



Numerik-Vorlesungen
Teil 2
Anfangs- und Randwertprobleme,
lineare Gleichungssysteme

Peter Szyler, Horst Hollatz

Letzte Änderung: 20. Januar 2006

Inhaltsverzeichnis

6. Anfangswertprobleme	1
6.1. Einführung	1
6.2. Einschrittverfahren	6
6.2.1. Grundlegende Begriffe und Euler-Verfahren	6
6.2.2. Runge-Kutta-Verfahren	11
6.2.3. Konvergenz von Einschrittverfahren	15
6.2.4. Rundungsfehlereinfluss	19
6.2.5. Schrittweitensteuerung	22
6.2.6. Steife Differentialgleichungen und implizite Verfahren	29
6.3. Mehrschrittverfahren	36
6.3.1. Prediktor-Korrektor-Verfahren	36
6.3.2. Konvergenz von Mehrschrittverfahren	41
6.4. Extrapolationsverfahren	52
6.5. Aufgaben	55
7. Randwertprobleme	59
7.1. Einführung	59
7.2. Das einfache Schießverfahren	62
7.3. Die Mehrzielmethode	71
7.4. Differenzenverfahren	72
7.5. Aufgaben	75
8. Lineare Gleichungssysteme	77
8.1. Allgemeine Grundlagen und Störungstheorie	77
8.1.1. Vektor- und Matrixnormen	77
8.1.2. Ordnungen und Beträge	83
8.1.3. Spezielle Transformationsmatrizen	85
8.1.4. Eigenwerte und Singulärwerte	89
8.1.5. Störungstheorie	93
8.2. Direkte Lösungsverfahren	106
8.2.1. Die LU-Zerlegung	106

8.2.2.	Rundungsfehleranalyse der LU-Zerlegung	115
8.2.3.	Pivotisierung und Skalierung	129
8.2.4.	Symmetrische Matrizen	140
8.2.5.	Orthogonalisierungsverfahren	147
8.3.	Iterative Verfahren	158
8.3.1.	Iterationsverfahren und ihre Konvergenz	158
8.3.2.	Das Jacobi- und das Gauß-Seidel-Verfahren	161
8.3.3.	Nachiteration	169
8.3.4.	Das cg-Verfahren von Hestenes und Stiefel	172
8.4.	Aufgaben	188
Index		196

Kapitel 6

Anfangswertprobleme

6.1. Einführung

Wir betrachten das folgende Problem:

Gegeben seien eine Teilmenge $D \subseteq \mathbb{R}$ und eine Funktion

$$f : D \times \mathbb{R}^d \longrightarrow \mathbb{R}^d.$$

Gesucht ist eine differenzierbare Funktion

$$y : D \longrightarrow \mathbb{R}^d,$$

die die Gleichung

$$y'(x) = f(x, y(x))$$

erfüllt.

Diese vektorielle Differentialgleichung erster Ordnung hat im allgemeinen unendlich viele Lösungen. Oft sucht man nach Lösungen, die gewisse Zusatzbedingungen erfüllen. Sind diese Bedingungen in der Form

$$y(x_0) = y_0$$

gegeben, so spricht man von einem **Anfangswertproblem**¹. Solche Probleme werden in diesem Kapitels untersucht. Allgemeiner als AWP sind Randwertprobleme², die wir im nächsten Kapitel behandeln werden. Bei ihnen sind zusätzliche Bedingungen der Form

$$r(y(a), y(b)) = 0$$

¹Wir werden statt Anfangswertproblem im weiteren die Abkürzung AWP verwenden.

²Für diese werden wir die Abkürzung RWP verwenden.

gegeben.

Bemerkung: Die Beschränkung auf Differentialgleichungen erster Ordnung ist keine Einschränkung. Jede Differentialgleichung m -ter Ordnung lässt sich formal in eine Differentialgleichung erster Ordnung überführen:

Wir definieren

$$z_1(x) = \mathbf{y}(x), \quad z_2(x) = \mathbf{y}'(x), \quad z_3(x) = \mathbf{y}''(x), \dots, z_m(x) = \mathbf{y}^{(m-1)}(x).$$

Dann gilt

$$z'_1 = z_2, \quad z'_2 = z_3, \quad z'_3 = z_4, \dots, z'_m = \mathbf{f}(x, z_1, z_2, \dots, z_m).$$

Mit

$$z(x) = \begin{pmatrix} z_1(x) \\ z_2(x) \\ \vdots \\ z_m(x) \end{pmatrix} : D \subseteq \mathbb{R} \longrightarrow \mathbb{R}^{m \cdot d}$$

und

$$\mathbf{F}(x, z) = \begin{pmatrix} z_2 \\ z_3 \\ \vdots \\ z_m \\ \mathbf{f}(x, z_1, \dots, z_m) \end{pmatrix} : D \times \mathbb{R}^{m \cdot d} \longrightarrow \mathbb{R}^{m \cdot d}$$

erhalten wir die Differentialgleichung erster Ordnung

$$z' = \mathbf{F}(x, z).$$

Die numerische Aufgabenstellung unterscheidet sich von der analytischen. Wir suchen nicht nach einem formelmäßigen Ausdruck für die Funktion \mathbf{y} , die das AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

löst, sondern wir werden versuchen, Näherungswerte

$$\boldsymbol{\eta}_i = \boldsymbol{\eta}(x_i)$$

an gewissen Stellen x_i zu bestimmen, die die exakten Werte

$$\mathbf{y}_i = \mathbf{y}(x_i)$$

möglichst gut approximieren.

Bevor wir zum Konstruieren und Untersuchen solcher Verfahren zum Berechnen der Näherungen $\boldsymbol{\eta}_i$ kommen, ist zu untersuchen, unter welchen Bedingungen ein AWP lösbar ist. Aus der Analysis ist dazu der folgende Satz bekannt.

6.1. Existenz- und Eindeigkeitssatz: Die Funktion f sei auf dem Streifen

$$S = [a, b] \times \mathbb{R}^d \quad (a \text{ und } b \text{ endlich})$$

definiert und stetig. Weiterhin existiere eine Konstante L , so dass für alle $x \in [a, b]$ und alle $\mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^n$

$$\|\mathbf{f}(x, \mathbf{y}_1) - \mathbf{f}(x, \mathbf{y}_2)\| \leq L \|\mathbf{y}_1 - \mathbf{y}_2\|$$

gelte. Dann existiert zu jedem $x_0 \in [a, b]$ und jedem $\mathbf{y}_0 \in \mathbb{R}^n$ genau eine Funktion \mathbf{y} mit

1. $\mathbf{y} \in F^1[a, b]$,
2. $\mathbf{y}' = \mathbf{f}(x, \mathbf{y})$ für alle $x \in [a, b]$ und
3. $\mathbf{y}(x_0) = \mathbf{y}_0$.

Für eine stetige Funktion f , die bezüglich \mathbf{y} zusätzlich lipschitzstetig ist, existiert also genau eine Lösung des AWP.

Weiterhin führen wir folgende Bezeichnung ein: Es sei $F^n[a, b]$ die Menge aller Funktionen

$$\mathbf{f} : [a, b] \times \mathbb{R}^d \longrightarrow \mathbb{R}^d,$$

für die alle partiellen Ableitungen bis zur Ordnung n auf dem Streifen $S = [a, b] \times \mathbb{R}^d$ existieren, dort stetig und beschränkt sind. Insbesondere erfüllen alle Funktionen $\mathbf{f} \in F^1[a, b]$ die Voraussetzungen von Satz 6.1. Aus numerischer Sicht ist die Abhängigkeit der Lösung von den Eingabedaten interessant.

6.2. Beispiel: Wir betrachten das AWP

$$y' = 10 \left(y - \frac{x^2}{1+x^2} \right) + \frac{2x}{(1+x^2)^2}, \quad y(0) = y_0 = 0.$$

Die exakte Lösung lautet

$$y(x) = \frac{x^2}{1+x^2}.$$

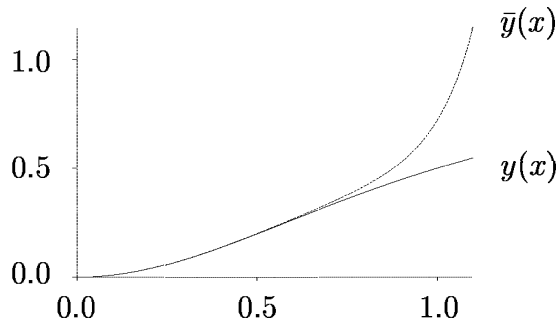
Ändern wir den Anfangswert leicht ab:

$$y(0) = \bar{y}_0 = \varepsilon,$$

so erhalten wir die Lösung

$$\bar{y}(x) = \varepsilon e^{10x} + \frac{x^2}{1+x^2}.$$

Für hinreichend großes x wird diese Lösung i. a. beliebig stark von der ersten Lösung abweichen. Das folgende Bild zeigt $y(x)$ und $\bar{y}(x)$ für $\varepsilon = 0.00001$.



Diese starke Abhängigkeit der Lösung von den Anfangswerten ist kein Einzelfall.

6.3. Satz: Die Funktion f erfülle die Bedingungen aus Satz 6.1. Dann gilt für die Lösung $\mathbf{y}(x; \mathbf{s})$ des AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0; \mathbf{s}) = \mathbf{s}$$

für ein beliebiges $x \in [a, b]$ die Abschätzung

$$\|\mathbf{y}(x; \mathbf{s}_1) - \mathbf{y}(x; \mathbf{s}_2)\| \leq e^{L|x-x_0|} \|\mathbf{s}_1 - \mathbf{s}_2\|.$$

Beweis: Es gilt

$$\mathbf{y}(x; \mathbf{s}) = \mathbf{s} + \int_{x_0}^x \mathbf{f}(t, \mathbf{y}(t; \mathbf{s})) dt.$$

Daraus folgt

$$\mathbf{y}(x; \mathbf{s}_1) - \mathbf{y}(x; \mathbf{s}_2) = \mathbf{s}_1 - \mathbf{s}_2 + \int_{x_0}^x [\mathbf{f}(t, \mathbf{y}(t; \mathbf{s}_1)) - \mathbf{f}(t, \mathbf{y}(t; \mathbf{s}_2))] dt,$$

und weiter

$$\begin{aligned} \|\mathbf{y}(x; \mathbf{s}_1) - \mathbf{y}(x; \mathbf{s}_2)\| &\leq \|\mathbf{s}_1 - \mathbf{s}_2\| + \left| \int_{x_0}^x \|\mathbf{f}(t, \mathbf{y}(t; \mathbf{s}_1)) - \mathbf{f}(t, \mathbf{y}(t; \mathbf{s}_2))\| dt \right| \\ &\leq \|\mathbf{s}_1 - \mathbf{s}_2\| + L \left| \int_{x_0}^x \|\mathbf{y}(t; \mathbf{s}_1) - \mathbf{y}(t; \mathbf{s}_2)\| dt \right|. \end{aligned}$$

Es sei

$$\Phi(x) = \int_{x_0}^x \|\mathbf{y}(t; \mathbf{s}_1) - \mathbf{y}(t; \mathbf{s}_2)\| dt$$

und damit

$$\Phi'(x) = \|\mathbf{y}(x; \mathbf{s}_1) - \mathbf{y}(x; \mathbf{s}_2)\|.$$

Wir betrachten den Fall $x \geq x_0$. Für die Funktion

$$\alpha(x) = \Phi'(x) - L\Phi(x)$$

gilt dann die Abschätzung

$$\alpha(x) \leq \|\mathbf{s}_1 - \mathbf{s}_2\|$$

und Φ ist Lösung des AWP

$$\Phi'(x) = \alpha(x) + L\Phi(x), \quad \Phi(x_0) = 0.$$

Die Lösung dieses AWP lässt sich in der Form

$$\Phi(x) = e^{L(x-x_0)} \int_{x_0}^x \alpha(t) e^{-L(t-x_0)} dt$$

schreiben. Wegen $\alpha(x) \leq \|\mathbf{s}_1 - \mathbf{s}_2\|$ folgt für $x \geq x_0$ die Abschätzung

$$\begin{aligned} 0 \leq \Phi(x) &\leq e^{L(x-x_0)} \|\mathbf{s}_1 - \mathbf{s}_2\| \int_{x_0}^x e^{-L(t-x_0)} dt \\ &= \frac{1}{L} e^{L(x-x_0)} \|\mathbf{s}_1 - \mathbf{s}_2\| \left(-e^{-L(x-x_0)} + 1 \right) \\ &= \frac{1}{L} \|\mathbf{s}_1 - \mathbf{s}_2\| \left(e^{L(x-x_0)} - 1 \right). \end{aligned}$$

Damit gilt

$$\begin{aligned} \|\mathbf{y}(x; \mathbf{s}_1) - \mathbf{y}(x; \mathbf{s}_2)\| &= \Phi'(x) = \alpha(x) + L\Phi(x) \\ &\leq \|\mathbf{s}_1 - \mathbf{s}_2\| + \|\mathbf{s}_1 - \mathbf{s}_2\| \left(e^{L(x-x_0)} - 1 \right) \\ &= \|\mathbf{s}_1 - \mathbf{s}_2\| e^{L(x-x_0)}. \end{aligned}$$

Im Falle $x < x_0$ folgt aus

$$\|\mathbf{y}(x; \mathbf{s}_1) - \mathbf{y}(x; \mathbf{s}_2)\| \leq \|\mathbf{s}_1 - \mathbf{s}_2\| - L \int_{x_0}^x \|\mathbf{y}(t; \mathbf{s}_1) - \mathbf{y}(t; \mathbf{s}_2)\| dt$$

und

$$\alpha(x) = \Phi'(x) + L\Phi(x)$$

wieder

$$\alpha(x) \leq \|\mathbf{s}_1 - \mathbf{s}_2\|$$

und

$$-\Phi(x) \leq \frac{1}{L} \|\mathbf{s}_1 - \mathbf{s}_2\| \left(e^{-L(x-x_0)} - 1 \right).$$

Damit gilt dann

$$\|\mathbf{y}(x; \mathbf{s}_1) - \mathbf{y}(x; \mathbf{s}_2)\| \leq \|\mathbf{s}_1 - \mathbf{s}_2\| \left(e^{-L(x-x_0)} - 1 \right) = \|\mathbf{s}_1 - \mathbf{s}_2\| \left(e^{L|x-x_0|} - 1 \right).$$

✱

Dieser Satz sagt aus, dass die Lösung eines AWP stetig von den Anfangswerten abhängt. Trotzdem unterscheiden sich Lösungen beliebig stark voneinander. Wie das Beispiel zeigt, ist die im Satz angegebene Abschätzung scharf.

6.2. Einschrittverfahren

6.2.1. Grundlegende Begriffe und Euler-Verfahren

Wir betrachten zunächst ein AWP im \mathbb{R}^1 :
Gesucht ist eine Funktion $y \in C^1[a, b]$ mit

$$y' = f(x, y), \quad y(x_0) = y_0.$$

Die Funktion $f : [a, b] \times \mathbb{R} \rightarrow \mathbb{R}$ erfülle alle Bedingungen aus Satz 6.1. Damit existiert genau eine Lösung des AWP. Als einfachste Methode zur numerischen Behandlung dieses Problems bietet es sich an, die Ableitung $y'(x)$ durch den Differenzenquotienten

$$\frac{y(x+h) - y(x)}{h}$$

mit einer geeigneten Schrittweite h zu ersetzen. Wir erhalten

$$\frac{y(x+h) - y(x)}{h} \approx f(x, y(x)),$$

$$y(x+h) \approx y(x) + hf(x, y(x)).$$

Gehen wir von den exakten Werten $y(x)$ zu Näherungswerten $\eta(x)$ über und ersetzen „ \approx “ durch „ $=$ “, so erhalten wir das folgende Verfahren.

6.4. EULERSches Polygonzugverfahren:

S0 Setze $\eta_0 = \mathbf{y}_0$ und $k = 0$.

S1 Berechne

$$\eta_{k+1} = \eta_k + hf(x_k, \eta_k),$$

$$x_{k+1} = x_k + h.$$

S2 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Das EULERSche Polygonzugverfahren ist ein typisches Einschrittverfahren³. Für das Berechnen der neuen Näherung η_{k+1} an der Stelle x_{k+1} wird nur **eine** alte Näherung η_k an der Stelle x_k benötigt. Ein ESV hat damit folgende Struktur:

6.5. Allgemeines Einschrittverfahren:

S0 Setze $\eta_0 = \mathbf{y}_0$ und $k = 0$.

S1 Berechne

$$\eta_{k+1} = \eta_k + h\Phi(x_k, \eta_k; h)$$

$$x_{k+1} = x_k + h.$$

S2 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Bemerkung: Im Schritt **S1** muss nicht jeweils die gleiche Schrittweite verwendet werden. Wir dürfen ihn also zu:

$$\text{Wähle Schrittweite } h_k \text{ und berechne } \eta_{k+1} = \eta_k + h_k \Phi(x_k, \eta_k; h_k) \text{ und}$$

$$x_{k+1} = x_k + h_k$$

abändern.

Bevor wir weitere Verfahren konstruieren, wollen wir uns mit zwei wichtigen Größen

³Wir werden Einschrittverfahren im weiteren Text durch ESV abkürzen.

beschäftigen, durch die wesentliche Eigenschaften von ESV charakterisiert sind. Es sei z die exakte Lösung des AWP

$$z'(t) = \mathbf{f}(t, z(t)), \quad z(x) = \mathbf{y}$$

und

$$\Delta(x, \mathbf{y}; h) = \begin{cases} \frac{z(x+h) - z(x)}{h} = \frac{z(x+h) - \mathbf{y}}{h} & \text{für } h \neq 0, \\ \mathbf{f}(x, \mathbf{y}) & \text{für } h = 0. \end{cases}$$

Φ sei die Iterationsfunktion eines allgemeinen ESV. Dann heißt die Größe

$$\tau(x, \mathbf{y}; h) = \Delta(x, \mathbf{y}; h) - \Phi(x, \mathbf{y}; h)$$

lokaler Diskretisierungsfehler an der Stelle (x, \mathbf{y}) zur Schrittweite h . Führt man ausgehend von einer Stelle $(\bar{x}, \bar{\mathbf{y}})$ einen Schritt eines ESV aus, so gibt die Größe $h\tau(\bar{x}, \bar{\mathbf{y}}; h)$ an, wie stark die berechnete Näherung

$$\boldsymbol{\eta}(\bar{x} + h) = \bar{\mathbf{y}} + h\Phi(\bar{x}, \bar{\mathbf{y}}; h)$$

von dem Wert abweicht, der sich bei exakter Lösung des AWP

$$z'(x) = \mathbf{f}(x, z(x)), \quad z(\bar{x}) = \bar{\mathbf{y}}$$

an der Stelle $\bar{x} + h$ ergeben würde. Oder anders gesagt: Der lokale Diskretisierungsfehler ist ein Maß dafür, wie gut die Gleichung des Näherungsverfahrens im Punkt $(\bar{x}, \bar{\mathbf{y}})$ durch die exakte Lösung des entsprechenden AWP erfüllt wird.

Eine erste Forderung an ein Näherungsverfahren wäre die, dass der lokale Diskretisierungsfehler für $h \rightarrow 0$ verschwinden soll. Es soll also

$$\lim_{h \rightarrow 0} \tau(x, \mathbf{y}; h) = \lim_{h \rightarrow 0} \Delta(x, \mathbf{y}; h) - \lim_{h \rightarrow 0} \Phi(x, \mathbf{y}; h) = 0$$

und damit

$$\lim_{h \rightarrow 0} \Phi(x, \mathbf{y}; h) = \mathbf{f}(x, \mathbf{y})$$

gelten. Ein Verfahren, das diese Bedingung für alle Funktionen $\mathbf{f} \in \mathbf{F}^1[a, b]$ erfüllt, heißt **konsistent**. Das EULERSche Polygonzugverfahren ist offensichtlich konsistent. Von praktisch größerem Interesse als der lokale Diskretisierungsfehler ist der globale Diskretisierungsfehler.

Es sei \mathbf{y} die exakte Lösung des AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0.$$

$\eta(x; h)$ bezeichne die mit einem ESV und der Schrittweite h berechnete Näherung für \mathbf{y} an der Stelle x . Dann heißt die Größe

$$e(x; h) = \eta(x; h) - \mathbf{y}(x)$$

globaler Diskretisierungsfehler an der Stelle x zur Schrittweite h .

Der globale Diskretisierungsfehler gibt also an, wie stark die berechnete Näherung von der exakten Lösung des AWP abweicht.

Hier wird man natürlich auch fordern, dass der globale Diskretisierungsfehler für $h \rightarrow 0$ verschwinden soll, dass also

$$\lim_{h \rightarrow 0} \eta(x; h) = \mathbf{y}(x)$$

gelte. Ein Verfahren, das diese Bedingung für alle Funktionen $\mathbf{f} \in \mathbf{F}^1[a, b]$ und alle $x \in [a, b]$ erfüllt, heißt **konvergent**. Wie wir später sehen werden hängen Konvergenz und Konsistenz von ESV eng zusammen.

Es ist zu erwarten, dass ein ESV um so bessere Näherungen liefert, je schneller der lokale Diskretisierungsfehler mit $h \rightarrow 0$ gegen Null konvergiert. Das führt zu folgender Definition. Das durch Φ erzeugte ESV zum Lösen des AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

hat die Konsistenzordnung p , falls

$$\|\tau(x, \mathbf{y}; h)\| = O(h^p)$$

für alle $(x, \mathbf{y}) \in S$ und alle Funktionen $\mathbf{f} \in \mathbf{F}^p[a, b]$ gilt.

Die Konsistenzordnung eines Verfahrens lässt sich durch TAYLOR-Entwicklung des lokalen Diskretisierungsfehlers bestimmen. Es gilt

$$\tau(x, \mathbf{y}; h) = \Delta(x, \mathbf{y}; h) - \Phi(x, \mathbf{y}; h) = \frac{\mathbf{z}(x+h) - \mathbf{y}}{h} - \Phi(x, \mathbf{y}; h)$$

mit

$$\mathbf{z}(x+h) = \mathbf{z}(x) + h\mathbf{z}'(x) + \frac{h^2}{2}\mathbf{z}''(x) + \frac{h^3}{6}\mathbf{z}'''(x) + \cdots + \frac{h^p}{p!}\mathbf{z}^{(p)}(x + \vartheta h), \quad \vartheta \in (0, 1).$$

Da \mathbf{z} Lösung des AWP

$$\mathbf{z}'(t) = \mathbf{f}(t, \mathbf{z}(t)), \quad \mathbf{z}(x) = \mathbf{y}$$

ist, gilt

$$\begin{aligned}
 z(x) &= \mathbf{y}, \\
 z'(x) &= \mathbf{f}(x, z(x)) = \mathbf{f}(x, \mathbf{y}) \\
 z''(x) &= \mathbf{f}_x(x, z(x)) + \mathbf{f}_y(x, \mathbf{y})z'(x) = \mathbf{f}_x(x, \mathbf{y}) + \mathbf{f}_y(x, \mathbf{y})\mathbf{f}(x, \mathbf{y}) \\
 z'''(x) &= \mathbf{f}_{xx}(x, \mathbf{y}) + 2\mathbf{f}_{xy}(x, \mathbf{y})\mathbf{f}(x, \mathbf{y}) + \mathbf{f}_{yy}(x, \mathbf{y})[\mathbf{f}(x, \mathbf{y})]^2 + \\
 &\quad + \mathbf{f}_y(x, \mathbf{y})[\mathbf{f}_x(x, \mathbf{y}) + \mathbf{f}_y(x, \mathbf{y})\mathbf{f}(x, \mathbf{y})] \\
 &\quad \vdots \quad \quad \quad \vdots
 \end{aligned}$$

Damit erhält man

$$\begin{aligned}
 \frac{z(x+h) - \mathbf{y}}{h} &= \mathbf{f} + \frac{h}{2}[\mathbf{f}_x + \mathbf{f}_y\mathbf{f}] \\
 &\quad + \frac{h^2}{6}[\mathbf{f}_{xx} + 2\mathbf{f}_{xy}\mathbf{f} + \mathbf{f}_{yy}\mathbf{f}^2 + \mathbf{f}_y(\mathbf{f}_x + \mathbf{f}_y\mathbf{f})] + \dots,
 \end{aligned}$$

wobei der besseren Lesbarkeit halber die Argumente der Funktion \mathbf{f} und ihrer Ableitungen fortgelassen wurden. Hat nun Φ die Entwicklung

$$\Phi(x, \mathbf{y}; h) = \mathbf{f}(x; \mathbf{y}) + h\varphi_1(x; \mathbf{y}) + h^2\varphi_2(x; \mathbf{y}) + h^3\varphi_3(x; \mathbf{y}) + \dots,$$

so ergibt sich

$$\begin{aligned}
 \tau(x, \mathbf{y}; h) &= \frac{h}{2}[\mathbf{f}_x + \mathbf{f}_y\mathbf{f} - 2\varphi_1] \\
 &\quad + \frac{h^2}{6}[\mathbf{f}_{xx} + 2\mathbf{f}_{xy}\mathbf{f} + \mathbf{f}_{yy}\mathbf{f}^2 + \mathbf{f}_y(\mathbf{f}_x + \mathbf{f}_y\mathbf{f}) - 6\varphi_2] + \dots
 \end{aligned}$$

Für das EULER-Verfahren gilt $\Phi(x, \mathbf{y}; h) = \mathbf{f}(x; \mathbf{y})$ und damit

$$\tau(x, \mathbf{y}; h) = \frac{h}{2}[\mathbf{f}_x + \mathbf{f}_y\mathbf{f}] + \dots,$$

also

$$\|\tau(x, \mathbf{y}; h)\| = O(h).$$

Das EULER-Verfahren hat die Konsistenzordnung $p = 1$. Um Verfahren höherer Ordnung zu bekommen, könnte man nun die Funktionen $\varphi_1(x; \mathbf{y})$, $\varphi_2(x; \mathbf{y})$, \dots , so wählen, dass die entsprechenden Glieder der Entwicklung des lokalen Diskretisierungsfehlers verschwinden. Man erhält

$$\Phi(x, \mathbf{y}; h) = \mathbf{f}(x; \mathbf{y}) + \frac{h}{2}[\mathbf{f}_x(x; \mathbf{y}) + \mathbf{f}_y(x; \mathbf{y})\mathbf{f}(x; \mathbf{y})]$$

oder

$$\begin{aligned}\Phi(x, \mathbf{y}; h) = & \mathbf{f}(x; \mathbf{y}) + \frac{h}{2} [\mathbf{f}_x(x; \mathbf{y}) + \mathbf{f}_y(x; \mathbf{y}) \mathbf{f}(x; \mathbf{y})] + \\ & + \frac{h^2}{6} [\mathbf{f}_{xx}(x; \mathbf{y}) + 2\mathbf{f}_{xy}(x; \mathbf{y}) \mathbf{f} + \mathbf{f}_{yy}(x; \mathbf{y}) (\mathbf{f}(x; \mathbf{y}))^2 + \\ & + \mathbf{f}_y(x; \mathbf{y}) [\mathbf{f}_x(x; \mathbf{y}) + \mathbf{f}_y(x; \mathbf{y}) \mathbf{f}(x; \mathbf{y})]].\end{aligned}$$

Diese Verfahren sind aber praktisch wertlos. In vielen Fällen ist für die Funktion \mathbf{f} kein analytischer Ausdruck bekannt, so dass das Berechnen der partiellen Ableitungen nur näherungsweise erfolgen könnte. Ist dagegen ein analytischer Ausdruck für die Funktion \mathbf{f} bekannt, so ist die Auswertung der partiellen Ableitungen oft aufwendiger als die Auswertung der Funktion \mathbf{f} selbst. Dazu kommt noch ein hoher Programmieraufwand zur Implementierung der Ableitungen von \mathbf{f} . Im nächsten Abschnitt werden wir eine Möglichkeit zum Konstruieren von ESV höherer Ordnung kennenlernen, die nur Funktionsberechnungen von \mathbf{f} benötigen.

6.2.2. Runge-Kutta-Verfahren

Wir machen für die Funktion Φ folgenden Ansatz:

$$\Phi(x, \mathbf{y}; h) = \sum_{k=0}^n \alpha_k \mathbf{f}(\xi_k, \boldsymbol{\eta}_k)$$

mit

$$\xi_0 = x, \quad \xi_k = x + \vartheta_k h, \quad k = 1, \dots, n$$

und

$$\boldsymbol{\eta}_0 = \mathbf{y}, \quad \boldsymbol{\eta}_k = \mathbf{y} + h \sum_{l=0}^{k-1} \beta_{kl} \mathbf{f}(\xi_l, \boldsymbol{\eta}_l), \quad k = 1, \dots, n.$$

Die Parameter $\alpha_0, \dots, \alpha_n$, $\vartheta_1, \dots, \vartheta_n$ und $\beta_{10}, \dots, \beta_{n,n-1}$ werden so bestimmt, dass in der Entwicklung des lokalen Diskretisierungsfehlers möglichst viele Glieder verschwinden. Betrachten wir die einfachsten Fälle $n = 0$ und $n = 1$.

$n = 0$ Hier gilt $\Phi(x, \mathbf{y}; h) = \alpha_0 \mathbf{f}(x, \mathbf{y})$, und es folgt sofort $\alpha_0 = 1$. Man erhält das EULER-Verfahren mit der Konsistenzordnung $p = 1$.

$n = 1$ Wir machen für $\Phi(x, \mathbf{y}; h)$ den Ansatz

$$\Phi(x, \mathbf{y}; h) = \alpha_0 \mathbf{f}(x, \mathbf{y}) + \alpha_1 \mathbf{f}(x + \vartheta_1 h, \mathbf{y} + h\beta_{10} \mathbf{f}(x, \mathbf{y})).$$

Es gilt für $h = 0$

$$\Phi(x, \mathbf{y}; 0) = (\alpha_0 + \alpha_1) \mathbf{f}(x, \mathbf{y}).$$

Differentiation nach h liefert

$$\begin{aligned} \frac{d}{dh} \Phi(x, \mathbf{y}; h) &= \alpha_1 \mathbf{f}_x(x + \vartheta_1 h, \mathbf{y} + h\beta_{10} \mathbf{f}(x, \mathbf{y})) \vartheta_1 + \\ &\quad + \alpha_1 \mathbf{f}_y(x + \vartheta_1 h, \mathbf{y} + h\beta_{10} \mathbf{f}(x, \mathbf{y})) \beta_{10} \mathbf{f}(x, \mathbf{y}) \end{aligned}$$

und an der Stelle $h = 0$

$$\left. \frac{d}{dh} \Phi(x, \mathbf{y}; h) \right|_{h=0} = \alpha_1 \vartheta_1 \mathbf{f}_x(x, \mathbf{y}) + \alpha_1 \beta_{10} \mathbf{f}_y(x, \mathbf{y}) \mathbf{f}(x, \mathbf{y}).$$

Durch Koeffizientenvergleich mit der TAYLOR-Entwicklung von $\Delta(x, \mathbf{y}; h)$ erhält man

$$\begin{aligned} h^0 &: \quad \alpha_0 + \alpha_1 &= & 1, \\ h^1 \mathbf{f}_x &: \quad \alpha_1 \vartheta_1 &= & \frac{1}{2}, \\ h^1 \mathbf{f}_y \mathbf{f} &: \quad \alpha_1 \beta_{10} &= & \frac{1}{2}. \end{aligned}$$

als Bestimmungsgleichungen für ein Verfahren der Konsistenzordnung $p = 2$. Aus den beiden letzten Gleichungen folgt $\vartheta_1 = \beta_{10}$. Ansonsten ist ein Parameter frei wählbar. Man erhält verschiedene Verfahren.

- $\alpha_0 = \alpha_1 = 1/2$ und $\beta_{10} = \vartheta_1 = 1$
Wir erhalten das Verfahren von HEUN(1900):

$$\Phi(x, \mathbf{y}; h) = \frac{1}{2} [\mathbf{f}(x, \mathbf{y}) + \mathbf{f}(x + h, \mathbf{y} + h \mathbf{f}(x, \mathbf{y}))].$$

Dieses Verfahren braucht zwei Auswertungen der Funktion \mathbf{f} pro Schritt.

- $\alpha_0 = 0, \alpha_1 = 1$ und $\beta_{10} = \vartheta_1 = 1/2$
Wir erhalten das modifizierte EULER-Verfahren (COLLATZ 1960):

$$\Phi(x, \mathbf{y}; h) = \mathbf{f} \left(x + \frac{h}{2}, \mathbf{y} + \frac{h}{2} \mathbf{f}(x, \mathbf{y}) \right).$$

Dieses Verfahren erfordert ebenfalls zwei Auswertungen der Funktion \mathbf{f} pro Schritt.

