

# Numerik IIB

— Numerische Verfahren für Differentialgleichungen —

Bernhard Schmitt

Sommersemester 1999

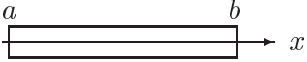
## Inhaltsverzeichnis

<b>1</b>	<b>Gewöhnliche Differentialgleichungen</b>	<b>5</b>
1.1	Theoretische Grundlagen . . . . .	5
1.2	Einschrittverfahren für Anfangswertprobleme . . . . .	9
1.2.1	Herleitung . . . . .	9
1.2.2	Konsistenz . . . . .	12
1.2.3	Stabilität . . . . .	13
1.2.4	Konvergenz . . . . .	14
1.2.5	Schrittweitensteuerung . . . . .	15
1.3	Mehrschrittverfahren . . . . .	21
1.3.1	Adams-Verfahren . . . . .	21
1.3.2	Lineare Mehrschrittverfahren und Stabilität . . . . .	24
1.4	Extrapolationsverfahren . . . . .	27
1.5	Steife Anfangswertprobleme . . . . .	30
1.6	Schießverfahren für Randwertprobleme . . . . .	36
1.7	Differenzenverfahren für Randwertprobleme . . . . .	42
<b>2</b>	<b>Gemeinsame Prinzipien von Diskretisierungsverfahren</b>	<b>46</b>
2.1	Konvergenz, Konsistenz, Stabilität . . . . .	46

<i>INHALTSVERZEICHNIS</i>	2
2.2 Allgemeine Verfahrens-Ansätze . . . . .	49
<b>3 Partielle Differentialgleichungen</b>	<b>52</b>
3.1 Allgemeine Eigenschaften . . . . .	52
3.2 Differenzenverfahren für elliptische Randwertprobleme . . . . .	54
3.2.1 Die Poissongleichung auf einfachen Gebieten . . . . .	54
3.2.2 Eigenschaften der Differenzen-Matrizen . . . . .	59
3.2.3 Diskretisierungsfehler, Konvergenz . . . . .	62
3.2.4 Allgemeinere Gebiete und Gleichungen . . . . .	66
3.3 Finite-Elemente-Verfahren für elliptische Probleme . . . . .	69
3.3.1 Variationsformulierung . . . . .	69
3.3.2 Rayleigh-Ritz-Galerkin-Verfahren . . . . .	71
3.3.3 Spline-Räume, <i>Finite Elemente</i> . . . . .	73
3.4 Partielle Anfangs-Randwert-Probleme . . . . .	77
<b>Index</b>	<b>82</b>



Einige durch Differentialgleichungen beschriebene Problemstellungen

<p>Einfache Modelle physik., biolog., chem... Prozesse; z.B.,</p> <p>a) <i>Massenbewegung unter Gravitation.</i>  Vereinfachende Annahmen: Vakuum, Massen starr, punktförmig, konzentriert,...</p> <p>a1) Wurf auf Erde: Flugbahn <math>u(t) = \begin{pmatrix} x(t) \\ y(t) \end{pmatrix}</math>, Newtonsches Gesetz:  <math display="block">m u''(t) = m \begin{pmatrix} x''(t) \\ y''(t) \end{pmatrix} = -m \begin{pmatrix} 0 \\ g \end{pmatrix}</math>. Explizite Lösung Wurfparabel  'Anfangswertproblem': Start-Ort &amp;-Geschwindigkeit gegeben <math>\Rightarrow \exists_1</math> Lösung!  'Randwertproblem': Start- &amp; Zielpunkt gegeben, Lösung?</p> <p>a2) Astronomie (Raumflug, <math>n</math>-Körper-Problem): nicht explizit lösbar, numerische Verfahren erforderlich (hoher Genauigkeit)</p> <p>b) <i>Wachsende Population in Nährmedium</i> (Bakterien, Protozoen,...)  b1) Kleiner Behälter, gute Durchmischung:  Population <math>p(t)</math> zur Zeit <math>t</math> wächst nach Gesetz  <math display="block">\frac{dp}{dt} = \alpha p - \beta p^2 = \alpha \left(1 - \frac{\beta}{\alpha} p\right) p</math>  Dabei: <math>\alpha</math> = Vermehrungsrate, <math>\beta p^2</math> = Wettbewerb. Lösung: logistische Kurve</p>	<p>Gewöhnliche Dgln</p>
<p>b2) Rohrförmiger Behälter, <math>p</math> abhängig von <math>x</math> und <math>t</math>: </p> <p>Wanderung (Diffusion) im Rohr prop. zu <math>p_{xx}</math>, für <math>p(t, x)</math> gilt  <math display="block">\frac{\partial p}{\partial t} = \alpha p - \beta p^2 + \gamma \frac{\partial^2 p}{\partial x^2}</math>  Zusatzinformation: Startpopulation <math>p(0, x)</math>, Rohrende: <math>p_x(t, a) = p_x(t, b) = 0</math>, d.h., Anfangswerte bzgl. <math>t</math>, Randwerte bzgl. <math>x</math>.  &lt;&lt;&lt;&lt; _____ &gt;&gt;&gt;&gt;</p> <p>Partielle Dgln 2.Ordnung, Einteilung nach Typen (hier in Orts-Gebiet <math>D \subset \mathbb{R}^2</math>)</p> <p>a) Parabolische Dgln: Ausgleichsvorgänge <math>u(t, x, y)</math>, z.B., Wärmeleitungs-Gl.  <math display="block">u_t = \alpha(u_{xx} + u_{yy})</math>  Sinnvolle Aufgabe: Anfangswerte in <math>t = 0</math>, Randwerte auf <math>\partial D</math></p> <p>b) Elliptische Dgln: Gleichgewichtszustände <math>u(x, y)</math>, z.B., Potentialgleichung  <math display="block">u_{xx} + u_{yy} = 0</math>  Sinnvolle Aufgabe: Randwerte auf <math>\partial D</math></p> <p>c) Hyperbolische Dgln: Schwingungen/Wellen <math>u(t, x, y)</math>, z.B., Wellen-Gl.  <math display="block">u_{tt} = \alpha(u_{xx} + u_{yy})</math>  Sinnvolle Aufgabe: Anfangswerte <math>u, u_t</math> in <math>t = 0</math>, Randwerte auf <math>\partial D</math></p>	<p>Partielle Dgln</p>

# 1 Gewöhnliche Differentialgleichungen

## 1.1 Theoretische Grundlagen

Im allgemeinsten Fall treten gewöhnliche Differentialgleichungen als Systeme von nichtlinearen Gleichungen auf, in denen verschiedene Ableitungen einer Vektor-Funktion  $u : \mathbf{R} \rightarrow \mathbf{R}^n$  miteinander verknüpft werden. Meist betrachtet man aber zur Standardisierung explizite Gleichungen der Form

$$u'(x) = \frac{du}{dx} = f(x, u(x)), \quad \text{kurz} \quad u' = f(x, u), \quad (1.1.1)$$

wobei  $f : \mathbf{R} \times \mathbf{R}^n \rightarrow \mathbf{R}^n$  eine gegebene stetige Funktion ist. Als Lösung einer solchen Differentialgleichung (Dgl) ist eine stetig differenzierbare Vektorfunktion  $u : [a, b] \rightarrow \mathbf{R}^n$ , also eine Raumkurve gesucht, die für jeden  $x$ -Wert aus  $[a, b]$  im Punkt  $u(x)$  den Ableitungswert  $f(x, u(x))$  annimmt, sich also in das durch  $f(x, y)$ ,  $x \in [a, b]$ ,  $y \in \mathbf{R}^n$ , gegebene *Richtungsfeld* einpaßt.

Die ausführliche Bezeichnung für (1.1.1) ist *System von Dgln erster Ordnung*, da nur die erste Ableitung  $u'$  der gesuchten Funktion auftaucht. Allerdings kann jede Dgl  $m$ -ter Ordnung,

$$z^{(m)}(x) = g(x, z(x), z'(x), \dots, z^{(m-1)}(x))$$

für eine Funktion  $z : \mathbf{R} \rightarrow \mathbf{R}$  (zur Vereinfachung) durch Einführung von Hilfsfunktionen

$$u_1(x) := z(x), \quad u_2(x) := z'(x), \quad \dots, \quad u_m(x) := z^{(m-1)}(x),$$

in die Standardform (1.1.1) mit  $n = m$  Komponenten gebracht werden, denn dann ist

$$u' = \begin{pmatrix} u_1' \\ \vdots \\ u_{m-1}' \\ u_m' \end{pmatrix} = \begin{pmatrix} u_2 \\ \vdots \\ u_m \\ g(x, u_1, \dots, u_m) \end{pmatrix}.$$

Aus diesem Grund werden im folgenden meist Systeme erster Ordnung behandelt, und nur in wenigen Fällen der häufig auftretende Fall von Dgln zweiter Ordnung.

Die Dgl (1.1.1) definiert unter geeigneten Voraussetzungen durch jeden Punkt  $(x, y) \in \mathbf{R} \times \mathbf{R}^n$  eine eigene Lösungskurve. Zur Festlegung einer eindeutigen Lösung sind weitere Bedingungen erforderlich. Dazu werden zwei Standardfälle behandelt. Beim *Anfangswertproblem* (AWP) wird der Funktionswert an einer Stelle  $x_0 \in [a, b]$  vollständig vorgeschrieben (im folgenden oBdA  $x_0 = a$ )

$$u(a) = u_0 \in \mathbf{R}^n. \quad (1.1.2)$$

Für die Lösung interessiert man sich dann in einem (evtl. unbeschränkten) Intervall  $[a, b]$ . Für dieses Problem gibt es recht allgemeine Existenz- und Eindeutigkeitsaussagen. Sowohl aus theoretischer als auch praktischer Sicht schwieriger ist das *Randwertproblem* (RWP), wo man Lösungen sucht, für die

eine bestimmte Beziehung der Funktionswerte am Anfangs- und Endpunkt eines endlichen Intervalls  $[a, b]$  gefordert wird, etwa in Form eines zusätzlichen (evtl. nichtlinearen) Gleichungssystems

$$r(u(a), u(b)) = 0, \quad r : \mathbf{R}^n \times \mathbf{R}^n \rightarrow \mathbf{R}^n. \quad (1.1.3)$$

Es folgen einige grundlegende theoretische Aussagen zu den genannten Problemen.

Für das Anfangswertproblem (1.1.1),(1.1.2) gibt es eine sehr einfache und allgemeine Existenzaussage im Satz von Picard-Lindelöf. Er beruht auf dem Übergang von der Dgl durch Integration zur äquivalenten *Integralgleichung*

$$u(x) = u_0 + \int_a^x f(t, u(t)) dt, \quad (1.1.4)$$

der auch für die Konstruktion numerischer Verfahren grundlegend ist.

**Satz 1.1.1** Die Funktion  $f$  sei auf dem Streifen  $\Omega := \{(x, y) : x \in [a, b], y \in \mathbf{R}^n\}$  über dem endlichen Intervall  $[a, b]$  definiert und stetig. Außerdem gelte dort die Lipschitzbedingung

$$\|f(x, u) - f(x, v)\| \leq L\|u - v\| \quad \forall x \in [a, b], u, v \in \mathbf{R}^n, \quad (1.1.5)$$

in einer  $\mathbf{R}^n$ -Norm  $\|\cdot\|$  mit einer Lipschitzkonstanten  $L \geq 0$ . Dann existiert zu jedem  $u_0 \in \mathbf{R}^n$  genau eine Lösung  $u \in (C^1[a, b])^n$  des Anfangswertproblems (1.1.1),(1.1.2). Diese Lösung hängt Lipschitz-stetig vom Anfangswert ab. Für zwei Lösungen  $u, v \in (C^1[a, b])^n$  mit

$$\begin{aligned} u' &= f(x, u), & v' &= f(x, v), \\ u(x_0) &= u_0, & v(x_0) &= v_0, \end{aligned}$$

gilt dabei die Schranke

$$\|u(x) - v(x)\| \leq e^{L|x-a|} \|u_0 - v_0\| \quad \forall x \in [a, b]. \quad (1.1.6)$$

**Beweis a)** Bei (1.1.4) handelt es sich um ein Fixpunktproblem  $u = Tu$  mit der Abbildung

$$T : \begin{cases} (C[a, b])^n \rightarrow (C[a, b])^n, \\ v \mapsto Tv, \end{cases} \quad \text{mit } (Tv)(x) := u_0 + \int_a^x f(t, v(t)) dt.$$

Für  $y \in (C[a, b])^n$  wird die gewichtete Norm

$$\|y\|_L := \max\{\|y(x)\| e^{-L|x-a|} : x \in [a, b]\}$$

eingeführt, mit der  $(C[a, b])^n$  ein Banachraum wird. Aus Voraussetzung (1.1.5) läßt sich direkt die Kontraktivität von  $T$  ableiten, denn für beliebige  $u, v \in (C[a, b])^n$  gilt

$$\begin{aligned} \|Tu - Tv\|_L &= \max_{x \in [a, b]} e^{-L|x-a|} \left\| \int_a^x (f(t, u(t)) - f(t, v(t))) dt \right\| \\ &\leq \max_{x \in [a, b]} e^{-L|x-a|} \int_a^x \|f(t, u(t)) - f(t, v(t))\| dt \\ &\leq L \max_{x \in [a, b]} e^{-L|x-a|} \int_a^x e^{L|t-a|} e^{-L|t-a|} \|u(t) - v(t)\| dt \\ &\leq L \max_{x \in [a, b]} e^{-L|x-a|} \int_a^x e^{L|t-a|} dt \|u - v\|_L \\ &\leq (1 - e^{-L|b-a|}) \|u - v\|_L. \end{aligned}$$

Da der Vorfaktor kleiner als eins ist, ist  $T$  eine Kontraktion und nach dem Banachschen Fixpunktsatz existiert ein eindeutiger Fixpunkt  $u = Tu \in (C[a, b])^n$ . Mit  $u$  ist auch die Funktion  $f(x, u(x)) = u'(x)$  stetig.

b) Aus der Differenz der beiden für  $u$  und  $v$  geltenden Gleichungen (1.1.4) folgt wie eben

$$\begin{aligned} \|u(x) - v(x)\| &= \|u_0 - v_0 + \int_a^x (f(t, u(t)) - f(t, v(t))) dt\| \\ &\leq \|u_0 - v_0\| + L \int_a^x \|u(t) - v(t)\| dt. \end{aligned}$$

Für die skalare Funktion  $\delta(x) := \|u(x) - v(x)\|$  gilt also die Ungleichung  $\delta(x) \leq \|u_0 - v_0\| + L \int_a^x \delta(t) dt$ . Mit dem folgenden Gronwall-Lemma ergibt sich daraus die Schranke (1.1.6). ■

Im Beweis wurde die folgende Aussage benutzt.

**Lemma 1.1.2 (Gronwall-Lemma)** Die Funktion  $\delta \in C[a, b]$  erfülle die Ungleichung

$$\delta(x) \leq \alpha + L \int_a^x \delta(t) dt \quad \forall x \in [a, b],$$

mit  $\alpha, L \geq 0$ . Dann folgt

$$\delta(x) \leq \alpha e^{L(x-a)} \quad \forall x \in [a, b].$$

**Beweis** Es sei  $\varepsilon > 0$ . Für die Funktion  $\gamma(x) := (\alpha + \varepsilon)e^{L(x-a)}$  gilt  $\gamma(x) = \alpha + \varepsilon + L \int_a^x \gamma(t) dt$ . Offensichtlich ist  $\delta(a) < \gamma(a)$ . Sei nun  $a < x_1 \leq b$  die erste Stelle mit  $\delta(x_1) = \gamma(x_1)$ , insbesondere gelte also  $\delta(x) \leq \gamma(x) \forall a \leq x \leq x_1$ . Dann folgt in  $x_1$  mit

$$\delta(x_1) \leq \alpha + L \int_a^{x_1} \delta(t) dt < \alpha + \varepsilon + L \int_a^{x_1} \gamma(t) dt = \gamma(x_1)$$

aber ein Widerspruch. Daher gilt  $\delta(x) < (\alpha + \varepsilon)e^{L(x-a)}$  für jedes  $\varepsilon > 0$ . ■

Besonders einfach ist das Studium von linearen Differentialgleichungen,

$$u'(x) = A(x)u(x) + g(x), \quad \text{mit } f(x, y) = A(x)y + g(x), \quad (1.1.7)$$

einer Matrixfunktion  $A(x) \in \mathbf{R}^{n \times n}$  und einer Vektorfunktion  $g(x) \in \mathbf{R}^n$ , da hier explizite Lösungsdarstellungen gelten. Mit Hilfe eines Fundamentalsystems der Dgl (1.1.7), d.h. einer regulären Matrixfunktion  $W(x) \in \mathbf{R}^{n \times n}$ , die die Dgl

$$W'(x) = A(x)W(x)$$

erfüllt, läßt sich die Lösung des Anfangswertproblems (1.1.7), (1.1.2) explizit angeben durch

$$u(x) = W(x)W(a)^{-1}u_0 + W(x) \int_a^x W(t)^{-1}g(t) dt.$$

Da mit  $W(x)$  auch  $W(x)M$  ( $M$  regulär) ein Fundamentalsystem ist, vereinfacht sich diese Darstellung bei Übergang zum ausgezeichneten System  $Y(x) := W(x)W(a)^{-1}$  mit  $Y(a) = I$ .

**Satz 1.1.3** Die Lösung des linearen Anfangswertproblems (1.1.7),(1.1.2) mit einer beschränkten, stetigen Matrixfunktion  $A(x)$  ist

$$u(x) = Y(x)u_0 + Y(x) \int_a^x Y(t)^{-1}g(t) dt, \quad (1.1.8)$$

wobei  $Y(x)$  das Fundamentalsystem ist mit  $Y' = AY$ ,  $Y(a) = I$ .

Bei konstanter Matrix  $A$  kann man das Fundamentalsystem explizit konstruieren. Über den Ansatz  $u(x) = e^{\lambda x}y$  läßt sich die homogene Dgl auf das Eigenwertproblem  $(A - \lambda I)y = 0$  zurückführen. Bei einer diagonalisierbaren Matrix  $A$  mit Eigen-Wert-/Vektor-Paaren  $(\lambda_j, y_j)$  folgt  $u(x) = \sum_{j=1}^n \alpha_j e^{\lambda_j x} y_j$ . Äquivalent dazu ist analog zum skalaren Fall ( $n = 1$ ) die Darstellung

$$Y(x) = e^{(x-a)A} := \sum_{k=0}^{\infty} \frac{(x-a)^k}{k!} A^k.$$

Dies vereinfacht die Lösungsformel (1.1.8) im Spezialfall weiter zu

$$u(x) = e^{(x-a)A}u_0 + \int_a^x e^{(x-t)A}g(t) dt.$$

Diese expliziten Formeln haben folgende Beziehung zu Satz 1.1.1. Die Lipschitzbedingung (1.1.5) gilt bei (1.1.7) mit  $L = \max_x \|A(x)\|$ , aus der für die Differenz von Lösungen folgt

$$\|u(x) - v(x)\| = \|Y(x)(u_0 - v_0)\| \leq e^{L|x-a|} \|u_0 - v_0\|,$$

also  $\|Y(x)\| \leq e^{L|x-a|}$ . Im Fall konstanter Koeffizienten reduziert sich dies auf die Ungleichung  $\|e^{tA}\| \leq e^{t\|A\|}$ ,  $t \geq 0$ , die sich auch direkt aus der Reihendarstellung ergibt.

Mit Hilfe der expliziten Lösungsdarstellungen soll kurz die Situation bei linearen Randwertproblemen diskutiert werden. Gesucht ist eine Lösung von

$$u'(x) = A(x)u(x) + g(x), \quad x \in (a, b), \quad R_0u(a) + R_1u(b) = d, \quad (1.1.9)$$

mit festen  $n \times n$ -Matrizen  $R_0, R_1$  und  $d \in \mathbf{R}^n$ . Dabei ist  $\text{rang}(R_0, R_1) = n$  vorzusetzen. Im Gegensatz zum Anfangswertproblem können hier alle Lösbarkeitsfälle auftreten. Die Lösung  $u$  wird durch die explizite Formel (1.1.8) mit einem unbekanntem Startvektor  $\eta \in \mathbf{R}^n$  dargestellt,

$$u(x) = Y(x)\eta + \gamma(x), \quad \gamma(x) := \int_a^x Y(x)Y(t)^{-1}g(t) dt.$$

Einsetzen in die Randbedingung führt auf das lineare Gleichungssystem

$$(R_0Y(a) + R_1Y(b))\eta = d - R_1\gamma(b) \quad (1.1.10)$$

für die Unbekannte  $\eta$ . Nach den Ergebnissen der Linearen Algebra treten dabei folgende Fälle auf:

$$\text{Wenn } R_0Y(a) + R_1Y(b) \begin{cases} \text{regulär:} & \text{es existiert eine eindeutige Lösung} \\ \text{singulär:} & \begin{cases} \text{es gibt keine Lösung} \\ \text{es gibt unendlich viele Lösungen} \end{cases} \end{cases}.$$

Bei Randwertproblemen kann es also keine zum Satz 1.1.1 vergleichbare allgemeine Existenz- oder Eindeutigkeitsaussage geben.

## 1.2 Einschrittverfahren für Anfangswertprobleme

### 1.2.1 Herleitung

Die Lösung der Dgl  $u' = f(x, u)$  ist diejenige Kurve  $(x, u(x))$ , die dem Richtungsfeld  $(x, y, f(x, y))$  folgt. Ein elementares Approximationsverfahren beruht darauf, daß man vom Startwert  $u_0 = u(x_0)$  'ein Stück weit' in Richtung der Ableitung  $u'(x_0) = f(x_0, u_0)$  geht und dann dort dieses Verfahren wiederholt (Eulerscher Polygonzug, s.u.). Einen allgemeinen Zugang zu solchen schrittweise vorgehenden numerischen Integrationsverfahren bekommt man über die Formel (1.1.4). Dazu wird im Punkt  $\bar{x}$  eine Schrittweite  $h > 0$  gewählt und die Dgl von  $\bar{x}$  bis  $\bar{x} + h$  integriert. Dies ergibt

$$\frac{1}{h}[u(\bar{x} + h) - u(\bar{x})] = \frac{1}{h} \int_{\bar{x}}^{\bar{x}+h} f(t, u(t)) dt = \int_0^1 f(\bar{x} + hs, u(\bar{x} + hs)) ds. \quad (1.2.1)$$

Zur Konstruktion numerischer Verfahren wird das Integral auf der rechten Seite approximiert, z.B., durch eine geeignete Quadraturformel. In diesem Abschnitt werden Näherungen der Form

$$\int_0^1 f(\bar{x} + hs, u(\bar{x} + hs)) ds = \underbrace{f_h(\bar{x}, u(\bar{x}))}_{\text{Näherung}} + \underbrace{T_h(\bar{x})}_{\text{Fehler}}. \quad (1.2.2)$$

betrachtet, in denen die Verfahrensfunktion  $f_h$  nur vom Wert  $(\bar{x}, u(\bar{x}))$  abhängt. Im Verfahren selbst wird der lokale Fehler  $T_h$  vernachlässigt, die Differentialgleichung (1.1.1) wird durch eine *Differenzgleichung* ersetzt. Somit ergeben sich Näherungswerte  $u_\Delta(x)$  an  $u(x)$  aus der Formel

$$\frac{1}{h}[u_\Delta(\bar{x} + h) - u_\Delta(\bar{x})] := f_h(\bar{x}, u_\Delta(\bar{x}))$$

nur an endlich vielen Stellen  $x_i$ , also auf einem *Gitter*  $\Delta$ .

**Definition 1.2.1** Ein Einschrittverfahren zur Lösung des Anfangswertproblems (1.1.1), (1.1.2) besteht aus der Wahl einer Verfahrensfunktion  $f_h(x, y)$  und eines Gitters

$$\Delta : a := x_0 < x_1 < \dots < x_N = b, h_i := x_{i+1} - x_i, H := \max_{i=0}^{N-1} h_i, |\Delta| := H.$$

Damit berechnet sich die Näherungslösung  $u_\Delta(x_i) = y_i$ ,  $i = 0, \dots, N$ , auf dem Gitter rekursiv

$$\begin{aligned} y_0 &:= u_0, \\ y_{i+1} &:= y_i + h_i f_{h_i}(x_i, y_i), \quad i = 0, \dots, N-1. \end{aligned} \quad (1.2.3)$$

Die Bezeichnung Einschrittverfahren bezieht sich auf die Tatsache, daß in (1.2.3) nur Näherungswerte  $y_i, y_{i+1}$  aus einem Schritt  $x_i \rightarrow x_{i+1}$  verknüpft werden. Es folgen vier einfache Beispiele.

**Verfahren 1**, Euler-Cauchy-Polygonzug: Das Integral in (1.2.2) wird durch den Funktionswert am linken Rand ersetzt,

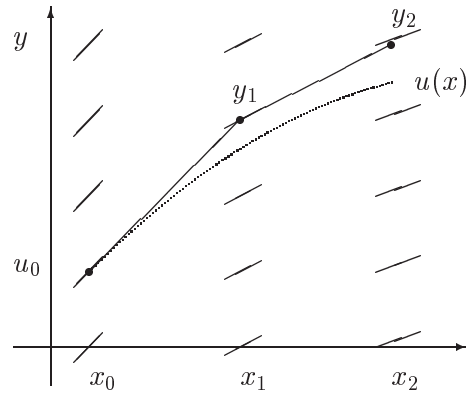
$$\int_0^1 f(\bar{x} + hs, u(\bar{x} + hs)) ds \cong f(\bar{x}, u(\bar{x})) =: f_h(\bar{x}, u(\bar{x})),$$

das Verfahren lautet also

$$y_{i+1} := y_i + h_i f(x_i, y_i),$$

$i = 0, 1, \dots$

Klartext: 'In  $x_i$  geht man einen Schritt der Länge  $h_i$  in Richtung des Richtungsfeldes'



**Verfahren 2**, von Runge: Integralapproximation durch die Rechteckregel,

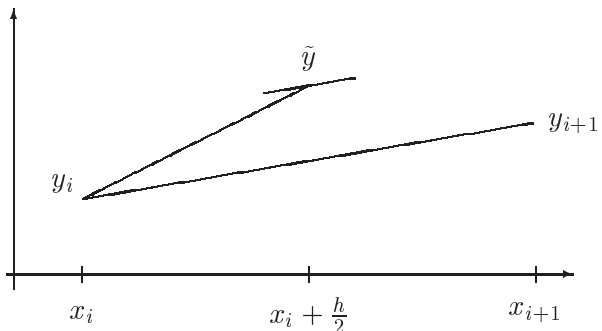
$$\int_0^1 f(\bar{x} + hs, u(\bar{x} + hs)) ds \cong f\left(\bar{x} + \frac{h}{2}, u\left(\bar{x} + \frac{h}{2}\right)\right).$$

Der unbekannte Wert  $u(\bar{x} + \frac{h}{2})$  wird durch Verfahren 1 angenähert,

$$u\left(\bar{x} + \frac{h}{2}\right) \cong u(\bar{x}) + \frac{h}{2} f(\bar{x}, u(\bar{x})) =: \tilde{y},$$

die Verfahrensfunktion ist hier also

$$f_h(x, y) := f\left(x + \frac{h}{2}, y + \frac{h}{2} f(x, y)\right).$$



**Verfahren 3**, von Heun: Integralapproximation durch die Trapezregel,

$$\int_0^1 f(\bar{x} + hs, u(\bar{x} + hs)) ds \cong \frac{1}{2} [f(\bar{x}, u(\bar{x})) + f(\bar{x} + h, u(\bar{x} + h))].$$

Wieder wird  $u(\bar{x} + h)$  durch Verfahren 1 angenähert. Jetzt ist

$$f_h(x, y) := \frac{1}{2} [f(x, y) + f(x + h, y + hf(x, y))].$$

**Verfahren 4**, klassisches Runge-Kutta-Verfahren: Integralapproximation durch eine modifizierte Simpsonregel,

$$f_h(x, y) := \frac{1}{6} [k_1 + 2k_2 + 2k_3 + k_4], \quad \text{mit}$$

$$k_1 := f(x, y) \cong u'(x)$$

$$k_2 := f\left(x + \frac{h}{2}, y + \frac{h}{2}k_1\right) \cong u'\left(x + \frac{h}{2}\right)$$

$$k_3 := f\left(x + \frac{h}{2}, y + \frac{h}{2}k_2\right) \cong u'\left(x + \frac{h}{2}\right)$$

$$k_4 := f(x + h, y + hk_3) \cong u'(x + h)$$

**Beispiel 1.2.2** Dgl  $u' = u =: f(x, u)$ ,  $x_0 = 0, u_0 = 1 \Rightarrow u(x) = e^x$ .

Verfahren	Näherung	lokal.Fehler
1	$y_1 = y_0 + hy_0 = 1 + h$	$h^2/2$
2	$y_{1/2} = 1 + \frac{h}{2}, y_1 = 1 + h(1 + \frac{h}{2}) = 1 + h + \frac{1}{2}h^2$	$h^3/6$
3	$\tilde{y}_1 = 1 + h, y_1 = 1 + \frac{h}{2}(1 + 1 + h) = 1 + h + \frac{1}{2}h^2$	$h^3/6$
4	$k_1 = 1, k_2 = 1 + \frac{h}{2}, k_3 = 1 + \frac{h}{2} + \frac{h^2}{4}, k_4 = 1 + h + \frac{h^2}{2} + \frac{h^3}{4} \Rightarrow$ $y_1 = 1 + h + \frac{1}{2}h^2 + \frac{1}{6}h^3 + \frac{1}{24}h^4$	$h^5/120$

Alle Verfahren erzeugen also den Anfang der Exponentialreihe unterschiedlich weit. Sie sind Spezialfälle einer allgemeinen Verfahrensklasse, der *Runge-Kutta-Verfahren*. Diese bestehen aus  $m \in \mathbb{N}$  Stufen und werden durch eine Tabelle von Koeffizienten definiert,

$$\begin{array}{c|c} c & A \\ \hline b & \end{array} = \begin{array}{c|ccc} c_1 & a_{11} & \cdots & a_{1m} \\ \vdots & \vdots & & \vdots \\ c_m & a_{m1} & \cdots & a_{mm} \\ \hline & b_1 & \cdots & b_m \end{array}$$

wobei  $c_i = \sum_{j=1}^m a_{ij}, i = 1, \dots, m$  gilt. Die Verfahrensfunktion beim allgemeinen Runge-Kutta-Verfahren ist damit

$$f_h(x, y) := \sum_{i=1}^m b_i k_i, \quad \text{mit} \tag{1.2.4}$$

$$k_i := f(x + hc_i, y + h \sum_{j=1}^m a_{ij} k_j), \quad i = 1, \dots, m.$$

Wenn die Koeffizientenmatrix  $A$  strikt untere Dreiecksgestalt besitzt, also für  $a_{ij} = 0 \forall j \geq i$ , können die Hilfsgrößen  $k_i$  der Reihe nach explizit aus diesen Gleichungen berechnet werden. Dann liegt ein *explizites Runge-Kutta-Verfahren* vor. Die Verfahren 1–4 gehören zu dieser Klasse, das klassische Runge-Kutta-Verfahren 4 etwa hat die Tabelle

0	0	.	.	.
$\frac{1}{2}$	$\frac{1}{2}$	0	.	.
$\frac{1}{2}$	0	$\frac{1}{2}$	0	.
1	0	0	1	0
	$\frac{1}{6}$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{6}$

Im allgemeinen Fall ( $\exists a_{ij} \neq 0, j \geq i$ ) ist (1.2.4) ein *implizites Runge-Kutta-Verfahren*, da dann zur Berechnung der  $k_i$  i.a. nichtlineare Gleichungssysteme zu lösen sind. Solche Verfahren sind aufwendig, können aber sehr gute Eigenschaften besitzen und werden daher später noch einmal betrachtet.

Die Güte der Verfahrensnäherungen  $y_i$  hängt vor allem von den Schrittweiten  $h_i$  ab. Eine Konvergenz dieser Näherungswerte bei den numerischen Verfahren (1.2.3) kann nur für feiner werdende Gitter erwartet werden in dem Sinn  $u_\Delta(x) \rightarrow u(x), |\Delta| \rightarrow 0, (x \text{ fest})$ . Der Nachweis dieser Tatsache erfolgt nach einem sehr allgemeinen Prinzip zweistufig in der Form

Konsistenz + Stabilität $\implies$ Konvergenz
---

Diese Begriffe werden im folgenden erläutert.

### 1.2.2 Konsistenz

Konsistenz besagt, daß die Verfahrensfunktion  $f_h$  so gewählt ist, daß die exakten Lösungen  $u$  der Differentialgleichung auch die Differenzengleichung (1.2.3) näherungsweise erfüllen.

**Definition 1.2.3** Das Verfahren (1.2.3) heißt konsistent, wenn mit jeder Lösung  $u$  der Dgl (1.1.1) bei stetig differenzierbarer rechter Seite  $f$  für den lokalen Fehler

$$T_h(x) := \frac{1}{h}[u(x+h) - u(x)] - f_h(x, u(x)), \quad x \in [a, b-h], \quad (1.2.5)$$

gilt  $T_h \rightarrow 0$  für  $h \rightarrow 0$ . Das Verfahren heißt konsistent mit Ordnung  $p > 0$ , wenn

$$T_h(x) = \mathcal{O}(h^p), \quad h \rightarrow 0, \quad x \in [a, b], \quad (1.2.6)$$

für genügend oft differenzierbare rechte Seiten  $f$ .

**Beispiel 1.2.4** Verfahren 2,  $f_h = f(x + \frac{h}{2}, y + \frac{h}{2}f(x, y))$ . Die Taylorentwicklung einer Lösung  $u$  in  $x$  ist

$$\frac{1}{h}[u(x+h) - u(x)] = u'(x) + \frac{h}{2}u''(x) + \mathcal{O}(h^2).$$

Für  $f \in C^2$  erhält man die Ableitungen der Lösung aus der Dgl,  $u' \equiv f(x, u)$ ,  $u'' \equiv f_x + f_y u' = f_x + f_y f$ . Daher stimmt die Taylorentwicklung von  $f_h$  (nach  $h$ ),

$$f(x + \frac{h}{2}, u + \frac{h}{2}f) = f(x, u) + \frac{h}{2}(f_x + f_y f) + \mathcal{O}(h^2)$$

mit der obigen in den ersten beiden Gliedern überein, das Verfahren hat also Konsistenzordnung  $p = 2$ .

Konsistenznachweise sind zwar elementar durchführbar, aber sehr arbeitsaufwendig. Standardmethode ist dabei die im Beispiel angewendete Taylorentwicklung. Umgekehrt können durch Vergleich der Taylorentwicklungen von  $u$  und  $f_h$  in (1.2.4) die erforderlichen Bedingungen an die Koeffizienten hergeleitet und damit geeignete Verfahren konstruiert werden. Für ein dreistufiges Verfahren mit Ordnung drei sind, z.B., folgende *Ordnungsbedingungen* zu erfüllen:

$$\begin{aligned} b_1 + b_2 + b_3 &= 1 \\ b_2 c_2 + b_3 c_3 &= \frac{1}{2} \\ b_2 c_2^2 + b_3 c_3^2 &= \frac{1}{3} \\ b_3 a_{32} c_2 &= \frac{1}{6} \end{aligned} \quad (1.2.7)$$

Dabei entsprechen die ersten drei Bedingungen denen an die Quadraturformel mit Stützstellen  $c_i$  und Gewichten  $b_i$  (Parabeln werden exakt integriert). Da das Verfahren 6 freie Parameter besitzt, gibt es also eine zweiparametrische Schar von dreistufigen Verfahren mit Ordnung drei. Analoges gilt bei 4-stufigen Verfahren (8 Bedingungen für 10 Parameter). Keines dieser Verfahren hat aber eine höhere Ordnung als drei bzw. vier. Die Zahl der Ordnungsbedingungen wächst zu höheren Ordnungen stark an, sodaß dann (im Unterschied zu den Quadraturformeln)  $m > p$  Stufen für Ordnung  $p$  erforderlich werden. Die mit

einer bestimmten Stufenzahl  $m$  erreichbare Ordnung  $p$  ist in der folgenden Tabelle aufgeführt:

$m =$	1	2	3	4	5	6	7	8	9	10	11	..	17
$p =$	1	2	3	4	4	5	6	6	7	7	8	..	10

Die vorher behandelten Verfahren 1–4 besitzen jeweils die höchsten erreichbaren Ordnungen. Die höchste bisher durch explizite Konstruktion erreichte Ordnung ist  $p = 10$  mit  $m = 17$  Stufen. Später werden zwar Einschrittverfahren beliebig hoher Ordnung auftreten, doch werden bei diesen noch sehr viel höhere Stufenzahlen benutzt.

Um die Konvergenz eines Verfahrens bzw. der durch es berechneten Näherungen nachzuweisen, reicht die Konsistenz i.a. nicht aus. Das Verfahren muß eine weitere zentrale Eigenschaft erfüllen.

### 1.2.3 Stabilität

Stabilität bedeutet im allgemeinen, daß kleine Störungen eines Verfahrens die berechneten Näherungen nur wenig (Lipschitz-stetig) verfälschen. Beim diskreten Anfangswertproblem (1.2.3) bietet sich wegen der einfachen Struktur aber eine speziellere Formulierung an.

**Definition 1.2.5** Ein Einschrittverfahren (1.2.3) heißt stabil, wenn Konstanten  $S, r, H_0$  existieren so, daß für je zwei gestörte Lösungen  $\{y_i\}, \{z_i\}$  mit

$$\left. \begin{aligned} \frac{1}{h_i}[y_{i+1} - y_i] - f_{h_i}(x_i, y_i) &= \eta_i \\ \frac{1}{h_i}[z_{i+1} - z_i] - f_{h_i}(x_i, z_i) &= \zeta_i \end{aligned} \right\} i = 0, 1, \dots, N-1, \quad (1.2.8)$$

gilt

$$\max_{i=0}^N \|y_i - z_i\| \leq S \left( \|y_0 - z_0\| + \max_{i=0}^{N-1} \|\eta_i - \zeta_i\| \right) \quad (1.2.9)$$

für alle Gitter  $\Delta$  mit  $|\Delta| \leq H_0$  und solange  $\|y_0 - z_0\| \leq r, \|\eta_i - \zeta_i\| \leq r, i = 0, \dots, N-1$ .

Im Gegensatz zu den später behandelten Mehrschrittverfahren gilt die Stabilität von Einschrittverfahren unter einfachen Voraussetzungen. In dieser Beziehung und auch in Bezug auf die verwendeten Beweistechniken liegen starke Ähnlichkeiten zu den Aussagen aus §1.1 vor. Umgekehrt kann übrigens auch der Existenzsatz 1.1.1 mit Hilfe des Euler-Cauchy-Verfahrens bewiesen werden. Als Beweishilfsmittel im folgenden Stabilitätssatz wird das diskrete Analogon zum Gronwall-Lemma benötigt.

**Lemma 1.2.6** Für Werte  $\delta_i \in \mathbf{R}_+$  gelte mit  $\alpha, \ell, h_i \geq 0$  rekursiv

$$\delta_{i+1} \leq (1 + h_i \ell) \delta_i + h_i \alpha, \quad i = 0, \dots, N-1.$$

Dann folgt (mit  $x_i := x_0 + h_0 + \dots + h_{i-1}$ ) die explizite Schranke

$$\delta_i \leq e^{\ell(x_i - x_0)} \left( \delta_0 + (x_i - x_0) \alpha \right), \quad i = 0, \dots, N.$$

**Beweis** Da  $1 + \ell h_i \leq e^{\ell h_i}$  ist und  $\delta_0$  die Behauptung erfüllt, folgt induktiv

$$\begin{aligned} \delta_{i+1} &\leq e^{\ell h_i} \delta_i + h_i \alpha \leq e^{\ell h_i} \left( e^{\ell(x_i - x_0)} [\delta_0 + (x_i - x_0) \alpha] \right) + h_i \alpha \\ &\leq e^{\ell(x_{i+1} - x_0)} [\delta_0 + (x_i - x_0) \alpha + h_i \alpha]. \quad \blacksquare \end{aligned}$$

Als wesentliche Voraussetzung wird beim Stabilitätsbeweis von Einschrittverfahren wieder eine Lipschitzbedingung benötigt.

**Satz 1.2.7 (Stabilitätssatz)** Die das Einschrittverfahren (1.2.3) definierende Funktion  $f_h$  erfülle in dem in Satz 1.1.1 definierten Streifen  $\Omega$  und für  $0 < h \leq H_0 \leq b - a$  eine Lipschitzbedingung

$$\|f_h(x, y) - f_h(x, z)\| \leq \ell \|y - z\| \quad \forall x \in [a, b], y, z \in \mathbf{R}^n. \quad (1.2.10)$$

Dann ist das Verfahren stabil, es gilt (1.2.9) mit der (von  $\Delta$  unabhängigen) Konstanten  $S = e^{(b-a)\ell} \max\{1, b-a\}$  sowie  $r = \infty$ .

**Beweis** Die Subtraktion der beiden Gleichungen in (1.2.8) ergibt

$$y_{i+1} - z_{i+1} = y_i - z_i + h_i [f_{h_i}(x_i, y_i) - f_{h_i}(x_i, z_i)] + h_i [\eta_i - \zeta_i].$$

Mit der Lipschitzbedingung folgt daraus die Ungleichungskette

$$\|y_{i+1} - z_{i+1}\| \leq \|y_i - z_i\| + h_i \ell \|y_i - z_i\| + h_i \|\eta_i - \zeta_i\| = (1 + h_i \ell) \|y_i - z_i\| + h_i \|\eta_i - \zeta_i\|.$$

Für die Differenzen  $\delta_i := \|y_i - z_i\|$  entspricht dies den Voraussetzungen von Lemma 1.2.6 mit  $\alpha = \max_j \|\eta_j - \zeta_j\|$  und liefert mit  $(x_i - x_0) e^{\ell(x_i - x_0)} \leq (b - a) e^{\ell(b-a)}$  die Behauptung.  $\blacksquare$

Die im Satz benötigte Lipschitzbedingung der Verfahrensfunktion  $f_h$  folgt bei expliziten Runge-Kutta-Verfahren aus derjenigen von  $f$ . Beim Euler-Cauchy-Verfahren ist dies trivialerweise der Fall, elementar ist auch

**Beispiel 1.2.8** Verfahren 2,  $f_h(x, y) = f(y + \frac{h}{2} f(y))$ , die Abhängigkeit von  $x$  ist unwesentlich. Aus der Bedingung  $\|f(y) - f(z)\| \leq L \|y - z\|$  folgt

$$\begin{aligned} \|f_h(y) - f_h(z)\| &= \left\| f\left(y + \frac{h}{2} f(y)\right) - f\left(z + \frac{h}{2} f(z)\right) \right\| \\ &\leq L \left\| y + \frac{h}{2} f(y) - z - \frac{h}{2} f(z) \right\| \leq L (\|y - z\| + \frac{h}{2} L \|y - z\|), \end{aligned}$$

d.h. (1.2.10) mit  $\ell = L + \frac{1}{2} H_0 L^2$ .

## 1.2.4 Konvergenz

Die doppelte Bezeichnung der Näherungen durch  $y_i$  und  $u_\Delta(x_i)$  wurde eingeführt, um einerseits die Schreibweise im Verfahren zu vereinfachen, andererseits aber auch die Konvergenzaussage klarer formulieren zu können. Konvergenz der Näherungslösungen kann natürlich nur bei feiner werdendem Gitter  $\Delta$ ,

$|\Delta| \rightarrow 0$ , also auch  $N \rightarrow \infty$  erwartet werden. Bezogen auf einen *festen* Punkt  $\bar{x}$  im Integrationsintervall, z.B.  $\bar{x} = b$ , betrifft daher die Aussage  $u_\Delta(\bar{x}) \rightarrow u(\bar{x})$  ( $|\Delta| \rightarrow 0$ ) Näherungswerte  $y$  mit wachsendem Index auf verschiedenen Gittern.

**Satz 1.2.9** *Ein konsistentes und stabiles Verfahren ist konvergent:*

Das Einschrittverfahren (1.2.3) sei konsistent mit der Ordnung  $p > 0$ , d.h. es gelte (1.2.6). Das Verfahren sei außerdem stabil, etwa aufgrund der Lipschitzbedingung (1.2.10) für  $f_h$ . Dann konvergieren die Näherungswerte  $\{y_i = u_\Delta(x_i)\}$  bei feiner werdendem Gitter  $\Delta$  ( $|\Delta| \rightarrow 0$ ) gegen die Lösung  $u$  des AWP's (1.1.1),(1.1.2). Die Konvergenzordnung ist mindestens  $p$ , d.h. es gibt eine Konstante  $K$  so, daß für alle  $|\Delta| \leq H_0$  gilt

$$\max_{x \in \Delta} \|u_\Delta(x) - u(x)\| = \max_{i=0}^{|\Delta|} \|y_i - u(x_i)\| \leq K|\Delta|^p. \quad (1.2.11)$$

**Beweis** Nach Definition erfüllen die Näherungen  $y_i$  die Beziehung

$$\frac{1}{h_i}(y_{i+1} - y_i) - f_{h_i}(x_i, y_i) = 0.$$

Für die Gitterwerte  $z_i := u(x_i)$  der Lösung der Dgl gilt nun die Konsistenzaussage (1.2.6)

$$\frac{1}{h_i}(z_{i+1} - z_i) - f_{h_i}(x_i, z_i) = T_{h_i}(x_i) = \mathcal{O}(h_i^p).$$

Diese Gleichungen entsprechen denen aus der Stabilitätsdefinition, (1.2.8), mit  $\eta_i = 0$  und  $\zeta_i = T_{h_i}(x_i)$ . Wegen  $y_0 = u_0$  und  $h_i^p \leq |\Delta|^p$  folgt daher die Schranke (1.2.11). ■

Der Beweis des Satzes ist natürlich trivial, nach erfolgter Begriffsbildung, da einfach zwei Definitionen zusammengefügt werden. Die eigentliche Arbeit liegt im Einzelnachweis von Konsistenz und Stabilität.

*Bemerkung:* Für stabile Verfahren gilt *Konvergenzordnung*  $\geq$  *Konsistenzordnung*.

Von theoretischer Bedeutung ist die Konvergenzaussage des Satzes, da sie zeigt, daß die Näherungslösung  $u_\Delta$  genügend schnell konvergiert. Bei einer Halbierung der Schrittweiten im Gitter kann, z.B., eine Verkleinerung des Fehlers um den Faktor  $2^{-p}$  erwartet werden. Dies legt die Verwendung von Verfahren möglichst hoher Ordnung nahe. Allerdings hängt die Fehlerkonstante  $K$  in (1.2.11) von immer höheren Ableitungen  $u^{(p)}$  ab und begrenzt dadurch die verwendbare Ordnung. Diese Konstante  $K$  ist nicht explizit bekannt und in der Praxis kaum berechenbar, daher kann aus der Fehlerschranke keine Angabe über die tatsächliche Größe des Fehlers  $u_\Delta$  bei gegebenem Gitter abgeleitet werden. Genausowenig läßt sich damit ein Gitter konstruieren, auf dem ein geforderter Fehler (möglichst effizient) eingehalten wird. Diese praktisch wichtige Frage wird jetzt behandelt.

### 1.2.5 Schrittweitensteuerung

In vielen realen Problemen kann die Gestalt der Lösungen in unterschiedlichen Bereichen des Integrationsintervalls  $[a, b]$  sehr unterschiedlich ausfallen (Beispiel Satellitenbahn: stark/schwach gekrümmt

nahe/entfernt von Himmelskörpern). Ein Verfahren kann hier nur dann effizient (und auch mit akzeptablen Rundungsfehlern, s.u.) arbeiten, wenn sich die Schrittweiten an diesen Verlauf anpassen. Durch eine Inspektion von Lemma 1.2.6 sieht man, daß die Fehlerschranke (1.2.11) eigentlich aus

$$\max_{x \in \Delta} \|u_{\Delta}(x) - u(x)\| \leq \tilde{K} \sum_{j=0}^{N-1} h_j \|T_{h_j}\|, \quad \tilde{K} = e^{\ell(b-a)} \quad (1.2.12)$$

hervorgeht. Ein Gesamtfehler der Größenordnung  $(b-a)\tilde{K}\varepsilon$  kann daher durch die Forderung

$$\|T_{h_i}(x_i)\| \stackrel{!}{\leq} \varepsilon \quad \forall i = 0, \dots, N-1 \quad (1.2.13)$$

erreicht werden, vgl. Numerik I, §5.4 (adaptive Integration). Da der exakte lokale Fehler  $T_h$  nicht bekannt ist, arbeitet man mit einer Schätzung für diesen und nutzt das bekannte Verhalten  $T_h \doteq \gamma h^p$  mit der Forderung (1.2.13) als Richtlinie zur lokalen, fortlaufenden Schrittweitensteuerung. Hat man nämlich nach einem mit Schrittweite  $h$  ausgeführten Integrationsschritt (eine Schätzung für) den aktuellen lokalen Fehler  $T_h$ , dann läßt sich die eigentlich für (1.2.13) erforderliche Schrittweite  $\hat{h}$  schätzen nach

$$T_{\hat{h}} \doteq \gamma \hat{h}^p \doteq \left(\frac{\hat{h}}{h}\right)^p T_h \stackrel{!}{=} \varepsilon \quad \Rightarrow \quad \hat{h} \doteq h \sqrt[p]{\frac{\varepsilon}{T_h}}.$$

Zur Schätzung des lokalen Fehlers können zwei Verfahren verschiedener Ordnung, etwa Ordnung  $p$  und  $p+1$ , verwendet werden, ihre Verfahrensfunktionen seien  $f_h$  und  $\bar{f}_h$ . Dazu berechnet man im  $i$ -ten Schritt zwei Näherungen

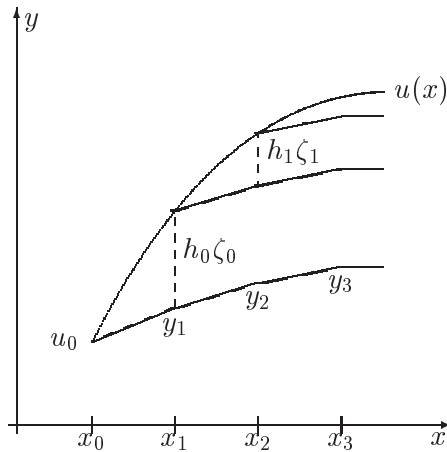
$$\left. \begin{aligned} y_{i+1} &= y_i + h f_h(x_i, y_i) \\ \bar{y}_{i+1} &= y_i + h \bar{f}_h(x_i, y_i) \end{aligned} \right\} \text{mit Fehler} \begin{cases} \mathcal{O}(h^p) \\ \mathcal{O}(h^{p+1}) \end{cases}. \quad (1.2.14)$$

Ist  $v_i(x)$  diejenige Lösung der Dgl, die die Anfangsbedingung  $v_i(x_i) = y_i$  erfüllt, dann gilt bezüglich dieser Lösung die Beziehung

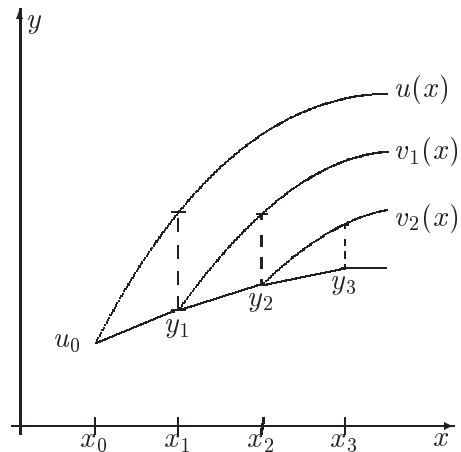
$$\begin{aligned} T_h(x_i) &= \frac{1}{h} [v_i(x_i + h) - y_i - h f_h(x_i, y_i)] = \frac{1}{h} [v_i(x_i + h) - y_{i+1}] \\ &= \frac{1}{h} [\bar{y}_{i+1} - y_{i+1}] + \mathcal{O}(h^{p+1}) = \bar{f}_h(x_i, y_i) - f_h(x_i, y_i) + \mathcal{O}(h^{p+1}). \end{aligned} \quad (1.2.15)$$

Daher ist (für genügend kleine Schrittweite) die Differenz  $\bar{f}_h(x_i, y_i) - f_h(x_i, y_i)$  eine Schätzung des lokalen Fehlers.

Die Unterschiede zwischen der Beweismethode zur Fehlerabschätzung in Satz 1.2.9 und der hier verwendeten zur Fehlerschätzung lassen sich graphisch mit den Lösungsfächern von Differentialgleichung (1.1.1) und Differenzgleichung (1.2.3) veranschaulichen. Die gestrichelten Fehlerbeiträge der beiden Ansätze sind dabei vergleichbar, aber nicht identisch.



In Satz 1.2.9: Lösungsfächer Differenzgl.  
mit  $\zeta_j = T_{h_j}(x_j)$ ,  $j = 0, 1, \dots$



In (1.2.15): Lösungsfächer Dgl

Mit einer solchen Schätzung des lokalen Fehlers kann auf sehr einfache Weise eine Schrittweitensteuerung formuliert werden:

**Algorithmus 1.2.10** Schrittweitensteuerung beim AWP:

```

x := a; y := u_0; h = ?, ε = ?
wiederhole:
  y0 := y + h f_h(x, y);
  y1 := y + h f̄_h(x, y);
  q := ε h / ||y1 - y0||;
  falls q > 1 dann {Schritt akzeptieren}
    [ x := x + h; y := y0; ]
  h := h * √[q];
  falls x + h > b dann h := b - x;
bis x ≥ b;

```

In der Praxis wird meist eine geringfügig kleinere Schrittweite  $h := 0.9 * h * \sqrt[q]}$  verwendet, da eine Schrittweiterholung wesentlich teurer wird als ein etwas zu kleiner Schritt. Außerdem ist es ratsam, einen zu großen Wechsel bei den Schrittweiten durch Zusatzabfragen zu verhindern. Zusätzlich wird meist auch mit der genaueren Lösung  $y_1$  weitergerechnet ( $y := y_1$  statt der unterstrichenen Anweisung  $y := y_0$ ), da die lokale Fehlerschätzung nur zur Steuerung verwendet wird. Zur Wahl der erforderlichen Toleranz  $\varepsilon$  für den lokalen Fehler  $T_h$  um einen bestimmten globalen Fehler  $\|u_\Delta - u\|$  einzuhalten kann allerdings keine allgemeine Aussage gemacht werden, da die Konstante  $\tilde{K}$  in (1.2.12) i.a. unbekannt ist.

Zur Konstruktion von Verfahrenspaaren (1.2.14) gibt es zwei übliche Zugänge unterschiedlicher Effizienz und Allgemeinheit:

a) *Eingebettete Runge-Kutta-Verfahren*: Analog zum Konstruktionsprinzip der Runge-Kutta-Verfahren versucht man, die gewünschte Näherung  $\bar{y}_{i+1}$  mit möglichst wenig Funktionsauswertungen zu berechnen. Dazu sucht man nach zwei Verfahren (1.2.4) mit  $m$  bzw.  $m + 1$  Stufen, bei denen das zweite die

Funktionswerte  $k_i$  des ersten mitverwendet,

$$\begin{aligned}
 y_{i+1} &= y_i + h_i \sum_{j=1}^m b_j k_j, \\
 \bar{y}_{i+1} &= y_i + h_i \sum_{j=1}^{m+1} \bar{b}_j k_j, \\
 k_j &= f\left(x + h_i c_j, y_i + h_i \sum_{l=1}^{j-1} a_{jl} k_l\right), \quad j = 1, \dots, m+1,
 \end{aligned}$$

$c_1$	0	...	0	0
$\vdots$		$\ddots$		$\vdots$
$c_{m+1}$	$a_{j1}$			0
	$b_1$	...	$b_m$	0
	$\bar{b}_1$	...	$\bar{b}_m$	$\bar{b}_{m+1}$

Die lokale Fehlerschätzung zum  $i$ -ten Schritt ist damit (und  $b_{m+1} := 0$ )

$$T_{h_i}(x_i) \cong \frac{1}{h_i} [\bar{y}_{i+1} - y_{i+1}] = \sum_{j=1}^{m+1} (\bar{b}_j - b_j) k_j.$$

Mit Hilfe der Ordnungsbedingungen (1.2.7) wird kurz die Konstruktion eines eingebetteten Verfahrens zu Verfahren 3 betrachtet. Hier sind  $c_2 = a_{21} = 1$  vorgegeben. Für die restlichen Parameter bleiben die Bedingungen

$$\bar{b}_1 + \bar{b}_2 + \bar{b}_3 = 1, \quad \bar{b}_2 + \bar{b}_3 c_3 = \frac{1}{2}, \quad \bar{b}_2 + \bar{b}_3 c_3^2 = \frac{1}{3}, \quad \bar{b}_3 a_{32} = \frac{1}{6}.$$

Für  $c_3 = \frac{1}{2}$  ergeben sich daraus die Koeffizienten in der folgenden rechten Tabelle. Die linke enthält ein (2,3)-Paar zu Verfahren 2.

(2,3)-Paar zu Verf. 2	(2,3)-Paar zu Verf. 3
0	0
$\frac{1}{2}$	1
1	$\frac{1}{2}$
	$\frac{1}{4}$
	$\frac{1}{4}$
	0
	$\frac{1}{2}$
	$\frac{1}{2}$
	0
	$\frac{1}{6}$
	$\frac{2}{3}$
	$\frac{1}{6}$

Bekannt sind Runge-Kutta-Paare mit

Ordnungen	(1,2)	(2,3)	(3,4)	(4,5)	(4,5)	..	(7,8)
Stufen	(1,2)	(2,3)	(3,5)	(4,6)	(5,6)	..	(11,13)

Bei den 3-stufigen Verfahren gibt es keine eingebetteten 4-stufigen Verfahren. Allerdings fand Fehlberg ein 5-stufiges Paar der Ordnung (3,4), bei dem der letzte Funktionswert  $f(x + c_5 h, \dots)$  als erster Funktionswert im nächsten Intervall,  $f(x + h + c_1 h, \dots)$ , verwendet wird, das also beinahe diese Anforderungen erfüllt. Auch zum klassischen Runge-Kutta-Verfahren gibt es kein eingebettetes mit 5 Stufen, da überhaupt keine 5-stufigen Verfahren mit Ordnung 5 existieren (s.o.). Somit werden sowieso 6 Stufen nötig und es ist günstiger, das Stufenpaar (5,6) zu verwenden. Sehr effiziente Verfahren (und Implementierungen) sind das (4,5)-Paar DOPRI5 bzw. (3,5,8)-Tripel DOP853 von Dormand und Prince.

In praktischen Anwendungen benötigt man öfters (z.B. Graphikausgabe) die Näherung  $y$  bzw.  $u_\Delta$  tatsächlich als Funktion von  $t$ . Zu diesem Zweck betrachtet man *stetige Erweiterungen*. Dabei sind die

Koeffizienten in (1.2.4) Polynome  $\bar{b}_i(t)$  mit  $\bar{b}_i(0) = 0$ ,  $\bar{b}_i(1) = b_i$ , und durch

$$u_{\Delta}(x_i + th_i) = y_i + h_i \sum_{j=1}^m b_j(t) k_j,$$

wird eine stetige Funktion  $u_{\Delta}$  definiert. Für das Verfahren 3 etwa ergibt sich mit

$$b_1(t) = t - \frac{1}{2}t^2, \quad b_2(t) = \frac{1}{2}t^2$$

eine stetige Erweiterung  $u_{\Delta}(x)$  der gleichmäßigen Konvergenzordnung 2. Im allgemeinen kann man aber für stetige Erweiterungen nicht die volle normale Ordnung der Endformel mit  $\bar{b}_i(1)$  erwarten. Für 3-stufige Verfahren ist die glm Ordnung 2 erreichbar, zu DOPRI5 existiert eine Erweiterung mit glm Ordnung 4.

b) *Richardson-Extrapolation*: Bei vielen Verfahren, z.B. auch den Runge-Kutta-Verfahren, existiert für äquidistante Gitter  $\Delta = \{x_i = a + ih\}$  eine asymptotische Entwicklung des Fehlers

$$u_{\Delta}(x) - u(x) = h^p g_p(x) + \mathcal{O}(h^{p+1}), \quad h \rightarrow 0, \quad (1.2.16)$$

vgl. Numerik I, §5.5. Durch Linearkombination von zwei Näherungen zu Schrittweiten  $h$  und  $h/2$  kann der dominierende Fehleranteil  $h^p g_p$  eliminiert ( $\rightarrow$  Richardson-Extrapolation) oder aber geschätzt werden, da für eine zweite Lösung  $u_{\Delta'}$ , zu  $\Delta' = \{x_i = a + i\frac{h}{2} : i = 0, \dots, 2N\}$ , gilt

$$\begin{aligned} u_{\Delta'}(x) - u(x) &= 2^{-p} h^p g_p(x) + \mathcal{O}(h^{p+1}) \quad \Rightarrow \\ (1 - 2^{-p}) h^p g_p(x) &= u_{\Delta}(x) - u_{\Delta'}(x) + \mathcal{O}(h^{p+1}). \end{aligned}$$

Daraus kann eine bessere Approximation  $\bar{u}_{\Delta}(x) := u_{\Delta'}(x) + \frac{1}{2^p - 1} [u_{\Delta'}(x) - u_{\Delta}(x)]$ ,  $x \in \Delta$ , höherer Ordnung oder die Fehlerschätzungen

$$u_{\Delta} - u \doteq \frac{2^p}{2^p - 1} [u_{\Delta} - u_{\Delta'}], \quad u_{\Delta'} - u \doteq \frac{1}{2^p - 1} [u_{\Delta} - u_{\Delta'}],$$

berechnet werden. Wie in (1.2.15) ergibt sich

$$\begin{aligned} T_h(x_i) &= \frac{1}{h} [v_i(x_i + h) - u_{\Delta}(x_i + h)] \doteq \frac{1}{h} [\bar{u}_{\Delta}(x_i + h) - u_{\Delta}(x_i + h)] \\ &= \frac{1}{1 - 2^{-p}} \frac{1}{h} [u_{\Delta'}(x_i + h) - u_{\Delta}(x_i + h)] \doteq -h^{p-1} g_p(x_i + h). \end{aligned} \quad (1.2.17)$$

Diese Größe ist im speziellen Fall von der Ordnung  $\mathcal{O}(h^p)$ . Da man sich hier in (1.2.16) auf (Näherungs-)Lösungen bezieht mit  $u_{\Delta}(x_i) = v_i(x_i) = y_i$ , gilt  $g_p(x_i) = 0$  und somit nach dem Mittelwertsatz  $h^{p-1} g_p(x_i + h) = \mathcal{O}(h^p)$ . Daher ist (1.2.17) eine Fehlerschätzung von der in Algorithmus 1.2.5 verwendeten Gestalt. Im Vergleich zu eingebetteten Verfahren werden hier aber wesentlich mehr Zusatzfunktionsauswertungen verwendet (bezogen auf die genauere Lösung  $u_{\Delta'}$  nach 2 Schritten mit  $h/2$  also  $m/2$  pro Schritt).

**Beispiel 1.2.11** Periodische Satellitenbahn um Erde und Mond: In der Ebene des rotierenden Erde-Mond-System bewegt sich ein Satellit mit den Koordinaten  $u(t) \in \mathbf{R}^2$  nach der Dgl

$$u'' = u + \begin{pmatrix} 0 & 2 \\ -2 & 0 \end{pmatrix} u' - \nu \frac{u - E}{\|u - E\|^3} - \mu \frac{u - M}{\|u - M\|^3},$$

wobei  $\mu = 1/82.45, \nu = 1 - \mu$  die Massenverhältnisse und  $E = (-\mu, 0)^T$  bzw.  $M = (\nu, 0)^T$  die Positionen von Erde und Mond sind. Es gibt eine periodische Bahn, die nahe an der Erde vorbeiführt und dort sehr kleine Schrittweiten erfordert. Ergebnis bei  $\varepsilon =_{10} -4$ :

	Anz.f-Werte	$h_{\min}$	$h_{\max}$	Endfehler
RuKuFe-2/3	35805	$2_{10}-6$	$8_{10}-2$	$2_{10}-5$
RuKuFe-4/5	972	$4_{10}-4$	$6_{10}-1$	$9_{10}-3$

Schrittweitensteuerung hatte zunächst nur das Ziel, die geforderte Genauigkeit der Näherung mit möglichst wenig Aufwand zu erreichen. Zusätzlich reduziert sie aber auch den Einfluß von *Rundungsfehlern*. Bei Maschinenrechnung werden in (1.2.3) gestörte Näherungen  $\tilde{y}_i$  berechnet, die pro Schritt anfallenden Rundungsfehler  $\varepsilon_i$  führen also zu einer verfälschten Rekursion

$$\tilde{y}_{i+1} = \tilde{y}_i + h_i f_{h_i}(x_i, \tilde{y}_i) + \varepsilon_{i+1}, \quad i = 0, \dots, N-1,$$

$\tilde{y}_0 = u_0$ . Daher kommen bei praktischer Rechnung zum Diskretisierungsfehler  $y_i - u(x_i)$  noch Rundungsfehler hinzu. Unter Voraussetzung (1.2.10) gilt dafür

$$\begin{aligned} \|\tilde{y}_{i+1} - y_{i+1}\| &\leq \|\tilde{y}_i - y_i\| + h_i \|f_{h_i}(x_i, \tilde{y}_i) - f_{h_i}(x_i, y_i)\| + \|\varepsilon_{i+1}\| \\ &\leq (1 + h_i \ell) \|\tilde{y}_i - y_i\| + \hat{\varepsilon}_{i+1}, \quad \hat{\varepsilon}_{i+1} := \|\varepsilon_{i+1}\|. \end{aligned}$$

Analog zu Lemma 1.2.6 ergibt sich mit  $\tilde{y}_0 = y_0$  die Abschätzung

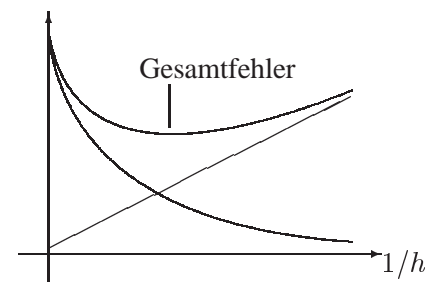
$$\|\tilde{y}_i - y_i\| \leq \sum_{j=1}^i e^{\ell(x_i - x_j)} \hat{\varepsilon}_j \leq S \sum_{j=1}^N \hat{\varepsilon}_j, \quad S := e^{\ell(b-a)}. \quad (1.2.18)$$

Bei einem äquidistanten Gitter führt dies für den Gesamtfehler auf die Schranke

$$\|\tilde{y}_i - u(x_i)\| \leq S(b-a) \frac{\varepsilon}{h} + Kh^p, \quad \varepsilon := \max_i \hat{\varepsilon}_i.$$

Bei äquidistantem Gitter muß die (zunächst unbekannte) Schrittweite  $h$  nach der ungünstigsten Stelle im Intervall gewählt werden und erzwingt daher evtl. eine sehr große Zahl  $N$  von Schritten. Da der Rundungsfehler nach (1.2.18) im wesentlichen von dieser Schrittzahl  $N$  abhängt, können

diese Fehler bei variablem Gitter reduziert werden, wenn man in glatten Passagen große, und daher insgesamt weniger, Schritte benutzt. Zusätzlich sei daran erinnert (vgl. Numerik I, §7.2), daß der kleinstmögliche Wert in der Schranke für  $\|\tilde{y}_i - u(x_i)\|$  von der Größe  $O(\varepsilon^{p/(p+1)})$  ist und mit wachsender Ordnung kleiner wird.



### 1.3 Mehrschrittverfahren

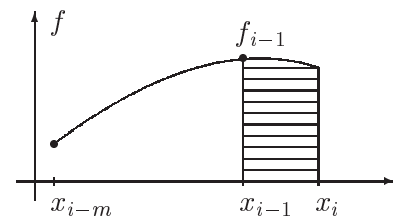
#### 1.3.1 Adams-Verfahren

Bei Einschrittverfahren der Ordnung  $p$  sind  $m \geq p$  Stufen erforderlich und daher zur Approximation des Integrals in (1.2.1),

$$\int_0^1 f(x_{i-1} + ht, u(x_{i-1} + ht)) dt = \frac{1}{h} [u(x_{i-1} + h) - u(x_{i-1})]$$

beim Schritt  $x_{i-1} \rightarrow x_i$  jeweils mindestens  $p$  Funktionswerte  $k_j \cong f(x_{i-1} + hc_j, u(x_{i-1} + hc_j))$  **neu** zu berechnen. Eine evtl. günstigere Alternative ist es, stattdessen schon bekannte, vor  $x_{i-1}$  gelegene Funktionswerte  $f(x_{i-j}, y_{i-j}) =: f_{i-j}, j = 1, 2, \dots, m$ , heranzuziehen.

Durch Integration des Interpolationspolynoms  $p_{m-1}(x)$  vom Grad  $m - 1$  zu den Funktionswerten  $f_{i-1}, \dots, f_{i-m}$ , das mit den Lagrangepolynomen  $L_j(x)$  zu den Stützstellen  $x_{i-1}, \dots, x_{i-m}$  dargestellt werden kann (vgl. Numerik I, §3.2), ergibt sich für dieses Integral die Approximation



$$f_h(x_{i-1}, y_{i-m}, \dots, y_{i-1}) := \sum_{j=-m}^{-1} f_{i+j} \underbrace{\int_0^1 L_j(x_{i-1} + ht) dt}_{=: \beta_{-j}}. \tag{1.3.1}$$

Diese Verfahrensfunktion hängt also von  $m$  früheren Lösungswerten  $y$  ab. Speziell für äquidistante Punkte (d.h.  $h_i \equiv h$ ) erhält man damit die *Methode von Adams-Bashforth* (A-B):

$$\left. \begin{aligned} \frac{1}{h}(y_i - y_{i-1}) &= f_h(x_{i-1}, y_{i-m}, \dots, y_{i-1}) \\ f_h(x_{i-1}, y_{i-m}, \dots, y_{i-1}) &:= \sum_{j=1}^m \beta_j f(x_{i-j}, y_{i-j}) \end{aligned} \right\} i = m, m + 1, \dots, N. \tag{1.3.2}$$

Dieses Verfahren ist ein explizites, lineares,  $m$ -Schrittverfahren. Pro Schritt fällt hier tatsächlich nur *eine* neue Funktionsauswertung,  $f_{i-1}$ , an. Explizit heißt das Verfahren wieder, da  $f_h$  nicht vom unbekanntem Wert  $y_i$  abhängt.

Der Konsistenzfehler entspricht dem Quadraturfehler bei der Funktion  $u'(x)$ , das Verfahren hat daher Ordnung  $p = m$ . Die Koeffizienten der ersten A-B-Verfahren enthält die Tabelle.

$m$	$j =$	1	2	3	4	Fehler
1	$\beta_j =$	1				$\frac{1}{2} h u''$
2	$2\beta_j =$	3	-1			$\frac{5}{12} h^2 u'''$
3	$12\beta_j =$	23	-16	5		$\frac{3}{8} h^3 u^{(4)}$
4	$24\beta_j =$	55	-59	37	-9	$\frac{251}{720} h^4 u^{(5)}$

Bei der Approximation (1.3.1) werden nur Daten auf einer Seite des Integrationsintervalls  $[x_{i-1}, x_i]$  benutzt. Ohne Erhöhung der Zahl  $m$  der Teilintervalle kann das A-B-Verfahren durch Berücksichtigung

des (noch unbekanntes) Wertes im Endpunkt  $x_i$  verbessert werden, man interpoliert also in den Stellen  $x_{i-m}, \dots, x_i$  mit einem Polynom vom nächsthöheren Grad  $m$ . Dies führt auf die *Methode von Adams-Moulton* (A-M)

$$\left. \begin{aligned} \frac{1}{h}(y_i - y_{i-1}) &= f_h(x_{i-1}, y_{i-m}, \dots, y_i) \\ f_h(x_{i-1}, y_{i-m}, \dots, y_i) &:= \sum_{j=0}^m \mu_j f(x_{i-j}, y_{i-j}) \end{aligned} \right\} i = m, m+1, \dots, N. \quad (1.3.3)$$

Hier liegt ein implizites, lineares,  $m$ -Schritt-Verfahren vor, da (1.3.3) für  $y_i$  ein (nichtlineares) Gleichungssystem darstellt. Die Koeffizienten der ersten A-M-Verfahren lauten

m	$j =$	0	1	2	3	Fehler
0	$\mu_j =$	1				$-\frac{1}{2}hu''$
1	$2\mu_j =$	1	1			$-\frac{1}{12}h^2u'''$
2	$12\mu_j =$	5	8	-1		$-\frac{1}{24}h^3u^{(4)}$
3	$24\mu_j =$	9	19	-5	1	$-\frac{19}{720}h^4u^{(5)}$

Die Konsistenz von beiden Verfahrensklassen folgt aus den Approximationseigenschaften der Interpolationspolynome. Beide Verfahrenstypen werden in einem Satz behandelt. Dazu wird das A-B-Verfahren als A-M-Verfahren mit Gewicht  $\beta_0 = 0$  interpretiert.

**Satz 1.3.1** *Der lokale Fehler bei den Adams-Bashforth-Verfahren (1.3.2),  $\gamma_0 := 0, \gamma_j := \beta_j, 1 \leq j \leq m$ , bzw. den Adams-Moulton-Verfahren (1.3.3),  $\gamma_j := \mu_j, 0 \leq j \leq m$ , ist gegeben durch*

$$T_h(x_{i-1}) := \frac{1}{h}[u(x_i) - u(x_{i-1})] - \sum_{j=0}^m \gamma_j u'(x_{i-j}).$$

Er hat die Form

$$T_h(x_{i-1}) = \begin{cases} c_m^B h^m u^{(m+1)}(x_{i-1}) + \mathcal{O}(h^{m+1}) & \text{bei A-B} \\ c_m^M h^{m+1} u^{(m+2)}(x_{i-1}) + \mathcal{O}(h^{m+2}) & \text{bei A-M} \end{cases}. \quad (1.3.4)$$

**Beweis** Beim A-B-Verfahren sei  $p$  das Interpolationspolynom vom Grad  $m-1$  zu den Werten  $(x_{i-j}, u'(x_{i-j}))$ ,  $j = 1, \dots, m$ . Dann folgt aus der Fehlerformel der Interpolation (Numerik I, 3.3.11) und durch Taylorentwicklung die Aussage

$$\begin{aligned} T_h(x_{i-1}) &= \int_0^1 [u'(x_{i-1} + ht) - p(x_{i-1} + ht)] dt \\ &= \frac{1}{m!} \int_0^1 (x_{i-1} + ht - x_{i-1}) \dots (x_{i-1} + ht - x_{i-m}) u^{(m+1)}(\xi(t)) dt \\ &= u^{(m+1)}(x_{i-1}) h^m \underbrace{\frac{1}{m!} \int_0^1 t(t+1) \dots (t+m-1) dt}_{=: c_m^B} + \mathcal{O}(h^{m+1}). \end{aligned}$$

Bei A-M betrachtet man das Polynom vom Grad  $m$  zu  $(x_{i-j}, u'(x_{i-j}))$ ,  $j = 0, \dots, m$ . ■

Die Fehlerkoeffizienten  $c_m$  sind in den obigen Tabellen angegeben. Das A-M-Verfahren mit  $m$  Stufen besitzt die gleiche Konvergenzordnung  $h^{m+1}$  wie das A-B-Verfahren mit  $m+1$  Stufen, hat aber kleinere

Fehlerkoeffizienten  $|c_m^M| < |c_{m+1}^B|$ ,  $m > 0$ . Daher ist das Adams-Moulton-Verfahren in der Regel vorzuziehen. Die Auflösung der hier auftretenden nichtlinearen Gleichung für  $y_i$  in (1.3.3) kann durch Kombination mit dem A-B-Verfahren umgangen werden, indem mit A-B zunächst eine Approximation  $\tilde{y}_i$  für den unbekanntem Wert berechnet wird. Dies ergibt das *Prädiktor-Korrektor-Verfahren* (P-K)

$$\left. \begin{aligned} \tilde{y}_i &:= y_{i-1} + h \sum_{j=1}^m \beta_j f(x_{i-j}, y_{i-j}) \\ y_i &:= y_{i-1} + h \mu_0 f(x_i, \tilde{y}_i) + h \sum_{j=1}^m \mu_j f(x_{i-j}, y_{i-j}) \end{aligned} \right\} i = m, m+1, \dots \quad (1.3.5)$$

Hier werden pro Schritt zwei Funktionsaufrufe benutzt. Der lokale Fehler dieses Verfahren ist tatsächlich von der Ordnung  $h^{m+1}$ . Denn nach Satz 1.3.1 gilt für Lipschitz-stetiges  $f$  mit  $\tilde{u}_i := u(x_{i-1}) + h \sum_{j=1}^m \beta_j f(x_{i-j}, u(x_{i-j}))$  hier

$$\begin{aligned} \|T_h^{PK}\| &= \|T_h^{AM} + \mu_0(u'(x_i) - f(x_i, \tilde{u}_i))\| \\ &\leq \|T_h^{AM}\| + |\mu_0|L\|u(x_i) - \tilde{u}_i\| = \|T_h^{AM}\| + |\mu_0|Lh\|T_h^{AB}\| \\ &= h^{m+1}(|c_m^M|\|u^{(m+2)}(x_{i-1})\| + |\mu_0|Lc_m^B\|u^{(m+1)}(x_{i-1})\|) + \mathcal{O}(h^{m+2}). \end{aligned}$$

Der Fehler bei Runge-Kutta-Verfahren besitzt keine so einfache Gestalt. Im allgemeinen ist er aber für ein Verfahren gleicher Ordnung kleiner als bei Mehrschrittverfahren. Für das Problem  $u' = u$ ,  $u(0) = 1$ , etwa ergeben sich bei den Verfahren vierter Ordnung nach einem Schritt

Verfahren:	Runge-Kutta	Adams-Bashf.	Adams-Moulton
lok. Fehler:	$\frac{1}{120}h^5$	$\frac{251}{720}h^5$	$\frac{19}{720}h^5$

Bei den Mehrschrittverfahren sind für gleiche Fehlerbeträge also kleinere Schrittweiten erforderlich. Außerdem ist die Programmierung aufwendiger, wie die folgende Diskussion zeigt.