Für größere n erhält man entsprechend kompliziertere Bestimmungsgleichungen für die Parameter. Wir geben einige Verfahren nach folgendem Schema an.

$$\begin{array}{c|ccc} \vartheta_1 & \beta_{10} & & \\ \vartheta_2 & \beta_{20} & \beta_{21} & \\ \vartheta_3 & \beta_{30} & \beta_{31} & \beta_{32} \\ \hline & \alpha_0 & \alpha_1 & \alpha_2 & \alpha_3 \end{array}$$

$n = 0$; EULER-Verfahren; $p = 1$

$$\begin{array}{c|c} & \\ \hline & 1 \end{array}$$

$n = 1$; HEUN-Verfahren; $p = 2$

$$\begin{array}{c|cc} 1 & 1 & \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array}$$

$n = 1$; mod. EULER-Verfahren; $p = 2$

$$\begin{array}{c|cc} 1 & 1 & \\ \frac{1}{2} & \frac{1}{2} & \\ \hline & 0 & 1 \end{array}$$

$n = 2$; einfache KUTTA-Regel; $p = 3$

$$\begin{array}{c|ccc} 1 & 1 & & \\ \frac{1}{2} & \frac{1}{2} & & \\ 1 & -1 & 2 & \\ \hline & \frac{1}{6} & \frac{4}{6} & \frac{1}{6} \end{array}$$

$n = 3$; RUNGE-KUTTA-Verfahren; $p = 4$

$$\begin{array}{c|cccc} 1 & 1 & & & \\ \frac{1}{2} & \frac{1}{2} & & & \\ \frac{1}{2} & 0 & \frac{1}{2} & & \\ 1 & 0 & 0 & 1 & \\ \hline & \frac{1}{6} & \frac{2}{6} & \frac{2}{6} & \frac{1}{6} \end{array}$$

$n = 3$; $\frac{3}{8}$ -Regel; $p = 4$

$$\begin{array}{c|cccc} 1 & 1 & & & \\ \frac{2}{3} & \frac{1}{3} & & & \\ \frac{2}{3} & -\frac{1}{3} & 1 & & \\ 1 & 1 & -1 & 1 & \\ \hline & \frac{1}{8} & \frac{3}{8} & \frac{3}{8} & \frac{1}{8} \end{array}$$

$n = 3$; vierstufige ENGLAND-Formel; $p = 4$;

$$\begin{array}{c|cccc} 1 & 1 & & & \\ \frac{1}{2} & \frac{1}{2} & & & \\ \frac{1}{2} & \frac{1}{4} & \frac{1}{4} & & \\ 1 & 0 & -1 & 2 & \\ \hline & \frac{1}{6} & 0 & \frac{2}{3} & \frac{1}{6} \end{array}$$

$n = 3$; GILL-Modifikation des RUNGE-KUTTA-Verfahrens; $p = 4$

$$\begin{array}{c|cccc} 1 & 1 & & & \\ \frac{1}{2} & \frac{1}{2} & & & \\ \frac{1}{2} & \frac{\sqrt{2}-1}{2} & \frac{2-\sqrt{2}}{2} & & \\ 1 & 0 & -\frac{\sqrt{2}}{2} & \frac{2+\sqrt{2}}{2} & \\ \hline & \frac{1}{6} & \frac{2-\sqrt{2}}{6} & \frac{2+\sqrt{2}}{6} & \frac{1}{6} \end{array}$$

$n = 3$; KUNTZMANN-Verfahren; $p = 4$

$$\begin{array}{r|rrrr}
 2 & 2 & & & \\
 5 & 5 & & & \\
 3 & 3 & 15 & & \\
 5 & -\frac{20}{5} & \frac{15}{20} & & \\
 1 & \frac{19}{44} & -\frac{15}{44} & \frac{40}{44} & \\
 \hline
 & \frac{11}{72} & 0 & \frac{25}{72} & \frac{11}{72}
 \end{array}$$

Bemerkungen: (i) Die GILL-Modifikation des RUNGE-KUTTA-Verfahrens besitzt ein günstigeres Rundungsfehlerverhalten als das RUNGE-KUTTA-Verfahren.

(ii) Beim KUNTZMANN-Verfahren fallen in der Entwicklung des lokalen Diskretisierungsfehlers auch noch einige Terme vierter Ordnung weg.

6.2.3. Konvergenz von Einschrittverfahren

Wir betrachten wieder ein AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

und ein ESV

$$\left. \begin{array}{l}
 \boldsymbol{\eta}_0 = \mathbf{y}_0, \\
 \boldsymbol{\eta}_{k+1} = \boldsymbol{\eta}_k + h\Phi(x_k, \boldsymbol{\eta}_k; h), \\
 x_{k+1} = x_0 + (k+1) \cdot h
 \end{array} \right\} \quad k = 0, 1, \dots$$

zur Lösungsberechnung. Wir werden nun untersuchen, wie sich der globale Diskretisierungsfehler

$$\mathbf{e}(x; h) = \boldsymbol{\eta}(x; h) - \mathbf{y}(x)$$

an einer Stelle x für $h \rightarrow 0$ verhält, wobei nur Schrittweiten

$$h \in H_x = \left\{ \frac{x - x_0}{n} \mid n = 1, 2, \dots \right\}$$

zulässig sind. Wir untersuchen

$$\lim_{n \rightarrow \infty} \mathbf{e} \left(x; \frac{x - x_0}{n} \right) = \lim_{n \rightarrow \infty} \mathbf{e}(x; h_n).$$

Wir werden zeigen, dass alle ESV der Konsistenzordnung $p > 0$ konvergent sind, und dass darüber hinaus die Konvergenzordnung mit der Konsistenzordnung übereinstimmt. Für den Beweis des Konvergenzsatzes benötigen wir den folgenden Hilfsatz.

6.6. Satz: *Gelten für die Zahlen ξ_i $i = 0, 1, \dots$ die Ungleichungen*

$$|\xi_{i+1}| \leq (1 + \delta)|\xi_i| + B, \quad i = 0, 1, \dots$$

mit gewissen Konstanten $\delta > 0$ und $B \geq 0$, so gilt

$$|\xi_n| \leq e^{n\delta}|\xi_0| + \frac{e^{n\delta} - 1}{\delta}B, \quad n = 0, 1, \dots$$

Beweis: Aus der Voraussetzung folgt

$$\begin{aligned} |\xi_1| &\leq (1 + \delta)|\xi_0| + B \\ |\xi_2| &\leq (1 + \delta)^2|\xi_0| + (1 + \delta)B + B \\ &= (1 + \delta)^2|\xi_0| + [(1 + \delta) + 1]B \\ |\xi_3| &\leq (1 + \delta)^3|\xi_0| + (1 + \delta)[(1 + \delta) + 1]B + B \\ &= (1 + \delta)^3|\xi_0| + [(1 + \delta)^2 + (1 + \delta) + 1]B \\ &\vdots \\ |\xi_n| &\leq (1 + \delta)^n|\xi_0| + [(1 + \delta)^{n-1} + \dots + (1 + \delta) + 1]B \\ &= (1 + \delta)^n|\xi_0| + \frac{(1 + \delta)^n - 1}{(1 + \delta) - 1}B \\ &= (1 + \delta)^n|\xi_0| + \frac{(1 + \delta)^n - 1}{\delta}B. \end{aligned}$$

Wegen $0 < 1 + \delta \leq e^\delta$ für $\delta > 0$ folgt daraus

$$|\xi_n| \leq (1 + \delta)^n|\xi_0| + \frac{e^{n\delta} - 1}{\delta}B.$$

✱

Nun beweisen wir den eigentlichen Konvergenzsatz.

6.7. Satz: *Gegeben sei für $x_0 \in [a, b]$ und $\mathbf{y}_0 \in \mathbb{R}^d$ das AWP*

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

mit der exakten Lösung $\mathbf{y}(x)$. Die Iterationsfunktion Φ eines ESV habe folgende Eigenschaften:

1. Φ ist stetig auf dem Gebiet

$$G = \left\{ (x, \mathbf{y}, h) \mid x \in [a, b], \|\mathbf{y} - \mathbf{y}(x)\| \leq \gamma, |h| \leq h_0 \right\}$$

mit Konstanten $\gamma > 0$ und $h_0 > 0$.

2. Es existiert eine Konstante $M > 0$, so dass

$$\|\Phi(x, \mathbf{y}_1; h) - \Phi(x, \mathbf{y}_2; h)\| \leq M \|\mathbf{y}_1 - \mathbf{y}_2\|$$

für alle $(x, \mathbf{y}_1, h), (x, \mathbf{y}_2, h) \in G$ gilt.

3. Es existiert eine Konstante $N > 0$, so dass

$$\|\boldsymbol{\tau}(x, \mathbf{y}(x); h)\| = \|\Delta(x, \mathbf{y}(x); h) - \Phi(x, \mathbf{y}(x); h)\| \leq N|h|^p$$

für alle $x \in [a, b]$ und alle $|h| \leq h_0$ gilt.

Dann existiert ein \bar{h} mit $0 < \bar{h} \leq h_0$, so dass

$$e(x, h_n) \leq N \frac{e^{M|x-x_0|} - 1}{M} |h_n|^p$$

für alle $x \in [a, b]$ und alle $h_n = (x - x_0)/n$ mit $|h_n| \leq \bar{h}$ gilt. Für $\gamma = \infty$ ist $\bar{h} = h_0$.

Beweis: Statt der Iterationsfunktion Φ betrachten wir die Funktion

$$\tilde{\Phi} = \begin{cases} \Phi(x, \mathbf{y}; h) & (x, \mathbf{y}, h) \in G, \\ \Phi\left(x, \mathbf{y}(x) + \gamma \frac{\mathbf{y} - \mathbf{y}(x)}{\|\mathbf{y} - \mathbf{y}(x)\|}; h\right) & x \in [a, b], \quad |h| \leq h_0, \quad \|\mathbf{y} - \mathbf{y}(x)\| > \gamma \end{cases}$$

und das durch sie definierte ESV. $\tilde{\Phi}$ ist offensichtlich stetig auf

$$\tilde{G} = \left\{ (x, \mathbf{y}, h) \mid x \in [a, b], \mathbf{y} \in \mathbb{R}^d, |h| \leq h_0 \right\}$$

und genügt der LIPSCHITZ-Bedingung

$$\|\tilde{\Phi}(x, \mathbf{y}_1; h) - \tilde{\Phi}(x, \mathbf{y}_2; h)\| \leq M \|\mathbf{y}_1 - \mathbf{y}_2\|$$

für alle $(x, \mathbf{y}_1, h), (x, \mathbf{y}_2, h) \in \tilde{G}$. Wegen $\tilde{\Phi}(x, \mathbf{y}(x); h) = \Phi(x, \mathbf{y}(x); h)$ gilt auch die dritte Voraussetzung des Satzes für das durch $\tilde{\Phi}$ definierte ESV:

$$\left. \begin{aligned} \tilde{\boldsymbol{\eta}}_0 &= \mathbf{y}_0, \\ \tilde{\boldsymbol{\eta}}_{k+1} &= \tilde{\boldsymbol{\eta}}_k + h\tilde{\Phi}(x_k, \tilde{\boldsymbol{\eta}}_k; h), \\ x_{k+1} &= x_k + h \end{aligned} \right\} \text{für } k = 0, 1, \dots$$

Wir wollen den globalen Diskretisierungsfehler $\tilde{e}(x; h)$ dieses Verfahrens abschätzen. Es gilt

$$\begin{aligned}\tilde{e}_{k+1} &= \tilde{\eta}_{k+1} - \mathbf{y}_{k+1} \\ &= \tilde{\eta}_k + h\tilde{\Phi}(x_k, \tilde{\eta}_k; h) - \mathbf{y}_{k+1} + \mathbf{y}_k - \mathbf{y}_k + h\tilde{\Phi}(x_k, \mathbf{y}_k; h) - h\tilde{\Phi}(x_k, \mathbf{y}_k; h) \\ &= \tilde{e}_k - h \left[\frac{\mathbf{y}_{k+1} - \mathbf{y}_k}{h} - \tilde{\Phi}(x_k, \mathbf{y}_k; h) \right] + h [\tilde{\Phi}(x_k, \tilde{\eta}_k; h) - \tilde{\Phi}(x_k, \mathbf{y}_k; h)] \\ &= \tilde{e}_k - h\boldsymbol{\tau}(x_k, \mathbf{y}_k; h) + h [\tilde{\Phi}(x_k, \tilde{\eta}_k; h) - \tilde{\Phi}(x_k, \mathbf{y}_k; h)].\end{aligned}$$

Damit erhalten wir für die Norm des globalen Diskretisierungsfehlers die Abschätzung

$$\begin{aligned}\|\tilde{e}_{k+1}\| &\leq \|\tilde{e}_k\| + |h| \|\boldsymbol{\tau}(x, \mathbf{y}_k; h)\| + |h| \|\tilde{\Phi}(x_k, \tilde{\eta}_k; h) - \tilde{\Phi}(x, \mathbf{y}_k; h)\| \\ &\leq \|\tilde{e}_k\| + |h|N|h|^p + |h|M\|\tilde{\eta}_k - \mathbf{y}_k\| \\ &= (1 + M|h|) \|\tilde{e}_k\| + N|h|^{p+1}.\end{aligned}$$

Mit Satz 6.6 und $\tilde{e}_0 = 0$ folgt daraus

$$\|\tilde{e}_n\| \leq \frac{e^{n|h|M} - 1}{M|h|} N|h|^{p+1} = \frac{e^{n|h|M} - 1}{M} N|h|^p.$$

Für ein festes $x \in [a, b]$, $x \neq x_0$, ist

$$h = h_n = \frac{x - x_0}{n}$$

und $\tilde{e}_n = \tilde{e}(x, h_n) = \tilde{e}(x_0 + nh_n, h_n)$. Damit folgt

$$\|\tilde{e}(x, h_n)\| \leq \frac{e^{M|x-x_0|} - 1}{M} N|h|^p \leq \frac{e^{M|b-a|} - 1}{M} N|h|^p.$$

Diese Abschätzung gilt für beliebige $x \in [a, b]$ und alle $|h_n| \leq h_0$. Ist nun

$$|h_n| \leq \bar{h} = \min \left\{ \left(\frac{\gamma M}{N(e^{M|b-a|} - 1)} \right)^{1/p}, h_0 \right\} > 0,$$

so gilt für alle h_n mit $|h_n| < \bar{h}$

$$\|\tilde{e}(x, h_n)\| \leq \frac{e^{M|b-a|} - 1}{M} N \frac{\gamma M}{N(e^{M|b-a|} - 1)} = \gamma.$$

Damit gilt aber für alle n $\tilde{\eta}_n \in G$ und $\tilde{\Phi}(x, \tilde{\eta}_n; h) = \Phi(x, \tilde{\eta}_n; h)$. Im Falle $|h_n| < \bar{h}$ geht das durch $\tilde{\Phi}$ erzeugte ESV in das durch Φ erzeugte ESV über. Alle Abschätzungen gelten dann auch für $\boldsymbol{\eta}(x; h_n)$. *

Der letzte Satz sagt aus, dass die Konvergenzordnung eines ESV gleich seiner Konsistenzordnung ist, falls die Iterationsfunktion Φ in einer Umgebung der exakten Lösung des AWP lipschitzstetig ist. Für die RUNGE-KUTTA-Verfahren ergibt sich aber die LIPSCHITZ-Stetigkeit von Φ aus der ohnehin geforderten LIPSCHITZ-Stetigkeit von f .

Man könnte nun das Ergebnis dieses Satzes dazu nutzen, um die Schrittweite zum Erreichen einer vorgegebenen Genauigkeit abzuschätzen. Dazu müsste man näherungsweise die Konstanten M und N kennen. Falls Φ hinreichend oft partiell differenzierbar ist, läuft dies auf die Abschätzung partieller Ableitungen von Φ bzw. f hinaus.

Für das EULER-Verfahren erhält man

$$M \approx \left\| \frac{\partial \Phi(x, \mathbf{y}; h)}{\partial \mathbf{y}} \right\| = \|\mathbf{f}_y(x, \mathbf{y})\|, \quad N \approx \frac{1}{2} \|\mathbf{f}_x(x, \mathbf{y}) + \mathbf{f}_y(x, \mathbf{y})\mathbf{f}(x, \mathbf{y})\|.$$

Für das HEUN-Verfahren erhält man

$$M \approx \left\| \frac{\partial \Phi(x, \mathbf{y}; h)}{\partial \mathbf{y}} \right\| = \|\mathbf{f}_y(x, \mathbf{y})\|$$

und

$$N \approx \frac{1}{12} \|\mathbf{f}_{xx} + 2\mathbf{f}_{xy}\mathbf{f} + \mathbf{f}_{yy}\mathbf{f}^2 - 2\mathbf{f}_y(\mathbf{f}_x + \mathbf{f}_y\mathbf{f})\|.$$

Praktisch ist das kaum durchführbar.

6.2.4. Rundungsfehlereinfluss

Wir betrachten das ESV

$$\left. \begin{aligned} \eta_0 &= \mathbf{y}_0, \\ \eta_{k+1} &= \eta_k + h\Phi(x_k, \eta_k; h), \\ x_{k+1} &= x_k + h \end{aligned} \right\}, \quad k = 0, 1, \dots$$

Auf einem Rechner treten Rundungsfehler auf, so dass das Verfahren in der Form

$$\left. \begin{aligned} \hat{\eta}_0 &= \mathbf{y}_0, \\ \hat{\eta}_{k+1} &= \hat{\eta}_k + h\Phi(x_k, \hat{\eta}_k; h) + \delta_{k+1}, \\ x_{k+1} &= x_k + h \end{aligned} \right\}, \quad k = 0, 1, \dots$$

realisiert wird; die Größe δ_k stellt dabei den pro Schritt erzeugten Rundungsfehler dar. Es gilt der folgende Satz.

6.8. Satz: *Es sei zum Lösen eines AWP obiges ESV gegeben. Die Funktion Φ genüge den Voraussetzungen aus Satz 6.7 und der pro Schritt erzeugte Rundungsfehler sei beschränkt:*

$$\|\delta_k\| \leq \Delta < \infty, \quad k = 1, 2, \dots$$

*Für hinreichend großes γ gilt dann:
Der gesamte Rundungsfehler*

$$\mathbf{r}(x_k; h) = \hat{\boldsymbol{\eta}}(x_k; h) - \boldsymbol{\eta}(x_k; h)$$

genügt der Abschätzung

$$\|\mathbf{r}(x_k; h)\| \leq \frac{\Delta}{|h|} \frac{e^{M|x_k - x_0|} - 1}{M}.$$

Für den gesamten absoluten Fehler

$$\mathbf{v}(x_k; h) = \hat{\boldsymbol{\eta}}(x_k; h) - \mathbf{y}(x_k)$$

gilt die Abschätzung

$$\|\mathbf{v}(x_k; h)\| \leq \left(N|h|^p + \frac{\Delta}{|h|} \right) \frac{e^{M|x_k - x_0|} - 1}{M}.$$

Beweis: Wir betrachten wieder, wie im Beweis von Satz 6.7, das durch die Funktion $\tilde{\Phi}$ erzeugte Verfahren

$$\left. \begin{aligned} \tilde{\boldsymbol{\eta}}_0 &= \mathbf{y}_0, \\ \tilde{\boldsymbol{\eta}}_{k+1} &= \tilde{\boldsymbol{\eta}}_k + h\tilde{\Phi}(x_k, \tilde{\boldsymbol{\eta}}_k; h), \\ x_{k+1} &= x_k + h \end{aligned} \right\}, \quad k = 0, 1, \dots$$

Die Berücksichtigung von Rundungsfehlern führt auf das Verfahren

$$\left. \begin{aligned} \hat{\boldsymbol{\eta}}_0 &= \mathbf{y}_0, \\ \hat{\boldsymbol{\eta}}_{k+1} &= \hat{\boldsymbol{\eta}}_k + h\tilde{\Phi}(x_k, \hat{\boldsymbol{\eta}}_k; h) + \tilde{\boldsymbol{\delta}}_{k+1}, \\ x_{k+1} &= x_k + h \end{aligned} \right\}, \quad k = 0, 1, \dots$$

mit

$$\|\tilde{\boldsymbol{\delta}}_k\| \leq \tilde{\Delta} < \infty, \quad k = 1, 2, \dots$$

Hieraus ergibt sich

$$\begin{aligned}\tilde{\mathbf{r}}_{k+1} &= \hat{\boldsymbol{\eta}}_{k+1} - \tilde{\boldsymbol{\eta}}_{k+1} \\ &= \hat{\boldsymbol{\eta}}_k + h\tilde{\Phi}(x_k, \hat{\boldsymbol{\eta}}_k; h) + \tilde{\boldsymbol{\delta}}_{k+1} - \tilde{\boldsymbol{\eta}}_k - h\tilde{\Phi}(x_k, \tilde{\boldsymbol{\eta}}_k; h) \\ &= \tilde{\mathbf{r}}_k + h[\tilde{\Phi}(x_k, \hat{\boldsymbol{\eta}}_k; h) - \tilde{\Phi}(x_k, \tilde{\boldsymbol{\eta}}_k; h)] + \tilde{\boldsymbol{\delta}}_{k+1}\end{aligned}$$

und

$$\begin{aligned}\|\tilde{\mathbf{r}}_{k+1}\| &\leq \|\tilde{\mathbf{r}}_k\| + |h| \|\tilde{\Phi}(x_k, \hat{\boldsymbol{\eta}}_k; h) - \tilde{\Phi}(x_k, \tilde{\boldsymbol{\eta}}_k; h)\| + \|\tilde{\boldsymbol{\delta}}_{k+1}\| \\ &\leq \|\tilde{\mathbf{r}}_k\| + |h|M \|\hat{\boldsymbol{\eta}}_k - \tilde{\boldsymbol{\eta}}_k\| + \|\tilde{\boldsymbol{\delta}}_{k+1}\| \\ &\leq (1 + M|h|)\|\tilde{\mathbf{r}}_k\| + \tilde{\Delta}.\end{aligned}$$

Mit Satz 6.6 und $\|\tilde{\mathbf{r}}_0\| = 0$ folgt die Abschätzung

$$\|\tilde{\mathbf{r}}_k\| \leq \tilde{\Delta} \frac{e^{kM|h|} - 1}{M|h|} = \frac{\tilde{\Delta}}{|h|} \frac{e^{M|x_k - x_0|} - 1}{M}.$$

Für den gesamten absoluten Fehler erhalten wir

$$\begin{aligned}\|\tilde{\mathbf{v}}_k\| &\leq \|\tilde{\mathbf{r}}_k\| + \|\tilde{\mathbf{e}}_k\| \\ &\leq \frac{\tilde{\Delta}}{|h|} \frac{e^{M|x_k - x_0|} - 1}{M} + N|h|^p \frac{e^{M|x_k - x_0|} - 1}{M} \\ &= \left(N|h|^p + \frac{\tilde{\Delta}}{|h|} \right) \frac{e^{M|x_k - x_0|} - 1}{M}\end{aligned}$$

Ist nun γ hinreichend groß, so dass Schrittweiten existieren, für die

$$\|\tilde{\mathbf{v}}_k\| = \|\tilde{\mathbf{v}}(x_k; h)\| \leq \gamma$$

für alle $x_k \in [a, b]$ gilt, so stimmt $\tilde{\Phi}$ für diese Schrittweiten mit Φ überein. Alle Abschätzungen gelten dann auch für das ursprüngliche Verfahren mit Δ statt $\tilde{\Delta}$. *

Betrachten wir nun noch den pro Schritt erzeugten Rundungsfehler genauer. Wir beschränken uns dabei auf den eindimensionalen Fall $y : D \subset \mathbb{R} \rightarrow \mathbb{R}$. Das Berechnen von $\hat{\eta}_{k+1}$ erfolgt in folgenden Schritten:

$$a_k = gl(\Phi(x_k, \hat{\eta}_k; h)) = \Phi(x_k, \hat{\eta}_k; h)(1 + \alpha_k)$$

$$b_k = gl(ha_k) = ha_k(1 + \beta_k)$$

$$\begin{aligned}
\hat{\eta}_{k+1} &= gl(\hat{\eta}_k + b_k) = (\hat{\eta}_k + b_k)(1 + \gamma_k) \\
&= (\hat{\eta}_k + ha_k(1 + \beta_k))(1 + \gamma_k) \\
&= (\hat{\eta}_k + h\Phi(x_k, \hat{\eta}_k; h)(1 + \alpha_k)(1 + \beta_k))(1 + \gamma_k) \\
&\doteq (\hat{\eta}_k + h\Phi(x_k, \hat{\eta}_k; h)(1 + \alpha_k + \beta_k))(1 + \gamma_k) \\
&\doteq \hat{\eta}_k + h\Phi(x_k, \hat{\eta}_k; h) + h\Phi(x_k, \hat{\eta}_k; h)(\alpha_k + \beta_k + \gamma_k) + \hat{\eta}_k\gamma_k.
\end{aligned}$$

Damit gilt

$$\delta_{k+1} \doteq h\Phi(x_k, \hat{\eta}_k; h)(\alpha_k + \beta_k + \gamma_k) + \hat{\eta}_k\gamma_k = h\Phi(x_k, \hat{\eta}_k; h)(\alpha_k + \beta_k) + \hat{\eta}_{k+1}\gamma_k.$$

Normalerweise ist die Schrittweite so klein, dass

$$|h\Phi(x_k, \hat{\eta}_k; h)y| \ll |\hat{\eta}_{k+1}|$$

gilt. Dann wird der Rundungsfehler im wesentlichen durch den Term $\hat{\eta}_{k+1}\gamma_k$, also den Fehler bei der Addition bestimmt. Es ist daher nützlich, die Addition mit einer höheren Genauigkeit durchzuführen.

6.2.5. Schrittweitensteuerung

Wie wir gesehen haben, ist die Schätzung einer günstigen Schrittweite zum Erreichen einer vorgegebenen Genauigkeit mittels Satz 6.7 praktisch nicht möglich. In diesem Abschnitt werden wir zwei andere Möglichkeiten zur Schätzung des globalen Diskretisierungsfehlers und damit zur Schrittweitensteuerung kennenlernen. Beide Verfahren beruhen darauf, dass man mit verschiedenen Methoden zwei Näherungen $\eta_I(x)$ und $\eta_{II}(x)$ berechnet und aus der Differenz $\eta_I(x) - \eta_{II}(x)$ auf die Größe des globalen Diskretisierungsfehlers schließt.

Betrachten wir ein ESV p -ter Ordnung. Es gilt

$$e(x; h) = \eta(x; h) - \mathbf{y}(x) = \mathbf{e}_p(x)h^p + \mathbf{O}(h^{p+1})$$

mit $\mathbf{e}_p(x_0) = \mathbf{o}$. Bei Vernachlässigung des Restgliedes folgt

$$\eta(x; h) \doteq \mathbf{y}(x) + \mathbf{e}_p(x)h^p.$$

Berechnet man mit der Schrittweite $h/2$ eine weitere Näherung, so gilt

$$\eta\left(x; \frac{h}{2}\right) \doteq \mathbf{y}(x) + \mathbf{e}_p(x)\left(\frac{h}{2}\right)^p.$$

Nun lässt sich $\mathbf{e}_p(x)$ näherungsweise berechnen. Es ergibt sich

$$\eta(x; h) - \eta\left(x; \frac{h}{2}\right) \doteq \mathbf{e}_p(x)\left(h^p - \left(\frac{h}{2}\right)^p\right)$$

und

$$e_p(x) \left(\frac{h}{2}\right)^p \doteq \frac{\eta(x; h) - \eta\left(x; \frac{h}{2}\right)}{2^p - 1}.$$

Diese Gleichung wird zunächst verwendet, um die berechneten Näherungen zu verbessern. Man erhält

$$y(x) \doteq \bar{\eta}(x) = \eta\left(x; \frac{h}{2}\right) - \frac{\eta(x; h) - \eta\left(x; \frac{h}{2}\right)}{2^p - 1} = \frac{2^p \eta\left(x; \frac{h}{2}\right) - \eta(x; h)}{2^p - 1}.$$

Dies entspricht dem Extrapolationsprinzip im ROMBERG-Verfahren.

Auch die beiden Näherungen $\eta(x; h)$ und $\eta(x; h/2)$ lassen sich verwenden, um eine bessere Schrittweite zu bestimmen. Wegen $e_p(x_0) = \mathbf{o}$ gilt

$$e_p(x) \doteq (x - x_0)e'_p(x_0).$$

Damit folgt weiter

$$e(x_0 + h; h) \doteq e_p(x_0 + h)h^p \doteq (x_0 + h - x_0)e'_p(x_0)h^p = e'_p(x_0)h^{p+1}$$

und

$$\begin{aligned} e\left(x_0 + h; \frac{h}{2}\right) &\doteq e_p(x_0 + h) \left(\frac{h}{2}\right)^p \\ &\doteq (x_0 + h - x_0)e'_p(x_0) \left(\frac{h}{2}\right)^p = e'_p(x_0)h \left(\frac{h}{2}\right)^p. \end{aligned}$$

Aus diesen beiden Gleichungen lässt sich $e'_p(x_0)$ näherungsweise berechnen. Man erhält

$$e'_p(x_0) \doteq \frac{2^p}{2^p - 1} \frac{\eta(x_0 + h; h) - \eta\left(x_0 + h; \frac{h}{2}\right)}{h^{p+1}}.$$

Von einer guten Schrittweite \bar{h} wird man fordern, dass der globale Diskretisierungsfehler unterhalb einer vorgegebenen Genauigkeitsschranke ε liegt:

$$\|e(x_0 + \bar{h}; \bar{h})\| \doteq \|e'_p(x_0)\| |\bar{h}|^{p+1} \leq \varepsilon.$$

Mit dem geschätzten $e'_p(x_0)$ liefert das folgende Schrittweiteschätzung

$$|\bar{h}|^{p+1} \leq \frac{2^p - 1}{2^p} \frac{\varepsilon}{\left\| \eta(x_0 + h; h) - \eta\left(x_0 + h; \frac{h}{2}\right) \right\|} |h|^{p+1}$$

oder

$$|\bar{h}| \leq |h|^{p+1} \sqrt{\frac{2^p - 1}{2^p} \frac{\varepsilon}{\left\| \boldsymbol{\eta}(x_0 + h; h) - \boldsymbol{\eta}\left(x_0 + h; \frac{h}{2}\right) \right\|}}.$$

Damit ergibt sich der folgende Algorithmus:

6.9. Schrittweitensteuerung bei ESV:

Es sei ein ESV mit der Konsistenzordnung p zum Lösen des AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

gegeben.

S0 Wähle Grundschriftweite H und eine zu erreichende Genauigkeit ε .

S1 Berechne mit dem gegebenen ESV die beiden Näherungen

$$\boldsymbol{\eta}(x_0 + H; H), \quad \boldsymbol{\eta}(x_0 + H; H/2).$$

S2 Berechne

$$\frac{H}{h} = \sqrt[p+1]{\frac{2^p}{2^p - 1} \frac{\left\| \boldsymbol{\eta}(x_0 + H; H) - \boldsymbol{\eta}\left(x_0 + H; \frac{H}{2}\right) \right\|}{\varepsilon}}.$$

S3 Ist $H/h \gg 2$, so setze $H = 2h$ und gehe zu Schritt **S1**.

S4 Setze

$$\begin{aligned} x_0 &= x_0 + H, \\ \mathbf{y}_0 &= \boldsymbol{\eta}\left(x_0 + h; \frac{H}{2}\right), \\ H &= 2h \end{aligned}$$

und gehe zu Schritt **S1**.

Bemerkungen: (i) In der Praxis wird man die Bedingung $H/h \gg 2$ in Schritt **S3** durch $H/h \geq 3..5$ ersetzen.

(ii) ε sollte man nicht kleiner als $K\varepsilon$ wählen, wobei K eine obere Schranke für den Betrag der Lösung im betrachteten Bereich ist. Also

$$K \approx \max \{ \|\mathbf{y}(x)\| \mid x \in [x_0, x_0 + H] \}.$$

Runge-Kutta-Fehlberg-Verfahren

Bei dem oben angegebenen Verfahren **6.9** zur Schrittweitensteuerung benötigt man pro Schritt mindestens drei Auswertungen der Iterationsfunktion Φ . Wir werden sehen, dass man auch mit geringerem Aufwand entsprechende Schätzungen für den globalen Diskretisierungsfehler und damit eine bessere Schrittweite erhält.