#### Fragen zum praktischen Einsatz:

- *Anlaufrechnung:* Die  $m$ -Schritt-Verfahren können, wie aus (1.3.2)-(1.3.5) ersichtlich, erstmals nach einer Anlaufphase von  $m$  Integrationsschritten eingesetzt werden. Daher sind Startwerte  $y_0, \dots, y_{m-1}$  ausreichender Genauigkeit auf andere Weise bereitzustellen, entweder mit Runge-Kutta-Verfahren oder Mehrschrittverfahren niedriger (wachsender) Ordnung und sehr kleinen Schrittweiten.
- *Schrittweitenwechsel:* Die Adams-Formeln (1.3.2)-(1.3.5) beziehen sich auf äquidistante Punkte  $x_{i-m}, \dots, x_i$  im konstanten Abstand  $h$ . Die Koeffizienten  $\beta_j, \mu_j$  ergaben sich dabei durch Integration der Lagrange-Darstellung des Interpolationspolynoms, bei Adams-Bashforth also von  $p(x) = \sum_{j=1}^m L_j(x) f_{i-j}$  in der Form

$$f_h^{AB} = \sum_{j=1}^m \beta_j f_{i-j} = \int_0^1 \sum_{j=1}^m L_j(x_{i-1} + ht) f_{i-j} dt.$$

Wird  $p$  dagegen als Newton-Polynom (Numerik I, §3.3) dargestellt, gilt

$$f_h^{AB} = \sum_{j=0}^{m-1} h^j \delta_j f[x_{i-j-1}, \dots, x_{i-1}].$$

Ein Wechsel der Schrittweite  $h \rightarrow \hat{h}$  ist in dieser Form jetzt einfach möglich durch Ersetzung der Faktoren  $h^j \rightarrow \hat{h}^j$ . Die schon berechneten Differenzen  $f[\dots, x_{i-1}]$  werden übernommen. In der Praxis sind dabei zusätzliche Vorsichtsmaßnahmen zu treffen.

- *Schrittweiten-, Ordnungs-Steuerung:* Da die A-B-Verfahren Ordnung  $m$  und die A-M-Verfahren Ordnung  $m + 1$  besitzen, sind die beiden Teile des Prädiktor-Korrektor-Verfahrens (1.3.5) als Verfahrenspaar zur Schätzung des lokalen Fehlers einsetzbar. Algorithmus 1.2.10 könnte damit also auch hier direkt eingesetzt werden (vgl. letzte Anmerkung). Darüberhinaus läßt sich aber auch noch die Ordnung des Verfahrens steuern, da es Adams-Methoden beliebiger Ordnung gibt. Eine einfache Strategie hierfür erhält man aus der Regel:

*Ist ein Schritt in  $[x_{i-1}, x_i]$  mit einem Verfahren der Ordnung  $m$  akzeptiert worden, so wähle für  $[x_i, x_{i+1}]$  dasjenige der drei Verfahren der Ordnungen  $m-1, m, m+1$ , das die größte Schrittweite  $\hat{h}$  gestattet.*

Die erwähnten Maßnahmen ermöglichen auch bei Mehrschrittverfahren einen selbststeuernden Einsatz. Allerdings ist der damit verbundene Verwaltungsaufwand erheblich höher als bei Einschrittverfahren und muß bei einem Vergleich berücksichtigt werden (vgl. §1.4).

### 1.3.2 Lineare Mehrschrittverfahren und Stabilität

In den Verfahren (1.3.2, 1.3.3) werden zwar die Funktionswerte  $f_{i-j}$  aus  $m$  zurückliegenden Punkten benutzt ('rechte Seite'), die Lösungswerte dagegen nur von zwei Stellen,  $y_i, y_{i-1}$ . Ohne Erhöhung der Schrittzahl  $m$  kann man daher eventuell eine Verbesserung durch Einführung zusätzlicher Parameter erwarten. Dazu werden jetzt allgemeinere Verfahren der Form

$$\sum_{j=0}^m \alpha_j y_{i-j} = h \sum_{j=0}^m \beta_j f_{i-j}, \quad i = m, m+1, \dots, \quad (1.3.6)$$

betrachtet mit  $x_j = a + jh$ ,  $f_j := f(x_j, y_j)$  und  $\alpha_0 = 1$ . Für  $\beta_0 = 0$  ist dies ein explizites, andernfalls ein implizites, *lineares  $m$ -Schrittverfahren*.

Eines der Ziele bei der Konstruktion von Verfahren ist eine möglichst hohe Ordnung des lokalen bzw. Konsistenzfehlers

$$T_h(x_{i-1}) = \frac{1}{h} \sum_{j=0}^m \alpha_j u(x_{i-j}) - \sum_{j=0}^m \beta_j u'(x_{i-j}). \quad (1.3.7)$$

Zur Überprüfung werden die Taylorentwicklungen in beiden Summen verglichen. Mit

$$\begin{aligned} u(x_i - jh) &= \sum_{k=0}^p \frac{(-jh)^k}{k!} u^{(k)}(x_i) + \mathcal{O}(h^{p+1}) \\ u'(x_i - jh) &= \sum_{k=0}^{p-1} \frac{(-jh)^k}{k!} u^{(k+1)}(x_i) + \mathcal{O}(h^p) \end{aligned}$$

in (1.3.7) ergibt sich

$$\begin{aligned}
 T_h(x_{i-1}) &= \sum_{j=0}^m \left[ \sum_{k=0}^p \alpha_j \frac{(-j)^k}{k!} h^{k-1} u^{(k)} - \sum_{k=0}^{p-1} \beta_j \frac{(-j)^k}{k!} h^k u^{(k+1)} \right] + \mathcal{O}(h^p) \\
 &= \frac{1}{h} u(x_i) \left( \sum_{j=0}^m \alpha_j \right) - \sum_{k=0}^{p-1} \frac{(-h)^k}{k!} u^{(k+1)}(x_i) \sum_{j=0}^m \left( \alpha_j \frac{j^{k+1}}{k+1} + \beta_j j^k \right) + \mathcal{O}(h^p).
 \end{aligned}$$

In dieser Entwicklung nach  $h$ -Potenzen liest man die Ordnungsbedingungen direkt ab:

**Satz 1.3.2** Das Verfahren (1.3.6) besitzt Konsistenzordnung  $p$  genau dann, wenn seine Koeffizienten mit  $\alpha_0 = 1$  das folgende lineare Gleichungssystem erfüllen,

$$\begin{aligned}
 \sum_{j=0}^m \alpha_j &= 0, \\
 \sum_{j=0}^m \left( \frac{j^{k+1}}{k+1} \alpha_j + j^k \beta_j \right) &= 0, \quad k = 0, \dots, p-1.
 \end{aligned} \tag{1.3.8}$$

Es ist hier also wesentlich leichter Ordnungsbedingungen aufzustellen und zugehörige Lösungen zu berechnen als bei den Einschrittverfahren. Ein naheliegendes Ziel wäre es, durch geeignete Wahl der  $2m(+1)$  freien Parameter  $\alpha_j, \beta_j$  eine möglichst hohe Ordnung zu erreichen.

**Beispiel 1.3.3** Konstruktion eines expliziten ( $\beta_0 = 0$ ) 2-Schrittverfahrens mit maximaler Ordnung  $p = 3$ .

$$\left. \begin{aligned}
 1 + \alpha_1 + \alpha_2 &= 0 \\
 \alpha_1 + 2\alpha_2 + \beta_1 + \beta_2 &= 0 \\
 \frac{1}{2}(\alpha_1 + 4\alpha_2) + \beta_1 + 2\beta_2 &= 0 \\
 \frac{1}{3}(\alpha_1 + 8\alpha_2) + \beta_1 + 4\beta_2 &= 0
 \end{aligned} \right\} \implies \begin{aligned}
 (\alpha_0, \alpha_1, \alpha_2) &= (1, 4, -5) \\
 (\beta_0, \beta_1, \beta_2) &= (0, 4, 2)
 \end{aligned}$$

Das Verfahren lautet also

$$y_i + 4y_{i-1} - 5y_{i-2} = h(4f_{i-1} + 2f_{i-2}), \quad i \geq 2.$$

Einfaches Zahlenbeispiel: Bei der Dgl  $u' = -u$ ,  $u(0) = 1$ , mit exakten Startwerten  $y_0 := 1, y_1 = e^{-h}$ , Berechnung der Näherungen  $y_i := -4(1+h)y_{i-1} + (5-2h)y_{i-2}$  mit  $h = 0.2$ .

Numerisches Ergebnis: Die Fehler in den Approximationen oszillieren und wachsen ungefähr um einen Faktor 5 pro Schritt.

$i$	$y_i$
1	0.8187
2	0.6701
3	0.5497
4	0.4438
5	0.3986
6	0.1283
7	1.2177
8	-5.2551
9	30.8262



Das Verfahren ist **instabil**, Störungen wachsen extrem an. Im Gegensatz zu den Einschrittverfahren ist die Stabilitätsforderung eine echte Einschränkung für Mehrschrittverfahren!

Es gibt sogar eine sogenannte *Ordnungsbarriere* (von Dahlquist):

Ein stabiles  $m$ -Schritt-Verfahren besitzt eine (Konsistenz-)

$$\text{Ordnung } p \leq \begin{cases} m + 1, & \text{wenn } m \text{ ungerade} \\ m + 2, & \text{wenn } m \text{ gerade} \end{cases}.$$

Gegenüber den Adams-Verfahren kann die Ordnung also höchstens um 1 erhöht werden, wenn  $m$  gerade ist. Die zugehörigen Verfahren besitzen eine den Adams-Verfahren ähnliche Struktur,

$$y_i - y_{i-k} = h \sum_{j=0}^m \beta_j f_{i-j}, \quad i \geq m, \quad 1 \leq k \leq m. \quad (1.3.9)$$

Daher beschränkt sich der folgende Konvergenzsatz auf solche Verfahren und enthält eine zu den Einschrittverfahren ähnliche Stabilitätsaussage. Auf die Stabilitätsforderungen an allgemeine lineare Mehrschrittverfahren wird noch einmal kurz in §1.4 eingegangen.

**Satz 1.3.4** Die rechte Seite  $f$  der Dgl erfülle eine Lipschitzbedingung (1.1.5) mit der Konstanten  $L$ . Die Schrittweite im Verfahren sei  $h \leq h_0 := 1/(2L|\beta_0|)$  ( $:= \infty$  für  $\beta_0 = 0$ ). Für die Startwerte gelte  $\|y_i - u(x_i)\| \leq \varphi(h)$ ,  $i = 0, \dots, m-1$ , und für die Konsistenzfehler  $\|T_h(x_{i-1})\| \leq \tau(h)$ ,  $i = m, \dots, N$ . Dann gilt die Fehlerschranke

$$\|y_i - u(x_i)\| \leq e^{\lambda(x_i - x_{m-1})} [\varphi(h) + 2(x_i - x_{m-1})\tau(h)], \quad i = m, \dots, N, \quad (1.3.10)$$

mit  $\lambda = L \sum_{j=0}^m |\beta_j| / (1 - hL|\beta_0|) \leq 2L \sum_{j=0}^m |\beta_j|$ . Für genügend oft differenzierbares  $f$  und  $\tau(h) = \mathcal{O}(h^p)$ ,  $\varphi(h) = \mathcal{O}(h^p)$  ist also insbesondere auch  $\max_{\Delta} \|u_{\Delta} - u\| = \mathcal{O}(h^p)$ ,  $h \rightarrow 0$ .

**Beweis** Für  $hL|\beta_0| < 1$  ist die Gleichung (1.3.9) nach  $y_i$  eindeutig auflösbar (Banachscher Fixpunktsatz). Wieder durch Subtraktion der Gleichungen (1.3.9) und (1.3.7),

$$\begin{aligned} y_i - y_{i-k} &= h \sum_{j=0}^m \beta_j f(x_{i-j}, y_{i-j}) \\ u(x_i) - u(x_{i-k}) &= h \sum_{j=0}^m \beta_j f(x_{i-j}, u(x_{i-j})) + hT_h(x_{i-1}) \end{aligned}$$

folgt für die Fehler  $\varepsilon_j := y_j - u(x_j)$  die Beziehung

$$\begin{aligned} \|\varepsilon_i - \varepsilon_{i-k}\| &= h \left\| \sum_{j=0}^m \beta_j [f(x_{i-j}, y_{i-j}) - f(x_{i-j}, u(x_{i-j}))] - hT_h(x_{i-1}) \right\| \\ &\leq hL \sum_{j=0}^m |\beta_j| \|\varepsilon_{i-j}\| + h\tau(h). \end{aligned}$$

Mit  $\|\varepsilon_i - \varepsilon_{i-k}\| \geq \|\varepsilon_i\| - \|\varepsilon_{i-k}\|$  erhält man daraus die Schranke

$$(1 - hL|\beta_0|) \|\varepsilon_i\| \leq \|\varepsilon_{i-k}\| + hL \sum_{j=1}^m |\beta_j| \max_{j=1}^m \|\varepsilon_{i-j}\| + h\tau(h).$$

Da auf der rechten Seite viele ältere Fehlerwerte auftreten, liegt es nahe deren Maximum

$$\delta_j := \max_{\nu \leq j} \|\varepsilon_\nu\|, \quad j = m-1, \dots, N,$$

zu betrachten. Nach Division durch  $1 - hL|\beta_0|$  und Berücksichtigung der Identität  $1/(1 - hL|\beta_0|) = 1 + hL|\beta_0|/(1 - hL|\beta_0|)$  ergibt sich so

$$\|\varepsilon_i\| \leq (1 + h\lambda)\delta_{i-1} + \frac{h}{1 - hL|\beta_0|}\tau(h), \quad i = m, \dots, N. \quad (1.3.11)$$

Da nun  $\delta_i = \max\{\|\varepsilon_i\|, \delta_{i-1}\}$  und die rechte Seite größer als  $\delta_{i-1}$  ist, kann die linke Seite in (1.3.11) durch  $\delta_i$  ersetzt werden. Dann liegt aber die in Lemma 1.2.6 behandelte Situation vor, die auf die Schranke (1.3.10) führt. Als Startpunkt ist dabei allerdings  $x_{m-1}$  zu nehmen. ■

Die stabilen Verfahren höchster Ordnung  $m+2$  bei geradem  $m$  ergeben sich mit der Wahl  $k = m$  in (1.3.9) aus den abgeschlossenen Newton-Cotes-Quadraturformeln zum Intervall  $[x_{i-m}, x_i]$ ,

$$\frac{1}{h}[y_i - y_{i-m}] = \sum_{j=0}^m \beta_j f_{i-j} \cong \frac{1}{h} \int_{x_{i-m}}^{x_i} f(x, u(x)) dx. \quad (1.3.12)$$

Bei diesen Verfahren ist allerdings, im Gegensatz zu den Adams-Verfahren, die mit  $x$  wachsende Stabilitätsschranke (1.3.10) scharf, selbst dann, wenn im Problem nur exponentiell fallende Lösungen auftreten können (vgl. §1.4). Daher werden in der Praxis doch meist die Adams-Verfahren verwendet.

Aufgrund ihrer Symmetrie bzgl.  $\frac{1}{2}(x_{i-m} + x_i)$  besitzen Verfahren der Form (1.3.12) aber noch eine weitere Eigenschaft, die von großer praktischer Bedeutung ist. Das einfachste, explizite Verfahren dieser Form zu offenen Newton-Cotes-Formeln ergibt sich aus der Rechteckregel zum Intervall  $[-2h, 0]$ , seine Symmetrie ist erkennbar in der Form

$$y_{i+1} - y_{i-1} = 2hf(x_i, y_i), \quad i = 1, 2, \dots$$

Diese *explizite Mittelpunkregel* hat nicht nur Ordnung 2, sondern besitzt sogar eine *asymptotische Entwicklung* nach Potenzen von  $h^2$ . Daher können mit Richardson-Extrapolation (vgl. Numerik I, §5.5) Näherungen beliebig hoher Ordnung berechnet werden. Das resultierende Gesamt-Verfahren ist vom Typus her aber wieder ein Einschrittverfahren und wird daher in einem eigenen Abschnitt behandelt.

## 1.4 Extrapolationsverfahren

Die explizite Mittelpunkregel benötigt zur Durchführung noch einen Näherungswert  $y_1$ . Dieser kann ohne Verlust der Gesamtordnung 2 mit dem Eulerverfahren (Verfahren 1 in §1.2.1) bestimmt werden. Das betrachtete Verfahren ist daher

$$\begin{aligned} y_1 - y_0 &= hf(x_0, y_0), & y_0 &= u_0, \\ y_{i+1} - y_{i-1} &= 2hf(x_i, y_i), & i &= 1, 2, \dots \end{aligned} \quad (1.4.1)$$

Der erste Schritt führt eine Asymmetrie ein, die von der Mittelpunktregel als oszillierende Störung weitergegeben wird. Daher besitzt die Näherungslösung  $u_\Delta$  eine ungewöhnliche asymptotische Entwicklung der Gestalt

$$u_\Delta(x_j) - u(x_j) = \sum_{k=1}^q h^{2k} \left[ g_{2k}(x_j) + (-1)^j \tilde{g}_{2k}(x_j) \right] + \mathcal{O}(h^{2q+2}), \quad h \rightarrow 0, \quad (1.4.2)$$

mit vom Gitter unabhängigen Koeffizientenfunktionen  $g_{2k}, \tilde{g}_{2k}$ . Insbesondere haben also Punkte  $x_j$  mit geradem und ungeradem Index verschiedene Entwicklungen und dürfen daher nicht gleichzeitig zur Extrapolation herangezogen werden.

**Beispiel 1.4.1** Beim Problem  $u' = \lambda u$ ,  $u(0) = 1$  lautet das Verfahren  $y_0 = 1$ ,  $y_1 = 1 + h\lambda$ ,  $y_{i+1} = y_{i-1} + 2h\lambda y_i$ . Wie bei den linearen Dgln führt der Ansatz  $y_j = c\xi^j$  in der Differenzgleichung  $y_{i+1} - 2h\lambda y_i - y_{i-1} = 0$  auf die quadratische Gleichung

$$\xi^2 - 2h\lambda\xi - 1 = 0 \Rightarrow \xi_{1/2} = h\lambda \pm \sqrt{1 + h^2\lambda^2} = \begin{cases} \sqrt{1 + h^2\lambda^2} + h\lambda > 0 \\ -(\sqrt{1 + h^2\lambda^2} - h\lambda) < 0 \end{cases}.$$

Die  $y_j$  besitzen daher die Darstellung  $y_j = C_1 \xi_1^j + C_2 \xi_2^j$ , wobei die Koeffizienten  $C_1, C_2$  aus den Anfangsbedingungen  $y_0 = 1$ ,  $y_1 = 1 + h\lambda$  zu bestimmen sind. Um den Bezug auf das Gitter herzustellen, wird der Gitterindex bei  $y_j = u_\Delta(jh)$  in der Form  $j = x/h$  geschrieben. So ergibt sich die Näherungslösung explizit in der Form

$$2u_\Delta(x) = \left(1 + \frac{1}{\sqrt{1 + h^2\lambda^2}}\right) \left(\sqrt{1 + h^2\lambda^2} + h\lambda\right)^{x/h} + (-1)^{x/h} \frac{h^2\lambda^2}{1 + h^2\lambda^2 + \sqrt{1 + h^2\lambda^2}} \left(\sqrt{1 + h^2\lambda^2} - h\lambda\right)^{x/h}.$$

Man sieht leicht, daß mit Ausnahme von  $(-1)^{x/h}$  alle Ausdrücke eine Talyorentwicklung nach  $h$  um  $h = 0$  besitzen. Daraus folgt hier für den Anfang der Entwicklung (1.4.2)

$$u_\Delta(x) = e^{\lambda x} + h^2 \left[ \underbrace{-\lambda^2 \left(\frac{1}{6}\lambda x + \frac{1}{4}\right) e^{\lambda x}}_{g_2} + (-1)^{x/h} \underbrace{\frac{1}{4}\lambda^2 e^{-\lambda x}}_{\tilde{g}_2} \right] + \mathcal{O}(h^4), \quad h \rightarrow 0.$$

Wenn die Lösung der Dgl fällt,  $\lambda < 0$ , kann der oszillierende *und wachsende* Störterm  $\tilde{g}_2$  die numerische Lösung erheblich beeinträchtigen. Dieses Problem läßt sich etwas durch einen zusätzlichen Glättungsschritt abschwächen. Bei diesem rechnet man im Lösungsschritt  $x \rightarrow x + H$  mit einer Teilschrittweite  $H_* = H/N_*$  über den Endpunkt  $x + H$  hinaus und bildet den Mittelwert ('Tiefpaß', dämpft alternierende Anteile)

$$\hat{u}_\Delta(x + H) := \frac{1}{4} \left( u_\Delta(x + H - H_*) + 2u_\Delta(x + H) + u_\Delta(x + H + H_*) \right).$$

Zur Extrapolation sind für einen Schritt  $x \rightarrow x + H$  mehrere Näherungen  $u_{\Delta_j}$  zu berechnen. Bei der Mittelpunktregel wählt man dazu eine wachsende Folge gerader Zahlen  $N_0 < \dots < N_q$  und bildet

Teilschrittweiten  $H_j := H/N_j$  sowie Gitter  $x_i^{[j]} := x_0 + iH_j$  ( $i = 0, \dots, N_j$ ). Dies ergibt das *Verfahren von Gragg-Bulirsch-Stoer*:

1	für $j = 0, \dots, q$ : $y_0^{[j]} := u_0, y_1^{[j]} := y_0 + H_j f(x_0, y_0),$	Startwerte	Mittelpunkt-	-Regel	(1.4.3)
2	für $i = 1, \dots, N_j$ : $y_{i+1}^{[j]} := y_{i-1}^{[j]} + 2H_j f(x_i^{[j]}, y_i^{[j]})$				
3	$P_{j,0} := \frac{1}{4}(y_{N_j-1}^{[j]} + 2y_{N_j}^{[j]} + y_{N_j+1}^{[j]}),$	Glättung	Extrapo-	-lation	
4	für $k = 1, \dots, j$ : $i := j - k$ ; $P_{i,k} := P_{i+1,k-1} + \frac{1}{(N_j/N_i)^2 - 1}(P_{i+1,k-1} - P_{i,k-1})$				

Im Verfahren (1.4.3) kann einfach die am langsamsten wachsende Folge  $\{N_j\} = \{2, 4, 6, 8, 10, \dots\}$  verwendet werden. Schritt 3 ist der erwähnte Glättungsschritt.

In Schritt 4 wird die Richardson-Extrapolation ausgeführt, die dem Neville-Algorithmus (I.3.3.5) zur Polynominterpolation entspricht und der Reihe nach jeweils den führenden Term der asymptotischen Entwicklung (1.4.2) eliminiert (Numerik I, §5.5).

Im  $j$ -ten Lauf durch Schritt 3-4 wird dabei mit der Näherung  $P_{j,0}$  vom  $j$ -ten Gitter eine zusätzliche Schrägzeile unten an das Extrapolations-Tableau angehängt.

$P_{0,0}$	$P_{0,1}$	$P_{0,2}$	$\dots$	$P_{0,j}$
$P_{1,0}$	$P_{1,1}$			
$P_{2,0}$	$\vdots$			
$\vdots$	$P_{j-1,1}$			
$P_{j,0}$				

Danach gilt die (lokale) Fehleraussage

$$\frac{1}{H}[P_{0,j} - u(x_0 + H)] = \mathcal{O}(H^{2j+2}), \quad j \leq q. \tag{1.4.4}$$

Dieses Endergebnis  $P_{0,j}$  gehört zu einem Einschrittverfahren der Ordnung  $2j + 2$  für den Schritt  $x \rightarrow x + H$ , da es nicht auf Werte vor der Stelle  $x$  zugreift. Damit wurde insbesondere folgende Existenzaussage zur Theorie der Einschrittverfahren geliefert, die sich durch Abzählung der insgesamt in (1.4.3) benötigten Funktionsauswertungen ergibt.

**Satz 1.4.2** Für jedes  $p \in 2\mathbb{N}$  gibt es ein Einschrittverfahren mit Ordnung  $p$  und  $\frac{1}{4}p^2 + 1$  Stufen.

Das Extrapolationsverfahren benutzt also mehr Funktionsauswertungen als die in §1.2 besprochenen Runge-Kutta-Verfahren. Aufgrund der einfachen und einheitlichen Konstruktion von Verfahren verschiedener Ordnung kann hier aber eine Schrittweiten- und Ordnungssteuerung wie bei den Mehrschrittverfahren durchgeführt werden. Insbesondere wird die aktuell verwendete Ordnung  $2j + 2 (\leq 2q + 2)$  in (1.4.4) in der Praxis nicht vorgegeben, sondern durch Inspektion des Tableaus  $P_{i,k}$  bestimmt.

**Demo-Beispiel:** Extrapolationsverfahren mit Ordnungs- und Schrittweitensteuerung (vgl. Ende §1.3.1) beim AWP  $u' = f(x, u), u(-\frac{1}{2}) = u_0$  auf  $[-\frac{1}{2}, 1]$  mit ( $a := 800, b = \ln(1 + a) - 3$ )

$$f(x, y) = \begin{cases} -2axy^2, & x \leq 0 \\ 2[-\frac{ax}{1+ax^2} + bx]y^2, & x > 0 \end{cases} \quad \text{mit } u(x) = \begin{cases} \frac{1}{1+ax^2}, & x \leq 0 \\ \frac{1}{1+\ln(1+ax^2) - bx^2}, & x > 0 \end{cases} .$$

**Vergleich der Verfahren:**

	Einschritt	Extrapolation	Mehrschritt
Konsistenz-Ordnung	aufwendig	einfach	einfach
Stabilität	immer	immer	Einschränkung
Aufwand	merkbar	hoch	gering
Fehlerschätzung	möglich	einfach	einfach
Schrittweitenänd.	immer	immer	aufwendig
Ordnungswechsel	nein	einfach	möglich

Die *praktische* Erfahrung mit den drei behandelten Verfahrensklassen ergibt ein differenziertes Bild, keines kann als Universalverfahren bezeichnet werden. Zwar schneiden in Bezug auf die Gesamtanzahl der zur Lösung benötigten Funktionsaufrufe von  $f$  die Prädiktor-Korrektor-Verfahren am günstigsten ab. Da sie im adaptiven Einsatz beim Schrittweitenwechsel aber den höchsten Verwaltungsaufwand haben, sind sie tatsächlich nur selten schneller (in Bezug auf die Rechenzeiten) als Einschrittverfahren, nämlich in Fällen, wo die numerische Auswertung  $f(x, y)$  sehr aufwendig ist. Im Vergleich der Einschrittverfahren, Runge-Kutta und Extrapolation, zeigt sich ein Vorteil der möglichen Ordnungssteuerung bei letzteren erst für (extrem) hohe Genauigkeitsanforderungen. Das Runge-Kutta-Tripel DOP853 wird erst für extreme Toleranzen ( $< 10^{-12}$ ) von Extrapolationsverfahren übertroffen.

**1.5 Steife Anfangswertprobleme**

In der Praxis trifft man oft auf Differentialgleichungen, deren Lösungen sehr unterschiedliche *Zeitskalen* besitzen, bei denen also sowohl sehr schnell, als auch sehr langsam veränderliche Lösungen auftreten können. Dies ist, z.B., bei chemischen Reaktionen der Fall. Von größerer Bedeutung im Rahmen dieser Vorlesung ist aber die Tatsache, daß bestimmte Approximationsverfahren für parabolische, partielle Differentialgleichungen (Linienmethode, vgl. §3.4) auf gewöhnliche Anfangswertprobleme dieser Art mit großen Dimensionen  $n$  führen. Die Behandlung solcher Probleme ist ein sehr aktuelles Forschungsgebiet, daher können nur die wichtigsten Themen angerissen werden. Eine ausführlichere Darstellung findet sich in den Büchern von Hairer/Wanner und Strehmel/Weiner. Zur Problematik:

**Beispiel 1.5.1**  $u(x) \in \mathbf{R}^2$ ,

$$u'(x) = Lu(x) := \begin{pmatrix} -500 & 499 \\ 499 & -500 \end{pmatrix} u(x), \quad u(0) = u_0 := \begin{pmatrix} 1 \\ 0 \end{pmatrix}. \quad (1.5.1)$$

Da die Eigenwerte der Koeffizientenmatrix  $L$  die Werte  $\lambda_1 = -1$  und  $\lambda_2 = -999$  sind, lautet die exakte Lösung dieses Problems

$$u(x) = \begin{pmatrix} C_1 e^{-x} + C_2 e^{\lambda_2 x} \\ C_1 e^{-x} - C_2 e^{\lambda_2 x} \end{pmatrix} \quad (1.5.2)$$

mit  $C_1 = C_2 = \frac{1}{2}$ . Bis auf eine sehr kurze Startphase  $0 \leq x \ll 1/1000$  verhält sich die Lösung wie die Funktion  $\frac{1}{2} \begin{pmatrix} 1 \\ 1 \end{pmatrix} e^{-x}$ , trivialerweise gilt auch  $u(x) \rightarrow 0$ ,  $x \rightarrow \infty$ . Wird das Problem, z.B., mit dem

(expliziten) Euler-Verfahren und konstantem  $h$  behandelt, ergibt sich die numerische Lösung aus

$$y_{i+1} = (I + hL)y_i \quad \Rightarrow \quad y_i = \begin{pmatrix} C_1(1-h)^i + C_2(1+h\lambda_2)^i \\ C_1(1-h)^i - C_2(1+h\lambda_2)^i \end{pmatrix}.$$

Hier gilt nur für  $|1 + h\lambda_2| < 1$  die Aussage  $y_i \rightarrow 0$  ( $i \rightarrow \infty$ ), d.h., wenn

$$h < \frac{2}{-\lambda_2} = \frac{1}{500}.$$

Im anderen Fall wächst  $y_i$  oszillierend stark an. Obwohl also (z.B.)  $e^{\lambda_2 x} \leq 10^{-8}$  für  $x \geq 0.02$  ist und die Lösung für solche Argumente vollkommen glatt aussieht, auch durch einen Polygonzug mit größeren Schrittweiten gut approximiert werden kann, erzwingt das Vorhandensein des Eigenwerts  $\lambda_2$  im Problem die Verwendung kleiner Schrittweiten *bei diesem Verfahren*.

Beim *impliziten* Euler-Verfahren wird das Integral  $\int_0^1 f(x_i + th, u(x_i + th)) dt$  in (1.2.1) durch den Funktionswert am *rechten* Rand ersetzt. Bei diesem Verfahren,

$$y_{i+1} = y_i + h_i f(x_{i+1}, y_{i+1}) \quad \Leftrightarrow \quad \begin{cases} k_1 &= f(x_i + h_i, y_i + h_i k_1) \\ y_{i+1} &= y_i + h_i k_1 \end{cases} \quad (1.5.3)$$

sind  $k_1$  bzw.  $y_{i+1}$  nur implizit beschrieben, müssen also durch Auflösung eines (nicht-)linearen Gleichungssystems bestimmt werden. Im linearen Beispiel (1.5.1) führt dies auf die Beziehungen

$$y_{i+1} = y_i + h_i L y_{i+1} \quad \Leftrightarrow \quad (I - h_i L) y_{i+1} = y_i \quad \Leftrightarrow \quad y_{i+1} = (I - h_i L)^{-1} y_i.$$

Mit Hilfe der Eigenwertzerlegung von  $L$  kann auch hier wieder das Ergebnis für konstante Schrittweite  $h$  explizit bestimmt werden,

$$y_i = \begin{pmatrix} C_1(1+h)^{-i} + C_2(1-h\lambda_2)^{-i} \\ C_1(1+h)^{-i} - C_2(1-h\lambda_2)^{-i} \end{pmatrix}.$$

Da jetzt  $1/(1+h) < 1$  und  $1/(1-h\lambda_2) < 1$  gilt, für alle  $h > 0$ , ist zumindestens das asymptotische Verhalten der Lösung  $u$  reproduziert,  $y_i \rightarrow 0$ ,  $i \rightarrow \infty$  für alle Schrittweiten.

Die beiden Beispielverfahren unterscheiden sich also ganz wesentlich in ihren Stabilitätseigenschaften. Zur Überprüfung solcher Eigenschaften von numerischen Lösungen verschiedener Verfahren betrachtet man deren Verhalten bei bestimmten *Testproblemen* im Vergleich zur exakten Lösung. Die einfachste Testgleichung ist die skalare, linear-homogene

$$u'(x) = \lambda u(x), \quad u(0) = 1, \quad \lambda \in \mathbf{C}, \quad \operatorname{Re} \lambda \leq 0, \quad (1.5.4)$$

die aber über die Eigenvektor-Zerlegung auch Aussagen über Systeme der Form (1.5.1) gestattet. Für deren Lösungen gilt  $|u(x+s)| = |e^{\lambda s} u(x)| \leq |u(x)| \forall s \geq 0$ . Runge-Kutta-Verfahren (1.2.4) führen hier auf die skalaren Gleichungen

$$\begin{aligned} y_{i+1} &= y_i + h \sum_{j=1}^m b_j k_j, \\ k_j &= \lambda y_i + h \lambda \sum_{l=1}^m a_{jl} k_l \quad \Leftrightarrow \quad (I - zA)k = \lambda y_i \mathbb{1} \end{aligned}$$

wobei  $z := h\lambda$ ,  $k := (k_1, \dots, k_m)^\top$ ,  $A = (a_{jl})$ ,  $\mathbf{1} := (1, \dots, 1)^\top$  gesetzt wurde. Analog sei  $b = (b_1, \dots, b_m)^\top$ . Aus dieser Beziehung ergibt sich die Rekursion  $y_i \rightarrow y_{i+1}$  explizit in der Form

$$y_{i+1} = y_i + hb^\top k = \left[ 1 + zb^\top (I - zA)^{-1} \mathbf{1} \right] y_i = \varphi(z) y_i, \tag{1.5.5}$$

$$\varphi(z) := 1 + zb^\top (I - zA)^{-1} \mathbf{1}. \tag{1.5.6}$$

Einschrittverfahren führen bei (1.5.4) immer auf eine solche Beziehung  $y_{i+1} = \varphi(h\lambda)y_i$ . In (1.5.5) ist  $\varphi(z)$  eine rationale Funktion der komplexen Variablen  $z = h\lambda$ , die sogenannte *Stabilitätsfunktion*. Für explizite Verfahren wurde diese schon in Beispiel 1.2.1 bzw. bei (1.5.3) berechnet:

Verfahren:	expl. Euler	Runge/Heun (Nr.2/3)	klass. Runge-Kutta (Nr.4)	impl. Euler
$\varphi(z) =$	$1 + z$	$1 + z + \frac{1}{2}z^2$	$1 + z + \frac{1}{2}z^2 + \frac{1}{6}z^3 + \frac{1}{24}z^4$	$\frac{1}{1-z}$

Offensichtlich sind die Stabilitätsfunktionen der expliziten Verfahren ganze rationale Funktionen, d.h., Polynome. Für diese ist ja  $A$  eine strikt untere Dreiecksmatrix, also nilpotent mit  $A^m = 0$ . Daher ist die Neumannreihe ein Polynom,  $\varphi(z) = 1 + zb^\top (I - zA)^{-1} \mathbf{1} = 1 + z \sum_{j=0}^{m-1} z^j b^\top A^j \mathbf{1}$ . Wenn ein Verfahren mit Ordnung  $p$  konvergiert, bildet die Stabilitätsfunktionen natürlich eine rationale Approximationen an die  $e$ -Funktion mit  $\varphi(z) = e^z + \mathcal{O}(z^{p+1})$ ,  $z \rightarrow 0$ .

Bei der Testgleichung (1.5.4), die nur beschränkte Lösungen besitzt, interessieren diejenigen Argumentwerte  $z = h\lambda$ , für die die numerische Lösung  $y$  ebenfalls nicht wächst.

**Definition 1.5.2** Ein Einschrittverfahren besitze die Stabilitätsfunktion  $\varphi$ , d.h. bei Anwendung auf (1.5.4) führe es auf die Beziehung  $y_{i+1} = \varphi(h\lambda)y_i$ . Das Verfahren heißt *absolut stabil* bei  $z \in \mathbf{C}$ , falls  $|\varphi(z)| \leq 1$  ist. Die Menge

$$S := \{z \in \mathbf{C} : |\varphi(z)| \leq 1\} \tag{1.5.7}$$

heißt *Bereich der absoluten Stabilität des Verfahrens*.

Zur stabilen Integration von (1.5.4) muß bei einem gegebenen Verfahren die Schrittweite  $h$  so gewählt werden, daß  $h\lambda \in S$  gilt, da sonst unsinnige, wachsende Lösungen auftreten können. Für die in  $\mathbf{R}$  liegenden *Stabilitätsintervalle*  $S \cap \mathbf{R}$  der obigen Beispielverfahren gilt

Verfahren:	expl. Euler	Runge/Heun (Nr.2/3)	klass. Runge-Kutta (Nr.4)	impl. Euler
Stab.Intervall:	$[-2, 0]$	$[-2, 0]$	$[-2.78, 0]$	$(-\infty, 0]$

Die in Beispiel 1.5 für das explizite Euler-Verfahren gefundene Schrittweiten-Einschränkung  $h \leq 2/(-\lambda)$  entspricht also der Bedingung  $h\lambda \in [-2, 0] = S \cap \mathbf{R}$ . Von allen erwähnten Verfahren kommt bei Anwendung auf Probleme der Form (1.5.1) nur das implizite Eulerverfahren ohne Schrittweitenbeschränkung aus, denn dessen Stabilitätsbereich umfaßt die gesamte negative reelle Achse. Da bei expliziten Verfahren die Stabilitätsfunktion  $\varphi$  ein Polynom ist und daher gilt  $|\varphi(z)| \rightarrow \infty$ ,  $z \rightarrow \infty$ , besitzen diese nur beschränkte Stabilitätsbereiche.

Daher läßt sich allgemein feststellen, daß zur Lösung *steifer* Probleme, mit möglicherweise schnell ausklingenden Lösungen, nur implizite Verfahren sinnvoll sind, nur diese müssen sich ausschließlich

nach Genauigkeitsforderungen richten. Bei expliziten Verfahren dagegen wird die Schrittweite nicht durch die gewünschte Genauigkeit bestimmt, sondern (sehr restriktiv) durch Stabilitätsschwächen begrenzt. Eine automatische Schrittweitensteuerung erzeugt hier eine um den Grenzfall schwankende Schrittweitenfolge (vgl. Demo-Beispiel). In Anlehnung an das Lösungsverhalten der Dgl (1.5.4) zeichnet man Verfahren mit analogen Dämpfungseigenschaften aus.

**Definition 1.5.3** Ein Einschrittverfahren heißt A-stabil, wenn gilt

$$S \supseteq \mathbf{C}_- := \{z \in \mathbf{C} : \operatorname{Re} z \leq 0\}, \quad \text{d.h.,} \quad |\varphi(z)| \leq 1 \quad \forall z \in \mathbf{C}_-.$$

Offensichtlich ist das implizite Eulerverfahren A-stabil, da die Bedingung  $|\varphi(z)| = \left|\frac{1}{1-z}\right| \leq 1 \iff 1 \leq |1-z|^2 = 1 - 2\operatorname{Re} z + |z|^2$  für  $\operatorname{Re} z \leq 0$  hier sicher erfüllt ist.