Wir betrachten zwei ESV. Diese seien durch die Iterationsfunktionen Φ_1 und Φ_2 gegeben. Die Konsistenzordnung des ersten Verfahrens sei p und die des zweiten Verfahrens sei $p+1$. Für die lokalen Diskretisierungsfehler gilt dann

$$\begin{aligned}\tau_1(x, \mathbf{y}; h) &= \Delta(x, \mathbf{y}; h) - \Phi_1(x, \mathbf{y}; h) = \mathbf{C}_1(x)h^p + \mathcal{O}(h^{p+1}), \\ \tau_2(x, \mathbf{y}; h) &= \Delta(x, \mathbf{y}; h) - \Phi_2(x, \mathbf{y}; h) = \mathbf{C}_2(x)h^{p+1} + \mathcal{O}(h^{p+2}).\end{aligned}$$

Ausgehend von einer Näherung $\boldsymbol{\eta}(x)$ berechnen wir mit beiden Verfahren Näherungen an der Stelle $x+H$:

$$\begin{aligned}\boldsymbol{\eta}_1(x+H; H) &= \boldsymbol{\eta}(x) + H\Phi_1(x, \boldsymbol{\eta}(x); H), \\ \boldsymbol{\eta}_2(x+H; H) &= \boldsymbol{\eta}(x) + H\Phi_2(x, \boldsymbol{\eta}(x); H).\end{aligned}$$

Für die Differenz der beiden Näherungen gilt

$$\begin{aligned}\boldsymbol{\eta}_1(x+H; H) - \boldsymbol{\eta}_2(x+H; H) &= H(\Phi_1(x, \boldsymbol{\eta}(x); H) - \Phi_2(x, \boldsymbol{\eta}(x); H)) \\ &= \mathbf{C}_1(x)H^{p+1} + \mathcal{O}(H^{p+2}).\end{aligned}$$

In erster Näherung gilt damit

$$\mathbf{C}_1(x) \doteq \frac{\boldsymbol{\eta}_1(x+H; H) - \boldsymbol{\eta}_2(x+H; H)}{H^{p+1}}.$$

Von einem erfolgreichen Schritt mit der Schrittweite h werden wir fordern, dass

$$\varepsilon \geq \|\boldsymbol{\eta}_1(x+h; h) - \boldsymbol{\eta}_2(x+h; h)\| \doteq \|\mathbf{C}_1(x)\| |h|^{p+1}$$

gilt. Mit der obigen Schätzung für $\mathbf{C}_1(x)$ liefert das die Schrittweitschätzung

$$h \approx H \sqrt[p+1]{\frac{\varepsilon}{\|\boldsymbol{\eta}_1(x+H; H) - \boldsymbol{\eta}_2(x+H; H)\|}}.$$

Es ergibt sich der folgende Algorithmus:

6.10. RUNGE-KUTTA-FEHLBERG-Verfahren:

Es seien zwei ESV mit den Konsistenzordnungen p und $p+1$ zum Lösen des AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

gegeben.

S0 Wähle Grundschriftweite H und eine zu erreichende Genauigkeit ε .

S1 Berechne mit den gegebenen ESV die Näherungen

$$\eta_1(x_0 + H; H), \quad \eta_2(x_0 + H; H).$$

S2 Berechne

$$\frac{H}{h} = \sqrt[p+1]{\frac{\|\eta_1(x_0 + H; H) - \eta_2(x_0 + H; H)\|}{\varepsilon}}.$$

S3 Ist $H/h \gg 2$, so setze $H = 2h$ und gehe zu Schritt **S1**.

S4 Setze

$$\begin{aligned} x_0 &= x_0 + H, \\ y_0 &= \eta_2(x_0 + H; H), \\ H &= 2h \end{aligned}$$

und gehe zu Schritt **S1**.

Bemerkungen: (i) Auf den ersten Blick scheint der Aufwand beim Algorithmus **6.10** nicht geringer als beim Algorithmus **6.9** zu sein. Rechnet man beim Algorithmus **6.9** k Auswertungen der Funktion f für einen Iterationsschritt, so benötigt dieser Algorithmus mindestens $3k$ Auswertungen der Funktion f pro Schritt. Rechnet man beim Algorithmus **6.10** k Auswertungen der Funktion f für das erste Verfahren und $k + 1$ Auswertungen der Funktion f für das zweite Verfahren, so benötigt dieser Algorithmus mindestens $2k + 1$ Auswertungen der Funktion f pro Schritt. Es lassen sich aber ESV der Ordnung $p + 1$ konstruieren, bei denen ein ESV der Ordnung p als Zwischenergebnis abfällt. So zum Beispiel

$$\begin{array}{r|l} \frac{1}{4} & \frac{1}{4} \\ \frac{27}{40} & \frac{189}{800} \quad \frac{729}{800} \\ 1 & \frac{214}{891} \quad \frac{27}{891} \quad \frac{650}{891} \\ \hline p=2 & \frac{214}{891} \quad \frac{27}{891} \quad \frac{650}{891} \\ p=3 & \frac{533}{2106} \quad 0 \quad \frac{1600}{2106} \quad -\frac{27}{2106} \end{array}.$$

Wir berechnen

$$\begin{aligned} \mathbf{k}_1 &= \mathbf{f}(x, \mathbf{y}), \\ \mathbf{k}_2 &= \mathbf{f}\left(x + \frac{1}{4}H, \mathbf{y} + \frac{1}{4}H\mathbf{k}_1\right), \\ \mathbf{k}_3 &= \mathbf{f}\left(x + \frac{27}{40}H, \mathbf{y} - \frac{189}{800}H\mathbf{k}_1 + \frac{79}{800}H\mathbf{k}_2\right), \\ \eta_1 &= \mathbf{y} + H\left(\frac{214}{891}\mathbf{k}_1 + \frac{27}{891}\mathbf{k}_2 + \frac{650}{891}\mathbf{k}_3\right), \\ \mathbf{k}_4 &= \mathbf{f}(x + H\eta_1), \\ \eta_2 &= \mathbf{y} + H\left(\frac{533}{2106}\mathbf{k}_1 + \frac{1600}{2106}\mathbf{k}_3 + \frac{27}{2106}\mathbf{k}_4\right). \end{aligned}$$

Der Aufwand des gesamten Verfahrens wird durch das Verfahren mit der höheren Konvergenzordnung bestimmt. Beim Vergleich der Algorithmen **6.9** und **6.10** würde man damit ein Verhältnis von $3k$ zu $k + 1$ Funktionsauswertungen pro Schritt erhalten. Damit sind die RUNGE-KUTTA-FEHLBERG-Verfahren bedeutend effektiver.

(ii) Man beachte die beiden Bemerkungen zum Algorithmus **6.9**.

Anwendung der Schrittweitensteuerung auf die numerische Integration

Die Aufgabe der numerischen Integration ist auch als AWP formulierbar. Das Berechnen von

$$I(f) = \int_a^b f(x) dx$$

ist dem Problem des Berechnens des Wertes $y(b)$ der Lösung des AWP

$$y'(x) = f(x), \quad y(a) = 0$$

äquivalent. Ist nun durch

$$Q_n(f) = (b-a) \sum_{i=0}^n w_i f(x_i), \quad x_i = a + \vartheta_i(b-a)$$

eine Quadraturformel zum näherungsweise Berechnen von $I(f)$ gegeben, so entspricht dem das ESV

$$\left. \begin{aligned} \eta_0 &= y_0, \\ \eta_{k+1} &= \eta_k + h\Phi(x_k, \eta_k; h), \\ x_{k+1} &= x_k + h \end{aligned} \right\}, \quad k = 0, 1, \dots$$

mit

$$\Phi(x_k, \eta_k; h) = \sum_{i=0}^n w_i f(x_k + \vartheta_i h).$$

Hat die Quadraturformel den Exaktheitsgrad q , so gilt

$$\int_0^1 x^l dx = Q_n(x^l) = \sum_{i=0}^n w_i \vartheta_i^l$$

für $l = 0, \dots, q$. Daraus ergibt sich

$$\sum_{i=0}^n w_i \vartheta_i^l = \frac{1}{l+1}, \quad l = 0, \dots, q.$$

Die ersten q Glieder der TAYLOR-Entwicklung

$$\Phi(x, y; h) = \varphi_0(x, y) + \varphi_1(x, y)h + \varphi_2(x, y)h^2 + \dots$$

von Φ lauten dann

$$\varphi_l(x, y)h^l = \frac{h^l}{l!} \frac{d^l}{dh^l} \Phi(x, y; h)|_{h=0} = \frac{h^l}{l!} \sum_{i=0}^n w_i f^{(l)}(x) \vartheta_i^l = \frac{h^l}{(l+1)!} f^{(l)}(x).$$

Damit ergibt sich für den lokalen Diskretisierungsfehler

$$\tau(x, y; h) = O(h^{q+1}).$$

Das ESV, das der Quadraturformel $Q_n(f)$ entspricht, hat daher die Konsistenzordnung $p = q + 1$. Damit ist das erste Verfahren zur Schrittweitensteuerung von ESV (Algorithmus 6.9) unmittelbar auf die numerische Integration anwendbar. Wir erhalten den folgenden Algorithmus.

6.11. Schrittweitensteuerung für Quadraturverfahren:

Es sei durch

$$Q_n(f) = (b-a) \sum_{i=0}^n w_i f(a + \vartheta_i(b-a))$$

ein Quadraturverfahren mit dem Exaktheitsgrad q zum näherungsweise Berechnen von

$$I(f) = \int_a^b f(x) dx$$

gegeben.

S0 Wähle Grundschriftweite $H \leq b - a$ und eine zu erreichende Genauigkeit ε .

Setze $x_0 = a$, $\eta_0 = 0$ und $k = 0$.

S1 Falls $x_k = b$ so setze $I(f) = \eta_k$. STOPP

S2 Berechne die Näherungen

$$\bar{\eta}_{k+1} = \eta_k + H \sum_{i=0}^n w_i f(a + \vartheta_i H)$$

und

$$\hat{\eta}_{k+1} = \eta_k + \frac{H}{2} \sum_{i=0}^n w_i f\left(a + \vartheta_i \frac{H}{2}\right) + \frac{H}{2} \sum_{i=0}^n w_i f\left(a + \frac{H}{2} + \vartheta_i \frac{H}{2}\right).$$

S3 Berechne

$$\frac{H}{h} = \sqrt[q+2]{\frac{2^{q+1}}{2^{q+1} - 1} \frac{\|\bar{\eta}_{k+1} - \hat{\eta}_{k+1}\|}{\varepsilon}}.$$

S4 Ist $H/h \gg 2$, so setze $H = 2h$ und gehe zu Schritt **S2**.

S5 Setze

$$\begin{aligned} x_{k+1} &= x_k + H, \\ \eta_{k+1} &= \hat{\eta}_{k+1}, \\ H &= \min\{2h, b - x_{k+1}\}, \\ k &= k + 1 \end{aligned}$$

und gehe zu Schritt **S1**.

6.2.6. Steife Differentialgleichungen und implizite Verfahren

Viele AWP aus der Praxis bereiten beim Lösen eigenartige Schwierigkeiten. Betrachten wir dazu als Beispiel das AWP

$$y' = \mu y, \quad y(0) = 1$$

mit der exakten Lösung $y(x) = e^{\mu x}$. Ist $\mu < 0$, so strebt die Lösung für $x \rightarrow \infty$ gegen Null. Für hinreichend große x ändert sich die Lösung nur noch unwesentlich. Man sollte annehmen, dass man in diesem Bereich mit großen Schrittweiten arbeiten könnte. Verwendet man nun das EULERSche Polygonzugverfahren, erhält man die Näherungen

$$\eta_{k+1} = \eta_k + \mu h \eta_k = (1 + \mu h) \eta_k,$$

folglich

$$\eta_k = (1 + \mu h)^k \eta_0 = (1 + \mu h)^k.$$

Die η_k konvergieren offensichtlich nur für $|1 + \mu h| < 1$ gegen Null. Positive Näherungslösungen erhält man nur für $1 + \mu h > 0$. Damit ergibt sich für „vernünftige“ Schrittweiten die Forderung

$$-1 < \mu h < 0,$$

und daraus wegen $\mu < 0$

$$0 < h < -\frac{1}{\mu}.$$

Je größer $|\mu|$, also je schneller die exakte Lösung gegen Null konvergiert, desto kleinere Schrittweiten muss man wählen, um vernünftige Näherungslösungen zu erhalten.

Auch für Verfahren höherer Ordnung ergibt sich die gleiche Situation. Das HEUN-Verfahren liefert für das obige AWP die Näherungslösungen

$$\eta_k = \left(1 + \mu h + \frac{\mu^2 h^2}{2}\right)^k.$$

Damit wieder $0 < \eta_k$ und $\eta_k \rightarrow 0$ für $k \rightarrow \infty$ gilt, muss die Schrittweite h die Ungleichungen

$$0 < 1 + \mu h + \frac{\mu^2 h^2}{2} < 1$$

erfüllen. Daraus ergibt sich

$$-1 < (1 + \mu h)^2 < 1$$

und weiter

$$|1 + \mu h| < 1,$$

also

$$h < -\frac{2}{\mu}.$$

Das ist nicht wesentlich besser als beim EULER-Verfahren. Ein Verfahren, das besser zum Lösen des obigen AWP geeignet ist, erhalten wir folgendermaßen. Wir integrieren die Differentialgleichung von x bis $x + h$ und erhalten

$$\int_x^{x+h} \mathbf{y}'(t) dt = \int_x^{x+h} \mathbf{f}(t, \mathbf{y}(t)) dt,$$

und daraus

$$\mathbf{y}(x+h) = \mathbf{y}(x) + \int_x^{x+h} \mathbf{f}(t, \mathbf{y}(t)) dt.$$

Wendet man die Rechteckregel mit der linken Intervallgrenze zur Approximation des Integrals an, erhält man das bekannte EULER-Verfahren

$$\boldsymbol{\eta}(x+h) = \boldsymbol{\eta}(x) + h\mathbf{f}(x, \boldsymbol{\eta}(x)).$$

Wendet man jedoch die Rechteckregel mit der rechten Intervallgrenze an, erhält man

$$\boldsymbol{\eta}(x+h) = \boldsymbol{\eta}(x) + h\mathbf{f}(x+h, \boldsymbol{\eta}(x+h)).$$

Hier ergibt sich die neue Näherung $\boldsymbol{\eta}(x+h)$ aus einem im allgemeinen nichtlinearen Gleichungssystem. Das Verfahren wird als **implizites EULER-Verfahren** bezeichnet. Für obiges eindimensionales AWP ($\mathbf{f}(x, \mathbf{y}) = \mu y$) lässt sich das Gleichungssystem nach $\eta(x+h)$ auflösen. Es ergibt sich

$$\eta(x+h) = \frac{1}{1-\mu h} \eta(x).$$

Damit hat die k -te Näherung die Darstellung

$$\eta_k = \eta(x_k) = \frac{1}{(1-\mu h)^k}.$$

Die Forderungen

$$\lim_{k \rightarrow \infty} \eta_k = 0, \quad \eta_k > 0$$

liefern dann im Falle $\mu < 0$

$$0 < \frac{1}{1-\mu h} < 1,$$

folglich

$$h > 0.$$

Das Verfahren liefert für beliebige positive Schrittweiten vernünftige Lösungen. Wir wollen nun Differentialgleichungen mit diesen Schwierigkeiten genauer charakterisieren. Das lineare Differentialgleichungssystem

$$\mathbf{y}' = \mathbf{A}\mathbf{y}, \quad \mathbf{A} \in \mathbb{R}^{d \times d}$$

heißt **steif**, falls alle Eigenwerte λ_i der Matrix \mathbf{A} einen negativen Realteil haben und zusätzlich

$$q = \frac{\max \{ |\Re(\lambda_i)| \mid i = 1, \dots, d \}}{\min \{ |\Re(\lambda_i)| \mid i = 1, \dots, d \}} \gg 1$$

gilt. Für $q \approx 10$ heißt das Differentialgleichungssystem **schwach steif**, für $q \geq 10$ **steif**.

Bemerkung: Der Begriff der Steifheit ist auf beliebige Differentialgleichungen übertragbar. Die Differentialgleichung

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y})$$

wird durch die Einführung der Funktionen $\bar{y}(x) = x$ und

$$\tilde{\mathbf{y}}(x) = \begin{pmatrix} \mathbf{y}(x) \\ \bar{y}(x) \end{pmatrix}$$

auf die Gestalt

$$\tilde{\mathbf{y}}' = \tilde{\mathbf{f}}(\tilde{\mathbf{y}}) = \begin{pmatrix} \mathbf{f}(x, \mathbf{y}) \\ 1 \end{pmatrix}$$

gebracht. Durch Linearisierung in der Nähe eines Punktes $\tilde{\mathbf{y}}_0$ erhält man dann eine Differentialgleichung der Form $\mathbf{y}' = \mathbf{A}\mathbf{y}$. Damit lässt sich für eine beliebige Differentialgleichung lokal der Steifheitsbegriff erklären.

Wendet man nun ein ESV zum Lösen eines AWP der Form

$$\mathbf{y}' = \mathbf{A}\mathbf{y}, \quad \mathbf{y}(x_0) = \mathbf{y}_0,$$

an, lässt sich für die Folge $\{\boldsymbol{\eta}_k\}_{k \in \mathbb{N}}$ der Näherungslösungen meist eine Rekursionsformel

$$\boldsymbol{\eta}_{k+1} = g(h\mathbf{A})\boldsymbol{\eta}_k$$

mit einer rationalen Funktion g angeben, die nur vom betrachteten Verfahren abhängt. Man erhält zum Beispiel:

- EULER-Verfahren:

$$\boldsymbol{\eta}_{k+1} = (\mathbf{I} + h\mathbf{A})\boldsymbol{\eta}_k, \quad g(z) = 1 + z,$$

- implizites EULER-Verfahren:

$$\boldsymbol{\eta}_{k+1} = (\mathbf{I} - h\mathbf{A})^{-1}\boldsymbol{\eta}_k, \quad g(z) = \frac{1}{1-z},$$

- HEUN-Verfahren:

$$\boldsymbol{\eta}_{k+1} = \left(\mathbf{I} + h\mathbf{A} + \frac{1}{2}(h\mathbf{A})^2 \right) \boldsymbol{\eta}_k, \quad g(z) = 1 + z + \frac{z^2}{2},$$

- RUNGE-KUTTA-Verfahren:

$$\boldsymbol{\eta}_{k+1} = \left(\mathbf{I} + h\mathbf{A} + \frac{1}{2}(h\mathbf{A})^2 + \frac{1}{6}(h\mathbf{A})^3 + \frac{1}{24}(h\mathbf{A})^4 \right) \boldsymbol{\eta}_k,$$

$$g(z) = 1 + z + \frac{z^2}{2} + \frac{z^3}{6} + \frac{z^4}{24}.$$

Genügen die Näherungslösungen eines ESV für das angegebene AWP einer derartigen Rekursionsformel, so gilt

$$\boldsymbol{\eta}_k = [g(h\mathbf{A})]^k \mathbf{y}_0$$

bzw.

$$\boldsymbol{\eta}(x; h) = [g(h\mathbf{A})]^{\frac{x-x_0}{h}} \mathbf{y}_0.$$

Falls die Eigenwerte der Matrix $g(h\mathbf{A})$ betragsmäßig größer als 1 sind, wachsen die Näherungen $\boldsymbol{\eta}_k$ unbeschränkt. Näherungslösungen, die für große x gegen Null konvergieren, ergeben sich nur für $|\lambda(g(h\mathbf{A}))| < 1$. Ist λ Eigenwert von \mathbf{A} , so ist $g(h\lambda)$ Eigenwert von $g(h\mathbf{A})$. Damit spielen gerade die Argumente z von g eine besondere Rolle, für die $|g(z)| < 1$ gilt;. Das lineare AWP

$$\mathbf{y}' = \mathbf{A}\mathbf{y}, \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

werde mit einem ESV gelöst. Für die Näherungslösungen existiere bei konstanter Schrittweite eine Rekursionsbeziehung der Form

$$\left. \begin{aligned} \boldsymbol{\eta}_0 &= \mathbf{y}_0, \\ \boldsymbol{\eta}_{k+1} &= g(h\mathbf{A})\boldsymbol{\eta}_k, \\ x_{k+1} &= x_k + h \end{aligned} \right\}, \quad k = 0, 1, \dots$$

Die Menge

$$M = \{ z \in \mathbb{C} \mid |g(z)| < 1 \}$$

heißt dann (**absolute**) **Stabilitätsgebiet** des ESV.

Steife Differentialgleichungen sind gerade dadurch gekennzeichnet, dass alle Eigenwerte der Matrix \mathbf{A} einen negativen Realteil haben. Ein ESV wird darum genau dann für jene Schrittweiten $h > 0$ vernünftige Näherungslösungen einer steifen Differentialgleichung liefern, für die $\lambda h \in M$ für alle Eigenwerte λ von \mathbf{A} gilt. Ist $h > 0$, so folgt aus

$$\lambda \in \mathbb{C} = \{ z \in \mathbb{C} \mid \Re(z) < 0 \}$$

$h\lambda \in \mathbb{C}$. Ein ESV ist darum um so besser zur Integration einer steifen Differentialgleichung geeignet, je größer der Durchschnitt zwischen absolutem Stabilitätsgebiet M und der linken komplexen Halbebene \mathbb{C} ist. Ein ESV zum Lösen eines AWP heißt **absolut stabil** bzw. **A-stabil**, falls $\mathbb{C} \subset M$ gilt. Das implizite EULER-Verfahren ist damit absolut stabil. Weitere absolut stabile Verfahren höherer Ordnung gewinnt man durch ähnliche Ansätze wie bei den RUNGE-KUTTA-Verfahren:

$$\Phi(x, \mathbf{y}; h) = \sum_{k=0}^n \alpha_k \mathbf{K}_k(x, \mathbf{y}; h)$$

mit

$$\mathbf{K}_k(x, \mathbf{y}; h) = \mathbf{f} \left(x + \vartheta_k h, \mathbf{y} + h \sum_{l=0}^n \beta_{kl} \mathbf{K}_l(x, \mathbf{y}; h) \right), \quad k = 0, \dots, n.$$

Die Größen \mathbf{K}_k ergeben sich hier als Lösungen von i.a. nichtlinearen Gleichungssystemen. Für den Fall $n = 1$ ergibt sich zum Beispiel

$$\begin{aligned} \Phi(x, \mathbf{y}; h) &= \alpha_0 \mathbf{K}_0 + \alpha_1 \mathbf{K}_1, \\ \mathbf{K}_0 &= \mathbf{f}(x + \vartheta_0 h, \mathbf{y} + \beta_{00} \mathbf{K}_0 h + \beta_{01} \mathbf{K}_1 h), \\ \mathbf{K}_1 &= \mathbf{f}(x + \vartheta_1 h, \mathbf{y} + \beta_{10} \mathbf{K}_0 h + \beta_{11} \mathbf{K}_1 h). \end{aligned}$$

Die Konstanten $\alpha_0, \dots, \alpha_n, \vartheta_0, \dots, \vartheta_n$ und $\beta_{00}, \dots, \beta_{nn}$ werden wieder so bestimmt, dass man Verfahren möglichst hoher Ordnung erhält. Wir wollen die Koeffizienten einiger Verfahren nach folgendem Schema angeben:

ϑ_0	β_{00}	β_{01}	β_{02}	β_{03}
ϑ_1	β_{10}	β_{11}	β_{12}	β_{13}
ϑ_2	β_{20}	β_{21}	β_{22}	β_{23}
ϑ_3	β_{30}	β_{31}	β_{32}	β_{33}
	α_0	α_1	α_2	α_3

Speziell sind die expliziten RUNGE-KUTTA-Formeln in diesem Schema enthalten ($\beta_{ij} = 0$ für $j \geq i$). Gilt $\beta_{ij} = 0$ für $j > i$, so heißen die Verfahren **diagonalimplizit**.

- GAUSS-LEGENDRE

$$\begin{array}{c|c} \frac{1}{2} & \frac{1}{2} \\ \hline & 1 \end{array} \quad p = 2$$

$$\begin{array}{c|cc} \frac{1}{2} - \alpha & \frac{1}{4} & \frac{1}{4} - \alpha \\ \frac{1}{2} + \alpha & \frac{1}{4} + \alpha & \frac{1}{4} \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array} \quad \alpha = \frac{1}{2\sqrt{3}} \quad p = 4$$

- RADAU IA

$$\begin{array}{c|c} 0 & 0 \\ \hline & 1 \end{array} \quad p = 1$$

$$\begin{array}{c|cc} 0 & \frac{1}{4} & -\frac{1}{4} \\ \frac{2}{3} & \frac{1}{4} & \frac{5}{12} \\ \hline & -\frac{1}{4} & \frac{3}{4} \end{array} \quad p = 3$$

- RADAU IIA

$$\begin{array}{c|c} 1 & 1 \\ \hline & 1 \end{array} \quad p = 1$$

$$\begin{array}{c|cc} \frac{1}{3} & \frac{5}{12} & -\frac{1}{12} \\ 1 & \frac{3}{4} & \frac{1}{4} \\ \hline & \frac{3}{4} & \frac{1}{4} \end{array} \quad p = 3$$

- LOBATTO IIIA

$$\begin{array}{c|cc} 0 & 0 & 0 \\ 1 & \frac{1}{2} & \frac{1}{2} \\ \hline & \frac{1}{2} & \frac{1}{2} \end{array} \quad p = 2$$

$$\begin{array}{c|ccc} 0 & 0 & 0 & 0 \\ \frac{1}{2} & \frac{5}{24} & \frac{1}{3} & -\frac{1}{24} \\ 1 & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} \\ \hline & \frac{1}{6} & \frac{2}{3} & \frac{1}{6} \end{array} \quad p = 4$$

Schrittweisensteuerungen sind analog zu expliziten Verfahren möglich. Es existieren auch Verfahren, die ähnlich wie die RUNGE-KUTTA-FEHLBERG-Verfahren funktionieren.

6.3. Mehrschrittverfahren

6.3.1. Prediktor-Korrektor-Verfahren

ESV waren dadurch charakterisiert, dass zum Berechnen einer Näherung η_{k+1} nur die Kenntnis der vorigen Näherung η_k benötigt wurde. Nun liegt es nahe, Verfahren zu konstruieren, die zum Berechnen der Näherung η_{k+1} mehrere vorhergehende Näherungen $\eta_k, \eta_{k-1}, \dots, \eta_{k-r+1}$ verwenden. Verfahren dieser Art heißen Mehrschrittverfahren.⁴ Wir wollen nun einfache MSV konstruieren. Es sei

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

das zu lösende AWP. Die Funktion \mathbf{f} erfülle die Voraussetzungen aus Satz 6.1. Integrieren wir die Differentialgleichung von \underline{x} bis \bar{x} , so erhalten wir

$$\int_{\underline{x}}^{\bar{x}} \mathbf{y}'(t) dt = \int_{\underline{x}}^{\bar{x}} \mathbf{f}(t, \mathbf{y}(t)) dt$$

und damit

$$\mathbf{y}(\bar{x}) - \mathbf{y}(\underline{x}) = \int_{\underline{x}}^{\bar{x}} \mathbf{f}(t, \mathbf{y}(t)) dt.$$

Kennt man einen Näherungswert für $\mathbf{y}(\underline{x})$, lässt sich eine Näherung für $\mathbf{y}(\bar{x})$ gewinnen, indem man das Integral auf der rechten Seite der obigen Gleichung in geeigneter Weise approximiert. Hier treten zwei Schwierigkeiten auf. Einerseits kennt man die Abhängigkeit des Integranden \mathbf{f} von t nicht, da die Funktion $\mathbf{y}(t)$ nicht bekannt ist. Andererseits ließe sich das Integral auch bei bekannter Abhängigkeit der Funktion \mathbf{f} von t in den meisten Fällen nur näherungsweise berechnen. Betrachten wir das Problem der Approximation des bestimmten Integrals bei bekanntem Integranden. Zum Lösen dieses Problems gehen wir ähnlich wie bei der Herleitung der NEWTON-COTES-Formeln vor. Wir ersetzen den Integranden durch ein geeignetes Interpolationspolynom und integrieren dieses. Zur Konstruktion des Interpolationspolynoms wählen wir äquidistante Stützstellen

$$x_i = x_0 + ih, \quad i = 0, 1, \dots$$

⁴Wir werden für Mehrschrittverfahren die Abkürzung MSV verwenden.

Weiterhin sei

$$\bar{x} = x_{p+k}, \quad \underline{x} = x_{p-j}.$$

Das gesuchte Interpolationspolynom $\mathbf{P}_q(t)$ sei durch

1. $\text{Grad} \mathbf{P}_q \leq q$ und
2. $\mathbf{P}_q(x_i) = \mathbf{f}(x_i, \mathbf{y}(x_i))$ für $i = p, p-1, \dots, p-q$

festgelegt. Man beachte, dass \mathbf{P}_q im Falle $\mathbf{y} : \mathbb{R} \rightarrow \mathbb{R}^d$ selbst eine Vektorfunktion ist, deren Komponenten Polynome vom Höchstgrad q sind. Mit der LAGRAN-GESEN Interpolationsformel ist \mathbf{P}_q sofort angebar. Wir erhalten

$$\mathbf{P}_q(t) = \sum_{i=0}^q \mathbf{f}(x_{p-i}, \mathbf{y}(x_{p-i})) L_i(t)$$

mit

$$L_i(t) = \prod_{\substack{l=0 \\ l \neq i}}^q \frac{t - x_{p-l}}{x_{p-i} - x_{p-l}}, \quad i = 0, 1, \dots, q.$$

Damit gilt

$$\begin{aligned} \mathbf{y}_{p+k} &\approx \mathbf{y}_{p-j} + \sum_{i=0}^q \mathbf{f}(x_{p-i}, \mathbf{y}(x_{p-i})) \int_{x_{p-j}}^{x_{p+k}} L_i(t) dt \\ &= \mathbf{y}_{p-j} + h \sum_{i=0}^q \beta_{qi} \mathbf{f}(x_{p-i}, \mathbf{y}(x_{p-i})) \end{aligned}$$

mit

$$\begin{aligned} \beta_{qi} &= \frac{1}{h} \int_{x_{p-j}}^{x_{p+k}} L_i(t) dt \\ &= \frac{1}{h} \int_{x_{p-j}}^{x_{p+k}} \prod_{\substack{l=0 \\ l \neq i}}^q \frac{t - x_{p-l}}{x_{p-i} - x_{p-l}} dt. \end{aligned}$$

Durch die Substitution $t = x_p + sh = x_0 + (p+s)h$ erhält man daraus

$$\beta_{qi} = \int_{-j}^k \prod_{\substack{l=0 \\ l \neq i}}^q \frac{s+l}{-j+l} ds.$$

Je nach Wahl von k, j und q ergeben sich verschiedene Verfahren der Form

$$\eta_{p+k} = \eta_{p-j} + h \sum_{i=0}^q \beta_{qi} \mathbf{f}(x_{p-i}, \eta_{p-i}).$$

Wir wollen die wichtigsten vier Verfahrensklassen angeben.

Adams-Bashforth-Verfahren

Wir setzen $k = 1, j = 0$ und $q = 0, 1, \dots$. Es gilt dann

$$\beta_{qi} = \int_0^1 \prod_{\substack{l=0 \\ l \neq i}}^q \frac{s+l}{-j+l} ds$$

und

$$\eta_{p+1} = \eta_p + h [\beta_{q0} \mathbf{f}_p + \beta_{q1} \mathbf{f}_{p-1} + \dots + \beta_{qq} \mathbf{f}_{p-q}]$$

mit

$$\mathbf{f}_i = \mathbf{f}(x_i, \eta_i), \quad i = 0, 1, \dots$$

q	$s\beta_{qi}$					s	Ordnung
0	1					1	1
1	3	-1				2	2
2	23	-16	5			12	3
3	55	-59	37	-9		24	4
4	1901	-2774	2616	-1274	251	720	5

Tabelle 6.1: ADAMS-BASHFORTH-Verfahren

Adams-Moulton-Verfahren

Wir setzen $k = 0, j = 1$ und $q = 0, 1, \dots$. Es gilt dann

$$\beta_{qi} = \int_{-1}^0 \prod_{\substack{l=0 \\ l \neq i}}^q \frac{s+l}{-j+l} ds$$

und

$$\boldsymbol{\eta}_p = \boldsymbol{\eta}_{p-1} + h [\beta_{q0} \mathbf{f}(x_p, \boldsymbol{\eta}_p) + \beta_{q1} \mathbf{f}_{p-1} + \cdots + \beta_{qq} \mathbf{f}_{p-q}].$$

Ersetzt man p durch $p+1$, so erhält man die übliche Form

$$\boldsymbol{\eta}_{p+1} = \boldsymbol{\eta}_p + h [\beta_{q0} \mathbf{f}(x_{p+1}, \boldsymbol{\eta}_{p+1}) + \beta_{q1} \mathbf{f}_p + \cdots + \beta_{qq} \mathbf{f}_{p-q+1}].$$

q	$s\beta_{qi}$					s	Ordnung
0	1					1	1
1	1	1				2	2
2	5	8	-1			12	3
3	9	19	-5	1		24	4
4	251	646	-264	106	19	720	5

Tabelle 6.2: ADAMS-MOULTON-Verfahren

Nyström-Verfahren

Wir setzen $k = 1$, $j = 1$ und $q = 0, 1, \dots$. Es gilt dann

$$\beta_{qi} = \int_{-1}^1 \prod_{\substack{l=0 \\ l \neq i}}^q \frac{s+l}{-j+l} ds$$

und

$$\boldsymbol{\eta}_{p+1} = \boldsymbol{\eta}_p + h [\beta_{q0} \mathbf{f}_p + \beta_{q1} \mathbf{f}_{p-1} + \cdots + \beta_{qq} \mathbf{f}_{p-q}].$$

Milne-Verfahren

Wir setzen $k = 0$, $j = 2$ und $q = 0, 1, \dots$. Es gilt dann

$$\beta_{qi} = \int_{-2}^0 \prod_{\substack{l=0 \\ l \neq i}}^q \frac{s+l}{-j+l} ds$$

q	$s\beta_{qi}$					s	Ordnung
0	2					1	2
1	2	0				1	2
2	7	-2	1			3	3
3	8	-5	4	-1		3	4
4	269	-266	294	-146	29	90	5

Tabelle 6.3: NYSTRÖM-Verfahren

und

$$\boldsymbol{\eta}_p = \boldsymbol{\eta}_{p-2} + h [\beta_{q0}\mathbf{f}(x_p, \boldsymbol{\eta}_p) + \beta_{q1}\mathbf{f}_{p-1} + \cdots + \beta_{qq}\mathbf{f}_{p-q}].$$

Ersetzen wir wieder p durch $p+1$, so erhalten wir die übliche Form

$$\boldsymbol{\eta}_{p+1} = \boldsymbol{\eta}_{p-1} + h [\beta_{q0}\mathbf{f}(x_{p+1}, \boldsymbol{\eta}_{p+1}) + \beta_{q1}\mathbf{f}_p + \cdots + \beta_{qq}\mathbf{f}_{p-q+1}].$$

Bei den Verfahren von ADAMS-BASHFORTH und NYSTRÖM lässt sich aus den be-

q	$s\beta_{qi}$					s	Ordnung
0	2					1	1
1	0	2				1	2
2	1	4	1			3	4
3	1	4	1	0		3	4
4	29	124	24	4	-1	90	5

Tabelle 6.4: MILNE-Verfahren

kannten Näherungen $\boldsymbol{\eta}_p, \dots, \boldsymbol{\eta}_{p-q}$ leicht die neue Näherung $\boldsymbol{\eta}_{p+1}$ bestimmen. Dazu ist nur eine Auswertung der Funktion \mathbf{f} notwendig, nämlich das Berechnen von $\mathbf{f}(x_p, \boldsymbol{\eta}_p)$. Bei den Verfahren von ADAMS-MOULTON und MILNE tritt die neue Näherung $\boldsymbol{\eta}_{p+1}$ jedoch auch auf der rechten Seite der Gleichung auf. Damit ergibt sich $\boldsymbol{\eta}_{p+1}$ als Lösung eines i.a. nichtlinearen Gleichungssystems. Die Verfahren von ADAMS-MOULTON und MILNE sind **implizite** Verfahren im Gegensatz zu den **expliziten** Verfahren von ADAMS-BASHFORTH und NYSTRÖM. Zum Lösen der nichtlinearen Gleichungen liegt es nahe, eine Fixpunktiteration folgender Form anzuwenden:

6.12. Fixpunktiteration für Korrektor-Verfahren:

S0 Wähle Startwert $\boldsymbol{\eta}_{p+1}^{(0)}$ und setze $\nu = 0$. Berechne für das

- ADAMS-MOULTON-Verfahren

$$\varrho = \boldsymbol{\eta}_p + h [\beta_{q1} \mathbf{f}_p + \beta_{q2} \mathbf{f}_{p-1} + \cdots + \beta_{qq} \mathbf{f}_{p-q}]$$

- MILNE-Verfahren

$$\varrho = \boldsymbol{\eta}_{p-1} + h [\beta_{q1} \mathbf{f}_p + \beta_{q2} \mathbf{f}_{p-1} + \cdots + \beta_{qq} \mathbf{f}_{p-q}]$$

S1 Berechne

$$\boldsymbol{\eta}_{p+1}^{(\nu+1)} = \varrho + h \beta_{q0} \mathbf{f}(x_{p+1}, \boldsymbol{\eta}_{p+1}^{(\nu)}).$$

S2 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Die Iteration konvergiert, falls die Funktion \mathbf{f} eine LIPSCHITZ-Bedingung

$$\|\mathbf{f}(x, \mathbf{y}_1) - \mathbf{f}(x, \mathbf{y}_2)\| \leq \|\mathbf{y}_1 - \mathbf{y}_2\|$$

für alle $x \in [a, b]$ und alle $\mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^d$ erfüllt, und falls die Schrittweite h hinreichend klein gewählt wird. Eine gute Startnäherung $\boldsymbol{\eta}_{p+1}^{(0)}$ beschafft man sich mit einem expliziten MSV. Darum heißen die expliziten MSV auch **Prediktor-Verfahren** und die impliziten MSV **Korrektor-Verfahren**. Die Kombination eines Prediktor-Verfahrens mit einem Korrektor-Verfahren wird **Prediktor-Korrektor-Verfahren** genannt. Für das Berechnen einer neuen Näherung reichen i. a. ein Prediktor-Schritt und ein oder zwei Korrektor-Schritte.

6.3.2. Konvergenz von Mehrschrittverfahren

Die bisher behandelten MSV und ESV sind von der Form

$$\boldsymbol{\eta}_{j+r} + \alpha_{r-1} \boldsymbol{\eta}_{j+r-1} + \cdots + \alpha_0 \boldsymbol{\eta}_j = h \mathbf{F}(x_j, \boldsymbol{\eta}_{j+r}, \dots, \boldsymbol{\eta}_j; h; \mathbf{f}).$$

Ein derartiges Verfahren wird auch als r -**Schritt-Verfahren** bezeichnet. Bei den ESV ist $r = 1$, $\alpha_0 = -1$ und $\mathbf{F} \equiv \Phi$. Bei den im vorigen Abschnitt behandelten MSV hing die Funktion \mathbf{F} linear von der Funktion \mathbf{f} ab:

$$\begin{aligned} \mathbf{F}(x_j, \boldsymbol{\eta}_{j+r}, \dots, \boldsymbol{\eta}_j; h; \mathbf{f}) &= \beta_r \mathbf{f}(x_{j+r}, \boldsymbol{\eta}_{j+r}) \\ &+ \beta_{r-1} \mathbf{f}(x_{j+r-1}, \boldsymbol{\eta}_{j+r-1}) + \cdots + \beta_0 \mathbf{f}(x_j, \boldsymbol{\eta}_j). \end{aligned}$$

Man spricht von einem linearen r -Schritt-Verfahren. Im Gegensatz dazu hing die Funktion Φ bei den ESV hochgradig nichtlinear von \mathbf{f} ab.

Wir wollen nun die Eigenschaften von allgemeinen MSV untersuchen. Dazu werden wir wieder lokale und globale Diskretisierungsfehler betrachten und Konsistenz und Konvergenz der Verfahren untersuchen.

Es seien $\mathbf{f} \in F^1[a, b]$, $\mathbf{z}(t)$ die exakte Lösung des AWP

$$\mathbf{z}'(t) = \mathbf{f}(t, \mathbf{z}(t)), \quad \mathbf{z}(x) = \mathbf{y}$$

und \mathbf{F} die Verfahrensfunktion eines zugeordneten r -Schritt-Verfahrens. Dann heißt die Größe

$$\tau(x, \mathbf{y}; h) = \frac{1}{h} \left[\mathbf{z}(x + rh) + \sum_{i=0}^{r-1} \alpha_i \mathbf{z}(x + ih) - h \mathbf{F}(x_i, \mathbf{z}(x + rh), \dots, \mathbf{z}(x + h), \mathbf{z}(x); h; \mathbf{f}) \right]$$

lokaler Diskretisierungsfehler des entsprechenden MSV an der Stelle (x, \mathbf{y}) .

Der lokale Diskretisierungsfehler gibt folglich an, wie gut die exakte Lösung die Gleichung des MSV erfüllt. Wie bei ESV definieren wir nun auch die Konsistenz von Verfahren. Ein MSV heißt **konsistent**, falls für alle $x \in [a, b]$, alle $\mathbf{y} \in \mathbb{R}^d$ und alle $\mathbf{f} \in F^1[a, b]$

$$\lim_{h \rightarrow 0} \tau(x, \mathbf{y}; h) = 0$$

gilt. Es ist natürlich wieder wünschenswert, dass der lokale Diskretisierungsfehler möglichst schnell mit h gegen Null konvergiert. Das führt wieder auf den Begriff der Konsistenzordnung. Ein MSV hat die Konsistenzordnung p , falls für alle $x \in [a, b]$, alle $\mathbf{y} \in \mathbb{R}^d$ und alle $\mathbf{f} \in F^p[a, b]$

$$\|\tau(x, \mathbf{y}; h)\| = O(|h|^p)$$

gilt. Die Konsistenzordnung eines Verfahrens lässt sich wieder bestimmen, indem man für den lokalen Diskretisierungsfehler die TAYLOR-Entwicklung bezüglich h verwendet. Die Ordnung des ersten nichtverschwindenden Gliedes gibt dann die Konsistenzordnung an.

6.13. Beispiel: Wir betrachten ein NYSTRÖM-Verfahren mit $q = 2$:

$$\boldsymbol{\eta}_{p+1} - \boldsymbol{\eta}_{p-1} = h \left[\frac{7}{3} \mathbf{f}(x_p, \boldsymbol{\eta}_p) - \frac{2}{3} \mathbf{f}(x_{p-1}, \boldsymbol{\eta}_{p-1}) + \frac{1}{3} \mathbf{f}(x_{p-2}, \boldsymbol{\eta}_{p-2}) \right].$$

Setzt man $p = j + 2$, so erhält man

$$\boldsymbol{\eta}_{j+3} - \boldsymbol{\eta}_{j+1} = h \left[\frac{7}{3} \mathbf{f}(x_{j+2}, \boldsymbol{\eta}_{j+2}) - \frac{2}{3} \mathbf{f}(x_{j+1}, \boldsymbol{\eta}_{j+1}) + \frac{1}{3} \mathbf{f}(x_j, \boldsymbol{\eta}_j) \right].$$

Damit gilt

$$\begin{aligned}\tau(x, \mathbf{z}; h) &= \frac{1}{h} [\mathbf{z}(x+3h) - \mathbf{z}(x+h)] - \frac{7}{3} \mathbf{f}(x+2h, \mathbf{z}(x+2h)) + \\ &\quad + \frac{2}{3} \mathbf{f}(x+h, \mathbf{z}(x+h)) - \frac{1}{3} \mathbf{f}(x, \mathbf{z}(x)) \\ &= \frac{1}{h} [\mathbf{z}(x+3h) - \mathbf{z}(x+h)] - \frac{7}{3} \mathbf{z}'(x+2h) + \frac{2}{3} \mathbf{z}'(x+h) - \frac{1}{3} \mathbf{z}'(x).\end{aligned}$$

Wir entwickeln die einzelnen Terme nach Potenzen von h :

$$\begin{aligned}\mathbf{z}(x+3h) &= \mathbf{z}(x) + 3h\mathbf{z}'(x) + \frac{9h^2}{2}\mathbf{z}''(x) + \frac{9h^3}{2}\mathbf{z}'''(x) + \frac{27h^4}{8}\mathbf{z}^{(4)}(x) + \mathbf{O}(h^5), \\ \mathbf{z}(x+h) &= \mathbf{z}(x) + h\mathbf{z}'(x) + \frac{h^2}{2}\mathbf{z}''(x) + \frac{h^3}{6}\mathbf{z}'''(x) + \frac{h^4}{24}\mathbf{z}^{(4)}(x) + \mathbf{O}(h^5), \\ \mathbf{z}'(x+2h) &= \mathbf{z}'(x) + 2h\mathbf{z}''(x) + 2h^2\mathbf{z}'''(x) + \frac{4h^3}{3}\mathbf{z}^{(4)}(x) + \mathbf{O}(h^4), \\ \mathbf{z}'(x+h) &= \mathbf{z}'(x) + h\mathbf{z}''(x) + \frac{h^2}{2}\mathbf{z}'''(x) + \frac{h^3}{6}\mathbf{z}^{(4)}(x) + \mathbf{O}(h^4).\end{aligned}$$

Setzt man diese Entwicklungen in die obige Darstellung von τ ein, so ergibt sich

$$\begin{aligned}\tau(x, \mathbf{z}; h) &= \frac{1}{h} \left[\mathbf{z}(x) + 3h\mathbf{z}'(x) + \frac{9h^2}{2}\mathbf{z}''(x) + \frac{9h^3}{2}\mathbf{z}'''(x) + \frac{27h^4}{8}\mathbf{z}^{(4)}(x) - \right. \\ &\quad \left. - \mathbf{z}(x) - h\mathbf{z}'(x) - \frac{h^2}{2}\mathbf{z}''(x) - \frac{h^3}{6}\mathbf{z}'''(x) - \frac{h^4}{24}\mathbf{z}^{(4)}(x) + \mathbf{O}(h^5) \right] - \\ &\quad - \frac{7}{3} \left[\mathbf{z}'(x) + 2h\mathbf{z}''(x) + 2h^2\mathbf{z}'''(x) + \frac{4h^3}{3}\mathbf{z}^{(4)}(x) + \mathbf{O}(h^4) \right] + \\ &\quad + \frac{2}{3} \left[\mathbf{z}'(x) + h\mathbf{z}''(x) + \frac{h^2}{2}\mathbf{z}'''(x) + \frac{h^3}{6}\mathbf{z}^{(4)}(x) + \mathbf{O}(h^4) \right] - \frac{1}{3} \mathbf{z}'(x) \\ &= \frac{1}{h} \left[2h\mathbf{z}'(x) + 4h^2\mathbf{z}''(x) + \frac{13h^3}{3}\mathbf{z}'''(x) + \frac{10h^4}{3}\mathbf{z}^{(4)}(x) + \mathbf{O}(h^5) \right] - \\ &\quad - 2\mathbf{z}'(x) - 4h\mathbf{z}''(x) - \frac{13h^2}{3}\mathbf{z}'''(x) - 3h^3\mathbf{z}^{(4)}(x) + \mathbf{O}(h^4) \\ &= \frac{h^3}{3}\mathbf{z}^{(4)}(x) + \mathbf{O}(h^4).\end{aligned}$$

Es handelt sich somit um ein Verfahren der Konsistenzordnung 3. ♡

Nun möchte man natürlich wieder Verfahren möglichst hoher Konsistenzordnung konstruieren. Für lineare MSV erscheint dies einfach zu sein. Der Ansatz

$$\boldsymbol{\eta}_{j+r} + \alpha_{r-1}\boldsymbol{\eta}_{j+r-1} + \cdots + \alpha_0\boldsymbol{\eta}_j = h [\beta_r \mathbf{f}(x_{j+r}, \boldsymbol{\eta}_{j+r}) + \cdots + \beta_0 \mathbf{f}(x_j, \boldsymbol{\eta}_j)]$$

führt über die entsprechende Entwicklung des lokalen Diskretisierungsfehlers auf ein lineares Gleichungssystem zur Bestimmung der Koeffizienten $\alpha_{r-1}, \dots, \alpha_0$ und β_r, \dots, β_0 . Bei der Konstruktion von ESV über einen entsprechenden Ansatz lieferte jede Lösung des dazugehörigen nichtlinearen Gleichungssystems ein konvergentes ESV. Bei MSV ist dem nicht so. Das zeigt schon das folgende Beispiel.

6.14. Beispiel: Wir wollen ein explizites lineares MSV mit $r = 2$ und maximaler Konsistenzordnung konstruieren. Dazu dient der folgende Ansatz:

$$\eta_{j+2} + \alpha_1 \eta_{j+1} + \alpha_0 \eta_j = h [\beta_1 \mathbf{f}(x_{j+1}, \eta_{j+1}) + \beta_0 \mathbf{f}(x_j, \eta_j)].$$

Für den lokalen Diskretisierungsfehler erhalten wir

$$\begin{aligned} \tau(x, \mathbf{y}; h) &= \frac{1}{h} [z(x+2h) + \alpha_1 z(x+h) + \alpha_0 z(x)] \\ &\quad - [\beta_1 \mathbf{f}(x+h, z(x+h)) + \beta_0 \mathbf{f}(x, z(x))]. \end{aligned}$$

Mit $\mathbf{z}'(t) = \mathbf{f}(t, z(t))$ und $z(x) = \mathbf{y}$ ergibt sich daraus

$$\tau(x, \mathbf{y}; h) = \frac{1}{h} [z(x+2h) + \alpha_1 z(x+h) + \alpha_0 z(x)] - [\beta_1 \mathbf{z}'(x+h) + \beta_0 \mathbf{z}'(x)].$$

Wir entwickeln die einzelnen Terme nach Potenzen von h :

$$\begin{aligned} z(x+2h) &= z(x) + 2hz'(x) + \frac{4h^2}{2}z''(x) + \frac{8h^3}{6}z'''(x) + \mathbf{O}(h^4), \\ z(x+h) &= z(x) + hz'(x) + \frac{h^2}{2}z''(x) + \frac{h^3}{6}z'''(x) + \mathbf{O}(h^4), \\ z'(x+h) &= z'(x) + hz''(x) + \frac{h^2}{2}z'''(x) + \mathbf{O}(h^3). \end{aligned}$$

Einsetzen dieser Entwicklungen liefert

$$\begin{aligned} \tau(x, \mathbf{y}; h) &= \frac{1}{h} \left[(1 + \alpha_1 + \alpha_0)z(x) + (2 + \alpha_1)hz'(x) + (4 + \alpha_1)\frac{h^2}{2}z''(x) + \right. \\ &\quad \left. + (8 + \alpha_1)\frac{h^3}{6}z'''(x) + \mathbf{O}(h^4) \right] - \\ &\quad - \left[(\beta_1 + \beta_0)z'(x) + \beta_1 hz''(x) + \beta_1 \frac{h^2}{2}z'''(x) + \mathbf{O}(h^3) \right] \\ &= (1 + \alpha_1 + \alpha_0)\frac{1}{h}z(x) + (2 + \alpha_1 - \beta_1 - \beta_0)z'(x) \\ &\quad + (4 + \alpha_1 - 2\beta_1)\frac{h}{2}z''(x) + (8 + \alpha_1 - 3\beta_1)\frac{h^2}{6}z'''(x) + \mathbf{O}(h^3). \end{aligned}$$

Um ein Verfahren möglichst hoher Konsistenzordnung zu erhalten, müssen die ersten Terme verschwinden. Es ergibt sich das folgende lineare Gleichungssystem.

$$\begin{array}{rcl} \alpha_1 + \alpha_0 & & = -1 \\ \alpha_1 & - \beta_1 - \beta_0 & = -2 \\ \alpha_1 & - 2\beta_1 & = -4 \\ \alpha_1 & - 3\beta_1 & = -8 \end{array}$$

Dieses besitzt die eindeutige Lösung $\alpha_0 = -5$, $\alpha_1 = 4$, $\beta_0 = 2$ und $\beta_1 = 4$. Wir erhalten folglich das Verfahren

$$\eta_{j+2} + 4\eta_{j+1} - 5\eta_j = h [4f(x_{j+1}, \eta_{j+1}) + 2f(x_j, \eta_j)].$$

mit der Konsistenzordnung 3. Wenden wir dieses Verfahren auf das AWP

$$y' = -y, \quad y(0) = 1$$

mit der exakten Lösung $y(x) = e^{-x}$ an, so lässt sich die Näherungslösung explizit darstellen. Sie erfüllt die homogenen Differenzengleichung

$$\eta_{j+2} + (4 + 4h)\eta_{j+1} + (-5 + 2h)\eta_j = 0.$$

Lösungen von homogenen Differenzengleichungen erhält man über den Ansatz

$$\eta_j = \lambda^j.$$

Einsetzen dieses Ansatzes in die Differenzengleichung liefert

$$\lambda^{j+2} + (4 + 4h)\lambda^{j+1} + (-5 + 2h)\lambda^j = 0.$$

Diese Gleichung ist für $\lambda = 0$ trivial erfüllt. Wesentliche Lösungen erhält man aus

$$\lambda^2 + (4 + 4h)\lambda - 5 + 2h = 0.$$

Das sind

$$\lambda_1(h) = -2 - 2h + \sqrt{1 + \frac{2}{3}h + \frac{4}{9}h^2}$$

und

$$\lambda_2(h) = -2 - 2h - \sqrt{1 + \frac{2}{3}h + \frac{4}{9}h^2}.$$

Die allgemeine Lösung der Differenzgleichung lässt sich als Linearkombination dieser speziellen Lösungen darstellen

$$\eta_j = \mu_1 \lambda_1^j + \mu_2 \lambda_2^j.$$

μ_1 und μ_2 werden aus den Anfangsbedingungen ermittelt. Es gelte

$$\begin{aligned} \eta_0 &= y_0 = 1, \\ \eta_1 &= y_1 = e^{-h} \quad (\text{exakter Wert}). \end{aligned}$$

Daraus folgt

$$\mu_1 = \frac{e^{-h} - \lambda_2}{\lambda_1 - \lambda_2}, \quad \mu_2 = -\frac{e^{-h} - \lambda_1}{\lambda_1 - \lambda_2}.$$

Entwickelt man nun $\lambda_1(h)$, $\lambda_2(h)$, $\mu_1(h)$ und $\mu_2(h)$ nach Potenzen von h und setzt alles in die Darstellung

$$\eta_j = \mu_1(h) [\lambda_1(h)]^j + \mu_2(h) [\lambda_2(h)]^j$$

ein, so erhält man für festes x und $h_n = x/n$

$$\begin{aligned} \eta_n &= \eta(x; h_n) \\ &= \left[1 + O\left(\frac{x}{n}\right)\right] \left[1 - \frac{x}{n} + O\left(\left(\frac{x}{n}\right)^2\right)\right]^n - \\ &\quad - \frac{1}{216} \left(\frac{x}{n}\right)^4 \left[1 + O\left(\frac{x}{n}\right)\right] \left[-5 - 3\frac{x}{n} + O\left(\left(\frac{x}{n}\right)^2\right)\right]^n. \end{aligned}$$

Für den ersten Term gilt offensichtlich

$$\lim_{n \rightarrow \infty} \left[1 + O\left(\frac{x}{n}\right)\right] \left[1 - \frac{x}{n} + O\left(\left(\frac{x}{n}\right)^2\right)\right]^n = e^{-x}.$$

Das ist der Anteil, der die exakte Lösung beschreibt. Für den zweiten Term gilt

$$\begin{aligned} &-\frac{1}{216} \left(\frac{x}{n}\right)^4 \left[1 + O\left(\frac{x}{n}\right)\right] \left[-5 - 3\frac{x}{n} + O\left(\left(\frac{x}{n}\right)^2\right)\right]^n = \\ &-\frac{(-5)^n}{216} \left(\frac{x}{n}\right)^4 \left[1 + O\left(\frac{x}{n}\right)\right] \left[1 + \frac{3x}{5n} + O\left(\left(\frac{x}{n}\right)^2\right)\right]^n. \end{aligned}$$

Wegen

$$\lim_{n \rightarrow \infty} \left[1 + \frac{3x}{5n} + O\left(\left(\frac{x}{n}\right)^2\right) \right]^n = e^{3x/5}$$

verhält sich der zweite Term für großes n wie

$$-\frac{(-5)^n}{216} \left(\frac{x}{n}\right)^4 e^{3x/5}.$$

Dieser Term oszilliert mit wachsendem n immer heftiger.

Genau dieses Verhalten beobachtet man auch bei Anwendung des obigen Verfahrens auf ein beliebiges AWP. Der Grund dafür ist darin zu suchen, dass die allgemeine Lösung der Differenzgleichung den Term

$$\lambda_2^j = [-5 + O(h)]^j$$

enthält, der für großes j beliebig groß wird. ♡

Das Beispiel lässt vermuten, dass die Lösungen der Differenzgleichung

$$\eta_{j+r} + \alpha_{r-1}\eta_{j+r-1} + \cdots + \alpha_0\eta_j = 0$$

entscheidenden Einfluss auf das Konvergenzverhalten des MSV haben. Genau dies ist der Fall. Ein MSV heißt **nullstabil**, falls das erste charakteristische Polynom

$$\Psi(\mu) = \mu^r + \alpha_{r-1}\mu^{r-1} + \cdots + \alpha_1\mu + \alpha_0$$

nur Nullstellen λ besitzt, für die $|\lambda| \leq 1$ gilt, und die im Falle $|\lambda| = 1$ einfach sind.

Bemerkung: Für die bisher behandelten MSV ist die Stabilitätsbedingung stets erfüllt. Man erhält für die ADAMS-BASHFORTH- und ADAMS-MOULTON-Verfahren die charakteristischen Polynome

$$\Psi(\mu) = \mu^r - \mu^{r-1} = (\mu - 1)\mu^{r-1}$$

und für die NYSTRÖM- und MILNE-Verfahren die charakteristischen Polynome

$$\Psi(\mu) = \mu^r - \mu^{r-2} = (\mu + 1)(\mu - 1)\mu^{r-2}.$$

Genauere Untersuchungen zeigen, dass es keine r -Schritt-Verfahren der Konsistenzordnung $2r$ gibt, die auch nullstabil sind. Es gilt:

6.15. DAHLQUIST'S 1. Ordnungsschranke: *Für die maximale Konsistenzordnung eines nullstabilen r -Schritt-Verfahrens gilt*

$$p = \begin{cases} r+1 & \text{für } r \text{ ungerade} \\ r+2 & \text{für } r \text{ gerade} \end{cases}.$$

Ein r -Schritt-Verfahren der Konsistenzordnung $r + 1$ ist z. B. durch ein Prediktor-Korrektor-Verfahren vom ADAMS-BASHFORTH-MOULTON-Typ gegeben.

Ein r -Schritt-Verfahren der Konsistenzordnung $r + 2$ ist das MILNE-Verfahren für $q = 2$:

$$\eta_{p+1} = \eta_{p-1} + \left[\frac{1}{3} \mathbf{f}(x_{p+1}, \eta_{p+1}) + \frac{4}{3} \mathbf{f}(x_p, \eta_p) + \frac{1}{3} \mathbf{f}(x_{p-1}, \eta_{p-1}) \right].$$

Löst man mit diesem Verfahren das AWP

$$y' = \kappa y, \quad y(0) = 1$$

mit der exakten Lösung $y(x) = e^{\kappa x}$, so erfüllen die Näherungslösungen $\eta_k = \eta(x_k; h)$ die homogene lineare Differenzgleichung

$$\left(1 - \frac{\kappa h}{3}\right) \eta_{k+1} - \frac{4\kappa h}{3} \eta_k - \left(1 + \frac{\kappa h}{3}\right) \eta_{k-1} = 0.$$

Die allgemeine Lösung dieser Differenzgleichung lässt sich wieder in der Form

$$\eta_k = \mu_1 \lambda_1^k + \mu_2 \lambda_2^k$$

angeben, wobei λ_1 und λ_2 Lösungen der quadratischen Gleichung

$$\left(1 - \frac{\kappa h}{3}\right) \lambda^2 - \frac{4\kappa h}{3} \lambda - \left(1 + \frac{\kappa h}{3}\right) = 0$$

sind. Es ergibt sich

$$\lambda_1 = \frac{1}{1 - \frac{\kappa h}{3}} \left(\frac{2\kappa h}{3} + \sqrt{1 + \frac{(\kappa h)^2}{3}} \right)$$

und

$$\lambda_2 = \frac{1}{1 - \frac{\kappa h}{3}} \left(\frac{2\kappa h}{3} - \sqrt{1 + \frac{(\kappa h)^2}{3}} \right).$$

Die Koeffizienten μ_1 und μ_2 sind aus den Anfangsbedingungen zu bestimmen.

Aus $y(0) = 0$ folgt $\mu_1 + \mu_2 = 1$. Damit gilt

$$\eta_1 = \mu_1 \lambda_1 + (1 - \mu_1) \lambda_2,$$

folglich

$$\mu_1 = \frac{\eta_1 - \lambda_2}{\lambda_1 - \lambda_2}.$$

Entwickelt man λ_1 und λ_2 nach Potenzen von h , so gilt in erster Näherung

$$\lambda_1 \doteq \left(1 + \frac{\kappa h}{3}\right) \left(\frac{2\kappa h}{3} + 1\right) \doteq 1 + \kappa h$$

und

$$\lambda_2 \doteq \left(1 + \frac{\kappa h}{3}\right) \left(\frac{2\kappa h}{3} - 1\right) \doteq -1 + \frac{\kappa h}{3}.$$

Für die Näherungslösung an der Stelle x erhalten wir so

$$\eta(x; h) = \mu_1 \left(1 + \kappa h + O(h^2)\right)^{x/h} + (1 - \mu_1) \left(-1 + \frac{\kappa h}{3} + O(h^2)\right)^{x/h}.$$

Der erste Term verhält sich wie $e^{\kappa x}$. Er entspricht der exakten Lösung. Der zweite Term dagegen verhält sich wie

$$(-1)^{x/h} e^{-\kappa x/3}.$$

Für $\kappa < 0$, also eine exponentiell abfallende Lösung, liefert dieser Term eine oszillierende exponentiell wachsende Störung. Genau dieses Verhalten beobachtet man in der Praxis bei allen r -Schritt-Verfahren der Konsistenzordnung $r + 2$, falls diese auf AWP mit exponentiell abklingenden Lösungen angewendet werden. Damit ist die Anwendbarkeit dieser Verfahren stark eingeschränkt.

Neben der Konsistenz von MSV interessiert auch wieder die Konvergenz. Bei ESV folgt aus der Konsistenz die Konvergenz, und Konsistenzordnung und Konvergenzordnung stimmen überein. Bei MSV sind die Verhältnisse so: Zunächst benötigt man, um ein r -Schritt-Verfahren zu starten, neben dem Anfangswert $\boldsymbol{\eta}_0 = \boldsymbol{y}_0$ weitere $r - 1$ Näherungen $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_{r-1}$ für die exakten Werte $\boldsymbol{y}_1, \dots, \boldsymbol{y}_{r-1}$. (Diese beschafft man sich zum Beispiel mit einem ESV.) Die Güte dieser Näherungen beeinflusst natürlich die Güte aller weiteren Näherungswerte. Wir nehmen an, dass sich diese Startnäherungen gemäß

$$\begin{aligned} \boldsymbol{\eta}_0 &= \boldsymbol{y}_0, \\ \boldsymbol{\eta}_1 &= \boldsymbol{y}_1 + \boldsymbol{\varepsilon}_1(h), \\ &\vdots \\ \boldsymbol{\eta}_{r-1} &= \boldsymbol{y}_{r-1} + \boldsymbol{\varepsilon}_{r-1}(h) \end{aligned}$$

von den exakten Werten unterscheiden. Die berechneten Näherungen $\boldsymbol{\eta}_r, \boldsymbol{\eta}_{r+1}, \dots$ hängen damit nicht nur von der Schrittweite und dem angewendeten MSV ab, sondern auch von der Güte der Anfangsnäherungen. Folglich ist für die Näherung $\boldsymbol{\eta}$ an der Stelle x zu schreiben

$$\boldsymbol{\eta}(x; \boldsymbol{\varepsilon}; h),$$

wobei ε eine Funktion bezeichnet, für die $\varepsilon_i = \varepsilon(x_i; h)$ für $i = 1, \dots, r-1$ gilt. Nun definieren wir den globalen Diskretisierungsfehler für MSV. Es sei \mathbf{y} die exakte Lösung des AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0.$$

$\boldsymbol{\eta}(x; \boldsymbol{\varepsilon}; h)$ bezeichne die mit einem MSV mit der Schrittweite h berechnete Näherung für $\mathbf{y}(x)$, wobei die Güte der Startnäherungen durch die Funktion $\boldsymbol{\varepsilon}$ beschrieben wird. Dann heißt die Größe

$$\mathbf{e}(x; \boldsymbol{\varepsilon}; h) = \boldsymbol{\eta}(x; \boldsymbol{\varepsilon}; h) - \mathbf{y}(x)$$

globaler Diskretisierungsfehler an der Stelle x zur Schrittweite h . Nun lässt sich auch sagen, was unter einem konvergenten MSV zu verstehen ist. Ein MSV zum Lösen von AWP der Form

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

heißt **konvergent**, falls

$$\lim_{n \rightarrow \infty} \mathbf{e}(x; \boldsymbol{\varepsilon}; h_n) = \mathbf{o}, \quad h_n = \frac{x - x_0}{n},$$

für alle $x \in [a, b]$, für alle Funktionen $\mathbf{f} \in F^1[a, b]$ und für alle Funktionen $\boldsymbol{\varepsilon}(z; h)$ mit

$$\lim_{n \rightarrow \infty} \|\boldsymbol{\varepsilon}(z; h_n)\| = 0, \quad z = x_0 + ih_n, \quad i = 1, \dots, r-1$$

gilt. Der folgende Satz, den wir ohne Beweis angeben, liefert Bedingungen für die Konvergenz von MSV.

6.16. Satz: *Es sei durch*

$$\boldsymbol{\eta}_{j+r} + \alpha_{r-1} \boldsymbol{\eta}_{j+r-1} + \dots + \alpha_0 \boldsymbol{\eta}_j = h \mathbf{F}(x_j, \boldsymbol{\eta}_{j+r}, \dots, \boldsymbol{\eta}_j; h; \mathbf{f})$$

ein konsistentes MSV zum Lösen eines AWP der Form

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

gegeben. Die Funktion \mathbf{F} erfülle folgende Bedingungen:

1. $\mathbf{F} \equiv \mathbf{o}$ für alle $x \in [a, b]$, alle $\boldsymbol{\eta}_i \in \mathbb{R}^d$ und alle $h \in \mathbb{R}$ falls $\mathbf{f} \equiv \mathbf{o}$.