### Implizite Runge-Kutta-Verfahren:

Bei den Impliziten RK-Verfahren gibt es eine Klasse, bei der gleichzeitig die höchstmöglichen Konvergenzordnungen und gute Stabilitätseigenschaften aufeinandertreffen. Diese Verfahren ergeben sich aus Anwendung der Gaußschen Quadraturformeln auf (1.2.1). Dazu werden die Stützstellen  $c_i$  in (1.2.4) als Gaußknoten, d.h. die Nullstellen des  $m$ -ten Legendre-Polynoms gewählt (vgl. Numerik I, §5.3). Die restlichen Koeffizienten bestimmen sich als Integrale der zugehörigen Lagrangepolynome

$$L_i(x) := \prod_{\substack{j=1 \\ j \neq i}}^m \frac{x - c_j}{c_i - c_j}, \quad a_{ij} := \int_0^{c_i} L_j(x) dx, \quad b_i := \int_0^1 L_i(x) dx, \quad i, j = 1, \dots, m. \quad (1.5.8)$$

**Satz 1.5.4** Für beliebiges  $m \in \mathbf{N}$  sind die Gauß-Runge-Kutta-Verfahren A-stabil. Die Konsistenzordnung des  $m$ -stufigen Verfahrens ist  $2m$ .

Die Gauß-Runge-Kutta-Verfahren besitzen also exzellente theoretische Eigenschaften. Es können sogar Stabilitätsaussagen für nichtlineare Probleme gezeigt werden. Allerdings sind bei diesen Verfahren in jedem Verfahrensschritt nichtlineare Gleichungssysteme der Größe  $mn$ , z.B. mit dem Newtonverfahren, zu lösen. Daher sind sie für den praktischen Einsatz meist zu teuer. Es gibt verschiedene Ansätze zur effektiveren Implementierung von IRK-Verfahren. Diese schränken aber entweder die Struktur ein, und damit leider auch die erreichbare Ordnung auf  $m + 1$ , oder besitzen andere Nachteile.

Eine effiziente Alternative sind *Linear-implizite Runge-Kutta-Verfahren*. Zur Motivation werde die Lösung der Gleichung (1.5.3) beim impliziten Euler-Verfahren betrachtet mit sog. autonomer rechter Seite  $f = f(u)$ . Der erste Schritt der Newton-Iteration mit Startwert  $y_{i+1}^{[0]} := y_i$  und  $L_i := \frac{\partial f}{\partial y}(y_i)$  ergibt

$$y_{i+1}^{[1]} = y_i + h(I - hL_i)^{-1} f(y_i).$$

Man kann nun diese Vorschrift als eigenständiges Verfahren studieren in Bezug auf Stabilität und Konsistenz. Bei mehrstufigen Verfahren sind dabei zusätzliche Parameter  $\gamma, \gamma_{ij}$  sinnvoll. Die Verfahrensvor-

schrift lautet dann mit  $L := \frac{\partial f}{\partial y}(y_i)$ ,

$$\begin{aligned} (I - h\gamma L)k_l &= f\left(y_i + h \sum_{j=1}^{l-1} \alpha_{lj}k_j\right) + hL \sum_{j=1}^{l-1} \gamma_{lj}k_j, \quad l = 1, \dots, m \\ y_{i+1} &= y_i + h \sum_{j=1}^m b_j k_j. \end{aligned} \quad (1.5.9)$$

Die Stabilitätsfunktion dieser Verfahren entspricht der eines IRK-Verfahrens mit Koeffizienten  $a_{lj} = \alpha_{lj} + \gamma_{lj}$ , daher gibt es A-stabile Verfahren. Da in allen Stufen die gleiche Matrix  $(I - h\gamma L)$  auftritt, ist pro Schritt  $y_i \rightarrow y_{i+1}$  nur eine LR- oder QR-Zerlegung erforderlich. Die Verfahren gehen auf Rosenbrock und Wanner zurück. Man kann auch hier eingebettete Verfahren studieren zur Schrittweitensteuerung, z.B., existieren verschiedene (3,4)-Paare. Diese sind für mittlere Genauigkeitsforderungen konkurrenzfähig zur folgenden Verfahrensklasse.

### Implizite Mehrschrittverfahren

Im praktischen Einsatz hat sich eine spezielle Verfahrensklasse bewährt, die auf den rückwärtigen Differenzenformeln beruht. Im Gegensatz zu den Adams-Verfahren wird in der allgemeinen Form des Mehrschrittverfahrens (1.3.6) nur ein einziger nichtverschwindender  $\beta$ -Koeffizient gewählt,  $\beta_0 \neq 0$ , während die Koeffizienten  $\alpha_j$  aus dem linksseitigen Differenzenquotienten abgeleitet werden mit

$$\frac{1}{h} \sum_{j=0}^m \alpha_j v(-jh) = \beta_0 v'(0) + \mathcal{O}(h^m), \quad h \rightarrow 0, \quad (1.5.10)$$

für  $v \in C^m[-m, 0]$ . Aufgrund von Satz 1.3.2 besitzt daher die  $m$ -stufige Version dieser BDF-Verfahren (backward difference formula) Ordnung  $m$ . Die Koeffizienten sind

$m$	$j =$	1	2	3	4	$\beta_0$	$\delta_m$
1	$\alpha_j =$	-1				1	$90^\circ$
2	$3\alpha_j =$	-4	1			$\frac{2}{3}$	$90^\circ$
3	$11\alpha_j =$	-18	9	-2		$\frac{6}{11}$	$86^\circ$
4	$25\alpha_j =$	-48	36	-16	3	$\frac{12}{25}$	$73^\circ$

Ihre Stabilitätseigenschaften sind allerdings weniger gut als bei Einschrittverfahren. Bei der Testgleichung (1.5.4) haben die BDF-Verfahren die Gestalt ( $z := h\lambda$ )

$$(1 - \beta_0 z)y_i + \sum_{j=1}^m \alpha_j y_{i-j} = 0, \quad i \geq m. \quad (1.5.11)$$

Über den Lösungsansatz  $y_i = C\xi^i$  kommt man wieder (vgl. Beispiel 1.4) auf die allgemeine Lösung  $y_i = \sum_{j=1}^m C_j \xi_j^i$ , wenn das *charakteristische Polynom*

$$p_m(\xi; z) = (1 - \beta_0 z)\xi^m + \alpha_1 \xi^{m-1} + \dots + \alpha_m$$

der Differenzgleichung (1.5.11)  $m$  verschiedene Nullstellen  $\xi_j$  besitzt. Bei mehrfachen Nullstellen ist die Darstellung durch Polynomanteile  $i^l \xi^i$  zu modifizieren. In Analogie zu den Einschrittverfahren

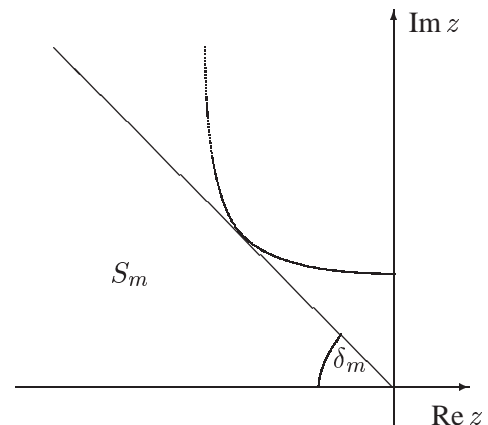
kann man die Menge dieser Nullstellen  $\{\xi_j\} = \{\xi_j(z)\}_{j=1}^m$  als (mehrwertige) Stabilitätsfunktion des Mehrschrittverfahrens interpretieren. Insbesondere nennt man ein Mehrschrittverfahren absolut stabil an der Stelle  $z \in \mathbf{C}$ , wenn gilt

$$p_m(\xi; z) = 0 \quad \Rightarrow \quad |\xi| < 1,$$

und A-stabil, wenn es für jedes  $z \in \mathbf{C}_-$  absolut stabil ist. Allerdings gibt es eine zweite *Ordnungsbarriere* von Dahlquist:

**Satz 1.5.5** Die Ordnung eines A-stabilen Mehrschrittverfahrens ist höchstens 2.

Daher können die Stabilitätsbereiche  $S_m$  der BDF-Verfahren höherer Ordnung  $m \geq 2$  nicht mehr die ganze Halbebene  $\mathbf{C}_-$  umfassen, in der obigen Tabelle ist der (halbe) Öffnungswinkel  $\delta_m$  des größten in  $S_m$  gelegenen Winkelraums angegeben. Für  $m > 6$  sind die Verfahren unbrauchbar, da dann nicht einmal mehr die negative reelle Achse vollständig in  $S$  enthalten ist. Auch für  $m = 6$  ist der Winkel  $\delta_6 = 17.8^\circ$  schon sehr klein (vgl. Diagramme).



**Demo-Beispiel** Das nichtlineare System von Dgln

$$\begin{aligned} u_j' &= -jdu_j, & j &\in \{1, n-1, n\} \\ u_j' &= -jdu_j + (u_{j-1}u_{j+2} - u_{j+1}u_j)^2, & j &= 2, \dots, n-2 \end{aligned}$$

besitzt die Lösung  $u_j(t) = e^{-jdt}$ . Für große  $n$  ist die Lipschitzkonstante der rechten Seite groß ( $\geq dn$ ). Die Stabilitätsprobleme expliziter Verfahren mit Schrittweitensteuerung sind daran erkennbar, daß Schrittweiten und Rechenzeit i.w. unabhängig von der geforderten Genauigkeit sind.

## 1.6 Schießverfahren für Randwertprobleme

Hier wird das standardisierte, nichtlineare Randwertproblem ('RWP', vgl. §1.1)

$$u'(x) = f(x, u(x)), \quad x \in (a, b), \quad r(u(a), u(b)) = 0, \quad (1.6.1)$$

betrachtet mit differenzierbaren Funktionen  $f(x, y)$ ,  $r(y_0, y_1)$ . Bei der numerischen Behandlung solcher gewöhnlicher RWPe kann man sowohl Lösungsverfahren von den Anfangswertproblemen übertragen, als auch neue Verfahren untersuchen, die auf partielle Randwertprobleme verallgemeinert werden können. Aus der Reihe der letzteren wird eine spezielle Klasse schon hier besprochen, da deren wesentliche Eigenschaften im gewöhnlichen Fall einfacher zu diskutieren sind. Zunächst werden aber Methoden behandelt, die auf AWP-Verfahren aufbauen. Dies hat den praktischen Vorteil, daß Standardprogramme aus dem AWP-Bereich einsetzbar sind.

### Schießverfahren

Aufbauend auf dem letzten Paragraphen wird vorausgesetzt, daß Anfangswertprobleme (mit genügender Genauigkeit) numerisch gelöst werden können. In §1.1 war das lineare Randwertproblem

$$u'(x) = A(x)u(x) + g(x), \quad x \in (a, b), \quad R_0u(a) + R_1u(b) = d, \quad (1.6.2)$$

mit Hilfe der Fundamentallösung  $Y(x)$  des Matrix-Anfangswertproblems  $Y'(x) = A(x)Y(x)$ ,  $Y(a) = I$ , gelöst worden. Dazu wurde in der Lösungsdarstellung  $u(x) = Y(x)\eta + \dots$  der unbekannte Anfangswert  $\eta = u(a) \in \mathbf{R}^n$  durch Einsetzen in die Randbedingung bestimmt:  $R_0\eta + R_1[Y(b)\eta + \dots] = d$ . Somit wurde also die Lösungskurve  $u(x)$  durch Wahl des Startwertes  $\eta = u(a)$  so angepaßt, daß ihr Wert  $u(b)$  am rechten Rand (zusammen mit  $u(a) = \eta$ ) die Randbedingung erfüllt (Bildlich *Schießen/Werfen*: Anpassung der Wurfrichtung so, daß Treffer erzielt wird).

Zur Verallgemeinerung dieses Prinzips muß im folgenden die Abhängigkeit der Lösung  $u$  vom Anfangswert betont werden. Daher bezeichne jetzt  $u(x; a, \eta)$  die Lösung des AWP

$$\frac{d}{dx}u(x; a, \eta) = f(x, u(x; a, \eta)), \quad x > a, \quad u(a; a, \eta) = \eta. \quad (1.6.3)$$

Die spezielle Lösung des RWPs (1.6.1) ist dann diejenige,  $u(x) = u(x; a, \hat{\eta})$ , deren Anfangswert  $\hat{\eta}$  das nichtlineare Gleichungssystem

$$F(\eta) = 0, \quad \text{mit } F : \begin{cases} \mathbf{R}^n & \rightarrow \mathbf{R}^n \\ \eta & \mapsto r(\eta, u(b; a, \eta)) \end{cases} \quad (1.6.4)$$

erfüllt. Zur Lösung solcher Gleichungssysteme kommt vor allem das Newton-Verfahren in Betracht, vgl. Numerik I, §6.2. Sein Einsatz setzt aber die genügende Differenzierbarkeit der Funktion  $F$  und die Kenntnis der Ableitung  $F'$  voraus. Die Lipschitz-stetige Abhängigkeit der Funktion  $u(x; a, \eta)$  vom Anfangswert  $\eta$  wurde schon in Satz 1.1.1 gezeigt. Darüber hinaus gilt aber auch

**Satz 1.6.1** Die rechte Seite  $f$  sei zweimal stetig differenzierbar, ihre Ableitung  $f_y = \frac{\partial f}{\partial y}$  beschränkt auf  $[a, b] \times \mathbf{R}^n$ . Dann hängt die Lösung  $u(x; a, \eta)$  von (1.6.3) differenzierbar vom Anfangswert  $\eta$  ab. Ihre Ableitung nach  $\eta$  an der Stelle  $(x, u(x; a, \eta))$  ist

$$\frac{\partial}{\partial \eta} u(x; a, \eta) = Y(x; a, \eta), \quad (1.6.5)$$

mit dem Fundamentalsystem  $Y$  zur linearisierten Differentialgleichung,

$$\frac{\partial}{\partial x} Y(x; a, \eta) = f_y(x, u(x; a, \eta)) \cdot Y(x; a, \eta), \quad Y(a; a, \eta) = I. \quad (1.6.6)$$

**Beweis** Nach Definition des Begriffs Differenzierbarkeit ist die Ableitung  $G'(\eta)$  einer Funktion  $G : \mathbf{R}^n \rightarrow \mathbf{R}^n$  diejenige lineare Abbildung, die bei Störung des Arguments  $\eta$  um  $s \in \mathbf{R}^n$  zum linearen Anteil der Differenz  $G(\eta + s) - G(\eta)$  gehört:  $G(\eta + s) - G(\eta) = G'(\eta)s + o(\|s\|)$ . Als Abbildung werde für festes  $x$  hier  $G : \eta \mapsto u(x; a, \eta)$  betrachtet, dazu werde  $v(x) := u(x; a, \eta + s) - u(x; a, \eta)$  definiert. Offensichtlich ist  $v(a) = \eta + s - \eta = s$  bei  $x = a$ . Nach Satz 1.1.1 gilt daher die Schranke  $\|v(x)\| \leq \exp(L(b-a))\|s\|$ . Da beide Funktionen  $u$  Lösungen der Dgl sind, gilt auch für  $x \in (a, b)$  (mit  $' = \partial/\partial x$ )

$$\begin{aligned} v'(x) &= u'(x; a, \eta + s) - u'(x; a, \eta) = f(x, u(x; a, \eta) + v(x)) - f(x, u(x; a, \eta)) \\ &= f_y(x, u(x; a, \eta))v(x) + \mathcal{O}(\|v(x)\|^2) = f_y(x, u(x; a, \eta))v(x) + \mathcal{O}(\|s\|^2). \end{aligned}$$

Die Funktion  $\hat{v}(x) := Y(x; a, \eta)s$  mit dem Fundamentalsystem  $Y$  aus (1.6.6) ist nun Lösung des linearen AWP's  $\hat{v}'(x) = f_y(x, u(x; a, \eta))\hat{v}(x)$ ,  $\hat{v}(a) = s$ . Da sich die rechten Seiten der Dgln für  $v$  und  $\hat{v}$  nur um  $\mathcal{O}(\|s\|^2)$  unterscheiden, gilt  $v(x) - \hat{v}(x) = \mathcal{O}(\|s\|^2)$  und somit tatsächlich

$$u(x; a, \eta + s) = u(x; a, \eta) + Y(x; a, \eta)s + \mathcal{O}(\|s\|^2), \quad s \rightarrow 0.$$

Damit erfüllt die Matrix  $Y(x; a, \eta)$  die Definition der Ableitung von  $u(x; a, \cdot)$  nach  $\eta$ , es gilt tatsächlich (1.6.5). ■

Die im Newtonverfahren zum System (1.6.4) auftretende Ableitung lautet nach der Kettenregel

$$\begin{aligned} F'(\eta) &= \frac{d}{d\eta} r(\eta, u(b; a, \eta)) = \frac{\partial r}{\partial y_0}(\eta, u(b; a, \eta)) + \frac{\partial r}{\partial y_1}(\eta, u(b; a, \eta)) \frac{\partial u}{\partial \eta}(b; a, \eta) \\ &= \frac{\partial r}{\partial y_0}(\eta, u(b; a, \eta)) + \frac{\partial r}{\partial y_1}(\eta, u(b; a, \eta)) Y(b; a, \eta). \end{aligned}$$

Die Bestimmung von  $F'$  ist somit durch Lösung des Matrix-AWP's (1.6.6), also durch (numerische) Berechnung von  $n$  Lösungen  $y_i(x) = Y(x)e_i, i = 1, \dots, n$ , der linearen AWP's

$$y_i' = f_y(x, u(x; a, \eta))y_i, \quad y_i(a) = e_i,$$

möglich ( $e_i \in \mathbf{R}^n$  :  $i$ -ter Einheitsvektor). Zur praktischen Durchführung können alle im §1.2-1.4 behandelten Verfahren herangezogen werden, es ist für die Diskussion aber einfacher von einer exakten Lösung der AWP's auszugehen. Damit kennt man die beim Newtonverfahren

$$\eta^{[k+1]} := \eta^{[k]} - \left(F'(\eta^{[k]})\right)^{-1} F(\eta^{[k]}), \quad k = 0, 1, \dots$$

auftretenden Größen. In der folgenden Formulierung des Gesamtverfahrens wird bei den Hilfsgrößen  $u(x; a, \eta^{[k]})$ ,  $Y(x; a, \eta^{[k]})$  der Index  $k$  und die Abhängigkeit vom Startwert unterdrückt.

**Algorithmus 1.6.2** Schießverfahren: Gegeben sei  $\eta^{[0]} \in \mathbf{R}^n$ . Für  $k = 0, 1, \dots$  löse man

$$\text{die AWPe } \begin{cases} u'(x) = f(x, u(x)), & u(a) = \eta^{[k]}, \\ Y'(x) = f_y(x, u(x))Y(x), & Y(a) = I, \end{cases} \quad (1.6.7)$$

$$\text{das LGS } (r_{y_0} + r_{y_1}Y(b))(\eta^{[k+1]} - \eta^{[k]}) = -r(\eta^{[k]}, u(b)). \quad (1.6.8)$$

Die partiellen Ableitungen von  $r$  sind ebenfalls an der Stelle  $(\eta^{[k]}, u(b))$  auszuwerten. In (1.6.7) sind also zunächst  $n + 1$  gekoppelte AWPe zu lösen, danach in (1.6.8) ein lineares  $n \times n$ -System.

Im linearen Spezialfall (1.6.2) ist  $r(y_0, y_1) = R_0 y_0 + R_1 y_1 - d$ , d.h.  $\frac{\partial r}{\partial y_0} = R_0$ ,  $\frac{\partial r}{\partial y_1} = R_1$ . Auch  $f_y(x, y) = A(x)$  ist unabhängig von  $u(x; a, \eta)$  also auch von  $\eta$ . Somit hängt auch  $Y$  nicht von  $\eta$  ab. Ein Schritt des Schießverfahrens mit Startwert  $\eta^{[0]} = 0$  liefert daher

$$\eta^{[1]} = -\left(R_0 + R_1 Y(b)\right)^{-1} (0 + R_1 u(b) - d).$$

Dies ist genau der Lösungswert, der bei der expliziten Lösung in (1.1.10) berechnet wurde, (1.6.8) konvergiert also bei linearen Problemen erwartungsgemäß in einem Schritt.

Wie auch sonst beim Newton-Verfahren kann die Konvergenz dieser Iteration für 'genügend genaue' Startwerte  $\eta^{[0]}$  gezeigt werden, wenn die zweite Ableitung  $f_{yy}$  beschränkt ist, und in der exakten Lösung  $u$  des RWP's die Randmatrix  $F'(u(a)) = r_{y_0} + r_{y_1}Y(b)$  regulär ist, mit  $Y(b) = Y(b; a, u(a))$ . In der Praxis kann diese Aussage aber wertlos sein, wenn die Lösungen  $u(x; a, \eta)$  sich sehr schnell von der gesuchten Lösung  $u(x; a, u(a))$  entfernen. Dann ist der Einzugsbereich der Iteration (1.6.8), eventuell sogar der Definitionsbereich der Funktion  $F$ , sehr klein.

**Beispiel 1.6.3** Diese Probleme können anhand der allgemeinen Lösung der skalaren Dgl

$$u' = \frac{1 + u^2}{1 + x^2} =: f(x, u) \quad \Rightarrow \quad u(x; a, \eta) = \frac{\eta - a + (1 + a\eta)x}{1 + a\eta - (\eta - a)x}$$

erläutert werden. Soll das RWP mit der Randbedingung  $u(0) + u(b) = b$ ,  $b > 0$ , (Lösung  $u(x) \equiv x$ ) mit dem Schießverfahren gelöst werden, ist ein Startwert  $\eta < 1/b$  zu wählen, damit die Lösung  $u(x; 0, \eta) = (\eta + x)/(1 - \eta x)$  überhaupt den rechten Intervallrand erreicht. Mit der explizit bekannten allgemeinen Lösung erhält die Randgleichung (1.6.4) die Form

$$F(\eta) = \eta + u(b; 0, \eta) - b = \eta - b + \frac{b + \eta}{1 - b\eta} \stackrel{!}{=} 0.$$

Die Ableitung  $F'(\eta)$  kann hieraus natürlich direkt bestimmt werden. Andererseits gilt auch

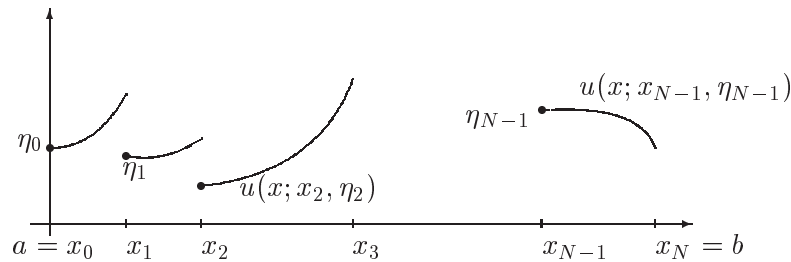
$$Y(x) = \frac{\partial u}{\partial \eta}(x; 0, \eta) = \frac{1 + x^2}{(1 - \eta x)^2} \quad \text{und} \quad Y'(x) = 2 \frac{\eta + x}{(1 - \eta x)(1 + x^2)} Y(x) = \underbrace{\frac{2u(x; 0, \eta)}{1 + x^2}}_{f_y(x, u)} Y(x).$$

Bei positiven Startwerten konvergiert das Newtonverfahren ebenfalls für  $\eta < 1/b$  (Ü-Aufgabe).

Die oben erwähnten Schwierigkeiten bei rasch divergierenden Lösungen des Anfangswertproblems können entschärft werden dadurch, daß man die Lösungen jeweils nur über kürzere Intervalle integriert, das Gesamtintervall also in mehrere Teile zerlegt. Dies führt auf die im folgenden beschriebene **Mehrzielmethode**:

Das Gesamtintervall  $[a, b]$  wird jetzt unterteilt durch ein Gitter

$$\Delta : a = x_0 < x_1 < \dots < x_N = b.$$



In jedem Teilintervall  $[x_i, x_{i+1}]$  wird das Schießverfahren getrennt angewendet. Dabei werden  $N$  Teil-Lösungen  $u(x; x_i, \eta_i)$  zu

$$u'(x; x_i, \eta_i) = f(x, u(x; x_i, \eta_i)), \quad u(x_i; x_i, \eta_i) = \eta_i, \quad i = 0, \dots, N-1, \quad (1.6.9)$$

berechnet. Eine Lösung  $u$  des Randwertproblems (1.6.1) kann nun aus diesen Stücken genau dann zusammengesetzt werden, wenn diese an den Gitterpunkten 'zusammenpassen' und außerdem die Randbedingung erfüllen, wenn also gilt

$$\begin{aligned} u(x_i; x_{i-1}, \eta_{i-1}) &= \eta_i && \equiv u(x_i; x_i, \eta_i), \quad i = 1, \dots, N-1, \\ r(\eta_0, u(x_N; x_{N-1}, \eta_{N-1})) &= 0. \end{aligned}$$

Wegen der Stetigkeit von  $f$  ist die so zusammengesetzte Lösung auch stetig differenzierbar. Im Unterschied zum Einfach-Schießverfahren sind nun  $N$  unbekannte Startwerte  $\eta_i$  zu bestimmen als eine Lösung des nichtlinearen Gleichungssystems für

$$\vec{\eta} := \begin{pmatrix} \eta_0 \\ \eta_1 \\ \vdots \\ \eta_{N-1} \end{pmatrix} \in \mathbf{R}^{nN} \text{ mit } F(\vec{\eta}) := \begin{pmatrix} u(x_i; x_{i-1}, \eta_{i-1}) - \eta_i, & i = 1, \dots, N-1, \\ r(\eta_0, u(x_N; x_{N-1}, \eta_{N-1})), & i = N \end{pmatrix} = 0 \quad (1.6.10)$$

Hierin ist das einfache Schießverfahren mit  $N = 1$  enthalten. Beim linearen RWP (1.6.2) können die Stück-Lösungen (1.6.9) wieder explizit dargestellt werden (vgl. Satz 1.1.3) durch

$$u(x; x_i, \eta_i) = Y(x)Y(x_i)^{-1}\eta_i + \int_{x_i}^x Y(x)Y(t)^{-1}g(t)dt.$$

Dabei sind  $Y(x)Y(x_i)^{-1} =: Y(x; x_i)$  diejenigen (hier von  $\eta_i$  unabhängigen) Fundamentallösungen der linearen Dgl mit  $Y(x_i; x_i) = I$ . Die Integrale entsprechen den speziellen Lösungen  $u(x; x_i, 0)$ . Das System (1.6.10) ist hier dann auch linear und lautet

$$\begin{aligned} Y(x_i; x_{i-1})\eta_{i-1} - \eta_i &= -u(x_i; x_{i-1}, 0), \quad i = 1, \dots, N-1, \\ R_0\eta_0 + R_1Y(b; x_{N-1})\eta_{N-1} &= d - R_1u(b; x_{N-1}, 0). \end{aligned}$$



Hilfe der ersten Gleichungen in (1.6.13):

$$\begin{aligned}\delta_1 &= Y_1\delta_0 + r_1, & \delta_2 &= Y_2\delta_1 + r_2 = Y_2Y_1\delta_0 + Y_2r_1 + r_2, & \dots \\ \delta_{N-1} &= Y_{N-1}Y_{N-2}\dots Y_1\delta_0 + \hat{r},\end{aligned}$$

$Y_i := Y(x_i; x_{i-1}, \eta_{i-1})$ . Mit der letzten Gleichung aus (1.6.13) ergibt dies also das System

$$\left(R_0 + R_1Y_NY_{N-1}\dots Y_2Y_1\right)\delta_0 = -r(\eta_0, \dots) - R_1Y_N\hat{r}, \quad (1.6.14)$$

das nur noch von der Größe  $n \times n$  ist. Diese Kondensation kann übrigens zusammen mit der Berechnung der  $Y_i$ , etc., erfolgen. Dann sind für einen Newtonschritt jeweils zwei Läufe durch das Intervall nötig, nämlich:

AWPe mit Kondensation, Lösung  $\delta_0$  aus (1.6.14), Berechnung von  $\delta_1, \dots, \delta_{N-1}$ .

Im linearen Fall und in der Nähe der Lösung  $u$  unterscheidet sich die Mehrzielmethode mit Kondensation nicht vom einfachen Schießverfahren. Denn, für die Matrix in (1.6.14) gilt mit exakten Startwerten  $\eta_i \equiv u(x_i)$  die Aussage

$$R_0 + R_1Y_N \dots Y_1 = R_0 + R_1Y(b; x_{N-1}, u(x_{N-1})) \dots Y(x_1; a, u(a)) = R_0 + R_1Y(b; a, u(a)).$$

Letzteres ist aber genau die Matrix aus (1.6.8), da in der Lösung die (Halb-)Gruppenidentität  $Y(x; s, u(s))Y(s; t, u(t)) = Y(x; t, u(t)) \forall t < s < x$  gilt. Die lokale Konvergenz (bei  $u$ ) des Newtonverfahrens hängt also insbesondere nur von der Regularität dieser Matrix ab. Der entscheidende Vorteil der Mehrzielmethode wirkt sich erst im nichtlinearen Fall aus durch (wesentliche) Vergrößerung des Einzugs-, d.h., Konvergenzbereichs beim Newtonverfahren. Die Wahl der Zwischenzielpunkte  $x_1, \dots, x_{N-1}$  kann dabei leicht während der Durchführung des Verfahrens modifiziert werden. Z. B. ist es sinnvoll, bei starkem Anwachsen einer Zwischenlösung  $u(\cdot; x_i, \eta_i)$  einen neuen Zielpunkt einzufügen. Die Auswirkungen der Mehrzielmethode kann anhand des einführenden Beispiels diskutiert werden.

**Beispiel 1.6.4** Eine Diskussion der Konvergenz ist hier naturgemäß viel schwieriger als im Fall  $N = 1$ . Als Mindestvoraussetzung kann aber die Existenz der Teillösungen in  $[x_i, x_{i+1}]$  überprüft werden. Beim Problem aus Beispiel 1.6 lautet der Lösungswert in  $x_{i+1}$ ,

$$u(x_{i+1}; x_i, \eta_i) \equiv \frac{\eta_i(1 + x_i x_{i+1}) + (x_{i+1} - x_i)}{1 + x_i^2 - (x_{i+1} - x_i)(\eta_i - x_i)}$$

und existiert insbesondere nur dann, wenn der Nenner positiv ist, also für

$$\eta_i - x_i = \eta_i - u(x_i) \leq \frac{1 + x_i^2}{x_{i+1} - x_i}, \quad i = 0, \dots, N - 1.$$

Diese Genauigkeitsforderungen an die lokalen Startwerte  $\eta_i$  sind schwächer als beim Einfach-Schießverfahren. Bei äquidistanter Unterteilung,  $x_i = ib/N$ , etwa ist für  $\eta_0$  nur zu fordern  $\eta_0 - u(x_0) \leq N/b$ , während in Beispiel 1.6  $\eta - u(x_0) < 1/b$  erforderlich war.

**Demo-Beispiel:** Mehrzielmethode beim Randwertproblem 2. Ordnung  $u'' = \frac{3}{2}u^2$ ,  $u(0) = 4$ ,  $u(1) = 1$ . Dieses besitzt 2 Lösungen.

## 1.7 Differenzenverfahren für Randwertprobleme

Die Konstruktion von numerischen Verfahren hatte bisher an der zur Dgl erster Ordnung äquivalenten Integralgleichung angesetzt. Bei allgemeineren Problemstellungen ist aber eine solche Umformung nicht einfach möglich bzw. unbequem, z.B. wenn ursprünglich Dgln höherer Ordnung vorliegen. Hier ist ein Zugang, der auf einer direkten Approximation der Ableitungen aufbaut, naheliegender, eine Standardmethode ist dabei die Approximation durch Differenzenausdrücke. Diese können explizit konstruiert oder systematisch durch Ableitung von Interpolationspolynomen bestimmt werden (vgl. Numerik I, §5.6). Eine besonders einfache Form besitzen dabei die symmetrischen Differenzenquotienten der Ordnung 2. Hier gilt

**Lemma 1.7.1** Für  $u \in C^{m+2}[-\epsilon, \epsilon]$ ,  $0 < \frac{m}{2}h \leq \epsilon$ , ist

$$\frac{1}{h^m} \sum_{j=0}^m \binom{m}{j} (-1)^j u\left(\left(j - \frac{m}{2}\right)h\right) = u^{(m)}(0) + C_m h^2 u^{(m+2)}(\xi), \quad \xi \in (-\epsilon, \epsilon). \quad (1.7.1)$$

**‘Beweis’:** Eine einfache Herleitung der Formel ist durch formale Rechnung mit Operatorreihen bei  $C^\infty(\mathbf{R})$ -Funktionen  $u$  (z.B. Polynomen) möglich. Mit  $h > 0$  werden Verschiebungs- und Ableitungsoperatoren punktweise definiert durch

$$(S_h u)(x) := u(x+h), \quad (Du)(x) := u'(x).$$

Für  $r \in \mathbf{R}$  ist dann  $(S_h^r u)(x) = (S_{rh} u)(x) = u(x+rh)$ . Der Satz von Taylor besagt nun

$$S_h u(x) = u(x+h) = \sum_{j=0}^{\infty} \frac{h^j}{j!} u^{(j)}(x) = \sum_{j=0}^{\infty} \frac{h^j}{j!} D^j u(x) = e^{hD} u(x), \quad \text{d.h., } \boxed{S_h = e^{hD}}.$$

Durch Differenzbildung erhält man Darstellungen für Differenzenoperatoren, etwa  $u(\cdot+h) - u(\cdot) = (S_h - I)u = (e^{hD} - I)u$ . Der symmetrische Differenzenausdruck in (1.7.1) folgt aus der binomischen Formel, denn

$$\begin{aligned} S_{h/2} - S_{-h/2} &= e^{\frac{h}{2}D} - e^{-\frac{h}{2}D} = 2 \sinh \frac{h}{2}D \quad \Rightarrow \quad (1.7.2) \\ \sum_{j=0}^m \binom{m}{j} (-1)^j S_{(j-m/2)h} &= (S_{h/2} - S_{-h/2})^m = \left(2 \sinh \frac{h}{2}D\right)^m \\ &= \left(hD + \frac{h^3}{24}D^3 + \dots\right)^m = h^m \left(D^m + \frac{m}{24}h^2 D^{m+2} + \dots\right). \end{aligned}$$

Im letzten Schritt wird die Reihenentwicklung von  $\sinh$  eingesetzt und liefert die Darstellung (1.7.1), wobei durch einige Zusatzüberlegungen das (Taylor-) Restglied vereinfacht wurde. ■

Von besonderem Interesse sind Differenzenverfahren bei (Systemen von) Dgln höherer Ordnung, da sich mit ihnen der Übergang zum äquivalenten System erster Ordnung, mit einer vielfachen Anzahl von Variablen, umgehen läßt. Als Vorbereitung auf die partiellen Dgln wird das häufig auftretende lineare Randwertproblem zweiter Ordnung betrachtet

$$-u''(x) + p(x)u'(x) + q(x)u(x) = g(x), \quad u(a) = \alpha, \quad u(b) = \beta. \quad (1.7.3)$$



**Satz 1.7.2** Für die Koeffizientenwerte in (1.7.5) gelte  $q_i \geq q^* > 0$ ,  $h|p_i| \leq 2$ ,  $i = 1, \dots, N$ . Dann ist die Matrix  $A$  streng diagonaldominant, also regulär, und es gilt (elementweise)

$$A^{-1} \geq 0 \quad \text{sowie} \quad \|A^{-1}\|_\infty \leq \frac{1}{q^*}.$$

**Beweis** Nach Voraussetzung sind die Nebendiagonalelemente von  $A = (a_{ij})$  nicht-positiv. Daher gilt  $A = D - B$  mit der Diagonalmatrix  $D = \text{diag}(a_{ii}) = \text{diag}(\frac{2}{h^2} + q_1, \dots, \frac{2}{h^2} + q_N)$  und einer nichtnegativen Matrix  $B$ . Die Zeilensummen von  $D^{-1}B$  für  $i = 2, \dots, N - 1$  sind

$$0 \leq -\frac{1}{a_{ii}}(a_{i,i-1} + a_{i,i+1}) = \frac{1 + \frac{h}{2}p_i + 1 - \frac{h}{2}p_i}{2 + h^2q_i} = \frac{2}{2 + h^2q_i} \leq \frac{2}{2 + h^2q^*} < 1.$$

Für  $i = 1, N$  gilt eine analoge Aussage. Damit folgt  $\|D^{-1}B\|_\infty < 1$  und die Neumann-Reihe

$$A^{-1} = (D - B)^{-1} = \sum_{j=0}^{\infty} \underbrace{(D^{-1}B)^j}_{\geq 0} D^{-1} \geq 0$$

konvergiert. Durch Übergang zu Normen in dieser Reihe folgt die zweite Aussage nach

$$\|A^{-1}\|_\infty \leq \sum_{j=0}^{\infty} \|D^{-1}B\|_\infty^j \|D^{-1}\|_\infty = \frac{\|D^{-1}\|_\infty}{1 - \|D^{-1}B\|_\infty} \leq \frac{h^2}{(2 + h^2q^*)(1 - \frac{2}{2 + h^2q^*})} = \frac{1}{q^*}. \quad \blacksquare$$

Der Satz erfordert ein positives  $q$ , hängt aber nicht von der Länge des Intervalls ab. Mit etwas mehr Aufwand läßt sich die Regularität von  $A$  unter Voraussetzungen zeigen, die nahe bei der Eindeutigkeitsbedingung an das RWP sind.

Der Satz ist auch von praktischem Interesse, da die Diagonaldominanz die Durchführbarkeit des Gaußalgorithmus ohne Zeilenvertauschung beim Tridiagonalsystem (1.7.6) garantiert mit  $\mathcal{O}(N)$  Operationen. Aus  $A^{-1} \geq 0$  folgt, daß die Lösung  $y$  monoton von der Inhomogenität  $g$  abhängt, und daß daher, z.B.,  $y \geq 0$  gilt, falls  $g \geq 0$ ,  $\alpha, \beta \geq 0$  sind. Die Normschränke  $\|A^{-1}\| \leq \frac{1}{q^*}$  schließlich ist eine Stabilitätsaussage, mit der sich aus der Konsistenz des Verfahrens direkt eine Fehlerschränke ableitet.