2. Zu jeder Funktion $\mathbf{F} \in F^1[a, b]$ gibt es Konstanten $h_0 > 0$ und M , so dass

$$\|\mathbf{F}(x, \mathbf{v}_r, \dots, \mathbf{v}_0; h; \mathbf{f}) - \mathbf{F}(x, \mathbf{w}_r, \dots, \mathbf{w}_0; h; \mathbf{f})\| \leq M \sum_{i=0}^r \|\mathbf{v}_i - \mathbf{w}_i\|$$

für alle $x \in [a, b]$, alle $\mathbf{v}_i, \mathbf{w}_i \in \mathbb{R}$, $i = 0, \dots, r$ und alle $|h| \leq h_0$ gilt.

Dann ist das gegebene MSV genau dann konvergent, wenn es nullstabil ist.

Bemerkungen: (i) Die erste Voraussetzung garantiert gemeinsam mit der Nullstabilität, dass das MSV das triviale AWP

$$\mathbf{y}' = \mathbf{o}, \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

exakt löst, falls $\varepsilon_1 = \varepsilon_2 = \dots = \varepsilon_{r-1} = 0$ gilt.

(ii) Die zweite Bedingung bedeutet die LIPSCHITZ-Stetigkeit von \mathbf{F} bezüglich der Näherungen $\boldsymbol{\eta}_j, \dots, \boldsymbol{\eta}_{j+r}$.

(iii) Für lineare MSV ist die erste Bedingung trivialerweise erfüllt. Die zweite Bedingung folgt sofort aus der LIPSCHITZ-Stetigkeit von \mathbf{f} , die ja schon im Existenz- und Eindeigkeitssatz gefordert wurde.

Damit folgt die Konvergenz der im Abschnitt 6.3.1. konstruierten Verfahren. Der folgende Satz macht noch eine Aussage über die Konvergenzordnung.

6.17. Satz: *Es sei durch*

$$\boldsymbol{\eta}_{j+r} + \alpha_{r-1} \boldsymbol{\eta}_{j+r-1} + \dots + \alpha_0 \boldsymbol{\eta}_j = h \mathbf{F}(x_j, \boldsymbol{\eta}_{j+r}, \dots, \boldsymbol{\eta}_j; h; \mathbf{f})$$

ein nullstabiles MSV der Konsistenzordnung p zum Lösen eines AWP der Form

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

gegeben. Die Funktion \mathbf{F} erfülle alle Voraussetzungen aus Satz 6.16.

Dann gilt für alle $\mathbf{f} \in F^p[a, b]$ und alle $x \in [a, b]$

$$\|\mathbf{e}(x; \boldsymbol{\varepsilon}; h_n)\| = O(h_n^p),$$

falls für die Güte der Startnäherungen

$$\|\boldsymbol{\varepsilon}_i(h^p)\| = O(h^p), \quad i = 1, \dots, r-1$$

gilt.

Um für ein MSV der Konsistenzordnung p auch die Konvergenzordnung p zu erreichen, ist es somit notwendig, die Startnäherungen ebenfalls mit der entsprechenden Güte zu berechnen, folglich zum Beispiel mit einem ESV der Ordnung p . Wenn wir ein Gesamtverfahren betrachten, das aus einem Startverfahren zum Berechnen der Näherungen $\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_{r-1}$ und einem Prediktor-Korrektor-Verfahren zum Berechnen der weiteren Näherungen besteht, so wird die Ordnung des gesamten Verfahrens durch das Teilverfahren mit der geringsten Ordnung bestimmt.

6.4. Extrapolationsverfahren

Alle bisher betrachteten Verfahren zum Lösen von AWP lieferten Näherungen, die von der verwendeten Schrittweite abhingen. Für die exakte Lösung $\mathbf{y}(x)$ an der Stelle x wird eine Näherungslösung $\boldsymbol{\eta}(x; h_n)$ mit $h_n = (x - x_0)/n$ berechnet. Gelingt es nun, für die Abhängigkeit der Näherungslösung $\boldsymbol{\eta}(x; h)$ von der Schrittweite h asymptotische Entwicklungen, analog der Entwicklung der Trapezsumme bei der numerischen Integration, herzuleiten, so lassen sich wie beim ROMBERG-Verfahren durch Extrapolation aus verschiedenen Näherungen, die mit verschiedenen Schrittweiten berechnet wurden, verbesserte Näherungen gewinnen. Leider ist nur in wenigen Spezialfällen die Existenz einer asymptotischen Entwicklung bekannt.

6.18. Satz von GRAGG: *Es sei $\mathbf{f} \in F^{2N+2}[a, b]$ und \mathbf{y} die exakte Lösung des AWP*

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0.$$

Für $x \in R_h = \{x_0 + ih \mid i = 0, 1, \dots\}$ sei $\boldsymbol{\eta}(x; h)$ definiert durch

$$\left. \begin{aligned} \boldsymbol{\eta}_0 &= \boldsymbol{\eta}(x_0; h) = \mathbf{y}_0, \\ \boldsymbol{\eta}_1 &= \boldsymbol{\eta}(x_1; h) = \mathbf{y}_0 + h\mathbf{f}(x_0, \mathbf{y}_0), \\ \boldsymbol{\eta}_{k+1} &= \boldsymbol{\eta}(x_{k+1}; h) = \boldsymbol{\eta}_{k-1} + 2h\mathbf{f}(x_k, \boldsymbol{\eta}_k), \\ x_{k+1} &= x_k + h \end{aligned} \right\}, \quad k = 1, 2, \dots$$

Dann besitzt $\boldsymbol{\eta}(x; h)$ eine Entwicklung der Form

$$\boldsymbol{\eta}(x; h) = \mathbf{y}(x) + \sum_{i=1}^N h^{2i} \left[\mathbf{u}_i(x) + (-1)^{\frac{x-x_0}{h}} \mathbf{v}_i(x) \right] + h^{2N+2} \mathbf{E}_{2N+2}(x; h),$$

die für alle $x \in (a, b)$ und alle $h = (x - x_0)/n$, $n = 1, 2, \dots$, gilt. Die Funktionen \mathbf{u}_i und \mathbf{v}_i , $i = 1, 2, \dots, N$, sind von h unabhängig. Das Restglied \mathbf{E}_{2N+2} bleibt bei festem x für alle $h = (x - x_0)/n$ beschränkt.

Die in diesem Satz angegebene Entwicklung berechtigt uns noch nicht, analog zum ROMBERG-Verfahren vorzugehen. Dazu benötigen wir eine Entwicklung der Form

$$\boldsymbol{\eta}(x; h) = \mathbf{y}(x) + \sum_{i=1}^N \boldsymbol{\tau}_i(x) h^{2i} + h^{2N+2} \mathbf{E}_{2N+2}(x; h),$$

mit von h unabhängigen Funktionen $\boldsymbol{\tau}_i(x)$. In der Entwicklung aus Satz 6.18 hängt der Ausdruck

$$\mathbf{u}_i(x) + (-1)^{\frac{x-x_0}{h}} \mathbf{v}_i(x)$$

über den Term $(-1)^{(x-x_0)/h}$ von der Schrittweite h ab. Diese Abhängigkeit lässt sich beseitigen, indem man die Schrittweite h so wählt, dass der Exponent $(x-x_0)/h$ immer gerade oder immer ungerade ist. Wir wählen die Schrittweiten folglich in der Form

$$h = \frac{x-x_0}{2n} \quad \text{oder} \quad h = \frac{x-x_0}{2n-1}.$$

Das liefert im ersten Falle Entwicklungen der Form

$$\eta(x; h) = \mathbf{y}(x) + \sum_{i=1}^N h^{2i} [\mathbf{u}_i(x) + \mathbf{v}_i(x)] + h^{2N+2} \mathbf{E}_{2N+2}(x; h)$$

und im zweiten Falle Entwicklungen der Form

$$\eta(x; h) = \mathbf{y}(x) + \sum_{i=1}^N h^{2i} [\mathbf{u}_i(x) - \mathbf{v}_i(x)] + h^{2N+2} \mathbf{E}_{2N+2}(x; h).$$

Dies sind wirkliche asymptotische Entwicklungen in h^2 . Eine weitere Verbesserung der Entwicklung erreicht man durch folgenden Trick (GRAGG), der den Oszillationsterm im ersten Glied prinzipiell beseitigt.

6.19. GRAGGSche Näherung:

Es ist für die Stelle x eine Näherungslösung des AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_0) = \mathbf{y}_0$$

zu berechnen.

S0 Wähle ein n , setze $h = (x-x_0)/n$ und

$$\begin{aligned} \boldsymbol{\eta}_0 &= \boldsymbol{\eta}(x_0; h) = \mathbf{y}_0, \\ \boldsymbol{\eta}_1 &= \boldsymbol{\eta}(x_1; h) = \mathbf{y}_0 + h\mathbf{f}(x_0, \mathbf{y}_0). \end{aligned}$$

S1 Für $k = 1, 2, \dots, n-1$ berechne

$$\boldsymbol{\eta}_{k+1} = \boldsymbol{\eta}_{k-1} + 2h\mathbf{f}(x_k, \boldsymbol{\eta}_k), \quad x_{k+1} = x_k + h.$$

S2 Berechne

$$\tilde{\boldsymbol{\eta}}(x; h) = \frac{1}{2} [\boldsymbol{\eta}_n + \boldsymbol{\eta}_{n-1} + h\mathbf{f}(x_n, \boldsymbol{\eta}_n)].$$

Für die GRAGGsche Näherung lässt sich unter den Voraussetzungen von Satz 6.18 leicht zeigen, dass sie eine Entwicklung der Form

$$\begin{aligned} \tilde{\eta}(x; h) &= \mathbf{y}(x) + h^2 \left[\mathbf{u}_1(x) + \frac{1}{4} \mathbf{y}''(x) \right] + \sum_{i=2}^N h^{2i} \left[\tilde{\mathbf{u}}_i(x) + (-1)^{\frac{x-x_0}{h}} \tilde{\mathbf{v}}_i(x) \right] \\ &\quad + h^{2N+2} \tilde{\mathbf{E}}_{2N+2}(x; h) \end{aligned}$$

besitzt.

Damit ergibt sich der folgende Algorithmus.

6.20. Extrapolationsverfahren zum Lösen von AWP:

S0 Wähle Grundschriftweite H und eine zu erreichende Genauigkeit ε .

Setze $x = x_0 + H$, $h_0 = H/2$ und $k = 0$.

S1 Berechne mit der Schrittweite h_k eine GRAGG'sche Näherung

$$\mathbf{T}_{k0} = \tilde{\eta}(x; h_k).$$

S2 Für $l = 1, \dots, k$ berechne

$$\mathbf{T}_{k,l} = \mathbf{T}_{k,l-1} + \frac{\mathbf{T}_{k,l-1} - \mathbf{T}_{k-1,l-1}}{\left(\frac{h_{k-l}}{h_k}\right)^2 - 1}.$$

S3 Gilt

$$\|\mathbf{T}_{k,k} - \mathbf{T}_{k,k-1}\| \leq \varepsilon,$$

so setze $\mathbf{y}_0 = \boldsymbol{\eta}(x) = \mathbf{T}_{k,k}$, $x_0 = x$ und gehe zu Schritt **S0**.

S4 Setze $h_{k+1} = h_k/2$ und $k = k + 1$. Gehe zu Schritt **S1**.

Bemerkungen: (i) Die Abbruchbedingung im Schritt **S3** lässt sich etwas verfeinern. Man bricht den Aufbau des Interpolationsschemas im Schritt **S2** ab, falls sich zwei benachbarte Werte im Rahmen der Genauigkeit ε nicht mehr unterscheiden. Um zufällige Effekte zu vermeiden, sollte man erst dann abbrechen, wenn diese Bedingung mehrmals hintereinander erfüllt war.

(ii) Man sollte in Schritt **S2** nur 5...7 Spalten des Interpolationsschemas berechnen, da die Polynominterpolation für Polynome höheren Grades ungenaue Resultate liefert.

(iii) Statt der Polynominterpolation darf wieder Rationale Interpolation verwendet werden. Die Rekursionsformel in Schritt **S2** ist dann gemäß

$$\mathbf{T}_{k,l} = \mathbf{T}_{k,l-1} + \frac{\mathbf{T}_{k,l-1} - \mathbf{T}_{k-1,l-1}}{\left(\frac{h_{k-l}}{h_k}\right)^2 \left[1 - \frac{\mathbf{T}_{k,l-1} - \mathbf{T}_{k-1,l-1}}{\mathbf{T}_{k,l-1} - \mathbf{T}_{k-1,l-2}}\right] - 1}$$

abzuändern.

(iv) Die Grundschriftweite H wird von Schritt zu Schritt geändert. Man sollte sie so wählen, dass das Interpolationsschema nach 5...7 Spalten abbricht.

Die Konvergenzaussagen zur ROMBERG-Integration gelten entsprechend auch für die mit dem Algorithmus 6.20 berechneten Werte. Insbesondere gilt also

$$\lim_{k \rightarrow \infty} \mathbf{T}_{k,l} = \mathbf{y}(x).$$

Somit konvergieren die Spalten des Interpolationsschemas gegen die exakte Lösung. Die Konvergenzordnung ist dabei $2k + 2$ falls $\mathbf{f} \in F^{2k+2}[a, b]$ und $k \leq N$.

6.5. Aufgaben

1. Man bestimme die Lösung des Anfangswertproblems

$$\mathbf{y}' = \mathbf{A}\mathbf{y}, \quad \mathbf{y}(0) = \mathbf{y}_0$$

mit $\mathbf{y}_0 \in \mathbb{R}^d$ und der $d \times d$ -Matrix

$$\mathbf{A} = \begin{bmatrix} \lambda & 1 & & & 0 \\ & \lambda & 1 & & \\ & & \cdot & \cdot & \\ & & & \cdot & \cdot \\ & & & & \lambda & 1 \\ 0 & & & & & \lambda \end{bmatrix}, \quad \lambda \in \mathbb{R}.$$

2. Gegeben sei das Anfangswertproblem

$$y' = x^2 + 2x^3 \quad y(0) = 0.$$

Zur Schrittweite h sollen mit dem EULER-Verfahren Näherungswerte $\eta(x_j; h)$ für $y(x_j)$, $x_j = jh$ berechnet werden. Man gebe $\eta(x_j; h)$ und den globalen Diskretisierungsfehler $e(x_j; h)$ explizit an und zeige, dass $e(x; h)$ bei festem x für $h = \frac{x}{n} \rightarrow 0$ gegen Null geht.

3. Es sei ein Einschrittverfahren durch

$$\Phi(x, y; h) = \sum_{i=0}^n \alpha_i f(\xi_i, \eta_i)$$

mit

$$\xi_0 = x, \quad \xi_i = x + \vartheta_i h, \quad i = 1, \dots, n$$

$$\eta_0 = y, \quad \eta_i = y + h \sum_{j=0}^{i-1} \beta_{ij} f(\xi_j, \eta_j), \quad i = 1, \dots, n$$

gegeben. Welche Gleichungen haben die Parameter α_i, ϑ_i und β_{ij} zu erfüllen, damit man für $n = 2$ ein Verfahren 3. Ordnung erhält? Man gebe eine Lösung dieser Gleichungen an.

4. Man schreibe ein Programm zum Berechnen der Iterationsfunktion $\Phi(x, y; h)$ für beliebiges n und beliebige Parameter; n und die Parameter sollen beim Aufruf des Programms übergeben werden.
5. Das Anfangswertproblem

$$y' = y^{\frac{1}{2}}, \quad y(0) = 0$$

hat die Lösung $y(x) = \frac{x^2}{4}$. Die Anwendung des EULER-Verfahrens liefert jedoch $\eta(x; h) = 0$ für alle x und $h = \frac{x}{n}$, $n = 1, 2, \dots$. Man begründe dieses Verhalten.

6. Gegeben sei das Anfangswertproblem

$$y' = \lambda y, \quad y(0) = 1.$$

Welche Schrittweitschätzung liefert die Schrittweitensteuerung

- (a) nach Algorithmus 6.9 mit dem EULER-Verfahren,
 - (b) nach Algorithmus 6.9 mit dem impliziten EULER-Verfahren,
 - (c) nach der RUNGE-KUTTA-FEHLBERG-Methode mit EULER- und HEUN-Verfahren?
7. Man zeige, dass das modifizierte EULER-Verfahren die exakte Lösung der Differentialgleichung

$$y' = -2ax$$

liefert.

8. Durch

$$F(x, y; h) = f(x, y) + h \frac{g(x + \frac{h}{3}, y + \frac{hf(x, y)}{3})}{2}$$

mit

$$g(x, y) = f_x(x, y) + f_y(x, y)f(x, y)$$

ist ein Einschrittverfahren gegeben. Man zeige, dass es sich um ein Verfahren 3. Ordnung handelt.

9. Bei den ADAMS-MOULTON-Verfahren wird, ausgehend von den Näherungswerten

$$\eta_{p-j}, \dots, \eta_p \text{ für } y(x_{p-j}), \dots, y(x_p)$$

ein Näherungswert η_{p+1} für $y(x_{p+1}), x \in [a, b]$, durch folgende Iterationsvorschrift berechnet:

$$\eta_{p+1}^{(0)} \text{ beliebig}$$

für $i = 0, 1, \dots$

$$\eta_{p+1}^{(i+1)} = \Psi(\eta_{p+1}^{(i)}) = \eta_p + h[\beta_{q0}f(x_{p+1}, \eta_{p+1}^{(i)}) + \beta_{q1}f_p + \dots + \beta_{qq}f_{p+1-q}].$$

Man zeige: Für einmal stetig partiell differenzierbare Funktionen f gibt es ein $h_0 > 0$, so dass für alle $|h| \leq h_0$ die Folge $\eta_{p+1}^{(i)}$ gegen ein $\eta_{p+1} = \Psi(\eta_{p+1})$ konvergiert.

10. Man prüfe, ob das lineare Mehrschrittverfahren

$$\eta_p - \eta_{p-4} = \frac{h}{3}(8f_{p-1} - 4f_{p-2} + 8f_{p-3})$$

konvergent ist.

11. Man bestimme α , β und γ so, dass das lineare Mehrschrittverfahren

$$\eta_{j+4} - \eta_{j+2} + \alpha(\eta_{j+3} - \eta_{j+1}) = h[\beta(f_{j+3} - f_{j+1}) + \gamma f_{j+2}]$$

die Ordnung 3 hat. Ist das so gewonnene Verfahren stabil?

12. Es werde das durch

$$\eta_{j+2} + a_1\eta_{j+1} + a_0\eta_j = h[b_0f(x_j, \eta_j) + b_1f(x_{j+1}, \eta_{j+1})]$$

gegebene PREDIKTOR-Verfahren betrachtet.

- (a) Man bestimme a_0 , b_0 und b_1 in Abhängigkeit von a_1 so, dass ein Verfahren mindestens 2. Ordnung entsteht.
- (b) Für welche a_1 -Werte ist das so gewonnene Verfahren stabil?
- (c) Welche speziellen Verfahren erhält man für $a_1 = 0$ und $a_1 = -1$?
- (d) lässt sich a_1 so wählen, dass man stets ein stabiles Verfahren 3. Ordnung erhält?

Kapitel 7

Randwertprobleme

7.1. Einführung

Allgemeiner als AWP sind Randwertprobleme (RWP):

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{r}(\mathbf{y}(a), \mathbf{y}(b)) = \mathbf{0}$$

mit

$$\begin{aligned} \mathbf{y} &: [a, b] \rightarrow \mathbb{R}^n, \\ \mathbf{f} &: [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n, \\ \mathbf{r} &: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^n. \end{aligned}$$

Oft liegen die Randbedingungen in linearer Form vor:

$$\mathbf{r}(\mathbf{y}(a), \mathbf{y}(b)) = \mathbf{A}\mathbf{y}(a) + \mathbf{B}\mathbf{y}(b) - \mathbf{c}, \quad \mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}, \quad \mathbf{c} \in \mathbb{R}^n$$

gilt. Man spricht von separierten Randbedingungen, falls

$$\mathbf{A}_1\mathbf{y}(a) = \mathbf{c}_1, \quad \mathbf{B}_2\mathbf{y}(b) = \mathbf{c}_2$$

gilt, also Permutationsmatrizen \mathbf{P}_1 und \mathbf{P}_2 existieren, so dass die Zeilen und Spalten der Matrix $(\mathbf{A}, \mathbf{B}, \mathbf{c})$ so vertauschbar sind, dass

$$\mathbf{P}_1(\mathbf{A}, \mathbf{B}, \mathbf{c})\mathbf{P}_2 = \begin{pmatrix} \mathbf{A}_1 & \mathbf{0} & \mathbf{c}_1 \\ \mathbf{0} & \mathbf{B}_2 & \mathbf{c}_2 \end{pmatrix}$$

gilt. AWP sind spezielle RWP mit $\mathbf{A} = \mathbf{I}$, $\mathbf{B} = \mathbf{0}$ und $\mathbf{c} = \mathbf{y}_0$.

Für RWP existiert kein Existenz- und Eindeutigkeitsatz wie für AWP, der unter schwachen Voraussetzungen die Existenz und Eindeutigkeit der Lösung sichert.

7.1. Beispiel: Es sei

$$\mathbf{y}' = \begin{pmatrix} y_1' \\ y_2' \end{pmatrix} = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix} \mathbf{y} = \mathbf{A}\mathbf{y}.$$

Die exakte Lösung ist

$$\mathbf{y}(x) = e^{\mathbf{A}x} \mathbf{c}, \quad \mathbf{c} \in \mathbb{R}^2.$$

Mit

$$\mathbf{A} = \frac{1}{\sqrt{2}} \begin{pmatrix} -i & i \\ 1 & 1 \end{pmatrix} \begin{pmatrix} i & 0 \\ 0 & -i \end{pmatrix} \frac{1}{\sqrt{2}} \begin{pmatrix} i & 1 \\ -i & 1 \end{pmatrix}$$

erhält man

$$\begin{aligned} \exp(\mathbf{A}x) &= \frac{1}{2} \begin{pmatrix} -i & i \\ 1 & 1 \end{pmatrix} \begin{pmatrix} e^{ix} & 0 \\ 0 & e^{-ix} \end{pmatrix} \begin{pmatrix} i & 1 \\ -i & 1 \end{pmatrix} \\ &= \frac{1}{2} \begin{pmatrix} -ie^{ix} & ie^{-ix} \\ e^{ix} & e^{-ix} \end{pmatrix} \begin{pmatrix} i & 1 \\ -i & 1 \end{pmatrix} \\ &= \frac{1}{2} \begin{pmatrix} e^{ix} + e^{-ix} & -ie^{ix} + ie^{-ix} \\ ie^{ix} - ie^{-ix} & e^{ix} + e^{-ix} \end{pmatrix} \\ &= \begin{pmatrix} \cos x & \sin x \\ -\sin x & \cos x \end{pmatrix}. \end{aligned}$$

Die allgemeine Lösung der Differentialgleichung lautet damit

$$\begin{aligned} \mathbf{y} &= \begin{pmatrix} \cos x & \sin x \\ -\sin x & \cos x \end{pmatrix} \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} \\ &= \begin{pmatrix} c_1 \cos x + c_2 \sin x \\ -c_1 \sin x + c_2 \cos x \end{pmatrix}. \end{aligned}$$

Betrachten wir nun verschiedene Randbedingungen.

1. $(1, 0)\mathbf{y}(0) = 0$ und $(1, 0)\mathbf{y}(\pi/2) = 1$:

Es folgt

$$(1, 0) \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = c_1 = 0$$

und

$$(1, 0) \begin{pmatrix} c_2 \\ -c_1 \end{pmatrix} = c_2 = 1.$$

Damit erhalten wir die eindeutige Lösung

$$\mathbf{y} = \begin{pmatrix} \sin x \\ \cos x \end{pmatrix}.$$

2. $(1, 0)\mathbf{y}(0) = 0$ und $(1, 0)\mathbf{y}(\pi) = 1$:

Es folgt

$$(1, 0) \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = c_1 = 0$$

und

$$(1, 0) \begin{pmatrix} -c_1 \\ -c_2 \end{pmatrix} = -c_1 = 1.$$

Das ist ein Widerspruch. Es existiert keine Lösung dieses RWP.

3. $(1, 0)\mathbf{y}(0) = 0$ und $(1, 0)\mathbf{y}(\pi) = 0$:

Es folgt

$$(1, 0) \begin{pmatrix} c_1 \\ c_2 \end{pmatrix} = c_1 = 0$$

und

$$(1, 0) \begin{pmatrix} -c_1 \\ -c_2 \end{pmatrix} = -c_1 = 0.$$

In diesem Falle ist c_2 beliebig wählbar. Es existieren unendlich viele Lösungen

$$\mathbf{y} = \begin{pmatrix} c_2 \sin x \\ c_2 \cos x \end{pmatrix}.$$



Unter starken Voraussetzungen lässt sich ein Existenz- und Eindeutigkeitssatz formulieren.

7.2. Satz: Für das RWP $\mathbf{y}' = \mathbf{f}(x, \mathbf{y})$, $r(\mathbf{y}(a), \mathbf{y}(b)) = \mathbf{o}$ gelte

1. \mathbf{f} und $\frac{\partial}{\partial \mathbf{y}} \mathbf{f}$ sind stetig auf dem Streifen $S = [a, b] \times \mathbb{R}^n$.

2. Es existiert eine Funktion $K \in C[a, b]$ mit

$$\left\| \frac{\partial}{\partial \mathbf{y}} \mathbf{f}(x, \mathbf{y}) \right\| \leq K(x) \quad \text{für alle } (x, \mathbf{y}) \in S.$$

3. Die Matrix

$$\mathbf{P}(\mathbf{u}, \mathbf{v}) = \frac{\partial}{\partial \mathbf{u}} \mathbf{r}(\mathbf{u}, \mathbf{v}) + \frac{\partial}{\partial \mathbf{v}} \mathbf{r}(\mathbf{u}, \mathbf{v})$$

besitzt für alle $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ eine Darstellung der Form

$$\mathbf{P}(\mathbf{u}, \mathbf{v}) = \mathbf{P}_0 (\mathbf{I} + \mathbf{M}(\mathbf{u}, \mathbf{v}))$$

mit einer konstanten regulären Matrix \mathbf{P}_0 und einer Matrix $\mathbf{M}(\mathbf{u}, \mathbf{v})$, für die

$$\left\| \mathbf{P}_0^{-1} \frac{\partial}{\partial \mathbf{v}} \mathbf{r}(\mathbf{u}, \mathbf{v}) \right\| \leq m$$

und

$$\|\mathbf{M}(\mathbf{u}, \mathbf{v})\| \leq \mu$$

für alle $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ mit Konstanten $m < \infty$ und $\mu < 1$ gilt.

4. Es existiert eine Zahl λ mit $0 < \lambda < 1 - \mu$, so dass

$$\int_a^b K(t) dt \leq \ln \left(1 + \frac{\lambda}{m} \right).$$

Dann besitzt das obige RWP genau eine Lösung $\mathbf{y}(x)$.

Wie das Beispiel zeigte, sind die Voraussetzungen dieses Satzes schon für einfache Fälle nicht erfüllt.

7.2. Das einfache Schießverfahren

Wir wollen das RWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{r}(\mathbf{y}(a), \mathbf{y}(b)) = \mathbf{o}$$

lösen. Dazu wenden wir die folgende Idee an. Wir betrachten das AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(a) = \mathbf{s}$$

mit dem frei wählbaren Parameter \mathbf{s} . Es sei $\mathbf{y}(x; \mathbf{s})$ die exakte Lösung dieses AWP. An der Stelle $x = b$ hat diese Lösung dann den Wert $\mathbf{y}(b; \mathbf{s})$. Damit die Lösung des AWP für ein gewisses \mathbf{s} auch Lösung des RWP ist, muss

$$\mathbf{r}(\mathbf{y}(a, \mathbf{s}), \mathbf{y}(b, \mathbf{s})) = \mathbf{r}(\mathbf{s}, \mathbf{y}(b, \mathbf{s})) = \mathbf{o}$$

gelten. Wir versuchen, den Parameter s so zu bestimmen, dass die Lösung $\mathbf{y}(x; s)$ des AWP auch die Randbedingungen erfüllt. Diese Vorgehensweise wird als einfaches Schießverfahren bezeichnet.

7.3. Einfaches Schießverfahren zum Lösen von RWP:

Gegeben sei das RWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{r}(\mathbf{y}(a), \mathbf{y}(b)) = \mathbf{o}.$$

S0 Bestimme den Parameter s so, dass die Lösung $\mathbf{y}(x; s)$ des AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(a) = \mathbf{s}$$

die Randbedingungen

$$\mathbf{r}(\mathbf{y}(a; s), \mathbf{y}(b; s)) = \mathbf{r}(\mathbf{s}, \mathbf{y}(b; s)) = \mathbf{o}.$$

erfüllt.

Wir haben ein nichtlineares Gleichungssystem

$$\mathbf{F}(s) = \mathbf{r}(s, \mathbf{y}(b; s)) = \mathbf{o}$$

zu lösen. Dazu ist prinzipiell jedes Verfahren anwendbar, das zum Berechnen von Nullstellen für Funktionen geeignet ist. Die Schwierigkeit besteht darin, dass bei jeder Funktionswertberechnung von \mathbf{F} ein AWP zu lösen ist. Im eindimensionalen Fall

$$y : [a, b] \rightarrow \mathbb{R}$$

lassen sich die Methoden aus Kapitel 5 anwenden, um eine Nullstelle der entsprechenden nichtlinearen Gleichung $F(s) = r(s, y(b; s)) = 0$ zu berechnen. Die einfachste Möglichkeit wäre, das Bisektionsverfahren anzuwenden.

7.4. Bisektionsverfahren zum Lösen eindimensionaler RWP:

Gegeben sei das RWP

$$y' = f(x, y), \quad r(y(a), y(b)) = 0.$$

Für einen gewissen Parameter s sei $y(x; s)$ die Lösung des AWP

$$y' = f(x, y), \quad y(a) = s.$$

Weiterhin sei eine Funktion F gemäß

$$F(s) = r(s, y(b; s))$$

definiert. Gegeben seien Werte s_0 und t_0 mit

$$F(s_0)F(t_0) < 0.$$

S0 Setze $k = 0$.

S1 Setze $c = (s_k + t_k)/2$ und berechne die Lösung des AWP

$$y' = f(x, y), \quad y(a) = c$$

an der Stelle $x = b$.

S2 Berechne $d = F(c)$.

S3 Für

$$d \cdot F(s_k) \begin{cases} > 0 & s_{k+1} = c, \quad t_{k+1} = t_k, \\ = 0 & s^* = c, \quad \text{STOPP}, \\ < 0 & s_{k+1} = s_k, \quad t_{k+1} = c. \end{cases}$$

S4 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Um schnellere Konvergenz zu erhalten, verwendet man das NEWTON-Verfahren:

$$s_{k+1} = s_k - \frac{F(s_k)}{F'(s_k)}.$$

Dazu ist aber die Kenntnis der ersten Ableitung von $F(s) = r(s, y(b; s))$ notwendig. Es gilt

$$F'(s) = \frac{d}{ds} r(s, y(b; s)) = r_u(s, y(b; s)) + r_v(s, y(b; s)) \frac{d}{ds} y(b; s).$$

Die Ableitung $\frac{d}{ds} y(b; s)$ wird aus einem weiteren AWP berechnet. Durch formale Integration erhält man aus dem AWP

$$y' = f(x, s), \quad y(a) = s$$

mit der exakten Lösung $y(x; s)$

$$\int_a^x y'(t; s) dt = \int_a^x f(t, y(t; s)) dt,$$

$$y(x; s) - y(a; s) = \int_a^x f(t, y(t; s)) dt,$$

$$y(x; s) = s + \int_a^x f(t, y(t; s)) dt,$$

$$\frac{\partial}{\partial s} y(x; s) = 1 + \int_a^x \frac{\partial}{\partial y} f(t, y(t; s)) \frac{\partial}{\partial s} y(t; s) dt.$$

Es sei nun für festes s

$$v(x) = \frac{\partial}{\partial s} y(x; s),$$

also

$$v(x) = 1 + \int_a^x f_y(t, y(t; s)) v(t) dt,$$

woraus $v' = f_y(x, y(x; s))v$, $v(a) = 1$ folgt. Damit ergibt sich der folgende Algorithmus.

7.5. Newton-Verfahren zum Lösen von eindimensionalen RWP:

S0 Wähle s_0 und setze $k = 0$.

S1 Berechne simultan die Lösungen $y(b; s_k)$ und $v(b; s_k)$ der AWP

$$y' = f(x, y), \quad y(a) = s_k$$

$$v' = f_y(x, y)v, \quad v(a) = 1$$

an der Stelle $x = b$.