**Satz 1.7.3** Für die Koeffizientenfunktionen in (1.7.3) gelte  $p, q, g \in C^2[a, b]$  und  $q(x) \geq q^* > 0$ . Die Schrittweite im Verfahren sei so, daß  $h|p(x)| \leq 2 \forall x \in \Delta$ . Dann gilt für den globalen Fehler der Näherung  $u_\Delta$  aus (1.7.6) mit einer von  $\Delta$  unabhängigen Konstanten  $C$  die Aussage

$$\max_{x \in \Delta} |u(x) - u_\Delta(x)| \leq Ch^2 (\|u^{(4)}\|_\infty + \|u^{(3)}\|_\infty).$$

**Beweis** Nach den einleitenden Bemerkungen erfüllt der Gittervektor  $\vec{u}$  der Lösung  $u$  das System  $A\vec{u} = \tilde{g} + \tau$ ,  $\tau = (T_h(x_1), \dots, T_h(x_N))^T$ , vgl. (1.7.4). Durch Subtraktion von (1.7.6) folgt

$$A(\vec{u} - y) = \tau \quad \Rightarrow \quad \|u - y\|_\infty \leq \|A^{-1}\|_\infty \|\tau\|_\infty \leq \frac{1}{q^*} \|\tau\|_\infty.$$

Unter den Voraussetzungen des Satzes ist  $u \in C^4[a, b]$ . Daher gilt für die Konsistenzfehler  $T_h$  nach (1.7.4) und Lemma 1.7.1

$$\begin{aligned} |T_h(x_i)| &= \left| \frac{1}{h^2}(-u_{i-1} + 2u_i - u_{i+1}) + \frac{p_i}{2h}(u_{i+1} - u_{i-1}) + q_i u_i - g_i \right| \\ &\leq \underbrace{\left| -u''(x_i) + p_i u'(x_i) - q_i u(x_i) - g_i \right|}_{=0} + h^2(C_2 |u^{(4)}(\xi_i)| + 4C_1 |p_i u^{(3)}(\zeta_i)|) \\ &\leq C' h^2 (|p_i| \|u^{(3)}\|_\infty + \|u^{(4)}\|_\infty). \quad \blacksquare \end{aligned}$$

Dieses einfache Differenzenverfahren erzeugt also Näherungen an  $u$  von der (geringen) Konvergenzordnung zwei. Bei gewöhnlichen Randwertproblemen ist dies im Vergleich zum Schießverfahren im Verbund mit Anfangswertverfahren höherer Ordnung uninteressant.

Durch Umkehrung der im Beweis von Lemma 1.7.1 gefundenen formalen Identität,  $\frac{1}{2}(S_{h/2} - S_{-h/2}) = \sinh \frac{h}{2} D$ , ergibt sich folgende Entwicklung der zweiten Ableitung durch Differenzen

$$D^2 = \left( \frac{2}{h} \operatorname{arsinh} \frac{1}{2}(S_{h/2} - S_{-h/2}) \right)^2 = \frac{1}{h^2} \left( (S_{h/2} - S_{-h/2})^2 - \frac{1}{12}(S_{h/2} - S_{-h/2})^4 + \dots \right).$$

Daraus leitet man Differenzenapproximationen höherer Ordnung her. Die beiden in der letzten Gleichung angegebenen Terme führen auf die (wieder symmetrische) Differenzenformel

$$\begin{aligned} &\frac{1}{h^2} \left( S_{-h} - 2I + S_h - \frac{1}{12}(S_{-2h} - 4S_{-h} + 6I - 4S_h + S_{2h}) \right) u(x) \\ &= \frac{1}{12h^2} \left( -u(x-2h) + 16u(x-h) - 30u(x) + 16u(x+h) - u(x+2h) \right) \\ &= u''(x) - \frac{h^4}{90} u^{(6)}(\xi). \end{aligned} \quad (1.7.7)$$

In den randnahen Punkten  $x_1, x_N$  sind allerdings andere, unsymmetrische Formeln zu verwenden. Analoge Approximationen für  $u'$  sind ebenfalls leicht herzuleiten. Da die Formel (1.7.7) Funktionswerte aus 5 Gitterpunkten benutzt, ergibt sich nach Approximation der Dgl anstelle der Tridiagonalmatrix (1.7.5) eine Matrix mit fünf besetzten Diagonalen. Auch gehen weitere Eigenschaften, wie Diagonaldominanz, verloren, sodaß auch die einfache Stabilitätsaussage von Satz 1.7.2 nicht mehr möglich ist.

Bei speziellen Gleichungen ohne Beteiligung der ersten Ableitung  $u'$  (d.h. bei  $p \equiv 0$ ) kann allerdings ein  $\mathcal{O}(h^4)$ -Verfahren ohne diese Nachteile konstruiert werden. Dabei wird ausgenutzt, daß hier bei Kenntnis von  $u(x)$  auch  $u''(x) = q(x)u(x) - g(x) =: f(x, u(x))$  bekannt ist. Daher ersetzt man den  $h^2$ -Korrekturterm  $(S_{h/2} - S_{-h/2})^4 u(x) \cong h^4 u^{(4)}(x_i)$  in (1.7.7) durch  $h^2 (S_{h/2} - S_{-h/2})^2 u''(x_i)$ . Dieses *Mehrstellenverfahren* führt nun auf das Gleichungssystem

$$\begin{aligned} \frac{1}{h^2}(-y_{i-1} + 2y_i - y_{i+1}) + \frac{1}{12} \left( f(x_{i-1}, y_{i-1}) + 10f(x_i, y_i) + f(x_{i+1}, y_{i+1}) \right) &= 0, & \iff \\ -\left(1 - \frac{h^2}{12} q_i\right) y_{i-1} + \left(2 + \frac{5h^2}{6} q_i\right) y_i - \left(1 - \frac{h^2}{12} q_i\right) y_{i+1} &= \frac{h^2}{12} (g_{i-1} + 10g_i + g_{i+1}), \end{aligned}$$

$i = 1, \dots, N$ , dessen Koeffizientenmatrix wieder tridiagonal und diagonaldominant ist für  $q^* > 0$ .

## 2 Gemeinsame Prinzipien von Diskretisierungsverfahren

### 2.1 Konvergenz, Konsistenz, Stabilität

Den Fehlerbetrachtungen bei allen bisher behandelten Verfahren liegen gemeinsame Prinzipien zugrunde, die hier noch einmal herausgearbeitet werden sollen. Alle behandelten Probleme können als Gleichungssysteme geschrieben werden, wobei die jeweilige Lösung  $u$  implizit

$$F(u) = 0, \quad u \in D, \quad \text{mit } F : D \subset V \rightarrow W \quad (2.1.1)$$

definiert ist mit Hilfe einer Abbildung  $F$  zwischen normierten Räumen  $V, W$ . Der Definitionsbereich  $D$  von  $F$  ist dabei i. a. eine echte Untermenge von  $V$ . Bei der Lösung von Differentialgleichungen ist  $V$  ein Raum von stetigen, differenzierbaren Funktionen, z.B.,  $(C^1[a, b])^n$ . Die Abbildung  $F$  besteht hauptsächlich aus einem Ausdruck, der Ableitungen der eingesetzten Funktion enthält, also einem Differentialoperator. Die zusätzlichen Anfangs- bzw. Randbedingungen können entweder in die Definition von  $F$  oder in die von  $V$  bzw.  $D$  eingebaut werden. Bei den Systemen von Dgln erster Ordnung läßt sich

$$\begin{aligned} V &= (C^1[a, b])^n & \text{mit Norm} & \quad \|v\|_V := \max\{\|v(x)\| : x \in [a, b]\}, \\ W &= \mathbf{R}^n \times (C[a, b])^n & \text{mit Norm} & \quad \left\| \begin{pmatrix} d \\ g \end{pmatrix} \right\|_W := \|d\| + \max\{\|g(x)\| : x \in [a, b]\} \end{aligned}$$

$D = V$ , wählen. Dabei sei die nicht-indizierte Norm eine beliebige des  $\mathbf{R}^n$ . Man beachte, daß der so gewählte Raum  $V$  nicht vollständig ist, da die Norm keine Ableitungen einschließt. Als Abbildung ergibt sich beim

$$\underline{\text{AWP}} \quad F : v \mapsto \begin{pmatrix} v(a) - u_0 \\ v' - f(\cdot, v) \end{pmatrix}, \quad \underline{\text{RWP}} \quad F : v \mapsto \begin{pmatrix} r(v(a), v(b)) \\ v' - f(\cdot, v) \end{pmatrix}. \quad (2.1.2)$$

Beim Randwertproblem (1.7.3) zweiter Ordnung können homogene Randbedingungen,  $\alpha = \beta = 0$ , besser in die Definition des Raums  $V$  eingebaut werden. Hier bieten sich an

$$V := \{v \in C^2[a, b] : v(a) = v(b) = 0\}, \quad W := C[a, b], \quad (2.1.3)$$

$$F : v \mapsto -v'' + pv' + qv - g. \quad (2.1.4)$$

In  $V$  und  $W$  kann jeweils die Maximumnorm von  $v$  verwendet werden.

Zur numerischen Approximation solcher Probleme wird das Originalproblem durch eine ersetzt. Das Ziel bei der Konstruktion solcher Verfahren ist, bei wachsender Dimension im diskreten Problem den Fehler zwischen diskreter und exakter Lösung immer kleiner zu machen. In den bisher behandelten Beispielen entsprachen die endlichdimensionalen Räume jeweils einer Sammlung von Funktionswerten auf einem Gitter  $\Delta \subset [a, b]$ , wobei die Raumdimension bei feiner werdendem Gitter anwächst. Stattdessen kommen aber auch etwa Räume von Polynomen oder Splinefunktionen wachsenden Grads in Frage.

*Folge von eff*

Ein Diskretisierungsverfahren besteht also aus einer unendlichen *Schar* von endlichdimensionalen Räumen  $V_\Delta, W_\Delta$ ,  $\dim V_\Delta = \dim W_\Delta$ , und Abbildungen  $F_\Delta : V_\Delta \rightarrow W_\Delta$ , mit denen diskrete Lösungen definiert werden durch (nicht-) lineare Gleichungssysteme

$$F_\Delta(u_\Delta) = 0. \quad (2.1.5)$$

Um diese Näherungen mit der Originallösung  $u$  vergleichen zu können, wird mindestens ein Transfer von  $u$  nach  $V_\Delta$  (Restriktion) oder von  $u_\Delta$  nach  $V$  (Prolongation) benötigt. Bisher trat nur der erste Fall auf. Es sei also eine Restriktion bekannt, d.h., eine (lineare) Abbildung

$$R_\Delta : V \rightarrow V_\Delta, \text{ womit } e_\Delta := u_\Delta - R_\Delta u \quad (2.1.6)$$

als Fehler der diskreten Lösung definiert wird. Normen in  $V_\Delta$  und  $W_\Delta$  werden beide mit  $\|\cdot\|_\Delta$  bezeichnet. Diese Normen sollten in Beziehung zu den Normen der Räume  $V, W$  stehen. Dazu wird angenommen, daß ein Maß  $|\Delta|$  für die *Feinheit* der Gitter  $\Delta$  existiert mit den Eigenschaften, daß für jedes  $v \in V$  gilt

$$|\Delta| \rightarrow 0 \quad \Rightarrow \quad \dim V_\Delta \rightarrow \infty, \quad \|R_\Delta v\|_\Delta \rightarrow \|v\|_V. \quad (2.1.7)$$

Für die einfachsten Gitter, die äquidistanten mit Schrittweite  $h = (b - a)/N$ , ist  $|\Delta| := h$  ein solches Maß. Hier liegt sogar eine abzählbare Folge von Gittern vor, die durch die Zahl der Gitterpunkte indiziert werden kann. Die Beziehungen können durch folgendes Diagramm zusammengefaßt werden:

$$\begin{array}{ccc} V & \xrightarrow{F} & W \\ R_\Delta \downarrow & & \\ V_\Delta & \xrightarrow{F_\Delta} & W_\Delta \end{array}$$

**Beispiel 2.1.1** Einschrittverfahren bei Anfangs- und Randwertproblemen. Mit  $\Delta = \{x_i := a + ih, i = 0, \dots, N\}$ ,  $h = (b - a)/N$ ,  $N \in \mathbf{N}$ ,  $|\Delta| := h$ , sei

$$V_\Delta = \left\{ y_\Delta = \begin{pmatrix} y_0 \\ \vdots \\ y_N \end{pmatrix} : y_i \in \mathbf{R}^n, i = 0, \dots, N \right\}, \quad R_\Delta v := \begin{pmatrix} v(x_0) \\ \vdots \\ v(x_N) \end{pmatrix}, \quad v \in V, \\ \|y_\Delta\|_\Delta := \max_{i=0}^N \|y_i\|, \quad \dim V_\Delta = n(N + 1).$$

Die Restriktion  $R_\Delta v$  besteht also aus den Gitterwerten von  $v$ . Der Bildraum  $W_\Delta$  kann mit  $V_\Delta$  identifiziert werden. Mit einem durch seine Verfahrensfunktion  $f_h$  definierten Einschrittverfahren erhält man beim *Anfangswertproblem* die diskrete Abbildung

$$y_\Delta \mapsto F_\Delta(y_\Delta) = \begin{pmatrix} y_0 - u_0 \\ \frac{1}{h}(y_1 - y_0) - f_h(x_0, y_0) \\ \vdots \\ \frac{1}{h}(y_N - y_{N-1}) - f_h(x_{N-1}, y_{N-1}) \end{pmatrix}. \quad (2.1.8)$$

Beim *Randwertproblem* ist nur die erste Komponente,  $y_0 - u_0$ , von  $F_\Delta$  durch  $r(y_0, y_N)$  zu ersetzen. Beim AWP läßt sich das Gleichungssystem  $F_\Delta(u_\Delta) = 0$  natürlich schrittweise nach  $y_0, y_1, \dots$  auflösen.

Beim RWP entspricht dieses Gleichungssystem dem der Mehrzielmethode (1.6.10), mit einer zyklischen Verschiebung der Zeilen, und kann im Rahmen der Newtoniteration auf ein  $n \times n$ -System kondensiert werden.

**Beispiel 2.1.2** Differenzenverfahren beim linearen RWP zweiter Ordnung mit homogenen Randbedingungen. Da die Randwerte durch die Randbedingungen bekannt sind, müssen sie nicht mehr als Unbekannte geführt werden. Mit dem Gitter aus Beispiel 2.1 wählt man hier

$$V_\Delta = \left\{ y_\Delta = \begin{pmatrix} y_1 \\ \vdots \\ y_{N-1} \end{pmatrix} : y_i \in \mathbf{R}, i = 1, \dots, N-1 \right\}, \quad R_\Delta v := \begin{pmatrix} v(x_1) \\ \vdots \\ v(x_{N-1}) \end{pmatrix}, \quad v \in V,$$

$$\|y_\Delta\|_\Delta := \max_{i=1}^{N-1} |y_i|, \quad \dim V_\Delta = N-1.$$

Das diskrete Problem hat damit die in §1.7 behandelte Form  $F_\Delta(u_\Delta) = Au_\Delta - \tilde{g}$ , mit der in (1.7.5) definierten Tridiagonalmatrix  $A$  (hier mit  $h = (b-a)/N$ , Dimensionen  $(N-1) \times (N-1)$ ).

Diese Bezeichnungen wurden vor allem eingeführt, um noch einmal die gemeinsame Struktur der alten Beweise hervorzuheben, insbesondere in Bezug auf das Zusammenspiel der Begriffe Konsistenz, Stabilität und Konvergenz, und deren Abhängigkeit von der geeigneten Wahl der Definitionen von  $F_\Delta$  und  $\|\cdot\|_\Delta$ . Im Beispiel 2.1 hätte man alle Komponenten von  $F_\Delta$  mit  $h$  multiplizieren können, ohne das Verfahren zu ändern. In diesem Fall wäre das Verfahren zwar noch konsistent, aber nicht mehr stabil im Sinne der folgenden Definition.

**Definition 2.1.3** Ein Diskretisierungsverfahren heißt konsistent (von der Ordnung  $p > 0$ ), wenn mit der Lösung  $u \in V$  für den Konsistenzfehler

$$T_\Delta := F_\Delta(R_\Delta u) \quad \text{gilt} \quad \|T_\Delta\|_\Delta \rightarrow 0, \quad \text{bzw.} \quad \|T_\Delta\|_\Delta = \mathcal{O}(|\Delta|^p), \quad |\Delta| \rightarrow 0.$$

Das Verfahren heißt stabil, wenn Konstanten  $r, S, H$  existieren so, daß für alle Gitter mit  $|\Delta| \leq H$  für je zwei Elemente  $v_\Delta^{[j]} \in V_\Delta$  mit  $F_\Delta(v_\Delta^{[j]}) = g_\Delta^{[j]}$  und  $\|g_\Delta^{[j]}\|_\Delta \leq r$ ,  $j = 1, 2$ , folgt

$$\|v_\Delta^{[1]} - v_\Delta^{[2]}\|_\Delta \leq S \|g_\Delta^{[1]} - g_\Delta^{[2]}\|_\Delta = S \|F_\Delta(v_\Delta^{[1]}) - F_\Delta(v_\Delta^{[2]})\|_\Delta. \quad (2.1.9)$$

Bei der Stabilität betrachtet man also nicht nur die Ausgangsgleichung (2.1.5), sondern um  $g_\Delta$  gestörte Versionen. Das wesentliche Merkmal der Stabilitätsforderung ist, daß die Schranke gleichmäßig gilt für die gesamte (unendliche) Folge von Gittern, die Stabilitätskonstante  $S$  also nicht von  $\Delta$  abhängt. Daher bedeutet Stabilität insbesondere auch die gleichmäßig beschränkte Invertierbarkeit der Ableitungen  $F'_\Delta$  für die ganze Folge der Gitter! Die Stabilitätsschwelle  $r$  beschränkt die Forderung auf eine nicht zu kleine Umgebung der diskreten Lösung  $u_\Delta$  zur rechten Seite  $g_\Delta = 0$ .

Diese Definition stimmt bis auf kleinere Modifikationen mit denen von früher überein, beim AWP war das Def. 1.2.3 und Def. 1.2.5, sowie Satz 1.2.7. Beim RWP zweiter Ordnung wurden die entsprechenden Aussagen direkt nachgewiesen, Konsistenz mit  $p = 2$  in Satz 1.7.3, Stabilität mit  $S = 1/q^*$  in Satz 1.7.2. Wie in der obigen Bemerkung gilt die Normschranke dieses Satzes gleichmäßig für die Inversen einer unendlichen Folge von Matrizen wachsender Dimension  $n \rightarrow \infty$ .

Die Konvergenz eines Verfahrens folgt nun wieder durch Kombination der beiden Eigenschaften Konsistenz und Stabilität.

**Satz 2.1.4** *Das diskrete Problem (2.1.5) sei stabil nach (2.1.9) und konsistent von der Ordnung  $p$  mit  $\|T_\Delta\|_\Delta \leq C_T |\Delta|^p$ . Dann konvergieren die Näherungslösungen mit Ordnung  $p$ , d.h., für alle Gitter mit  $|\Delta| \leq H_0 := \min\{H, (r/C_T)^{1/p}\}$ , gilt*

$$\|u_\Delta - R_\Delta u\|_\Delta \leq SC_T |\Delta|^p.$$

**Beweis** In (2.1.9) wird gewählt  $v_\Delta^{[1]} := u_\Delta$  mit  $g_\Delta^{[1]} = 0$  und  $v_\Delta^{[2]} := R_\Delta u$  mit  $g_\Delta^{[2]} = T_\Delta$ . Nach Voraussetzung ist  $|\Delta| \leq H$ ,  $\|g_\Delta^{[j]}\|_\Delta \leq r$ ,  $j = 1, 2$ , und die Behauptung folgt aus (2.1.9):

$$\|u_\Delta - R_\Delta u\|_\Delta \leq S \|F_\Delta(u_\Delta) - F_\Delta(R_\Delta u)\|_\Delta = S \|T_\Delta\|_\Delta. \quad \blacksquare$$

## 2.2 Allgemeine Verfahrens-Ansätze

In diesem allgemeinen Rahmen sollen noch kurz zwei Prinzipien skizziert werden, die jeweils einen allgemeinen Ansatz zur Konstruktion von Verfahren für unterschiedliche Problembereiche darstellen.

**Projektions-Verfahren** Näherungen für eine isolierte Lösung von  $F(u) = 0$ ,  $u \in V$ , werden hier in einem endlichdimensionalen Unterraum  $V_\Delta \subseteq V$  bestimmt, z.B., einem Raum von Polynomen oder Spline-Funktionen (vgl. Numerik I, §4). Für  $v_\Delta \in V_\Delta$  kann dann  $F(v_\Delta) \in W$  gebildet werden, allerdings ist das Gleichungssystem  $F(u_\Delta) = 0$  in der Regel unlösbar. Hier wird nun ein Projektor  $P_\Delta = P_\Delta^2$  eingeführt, und die Näherung  $u_\Delta \in V_\Delta$  definiert durch

$$P_\Delta F(u_\Delta) = 0, \quad P_\Delta : W \rightarrow W_\Delta. \quad (2.2.1)$$

Bei geeigneter Wahl von  $P_\Delta$  und  $\dim W_\Delta = \dim V_\Delta$  ist dieses System (lokal) eindeutig lösbar, Stabilitätseigenschaften der Abbildung  $P_\Delta F : V_\Delta \rightarrow W_\Delta$  müssen für das konkrete Verfahren untersucht werden. Aufgrund der speziellen Struktur läßt sich aber die Konvergenz nach einfachen Prinzipien analysieren. Im linearen Fall,  $F(u) = Lu - r$ , kann dazu das Hilfsproblem  $Lv = LR_\Delta u$  betrachtet werden. Es hat die Eigenschaft, daß exakte und numerische Lösung übereinstimmen, denn  $v = v_\Delta = R_\Delta u$ . Durch Subtraktion der Beziehungen  $P_\Delta L u_\Delta = P_\Delta r = P_\Delta L u$  und  $P_\Delta(Lv_\Delta - LR_\Delta u) = 0$  folgt die Gleichung  $P_\Delta L(u_\Delta - v_\Delta) = P_\Delta L(I - R_\Delta)u$ . Damit gilt

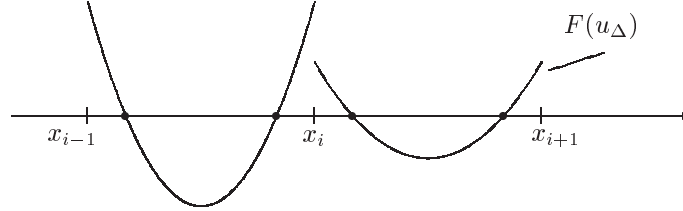
$$\|u - u_\Delta\| = \|u - v - (u_\Delta - v_\Delta)\| \leq (1 + K) \|(I - R_\Delta)u\|.$$

Der globale Fehler kann also bis auf die Konstante  $1 + K$  durch den Approximationsfehler  $(I - R_\Delta)u$  bei der Lösung  $u$  abgeschätzt werden. Allerdings muß dazu wieder im konkreten Fall die Konstante  $K$  als Norm der Abbildung  $(P_\Delta L|_{V_\Delta})^{-1} P_\Delta L : V \rightarrow V_\Delta$  ( $\neq \text{Id}$ !) bestimmt werden.

Spezielle Projektionsverfahren sind die in §3.3 behandelten Finite-Element-Verfahren und die *Kollokationsverfahren*. Ihr Name kommt von der Interpretation, daß man (bei Funktionen  $u = u(x)$ ,  $x \in \Omega$ ) für  $u_\Delta(x)$  die Gleichung  $F(u_\Delta(x)) = 0$  nicht mehr überall fordert,  $\forall x \in \Omega$ , sondern nur noch in endlich vielen Punkten  $x \in \bar{\Delta} \subseteq \Omega$ . Äquivalent dazu ist, daß man  $P_\Delta$  als Interpolationsabbildung mit den  $x \in \bar{\Delta}$  als Stützstellen definiert und (2.2.1) fordert. Als Beispiel werde eine Dgl  $k$ -ter Ordnung betrachtet, etwa (1.7.3) bei  $k = 2$ . Dazu wählt man ein Grundgitter  $\Delta : a = x_0 < x_1 < \dots < x_N = b$  und als Ansatzraum einen Raum von Spline-Funktionen

$$V_\Delta := S_{\Delta, p+k}^{k-1} = \{g \in C^{k-1}[a, b] : g|_{[x_i, x_{i+1}]} \in \Pi_{p+k}\},$$

die stückweise Polynome vom Grad  $p + k$  sind. Dabei ist nur  $g \in C^{m-1}$  wirklich erforderlich, vgl. Satz von Picard-Lindelöf ( $T : C^0 \rightarrow C^0$ ). Als Kollokationspunkte werden in jedem Teilintervall  $[x_i, x_{i+1}]$  jeweils  $p$  affin ähnliche Punkte verwendet  $\bar{\Delta} := \{x_i + h_i \xi_j : j = 1, \dots, p, i = 0, \dots, N-1\}$ .  $P_\Delta$  entspricht der Polynom-Interpolation unabhängig in allen Teilintervallen, Funktionen  $P_\Delta v$  sind also nur stückweise stetig. In der Skizze sind für  $p = 2$  die Kollokationspunkte  $\bar{\Delta}$  durch Kreise markiert, der Defekt  $F(u_\Delta)$  verschwindet nur dort.



Das Verfahren soll nicht weiter diskutiert werden, es sei jedoch erwähnt, daß bei Wahl der Kollokationsstellen  $\xi_1, \dots, \xi_p \in (0, 1)$  wie bei der Gauß-Quadratur (Numerik I, §5.3) sich eine *Superkonvergenz* auf den Gitterpunkten  $x_i$  einstellt, d.h., eine Konvergenzordnung, die über der Konsistenzordnung liegt. Dies war auch der Hintergrund für die hohe Ordnung in Satz 1.5.4, die Gauß-Runge-Kutta-Verfahren sind nämlich äquivalent zu Kollokationsverfahren in  $S_m^0$ .

**Defektkorrektur-Verfahren** Bei Verfahren höherer Ordnung reicht die Kopplung zwischen Unbekannten (Lösungs-Werten, -Koeffizienten) in der Regel weiter als bei solchen niedriger Ordnung. Daher ist die Struktur der (bei RWPen) auftretenden Gleichungssysteme komplexer und die Auflösung aufwendiger. Bei den Defektkorrekturverfahren wird dies umgangen dadurch, daß ausgehend von einer Lösung niedriger Ordnung schrittweise Näherungen wachsender Ordnung berechnet werden. Wegen der damit auch verfügbaren Fehlerschätzung bieten diese Verfahren als Bonus die Möglichkeit zu Gitter- und Ordnungssteuerung. Als Ansatzpunkt kann eine Newton-artige Iteration

$$L_\Delta(u_\Delta^{[k+1]} - u_\Delta^{[k]}) = -G_\Delta(u_\Delta^{[k]})$$

für eine nichtlineare Gleichung  $G_\Delta(u_\Delta) = 0$  dienen.  $L_\Delta$  ist dabei eine Approximation an  $G'_\Delta$ , da für die (lokale) Konvergenz die Bedingung  $\|I - L_\Delta^{-1}G'_\Delta(u_\Delta)\| = \|L_\Delta^{-1}(L_\Delta - G'_\Delta(u_\Delta))\| < 1$  erforderlich ist. Diese bietet aber auch den Spielraum merkliche Abweichungen  $L_\Delta - G'_\Delta$  zuzulassen, wenn  $L_\Delta$  dafür eine einfachere Struktur hat als  $G'_\Delta$ . Z.B., kann  $L_\Delta$  zu einem Verfahren niedriger Ordnung gehören als  $G_\Delta$ . Daher werde jetzt konkret folgende Iteration betrachtet,

$$F'_\Delta(u_\Delta^{[0]})(u_\Delta^{[k+1]} - u_\Delta^{[k]}) = -\hat{F}_\Delta(u_\Delta^{[k]}), \quad k = 0, 1, \dots \quad (2.2.2)$$

Dabei soll  $\hat{F}_\Delta$  ein Verfahren höherer Ordnung realisieren, während  $F'_\Delta$  die Ableitung zu einem Verfahren niedriger Ordnung darstellt. Als Startwert kann die Lösung  $u^{[0]}$  dieses *Basisverfahrens* dienen, definiert durch  $F_\Delta(u^{[0]}) = 0$ . Dann kann in allen Schritten (2.2.2) die LR-Zerlegung von  $F'_\Delta(u_\Delta^{[0]})$  wiederverwendet werden.

Zur Illustration werden Differenzenverfahren beim linearen RWP (1.7.3) mit  $p \equiv 0$  betrachtet. Zum Basisverfahren  $F_\Delta$  der Ordnung 2 gehört die Tridiagonalmatrix (1.7.5). Die Abbildung  $\hat{F}_\Delta$  zum Verfahren 4-ter Ordnung mit (1.7.7) ist für  $i = 2, \dots, N-2$

$$(\hat{F}_\Delta u_\Delta)(x_i) = q(x_i)y_i + \frac{1}{h^2}(-y_{i-1} + 2y_i - y_{i+1}) + \frac{1}{12h^2}(y_{i-2} - 4y_{i-1} + 6y_i - 4y_{i+1} + y_{i+2}) - g(x_i),$$

( $y_i := u_\Delta(x_i)$ ), mit Modifikationen in den randnahen Punkten  $x_1, x_{N-1}$ . Die Ableitungsmatrix  $\hat{F}'$  hat  $\geq 5$  Diagonalen. Es sei  $\hat{u}_\Delta$  die genauere Lösung mit  $\hat{F}'(\hat{u}_\Delta) = 0$ . Betrachtet man nun den ersten Schritt der Iteration (2.2.2), ergibt sich für die Fehler dazu die Beziehung

$$F'_\Delta \cdot (u_\Delta^{[1]} - \hat{u}_\Delta) = (F'_\Delta - \hat{F}'_\Delta) \cdot (u_\Delta^{[0]} - \hat{u}_\Delta), \quad (2.2.3)$$

da  $\hat{F}_\Delta(u_\Delta^{[0]}) = F_\Delta(u_\Delta^{[0]}) - \hat{F}_\Delta(\hat{u}_\Delta) = \hat{F}'_\Delta \cdot (u_\Delta^{[0]} - \hat{u}_\Delta)$ . Wegen der Linearität aller Abbildungen ist bei diesen kein Argument angegeben. Die Differenz  $(F'_\Delta - \hat{F}'_\Delta)v$  besteht jetzt nur noch aus dem Differenzenausdruck 4-ter Ordnung

$$\frac{1}{12h^2}(v_{i-2} - 4v_{i-1} + 6v_i - 4v_{i+1} + v_{i+2}) = 2h^2v[x_{i-2}, \dots, x_{i+2}].$$

In (2.2.3) ist  $v = u_\Delta^{[0]} - \hat{u}_\Delta = O(h^2)$ . Für genügend oft differenzierbare Koeffizienten der Dgl kann man (mit einigem Aufwand) nachweisen, daß dann auch für die dividierte Differenz  $v[x_{i-2}, \dots, x_{i+2}] = O(h^2)$  gilt, sodaß insgesamt die rechte Seite von (2.2.3) die Größenordnung  $O(h^4)$  besitzt, d.h., schon ein Defektkorrektur-Schritt erzeugt eine  $h^4$ -Näherung  $u_\Delta^{[1]}$ .

Wesentlich an dem Verfahren ist, daß der Defekt  $\hat{F}_\Delta(u_\Delta^{[k]})$  auf der rechten Seite der Iteration (2.2.2) mit einem Verfahren höherer Ordnung gebildet wird. Wenn allgemein  $\hat{F}_\Delta$  Konsistenz-Ordnung  $p$  hat und  $F_\Delta$  die Ordnung 2, ergibt sich bei passender Konstruktion der Fehler  $u_\Delta^{[k]} - R_\Delta u = O(h^{\min\{2k+2, p\}})$ .

### 3 Partielle Differentialgleichungen

#### 3.1 Allgemeine Eigenschaften

Im Unterschied zu den gewöhnlichen Differentialgleichungen werden jetzt Funktionen  $u(x_1, \dots, x_n)$  von mehr als einer Veränderlichen betrachtet. Da partielle Dgln ein sehr umfangreiches Problemgebiet darstellen, das hier nur ansatzweise diskutiert werden kann, werden hier nur lineare Probleme in zwei (Raum-) Dimensionen diskutiert für Funktionen  $u(x, y)$ . Partielle Differentialgleichungen verknüpfen Ableitungen der gesuchten Funktion nach verschiedenen Variablen und können daher auch i.a. nicht auf eine Standardform wie (1.1.1) gebracht werden. Geht man von Systemen erster Ordnung für Funktionen  $u(x, y) := (u_1(x, y), \dots, u_n(x, y))^T$  aus, dann ist ein recht allgemeine Form

$$u_x = A(x, y)u_y + B(x, y)u + g(x, y). \quad (3.1.1)$$

Charakteristische Eigenschaften dieser Dgl, nach denen sich auch die Wahl der zusätzlichen Randbedingungen zu richten hat, werden vor allem durch die  $n \times n$ -Koeffizientenmatrix  $A$  bestimmt. Zwei Sonderfälle erhalten eine Bezeichnung.

**Definition 3.1.1** a) Das System (3.1.1) heißt elliptisch im Punkt  $(x, y)$ , wenn  $A(x, y)$  keine reellen Eigenwerte besitzt.

b) Das System (3.1.1) heißt hyperbolisch in  $(x, y)$ , wenn  $A(x, y)$  reell diagonalisierbar ist.

Für  $n = 1$  ist jede reelle Gleichung hyperbolisch. Für  $n > 2$  sind diese beiden Fälle nicht erschöpfend. Die wichtigsten Beispiele in Dimension  $n = 2$  für  $u = (v, w)^T$  sind folgende

**Beispiel 3.1.2** Die Cauchy-Riemann-Dgln

$$v_x + w_y = 0, \quad w_x - v_y = 0 \quad \Longleftrightarrow \quad u_x = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} u_y$$

sind überall elliptisch (EWe  $\pm i$ ). Wenn die Lösung sogar zweimal differenzierbar ist, führt einmalige Differentiation der Gleichungen auf  $v_{xx} + w_{yx} = 0$ ,  $w_{xy} - v_{yy} = 0$ , und man kann das System auf die Potentialgleichung  $v_{xx} + v_{yy} = 0$  für  $v$  (analog auch für  $w$ ) reduzieren. Dies ist eine einzelne partielle Dgl zweiter Ordnung.

**Beispiel 3.1.3** Aus dem hyperbolischen System (EWe  $\pm 1$ )

$$v_x + w_y = 0, \quad w_x + v_y = 0 \quad \Longleftrightarrow \quad u_x = \begin{pmatrix} 0 & -1 \\ -1 & 0 \end{pmatrix} u_y$$

wird analog die Wellengleichung  $v_{xx} - v_{yy} = 0$  für  $v$  und  $w$ .

**Beispiel 3.1.4** Das System

$$v_x + w_y = 0, \quad w_x + v = 0 \quad \Longleftrightarrow \quad u_x = - \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} u_y + \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} u$$

führt auf die Wärmeleitungsgleichungen  $v_{xx} - v_y = 0$  und  $w_{xx} - w_y = 0$ .

Das letzte Beispiel ist durch die Fälle aus Definition 3.1.1 nicht abgedeckt, sein Typ ist *parabolisch*. Da viele Probleme direkt als Gleichung zweiter Ordnung auftreten, wird eine vollständige Typeneinteilung nur dafür formuliert.

**Definition 3.1.5** Die allgemeine lineare Dgl in zwei Variablen  $x, y$

$$a(x, y)u_{xx} + 2b(x, y)u_{xy} + d(x, y)u_{yy} + e(x, y)u_x + f(x, y)u_y = g(x, y)$$

heißt im Punkt  $(x, y)$

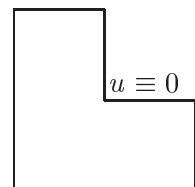
- a) elliptisch, wenn  $a(x, y)d(x, y) - b^2(x, y) > 0$ ,
- b) hyperbolisch, wenn  $a(x, y)d(x, y) - b^2(x, y) < 0$ ,
- c) parabolisch, wenn  $a(x, y)d(x, y) - b^2(x, y) = 0$  und  $\text{Rang} \begin{pmatrix} a & b & e \\ b & d & f \end{pmatrix} = 2$ .

Abhängig vom Typ sind nur bestimmte Randbedingungen sinnvoll. Bei einer (überall) elliptischen Dgl, die auf einem (evtl. unbeschränkten) Gebiet  $\Omega \subseteq \mathbf{R}^2$  gilt, sind auf dem ganzen Rand von  $\Omega$  Bedingungen an  $u$  und/oder  $u_x, u_y$  vorzuschreiben. Bei parabolischen und hyperbolischen dagegen nur auf einem Teil davon. Diese Unterschiede werden in den folgenden Kapiteln angesprochen. Praktische Auswirkung auf die Einsetzbarkeit numerischer Verfahren hat die in der Regel geringere Regularität der Lösungen partieller Gleichungen.

**Beispiel 3.1.6** Auf dem L-Gebiet  $\Omega := (-1, 1)^2 \setminus [0, 1]^2$  werde in Polarkoordinaten  $(x, y) = (r \cos \phi, r \sin \phi)$  die Funktion

$$u(r, \phi) = r^{2/3} \sin \frac{2\phi - \pi}{3}, \quad \frac{\pi}{2} \leq \phi \leq 2\pi,$$

definiert. Dann erfüllt  $u$  die Potentialgleichung  $u_{xx} + u_{yy} = 0$  auf  $\Omega$  und besitzt glatte Randwerte, z.B.  $u(x, y) = 0$  für  $x = 0, y \geq 0$  sowie  $x \geq 0, y = 0$ . Die ersten Ableitungen von  $u$  sind jedoch nicht beschränkt in der einspringenden Ecke im Nullpunkt  $r \rightarrow 0$ . In den übrigen Ecken des L treten etwas schwächere Singularitäten auf.



Diese Eigenschaft weicht von der Situation bei den gewöhnlichen Dgln ab und führt auf eine neue Rahmenbedingung bei der Diskussion numerischer Verfahren. Wenn eine gewöhnliche Dgl genügend oft differenzierbare Koeffizienten besitzt, kann durch Ableitung der Dgl auf eine entsprechende Differenzierbarkeit der Lösung geschlossen werden. Dies machte den Einsatz von Verfahren hoher Ordnung erst sinnvoll. Bei den partiellen Dgln dagegen kann die Lösung schon bei ganz einfachen Problemen Singularitäten (von Ableitungen) besitzen. Aus diesem und auch anderen, praktischen Gründen macht bei partiellen Dgln der Einsatz numerischer Verfahren hoher Ordnung Schwierigkeiten. Daher werden v.a. Verfahren der Ordnung zwei oder vier behandelt.