S2 Berechne $F(s_k) = r(s_k, y(b; s_k))$.

S3 Berechne

$$F'(s_k) = r_u(s_k, y(b; s_k)) + r_v(s_k, y(b; s_k))v(b; s_k).$$

S4 Berechne

$$s_{k+1} = s_k - \frac{F(s_k)}{F'(s_k)}.$$

Setze $k = k + 1$ und gehe zu Schritt **S1**.

Wegen der partiellen Ableitung von f nach y ist das zusätzliche AWP im allgemeinen wesentlich komplizierter. Damit wird die Anwendung des NEWTON-Verfahrens aufwendig. Um das Lösen dieses komplizierten AWP zu vermeiden, wendet man ein Quasi-NEWTON-Verfahren an. Hier wird $F'(s)$ durch einen Differenzenquotienten

$$F'(s) \approx \frac{F(s + \Delta s) - F(s)}{\Delta s}$$

approximiert. Zum Berechnen dieses Differenzenquotienten ist ein weiteres AWP mit der Anfangsbedingung $y(a) = s + \Delta s$ zu lösen. In diesem Falle erhalten wir folgendes Verfahren.

7.6. Quasi-Newton-Verfahren zum Lösen von eindimensionalen RWP:**S0** Wähle s_0 und setze $k = 0$.**S1** Berechne die Lösung $y(b; s_k)$ des AWP

$$y' = f(x, y), \quad y(a) = s_k$$

an der Stelle $x = b$.**S2** Wähle ein Δs_k und berechne die Lösung $y(b; s_k + \Delta s_k)$ des AWP

$$y' = f(x, y), \quad y(a) = s_k + \Delta s_k$$

an der Stelle $x = b$.**S3** Berechne $F(s_k) = r(s_k, y(b; s_k))$ und

$$F(s_k + \Delta s_k) = r(s_k + \Delta s_k, y(b; s_k + \Delta s_k)).$$

S4 Berechne

$$D_{\Delta s_k} F(s_k) = \frac{F(s_k + \Delta s_k) - F(s_k)}{\Delta s_k}.$$

S5 Berechne

$$s_{k+1} = s_k - \frac{F(s_k)}{D_{\Delta s_k} F(s_k)}$$

S6 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Die Konvergenz des Verfahrens hängt dabei stark von der Wahl der Δs_k ab. Wählt man Δs_k zu groß, so ist $D_{\Delta s_k} F(s_k)$ eine schlechte Approximation für $F'(s_k)$ und die Iteration konvergiert wesentlich schlechter als das NEWTON-Verfahren. Wählt man Δs_k zu klein, so gilt $F(s_k + \Delta s_k) \approx F(s_k)$ und beim Berechnen von $D_{\Delta s_k} F(s_k)$ tritt Auslöschung auf, was zu großen Fehlern in der Lösung führt.

Man hat folglich, um mit diesem Verfahren gute Resultate zu erzielen, zum einen in jedem Schritt Δs_k vernünftig wählen, und zum anderen die AWP zum Berechnen von $F(s_k + \Delta s_k)$ und $F(s_k)$ genau lösen. (Der relative Fehler von $F(s_k + \Delta s_k)$ und $F(s_k)$ sollte nur in der Größenordnung von eps liegen.) Diese Genauigkeit lässt sich mit Extrapolationsverfahren erreichen. Als vernünftige Wahl von Δs_k hat sich

$$\Delta s_k = \sqrt{\text{eps}} s_k$$

erwiesen. Hier wird erreicht, dass sich bei der Differenzbildung $F(s_k + \Delta s_k) - F(s_k)$ ungefähr die Hälfte der Mantissenstellen auslöscht.

Betrachten wir nun noch den allgemeinen Fall:

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{r}(\mathbf{y}(a), \mathbf{y}(b)) = \mathbf{o}$$

mit

$$\begin{aligned} \mathbf{y} &: [a, b] \rightarrow \mathbb{R}^n \\ \mathbf{f} &: [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n \\ \mathbf{r} &: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}^n. \end{aligned}$$

In diesem Falle ist das nichtlineare Gleichungssystem

$$\mathbf{F}(\mathbf{s}) = \mathbf{r}(\mathbf{s}, \mathbf{y}(b; \mathbf{s})) = \mathbf{0}$$

zu lösen. Man wendet dazu das allgemeine NEWTON-Verfahren

$$\mathbf{s}^{(k+1)} = \mathbf{s}^{(k)} - \left[\mathbf{F}'(\mathbf{s}^{(k)}) \right]^{-1} \mathbf{F}(\mathbf{s}^{(k)})$$

an. Hierbei bezeichnet \mathbf{F}' die JACOBI-Matrix von \mathbf{F} :

$$\mathbf{F}'(\mathbf{s}) = \mathbf{r}_u(\mathbf{s}, \mathbf{y}(b; \mathbf{s})) + \mathbf{r}_v(\mathbf{s}, \mathbf{y}(b; \mathbf{s})) \mathbf{y}_s(b; \mathbf{s})$$

mit

$$\begin{aligned} \mathbf{r}_u(\mathbf{u}, \mathbf{v}) &= \left(\frac{\partial \mathbf{r}_i(\mathbf{u}, \mathbf{v})}{\partial \mathbf{u}_j} \right), \\ \mathbf{r}_v(\mathbf{u}, \mathbf{v}) &= \left(\frac{\partial \mathbf{r}_i(\mathbf{u}, \mathbf{v})}{\partial \mathbf{v}_j} \right), \\ \mathbf{y}_s(b; \mathbf{s}) &= \left(\frac{\partial \mathbf{y}_i(b; \mathbf{s})}{\partial \mathbf{s}_j} \right) = \mathbf{Z}(b; \mathbf{s}). \end{aligned}$$

$\mathbf{Z}(b; \mathbf{s})$ ergibt sich für festes \mathbf{s} wieder als Lösung eines zusätzlichen AWP

$$\mathbf{Z}'(x; \mathbf{s}) = \mathbf{f}_y(x, \mathbf{y}(x; \mathbf{s})), \quad \mathbf{Z}(a; \mathbf{s}) = \mathbf{I}.$$

Dies ist ein Differentialgleichungssystem der Dimension n^2 ! Das Lösen dieses komplizierten Differentialgleichungssystems wird umgangen, indem man \mathbf{F}' durch Differenzenquotienten approximiert

$$\mathbf{F}'(\mathbf{s}) \approx \mathbf{D}_{\Delta \mathbf{s}} \mathbf{F}(\mathbf{s}) = (\mathbf{D}_{\Delta s_1} \mathbf{F}(\mathbf{s}), \mathbf{D}_{\Delta s_2} \mathbf{F}(\mathbf{s}), \dots, \mathbf{D}_{\Delta s_n} \mathbf{F}(\mathbf{s}))$$

mit

$$\mathbf{D}_{\Delta s_i} \mathbf{F}(\mathbf{s}) = \frac{\mathbf{F}(s_1, \dots, s_{i-1}, s_i + \Delta s_i, s_{i+1}, \dots, s_n) - \mathbf{F}(\mathbf{s})}{\Delta s_i}, \quad i = 1, \dots, n.$$

Zum Berechnen von Approximationen der partiellen Ableitungen ist jeweils wieder ein zusätzliches AWP zu lösen. Wir erhalten damit das folgende Verfahren.

7.7. Quasi-Newton-Verfahren zum Lösen von allgemeinen RWP:**S0** Wähle $\mathbf{s}^{(0)}$ und setze $k = 0$.**S1** Berechne die Lösung $\mathbf{y}(b; \mathbf{s}^{(k)})$ des AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(a) = \mathbf{s}^{(k)}$$

an der Stelle $x = b$.**S2** Wähle ein

$$\Delta \mathbf{s}^{(k)} = (\Delta s_1^{(k)}, \dots, \Delta s_n^{(k)})^T$$

und berechne für $i = 1, \dots, n$ die Lösungen

$$\mathbf{y}^{(i)} = \mathbf{y}(b; \mathbf{s}^{(k)} + \Delta s_i^{(k)} \mathbf{e}_i)$$

der AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(a) = \mathbf{s}^{(k)} + \Delta s_i^{(k)} \mathbf{e}_i$$

an der Stelle $x = b$.**S3** Berechne

$$\mathbf{F}(\mathbf{s}^{(k)}) = \mathbf{r}(\mathbf{s}^{(k)}, \mathbf{y}(b; \mathbf{s}^{(k)}))$$

und für $i = 1, \dots, n$

$$\mathbf{F}(\mathbf{s}^{(k)} + \Delta s_i^{(k)} \mathbf{e}_i) = \mathbf{r}(\mathbf{s}^{(k)} + \Delta s_i^{(k)} \mathbf{e}_i, \mathbf{y}^{(i)}).$$

S4 Berechne für $i = 1, \dots, n$

$$D_{\Delta s_i^{(k)}} \mathbf{F}(\mathbf{s}^{(k)}) = \frac{\mathbf{F}(\mathbf{s}^{(k)} + \Delta s_i^{(k)} \mathbf{e}_i) - \mathbf{F}(\mathbf{s}^{(k)})}{\Delta s_i^{(k)}}.$$

S5 Setze

$$D_{\Delta \mathbf{s}^{(k)}} \mathbf{F}(\mathbf{s}^{(k)}) = \left(D_{\Delta s_1^{(k)}} \mathbf{F}(\mathbf{s}^{(k)}), \dots, D_{\Delta s_n^{(k)}} \mathbf{F}(\mathbf{s}^{(k)}) \right).$$

S6 Löse das lineare Gleichungssystem

$$D_{\Delta \mathbf{s}^{(k)}} \mathbf{F}(\mathbf{s}^{(k)}) \mathbf{d}^{(k)} = \mathbf{F}(\mathbf{s}^{(k)}).$$

S7 Setze $\mathbf{s}^{(k+1)} = \mathbf{s}^{(k)} - \mathbf{d}^{(k)}$.

S8 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Pro Schritt sind hier $n + 1$ AWP zu lösen. Die Konvergenz des Verfahrens ist lokal und hängt von der Wahl des Startwertes $\mathbf{s}^{(0)}$ und von der Wahl der Größen $\Delta s_i^{(k)}$, $i = 1, \dots, n$, ab. Praktisch wendet man besser ein modifiziertes NEWTON-Verfahren an, bei dem Schritt **S6** durch

S6' Wähle eine Schrittweite λ_k und setze

$$\mathbf{s}^{(k+1)} = \mathbf{s}^{(k)} - \lambda_k \mathbf{d}^{(k)}.$$

ersetzt wird. Durch geschickte Wahl der Schrittweiten λ_k erreicht man eine bessere Konvergenz.

Für lineare RWP

$$\mathbf{y}' = \mathbf{T}(x)\mathbf{y} + \mathbf{g}(x), \quad \mathbf{A}\mathbf{y}(a) + \mathbf{B}\mathbf{y}(b) = \mathbf{c}$$

lässt sich zeigen, dass das obige Verfahren bei beliebigem Anfangswert $\mathbf{s}^{(0)}$ in einem Schritt die Lösung liefert.

Probleme beim einfachen Schießverfahren

Zum Lösen des RWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{r}(\mathbf{y}(a), \mathbf{y}(b)) = \mathbf{o}$$

wird beim einfachen Schießverfahren ein Wert \mathbf{s} so bestimmt, dass die Lösung des AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(a) = \mathbf{s}$$

auch die Randbedingungen erfüllt. Eigentlich interessieren uns aber die Werte der Lösung $\mathbf{y}(x)$ des RWP im vorgegebenen Intervall $[a, b]$. Prinzipiell sind Werte dieser Lösung berechenbar, indem das AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(a) = \mathbf{s}$$

näherungsweise für verschiedenen $x \in [a, b]$ gelöst wird. Hierbei tritt aber ein Problem auf. Nach Satz 6.3 hängt die Lösung eines AWP gemäß

$$\|\mathbf{y}(x; \mathbf{s}_1) - \mathbf{y}(x; \mathbf{s}_2)\| \leq \|\mathbf{s}_1 - \mathbf{s}_2\| e^{L|x-a|}$$

von den Anfangswerten ab. Das bedeutet aber, dass auch bei genauer Bestimmung des Anfangswertes s , für den die Lösung des entsprechenden AWP auch die Randbedingung erfüllt, die Lösung $\mathbf{y}(x; s)$ beliebig stark von der exakten Lösung des RWP abweicht (falls nur der Term $e^{L|x-a|}$ hinreichend groß ist). Das einfache Berechnen eines geeigneten Startwertes s genügt daher nicht, um die Lösung eines RWP mit hinreichender Genauigkeit zu berechnen.

7.8. Beispiel: Wir betrachten das RWP

$$\begin{aligned} y_1' &= y_2 & , & & y_1(0) &= 1, \\ y_2' &= 110y_1 + y_2 & , & & y_1(10) &= 1 \end{aligned}$$

und dazu das AWP

$$\begin{aligned} y_1' &= y_2 & , & & y_1(0) &= 1, \\ y_2' &= 110y_1 + y_2 & , & & y_2(0) &= s. \end{aligned}$$

Die allgemeine Lösung dieses AWP lautet

$$\begin{aligned} y_1(x; s) &= \frac{11-s}{21}e^{-10x} + \frac{10+s}{21}e^{11x}, \\ y_2(x; s) &= -10\frac{11-s}{21}e^{-10x} + 11\frac{10+s}{21}e^{11x}. \end{aligned}$$

Beim einfachen Schießverfahren müßte der Parameter s die Gleichung $y_1(10; s) = 1$ erfüllen. Es ergibt sich

$$\frac{11-s}{21}e^{-100} + \frac{10+s}{21}e^{110} = 1$$

und daraus

$$s = -10 + 21 \frac{1 - e^{-100}}{e^{110} - e^{-100}} \approx -10 + 3.5 \cdot 10^{-47} \approx -10.$$

Beim Lösen dieses Problems auf einem realen Rechner darf man höchstens erwarten, dass s mit relativer Maschinengenauigkeit berechnet wird. Man erhält eine Näherung $\bar{s} = -10(1 + \varepsilon)$ mit $|\varepsilon| \leq \text{eps}$. Schon für einen Fehler $\varepsilon = -10^{-10}$ erhält man aber

$$y_1(10; \bar{s}) \approx 2.8 \cdot 10^{37}.$$

Selbst beim Berechnen von s mit voller Maschinengenauigkeit weichen die damit berechneten Lösungen des AWP stark von der exakten Lösung des RWP ab. ♡

Aus der Abschätzung

$$\|\mathbf{y}(x; \mathbf{s}_1) - \mathbf{y}(x; \mathbf{s}_2)\| \leq \|\mathbf{s}_1 - \mathbf{s}_2\| e^{L|x-a|}$$

folgt aber auch, dass man den Einfluss der ungenauen Anfangsbedingungen durch Verkleinern des Intervalls $[a, x]$ beliebig verkleinert. Das führt auf die Mehrzielmethode.

7.3. Die Mehrzielmethode

Zum Lösen des RWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{r}(\mathbf{y}(a), \mathbf{y}(b)) = \mathbf{o}$$

gehen wir folgendermaßen vor:

Es sei $\Delta_n : \{x_1, x_2, \dots, x_m\}$ mit $a = x_1 < x_2 < \dots < x_m = b$ eine Unterteilung des Intervalls $[a, b]$. Wir wählen zu jedem $x_i, i = 1, \dots, m-1$, einen Anfangswert \mathbf{s}_i und lösen die AWP

$$\mathbf{y}' = \mathbf{f}(x, \mathbf{y}), \quad \mathbf{y}(x_i) = \mathbf{s}_i, \quad i = 1, \dots, m-1.$$

Es sei $\mathbf{y}(x; x_i, \mathbf{s}_i)$ die Lösungsfunktion, die wir für das Intervall $[x_i, x_{i+1})$ erhalten. Sind wir nun in der Lage, die Anfangswerte so zu wählen, dass die aus den Funktionen $\mathbf{y}(x; x_i, \mathbf{s}_i)$ stückweise zusammengesetzte Funktion stetig ist und noch die Randbedingungen erfüllt, so akzeptieren wir diese Funktion als Lösungsfunktion.

Es sei also $\mathbf{s} = (\mathbf{s}_1^T, \dots, \mathbf{s}_m^T)^T$ und

$$\mathbf{y}(x; \mathbf{s}) = \begin{cases} \mathbf{y}(x; x_i, \mathbf{s}_i) & \text{für } x \in [x_i, x_{i+1}), \quad i = 1, \dots, m-1 \\ \mathbf{s}_m & \text{für } x = x_m = b. \end{cases}$$

Die Stetigkeitsforderung an $\mathbf{y}(x; \mathbf{s})$ und die zu erfüllenden Randbedingungen führen auf ein nichtlineares Gleichungssystem für die Parameter $\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_m$:

$$\mathbf{F}(\mathbf{s}) = \begin{pmatrix} \mathbf{F}_1(\mathbf{s}_1, \mathbf{s}_2) \\ \mathbf{F}_2(\mathbf{s}_2, \mathbf{s}_3) \\ \vdots \\ \mathbf{F}_{m-1}(\mathbf{s}_{m-1}, \mathbf{s}_m) \\ \mathbf{F}_m(\mathbf{s}_1, \mathbf{s}_m) \end{pmatrix} = \begin{pmatrix} \mathbf{y}(x_2; x_1, \mathbf{s}_1) - \mathbf{s}_2 \\ \mathbf{y}(x_3; x_2, \mathbf{s}_2) - \mathbf{s}_3 \\ \vdots \\ \mathbf{y}(x_m; x_{m-1}, \mathbf{s}_{m-1}) - \mathbf{s}_m \\ \mathbf{r}(\mathbf{s}_1, \mathbf{s}_m) \end{pmatrix} = \mathbf{o}.$$

Zum Lösen wird wieder ein NEWTON-Verfahren der Form

$$\mathbf{s}^{(k+1)} = \mathbf{s}^{(k)} - \left[\mathbf{F}'(\mathbf{s}^{(k)}) \right]^{-1} \mathbf{F}(\mathbf{s}^{(k)})$$

oder ein modifiziertes NEWTON-Verfahren

$$\mathbf{s}^{(k+1)} = \mathbf{s}^{(k)} - \lambda_k \left[\mathbf{F}'(\mathbf{s}^{(k)}) \right]^{-1} \mathbf{F}(\mathbf{s}^{(k)})$$

angewendet. Man beachte, dass die \mathbf{s}_i Vektoren aus dem \mathbb{R}^n sind. Die Gleichung $\mathbf{F}(\mathbf{s}) = \mathbf{o}$ stellt damit ein nichtlineares Gleichungssystem mit $n \cdot m$ Gleichungen für

$n \cdot m$ Variable dar. Für die JACOBI-Matrix $F'(s)$ erhält man

$$F'(s) = \begin{pmatrix} G_1 & -I & O & \cdots & O & O & O \\ O & G_2 & -I & \ddots & O & O & O \\ O & O & G_3 & \ddots & O & O & O \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots & \vdots \\ O & O & O & \ddots & G_{m-2} & -I & O \\ O & O & O & \cdots & O & G_{m-1} & -I \\ A & O & O & \cdots & O & O & B \end{pmatrix}$$

mit

$$G_i = \left(\frac{\partial F_i(s)}{\partial s_i} \right), \quad A = \left(\frac{\partial r(s_1, s_m)}{\partial s_1} \right), \quad B = \left(\frac{\partial r(s_1, s_m)}{\partial s_m} \right).$$

Man ersetzt wieder die Differentialquotienten durch entsprechende Differenzenquotienten. Zu deren Berechnung ist wieder das Lösen von weiteren $(m-1)n$ AWP notwendig.

Die Berechnung von

$$\left[F'(s^{(k)}) \right]^{-1} F(s^{(k)}),$$

das bedeutet das Lösen des linearen Gleichungssystems

$$F'(s^{(k)})d = F(s^{(k)}),$$

lässt sich wegen der speziellen Struktur der JACOBI-Matrix auf die Behandlung von mehreren linearen Gleichungssystemen geringerer Dimension zurückführen.

Die Voraussetzungen an die Durchführbarkeit der Mehrzielmethode sind bedeutend schwächer als die Voraussetzungen an die Durchführbarkeit des einfachen Schießverfahrens. Besonders die Güte der Startwerte spielt bei der Mehrzielmethode eine geringere Rolle als beim einfachen Schießverfahren.

7.4. Differenzenverfahren

Eine weitere Möglichkeit Verfahren zum Lösen von RWP zu gewinnen, besteht darin, die in der Differentialgleichung auftretenden Differentialquotienten durch Differenzenquotienten zu ersetzen. Wir erläutern das Vorgehen am Beispiel des folgenden linearen RWP 2.Ordnung mit linearen separierten Randbedingungen.

$$\begin{aligned} -y'' + p(x)y' + q(x)y &= g(x), \quad x \in [a, b], \\ \alpha_1 y(a) + \alpha_2 y'(a) &= \alpha_0, \\ \beta_1 y(b) + \beta_2 y'(b) &= \beta_0. \end{aligned}$$

Weiterhin nehmen wir an, dass $y \in C^4[a, b]$ gilt.

Nun ersetzen wir in geeigneter Weise y'' und y' durch Differenzenquotienten. Dazu verwenden wir eine äquidistante Unterteilung des Intervalls $[a, b]$:

$$x_i = a + ih, \quad h = \frac{b-a}{N}.$$

Wir ersetzen $y''(x_i)$ durch

$$\frac{y(x_{i+1}) - 2y(x_i) + y(x_{i-1}))}{h^2}, \quad i = 1, \dots, N-1$$

und $y'(x_i)$ durch

$$\frac{y(x_{i+1}) - y(x_{i-1}))}{2h}, \quad i = 1, \dots, N-1.$$

Durch TAYLOR-Entwicklung von y bestätigt man leicht, dass

$$y''(x_i) = \frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} + O(h^2)$$

und

$$y'(x_i) = \frac{y_{i+1} - y_{i-1}}{2h} + O(h^2)$$

gilt. (Dabei ist $y_i = y(x_i)$ für $i = 0, 1, \dots, N$.) Setzen wir dies in die Differentialgleichung ein, so erhalten wir

$$-\frac{y_{i+1} - 2y_i + y_{i-1}}{h^2} + p(x_i) \frac{y_{i+1} - y_{i-1}}{2h} + q(x_i)y_i = g(x_i) + O(h^2), \quad i = 1, \dots, N-1.$$

In den Randbedingungen dürfen wir die ersten Ableitungen nicht durch zentrale Differenzenquotienten ersetzen. Wir verwenden hier einfache Differenzenquotienten

$$\begin{aligned} y'(x_0) &= y'(a) = \frac{y_1 - y_0}{h} + O(h), \\ y'(x_N) &= y'(b) = \frac{y_N - y_{N-1}}{h} + O(h). \end{aligned}$$

Damit ergibt sich

$$\begin{aligned} \alpha_1 y_0 + \alpha_2 \frac{y_1 - y_0}{h} &= \alpha_0 + O(h), \\ \beta_1 y_N + \beta_2 \frac{y_N - y_{N-1}}{h} &= \beta_0 + O(h). \end{aligned}$$

Ersetzen wir nun die y_i durch die Näherungen η_i und lassen die Restglieder in den Gleichungen fort, so ergibt sich zur Bestimmung der η_i das folgende lineare Gleichungssystem (mit $p_i = p(x_i)$, $q_i = q(x_i)$, $g_i = g(x_i)$).

$$\begin{aligned} (-\alpha_2 + h\alpha_1)\eta_0 + \alpha_2\eta_1 &= h\alpha_0, \\ \left(-1 - \frac{h}{2}p_i\right)\eta_{i-1} + \left(2 + h^2q_i\right)\eta_i + \left(-1 + \frac{h}{2}p_i\right)\eta_{i+1} &= h^2g_i, \quad i = 1, \dots, N-1, \\ -\beta_2\eta_{N-1} + (\beta_2 + h\beta_1)\eta_N &= h\beta_0. \end{aligned}$$

In Matrixschreibweise erhalten wir $\mathbf{A}(h)\boldsymbol{\eta} = \mathbf{b}(h)$ mit

$$\mathbf{A}(h) = \begin{pmatrix} -\alpha_2 + h\alpha_1 & \alpha_2 & 0 & \cdots & 0 & 0 \\ -1 - \frac{h}{2}p_1 & 2 + h^2q_1 & -1 + \frac{h}{2}p_1 & \ddots & 0 & 0 \\ 0 & -1 - \frac{h}{2}p_2 & 2 + h^2q_2 & \ddots & 0 & 0 \\ \vdots & & & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 2 + h^2q_{N-1} & -1 + \frac{h}{2}p_{N-1} \\ 0 & 0 & 0 & \cdots & -\beta_2 & \beta_2 + h\beta_1 \end{pmatrix}$$

und

$$\boldsymbol{\eta} = \begin{pmatrix} \eta_0 \\ \eta_1 \\ \vdots \\ \eta_{N-1} \\ \eta_N \end{pmatrix}, \quad \mathbf{b}(h) = \begin{pmatrix} h\alpha_0 \\ h^2g_1 \\ \vdots \\ h^2g_{N-1} \\ h\beta_0 \end{pmatrix}.$$

Wir erhalten genau dann eine eindeutige Lösung $\boldsymbol{\eta}$, wenn $\mathbf{A}(h)$ regulär ist. Die Regularität von $\mathbf{A}(h)$ ist natürlich in jedem Falle gesondert nachzuprüfen. Für bestimmte Parameterkombinationen lassen sich Aussagen über die Regularität von $\mathbf{A}(h)$ machen. Es gilt der folgende Satz (ohne Beweis).

7.9. Satz: Für das RWP

$$\begin{aligned} -y'' + p(x)y' + q(x)y &= g(x), \quad x \in [a, b], \\ \alpha_1 y(a) + \alpha_2 y'(a) &= \alpha_0, \\ \beta_1 y(b) + \beta_2 y'(b) &= \beta_0 \end{aligned}$$

gelte

1. $\alpha_1 > 0$, $\alpha_2 \leq 0$, $\beta_1 > 0$ und $\beta_2 \geq 0$,
2. $q(x) \geq 0$ für alle $x \in [a, b]$.

Dann ist die Matrix $\mathbf{A}(h)$ regulär für alle Schrittweiten h mit

$$0 < h < h_0 = \frac{2}{\sup_{x \in [a,b]} p(x)}.$$

Die Matrix $\mathbf{A}(h)$ wird für kleines h (großes N) fast singulär sein, so dass das numerische Lösen des Gleichungssystems $\mathbf{A}(h)\boldsymbol{\eta} = \mathbf{b}(h)$ schwierig ist.

Für die Güte der Näherungen lässt sich im Falle $\alpha_2 = \beta_2 = 0$ und $p(x) \equiv 0$ die folgende Aussage machen.

7.10. Satz: *Das RWP*

$$\begin{aligned} -y'' + q(x)y &= g(x), & x \in [a, b], \\ y(a) &= \alpha_0, \\ y(b) &= \beta_0 \end{aligned}$$

mit $q(x) \geq 0$ für alle $x \in [a, b]$ habe eine Lösung $y \in C^4[a, b]$. Es seien $|y^{(4)}| \leq M$ für alle $x \in [a, b]$ und $\boldsymbol{\eta}$ die mit dem Differenzenverfahren erzeugte Näherungslösung. Dann gilt

$$|y(x_i) - \eta_i| \leq \frac{Mh^2}{24}(x_i - a)(b - x_i).$$

Für dieses spezielle RWP erhalten wir so ein Verfahren 2. Ordnung.

Das obige Differenzenverfahren ist auch zur Behandlung von nichtlinearen RWP anwendbar. Dann erhält man jedoch zum Berechnen der Näherungen η_0, \dots, η_N nichtlineare Gleichungssysteme, die selbst mit Hilfe von Iterationsverfahren nur näherungsweise lösbar sind. Um mit Differenzenverfahren hohe Genauigkeiten zu erzielen, ist die Schrittweite h klein zu wählen. Das führt aber auf große Gleichungssysteme, die eventuell noch schlecht konditioniert sind. Differenzenverfahren spielen deshalb nur bei geringen Genauigkeitsanforderungen eine Rolle. Bei hohen Genauigkeitsanforderungen sollte man besser die Mehrzielmethode anwenden.

7.5. Aufgaben

1. Man löse das Randwertproblem

$$y'' - xy' + y = 1 \quad x \in [a, b] \quad y(0) = y'(1) = 1$$

numerisch mit einem Differenzenverfahren. Als Schrittweite wähle man $h = 0,2$. Anschließend bestimme man die exakte Lösung und vergleiche mit den erhaltenen Näherungen.

Kapitel 8

Lineare Gleichungssysteme

8.1. Allgemeine Grundlagen und Störungstheorie

Wir betrachten in diesem Kapitel lineare Gleichungssysteme

$$\begin{aligned} a_{11}x_1 + a_{12}x_2 + \cdots + a_{1n}x_n &= b_1, \\ a_{21}x_1 + a_{22}x_2 + \cdots + a_{2n}x_n &= b_2, \\ &\vdots \\ a_{n1}x_1 + a_{n2}x_2 + \cdots + a_{nn}x_n &= b_n, \end{aligned}$$

oder in Matrixschreibweise

$$\mathbf{Ax} = \mathbf{b}$$

mit

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^n.$$

Bevor wir uns mit der Existenz und Eindeutigkeit von Lösungen und Genauigkeitsfragen beschäftigen, sind noch einige Begriffe bereitzustellen.

8.1.1. Vektor- und Matrixnormen

Eine Abbildung

$$\|\circ\| : \mathbb{R}^n \longrightarrow \mathbb{R}_+ = \left\{ \alpha \in \mathbb{R} \mid \alpha \geq 0 \right\}$$

heißt **Vektornorm**, falls sie folgende Bedingungen erfüllt.

1. $\|\mathbf{x}\| \geq 0$, $\|\mathbf{x}\| = 0 \iff \mathbf{x} = \mathbf{o}$.
2. Für alle $\alpha \in \mathbb{R}$ und alle $\mathbf{x} \in \mathbb{R}^n$ gilt $\|\alpha\mathbf{x}\| = |\alpha| \|\mathbf{x}\|$.
3. Für alle $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ gilt die Dreiecksungleichung $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$.

Bemerkung: Die Dreiecksungleichung ist äquivalent zur Ungleichung

$$\|\mathbf{x} - \mathbf{y}\| \geq \left| \|\mathbf{x}\| - \|\mathbf{y}\| \right|.$$

Beispiele für Vektornormen

- Euklidische Norm:

$$\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}} = \sqrt{\sum_{i=1}^n x_i^2}.$$

- Betragssummennorm:

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|.$$

- Maximumnorm:

$$\|\mathbf{x}\|_\infty = \max_{i=1, \dots, n} |x_i|.$$

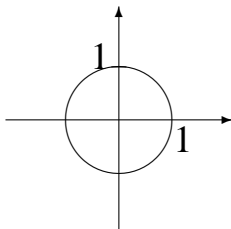
- p -Norm:

$$\|\mathbf{x}\|_p = \sqrt[p]{\sum_{i=1}^n |x_i|^p}.$$

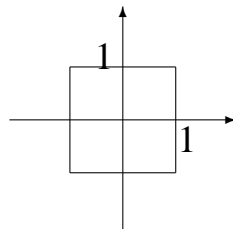
In den folgenden Bildern sind die sogenannten Einheitskugeln, die Mengen

$$\{ \mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| = 1 \},$$

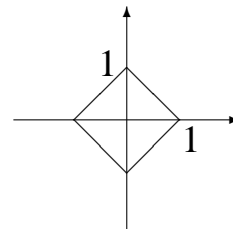
im \mathbb{R}^2 dargestellt.



$$\|\mathbf{x}\|_2 = 1$$



$$\|\mathbf{x}\|_\infty = 1$$



$$\|\mathbf{x}\|_1 = 1$$

Zwischen den Vektornormen gibt es gewisse Beziehungen. Die wichtigste kommt in folgendem Satz zum Ausdruck.

8.1. Satz: *Alle Vektornormen im \mathbb{R}^n sind in folgendem Sinne äquivalent: Für je zwei Normen $\|\circ\|^{(1)}$ und $\|\circ\|^{(2)}$ gibt es Konstanten $M \geq m > 0$, so dass für alle $\mathbf{x} \in \mathbb{R}^n$*

$$m\|\mathbf{x}\|^{(2)} \leq \|\mathbf{x}\|^{(1)} \leq M\|\mathbf{x}\|^{(2)}$$

gilt.