## 3.2 Differenzenverfahren für elliptische Randwertprobleme

### 3.2.1 Die Poissongleichung auf einfachen Gebieten

Es sei  $\Omega \subseteq \mathbf{R}^2$  ein beschränktes Gebiet, d.h. eine offene und zusammenhängende Menge und  $\Gamma = \partial\Omega$  sein Rand, der aus stückweise stetig differenzierbaren Kurven bestehen soll. Aus praktischen Gründen werden in diesem Paragraphen zunächst nur Polygone mit achsenparallelen oder diagonalen Kanten behandelt.  $\bar{\Omega} = \Omega \cup \Gamma$  ist der Abschluß von  $\Omega$ . Mit einer stetigen Funktion  $g(x, y)$  lautet die *Poissongleichung*

$$-u_{xx} - u_{yy} = g \quad \text{für } (x, y) \in \Omega. \quad (3.2.1)$$

Diese Gleichung ist elliptisch, sie stellt sogar die Normalform einer elliptischen Dgl dar. Physikalisch beschreibt sie, u.a., die Temperaturverteilung in einem homogenen Körper mit Wärmequellen  $g$ . Die homogene Gleichung ( $g \equiv 0$ ) heißt *Potentialgleichung*, ihre Lösungen harmonische Funktionen ( $\rightarrow$  Funktionentheorie). Zur Definition einer sinnvollen Problemstellung sind bei dieser Dgl auf ganz  $\Gamma$  Randwerte vorzugeben. Historisch unterscheidet man drei Arten von Randbedingungen, die den Funktionswert oder die Ableitung  $\partial u / \partial \mathbf{n}$  in Richtung der äußeren Normalen  $\mathbf{n}$  enthalten. Dazu wird der Rand disjunkt zerlegt in  $\Gamma = \Gamma_d \cup \Gamma_n \cup \Gamma_c$  und auf jedem Teil eine der folgenden Bedingungen gefordert mit stetigen Funktionen  $\phi, \psi$ .

$$u(x, y) = \phi_d(x, y), \quad (x, y) \in \Gamma_d \quad (\text{Dirichlet - RB}) \quad (3.2.2)$$

$$\frac{\partial u}{\partial \mathbf{n}}(x, y) = \phi_n(x, y), \quad (x, y) \in \Gamma_n \quad (\text{Neumann - RB}) \quad (3.2.3)$$

$$\frac{\partial u}{\partial \mathbf{n}}(x, y) + \psi(x, y)u(x, y) = \phi_c(x, y), \quad (x, y) \in \Gamma_c \quad (\text{Cauchy - RB}). \quad (3.2.4)$$

In der obigen physikalischen Interpretation entspricht die Dirichlet-Randbedingung einer Temperaturvorgabe auf  $\Gamma_d$ , die Neumann-RB (mit  $\phi_n \equiv 0$ ) einer isolierten Wand  $\Gamma_n$ , während die Cauchy-Bedingung sich als Wärmeabstrahlung durch  $\Gamma_c$  interpretieren läßt. In der Praxis kann natürlich auch jeweils nur eine der Randbedingungen auftreten. In diesen Fällen redet man von dem Dirichletschen ( $\Gamma = \Gamma_d$ ) bzw. Neumannschen ( $\Gamma = \Gamma_n$ ) Randwertproblem.

**Definition 3.2.1** Eine Funktion  $u(x, y)$  heißt (klassische) Lösung des Randwertproblems (3.2.1), (3.2.2-3.2.4), wenn  $u \in C^2(\Omega) \cap C(\bar{\Omega}) \cap C^1(\Omega \cup \Gamma_n \cup \Gamma_c)$  gilt und die Gleichungen (3.2.1), (3.2.2-3.2.4) erfüllt sind.

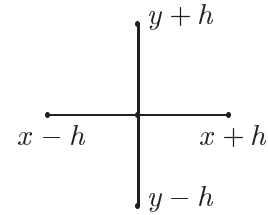
Während beim Dirichletproblem eine eindeutige Lösung existiert, besitzt das Neumannproblem Lösungen nur unter der Voraussetzung  $\int_{\Gamma} \phi_n ds = \int \int_{\Omega} g(x, y) dx dy$ . Die Lösung ist dann auch nur bis auf eine Konstante eindeutig.

Wie im eindimensionalen Fall (§1.7) lassen sich Näherungsverfahren gewinnen, indem man die verschiedenen Ableitungen in der Dgl durch geeignete Differenzenquotienten ersetzt. Mit Schrittweiten  $h_x, h_y > 0$  gilt für  $u \in C^4$  nach Lemma 1.7.1

$$u_{xx}(x, y) = \frac{1}{h_x^2} [u(x - h_x, y) - 2u(x, y) + u(x + h_x, y)] + \mathcal{O}(h_x^2)$$

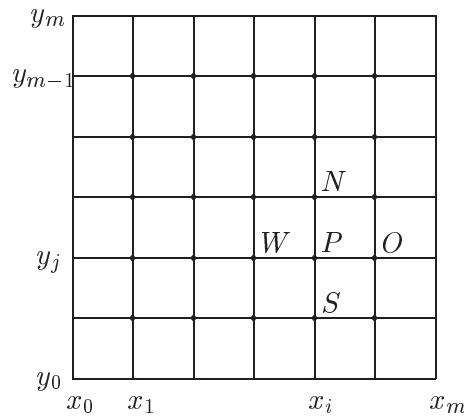
$$u_{yy}(x, y) = \frac{1}{h_y^2} [u(x, y - h_y) - 2u(x, y) + u(x, y + h_y)] + \mathcal{O}(h_y^2)$$

Bei Addition dieser beiden Näherungen werden somit zur Approximation der Poissongleichung in  $(x, y)$  fünf sternförmig verteilte Funktionswerte miteinander verknüpft. Für den Standardfall  $h_x = h_y = h$  haben die vier Strahlen des Sterns alle die Länge  $h$ .



Als erstes Modellproblem wird das Einheitsquadrat  $\Omega = (0, 1) \times (0, 1)$  betrachtet. Wie im eindimensionalen lassen sich mit einer Differenzenapproximation am einfachsten Funktionswerte eines äquidistanten Gitters verknüpfen. Dazu sei mit  $h := 1/m$ ,  $m \in \mathbb{N}$ , das Gitter

$$\Delta_h := \{(x_i, y_j), i, j = 0, \dots, m\}, \text{ mit } x_i = y_i := ih, i = 0, \dots, m. \tag{3.2.5}$$



Da die Handhabung im zweidimensionalen schwieriger ist als früher, wird diese Indizierung der Gitterpunkte mit zwei Indizes auch für die Näherungswerte  $v_{ij} = u_{\Delta}(x_i, y_j) \cong u(x_i, y_j)$  benutzt. Zur Beschreibung der Nachbarwerte eines speziellen Punktes  $P = (x_i, y_j)$  ist es jedoch einprägsamer, Himmelsrichtungen zu verwenden. Also gelte, zu

$$v_P = v_{i,j} \text{ sei } v_N := v_{i,j+1}, v_O := v_{i+1,j}, v_S := v_{i,j-1}, v_W := v_{i-1,j}, \tag{3.2.6}$$

wie in der obigen Skizze angedeutet. Die Poissongleichung im Punkt  $P$ ,  $-u_{xx}(P) - u_{yy}(P) = g(P)$  wird nun approximiert durch die Differenzengleichung  $-\frac{1}{h^2}[v_W - 2v_P + v_O] - \frac{1}{h^2}[v_S - 2v_P + v_N] = g_P$  mit  $g_P := g(x_i, y_j)$ , d.h.,

$$\begin{array}{c} -1 \\ -1 \quad | \quad 4 \quad -1 \\ -1 \quad | \quad -1 \end{array} \cdot \frac{1}{h^2} \qquad \frac{1}{h^2} [4v_P - v_N - v_O - v_S - v_W] = g_P. \tag{3.2.7}$$

Die Linearkombination in eckigen Klammern wird gerne kompakt durch den links dargestellten Fünf-Punkte-Stern beschrieben. Bei der (homogenen) Potentialgleichung kann diese Beziehung auch in der Form

$$v_P = \frac{1}{4} [v_N + v_O + v_S + v_W]$$

ausgedrückt werden, die aussagt, daß  $v_P$  der Mittelwert der Funktionswerte in den vier Nachbarn ist (Interpretation: 'aufgespanntes Gummnetz').

**Dirichletproblem:** Hier sind die Funktionswerte der Lösung auf dem ganzen Rand bekannt. Daher können die auf den Rand fallenden Variablen in (3.2.7) durch die entsprechenden Werte von  $\phi_d$  ersetzt werden. Aufgrund der einfachen Gestalt des Einheitsquadrats  $\Omega$  ergeben sich so für die Werte auf den im Bild bei (3.2.5) durch einen Punkt markierten inneren  $(m-1)^2$  Gitterpunkten gerade  $(m-1)^2$  Gleichungen (3.2.7). In Anlehnung an die Formulierung des RWPs (3.2.1),(3.2.2), kann man durch Einführung der Bezeichnungen

$$\begin{aligned}\Omega_h &:= \{(x_i, y_j) : i, j = 1, \dots, m-1\}, \\ \Gamma_h &:= \{(x_i, y_j) : i = 0 \vee i = m \vee j = 0 \vee j = m\},\end{aligned}\tag{3.2.8}$$

d.h.,  $\Omega_h \cup \Gamma_h = \Delta_h$ , dieses Gleichungssystem ebenfalls in kompakter Form schreiben:

$$\begin{aligned}\frac{1}{h^2}[4v_{ij} - v_{i,j+1} - v_{i+1,j} - v_{i,j-1} - v_{i-1,j}] &= g_{ij}, & \text{für } (x_i, y_j) \in \Omega_h, \\ v_{ij} &= \phi_d(x_i, y_j) & \text{für } (x_i, y_j) \in \Gamma_h.\end{aligned}\tag{3.2.9}$$

Mit den trivialen Randgleichungen werden die zugehörigen Unbekannten aus den Gleichungen auf  $\Omega_h$  eliminiert (vgl. auch (1.7.6)). Dann liegt ein lineares Gleichungssystem aus  $(m-1)^2$  Gleichungen für genauso viele Unbekannte vor. Bei den partiellen Dgln ist es angebracht, in der Formulierung wieder stärker zwischen der mit dem Gleichungssystem zusammenhängenden linearen Abbildung  $L_\Delta$  von Gitterfunktionen und den zugehörigen Matrizen zu unterscheiden. Wird zu Gitterfunktionen  $u_\Delta : \Omega_h \rightarrow \mathbf{R}$  die lineare Abbildung  $L_\Delta$  punktweise definiert durch

$$(L_\Delta u_\Delta)(x, y) := \frac{1}{h^2}[4u_\Delta(x, y) - u_\Delta(x, y+h) - u_\Delta(x, y-h) - \dots], \quad (x, y) \in \Omega_h,$$

mit Modifikationen in den randnahen Punkten, entspricht das Gleichungssystem (3.2.9) einem Problem der Form

$$L_\Delta u_\Delta = g_\Delta, \quad \text{auf } \Omega_h,$$

mit der schon verwendeten Abkürzung  $v_{ij} = u_\Delta(x_i, y_j)$ . 'Die' Matrix dieses linearen Systems hängt zumindest von einer Numerierung der bisher doppelt indizierten Unbekannten ab.

Bei *zeilenweiser Numerierung* (lexikographischer) des Gitters erhalten die drei Unbekannten  $v_W, v_P, v_O$  aufeinanderfolgende Nummern bzw. Indizes, während die Nummern von  $v_S, v_N$  genau um  $m-1$  kleiner bzw. größer sind als die von  $v_P$ . Die Nummern der Gleichungen (3.2.7)/(3.2.9) seien die von  $v_P$ , sodaß also das Hauptdiagonalelement immer  $4/h^2$  ist. Dann ist die zu dieser Numerierung gehörende Matrix  $A$  sehr regelmäßig aufgebaut, sie hat eine Block- und Bandstruktur mit Bandbreite  $2m-1$  in der Form

$$A = \frac{1}{h^2} \begin{pmatrix} T & -I & & & \\ -I & T & -I & & \\ & -I & T & -I & \\ & & & \ddots & \\ & & & & -I & T \end{pmatrix} = \begin{pmatrix} \diagup & & & & \\ & \diagup & & & \\ & & \diagup & & \\ & & & \diagup & \\ & & & & \diagup \end{pmatrix}, \tag{3.2.10}$$

$$T := \begin{pmatrix} 4 & -1 & & & \\ -1 & 4 & -1 & & \\ & -1 & 4 & -1 & \\ & & & \ddots & \\ -1 & & & & 4 \end{pmatrix}$$

Dabei sind die Matrizen  $I$  und  $T$  jeweils von der Größe  $(m-1) \times (m-1)$ , die Gesamtmatrix  $A$  von der Größe  $(m-1)^2 \times (m-1)^2$ .

Stellt man sich andererseits vor, daß die Gitterpunkte wie ein *Schachbrett* schwarz-weiß eingefärbt sind, erkennt man, daß in (3.2.9) die Nachbarpunkte  $N, O, S, W$  immer die zum Punkt  $P$  entgegengesetzte Farbe besitzen. Wenn daher zunächst weiße und dann schwarze Punkte numeriert werden, ergibt sich im Großen eine Blockstruktur der Form

$$\frac{1}{h^2} \begin{pmatrix} 4I & -M \\ -M^T & 4I \end{pmatrix},$$

wobei  $M$  in jeder Zeile höchstens vier Einsen und sonst Nullen enthält. Wenn  $m-1$  ungerade ist, haben die Blöcke  $4I$  geringfügig unterschiedliche Größen. Weitere Eigenschaften dieser Matrizen werden später diskutiert. Lösungsverfahren für solche Gleichungssysteme, die die geringe Anzahl der nichttrivialen Matrixelemente nutzen, sind Thema der Vorlesung Numerik IIA.

**Neumann-Randbedingung:** Hier ist der Wert der Normalenableitung vorgeschrieben. Daher muß auch bei der Randbedingung eine Approximation durch Differenzenquotienten vorgenommen werden. Beim Einheitsquadrat sind die Normalen parallel zu den Achsen. Daher gilt entweder  $\partial u / \partial \mathbf{n} = \pm u_x$  (auf der Ost- und Westseite), oder  $\partial u / \partial \mathbf{n} = \pm u_y$  (auf der Nord- und Südseite). Als Beispiel werde die Approximation der Randbedingung  $u_x(1, y) = \phi_n(1, y)$ ,  $y \in (0, 1)$  betrachtet. Dabei bieten sich folgende Möglichkeiten an:

1. Verwendung des Gitters (3.2.5), Approximation der Randbedingung durch den einseitigen Differenzenquotienten mit Schrittweite  $h$ : Im Beispiel sei  $P = (1-h, y_j)$  ein dem Rand benachbarter Gitterpunkt. Die Randbedingung wird also ersetzt durch  $\frac{1}{h}(v_O - v_P) = \phi_n(O)$ . Da die Variable  $v_O$  sonst nur noch im Differenzenstern von  $P$  auftaucht, kann die Bedingung dazu benutzt werden, diese Unbekannte sofort aus dem Stern zu eliminieren. Dann hat das Gleichungssystem wieder die gleiche Struktur wie beim Dirichletproblem, z.B. (3.2.10), nur die Sterne in den Randnachbarpunkten sind abgeändert, im Beispiel ergibt sich statt (3.2.7) jetzt

$$\frac{1}{h^2}(3v_P - v_N - v_S - v_W) = g_P + \frac{1}{h}\phi_n(O).$$

Nachteilig ist hierbei die schlechtere Approximation der Randbedingung mit einem  $\mathcal{O}(h)$ -Fehler wegen der Verwendung unsymmetrischer Differenzenquotienten.

2. Der symmetrische Differenzenquotient  $\frac{1}{2h}(u(x+h, y) - u(x-h, y))$  ist eine  $\mathcal{O}(h^2)$ -Approximation für den Wert  $u_x(x, y)$ . Um ihn einsetzen zu können, muß der entsprechende Randpunkt mit ins Berechnungsgitter  $\Omega_h$  aufgenommen werden, dies sei jetzt  $P = (1, y_j) \in \Gamma_n$ . Die approximierte Randbedingung

$$\frac{1}{2h}(v_O - v_W) = \phi_n(P) \quad (3.2.11)$$

verwendet dann den Punkt  $O$  außerhalb von  $\Omega$ . Die Unbekannte  $v_O$  kann aber, wie vorher, aus dem Fünfpunktstern im Randpunkt  $P$  eliminiert werden. Dies führt in  $P$  auf den einseitigen Stern

$$\begin{array}{c}
 -1 \\
 \hline
 -2 \quad \left| \quad 4 \cdot \frac{1}{h^2} \right. \\
 \hline
 -1
 \end{array}
 \quad
 \frac{1}{h^2} (4v_P - 2v_W - v_N - v_S) = g_P + \frac{2}{h} \phi_n(P).$$

Diese Methode ist besonders naheliegend, da (homogene) Neumann-Randbedingungen oft aufgrund von Symmetrieüberlegungen beim Dirchletproblem auftreten. Zum Beispiel gilt für die Lösung des Problems

$$-u_{xx} - u_{yy} = 1 \text{ auf } (-1, 1) \times (-1, 1), \quad u = 0 \text{ auf dem Rand,}$$

die Aussage  $u(x, y) = u(\pm x, \pm y)$ , also  $u_x(0, y) \equiv 0, u_y(x, 0) \equiv 0$ . Daher kann dieses Problem auf ein Viertel des ursprünglichen Quadrats eingeschränkt werden, etwa das Einheitsquadrat, wobei auf dessen linken und unteren Rand die Neumann-Bedingung (3.2.3) mit  $\phi_n \equiv 0$  gilt. Wird diese mit Hilfe von (3.2.11) diskretisiert, ergibt sich wieder eine Symmetriebedingung für  $v = u_\Delta$ , nämlich in  $P = (0, y_j)$  die Gleichung  $\frac{1}{h}(v_O - v_W) = 0 \iff v_O = v_W$  und in  $P = (x_i, 0)$  analog  $v_N = v_S$ . Diese Approximation der Randbedingung führt daher auf das gleiche Ergebnis, das man durch Berücksichtigung der Symmetrie beim diskreten Problem (3.2.9) erhalten würde. Die Symmetriebetrachtung reduziert dabei die Größe des Gleichungssystems auf ein Viertel.

Die Struktur der Matrix zu diesem reduzierten Problem

$$-u_{xx} - u_{yy} = g \quad \text{in} \quad \Omega = (0, 1) \times (0, 1), \tag{3.2.12}$$

$$u = 0 \quad \text{in} \quad \Gamma_d = \{(x, y) \in \bar{\Omega} : x = 1 \vee y = 1\}, \tag{3.2.13}$$

$$\frac{\partial}{\partial \mathbf{n}} u = 0 \quad \text{in} \quad \Gamma_n = \{(x, y) : x = 0, y \in [0, 1] \vee y = 0, x \in [0, 1]\} \tag{3.2.14}$$

ist ähnlich zu (3.2.10), wobei allerdings in den Nebendiagonalen auch der Wert  $-2/h^2$  auftritt. Dies zerstört die Symmetrie der Matrix. Durch Skalierung der Gleichungen (Multiplikation mit  $\frac{1}{2}, \frac{1}{4}$ ) kann die Symmetrie wiederhergestellt werden. Dies entspricht der Verwendung folgender Differenzensterne auf  $\Gamma_n$  (Vorfaktor  $1/h^2$ ):

$$\begin{array}{c}
 -\frac{1}{2} \\
 \hline
 2 \quad \left| \quad -1 \right. \\
 \hline
 -\frac{1}{2}
 \end{array}
 \quad
 \begin{array}{c}
 -1 \\
 \hline
 -\frac{1}{2} \quad \left| \quad -\frac{1}{2} \right. \\
 \hline
 2
 \end{array}
 \quad
 \begin{array}{c}
 -\frac{1}{2} \\
 \hline
 1 \quad \left| \quad -\frac{1}{2} \right. \\
 \hline
 -\frac{1}{2}
 \end{array}$$

Die Matrix ist auch etwas größer, da jetzt unbekannte Werte auf dem größeren Berechnungsgitter  $\Omega_h := \{(x_i, y_j) : i, j = 0, \dots, m-1\}$  auftreten. Bei zeilenweiser Numerierung dieser Unbekannten  $v_{ij} \cong u(x_i, y_j)$ , in der Reihenfolge  $v_{00}, \dots, v_{m-1,0}, v_{0,1}, \dots, v_{m-1,1}, \dots$ , hat die Systemmatrix zu (3.2.12) die Gestalt

$$A = \frac{1}{h^2} \begin{pmatrix}
 \frac{1}{2}T & -E & & & \\
 -E & T & -E & & \\
 & -E & T & -E & \\
 & & & \ddots & \\
 & & & & -E & T
 \end{pmatrix}, \tag{3.2.15}$$

$$T := \begin{pmatrix} 2 & -1 & & & \\ -1 & 4 & -1 & & \\ & -1 & 4 & -1 & \\ & & & \ddots & \\ & & & -1 & 4 \end{pmatrix}, \quad E := \begin{pmatrix} \frac{1}{2} & & & & \\ & 1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{pmatrix}.$$

Die Größe von  $T$  und  $E$  ist  $m \times m$ , die der Gesamtmatrix  $A$  ist  $m^2 \times m^2$ .

Die in (3.2.10) und (3.2.15) beobachtete Symmetrie der Matrizen macht den Einsatz effizienterer Verfahren möglich und ist auch bei der Analyse (Stabilität) hilfreich.

### 3.2.2 Eigenschaften der Differenzen-Matrizen

Zwei Matrix-Eigenschaften und damit verbundene Beweistechniken spielen bei der Behandlung elliptischer Probleme eine wichtige Rolle. Die erste Eigenschaft ist Symmetrie und Definitheit.

**Definition 3.2.2** Eine Matrix  $A = (a_{ij}) \in \mathbf{R}^{n \times n}$  heißt symmetrisch, wenn  $A = A^T$  gilt, d.h.,  $a_{ij} = a_{ji} \forall i, j = 1, \dots, n$ . Sie heißt positiv definit bzw. semi-definit, wenn sie symmetrisch ist und ein  $\alpha > 0$  bzw.  $\alpha \geq 0$  existiert so, daß gilt

$$v^T A v \geq \alpha v^T v \quad \forall v \in \mathbf{R}^n. \quad (3.2.16)$$

Bei einer symmetrisch, positiv definiten ('SPD-') Matrix sind alle Eigenwerte reell positiv, insbesondere ist (3.2.16) mit dem kleinsten Eigenwert,  $\alpha = \min \lambda(A)$ , erfüllt. Hier ist die Verwendung der Euklidnorm und der dazu gehörigen Matrixnorm,

$$\|v\|_2 = \sqrt{v^T v}, \quad v \in \mathbf{R}^n, \quad \|A\|_2 := \sqrt{\varrho(A^T A)},$$

der Spektralnorm angebracht. Dabei ist  $\varrho(M)$  der Spektralradius der Matrix  $M$ , d. h. der Maximalbetrag ihrer Eigenwerte. Da jede symmetrische Matrix unitär ähnlich zu einer reellen Diagonalmatrix ist, gilt  $\|A\|_2 = \varrho(A)$ . Zur Abschätzung der Lösungen von linearen Gleichungssystemen gilt bei Definitheit die einfache Aussage,

$$A \text{ SPD} \Rightarrow A^{-1} \text{ SPD}, \quad \|A^{-1}\|_2 = \varrho(A^{-1}) \leq \frac{1}{\alpha}, \quad \alpha \text{ aus (3.2.16)}.$$

Diese läßt sich auch auf allgemeinere Matrizen übertragen.

**Satz 3.2.3** Für die (i.a. nichtsymmetrischen) Matrizen  $A_j \in \mathbf{R}^{n \times n}$ ,  $j = 1, \dots, k$ , seien mit  $\alpha_j \in \mathbf{R}$  die Ungleichungen  $v^T A_j v \geq \alpha_j v^T v \quad \forall v \in \mathbf{R}^n$  erfüllt. Dann ist  $A = A_1 + \dots + A_k$  invertierbar für  $\alpha_1 + \dots + \alpha_k > 0$  und es gilt

$$\|A^{-1}\|_2 = \|(A_1 + \dots + A_k)^{-1}\|_2 \leq \frac{1}{\alpha_1 + \dots + \alpha_k}.$$

**Beweis** Es sei  $0 \neq v \in \mathbf{R}^n$  beliebig. Aus der Voraussetzung folgt durch Addition  $v^\top Av \geq \alpha v^\top v$  wobei  $\alpha = \alpha_1 + \dots + \alpha_k > 0$ . Mit der Cauchy-Schwarz-Ungleichung führt dies auf die Aussage  $\alpha \|v\|^2 \leq v^\top Av \leq \|v\|_2 \|Av\|_2 \Rightarrow \|Av\|_2 \geq \alpha \|v\|_2$ . Für  $v \neq 0$  ist insbesondere immer  $Av \neq 0$ . Daher ist  $A$  regulär und mit  $v = A^{-1}w$  folgt  $\|A^{-1}w\|_2 \leq \frac{1}{\alpha} \|w\|_2$ . ■

In der reellen quadratischen Form  $v^\top Av$  im Satz spielt insbesondere nur der symmetrische Anteil

$$S(A) := \frac{1}{2}(A + A^\top)$$

der Matrix  $A$  eine Rolle, da  $v^\top Av = (v^\top Av)^\top = v^\top A^\top v = v^\top S(A)v$ . Im Satzes wird also gerade die positive Definitheit des symmetrischen Anteils  $S(A) = S(A_1) + \dots + S(A_k)$  von  $A$  vorausgesetzt. Für die Matrizen des letzten Paragraphen folgt diese einfach aus der Tatsache, daß die Matrix  $A$  des Systems (3.2.9) nach Herleitung eine Summe  $A = A_1 + A_2$  ist, bei der  $A_1$  und  $A_2$  die Differenzenquotienten in  $x$ - bzw.  $y$ -Richtung enthalten. Beim diskreten Dirichlet-Problem (3.2.9) etwa zerfällt daher die quadratische Form  $v^\top Av$ ,  $v \neq 0$ , mit  $A$  aus (3.2.10) in zwei positive Summen,

$$\begin{aligned} & \frac{1}{h^2} \sum_{i,j=1}^{m-1} v_{ij} \left( 4v_{ij} - v_{i-1,j} - v_{i+1,j} - v_{i,j-1} - v_{i,j+1} \right) \\ &= \frac{1}{h^2} \sum_{j=1}^{m-1} \sum_{i=1}^{m-1} v_{ij} (-v_{i-1,j} + 2v_{ij} - v_{i+1,j}) + \frac{1}{h^2} \sum_{i=1}^{m-1} \sum_{j=1}^{m-1} v_{ij} (-v_{i,j-1} + 2v_{ij} - v_{i,j+1}). \end{aligned}$$

Randwerte bei  $v$  sind hier null zu setzen. Zu den einzelnen Summen gilt

**Lemma 3.2.4** Gegeben seien  $v_1, \dots, v_m \in \mathbf{R}$  mit  $v_m = 0$ . Dann gilt für  $a \in \{1, 2\}$

$$v_1(av_1 - v_2) + \sum_{j=2}^{m-1} v_j(-v_{j-1} + 2v_j - v_{j+1}) = (a-1)v_1^2 + \sum_{j=2}^{m-1} (v_j - v_{j-1})^2 + v_{m-1}^2 \geq C_a \sum_{j=1}^{m-1} v_j^2,$$

mit  $C_2 = 4 \sin^2 \frac{\pi}{2m}$ ,  $C_1 > 0$ .

**Beweis** Links steht die quadratische Form  $v^\top T v$  der symmetrischen Tridiagonalmatrix

$$T = \begin{pmatrix} a & -1 & & & \\ -1 & 2 & -1 & & \\ & & \ddots & & \\ & & & -1 & 2 \end{pmatrix} \in \mathbf{R}^{(m-1) \times (m-1)}.$$

Die erste Identität ist elementar, die Ungleichung gilt, da Quadrate summiert werden. Für  $a = 2$  ist die Konstante bekannt, da die Eigenwerte von  $T$  gleich  $4 \sin^2 \frac{j\pi}{2m}$  sind. ■

**Satz 3.2.5** Die Matrix  $A$  (3.2.10) des Gleichungssystems (3.2.9) ist positiv definit, ihre Inverse gleichmäßig beschränkt für  $h \in (0, 1/2]$  mit  $\|A^{-1}\|_2 \leq C_h$ ,  $C_h = \frac{h^2}{8} / \sin^2(h\pi/2) \leq \frac{1}{16}$ . Insbesondere gilt für die Lösung des diskreten Dirichletproblems (3.2.9) mit  $\phi_d = 0$  die Aussage

$$\sum_{i,j=1}^{m-1} v_{ij}^2 \leq C_h^2 \sum_{i,j=1}^{m-1} g_{ij}^2.$$

Die Symmetrie und Definitheit hat auch praktische Bedeutung, da deswegen zur Lösung der linearen Gleichungssysteme effizientere Verfahren als im unsymmetrischen Fall eingesetzt werden können (Cholesky-Zerlegung, spezielle Iterationsverfahren). Ihre Bedeutung im Theoretischen ist nicht so entscheidend, da die Euklidnorm zwar hier die schärfsten Schranken liefert, diese Norm aber bei (diskreten) Funktionen weniger aussagekräftig ist als die Maximumnorm. Für Abschätzungen in dieser Norm ist eine andere Eigenschaft entscheidend, die mit den Vorzeichen der Matrixelemente zu tun hat. Im folgenden sind Ungleichungszeichen ( $\leq, \geq, >, <$ ) zwischen Vektoren und Matrizen sowie Beträge elementweise zu verstehen,  $|v| = (|v_i|)$ ,  $|A| = (|a_{ij}|)$ .

**Definition 3.2.6** Eine Matrix  $A = (a_{ij}) \in \mathbf{R}^{n \times n}$  heißt M-Matrix, wenn gilt

$$a_{ii} > 0 \quad \forall i = 1, \dots, n, \quad a_{ij} \leq 0 \quad \forall i \neq j, \quad (3.2.17)$$

$$A \text{ ist regulär und } A^{-1} \geq 0. \quad (3.2.18)$$

Nach Zerlegung der Matrix  $A = D - N$  in Diagonale  $D$  und Nebendiagonalen  $-N$  lautet die Vorzeichenbedingung (3.2.17) kürzer  $N \geq 0, a_{ii} > 0 \forall i$ . Bei Normbetrachtungen in der Euklidnorm sind Eigenvektoren extremal. Diese Rolle wird hier von positiven Vektoren übernommen. Wichtige Schlußweisen für nichtnegative Vektoren  $v \geq 0$  bzw. Matrizen  $B \geq 0$  sind  $\|v\|_\infty \leq a \iff v \leq a\mathbf{1}$  mit  $\mathbf{1} = (1, \dots, 1)^\top$  und  $\|B\|_\infty \leq b \iff B\mathbf{1} \leq b\mathbf{1}$ . Außerdem ist das Produkt mit  $B$  monoton:  $v \leq w \Rightarrow Bv \leq Bw$ . Kriterien für M-Matrizen enthält

**Lemma 3.2.7** a) Eine Matrix  $A = D - N$  mit der Vorzeichenverteilung (3.2.17) ist genau dann M-Matrix, wenn  $\rho(D^{-1}N) < 1$  gilt.

b) Eine Matrix  $A$  mit der Vorzeichenverteilung (3.2.17) ist (genau dann) M-Matrix, wenn ein Vektor  $z > 0$  existiert mit  $w := Az > 0$ . Damit gilt

$$\|A^{-1}\|_\infty \leq \frac{\max_i z_i}{\min_i w_i}. \quad (3.2.19)$$

**Beweis** a) Für  $\rho(D^{-1}N) < 1$  konvergiert die Neumannreihe

$$A^{-1} = (I - D^{-1}N)^{-1}D^{-1} = \sum_{j=0}^{\infty} \underbrace{(D^{-1}N)^j}_{\geq 0} D^{-1} \geq 0,$$

die Inverse existiert und ist nichtnegativ. Sei umgekehrt  $A$  M-Matrix und  $v \neq 0$  ein Eigenvektor von  $D^{-1}N$  zum Eigenwert  $\lambda$ . Wegen  $N \geq 0$  gilt  $|Nv| \leq N|v|$ . Aus  $\lambda Dv = Nv$  folgt daher

$$|\lambda Dv| = |\lambda|D|v| \leq N|v| \quad \Rightarrow \quad (D - N)|v| \leq (1 - |\lambda|)D|v|.$$

Multiplikation mit  $A^{-1} \geq 0$  führt dann auf die Ungleichung  $0 \leq |v| \leq (1 - |\lambda|)A^{-1}D|v|$ , die für  $v \neq 0$  wegen  $A^{-1}D|v| \geq 0$  die Folgerung  $1 - |\lambda| > 0$  liefert.

b) Zu  $z > 0$  wird die reguläre Matrix  $Z := \text{diag}(z_i)$  definiert. Damit ist

$$Az = Dz - Nz > 0 \iff NZ\mathbf{1} < DZ\mathbf{1} \iff Z^{-1}D^{-1}NZ\mathbf{1} < \mathbf{1}.$$

Die letzte Ungleichung für diese nichtnegative Matrix bedeutet aber gerade  $\|Z^{-1}D^{-1}NZ\|_\infty < 1$ , woraus insbesondere  $\rho(Z^{-1}D^{-1}NZ) = \rho(D^{-1}N) < 1$  folgt. Nach Teil a) ist  $A$  daher M-Matrix. Mit  $\alpha := \min_i w_i$  gilt laut Voraussetzung auch  $Az = w \geq \alpha \mathbf{1}$ . Aus  $A^{-1} \geq 0$  folgt damit  $\alpha A^{-1} \mathbf{1} \leq z \leq \|z\|_\infty \mathbf{1}$ . Dies ist äquivalent zur Behauptung (3.2.19). ■

Für Matrizen mit (3.2.17) ist die strenge Diagonaldominanz

$$a_{ii} > \sum_{j \neq i} |a_{ij}| \quad \forall i,$$

ein spezielles M-Matrix-Kriterium und mit  $z = \mathbf{1}$  in Lemma 3.2.7 enthalten. Abschwächungen der Diagonaldominanz sind aufwendiger in der Formulierung und liefern nur die Regularitätsaussage, aber keine Schranke wie (3.2.19). Der in Lemma 3.2.7 benötigte Testvektor  $z$  kann bei  $L_\Delta$  einfach angegeben werden. Dies führt direkt auf eine Stabilitätsaussage für die betrachteten Modellprobleme auf dem Einheitsquadrat.

**Satz 3.2.8** Für das diskrete Dirichletproblem (3.2.9) auf dem Einheitsquadrat mit homogenen Randbedingungen gilt  $\|L_\Delta^{-1}\|_\infty \leq \frac{1}{8}$ , d.h., die Lösung  $v$  erfüllt

$$\max_{i,j=1}^{m-1} |v_{ij}| \leq \frac{1}{8} \max_{i,j=1}^{m-1} |g_{ij}|.$$

**Beweis** Für  $z_\Delta(x, y) := r - (x - \frac{1}{2})^2 - (y - \frac{1}{2})^2$  ergibt sich sofort  $L_\Delta z_\Delta(x, y) \geq 4$  mit Gleichheit in allen randfernen Punkten. Bei Wahl von  $r = \frac{1}{2}$  wird  $z_\Delta$  positiv,  $0 < z_\Delta \leq r$  in  $\Omega_h$ . Aus Lemma 3.2.7 folgt dann die Aussage. ■

Für das Problem (3.2.12) mit Neumann-Bedingung auf zwei Seiten des Quadrats gilt eine analoge Aussage mit der größeren Konstanten 1, beim vollständigen Neumann-Problem ist die Matrix natürlich singulär.

### 3.2.3 Diskretisierungsfehler, Konvergenz

Da die Konsistenz des Differenzenverfahrens (für  $C^4$ -Lösungen) aus Lemma 1.7.1 folgt und die Stabilität im letzten Abschnitt geklärt wurde, kann daraus sofort die Konvergenz der Näherungen abgeleitet werden, wenn die Lösung des Randwertproblems hinreichend glatt ist. Allerdings ist dies nicht immer der Fall, wie in Beispiel 3.1 gezeigt wurde. Daher wird zunächst eine abgeschwächte Konsistenzaussage behandelt, danach Möglichkeiten zur Beeinflussung des Fehlers durch Gitteranpassung oder Verbesserung des Verfahrens. Im folgenden Satz wird der eigentlich nur für Gitterfunktionen ( $\Omega_h \subseteq \Omega$ ) definierte Differenzenoperator  $L_\Delta$  auch auf normale Funktionen angewendet (formal wäre eine Gitter-Restriktion  $R_\Delta$  erforderlich:  $L_\Delta R_\Delta u$ , vgl. §2).

**Satz 3.2.9** Es sei  $u \in C^k(\Omega)$ ,  $2 \leq k \leq 4$ . Dann gilt für  $(x, y) \in \Omega_h$  die Konsistenz-Aussage

$$\begin{aligned} |(L_\Delta u)(x, y) + u_{xx}(x, y) + u_{yy}(x, y)| &\leq \\ Ch^{k-2} \left( \max\left\{ \left| \frac{\partial^k}{\partial x^k} u(\xi, y) \right| : |\xi - x| < h \right\} + \max\left\{ \left| \frac{\partial^k}{\partial y^k} u(x, \eta) \right| : |\eta - y| < h \right\} \right). \end{aligned} \quad (3.2.20)$$

**Beweis** Taylorentwicklung von  $u$  um  $(x, y)$ . ■

Wird im Satz die Lösung  $u$  der Poissongleichung eingesetzt, ergibt sich für den lokalen Fehler  $T_\Delta := L_\Delta u - g_\Delta$  die Größenordnung  $T_\Delta = \mathcal{O}(h^2)$ , wenn  $u \in C^4$  gilt. Wegen der Linearität von  $L_\Delta$  führt die Subtraktion von (3.2.9), auf eine Gleichung für den Fehler  $u_\Delta - u$

$$L_\Delta u_\Delta = g_\Delta \quad \Rightarrow \quad -T_\Delta = g_\Delta - L_\Delta u = L_\Delta u_\Delta - L_\Delta u = L_\Delta(u_\Delta - u). \quad (3.2.21)$$

Die tatsächliche Größe dieses globalen Fehlers kann dann mit den Stabilitäts-Sätzen 3.2.5 (Euklidnorm) und 3.2.8 (Maximumnorm) abgeschätzt werden. Bei der *Quadratsummennorm* muß allerdings auf die richtige Skalierung geachtet werden. Denn bei Gittervektoren  $g_{ij}$  ist die reine Summe  $\sum_{i,j} g_{ij}^2$  kein gutes Maß für die Größe von  $g$ , da sich wegen der  $(m-1)^2$  Summanden bei einer konstanten Funktion der Wert  $(m-1)^2 g^2$  ergibt. Daher werden nun für Gitterfunktionen  $g_\Delta$  folgende skalierten Normen definiert,

$$\|g_\Delta\|_{\Delta,2} = \left( h^2 \sum_{(x,y) \in \Omega_h} g_\Delta^2(x,y) \right)^{1/2}, \quad \|g_\Delta\|_{\Delta,\infty} = \max_{(x,y) \in \Omega_h} |g_\Delta(x,y)|. \quad (3.2.22)$$

Beide Normen sind Approximationen für die jeweiligen Normen von Funktionen,

$$\|g\|_2 = \left( \int_0^1 \int_0^1 g(x,y)^2 dx dy \right)^{1/2}, \quad \|g\|_\infty = \sup_{(x,y) \in \Omega} |g(x,y)|,$$

insbesondere gilt so  $\|g_\Delta\|_{\Delta,2} \leq \|g_\Delta\|_{\Delta,\infty}$ . Durch Kombination der Ergebnisse aus Formel (3.2.21) und den Sätzen 3.2.5, 3.2.8, 3.2.9, folgt die globale Fehleraussage:

**Satz 3.2.10** Die Lösung  $u$  des Dirichletproblems auf dem Einheitsquadrat erfülle  $u \in C^4(\bar{\Omega})$ . Dann gilt für den Fehler der Näherungslösung  $u_\Delta$  aus dem Differenzenverfahren (3.2.9) mit Schrittweite  $h$  in jeder der Normen aus (3.2.22) die Abschätzung

$$\|u_\Delta - u\|_{\Delta,p} \leq C_p h^2 (\|u_{xxxx}\|_\infty + \|u_{yyyy}\|_\infty),$$

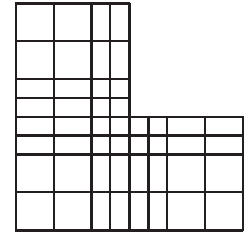
$p \in \{2, \infty\}$ , mit  $C_2 = 1/192$  ( $h \leq \frac{1}{2}$ ),  $C_\infty = 1/96$ .

**Beweis** Da die Werte von  $u$  und  $u_\Delta$  auf dem Rand  $\Gamma_h$  übereinstimmen, treten in der Fehlergleichung (3.2.21) nur Beiträge aus dem Konsistenzfehler  $T_\Delta$  auf, der mit Satz 3.2.9 (Konstante  $C = 1/12$ ) abgeschätzt werden kann. Daher sind die Sätze 3.2.5 und 3.2.8 anwendbar (homogene Randbedingungen) und ergeben die Behauptung. ■

Für das gemischte Randwertproblem (3.2.12) erhält man nach dem selben Prinzip eine ähnliche Aussage mit etwas größeren Konstanten.

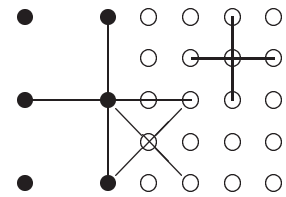
Nur auf den ersten Blick ist diese Fehleraussage zufriedenstellend. Denn in wichtigen Fällen erfüllt die Lösung  $u$  nicht die Voraussetzungen des Satzes, da dort von  $g \in C^2(\bar{\Omega})$  nur auf  $u \in C^4(\Omega)$ , ( $\Omega$  offen!) geschlossen werden kann (vgl. Beispiel 3.1). In diesem Beispiel wächst der lokale Konsistenzfehler (3.2.20) zum Rand hin stark an. Bei den gewöhnlichen Differentialgleichungen wurden solche Effekte durch Verkleinerung der lokalen Schrittweite kompensiert.

Überträgt man dieses Prinzip auf ein Produkt-Gitter  $(x_i, y_j)$ , indem man in bestimmten Teilen des  $x$ - und  $y$ -Intervalls die Punkte  $x_i, y_j$  dichter wählt, kann zwar an einer beliebigen Stelle im Gebiet ein feineres zweidimensionales Gitter erzeugt werden. Allerdings ergeben sich dann auch außerhalb dieses Bereichs Regionen mit kleinen Schrittweiten, die in  $x$ - und  $y$ -Richtung außerdem sehr unterschiedlich sind.



Auf einfache Weise läßt sich eine lokale Verfeinerung in einem kleinen Bereich dadurch erreichen, daß man ein Gitter mit halber Schrittweite  $h/2$  in das grobe Gitter 'einhängt'. An einer Übergangsstelle kann das folgendermaßen skizziert werden:

Im groben Gitter (ausgefüllte Punkte) wird der normale Fünfpunktstern mit Schrittweite  $h$  verwendet, der an den Übergangsstellen den ersten Punkt des feinen Gitters überspringt. In den Punkten des feinen Gitters (offene Kreise) wird der Stern mit Schrittweite  $h/2$  eingesetzt. In den beiden Punkten des feinen Gitters, die keinen westlichen Nachbarn haben, kann ein diagonaler



Fünf-Punkte-Stern mit Schrittweite  $h/\sqrt{2}$  benutzt werden (vgl. Lemma 3.2.12). Alle diese Sterne haben die Approximationsordnung 2. Das Verfahren führt nicht mehr auf eine symmetrische Matrix, allerdings läßt sich die Inverse in der  $\infty$ -Norm wieder mit Lemma 3.2.7 und dem Testvektor aus Satz 3.2.8 abschätzen. Die lokale Gitterverfeinerung kann in der Nähe einer Problemstelle natürlich mehrfach angewendet werden.

**Beispiel 3.2.11** Auf dem L-Gebiet  $\Omega = (0, 1) \times (0, 1) \setminus [\frac{1}{2}, 1) \times [\frac{1}{2}, 1)$  gelte  $-u_{xx} - u_{yy} = 25$ , auf dem Rand  $u = 0$ . Einsatz ineinandergeschachtelte Gitter mit  $h = 1/8, 1/16, 1/32$ . Der steile Anstieg in der einspringenden Ecke wird erst durch die lokal feinere Unterteilung sichtbar.

Wenn die (höheren) Ableitungen der Lösung zu groß werden, kann dies also durch Gitterverfeinerung kompensiert werden. Umgekehrt ist es natürlich sinnvoll, den Aufwand durch Verfahren höherer Konvergenzordnung zu senken, wenn die höheren Ableitungen der Funktion gutartig sind. Die Kenntnis der genauen Fehlerstruktur des oben verwendeten diagonalen Sterns kann dazu verwendet werden.

**Lemma 3.2.12** Für eine Funktion  $u \in C^6(\bar{\Omega})$  gilt

$$\begin{aligned} & \frac{1}{2h^2} \left( 4u(x, y) - u(x + h, y + h) - u(x + h, y - h) - u(x - h, y + h) - u(x - h, y - h) \right) \\ &= - \left( u_{xx} + u_{yy} + \frac{h^2}{12} (u_{xxxx} + 6u_{xxyy} + u_{yyyy}) \right) |_{(x,y)} + \mathcal{O}(h^4), \end{aligned}$$

**Beweis** Übungsaufgabe.

Zum Vergleich lautet der Fehler des normalen Fünf-Punkte-Sterns (Satz 3.2.9)

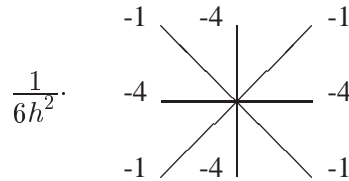
$$\frac{h^2}{12} (u_{xxxx}(x, y) + u_{yyyy}(x, y)) + \mathcal{O}(h^4).$$

Bei der Poissongleichung kann nun durch eine geeignete Linearkombination dieser beiden Sterne ein  $h^2$ -Fehlerterm erzeugt werden, der sich mit Hilfe der Dgl wieder approximieren läßt. Wendet man nämlich

den Laplace-Operator noch einmal auf die Poissongleichung an, folgt

$$-g_{xx} - g_{yy} = u_{xxxx} + 2u_{xxyy} + u_{yyyy}.$$

Bei einer Linearkombination, die den einfachen Fünfpunktstern mit  $2/3$  und den diagonalen mit  $1/3$  gewichtet, ergänzt sich der  $h^2$ -Fehleranteil gerade zu dem Ausdruck  $-g_{xx} - g_{yy}$ . Dies führt auf den *Neun-Punkte-Stern*

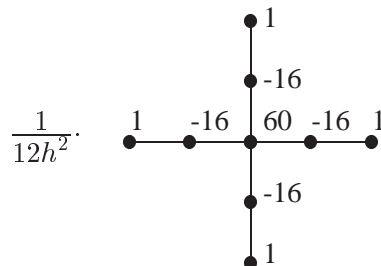


Das Gewicht im Mittelpunkt ist  $20/(6h^2)$ . Der Ausdruck  $\frac{h^2}{12}(g_{xx} + g_{yy})$  im führenden Fehlerbeitrag wird natürlich wieder durch den Fünfpunktstern approximiert. Dies führt auf das folgende *Mehrstellenverfahren*. Eine einzelne Gleichung desselben im Punkt  $(x, y)$  wird der Übersicht halber durch Angabe der Gewichte in den Sternen abgekürzt:

$$\frac{1}{6h^2} \begin{bmatrix} -1 & -4 & -1 \\ -4 & 20 & -4 \\ -1 & -4 & -1 \end{bmatrix} u_{\Delta}(x, y) = \frac{1}{12} \begin{bmatrix} & 1 & \\ 1 & 8 & 1 \\ & 1 & \end{bmatrix} g(x, y), \quad (x, y) \in \Omega_h. \quad (3.2.23)$$

Der lokale Fehler hat die Form  $\mathcal{O}(h^4)$  für  $u \in C^6$ , bei der Potentialgleichung mit  $g \equiv 0$  besitzt das Verfahren sogar die Ordnung 6. Die Stabilität beim Dirichletproblem kann wieder mit Lemma 3.2.7 wie in Satz 3.2.8 gezeigt werden mit dem gleichen Ergebnis,  $\|L_{\Delta}^{-1}\|_{\infty} \leq \frac{1}{8}$ . Der wesentliche Vorteil des Verfahrens ist die Kompaktheit seines Differenzsterns, der bei Anwendung nur geringe Modifikationen gegenüber dem Fünfpunktstern erfordert (Randbedingungen, Matrixstruktur). Da er aber auf der speziellen Form der Poissongleichung beruht, kann er nur in diesem Spezialfall verwendet werden.

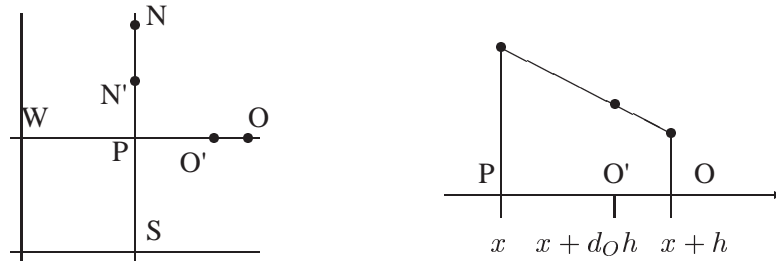
Eine auch bei allgemeineren Dgln verwendbare Methode zur Erhöhung der Konsistenzordnung ist der Einsatz von Differenzenausdrücken mit mehr als drei Punkten in jeder Ortsrichtung. In (1.7.7) wurde eine symmetrische Approximation für  $u_{xx}$  mit fünf Funktionswerten angegeben. Die Kombination mit der entsprechenden Approximation für  $u_{yy}$  führt bei der Poissongleichung auf den Neun-Punkte-Stern



mit einem  $\mathcal{O}(h^4)$ -Fehler. Da hier jetzt Werte verknüpft werden mit einem maximalen Abstand von 4 Gitterpunkten, ergeben sich erhebliche Probleme am Rand. Insbesondere bei nicht achsenparallelen Ecken können diese Schwierigkeiten auch nicht immer durch Verwendung unsymmetrischer Approximationen behoben werden. Dieses Verfahren wird nicht weiter betrachtet.

### 3.2.4 Allgemeinere Gebiete und Gleichungen

Der Fünfpunktstern kann in allen Gebieten verwendet werden, dessen Randpunkte auf ein geeignetes quadratisches oder rechteckiges ( $h_x \neq h_y$ ) Gitter  $\Delta_h$  fallen. Für nichtpolygonale Ränder ('krummlinige') läßt es sich aber meist nicht mehr erreichen, daß alle Nachbarn eines inneren Gitterpunkts entweder wieder Gitter- oder Randpunkte sind. Daher ist nun der Fall zu betrachten, daß ein oder mehrere Nachbarn eines (randnahen) Gitterpunkts  $P$  außerhalb von  $\Omega$  liegen wie in der gezeigten Skizze.



Der Rand  $\Gamma$  verlaufe also zwischen  $P$  und  $O$  durch den Punkt  $O'$  und (oder) zwischen  $P$  und  $N$  durch  $N'$ . Als allgemeinen Zugang kann man hier durch Taylor-Abgleich die Gewichte eines Fünfpunktsterns zu den Punkten  $P, N', O', S, W$  als Approximation an  $-(u_{xx} + u_{yy})$  bestimmen. Auf das gleiche Ergebnis kommt man, wenn Gitterwerte außerhalb von  $\Omega$ , also  $v_N, v_O$  durch lineare Extrapolation des inneren Wertes  $v_P$  und der bekannten Randwerte  $\phi_d(N'), \phi_d(O')$  (beim Dirichletproblem) approximiert werden. Wenn  $d_N h, d_O h, 0 \leq d_N, d_O \leq 1$  die Abstände der Randpunkte von  $P$  sind, dann ist

$$\phi_d(O') = d_O v_O + (1 - d_O) v_P \Rightarrow v_O = [\phi_d(O') - (1 - d_O) v_P] / d_O.$$

Durch Einsetzen in die zweite Differenz in  $x$ -Richtung, Anwendung auf  $u$  und Taylorentwicklung führt dies auf

$$\begin{aligned} & -u(x-h, y) + \left(1 + \frac{1}{d_O}\right)u(x, y) - \frac{1}{d_O}u(x + d_O h, y) = \\ & -u(W) + \left(1 + \frac{1}{d_O}\right)u(P) - \frac{1}{d_O}u(O') = -\frac{h^2}{2}(1 + d_O)u_{xx} + \frac{h^3}{6}(1 - d_O^2)u_{xxx} + \dots \end{aligned} \quad (3.2.24)$$

Division durch  $h^2(1 + d_O)/2$  liefert eine Approximation der zweiten Ableitung. Analoges gilt in  $y$ -Richtung. Im oben skizzierten Beispiel ergibt sich so im Punkt  $P$  die Gleichung

$$\frac{1}{h^2} \left[ \left( \frac{2}{d_O} + \frac{2}{d_N} \right) v_P - \frac{2}{1 + d_O} v_W - \frac{2}{d_O(1 + d_O)} v_{O'} - \frac{2}{1 + d_N} v_S - \frac{2}{d_N(1 + d_N)} v_{N'} \right] = g_P. \quad (3.2.25)$$

Da die Gewichte bei den Gitternachbarn  $v_S, v_W$  von den  $d$ -Werten abhängen, führt diese Approximation nicht mehr auf eine symmetrische, aber immer noch auf eine M-Matrix, da die Gleichung (3.2.25) diagonaldominant ist. Die Normabschätzung in der  $\infty$ -Norm aus §3.2.2 kann hier wiederholt werden, wobei die Konstanten vom Gebiet abhängen.

Bei der Diskussion des Gesamtfehlers ist allerdings eine ausführlichere Argumentation notwendig, da die (Konsistenz-) Ordnung des jetzt asymmetrischen Sterns (3.2.25) nur noch eins ist, wie man an der Taylorentwicklung (3.2.24) erkennt. Dennoch kann weiter  $\mathcal{O}(h^2)$ -Konvergenz des Gesamtverfahrens gezeigt werden, da der größere lokale Fehler nur am Rand auftritt. Hierbei spielt die M-Eigenschaft der Matrix wieder die entscheidende Rolle, da mit ihr Vergleiche zwischen Vektoren übertragen werden können. Es bezeichne  $L_\Delta$  wieder die lineare Abbildung, die sich beim Dirichletproblem bei Anwendung des normalen Fünf-Punkte-Sterns in inneren Gitterpunkten und einer Gleichung der Form (3.2.25) in randnahen Punkten ergibt. Bei Anwendung dieser Abbildung auf die konstante Funktion  $e_\Delta(x, y) \equiv 1$  liefert das Bild  $L_\Delta e_\Delta$  (die Zeilensumme der Matrix) nur in den randnahen Punkten einen nichtverschwindenden Beitrag, da die Dirichletrandwerte im Gleichungssystem  $L_\Delta u_\Delta = g_\Delta$  auf der rechten Seite berücksichtigt werden (vgl.(3.2.10)). Speziell in (3.2.25) ergibt dies

$$(L_\Delta e_\Delta)(P) = \frac{1}{h^2} \left[ \frac{2}{d_O} + \frac{2}{d_N} - \frac{2}{1 + d_O} - \frac{2}{1 + d_N} \right] = \frac{2}{h^2} \left[ \frac{1}{d_O(1 + d_O)} + \frac{1}{d_N(1 + d_N)} \right] > 0. \quad (3.2.26)$$

Bei einer anderen Anzahl von Randpunkten im Stern von  $P$  ist dieser Ausdruck zu modifizieren. Der Konsistenzfehler des Verfahrens hat hier die Form  $T_\Delta = L_\Delta u - g_\Delta = hr_\Delta + h^2 t_\Delta$ , wobei der  $h^2$ -Anteil dem des normalen Sterns aus Satz 3.2.9 entspricht, während  $hr_\Delta$  nur Randfehler enthält. Im Randpunkt  $P$  liefert Gleichung (3.2.25) nach (3.2.24) hierzu den Beitrag

$$hr_\Delta(P) = \frac{h}{3} \left[ (1 - d_O)u_{xxx}(x + \xi h, y) + (1 - d_N)u_{yyy}(x, y + \eta h) \right], \quad \xi, \eta \in (0, 1).$$

Mit  $K_j := \|\partial^j u / \partial x^j\|_\infty + \|\partial^j u / \partial y^j\|_\infty$ ,  $j = 3, 4$ , kann dieser Beitrag mit (3.2.26) verglichen werden, dabei gilt

$$h|r_\Delta(P)| \leq K_3 \frac{h}{3} [1 - d_O + 1 - d_N] \leq K_3 \frac{h}{3} \left[ \frac{1}{1 + d_O} + \frac{1}{1 + d_N} \right] \leq K_3 \frac{h^3}{6} (L_\Delta e_\Delta)(P).$$

Aus der üblichen Fehlergleichung (3.2.21)  $L_\Delta(u_\Delta - u) = -T_\Delta = -(hr_\Delta + h^2t_\Delta)$  kann man damit wegen  $L_\Delta^{-1} \geq 0$  komponentenweise Abschätzungen herleiten. In jedem Gitterpunkt gilt

$$\begin{aligned} |u_\Delta - u| &= |L_\Delta^{-1}(hr_\Delta + h^2t_\Delta)| \leq L_\Delta^{-1}(h|r_\Delta| + h^2|t_\Delta|) \leq L_\Delta^{-1}(K_3 \frac{h^3}{6} L_\Delta e_\Delta + h^2|t_\Delta|) \\ &\leq (K_3 \frac{h^3}{6} + K_4 \frac{h^2}{12} \|L_\Delta^{-1}\|_\infty) e_\Delta, \end{aligned}$$

also  $\|u_\Delta - u\|_\infty \leq Ch^2$  für den globalen Fehler.

**Beispiel 3.2.13** Auf dem Viertelkreis  $\Omega = \{(x, y) : x^2 + y^2 < 1, x > 0, y > 0\}$  sei  $-u_{xx} - u_{yy} = 4$  sowie  $u = 0$  auf dem Kreis und  $u_x(0, y) = 0, u_y(x, 0) = 0$ . Da die Lösung  $u(x, y) = 1 - x^2 - y^2$  ist, arbeitet das Verfahren exakt.

Bei einer *allgemeinen linearen, elliptischen Dgl 2.Ordnung*

$$au_{xx} + 2bu_{xy} + du_{yy} + eu_x + fu_y + qu = g \tag{3.2.27}$$

( $a, d > 0, ad - b^2 > 0$ ) ist gegenüber den schon aus §1.7 und §3.2.1 bekannten Approximationen nur noch die der gemischten Ableitung  $u_{xy}$  nachzutragen. Diese kann aus dem Produkt der ersten Differenzen in  $x$ - und  $y$ -Richtung mit Schrittweite  $2t$  zusammengesetzt werden. Dieser einfache Stern für  $u_{xy}$  hat die Form

$$\frac{1}{4t^2} \begin{matrix} & -1 & & 1 \\ & \diagdown & & / \\ & 1 & & -1 \end{matrix} u(x, y) := \frac{1}{4t^2} [u(x+t, y+t) + u(x-t, y-t) - u(x+t, y-t) - u(x-t, y+t)]$$

Er ist als symmetrischer Stern mit Zentrum  $(x, y)$  und Schrittweite  $2t = 2h$  eine  $O(h^2)$ -Approximation an  $u_{xy}(x, y)$ . Beim Einsatz dieses Sterns besitzen dann aber die Nebendiagonalelemente der entstehenden Matrix (die Gewichte der Sternstrahlen) nicht mehr alle das gleiche Vorzeichen, was für die M-Matrix-Eigenschaft entscheidend war. Einheitliche Vorzeichen auf den Raumdiagonalen können durch die Verwendung von zwei anderen Sternen in Abhängigkeit vom Vorzeichen von  $b$  erreicht werden. Dazu setzt man jeweils zwei der vier Sterne mit halber Achsenlänge  $2t = h$ , die um  $(x \pm h/2, y \pm h/2)$  zentriert sind, durch Mittelung zu einem um  $(x, y)$  zentrierten zusammen, der einen Fehler  $O(h^2)$  besitzt:

$$\frac{1}{2h^2} \begin{matrix} & -1 & & 1 \\ & \diagdown & & / \\ & 1 & & -1 \end{matrix} \quad \text{oder} \quad \frac{1}{2h^2} \begin{matrix} & -1 & & 1 \\ & \diagdown & & / \\ & 1 & & -1 \end{matrix}$$

Von diesen Sternen wird der linke verwendet in Punkten, wo  $b > 0$ , der rechte, wenn  $b < 0$  ist:

$$u_{xy} \cong \frac{1}{2h^2} \begin{matrix} & -1 & & 1 \\ & \diagdown & & / \\ & 1 & & -1 \end{matrix} \cdot u, \quad u_{xy} \cong \frac{1}{2h^2} \begin{matrix} & -1 & & 1 \\ & \diagdown & & / \\ & 1 & & -1 \end{matrix} \cdot u.$$

Gemeinsam führen diese Approximationen für die allgemeine Dgl (3.2.27) in jedem Gitterpunkt  $(x, y)$  auf folgende Differenzgleichung, die wieder in Stern/Matrixform geschrieben wird. Alle Koeffizienten sind ebenfalls in  $(x, y)$  auszuwerten. Dabei bedeutet  $b^+ := \max\{b, 0\} \geq 0$ ,  $b^- := \min\{b, 0\} \leq 0$ .

$$\frac{1}{h^2} \begin{bmatrix} -b^- & d - |b| & b^+ \\ a - |b| & -2(a + d - |b|) & a - |b| \\ b^+ & d - |b| & -b^- \end{bmatrix} u + \frac{1}{2h} \begin{bmatrix} f & & \\ -e & & e \\ & -f & \end{bmatrix} u + qu = g. \quad (3.2.28)$$

Insbesondere reicht also bei allgemeineren Gleichungen ein Fünfpunktstern nicht mehr aus. Unter einigen Einschränkungen an die Koeffizienten, die einerseits eine Verschärfung der Elliptizität der Dgl darstellen (vgl. Def. 3.1.5) bzw. der Schrittweitenbeschränkung aus Satz 1.7.2 entsprechen, läßt sich auch hier wieder zeigen, daß der Stern (3.2.28) eine M-Matrix erzeugt.

**Satz 3.2.14** Für die Koeffizienten der Dgl (3.2.27) gelte in jedem Gitterpunkt die Einschränkung

$$|b| + \frac{h}{2}|e| < a, \quad |b| + \frac{h}{2}|f| < d, \quad q \geq 0.$$

Für Gitter, bei denen alle in den Sternen (3.2.28) auftretenden Gitterpunkte im Innern oder auf dem Rand von  $\Omega$  liegen, ist beim Dirichletproblem die zugehörige negative Koeffizientenmatrix eine M-Matrix.

Bei allgemeineren Gebieten oder Randbedingungen sind außerdem in jedem randnahen Punkt spezielle Randsterne analog zu (3.2.25) zu erstellen. Das im nächsten Paragraphen besprochene Verfahren ist im Vergleich dazu wesentlich flexibler in Bezug auf die Geometrie des Problems, allerdings nicht in Bezug auf die Gestalt der Dgl (Ableitungen, Nichtlinearität).

### 3.3 Finite-Elemente-Verfahren für elliptische Probleme

#### 3.3.1 Variationsformulierung

Ein fundamentales physikalisches Prinzip besagt, daß jedes physikalische System den Zustand einnimmt, in dem seine Gesamtenergie minimal ist. Dieser Zustand kann dann *auch* durch eine Differentialgleichung beschrieben werden. Im eindimensionalen sei  $u \in C^2[0, 1]$  die Lösung des folgenden Randwertproblems

$$-(p(x)u'(x))' + q(x)u(x) = g(x), \quad x \in (0, 1), \quad u(0) = u(1) = 0, \quad (3.3.1)$$

mit Koeffizienten-Funktionen  $p \in C^1[0, 1]$ ,  $q \in C[0, 1]$ ,  $p(x) \geq p^* > 0$ . Die Randbedingungen wurden zur Vereinfachung homogen gewählt, da dann die Räume

$$C_0^k[0, 1] := \{v \in C^k[0, 1] : v(0) = v(1) = 0\}, \quad k = 0, 1, 2,$$

linear sind. In (3.3.1) ist also  $u \in C_0^2[0, 1]$  gesucht. Wenn auf beiden Seiten der Dgl in (3.3.1) das Innenprodukt mit einem beliebigen anderen Element  $v \in C_0^2[0, 1]$  gebildet wird, folgt durch partielle Integration

$$(g, v)_2 := \int_0^1 g(x)v(x)dx = \int_0^1 [-(pu')'v + quv]dx = \underbrace{-[pu'v]_0^1}_{=0} + \int_0^1 [pu'v' + quv]dx.$$

Der Ausdruck auf der rechten Seite ist schon für einmal differenzierbare Funktionen erklärt, damit läßt sich folgende symmetrische *Bilinearform*  $a : C_0^1[0, 1] \times C_0^1[0, 1] \rightarrow \mathbf{R}$  definieren

$$a(u, v) := \int_0^1 [p(x)u'(x)v'(x) + q(x)u(x)v(x)] dx, \quad u, v \in C_0^1[0, 1]. \quad (3.3.2)$$

Das Innenprodukt  $(g, v)_2$  ist ein beschränktes lineares Funktional, das mit  $f$  bezeichnet wird,  $f : C[0, 1] \rightarrow \mathbf{R}$ ,  $v \mapsto (g, v)_2$ . Aus der obigen Umformung folgt also, daß die Lösung  $u$  des RWP's (3.3.1) auch folgende Gleichung erfüllt

$$a(u, v) = f(v) \quad \forall v \in C_0^1[0, 1]. \quad (3.3.3)$$

Für  $p > 0, q \geq 0$  ist die Bilinearform  $a(u, v)$  definit, stellt auf  $C_0^1[0, 1]$  also ein Innenprodukt dar und  $\sqrt{a(v, v)}$  eine Norm. Aufgrund dieser Tatsache ist die Lösung  $u$  auch Minimalstelle des *konvexen Funktionals*

$$J(v) := \frac{1}{2}a(v, v) - f(v), \text{ d.h., } J(u) = \inf\{J(v) : v \in C_0^1[0, 1]\}. \quad (3.3.4)$$

Mit beliebigem  $0 \neq v \in C_0^1[0, 1], 0 \neq t \in \mathbf{R}$ , gilt für das Element  $u + tv$  nämlich

$$\begin{aligned} J(u + tv) &= \frac{1}{2}a(u + tv, u + tv) - f(u + tv) \\ &= \frac{1}{2}a(u, u) - f(u) + t[a(u, v) - f(v)] + \frac{t^2}{2}a(v, v) = J(u) + \frac{t^2}{2}a(v, v) \\ &> J(u). \end{aligned} \quad (3.3.5)$$

Umgekehrt wäre für  $a(u, v) - f(v) \neq 0$  auch  $u$  keine Minimalstelle. Daher stimmen die Lösungen des RWP's (3.3.1), von (3.3.3) und (3.3.4) überein, wenn sie in  $C_0^2$  liegen. Es muß aber betont werden, daß bei weniger glatten Koeffizienten  $p, q, g$  das Minimum bei (3.3.4) in einer Funktion  $u \notin C^2[0, 1]$  angenommen werden kann. Daher heißen Funktionen, die (3.3.3) oder (3.3.4) erfüllen, verallgemeinerte oder *schwache* Lösungen. Für eine korrekte Formulierung müssen dann andere Räume zugrundegelegt werden, vgl. (3.3.10).

Auch die Lösungen von partiellen, elliptischen Randwertproblemen sind Minimalstellen *konvexer Funktionale*. Es sei wieder  $\Omega \subseteq \mathbf{R}^2$  ein beschränktes Gebiet mit einem Rand  $\Gamma$  aus stückweise stetig differenzierbaren Kurven. Nach der Greenschen Formel gilt hier die Beziehung

$$\int \int_{\Omega} [u_x v_x + u_y v_y] dx dy = - \int \int_{\Omega} [u_{xx} + u_{yy}] v dx dy + \oint_{\Gamma} \frac{\partial u}{\partial \mathbf{n}} v ds \quad (3.3.6)$$

für  $u \in C^2(\bar{\Omega})$  und  $v \in C^1(\bar{\Omega})$ . Diese Relation spielt wieder eine Rolle, wenn man Störungen  $u + tv$  einer Lösung  $u$  betrachtet. Speziell beim Dirichletproblem müssen  $u$  und  $u + tv$  die gleichen Randwerte besitzen, also  $v$  auf dem Rand verschwinden. In diesem Fall entfällt daher das Randintegral in (3.3.6). Zur Vereinfachung wird wieder das Dirichletproblem bei der (etwas erweiterten) Poissongleichung mit  $q \geq 0$  diskutiert,

$$-u_{xx} - u_{yy} + qu = g \text{ in } \Omega, \quad u = 0 \text{ auf } \Gamma. \quad (3.3.7)$$

Zu diesem Randwertproblem kann schon für Funktionen aus dem Raum  $C_0^1(\Omega)$  mit

$$C_0^k(\Omega) := \{u \in C^k(\Omega) : u|_{\Gamma} \equiv 0\}, k = 0, 1,$$

die Bilinearform  $a : C_0^1(\Omega) \times C_0^1(\Omega) \rightarrow \mathbf{R}$  und die Linearform  $f : C_0^1(\Omega) \rightarrow \mathbf{R}$  gebildet werden,

$$\begin{aligned} a(u, v) &:= \int \int_{\Omega} [u_x v_x + u_y v_y + q u v] dx dy, \\ f(v) &:= (g, v)_2 = \int \int_{\Omega} g v dx dy, \end{aligned} \quad u, v \in C_0^1(\Omega). \quad (3.3.8)$$

Aufgrund der Greenschen Formel gilt also wieder:

$$u \in C_0^2(\Omega) \text{ löst (3.3.7)} \Rightarrow a(u, v) = f(v) \quad \forall v \in C_0^1(\Omega).$$

Auch Neumann- oder Cauchy-Randbedingungen können in einer ähnlichen Formulierung berücksichtigt werden. Dazu werden die Restriktionen an  $u$  in der Definition des Raums  $C_0^1$  nur auf  $\Gamma_d$  vorgenommen, für  $\Gamma_n, \Gamma_c$  werden dafür Randintegrale mit  $uv, u$  in die Funktionale  $a$  bzw.  $f$  aufgenommen. Im folgenden wird aber weiter nur das Dirichletproblem behandelt. Es läßt sich zeigen, daß

$$\|v\|_{1,2} := \left( \int \int_{\Omega} [v_x^2 + v_y^2] dx dy \right)^{1/2} \quad (3.3.9)$$

auf  $C_0^1(\Omega)$  eine Norm darstellt. Für  $q \geq 0$  ist daher auch  $a(v, v)$  wieder definit. Dies war die entscheidende Eigenschaft in der Entwicklung (3.3.5), die zeigt, daß auch die Lösung  $u$  von (3.3.7) als Minimalstelle (3.3.4) von  $J(v) = \frac{1}{2}a(v, v) - f(v)$  charakterisiert werden kann mit den Bi-Linearformen aus (3.3.8). Die Tatsache, daß dabei nur einmalige Differenzierbarkeit von  $u$  erforderlich ist, ist nach den Bemerkung zu Beginn über die Herkunft der Dgl zweiter Ordnung aus einem System 1. Ordnung, nicht so überraschend.

In diesem Minimalproblem dürfen im Prinzip alle Funktionen  $v$  betrachtet werden, für die das Funktional  $J$  sinnvoll definierbar ist, also auch solche, bei der die ersten Ableitungen nur quadrat-integrierbar sind. Dazu gehören bei (3.3.8) die Elemente des Sobolev-Raums

$$W_0^1(\Omega) := \{u \in L_2(\Omega) : u_x, u_y \in L_2(\Omega), u|_{\Gamma} = 0\}. \quad (3.3.10)$$

Jede klassische Lösung des RWP's (vgl. Def. 3.2.1) ist auch globale Minimalstelle von  $J$  in  $W_0^1(\Omega)$ . Diese Erweiterung des Grundraums von  $C_0^1$  auf  $W_0^1$  hat große Bedeutung für die Konstruktion von Verfahren, da  $W_0^1$  auch Funktionen enthält, die (nur) stückweise stetig differenzierbar sind. Dies reduziert die Anforderungen an die Ansatzfunktionen bei dem folgenden Verfahren.

### 3.3.2 Rayleigh-Ritz-Galerkin-Verfahren

Mit der im letzten Abschnitt hergeleiteten Minimaleigenschaft der Lösung bekommt man in elementarer Weise Näherungsverfahren dadurch, daß das Minimum des Energie-Funktional  $J$  nicht mehr im Gesamttraum  $V (= W_0^1[0, 1]$  bzw.  $= W_0^1(\Omega))$  gesucht wird, sondern nur noch in einem *endlichdimensionalen* linearen Unterraum  $V_{\Delta} \subseteq V$ . Man kann bei  $V_{\Delta}$  zunächst an die Menge aller Polynome eines bestimmten Maximalgrads denken, flexibler sind aber die im nächsten Abschnitt besprochenen Spline-Funktionen, die nur noch stückweise stetig differenzierbar sind. Betrachtet werde nun also ein linearer Unterraum

$$V_{\Delta} \subseteq V, \quad \dim V_{\Delta} = n \in \mathbf{N},$$

mit einer Basis  $\{B_i\}_{i=1}^n$ . Die *Rayleigh-Ritz-Galerkin-Lösung* (RRG-Lösung)  $u_\Delta \in V_\Delta$  des Minimalproblems (3.3.4) ist definiert durch

$$J(u_\Delta) = \min_{v \in V_\Delta} J(v) = \min_{v \in V_\Delta} \left( \frac{1}{2} a(v, v) - f(v) \right). \quad (3.3.11)$$

Für die Analyse dieses Verfahrens gibt es ein allgemeines Prinzip. Als wesentliche Voraussetzung wird dabei die Vergleichbarkeit der von der Bilinearform  $a$  induzierten Norm und der Standard-(Sobolev-) Normen benötigt. Daher wird jetzt angenommen, daß Konstanten  $\mu_0, \mu_1, M_1 > 0$  existieren so, daß gilt

$$\mu_0 \|v\|_2^2 \leq \mu_1 \|v\|_{1,2}^2 \leq a(v, v) \leq M_1 \|v\|_{1,2}^2 \quad \forall v \in V. \quad (3.3.12)$$

Diese Voraussetzung ist in den betrachteten Beispielen (3.3.2) bzw. (3.3.8) erfüllt, wenn die Koeffizienten beschränkt sind und  $p(\cdot) \geq p^* > 0$ ,  $q \geq 0$  gilt.

**Satz 3.3.1** *Unter der Voraussetzung (3.3.12) existiert eine eindeutige RRG-Lösung  $u_\Delta$  von (3.3.11). Der Vektor  $\eta = (\eta_j)_{j=1}^n$  der Koeffizienten von  $u_\Delta = \sum_{j=1}^n \eta_j B_j$  löst das lineare System*

$$A\eta = \gamma, \quad \text{mit } A = \left( a(B_i, B_j) \right)_{i,j=1}^n, \quad \gamma = \left( f(B_i) \right)_{i=1}^n. \quad (3.3.13)$$

Die Matrix  $A = (a_{ij})$  ist symmetrisch, positiv definit. Äquivalent zu (3.3.13) ist die Aussage

$$a(u_\Delta, v) = f(v) \quad \forall v \in V_\Delta. \quad (3.3.14)$$

*Bemerkung:* Die Bedingung (3.3.14) heißt *Galerkin-Bedingung*. Sie definiert auch dann brauchbare Lösungen, wenn die Bilinearform  $a$  nicht mehr definit ist, eine Minimallösung (3.3.11) also nicht mehr existiert.

**Beweis** Es sei  $v = \sum_{j=1}^n \xi_j B_j \in V_\Delta$ ,  $\xi := (\xi_i)$ , beliebig. Dann gilt wegen der Bilinearität von  $a(\cdot, \cdot)$  und wegen (3.3.12) zunächst

$$\xi^\top A \xi = \sum_{i,j=1}^n \xi_i \xi_j a(B_i, B_j) = a \left( \sum_{i=1}^n \xi_i B_i, \sum_{j=1}^n \xi_j B_j \right) = a(v, v) \geq \mu_0 \|v\|_2^2 > 0,$$

für  $v \neq 0$ , d.h.,  $\xi \neq 0$ . Also ist  $A$  definit und somit invertierbar. Daraus folgt

$$\begin{aligned} J(v) &= \frac{1}{2} a \left( \sum_{i=1}^n \xi_i B_i, \sum_{j=1}^n \xi_j B_j \right) - f \left( \sum_{i=1}^n \xi_i B_i \right) = \frac{1}{2} \xi^\top A \xi - \gamma^\top \xi \\ &= \frac{1}{2} (\xi - A^{-1} \gamma)^\top A (\xi - A^{-1} \gamma) - \frac{1}{2} \gamma^\top A^{-1} \gamma. \end{aligned}$$

Dieser Ausdruck wird minimal genau dann, wenn  $\xi = A^{-1} \gamma = \eta$  ist, also  $v = u_\Delta$  gilt. Nach Definition entspricht (3.3.13) der Bedingung  $a(u_\Delta, B_i) = f(B_i)$ ,  $i = 1, \dots, n$ . Da  $\{B_i\}$  eine Basis von  $V_\Delta$  darstellt (n.V.) ergibt sich (3.3.14). ■

Die Lösung  $u$  erfüllt (3.3.3), also  $f(v) = a(u, v)$  auch mit  $v = u$ . Für beliebige  $v \in V$  folgt daher

$$\begin{aligned} J(v) - J(u) &= \frac{1}{2} a(v, v) - f(v) - \frac{1}{2} a(u, u) + f(u) = \frac{1}{2} a(v, v) - a(u, v) + \frac{1}{2} a(u, u) \\ &= \frac{1}{2} a(v - u, v - u). \end{aligned}$$

Da die Minimalstellen von  $J(v) - J(u)$  und  $J(v)$  übereinstimmen (auch in  $V_\Delta$ ), bekommt man daraus eine andere Charakterisierung der RRG-Lösung, die direkt auf Fehlerschranken führen. Danach ist  $u_\Delta$  die Bestapproximierende an  $u$  im Raum  $V_\Delta$ , gemessen in der 'a-Norm'.

**Satz 3.3.2** Für die RRG-Lösung  $u_\Delta$  aus (3.3.11) gilt

$$\mu_1 \|u_\Delta - u\|_{1,2}^2 \leq a(u_\Delta - u, u_\Delta - u) = \inf_{v \in V_\Delta} a(v - u, v - u) \leq M_1 \inf_{v \in V_\Delta} \|v - u\|_{1,2}^2.$$

Durch diesen Satz erhält man sofort eine globale Fehlerschranke in der 1, 2-Norm, also eine Schranke für den Fehler in der Ableitung von  $u_\Delta$ , wenn man rechts eine beliebige, geeignete Funktion  $v \in V_\Delta$  einsetzt, z.B., eine Interpolierende der Lösung  $u$ . Schranken für den Abstand  $\text{dist}_{1,2}(u, V_\Delta) = \inf_{v \in V_\Delta} \|v - u\|_{1,2}$  werden später für konkrete Räume  $V_\Delta$  hergeleitet.