So gilt zum Beispiel

$$\begin{aligned} \|\mathbf{x}\|_\infty &\leq \|\mathbf{x}\|_2 \leq \sqrt{n}\|\mathbf{x}\|_\infty, \\ \frac{1}{\sqrt{n}}\|\mathbf{x}\|_1 &\leq \|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1, \\ \|\mathbf{x}\|_\infty &\leq \|\mathbf{x}\|_1 \leq n\|\mathbf{x}\|_\infty. \end{aligned}$$

Auch für Matrizen lassen sich Normen einführen. Eine Abbildung

$$\|\circ\| : \mathbb{R}^{m \times n} \longrightarrow \mathbb{R}_+$$

heißt (m, n) -**Matrixnorm**, falls sie folgende Bedingungen erfüllt:

1. $\|\mathbf{A}\| = 0 \iff \mathbf{A} = \mathbf{O}$.
2. Für alle $\alpha \in \mathbb{R}$ und alle (m, n) -Matrizen \mathbf{A} gilt

$$\|\alpha\mathbf{A}\| = |\alpha|\|\mathbf{A}\|.$$

3. Für alle (m, n) -Matrizen \mathbf{A}, \mathbf{B} gilt die Dreiecksungleichung

$$\|\mathbf{A} + \mathbf{B}\| \leq \|\mathbf{A}\| + \|\mathbf{B}\|.$$

Bemerkung: Die Dreiecksungleichung lässt sich wieder durch

$$\|\mathbf{A} - \mathbf{B}\| \geq \left| \|\mathbf{A}\| - \|\mathbf{B}\| \right|$$

ersetzen.

Bei Fehlerabschätzungen für lineare Gleichungssysteme treten oft Matrix- und Vektornormen gemeinsam auf. Hier ist darauf zu achten, dass die verwendeten Normen in gewisser Weise verträglich sind. Diese Eigenschaft wird in der folgenden Definition genauer beschrieben.

Eine (m, n) -Matrixnorm $\|\circ\|$ heißt mit den Vektornormen $\|\circ\|^{(m)}$ auf dem \mathbb{R}^m und $\|\circ\|^{(n)}$ auf dem \mathbb{R}^n **verträglich**, falls für alle $\mathbf{x} \in \mathbb{R}^n$ und alle (m, n) -Matrizen \mathbf{A}

$$\|\mathbf{Ax}\|^{(m)} \leq \|\mathbf{A}\| \cdot \|\mathbf{x}\|^{(n)}$$

gilt.

Die (n, n) -Matrixnorm $\|\circ\|$ heißt **submultiplikativ**, falls für alle (n, n) -Matrizen \mathbf{A}, \mathbf{B}

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{B}\|$$

gilt. Auch für Matrixnormen gilt ein zu Satz 8.1 analoger Satz.

8.2. Satz: *Alle (m, n) -Matrixnormen sind in folgendem Sinne äquivalent:*

Für je zwei Normen $\|\circ\|^{(1)}$ und $\|\circ\|^{(2)}$ gibt es Konstanten $L \geq l > 0$, so dass für alle (m, n) -Matrizen \mathbf{A}

$$l\|\mathbf{A}\|^{(2)} \leq \|\mathbf{A}\|^{(1)} \leq L\|\mathbf{A}\|^{(2)}$$

gilt.

Beispiele für Matrixnormen

- Zeilensummennorm:

$$\|\mathbf{A}\|_{\infty} = \max_{i=1, \dots, m} \sum_{k=1}^n |a_{ik}|.$$

- Spaltensummennorm:

$$\|\mathbf{A}\|_1 = \max_{j=1, \dots, n} \sum_{k=1}^m |a_{kj}|.$$

- FROBENIUS-Norm:

$$\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2}.$$

- Maximumnorm:

$$\|\mathbf{A}\|_M = \max_{i=1, \dots, m} \max_{j=1, \dots, n} |a_{ij}|.$$

Die Zeilensummennorm, die Spaltensummennorm und die FROBENIUS-Norm sind für quadratische Matrizen submultiplikativ. Die Zeilensummennorm ist mit der Maximumnorm für Vektoren verträglich, die Spaltensummennorm ist mit der Betragsnorm für Vektoren verträglich, und die FROBENIUS-Norm ist mit der euklidischen Vektornorm verträglich.

Unter allen (n, n) -Matrixnormen, die mit einer gegebenen Vektornorm $\|\circ\|$ auf dem \mathbb{R}^n verträglich sind, gibt es genau eine, die dieser Vektornorm besonders gut angepasst ist. Die Matrixnorm

$$\text{lub}(\mathbf{A}) = \max_{\substack{\mathbf{x} \in \mathbb{R}^n \\ \mathbf{x} \neq \mathbf{o}}} \frac{\|\mathbf{Ax}\|}{\|\mathbf{x}\|} = \max_{\substack{\mathbf{y} \in \mathbb{R}^n \\ \|\mathbf{y}\|=1}} \|\mathbf{Ay}\|$$

heißt **Matrixgrenznorm**¹ oder kurz Grenznorm zur Vektornorm $\|\circ\|$. Die besondere Eigenschaft von Grenznormen kommt in folgendem Satz zum Ausdruck.

8.3. Satz: *Es sei $\|\circ\|$ eine beliebige (n, n) -Matrixnorm, die mit der Vektornorm $\|\circ\|$ auf dem \mathbb{R}^n verträglich ist. Mit der zugehörigen Matrixgrenznorm gilt dann für alle (n, n) -Matrizen \mathbf{A}*

$$\text{lub}(\mathbf{A}) \leq \|\mathbf{A}\|.$$

Der Beweis folgt sofort aus der Definition der Matrixgrenznorm.

Weiterhin ist jede Matrixgrenznorm submultiplikativ und es gilt $\text{lub}(\mathbf{I}) = 1$.

Geometrische Interpretation der Matrixgrenznorm

Es sei φ die durch die Matrix \mathbf{A} vermittelte lineare Abbildung des \mathbb{R}^n in sich. Wir betrachten nun das Bild einer Einheitskugel:

$$\Omega = \{ \mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\| = 1 \}$$

und

$$\varphi(\Omega) = \{ \mathbf{Ax} \mid \mathbf{x} \in \Omega \}.$$

Dann folgt

$$\text{lub}(\mathbf{A}) = \max_{\mathbf{y} \in \varphi(\Omega)} \{\|\mathbf{y}\|\}.$$

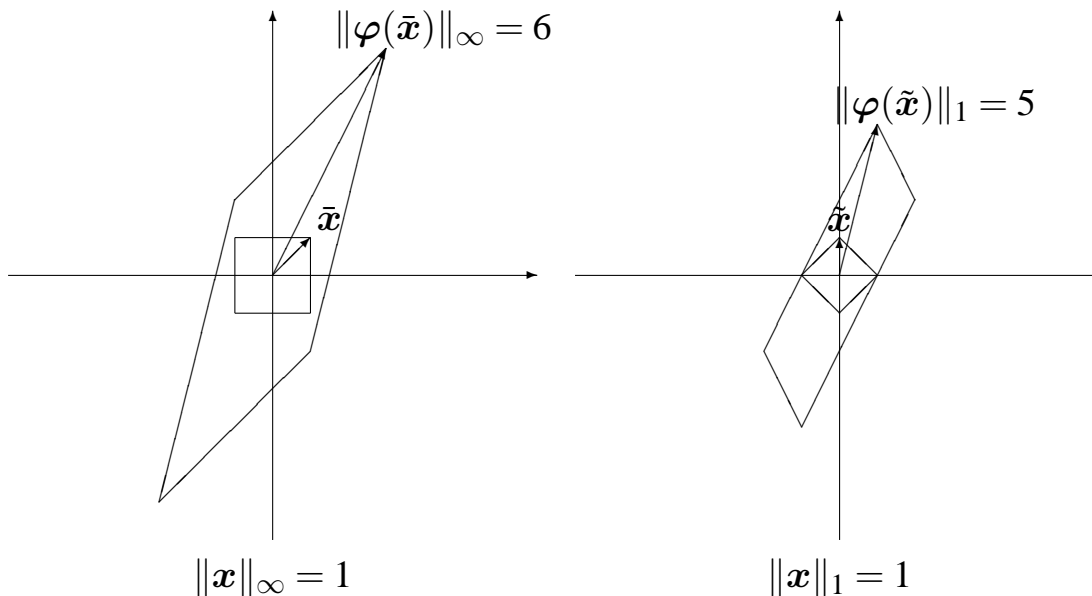
Die Größe $\text{lub}(\mathbf{A})$ stellt die größte „Verzerrung“ eines Bildpunktes gegenüber dem Originalpunkt dar.

¹Die Bezeichnung lub ist die Abkürzung von **l**owest **u**pper **b**ound.

8.4. Beispiel: Wir betrachten die Matrix

$$\mathbf{A} = \begin{pmatrix} 1 & 2 \\ 4 & 2 \end{pmatrix}.$$

Im Bild sind die Einheitssphären und ihre Bilder für die Maximumnorm und die Betragssummennorm dargestellt.



Zu den wichtigsten Vektornormen erhält man die folgenden Grenznormen.

- Betragssummennorm:

$$\|\mathbf{A}\|_1 = \text{lub}_1(\mathbf{A}) = \max_{\mathbf{x} \neq \mathbf{o}} \frac{\|\mathbf{A}\mathbf{x}\|_1}{\|\mathbf{x}\|_1} = \max_{j=1, \dots, n} \sum_{k=1}^n |a_{kj}|.$$

Das ist die Spaltensummennorm.

- Maximumnorm:

$$\|\mathbf{A}\|_\infty = \text{lub}_\infty(\mathbf{A}) = \max_{\mathbf{x} \neq \mathbf{o}} \frac{\|\mathbf{A}\mathbf{x}\|_\infty}{\|\mathbf{x}\|_\infty} = \max_{i=1, \dots, n} \sum_{k=1}^n |a_{ik}|.$$

Das ist die Zeilensummennorm.

- Euklidische Norm:

$$\|\mathbf{A}\|_2 = \text{lub}_2(\mathbf{A}) = \max_{\mathbf{x} \neq \mathbf{o}} \frac{\|\mathbf{A}\mathbf{x}\|_2}{\|\mathbf{x}\|_2} = \max_{\mathbf{x} \neq \mathbf{o}} \frac{\sqrt{(\mathbf{A}\mathbf{x})^T (\mathbf{A}\mathbf{x})}}{\sqrt{\mathbf{x}^T \mathbf{x}}} = \sqrt{\lambda_{\max}(\mathbf{A}^T \mathbf{A})}.$$

Das ist die Spektralnorm. $\lambda_{max}(\mathbf{A}^T \mathbf{A})$ bezeichnet dabei den größten Eigenwert der Matrix $\mathbf{A}^T \mathbf{A}$. Die Spektralnorm einer Matrix darf nicht mit dem Spektralradius verwechselt werden. Dieser ist der Betrag des betragsgrößten Eigenwertes von \mathbf{A} :

$$\varrho(\mathbf{A}) = \max_{i=1, \dots, n} \{|\lambda_i(\mathbf{A})|\}.$$

8.1.2. Ordnungen und Beträge

Die in Abschnitt 8.1.1. eingeführten Normen werden dazu verwendet, Vektoren und Matrizen in ihrer Gesamtheit miteinander zu vergleichen. Oft benötigt man aber auch elementweise Abschätzungen. Dazu führen wir eine natürliche Halbordnung für Vektoren aus dem \mathbb{R}^n und für (m, n) -Matrizen ein. Für je zwei Vektoren $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ gilt $\mathbf{x} \leq \mathbf{y}$ genau dann, wenn $x_i \leq y_i$ für $i = 1, \dots, n$ gilt. Für je zwei (m, n) -Matrizen \mathbf{A}, \mathbf{B} gilt $\mathbf{A} \leq \mathbf{B}$ genau dann, wenn $a_{ij} \leq b_{ij}$ für $i = 1, \dots, m$ und $j = 1, \dots, n$ gilt. Die so eingeführte Relation ist reflexiv, symmetrisch und transitiv. Aber es gibt Vektoren (Matrizen) für die weder $\mathbf{x} \leq \mathbf{y}$ ($\mathbf{A} \leq \mathbf{B}$) noch $\mathbf{y} \leq \mathbf{x}$ ($\mathbf{B} \leq \mathbf{A}$) gilt. Weiterhin führen wir Betragsvektoren und Betragsmatrizen ein. Für $\mathbf{x} \in \mathbb{R}^n$ sei

$$|\mathbf{x}| = \begin{pmatrix} |x_1| \\ |x_2| \\ \vdots \\ |x_n| \end{pmatrix}.$$

Für eine (m, n) -Matrix \mathbf{A} sei

$$|\mathbf{A}| = \begin{pmatrix} |a_{11}| & |a_{12}| & \dots & |a_{1n}| \\ |a_{21}| & |a_{22}| & \dots & |a_{2n}| \\ \vdots & \vdots & \ddots & \vdots \\ |a_{m1}| & |a_{m2}| & \dots & |a_{mn}| \end{pmatrix}.$$

Offensichtlich gilt für alle $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

- $|\mathbf{x}| = \mathbf{0} \iff \mathbf{x} = \mathbf{o}$.
- $|\alpha \mathbf{x}| = |\alpha| |\mathbf{x}|$ für alle $\alpha \in \mathbb{R}$.
- $|\mathbf{x} + \mathbf{y}| \leq |\mathbf{x}| + |\mathbf{y}|$.

Achtung: $|\circ|$ ist trotzdem keine Norm, denn $|\circ| : \mathbb{R}^n \rightarrow \mathbb{R}^n$.

Für beliebige (m, n) -Matrizen \mathbf{A}, \mathbf{B} und (n, p) -Matrizen \mathbf{C} gilt

- $|\mathbf{A}| = \mathbf{0} \iff \mathbf{A} = \mathbf{O}$.

- $|\alpha \mathbf{A}| = |\alpha| |\mathbf{A}|$ für alle $\alpha \in \mathbb{R}$.
- $|\mathbf{A} + \mathbf{B}| \leq |\mathbf{A}| + |\mathbf{B}|$.
- $|\mathbf{AC}| \leq |\mathbf{A}| |\mathbf{C}|$.

Insbesondere gilt für alle (m, n) -Matrizen \mathbf{A} und alle Vektoren $\mathbf{x} \in \mathbb{R}^n$

$$|\mathbf{Ax}| \leq |\mathbf{A}| |\mathbf{x}|.$$

Die Ordnungsrelationen und Beträge sind mit den Vektornormen $\|\circ\|_p$ in folgendem Sinne verträglich.

(*) Für alle Vektoren $\mathbf{x} \in \mathbb{R}^n$ gilt $\|\|\mathbf{x}\|\|_p = \|\mathbf{x}\|_p$.

(**) Für alle Vektoren $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ gilt

$$|\mathbf{x}| \leq |\mathbf{y}| \implies \|\mathbf{x}\|_p \leq \|\mathbf{y}\|_p.$$

Eine Norm mit der Eigenschaft (*) heißt **absolut**.

Eine Norm mit der Eigenschaft (**) heißt **monoton**. Die Vektornormen $\|\circ\|_p$ sind absolut und monoton. Bei den Matrixnormen sind die Normen $\|\mathbf{A}\|_1$, $\|\mathbf{A}\|_\infty$ und $\|\mathbf{A}\|_F$ absolut und monoton. Die Spektralnorm $\|\mathbf{A}\|_2$ ist weder absolut noch monoton. Es gilt jedoch

$$\|\mathbf{A}\|_2 \leq \|\|\mathbf{A}\|\|_2$$

und

$$|\mathbf{A}| \leq |\mathbf{B}| \implies \|\|\mathbf{A}\|\|_2 \leq \|\|\mathbf{B}\|\|_2.$$

Die Spektralnorm besitzt dafür (genau wie die Euklidische Vektornorm) eine andere wichtige Eigenschaft:

Für eine beliebige (m, n) -Matrix \mathbf{A} und beliebige orthogonale (m, m) -Matrizen \mathbf{U} und (n, n) -Matrizen \mathbf{V} gilt

$$\|\mathbf{UAV}\|_2 = \|\mathbf{A}\|_2.$$

Wir werden im weiteren nur submultiplikative Matrixnormen verwenden. Im wesentlichen werden das die Normen $\|\mathbf{A}\|_1$, $\|\mathbf{A}\|_2$, $\|\mathbf{A}\|_\infty$ und $\|\mathbf{A}\|_F$ sein.

8.1.3. Spezielle Transformationsmatrizen

Bei der Entwicklung von Verfahren zum Lösen von linearen Gleichungssystemen werden wir verschiedene Transformationsmatrizen verwenden. Die wichtigsten von ihnen mit ihren speziellen Eigenschaften wollen wir in diesem Abschnitt angeben.

Nichtorthogonale Transformationen

Eine (n, n) -Matrix M der Form

$$M = M_k(m) = I + m e_k^T$$

mit $m \in \mathbb{R}^n$ und $e_k^T m = 0$ heißt **NT-Matrix**². Eine (n, n) -Matrix L der Form

$$L = L_k(l) = I + l e_k^T$$

mit $l \in \mathbb{R}^n$ und $e_i^T l = 0$ für $i = 1, \dots, k$ heißt **LNT-Matrix**³. Eine NT-Matrix hat damit folgendes Aussehen:

$$M_k(m) = \begin{pmatrix} 1 & & m_1 & & & \\ & \ddots & \vdots & & & \\ & & 1 & m_{k-1} & & \\ & & & 1 & & \\ & & & m_{k+1} & 1 & \\ \mathbf{O} & & \vdots & & \ddots & \\ & & m_n & & & 1 \end{pmatrix}.$$

Eine LNT-Matrix hat dagegen das Aussehen:

$$L_k(l) = \begin{pmatrix} 1 & & & & & \\ & \ddots & & & & \\ & & 1 & & & \\ & & & 1 & & \\ & & & l_{k+1} & 1 & \\ \mathbf{O} & & \vdots & & \ddots & \\ & & l_n & & & 1 \end{pmatrix}.$$

Bemerkung: Eine Matrix der Form $I + m e_k^T$ mit beliebigem Vektor $m \in \mathbb{R}^n$ heißt **FROBENIUS-Matrix**.

Für NT-Matrizen und damit auch für LNT-Matrizen gilt der folgende Satz.

²Abkürzung für nichtorthogonale Transformationsmatrix

³Abkürzung für lower triangular NT-Matrix

8.5. Satz: Jede NT-Matrix $M_k(\mathbf{m})$ ist regulär, und es gilt

$$M_k(\mathbf{m})^{-1} = M_k(-\mathbf{m}).$$

Beweis:

$$\begin{aligned} M_k(\mathbf{m})M_k(-\mathbf{m}) &= (I + \mathbf{m}e_k^T)(I - \mathbf{m}e_k^T) \\ &= I + \mathbf{m}e_k^T - \mathbf{m}e_k^T - \mathbf{m}e_k^T\mathbf{m}e_k^T \\ &= I - (e_k^T\mathbf{m})\mathbf{m}e_k^T = I. \end{aligned}$$

✱

Die Inverse einer NT-Matrix ist wieder eine NT-Matrix. Speziell ist die Inverse einer LNT-Matrix wieder eine LNT-Matrix.

Eine (n, n) -Matrix L mit $l_{ij} = 0$ für $i < j$ heißt **untere Dreiecksmatrix**. Gilt zusätzlich $l_{ii} = 1$ für $i = 1, \dots, n$, so heißt die Matrix **untere Einsdreiecksmatrix**.

Diese Matrizen haben daher folgende Struktur:

$$L = \begin{pmatrix} l_{11} & & & & \\ l_{21} & l_{22} & & & \mathbf{O} \\ \vdots & \vdots & \ddots & & \\ l_{n-1,1} & l_{n-1,2} & \cdots & l_{n-1,n-1} & \\ l_{n,1} & l_{n,2} & \cdots & l_{n,n-1} & l_{nn} \end{pmatrix}$$

bzw.

$$L = \begin{pmatrix} 1 & & & & \\ l_{21} & 1 & & & \mathbf{O} \\ \vdots & \vdots & \ddots & & \\ l_{n-1,1} & l_{n-1,2} & \cdots & 1 & \\ l_{n,1} & l_{n,2} & \cdots & l_{n,n-1} & 1 \end{pmatrix}.$$

Eine (n, n) -Matrix U mit $u_{ij} = 0$ für $i > j$ heißt **obere Dreiecksmatrix**. Gilt zusätzlich $u_{ii} = 1$ für $i = 1, \dots, n$, so heißt die Matrix **obere Einsdreiecksmatrix**.

Eine obere Dreiecksmatrix hat folgende Struktur

$$U = \begin{pmatrix} u_{11} & u_{12} & \cdots & u_{1,n-1} & u_{1n} \\ & u_{22} & \cdots & u_{2,n-1} & u_{2n} \\ & & \ddots & \vdots & \vdots \\ & \mathbf{O} & & u_{n-1,n-1} & u_{n-1,n} \\ & & & & u_{nn} \end{pmatrix}.$$

$(\mathbf{G})_{pp}$, $(\mathbf{G})_{pq}$, $(\mathbf{G})_{qp}$ und $(\mathbf{G})_{qq}$. Die Orthogonalität ist offensichtlich. Geometrisch stellt die Transformation

$$\mathbf{x} \longrightarrow \mathbf{y} = \mathbf{G}_{pq}(c, s)\mathbf{x}$$

eine Drehung in der (pq) -Ebene um den Winkel $-\varphi$ mit $c = \cos \varphi$ und $s = \sin \varphi$ dar. Von großer Bedeutung sind auch die HOUSEHOLDER-Spiegelungen. Eine (n, n) -Matrix $\mathbf{H} = \mathbf{H}(\mathbf{u})$ der Form

$$\mathbf{H} = \mathbf{I} - 2\mathbf{u}\mathbf{u}^T, \quad \mathbf{u} \in \mathbb{R}^n, \quad \mathbf{u}^T\mathbf{u} = 1$$

heißt HOUSEHOLDER-**Spiegelung** oder HOUSEHOLDER-**Matrix**. Offensichtlich ist \mathbf{H} symmetrisch. Auch die Orthogonalität erkennt man leicht. Es gilt

$$\begin{aligned} \mathbf{H}\mathbf{H}^T &= \mathbf{H}\mathbf{H} = (\mathbf{I} - 2\mathbf{u}\mathbf{u}^T)(\mathbf{I} - 2\mathbf{u}\mathbf{u}^T) \\ &= \mathbf{I} - 2\mathbf{u}\mathbf{u}^T - 2\mathbf{u}\mathbf{u}^T + 4\mathbf{u}\mathbf{u}^T\mathbf{u}\mathbf{u}^T = \mathbf{I} - 4\mathbf{u}\mathbf{u}^T + 4\mathbf{u}\mathbf{u}^T = \mathbf{I}. \end{aligned}$$

Geometrisch lässt sich die Abbildung

$$\mathbf{x} \longrightarrow \mathbf{y} = \mathbf{H}(\mathbf{u})\mathbf{x}$$

als Spiegelung des Vektors \mathbf{x} an einer Ebene mit dem Normalenvektor \mathbf{u} deuten. Verzichtet man auf die Normierung $\mathbf{u}^T\mathbf{u} = 1$, so stellt sich eine HOUSEHOLDER-Matrix in der Form

$$\mathbf{H} = \mathbf{I} - \frac{\mathbf{v}\mathbf{v}^T}{\gamma}, \quad \gamma = \frac{\mathbf{v}^T\mathbf{v}}{2}$$

dar.

8.1.4. Eigenwerte und Singulärwerte

Bekanntlich lassen sich für quadratische Matrizen sogenannte Eigenwerte und Eigenvektoren definieren. Eine Zahl $\lambda \in \mathbb{C}$ heißt **Eigenwert** der (n, n) -Matrix \mathbf{A} , falls es einen Vektor $\mathbf{x} \in \mathbb{C}^n$, $\mathbf{x} \neq \mathbf{o}$, gibt, so dass

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$$

gilt. Jeder derartige Vektor \mathbf{x} heißt **Eigenvektor** zum Eigenwert λ . Eine notwendige und hinreichende Bedingung dafür, dass λ Eigenwert von \mathbf{A} ist, ist die Säkulargleichung

$$\det(\mathbf{A} - \lambda\mathbf{I}) = 0.$$

Das Polynom

$$\varphi(\mu) = \det(\mathbf{A} - \mu\mathbf{I})$$

heißt **charakteristisches Polynom** der (n, n) -Matrix \mathbf{A} . Es gilt $\varphi \in \Pi_n$. Die Eigenwerte einer Matrix sind damit die Nullstellen des charakteristischen Polynoms. Im allgemeinen sind die Eigenwerte und damit auch die Eigenvektoren einer beliebigen reellen Matrix komplex. Es gilt aber der folgende Satz.

8.6. Satz: *Die Eigenwerte einer symmetrischen (n, n) -Matrix \mathbf{A} sind sämtlich reell. Zu verschiedenen Eigenwerten gehörende Eigenvektoren sind paarweise orthogonal. Es existiert eine orthogonale (n, n) -Matrix \mathbf{U} , so dass*

$$\mathbf{U}^T \mathbf{A} \mathbf{U} = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

gilt. $\lambda_1, \dots, \lambda_n$ sind dabei die Eigenwerte von \mathbf{A} , die Spalten von \mathbf{U} sind zugehörige normierte Eigenvektoren.

Eine Abschätzung der Eigenwerte einer Matrix liefert der folgende Satz.

8.7. Satz: *Für alle Eigenwerte λ einer (n, n) -Matrix \mathbf{A} gilt*

$$|\lambda| \leq \text{lub}(\mathbf{A})$$

bezüglich einer beliebigen Matrixnorm.

Beweis: Ist $\mathbf{x} \neq \mathbf{o}$ Eigenvektor zum Eigenwert λ , so gilt $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$. Daraus folgt

$$|\lambda| \|\mathbf{x}\| = \|\mathbf{A}\mathbf{x}\| \leq \text{lub}(\mathbf{A}) \|\mathbf{x}\|,$$

und nach Division durch $\|\mathbf{x}\|$

$$|\lambda| \leq \text{lub}(\mathbf{A}).$$

✱

Bemerkung: Ist $\|\mathbf{A}\|$ eine Matrixnorm, die mit irgendeiner Vektornorm verträglich ist, so folgt aus diesem Satz und Satz 8.3 sofort $|\lambda| \leq \|\mathbf{A}\|$ für einen beliebigen Eigenwert λ von \mathbf{A} .

Beliebige Matrizen lassen sich durch orthogonale Transformationen nicht diagonalisieren. Für sie gilt aber der folgende Satz.

8.8. Singulärwertzerlegung: Zu jeder (m, n) -Matrix \mathbf{A} gibt es eine orthogonale (m, m) -Matrix \mathbf{U} und eine orthogonale (n, n) -Matrix \mathbf{V} , so dass

$$\mathbf{U}^T \mathbf{A} \mathbf{V} = \Sigma = \text{diag}(\sigma_1, \dots, \sigma_l)$$

mit $l = \min\{m, n\}$ und $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_l \geq 0$ gilt. Die Zahlen $\sigma_1, \dots, \sigma_l$ sind eindeutig bestimmt und heißen Singulärwerte der Matrix \mathbf{A} .

Bemerkung: Die Matrix Σ hat folgende Gestalt

$$\Sigma = \begin{pmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \\ & & & \mathbf{O} \end{pmatrix}, \quad m > n = l$$

$$\Sigma = \begin{pmatrix} \sigma_1 & & 0 \\ & \ddots & & \mathbf{O} \\ 0 & & \sigma_m & \\ & & & \end{pmatrix}, \quad n > m = l$$

$$\Sigma = \begin{pmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_n \end{pmatrix}, \quad m = n = l$$

Beweis: Die (n, n) -Matrix $\mathbf{A}^T \mathbf{A}$ ist symmetrisch und wegen

$$\mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x} = (\mathbf{A} \mathbf{x})^T (\mathbf{A} \mathbf{x}) = \|\mathbf{A} \mathbf{x}\|_2^2 \geq 0$$

positiv semidefinit. Nach Satz 8.6 existiert damit eine orthogonale (n, n) -Matrix \mathbf{V} , so dass

$$\mathbf{V}^T (\mathbf{A}^T \mathbf{A}) \mathbf{V} = \text{diag}(\lambda_1, \dots, \lambda_n) \geq \mathbf{O}$$

gilt. O.B.d.A. sei

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_r > \lambda_{r+1} = \dots = \lambda_n = 0$$

angenommen. Weiterhin sei

$$\mathbf{V} = (\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n)}).$$

Dann gilt

$$\mathbf{A}^T \mathbf{A} \mathbf{v}^{(j)} = \lambda_j \mathbf{v}^{(j)}, \quad j = 1, \dots, n,$$

und damit

$$\lambda_j = \mathbf{v}^{(j)T} \mathbf{A}^T \mathbf{A} \mathbf{v}^{(j)} = \|\mathbf{A} \mathbf{v}^{(j)}\|_2^2.$$

Es ist also $\mathbf{A} \mathbf{v}^{(j)} \neq \mathbf{o}$ für $j = 1, \dots, r$ und $\mathbf{A} \mathbf{v}^{(j)} = \mathbf{o}$ für $j = r+1, \dots, n$. Wir definieren nun

$$\sigma_j = \sqrt{\lambda_j}, \quad j = 1, \dots, r$$

und

$$\mathbf{u}^{(j)} = \frac{\mathbf{A} \mathbf{v}^{(j)}}{\sigma_j}, \quad j = 1, \dots, r.$$

Für $i, j = 1, \dots, r$ gilt dann

$$\mathbf{u}^{(i)T} \mathbf{u}^{(j)} = \frac{\mathbf{v}^{(i)T} \mathbf{A}^T \mathbf{A} \mathbf{v}^{(j)}}{\sigma_i \sigma_j} = \frac{\lambda_j}{\sigma_i \sigma_j} \mathbf{v}^{(i)T} \mathbf{v}^{(j)} = \delta_{ij}.$$

Die Vektoren $\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(r)}$ sind paarweise orthogonal und auf 1 normiert. Durch Hinzunahme weiterer Vektoren $\mathbf{u}^{(r+1)}, \dots, \mathbf{u}^{(m)}$ sind sie zu einer orthonormierten Basis des \mathbb{R}^m ergänzt. Dann ist die Matrix

$$\mathbf{U} = \left(\mathbf{u}^{(1)}, \dots, \mathbf{u}^{(m)} \right)$$

orthogonal, und es gilt

$$\left(\mathbf{U}^T \mathbf{A} \mathbf{V} \right)_{ij} = \mathbf{u}^{(i)T} \mathbf{A} \mathbf{v}^{(j)} = \begin{cases} \sigma_j \mathbf{u}^{(i)T} \mathbf{u}^{(j)} & \text{für } j = 1, \dots, r, \\ 0 & \text{für } j = r+1, \dots, n. \end{cases}$$

Setzen wir noch $\sigma_{r+1} = \dots = \sigma_l = 0$, so gilt

$$\mathbf{U}^T \mathbf{A} \mathbf{V} = \mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_l).$$

✱

Die Singulärwerte einer Matrix \mathbf{A} sind somit die Wurzeln aus den Eigenwerten der Matrix $\mathbf{A}^T \mathbf{A}$ (oder $\mathbf{A} \mathbf{A}^T$). Weiterhin gilt

- $\|\mathbf{A}\|_2 = \sigma_1,$
- $\|\mathbf{A}\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2},$

- $\text{rg}(\mathbf{A}) = r$.

Für den Zusammenhang zwischen Eigenwerten und Singulärwerten einer quadratischen Matrix gilt

- $\det(\mathbf{A}) = \prod_{i=1}^n \lambda_i$ und $|\det(\mathbf{A})| = \prod_{i=1}^n \sigma_i$,
- $\sigma_1 \geq |\lambda_i| \geq \sigma_n$ für $i = 1, \dots, n$,
- $\sigma_i = |\lambda_i|$ für $i = 1, \dots, n$ falls \mathbf{A} symmetrisch ist.

8.1.5. Störungstheorie

Wir betrachten ein lineares Gleichungssystem $\mathbf{Ax} = \mathbf{b}$ mit einer (n, n) -Matrix \mathbf{A} und $\mathbf{b} \in \mathbb{R}^n$. Bekanntlich gilt dann der folgende Satz.

8.9. Satz: *Das lineare Gleichungssystem $\mathbf{Ax} = \mathbf{b}$ ist genau dann eindeutig lösbar, wenn die Matrix \mathbf{A} regulär ist.*

Dieser Satz ist in der Praxis mit Vorsicht zu genießen, wie das folgende Beispiel zeigt.

8.10. Beispiel: Es sei

$$\mathbf{A} = \begin{pmatrix} 1.00 & 0.99 \\ 0.99 & 0.98 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 1.99 \\ 1.97 \end{pmatrix}.$$

Die exakte Lösung des Gleichungssystems $\mathbf{Ax} = \mathbf{b}$ ist $\mathbf{x} = (1, 1)^T$. Die Matrix \mathbf{A} ist regulär, denn es gilt $\det(\mathbf{A}) = 1 \cdot 0.98 - 0.99^2 = -0.0001 \neq 0$. Behandelt man dieses Gleichungssystem auf einem Rechner mit dem Maschinenzahlbereich $\mathbb{M}(10, 2, \dots)$, so ergibt sich

$$\text{gl}(\det(\mathbf{A})) = \text{gl}(0.98 - 0.99^2) = 0.98 - 0.98 = 0.$$

Auf diesem Rechner würde die Matrix als singulär betrachtet werden. ♡

Der Begriff der Regularität der Koeffizientenmatrix reicht daher in der Numerik nicht aus, um die (praktische) Lösbarkeit eines Gleichungssystems zu charakterisieren. Eine (n, n) -Matrix \mathbf{A} heißt **numerisch regulär**, falls alle (n, n) -Matrizen $\tilde{\mathbf{A}}$ aus einer Umgebung $U(\mathbf{A})$ regulär sind. Existiert eine singuläre Matrix $\tilde{\mathbf{A}} \in U(\mathbf{A})$, so heißt die Matrix \mathbf{A} **numerisch singulär**. Die Größe der Umgebung $U(\mathbf{A})$ hängt dabei sowohl von der verwendeten Hardware (relative Maschinengenauigkeit), als auch von der verwendeten Software (zur Behandlung der Matrix) ab. Zurück zum Beispiel 8.10.