### 3.3.3 Spline-Räume, Finite Elemente

Die Effizienz und Flexibilität des RRG-Verfahrens hängt wesentlich von der Wahl des endlichdimensionalen Raums  $V_\Delta$  ab. Man sieht schnell, daß ein Raum von Polynomen bei partiellen Problemen sehr unhandlich wird. Gerade bei der Behandlung unregelmäßiger Gebiete ist dagegen der folgende Zugang sehr anpassungsfähig, bei dem das Grundgebiet in kleine, einfache *Elemente* zerlegt wird. Hierbei ist es sehr hilfreich, daß das Minimalproblem (3.3.4) schon für stückweise differenzierbare Funktionen erklärt ist.

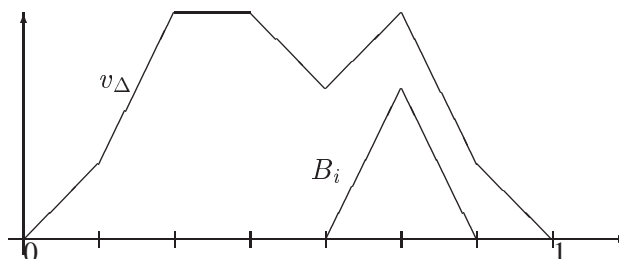
Zunächst wird das gewöhnliche Randwertproblem (3.3.1) auf  $[0, 1]$  betrachtet. Eine Unterteilung des Intervalls führt auf ein Gitter  $\Delta : 0 = x_0 < x_1 < \dots < x_m = 1$ , mit Schrittweiten  $h_i = x_{i+1} - x_i$ . Die einfachste Funktionenklasse in  $W_0^1[0, 1]$  sind die stückweise linearen Funktionen (lineare Splines, vgl. Numerik I, §4). Unter Berücksichtigung der Randbedingungen sei

$$V_\Delta := \{s \in C_0^0[0, 1] : s|_{[x_i, x_{i+1}]} \in \Pi_1, i = 0, \dots, m-1\} \quad (3.3.15)$$

der Raum aller stetigen Funktionen mit Randbedingung null, die in jedem Teilintervall Polynome vom Grad 1 ( $\in \Pi_1$ ) sind. Eine einfache Basis dieses Raums ist durch die *Dach-Funktionen*

$$B_i(x) := \begin{cases} (x - x_{i-1})/h_{i-1}, & x \in [x_{i-1}, x_i) \\ (x_{i+1} - x)/h_i, & x \in [x_i, x_{i+1}) \\ 0, & \text{sonst} \end{cases}, \quad (3.3.16)$$

$i = 1, \dots, m-1$ , gegeben, vgl. die Skizze.



Diese  $B_i$  bilden eine Kardinalbasis, da für  $x = x_k$  gilt  $u_\Delta(x_k) = \sum \eta_j B_j(x_k) = \eta_k$ . Der wichtigste praktische Vorteil der genannten Basis ist aber ihre Lokalität, die dazu führt, daß nur wenige der Matrixelemente  $a(B_i, B_j)$  in (3.3.13) von Null verschieden sind. Dies reduziert den Rechenaufwand sowohl bei der Aufstellung des Gleichungssystems, als auch bei seiner Auflösung.

Als Beispiel werde die Dgl  $-u'' = g$  mit  $p \equiv 1, q \equiv 0$ , also die Bilinearform  $a(u, v) = \int_0^1 u'(x)v'(x)dx$  betrachtet. Für die Basisfunktionen gilt nach (3.3.16)

$$B_i'(x) := \begin{cases} 1/h_{i-1}, & x \in [x_{i-1}, x_i) \\ -1/h_i, & x \in [x_i, x_{i+1}) \\ 0, & \text{sonst} \end{cases} .$$

Daher ist  $a(B_i, B_j) = 0$  für  $|i - j| > 1$ . Für  $j \in \{i - 1, i, i + 1\}$  gilt

$$a(B_i, B_j) = \int_{x_{i-1}}^{x_i} \frac{1}{h_{i-1}} B_j' dx - \int_{x_i}^{x_{i+1}} \frac{1}{h_i} B_j' dx = \begin{cases} -\frac{1}{h_{i-1}}, & j = i - 1 \\ \frac{1}{h_{i-1}} + \frac{1}{h_i}, & j = i \\ -\frac{1}{h_i}, & j = i + 1 \end{cases} .$$

Das Verfahren arbeitet daher, zumindestens bei der Approximation der zweiten Ableitung, ähnlich wie das Differenzenverfahren, denn es führt auf das Gleichungssystem ( $\eta_i = u_\Delta(x_i)$ , s.o.)

$$-\frac{1}{h_i}(\eta_{i+1} - \eta_i) + \frac{1}{h_{i-1}}(\eta_i - \eta_{i-1}) = \int_{x_{i-1}}^{x_{i+1}} g(x) B_i(x) dx, \quad i = 1, \dots, m - 1.$$

Für das RRG-Verfahren im Raum (3.3.15) sollte man daher auch einen Fehler der Größenordnung  $h^2$  erwarten. Dieser läßt sich mit Hilfe von Satz 3.3.2 auf den Approximationsfehler  $\inf\{\|u - v_\Delta\| : v_\Delta \in V_\Delta\}$  zurückführen. Bei der Interpolierenden in  $V_\Delta$  kann der Fehler unabhängig auf den einzelnen Teilintervallen untersucht werden. Die Interpolierende einer Funktion  $v \in V$  ist eine Restriktion  $R_\Delta v$  im Sinne von §2.

**Lemma 3.3.3** *Zu einer Funktion  $v \in C_0^2[0, 1]$  sei  $R_\Delta v := \sum_{i=1}^{m-1} v(x_i) B_i$  die Interpolierende aus  $V_\Delta$ . Dann gilt mit  $H := \max_i h_i$  die Fehlerschranke*

$$\|(R_\Delta v - v)^{(k)}\|_2 \leq \left(\frac{H}{\pi}\right)^{2-k} \|v''\|_2, \quad k = 0, 1.$$

Auf diese Aussage kommt man direkt über die Standarddarstellung des Interpolationsfehlers für Polynome (Numerik I, (3.3.11)), zur Verifikation der angegebenen Konstanten ist allerdings eine diffizilere Argumentation nötig. Das Infimum in Satz 3.3.2 kann nun durch Einsetzen von  $v = R_\Delta u$  abgeschätzt werden und ergibt mit dem letzten Lemma die Fehlerschranke

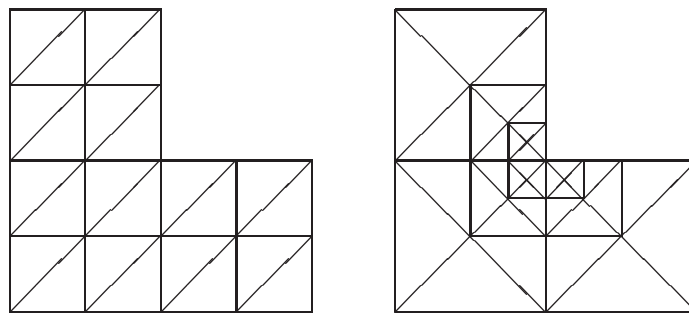
$$\|(u_\Delta - u)'\|_2 \leq \sqrt{\frac{M_1}{\mu_1}} \frac{H}{\pi} \|u''\|_2.$$

Mit einem zusätzlichen Argument ('Nitsche-Trick') kann aus dieser  $\mathcal{O}(H)$ -Schranke für die Ableitung bei vielen Randwertproblemen die globale Fehlerschranke

$$\|u_\Delta - u\|_2 \leq CH^2 \|u''\|_2 \leq \tilde{C} H^2 \|g\|_2$$

für die Funktionswerte hergeleitet werden. Im Vergleich zu den entsprechenden Aussagen bei Differenzenverfahren ist dabei hervorzuheben, daß hier nur die  $L_2$ -Integrierbarkeit von  $u''$  ausgenutzt wurde im Vergleich zur Voraussetzung  $u \in C^4$  bei den Differenzenverfahren. RRG-Verfahren höherer Ordnung erhält man durch Verwendung von Splines höherer Ordnung, wobei die Funktionen in  $V_\Delta$  stückweise aus Polynomen höheren Grads zusammengesetzt sind.

Von größter praktischer Bedeutung ist die Möglichkeit, beim RRG-Verfahren mit stückweise linearen Ansatzfunktionen arbeiten zu können, bei den partiellen Problemen. Sehr einfach sind Zerlegungen bzw. Approximationen des Gebiets  $\Omega$  durch *Dreieckgitter*. Daher ist das RRG-Verfahren wesentlich flexibler als Differenzenverfahren bei der Approximation komplizierterer Gebiete und bei der lokalen Gitterverfeinerung. Beispiele für regelmäßige bzw. lokal verfeinernde Triangulierungen beim L-Gebiet sind:



Dabei sind nur Triangulierungen aus offenen, paarweise disjunkten Dreiecken  $T_i$  zulässig,

$$\Delta = \{T_i\}_{i=1}^m, \quad \bar{\Omega} = \bigcup_{i=1}^m \bar{T}_i,$$

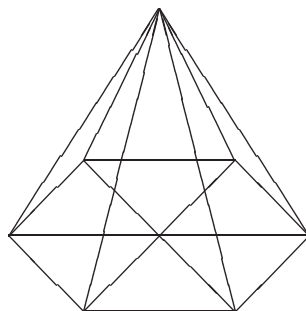
bei denen alle Ecken eines Dreiecks auf Ecken seiner Nachbardreiecke oder den Gebietsrand fallen. Der Raum aller stetigen, stückweise linearen Funktionen auf  $\Delta$ ,

$$V_\Delta := \{s \in C_0(\Omega) : s|_{T_i} \text{ linear}\}, \tag{3.3.17}$$

kann wieder durch eine Basis von Dachfunktionen aufgespannt werden. Es sei  $\{P_i\}_{i=1}^n \subseteq \Omega$  die Menge aller Eckpunkte von Dreiecken (Knoten) im Innern des Gebiets in einer geeigneten Numerierung. Dann wird durch

$$B_i \in V_\Delta, \quad B_i(P_j) = \delta_{ij}, \quad j = 1, \dots, n,$$

$i = 1, \dots, n$ , eine Basis von  $V_\Delta$  definiert. Ein Beispiel einer solchen Basisfunktion bei sechs in einer Ecke zusammenstoßenden Dreiecken ist



Wesentliche Merkmale dieser Basis sind wieder die Kardinalität,  $u_\Delta = \sum_{i=1}^n u_\Delta(P_i)B_i$ , und die Lokalität der Basis, da  $B_i$  nur auf den Dreiecken  $T_j$  mit  $P_i \in \bar{T}_j$  nicht verschwindet.

Diese Eigenschaft ist gerade bei mehrdimensionalen Problemen von großer Bedeutung, da dann nur wenige der Integrale  $a(B_i, B_j)$  in der Matrix  $A$  von Null verschieden sind, und bei den restlichen für  $i \neq j$  nur über zwei Dreiecke integriert werden muß.

$$a(B_i, B_j) = \sum_{P_i \in T_k \wedge P_j \in T_k} \int \int_{T_k} [(B_i)_x (B_j)_x + (B_i)_y (B_j)_y + q B_i B_j] dx dy$$

Bei der praktischen Berechnung der Integrale  $a_{ij} = a(B_i, B_j)$ ,  $\gamma_i = f(B_i)$ , ist es am einfachsten, alle Dreiecke  $T_i$  abzarbeiten und die einzelnen Teilintegrale in  $a_{ij}, \gamma_i$  aufzuaddieren. Die Matrix  $A$  ist hier wieder dünn besetzt und symmetrisch, pos. definit. Daher können sehr effiziente (Iterations-) Verfahren zur Lösung des Gleichungssystems (3.3.13) eingesetzt werden.

Zur Fehlerabschätzung der zugehörigen RRG-Lösung nach Satz 3.3.2 läßt sich eine Schranke für den Interpolationsfehler wieder lokal herleiten. Dabei ergibt sich eine verwertbare Aussage allerdings nur, wenn die Dreiecke  $T_k$  nicht zu spitz werden.

**Lemma 3.3.4** *Zu einer Funktion  $v \in C_0^2(\Omega)$  sei  $R_\Delta v := \sum_{i=1}^n v(P_i)B_i$  die Interpolierende aus  $V_\Delta$ , (3.3.17). In der Dreieckzerlegung des Gebiets sei  $h$  die längste Dreieckseite und  $\alpha$  der kleinste Innenwinkel eines Dreiecks  $T_i$ ,  $i = 1, \dots, m$ . Dann gilt die Fehlerschranke*

$$\|R_\Delta v - v\|_{1,2} \leq C \frac{h}{\sin \alpha} \|v\|_{2,2},$$

mit  $\|v\|_{2,2}^2 := \sum_{j+k=2} \|\frac{\partial^2}{\partial x^j \partial y^k} v\|_2^2$ .

Damit kann man wie im gewöhnlichen Fall in zwei Schritten einen Fehler  $\mathcal{O}(h/\sin \alpha)\|u\|_{2,2}$  in den ersten Ableitungen der RRG-Näherung  $u_\Delta$  und eine Schranke  $\mathcal{O}(h^2/\sin \alpha)\|u\|_{2,2}$  für den Fehler ihrer Funktionswerte herleiten. Zusätzlich kann unter einfachen Annahmen an die Koeffizienten zumindest für konvexe Gebiete die Abschätzung  $\|u\|_{2,2} \leq C\|g\|_2$  gezeigt werden. Daher hat man auch im partiellen Fall die  $h^2$ -Konvergenz unter wesentlich schwächeren Voraussetzungen als bei den Differenzenverfahren.

Allerdings ist bei diesen starken Aussagen zu beachten, daß sie nur für das exakte RRG-Verfahren zutreffen, bei dem alle Integrale  $a(B_i, B_j)$ ,  $f(B_i) = (g, B_i)_2$  ohne Fehler bestimmt werden. In der Praxis berechnet man jedoch zumindestens die Integralanteile, in denen allgemeine Koeffizienten (hier  $q, g$ ) auftreten, mit Hilfe von Quadraturformeln. Um auch für diese approximierten RRG-Verfahren einen  $\mathcal{O}(h^2)$ -Fehler garantieren zu können, müssen doch wieder stärkere Annahmen an die Koeffizienten, also auch an die Regularität der Lösung gemacht werden. Verfahren höherer Konvergenzordnung ergeben sich bei Verwendung von stückweise stetigen Polynomen höheren Grades auf den Dreieckzerlegungen.

### 3.4 Partielle Anfangs-Randwert-Probleme

Die nicht-elliptischen partiellen Differentialgleichungen 2. Ordnung, die parabolischen und hyperbolischen, besitzen eine gemeinsame Eigenschaft. Schreibt man den jeweiligen einfachsten Vertreter (Wärmeleitungs-, Wellengleichung) in der Form

$$u_{xx} - u_y = 0, \quad \text{bzw.} \quad u_{xx} - u_{yy} = 0,$$

so lassen sich damit sinnvolle Problemstellungen formulieren, wenn Randbedingungen nur in  $x$ -Richtung vorgegeben werden, in  $y$ -Richtung dagegen nur einseitige, also Anfangsbedingungen. Im physikalischen Zusammenhang hat die erste Variable,  $x$ , meist die Bedeutung einer Ortskoordinate, während die zweite,  $y$ , einen Zeitablauf beschreibt. Daher wird sie im folgenden mit  $t$  bezeichnet. Da die Geometrie des Problems sich i.a. nicht mit der Zeit ändert, ist das Grundgebiet meist ein Zylinder, z. B.,  $[0, 1] \times [0, \infty)$ . Mit stetigen Koeffizienten  $p, r, g$ ,  $p > 0, p_x$  stetig, werden also Anfangs-Randwertprobleme betrachtet mit einer der Dgln

$$u_t = [p(x, t)u_x]_x + q(x, t)u + g(x, t), \quad x \in (0, 1), \quad t > 0, \quad (\text{parabolisch}) \quad (3.4.1)$$

$$u_{tt} = [p(x, t)u_x]_x + q(x, t)u + g(x, t), \quad x \in (0, 1), \quad t > 0, \quad (\text{hyperbolisch}) \quad (3.4.2)$$

und den Rand- und Anfangsbedingungen

$$\begin{aligned} u(x, 0) &= \phi(x), \quad x \in (0, 1), \quad u(0, t) = u(1, t) = 0, \quad t \geq 0, \\ u_t(x, 0) &= \psi(x), \quad x \in (0, 1) \quad \text{nur bei (3.4.2)}. \end{aligned} \quad (3.4.3)$$

Die Probleme (3.4.1, 3.4.2), (3.4.3) verhalten sich also in Ortsrichtung wie (elliptische) Randwertprobleme, in Zeitrichtung wie Anfangswertprobleme erster bzw. zweiter Ordnung. Diese Tatsache kann bei der numerischen Approximation auch ausgenutzt werden durch Semidiskretisierung in der sogenannten *Linienmethode*. Dabei wird jeweils eines der von früher bekannten Verfahren zur Diskretisierung von Anfangs- oder Randwertproblemen getrennt auf die Zeit- bzw. Ortsrichtung angewendet. Allerdings muß man sich dabei über die Eigenschaften der nach dem ersten Diskretisierungsschritt entstehenden gewöhnlichen Probleme klar werden. Bei der Zeitlinienmethode, z.B., wird das Problem nur in der Ortsvariablen diskretisiert, das AWP in Zeitrichtung (zunächst) beibehalten. Dabei werde zunächst das Differenzenverfahren betrachtet mit einem äquidistanten Gitter  $\Delta = \{x_i = ih : i = 0, \dots, n+1\}$ ,  $h = 1/(n+1)$ , wie in §1.7. Der Ausdruck  $[pu_x]_x$  wird dabei durch zweifache Differenzenquotienten approximiert,

$$[p(x, t)u_x(x, t)]_x = \frac{p(x_{i+1/2}, t)(u(x_{i+1}, t) - u(x_i, t)) - p(x_{i-1/2}, t)(u(x_i, t) - u(x_{i-1}, t))}{h^2} + \mathcal{O}(h^2),$$

$x_{i\pm 1/2} := x_i \pm h/2$ . Unter Vernachlässigung des  $\mathcal{O}(h^2)$ -Fehlerterms ergibt sich für die Gitterwerte  $u_\Delta(x_i, t)$  einer Näherung  $u_\Delta$ , die immer noch Funktionen der Zeit sind, ein System von gewöhnlichen Anfangswertproblemen. Faßt man nämlich diese Funktionswerte als Komponenten einer Vektorfunktion auf,  $w(t) = (w_i(t))_{i=1}^n$ ,  $w_i(t) := u_\Delta(x_i, t)$  (und  $w_0 = w_{n+1} \equiv 0$ ), dann führt die obige Differenzenapproximation im parabolischen Fall (3.4.1) auf das AWP erster Ordnung

$$\begin{aligned} w'(t) &= f(t, w(t)), \quad t > 0 \\ w_i(0) &= \phi(x_i), \quad i = 1, \dots, n. \end{aligned} \quad (3.4.4)$$

mit

$$f_i(t, y) = \frac{1}{h^2} [p_{i+1/2}(y_{i+1} - y_i) - p_{i-1/2}(y_i - y_{i-1})] + q_i y_i + g_i, \quad i = 1, \dots, n.$$

Die Koeffizientenwerte sind, wie früher,  $q_i = q(x_i, t)$ ,  $g_i = g(x_i, t)$  und  $p_{i\pm 1/2} = p((i \pm \frac{1}{2})h, t)$ . Beim hyperbolischen Problem (3.4.2) ergibt sich das entsprechende AWP mit der Dgl zweiter Ordnung  $w'' = f(t, w(t))$  und der zusätzlichen Anfangsbedingung  $w'_i(0) = \psi(x_i)$ . Im betrachteten Beispiel ist die Funktion  $f$  affin linear,  $f(t, y) = A(t)y + \vec{g}(t)$ . Die zeitabhängige Matrix  $A(t)$  ist tridiagonal (vgl. (1.7.5)) und sogar symmetrisch mit den Elementen

$$a_{ii}(t) = -\frac{1}{h^2} [p_{i+1/2} + p_{i-1/2}] + q_i, \quad a_{i,i+1}(t) = a_{i+1,i}(t) = \frac{1}{h^2} p_{i+1/2}. \quad (3.4.5)$$

Wie in Lemma 3.2.4 kann für  $p > 0$  bei dieser Matrix die Aussage  $y^T A(t)y \leq q^* y^T y$ ,  $q^* = \max_i q_i$ , gezeigt werden. Daher sind die Eigenwerte von  $A$  reell und durch  $q^*$  nach oben beschränkt, also insbesondere negativ für  $q^* \leq 0$ . Allerdings wird die Norm  $\|A(t)\| \cong \frac{4}{h^2} \|p\|_\infty$  für kleine Schrittweiten  $h$  sehr groß. Tatsächlich läßt sich zeigen, daß sehr kleine Eigenwerte von  $A$  (mit großem Betrag)  $\lambda \ll 0$  auftreten können. Bei (3.4.4) handelt es sich somit um (eines der Standardbeispiele für) ein *steifes Anfangswertproblem*. Die Diskussion aus §1.5 ist also hier anwendbar. Die wichtigste Erkenntnis in diesem Zusammenhang war, daß zur Integration des AWP's (3.4.4) nur implizite Verfahren sinnvoll sind. Das läßt sich auch anhand der einfachsten Verfahren bei der Wärmeleitungsgleichung nachvollziehen. Bei dieser ( $p \equiv 1$ ,  $q = g \equiv 0$ ) wird das einfache Einschrittverfahren ( $w(t_k) \cong w^{[k]} \in \mathbf{R}^n$ )

$$\frac{1}{\tau} (w^{[j+1]} - w^{[j]}) = (1 - \vartheta)Aw^{[j]} + \vartheta Aw^{[j+1]}, \quad j = 0, 1, \dots, \quad (3.4.6)$$

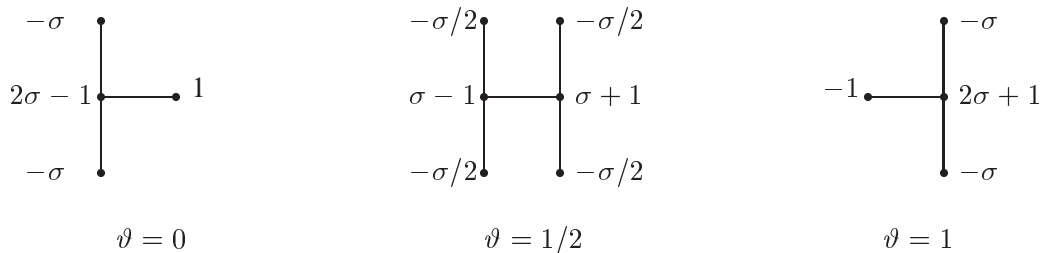
mit Zeitschrittweite  $\tau$  und einem Parameter  $\vartheta \in [0, 1]$  betrachtet. Die Matrix  $A$  (3.4.5) ist konstant. Das Verfahren stimmt für  $\vartheta = 0$  mit dem expliziten und für  $\vartheta = 1$  mit dem impliziten Eulerverfahren, überein, für  $\vartheta = \frac{1}{2}$  ist es die implizite Mittelpunkregel. Es führt auf eine Serie von Gleichungssystemen ( $\sigma := \tau/h^2$ ,  $\vartheta' := 1 - \vartheta$ ),

$$-\vartheta\sigma w_{i-1}^{[j+1]} + (1 + 2\vartheta\sigma)w_i^{[j+1]} - \vartheta\sigma w_{i+1}^{[j+1]} = \vartheta'\sigma w_{i-1}^{[j]} + (1 - 2\vartheta'\sigma)w_i^{[j]} + \vartheta'\sigma w_{i+1}^{[j]}, \quad (3.4.7)$$

$i = 1, \dots, n$ . Beim expliziten Euler-Verfahren,  $\vartheta = 0$ , reduziert sich dies auf

$$w_i^{[j+1]} = \sigma w_{i-1}^{[j]} + (1 - 2\sigma)w_i^{[j]} + \sigma w_{i+1}^{[j]}, \quad i = 1, \dots, n. \quad (3.4.8)$$

Wie im letzten Paragraphen können diese Beziehungen kompakt durch Differenzensterne beschrieben werden. Zu den verschiedenen Parameterwerten  $\vartheta$  gehören die Sterne



wobei die Zeitachse  $t$  nach rechts gezeichnet ist. Für  $\sigma \leq \frac{1}{2}$  sind alle Gewichte in (3.4.8) nichtnegativ,

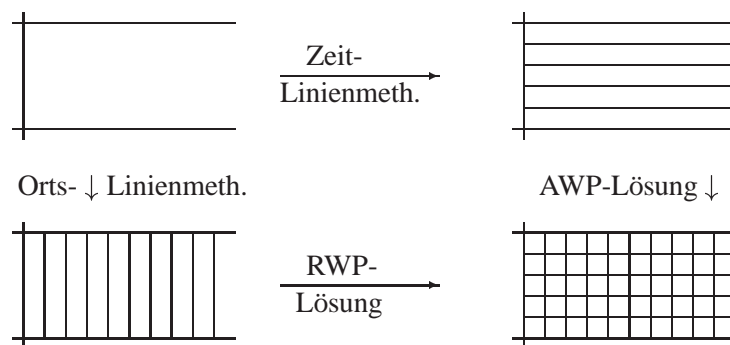
es findet also eine fortgesetzte Mittelwertbildung in den Zeitschritten statt. Dies ist ein gutes Modell für die durch die Wärmeleitungsgleichung beschriebenen Ausgleichsvorgänge. Allerdings muß dabei die Zeitschrittweite  $\tau$  aus Stabilitätsgründen inakzeptabel klein gewählt werden,  $\tau = h^2\sigma \leq h^2/2$ . Für  $\tau > h^2/2$ , d.h.  $\sigma > 1/2$  wird das explizite Verfahren (3.4.8) nämlich instabil, denn dann tritt in der (scharfen) Abschätzung

$$\|w^{[j+1]}\|_\infty = \max_i |w_i^{[j+1]}| \leq (2\sigma + |1 - 2\sigma|) \|w^{[j]}\|_\infty,$$

der Vorfaktor  $4\sigma - 1 = 4\tau/h^2 - 1 > 1$  auf. Dies paßt zur Diskussion von §1.5. Denn, da für die Norm  $\|A\| = 4/h^2$  gilt und die Stabilitätsabszisse des expliziten Eulerverfahren 2 ist, entspricht dieser Fall  $\sigma = \tau/h^2 > 1/2$  der instabilen Situation  $\tau\|A\| > 2$ .

Beim impliziten Eulerverfahren ( $\vartheta = 1$ ) dagegen ist das Tridiagonalsystem (3.4.7) in jedem Zeitschritt diagonaldominant (mit einer M-Matrix) und hier kann ohne jede Schrittweitenbeschränkung die Schranke  $\|w^{[j+1]}\|_\infty \leq \|w^{[j]}\|_\infty$  gezeigt werden. Das allgemeine Verfahren (3.4.6) ist für alle  $\vartheta \in [\frac{1}{2}, 1]$  A-stabil. Bei diesen Verfahren muß sich die Schrittweitenwahl nur nach den Genauigkeitsanforderungen richten. Dabei ist zu beachten, daß im AWP (3.4.4) auch schon ein Diskretisierungsfehler steckt, der Fehler beim Gesamtverfahren (3.4.6) setzt sich daher aus einem Orts- und einem Zeitfehler zusammen, die jeweils getrennt durch Wahl der Ortsschrittweite  $h$  und der Zeitschrittweite  $\tau$  beeinflusst werden können. Für die Mittelpunkregel ( $\vartheta = 1/2$ ) hat dieser globale Fehler die Form  $C_x h^2 + C_t \tau^2$ , für alle anderen  $\vartheta \neq 0$  nur  $C_x h^2 + C_t \tau$ . Das Verfahren arbeitet dann am effizientestens, wenn beide Fehler ungefähr gleich groß sind. Bei der Zeitlinienmethode, mit festem Ortsgitter, kann nur der Zeitfehler beeinflusst werden. Die Zeitschrittweite  $\tau$  ist daher so zu steuern, daß dieser ungefähr die Größe des Ortsfehlers besitzt.

Eine größere Anpassungsfähigkeit ist mit der Orts-Linienmethode (Rothe-Methode) möglich. Dabei werden die beiden Diskretisierungsschritte vertauscht. Die erste Diskretisierung erfolgt in Zeitrichtung mit einem impliziten Verfahren (z.B. dem impliziten Euler-Verfahren). Dies führt auf eine Serie von gewöhnlichen RWPen in jeder Zeitstufe. Das folgende Diagramm veranschaulicht die auftretenden Gitter für beide Möglichkeiten:



Wird bei den RWPen das Differenzenverfahren mit einer festen, konstanten Schrittweite  $h$  verwendet, ergibt sich wieder das Verfahren (3.4.6) und das im Bild gezeigte Endgitter. In dieser zweiten Interpretation ist es aber auch möglich, bei der Lösung der RWPe in jedem einzelnen Zeitschritt das Ortsgitter

neu (adaptiv) zu wählen. Die im Verfahren verwendeten Werte der vorhergehenden Zeitstufe können dabei interpoliert werden. Hiermit läßt sich sowohl die Zeit- als auch die Ortsunterteilung leichter an den Lösungsverlauf anpassen, die Gitter sind dann irregulär.

Auch bei den **hyperbolischen** Problemen (3.4.2) kann, im Prinzip, analog zum parabolischen Fall vorgegangen werden. Zur Diskretisierung der zweiten Zeitableitung  $u_{tt}$  bzw.  $w_{tt}$  (in der Zeitlinienmethode) müssen dabei allerdings drei Zeitstufen herangezogen werden. Ein explizites Verfahren ergibt sich direkt aus der Standardapproximation von  $u_{tt}$  mit Hilfe des zweiten Differenzenquotienten, für die Gitterwerte (wie oben)  $w_i^{[j]} = u_{\Delta}(ih, j\tau)$ ,  $i = 1, \dots, n$ , lautet das Gesamtverfahren

$$w_i^{[j+1]} = s^2 w_{i-1}^{[j]} + 2(1 - s^2) w_i^{[j]} + s^2 w_{i+1}^{[j]} - w_i^{[j-1]}, \quad i = 1, \dots, n, \quad (3.4.9)$$

mit  $s = \tau/h$ . Diese Rekursion kann im Zeitschritt  $j = 0$  gestartet werden, wenn die zweite Anfangsbedingung aus (3.4.3) mit der symmetrischen ersten Differenz approximiert wird,  $w_i^{[1]} - w_i^{[-1]} = 2\tau\psi(x_i)$ , die eine künstliche Zeitebene  $-\tau$  einbezieht. Dies entspricht dem Vorgehen bei Neumann-Randbedingungen in §3.2.1. Durch eine etwas aufwendigere Analyse als bei den parabolischen Problemen läßt sich zeigen, daß das Verfahren ebenfalls nur stabil ist unter der (schwächeren) Schrittweitereinschränkung  $s = \tau/h \leq 1$ . Es können auch einfache implizite Verfahren angegeben werden, die ohne Schrittweitenbeschränkung stabil sind. Der Vergleich mit den expliziten Verfahren fällt dennoch anders aus als bei den parabolischen Problemen, da sich die Lösungen in den beiden Fällen mit der Zeit unterschiedlich entwickeln. Dies wird jetzt kurz diskutiert.

Bei der parabolischen Gleichung (3.4.1) werden die Lösungen für wachsendes  $t > 0$  immer *glatter*. Daher kann nach einer kurzen 'transienten' Phase die Lösung in der Regel mit größeren Zeitschrittweiten approximiert werden, allerdings nur bei impliziten Ein- oder Mehrschritt-Verfahren, die die erforderlichen Stabilitätseigenschaften für steife Probleme besitzen. Die Glättungseigenschaft sieht man, z.B., bei der Wärmeleitungsgleichung  $u_t = u_{xx}$  mit Randbedingung  $u(0, t) = u(1, t) = 0$  daran, daß die Lösung aus einer Linearkombination von Funktionen der Form

$$e^{-k^2\pi^2 t} \sin k\pi x, \quad k \in \mathbf{N},$$

besteht. Diese verschwinden in  $t$ -Richtung umso schneller, je stärker sie in  $x$ -Richtung variieren, nach kurzer Zeit sind nur langsam variierende Anteile übrig. Außerdem hängt der Lösungswert  $u(x, t)$  von allen Daten im Bereich  $[0, 1] \times [0, t)$  ab.

Beim hyperbolischen Problem ergibt sich dagegen ein anderes Bild. Die allgemeine Lösung der Wellengleichung kann explizit konstruiert werden. Jede Funktion der Form

$$u(x, t) = a(x + t) + b(x - t) \quad \text{löst} \quad u_{tt} = u_{xx} \quad (3.4.10)$$

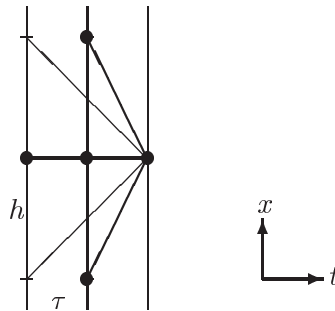
mit  $a, b \in C^2(\mathbf{R})$ . Wenn man die Randbedingungen in  $x$ -Richtung zur Vereinfachung außer Acht läßt, hat diese Lösung die Anfangsbedingungen

$$a(x) + b(x) = u(x, 0) = \phi(x), \quad a'(x) - b'(x) = u_t(x, 0) = \psi(x)$$

zu erfüllen. Durch Differentiation der ersten Gleichung können  $a$  und  $b$  bis auf Konstanten bestimmt werden. Daraus folgt die explizite Darstellung des AWP's bei der Wellengleichung,

$$u(x, t) = \frac{1}{2} \left( \phi(x+t) + \phi(x-t) + \int_{x-t}^{x+t} \psi(\xi) d\xi \right). \quad (3.4.11)$$

Zwei wichtige Eigenschaften können an dieser Lösungsdarstellung abgelesen werden. Erstens zeigt schon (3.4.10), daß die Eigenschaften der Lösung vom Startzeitpunkt praktisch unverändert fortgepflanzt werden, zweitens, daß die Lösung an einem Punkt  $(x, t)$  nur von den Startwerten im Intervall  $[x-t, x+t]$ , dem sogenannten Abhängigkeitsintervall, abhängt. Daher pflanzen sich Störungen (Sprünge, Schockwellen) längs der *Wellenfronten*  $x \pm t = \text{const}$  ungedämpft fort. Diese Linien nennt man *Charakteristiken*. Die aufgezählten Eigenschaften sollten bei den numerischen Verfahren berücksichtigt werden. Dies bedeutet, daß implizite Verfahren (mit großen Schrittweiten) nur angebracht sind, wenn die Lösung von Anfang an sehr glatt ist und daß andernfalls bei expliziten Verfahren die diskreten Abhängigkeitsgebiete an die der Dgl angepaßt werden sollten. In folgender Skizze



ist zu erkennen, daß für  $\tau < h$  die 'diskreten' Charakteristiken (dicke Linien) des Verfahrens (3.4.9) eine größere Steigung besitzen als die exakten (Steigung  $\pm 1$ ), Störungen sich in der Näherung in  $x$ -Richtung also schneller fortpflanzen als in der Realität. Im Verfahren (3.4.9) ist daher das Schrittweitenverhältnis  $s = 1$  am günstigsten. Bei allgemeinen Dgln (3.4.2) hängen die Charakteristiken von den Koeffizienten ab, sind im allgemeinen also keine geraden Linien mehr. Hier wird die Diskussion sehr viel schwieriger.

## Index

- Abhängigkeitsintervall, 81
- Adams-
  - Bashforth-Verfahren, 21
  - Moulton-Verfahren, 22
- Anfangswert, 36
- Anfangswertproblem, 5, 9, 47
  - steif, 30, 78
- Anlaufrechnung, 23
  
- BDF-Verfahren, 34
- Bilinearform, 70, 71
  
- Charakteristiken, 81
  
- Dahlquist, 26, 35
- Diagonaldominanz, 62
- Differenzen
  - Gleichung, 9, 28
  - Quotienten, 34, 42, 55, 77
- Dirichlet-Problem, 54, 56, 63, 70
- Dreieck-Gitter, 75
  
- Einschritt-Verfahren, 9, 13, 78
- Entwicklung
  - asymptotische, 19, 27
- Euler-Verfahren, 9, 78
  
- Fehler
  - Schätzung, 16, 19, 50
  - globaler, 15, 26, 44, 49, 63, 68, 73, 74, 79
  - lokaler, 12, 22, 23, 29, 43, 48, 63, 65, 67
- Fundamentalsystem, 7, 36, 37
- Funktional
  - konvexes, 70
  - lineares, 70
  
- Galerkin-Bedingung, 72
- Gauß-
  - Algorithmus, 40, 44
  - Quadraturformel, 33
  
- Gleichungssystem
  - nichtlinear, 36
- Gragg-Bulirsch-Stoer-Verfahren, 29
- Greensche Formel, 70
- Gronwall-Lemma, 7
  - diskret, 13
  
- Integralgleichung, 6
  
- Kollokationsverfahren, 49
- Kondensation, 40
- Konsistenz, 12, 22, 25, 48, 62
  - Ordnung, 12
  
- Linienmethode, 77
- Lipschitzbedingung, 6, 8, 14, 26
  
- M-Matrix, 61, 67, 69, 79
- Mehrschrittverfahren, 21
  - implizit, 34
  - lineare, 24
- Mehrstellenverfahren, 45, 65
- Mehrzielmethode, 39
- Mittelpunktregel
  - explizite, 27
  - implizite, 78
  
- Neumann-Randbedingung, 57
- Neumann-Reihe, 32, 44, 61
- Neville-Algorithmus, 29
- Newton-Verfahren, 37, 40
- Numerierung
  - lexikographisch, 56
  - Schachbrett-, 57
  
- Ordnungs-
  - Barriere, 26, 35
  - Bedingungen, 12, 25
  - Steuerung, 24, 29
  
- Partielle Differentialgleichung

- elliptische, 52–54, 70
- hyperbolische, 52, 53, 77, 80
- parabolische, 53, 77
- Picard-Lindelöf, Satz von, 6
- Poissongleichung, 54, 70
- Polynom
  - charakteristisches, 34
  - Lagrange-, 23
  - Legendre-, 33
  - Newton-, 23
- Prädiktor-Korrektor-Verfahren, 23
- Prolongation, 47
  
- Randwertproblem, 5, 36, 42, 47, 48, 54
- Restriktion, 47, 62, 74
- Richardson-Extrapolation, 19, 29
- Rosenbrock-Wanner-Verfahren, 34
- Rundungsfehler, 20
- Runge-Kutta-Verfahren, 50
  - allgemeine, 11
  - eingebettete, 17
  - explizite, 11
  - implizite, 33
  - klassisches, 10
  - linear-implizite, 33
  - stetige, 18
  
- Schrittweiten-Steuerung, 15, 17, 24, 29, 34
- Spline-Funktionen, 49, 75
- Stabilität, 13, 14, 26, 48, 62, 65, 79, 80
  - A-, 33
  - absolute, 32, 35
- Stabilitäts-
  - Bereich, 32
  - Funktion, 32, 35
- Stern
  - Fünf-Punkte-, 55, 64
  - Neun-Punkte-, 65
  
- Treppenmatrix, 40
- Triangulierung, 75
  
- Verfahren
  - explizite, 21, 78
  - implizite, 22, 31, 79
  
- Wärmeleitungsgleichung, 53, 78, 80
- Wellengleichung, 52, 80