8.11. Beispiel: Wir betrachten folgende Umgebung der Matrix \mathbf{A} :

$$U(\mathbf{A}) = \left\{ \bar{\mathbf{A}} \in \mathbb{R}^{2 \times 2} \mid |\bar{\mathbf{A}} - \mathbf{A}| \leq 0.5 \cdot 10^{-2} \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} \right\}.$$

Das ist gerade die Menge von Matrizen, die bei der Konvertierung in den Maschinenzahlbereich $\mathbb{M}(10, 2, \dots)$ auf die Matrix \mathbf{A} abgebildet werden. Wir sehen daher \mathbf{A} als Repräsentanten aller Matrizen $\bar{\mathbf{A}} \in U(\mathbf{A})$ an. Für die Matrix

$$\tilde{\mathbf{A}} = \begin{pmatrix} 1.000 & 0.9900 \\ 0.9900 & 0.9801 \end{pmatrix} \in U(\mathbf{A})$$

gilt aber $\det(\tilde{\mathbf{A}}) = 0$ und $\text{rg}(\tilde{\mathbf{A}}) = 1$. Folglich enthält $U(\mathbf{A})$ singuläre Matrizen. Die Matrix \mathbf{A} ist also in diesem Maschinenzahlbereich als numerisch singulär anzusehen. ♡

Bemerkungen: (i) Hinter der Aussage von Satz 8.9 steckt bei der Prüfung der Regularität von \mathbf{A} der Test einer REAL-Zahl auf Gleichheit mit Null. Solche Tests sollte man in Programmen vermeiden. Dafür ist ein Test der Form $|x| \leq \varepsilon$ mit einer vorgegebenen Genauigkeit ε besser.

(ii) Die Umgebung $U(\mathbf{A})$ lässt sich meist in der Form

$$U(\mathbf{A}) = \left\{ \bar{\mathbf{A}} \in \mathbb{R}^{n \times n} \mid \frac{|\bar{a}_{ij} - a_{ij}|}{1 + |a_{ij}|} \leq K \text{eps}; i, j = 1, \dots, n \right\}$$

darstellen. Dabei ist eps die relative Maschinengenauigkeit und K eine Konstante, die von der verwendeten Software abhängt.

Schon das letzte kleine Beispiel zeigt, dass Störungen die Aufgabe mehr oder weniger stark verändern. Den Einfluss dieser Störungen wollen wir nun genauer untersuchen.

Neben dem Gleichungssystem

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

mit der Lösung \mathbf{x} betrachten wir ein gestörtes System

$$(\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b} + \delta\mathbf{b}.$$

Die Lösung $\mathbf{x} + \delta\mathbf{x}$ des gestörten Systems weicht um $\delta\mathbf{x}$ von der Lösung des ursprünglichen Systems ab. Der Zusammenhang zwischen dieser Störung $\delta\mathbf{x}$ und den Eingabefeldern $\delta\mathbf{A}$ in der Matrix und $\delta\mathbf{b}$ in der rechten Seite soll nun untersucht werden. Wir stellen uns dabei zwei Fragen.

1. Wann besitzt das gestörte System eine eindeutige Lösung?
2. Wie lässt sich $\|\delta x\|$ durch $\|\delta A\|$ und $\|\delta b\|$ abschätzen?
Wie lässt sich $|\delta x|$ durch $|\delta A|$ und $|\delta b|$ abschätzen?

Die erste Frage ist gerade die Frage nach der numerischen Regularität von A . Wir fragen danach, wie groß eine Umgebung der regulären Matrix A maximal ist, so dass sie nur reguläre Matrizen enthält. Für die Einheitsmatrix wird diese Frage mit folgendem Satz beantwortet.

8.12. Satz: Für die (n, n) -Matrix P gelte $\|P\| < 1$. Dann ist die Matrix $I - P$ regulär, und es gilt die Abschätzung

$$\|(I - P)^{-1}\| \leq \frac{1}{1 - \|P\|}.$$

Beweis: Für jedes $x \in \mathbb{R}^n$ gilt

$$\begin{aligned} \|(I - P)x\| &= \|x - Px\| \geq |\|x\| - \|Px\|| \geq \|x\| - \|Px\| \\ &\geq \|x\| - \|P\|\|x\| = (1 - \|P\|)\|x\| > 0, \quad x \neq o. \end{aligned}$$

Damit ist $I - P$ regulär, und es existiert die Inverse $C = (I - P)^{-1}$. Dann gilt

$$\begin{aligned} 1 = \|I\| &= \|(I - P)C\| = \|C - PC\| \geq \|C\| - \|PC\| \\ &\geq \|C\| - \|P\|\|C\| = (1 - \|P\|)\|C\|. \end{aligned}$$

Wegen $1 - \|P\| > 0$ folgt daraus

$$\|C\| = \|(I - P)^{-1}\| \leq \frac{1}{1 - \|P\|}.$$

✱

Bemerkung: Die Bedingung $\|P\| < 1$ ist bezüglich einer Grenznorm auch notwendig. Das erkennt man, falls man

$$P = \begin{pmatrix} 1 & & & \\ & 0 & & O \\ & & \ddots & \\ & O & & 0 \end{pmatrix}$$

betrachtet. Hier gilt $\text{rg}(I - P) = n - 1$ und $\|P\| = 1$.

8.13. Störungslemma: *Gegeben seien eine reguläre (n, n) -Matrix \mathbf{A} und eine Störungsmatrix $\delta\mathbf{A}$, die der Bedingung*

$$\kappa = \|\mathbf{A}^{-1}\| \|\delta\mathbf{A}\| < 1$$

genügt. Dann ist auch die Matrix $\mathbf{A} + \delta\mathbf{A}$ regulär und es gilt

1.

$$\|(\mathbf{A} + \delta\mathbf{A})^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \kappa},$$

2.

$$\|(\mathbf{A} + \delta\mathbf{A})^{-1} - \mathbf{A}^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|^2}{1 - \kappa} \|\delta\mathbf{A}\|.$$

Beweis: Es sei $\mathbf{P} = -\mathbf{A}^{-1}\delta\mathbf{A}$. Dann gilt

$$\|\mathbf{P}\| = \|\mathbf{A}^{-1}\delta\mathbf{A}\| \leq \|\mathbf{A}^{-1}\| \|\delta\mathbf{A}\| = \kappa < 1.$$

Nach Satz 8.12 ist dann die Matrix $\mathbf{I} - \mathbf{P} = \mathbf{I} + \mathbf{A}^{-1}\delta\mathbf{A}$ regulär. Wegen der Regularität von \mathbf{A} ist damit aber auch $\mathbf{A} + \delta\mathbf{A}$ regulär. Für die erste Abschätzung folgt

$$\begin{aligned} \|(\mathbf{A} + \delta\mathbf{A})^{-1}\| &= \left\| \left[\mathbf{A} (\mathbf{I} + \mathbf{A}^{-1}\delta\mathbf{A}) \right]^{-1} \right\| = \left\| (\mathbf{I} + \mathbf{A}^{-1}\delta\mathbf{A})^{-1} \mathbf{A}^{-1} \right\| \\ &\leq \left\| (\mathbf{I} + \mathbf{A}^{-1}\delta\mathbf{A})^{-1} \right\| \|\mathbf{A}^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \|\mathbf{A}^{-1}\delta\mathbf{A}\|} \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \kappa}. \end{aligned}$$

Für die zweite Abschätzung erhalten wir

$$\begin{aligned} (\mathbf{A} + \delta\mathbf{A})^{-1} - \mathbf{A}^{-1} &= (\mathbf{A} + \delta\mathbf{A})^{-1} \left[\mathbf{I} - (\mathbf{A} + \delta\mathbf{A})\mathbf{A}^{-1} \right] \\ &= (\mathbf{A} + \delta\mathbf{A})^{-1} [\mathbf{A} - (\mathbf{A} + \delta\mathbf{A})] \mathbf{A}^{-1} \\ &= -(\mathbf{A} + \delta\mathbf{A})^{-1} \delta\mathbf{A} \mathbf{A}^{-1}. \end{aligned}$$

Damit erhält man

$$\|(\mathbf{A} + \delta\mathbf{A})^{-1} - \mathbf{A}^{-1}\| \leq \|(\mathbf{A} + \delta\mathbf{A})^{-1}\| \|\delta\mathbf{A}\| \|\mathbf{A}^{-1}\| \leq \frac{\|\mathbf{A}^{-1}\|^2}{1 - \kappa} \|\delta\mathbf{A}\|.$$

Bemerkungen: (i) Durch $\kappa = \|\mathbf{A}^{-1}\| \|\delta\mathbf{A}\| < 1$ wird eine Umgebung der regulären Matrix \mathbf{A} festgelegt, die nur reguläre Matrizen enthält:

$$U(\mathbf{A}) = \left\{ \bar{\mathbf{A}} \in \mathbb{R}^{n \times n} \mid \|\mathbf{A} - \bar{\mathbf{A}}\| \leq \Delta\mathbf{A} < \frac{1}{\|\mathbf{A}^{-1}\|} \right\}.$$

Die Größe $\|\mathbf{A}^{-1}\|$ ist damit ein Maß für die Regularität von \mathbf{A} .

(ii) Betrachtet man die Inversenbildung als Funktion $\Phi: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^{n \times n}$, so erkennt man aus der zweiten Abschätzung in Satz 8.13, dass für ein $\bar{\mathbf{A}} \in U(\mathbf{A})$

$$\begin{aligned} \|\bar{\mathbf{A}}^{-1} - \mathbf{A}^{-1}\| &= \|\Phi(\bar{\mathbf{A}}) - \Phi(\mathbf{A})\| \\ &\leq \frac{\|\mathbf{A}^{-1}\|^2}{1 - \|\bar{\mathbf{A}} - \mathbf{A}\| \|\mathbf{A}^{-1}\|} \|\bar{\mathbf{A}} - \mathbf{A}\| \\ &\leq \frac{\|\mathbf{A}^{-1}\|^2}{1 - \Delta\mathbf{A} \|\mathbf{A}^{-1}\|} \|\bar{\mathbf{A}} - \mathbf{A}\| \end{aligned}$$

und damit

$$\|\Phi(\bar{\mathbf{A}}) - \Phi(\mathbf{A})\| \leq L(\mathbf{A}, \Delta\mathbf{A}) \|\bar{\mathbf{A}} - \mathbf{A}\|$$

gilt. Die Inversenbildung ist daher für eine reguläre Matrix lokal lipschitzstetig mit der LIPSCHITZ-Konstanten

$$L(\mathbf{A}, \Delta\mathbf{A}) = \frac{\|\mathbf{A}^{-1}\|^2}{1 - \Delta\mathbf{A} \|\mathbf{A}^{-1}\|} \|\bar{\mathbf{A}} - \mathbf{A}\|.$$

Nach diesen Vorbereitungen ist der Einfluss von Störungen in einem linearen Gleichungssystem abschätzbar. Es sei durch

$$(\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x}) = \mathbf{b} + \delta\mathbf{b}$$

ein gestörtes lineares Gleichungssystem gegeben. Die Störung $\delta\mathbf{A}$ genüge der Bedingung

$$\kappa = \|\mathbf{A}^{-1}\| \|\delta\mathbf{A}\| < 1.$$

Nach Satz 8.13 ist dann die Matrix $\mathbf{A} + \delta\mathbf{A}$ regulär, und das gestörte System besitzt eine eindeutige Lösung

$$\mathbf{x} + \delta\mathbf{x} = (\mathbf{A} + \delta\mathbf{A})^{-1}(\mathbf{b} + \delta\mathbf{b}).$$

Dann gilt

$$\begin{aligned}
 \delta x &= (\mathbf{A} + \delta \mathbf{A})^{-1}(\mathbf{b} + \delta \mathbf{b}) - \mathbf{A}^{-1}\mathbf{b} \\
 &= (\mathbf{A} + \delta \mathbf{A})^{-1} \left[\mathbf{b} + \delta \mathbf{b} - (\mathbf{A} + \delta \mathbf{A})\mathbf{A}^{-1}\mathbf{b} \right] \\
 &= (\mathbf{A} + \delta \mathbf{A})^{-1} \left[\mathbf{b} + \delta \mathbf{b} - \mathbf{b} - \delta \mathbf{A} \cdot \mathbf{A}^{-1}\mathbf{b} \right] \\
 &= (\mathbf{A} + \delta \mathbf{A})^{-1} [-\delta \mathbf{A}\mathbf{x} + \delta \mathbf{b}] \\
 &= \mathbf{A}^{-1} [-\delta \mathbf{A}\mathbf{x} + \delta \mathbf{b}] + \left[(\mathbf{A} + \delta \mathbf{A})^{-1} - \mathbf{A}^{-1} \right] [-\delta \mathbf{A}\mathbf{x} + \delta \mathbf{b}] \\
 &= \delta \mathbf{x}' + O(\|\delta \mathbf{A}\|(\|\delta \mathbf{A}\| + \|\delta \mathbf{b}\|)).
 \end{aligned}$$

$\delta \mathbf{x}' = \mathbf{A}^{-1} [-\delta \mathbf{A}\mathbf{x} + \delta \mathbf{b}]$ stellt dabei den bezüglich $\delta \mathbf{A}$ und $\delta \mathbf{b}$ linearen Teil des Fehlers $\delta \mathbf{x}$ dar. Für $\delta \mathbf{x}$ erhält man die Abschätzung

$$\|\delta \mathbf{x}\| \leq \left\| (\mathbf{A} + \delta \mathbf{A})^{-1} \right\| \|\delta \mathbf{A}\mathbf{x} + \delta \mathbf{b}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \kappa} (\|\delta \mathbf{A}\| \|\mathbf{x}\| + \|\delta \mathbf{b}\|).$$

Für ein inhomogenes Gleichungssystem ($\mathbf{x} \neq \mathbf{o}$) lässt sich auch der relative Fehler der Lösung abschätzen. Man erhält

$$\begin{aligned}
 \frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} &\leq \frac{\|\mathbf{A}^{-1}\|}{1 - \kappa} \left(\|\delta \mathbf{A}\| + \frac{\|\delta \mathbf{b}\|}{\|\mathbf{x}\|} \right) \\
 &= \frac{1}{1 - \kappa} \left(\|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|} + \frac{\|\mathbf{A}^{-1}\| \|\mathbf{b}\|}{\|\mathbf{x}\|} \cdot \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|} \right).
 \end{aligned}$$

Berücksichtigt man noch

$$\|\mathbf{b}\| = \|\mathbf{A}\mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{x}\|,$$

so ergibt sich

$$\begin{aligned}
 \frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} &\leq \frac{1}{1 - \kappa} \left(\|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|} + \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|} \right) \\
 &= \frac{\text{cond}(\mathbf{A})}{1 - \kappa} \left(\frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|} + \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|} \right)
 \end{aligned}$$

mit $\text{cond}(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$.

Wir fassen die Ergebnisse in einem Satz zusammen.

8.14. Satz: Gegeben sei das lineare Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit einer regulären (n, n) -Matrix \mathbf{A} und $\mathbf{b} \in \mathbb{R}^n$; ferner sei

$$(\mathbf{A} + \delta \mathbf{A})(\mathbf{x} + \delta \mathbf{x}) = \mathbf{b} + \delta \mathbf{b}$$

ein gestörtes Gleichungssystem, wobei die Störung $\delta\mathbf{A}$ der Bedingung

$$\kappa = \left\| \mathbf{A}^{-1} \right\| \left\| \delta\mathbf{A} \right\| < 1$$

genüge. Dann ist die Matrix $\mathbf{A} + \delta\mathbf{A}$ regulär, und das gestörte Gleichungssystem besitzt die eindeutige Lösung

$$\mathbf{x} + \delta\mathbf{x} = (\mathbf{A} + \delta\mathbf{A})^{-1}(\mathbf{b} + \delta\mathbf{b}).$$

Für den Fehler $\delta\mathbf{x}$ der Lösung gilt:

1.

$$\delta\mathbf{x} = \delta\mathbf{x}' + O(\|\delta\mathbf{A}\|(\|\delta\mathbf{A}\| + \|\delta\mathbf{b}\|)),$$

wobei

$$\delta\mathbf{x}' = \mathbf{A}^{-1}[-\delta\mathbf{A}\mathbf{x} + \delta\mathbf{b}]$$

den bezüglich $\delta\mathbf{A}$ und $\delta\mathbf{b}$ linearen Teil des Fehlers $\delta\mathbf{x}$ darstellt.

2.

$$\|\delta\mathbf{x}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \kappa} (\|\mathbf{x}\| \|\delta\mathbf{A}\| + \|\delta\mathbf{b}\|).$$

$\|\mathbf{A}^{-1}\| \|\mathbf{x}\|$ bzw. $\|\mathbf{A}^{-1}\|$ sind die absoluten partiellen Konditionszahlen bezüglich \mathbf{A} und \mathbf{b} .

3. Für $\mathbf{b} \neq \mathbf{o}$ ($\mathbf{x} \neq \mathbf{o}$) gilt

$$\begin{aligned} \frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} &\leq \frac{1}{1 - \kappa} \left(\|\mathbf{A}\| \|\mathbf{A}^{-1}\| \frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} + \frac{\|\mathbf{A}^{-1}\| \|\mathbf{b}\| \|\delta\mathbf{b}\|}{\|\mathbf{x}\| \|\mathbf{b}\|} \right) \\ &\leq \frac{\text{cond}(\mathbf{A})}{1 - \kappa} \left(\frac{\|\delta\mathbf{A}\|}{\|\mathbf{A}\|} + \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|} \right). \end{aligned}$$

$\|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ und $\|\mathbf{A}^{-1}\| \|\mathbf{b}\|/\|\mathbf{x}\|$ sind die relativen partiellen Konditionszahlen bezüglich \mathbf{A} und \mathbf{b} .

Die Größe $\text{cond}(\mathbf{A}) = \|\mathbf{A}\| \|\mathbf{A}^{-1}\|$ wird schlechthin als **Kondition** der Matrix \mathbf{A} bezeichnet.

Bemerkungen: (i) Die Abschätzungen aus Satz 8.14 verkörpern den schlechtesten Fall. Im allgemeinen liefern sie zu pessimistische Schranken. Betrachten wir dazu wieder unser Beispiel. Es ist

$$\mathbf{A} = \begin{pmatrix} 1.00 & 0.99 \\ 0.99 & 0.98 \end{pmatrix}, \quad \mathbf{A}^{-1} = \begin{pmatrix} -9800 & 9900 \\ 9900 & -10000 \end{pmatrix}.$$

In der Spaltensummennorm gilt $\|\mathbf{A}\|_1 = 1.99$ und $\|\mathbf{A}^{-1}\|_1 = 19900$. Damit folgt $\text{cond}_1(\mathbf{A}) = 39601$. Die Voraussetzung von Satz 8.14 wäre schon scharf. Wir dürften nur Störungen in der Matrix zulassen, für die

$$\|\delta\mathbf{A}\|_1 < \frac{1}{\|\mathbf{A}^{-1}\|_1} = \frac{1}{19900} \approx 5.025 \cdot 10^{-5}$$

gilt. Schon für die kleine Störung

$$\delta\mathbf{A} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} 10^{-4}$$

wäre der Satz nicht anwendbar. Betrachten wir darum zunächst nur Störungen der rechten Seite, also den Fall $\delta\mathbf{A} = \mathbf{O}$. Für $\mathbf{b} \neq \mathbf{o}$ gilt dann

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq \text{cond}(\mathbf{A}) \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|}.$$

Für verschiedene rechte Seiten und verschiedene Störungen ergibt sich:

\mathbf{b}	$\delta\mathbf{b}$	$\frac{\ \delta\mathbf{b}\ }{\ \mathbf{b}\ }$	\mathbf{x}	$\delta\mathbf{x}$	$\frac{\ \delta\mathbf{x}\ }{\ \mathbf{x}\ }$	$\frac{\ \delta\mathbf{x}\ \ \mathbf{b}\ }{\ \mathbf{x}\ \ \delta\mathbf{b}\ }$
$\begin{pmatrix} 1.00 \\ 0.99 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 10^{-5} \end{pmatrix}$	$\approx 5 \cdot 10^{-6}$	$\begin{pmatrix} 1 \\ 0 \end{pmatrix}$	$\begin{pmatrix} 0.099 \\ 0.100 \end{pmatrix}$	≈ 0.2	39601
$\begin{pmatrix} -1 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 10^{-5} \\ 0 \end{pmatrix}$	$\approx 5 \cdot 10^{-6}$	$\begin{pmatrix} 19700 \\ -19900 \end{pmatrix}$	$\begin{pmatrix} -0.0998 \\ 0.0990 \end{pmatrix}$	$\approx 5 \cdot 10^{-6}$	0.995

In der letzten Spalte ist jeweils der wahre Verstärkungsfaktor, mit dem die relativen Fehler in der rechten Seite in die Lösung eingehen, angegeben.

Betrachten wir nun auch Störungen in der Matrix. Es sei

$$\delta\mathbf{A} = \begin{pmatrix} 0 & 0 \\ 0 & 2 \cdot 10^{-5} \end{pmatrix}, \quad \|\delta\mathbf{A}\|_1 = 2 \cdot 10^{-5}, \quad \frac{\text{cond}_1(\mathbf{A})}{1 - \|\delta\mathbf{A}\|_1 \|\mathbf{A}^{-1}\|_1} \approx 65800.$$

Für verschiedene rechte Seiten erhalten wir folgende Ergebnisse:

\mathbf{b}	\mathbf{x}	$\delta \mathbf{x}$	$\frac{\ \delta \mathbf{x}\ }{\ \mathbf{x}\ }$	$\frac{\ \delta \mathbf{x}\ \ \mathbf{b}\ }{\ \mathbf{x}\ \ \delta \mathbf{b}\ }$
$\begin{pmatrix} 1 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 100 \\ -100 \end{pmatrix}$	$\begin{pmatrix} 24.75 \\ -25.00 \end{pmatrix}$	≈ 0.24875	≈ 25000
$\begin{pmatrix} -1 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 19700 \\ -19900 \end{pmatrix}$	$\begin{pmatrix} 4925.25 \\ -4975.00 \end{pmatrix}$	≈ 0.2500	≈ 25000

Für die Störung

$$\delta \mathbf{A} = \begin{pmatrix} -10^{-5} & 10^{-5} \\ -10^{-5} & 10^{-5} \end{pmatrix}, \quad \|\delta \mathbf{A}\|_1 = 2 \cdot 10^{-5}, \quad \frac{\text{cond}_1(\mathbf{A})}{1 - \|\delta \mathbf{A}\|_1 \|\mathbf{A}^{-1}\|_1} \approx 65800$$

ergibt sich

\mathbf{b}	\mathbf{x}	$\delta \mathbf{x}$	$\frac{\ \delta \mathbf{x}\ }{\ \mathbf{x}\ }$	$\frac{\ \delta \mathbf{x}\ \ \mathbf{b}\ }{\ \mathbf{x}\ \ \delta \mathbf{b}\ }$
$\begin{pmatrix} 1 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 100 \\ -100 \end{pmatrix}$	$\begin{pmatrix} 0 \\ 0 \end{pmatrix}$	0	0
$\begin{pmatrix} -1 \\ 1 \end{pmatrix}$	$\begin{pmatrix} 19700 \\ -19900 \end{pmatrix}$	$\begin{pmatrix} -0.2 \\ 0.2 \end{pmatrix}$	$\approx 10^{-5}$	≈ 0.5

Wir sehen, dass für verschiedene Abschätzungen Fälle existieren, bei denen die jeweiligen maximalen Fehlerverstärkungen voll wirksam werden. Es existieren aber auch Störungen vom gleichen Störungsniveau, die sich fast gar nicht auf die Lösung auswirken. Störungen in der Matrix können sich dabei bedeutend stärker bemerkbar machen, da die Abhängigkeit der Lösung von der Matrix Unstetigkeitsstellen besitzt.

(ii) Verwendet man die euklidische Norm zur Fehlerabschätzung, so lässt sich die dritte Ungleichung in Satz 8.14 etwas verschärfen. Es gilt

$$\frac{\|\delta \mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \frac{1}{1 - \|\delta \mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2} \left(\text{cond}_2(\mathbf{A}) \frac{\|\delta \mathbf{A}\|_2}{\|\mathbf{A}\|_2} + \|\mathbf{A}^{-1}\|_2 \frac{\|\mathbf{A}^T \mathbf{b}\|_2}{\|\mathbf{b}\|_2} \frac{\|\delta \mathbf{b}\|_2}{\|\mathbf{b}\|_2} \right).$$

Hier wird die Abhängigkeit des Verstärkungsfaktors von der rechten Seite besser berücksichtigt. Für die ersten Beispiele ($\delta\mathbf{A} = \mathbf{O}$) liefert das die Abschätzungen

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq 36150 \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|}, \quad \mathbf{b} = \begin{pmatrix} 1 \\ 1 \end{pmatrix},$$

$$\frac{\|\delta\mathbf{x}\|}{\|\mathbf{x}\|} \leq 183 \frac{\|\delta\mathbf{b}\|}{\|\mathbf{b}\|}, \quad \mathbf{b} = \begin{pmatrix} -1 \\ 1 \end{pmatrix}.$$

(iii) In den Abschätzungen tauchen die Normen verschiedener Größen auf. Dabei sind $\|\mathbf{b}\|$ und $\|\mathbf{A}\|$ meist leicht zu berechnen. Für $\|\delta\mathbf{b}\|$ und $\|\delta\mathbf{A}\|$ kennt man oft Schranken. Das Berechnen von $\|\mathbf{A}^{-1}\|$ bereitet Schwierigkeiten, da die Inverse \mathbf{A}^{-1} im allgemeinen natürlich nicht bekannt ist.

Die Aussagen von Satz 8.14 stellen eine Abschätzung des unvermeidbaren Fehlers dar. Mit Hilfe dieses Satzes ist nicht die Güte einer berechneten Lösung bewertbar. Im Sinne der Rückwärtsanalyse müssten wir eine berechnete Lösung $\bar{\mathbf{x}}$ des linearen Gleichungssystems $\mathbf{A}\mathbf{x} = \mathbf{b}$ akzeptieren, falls sie als Lösung eines benachbarten Systems $\bar{\mathbf{A}}\bar{\mathbf{x}} = \bar{\mathbf{b}}$ interpretierbar ist, wobei sich die Matrix $\bar{\mathbf{A}}$ und der Vektor $\bar{\mathbf{b}}$ von \mathbf{A} und \mathbf{b} nur in der Größenordnung des Eingabefehlers unterscheiden. Mit dem folgenden Satz lässt sich für eine berechnete Lösung entscheiden, ob sie Lösung eines derartigen benachbarten Systems ist.

8.15. PRAGER-OETTLI: *Durch $\mathbf{A}\mathbf{x} = \mathbf{b}$ sei ein lineares Gleichungssystem gegeben. Mit der Matrix $\Delta\mathbf{A} \geq \mathbf{O}$ und dem Vektor $\Delta\mathbf{b} \geq \mathbf{o}$ seien folgende Mengen definiert:*

$$\mathcal{A} = \left\{ \hat{\mathbf{A}} \in \mathbb{R}^{n \times n} \mid |\hat{\mathbf{A}} - \mathbf{A}| \leq \Delta\mathbf{A} \right\}$$

und

$$\mathcal{B} = \left\{ \hat{\mathbf{b}} \in \mathbb{R}^n \mid |\hat{\mathbf{b}} - \mathbf{b}| \leq \Delta\mathbf{b} \right\}.$$

Der Vektor $\bar{\mathbf{x}}$ ist genau dann exakte Lösung eines Gleichungssystems $\bar{\mathbf{A}}\bar{\mathbf{x}} = \bar{\mathbf{b}}$ mit $\bar{\mathbf{A}} \in \mathcal{A}$ und $\bar{\mathbf{b}} \in \mathcal{B}$, wenn

$$|\mathbf{r}(\bar{\mathbf{x}})| \leq \Delta\mathbf{A}|\bar{\mathbf{x}}| + \Delta\mathbf{b}$$

gilt. Dabei bezeichnet $\mathbf{r}(\bar{\mathbf{x}}) = \mathbf{b} - \mathbf{A}\bar{\mathbf{x}}$ das Residuum von $\bar{\mathbf{x}}$ bezüglich des ursprünglichen Gleichungssystems.

Beweis: (\Rightarrow) Wir nehmen an, dass eine Matrix $\bar{\mathbf{A}} \in \mathcal{A}$ und ein Vektor $\bar{\mathbf{b}} \in \mathcal{B}$ existieren, so dass $\bar{\mathbf{A}}\bar{\mathbf{x}} = \bar{\mathbf{b}}$ gilt. Dann gilt

$$\bar{\mathbf{A}} = \mathbf{A} + \delta\mathbf{A}, \quad |\delta\mathbf{A}| \leq \Delta\mathbf{A}$$

und

$$\bar{\mathbf{b}} = \mathbf{b} + \delta\mathbf{b}, \quad |\delta\mathbf{b}| \leq \Delta\mathbf{b}.$$

Damit folgt

$$\begin{aligned} |\mathbf{r}(\bar{\mathbf{x}})| &= |\mathbf{b} - \mathbf{A}\bar{\mathbf{x}}| = |\bar{\mathbf{b}} - \delta\mathbf{b} - (\bar{\mathbf{A}} - \delta\mathbf{A})\bar{\mathbf{x}}| = |\bar{\mathbf{b}} - \bar{\mathbf{A}}\bar{\mathbf{x}} - \delta\mathbf{b} + \delta\mathbf{A}\bar{\mathbf{x}}| \\ &= |-\delta\mathbf{b} + \delta\mathbf{A}\bar{\mathbf{x}}| \leq |\delta\mathbf{b}| + |\delta\mathbf{A}||\bar{\mathbf{x}}| \leq \Delta\mathbf{b} + \Delta\mathbf{A}|\bar{\mathbf{x}}|. \end{aligned}$$

(\Leftarrow) Es gelte

$$|\mathbf{r}(\bar{\mathbf{x}})| \leq \Delta\mathbf{b} + \Delta\mathbf{A}|\bar{\mathbf{x}}|.$$

Es sei weiter

$$\mathbf{r}(\bar{\mathbf{x}}) = \begin{pmatrix} r_1 \\ r_2 \\ \vdots \\ r_n \end{pmatrix}, \quad \mathbf{s} = \Delta\mathbf{b} + \Delta\mathbf{A}|\bar{\mathbf{x}}| = \begin{pmatrix} s_1 \\ s_2 \\ \vdots \\ s_n \end{pmatrix} \geq \mathbf{o}, \quad \Delta\mathbf{b} = \begin{pmatrix} \Delta b_1 \\ \Delta b_2 \\ \vdots \\ \Delta b_n \end{pmatrix}, \quad \bar{\mathbf{x}} = \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_n \end{pmatrix}$$

und

$$\Delta\mathbf{A} = \begin{pmatrix} \Delta a_{11} & \Delta a_{12} & \cdots & \Delta a_{1n} \\ \Delta a_{21} & \Delta a_{22} & \cdots & \Delta a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \Delta a_{n1} & \Delta a_{n2} & \cdots & \Delta a_{nn} \end{pmatrix}.$$

Für die unbekanntenen Größen $\bar{\mathbf{A}}$, $\delta\mathbf{A}$, $\bar{\mathbf{b}}$, $\delta\mathbf{b}$ und $\bar{\mathbf{x}}$ machen wir die Ansätze

$$\bar{\mathbf{A}} = \begin{pmatrix} \bar{a}_{11} & \bar{a}_{12} & \cdots & \bar{a}_{1n} \\ \bar{a}_{21} & \bar{a}_{22} & \cdots & \bar{a}_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \bar{a}_{n1} & \bar{a}_{n2} & \cdots & \bar{a}_{nn} \end{pmatrix}, \quad \delta\mathbf{A} = \begin{pmatrix} \delta a_{11} & \delta a_{12} & \cdots & \delta a_{1n} \\ \delta a_{21} & \delta a_{22} & \cdots & \delta a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \delta a_{n1} & \delta a_{n2} & \cdots & \delta a_{nn} \end{pmatrix},$$

$$\bar{\mathbf{b}} = \begin{pmatrix} \bar{b}_1 \\ \bar{b}_2 \\ \vdots \\ \bar{b}_n \end{pmatrix}, \quad \delta\mathbf{b} = \begin{pmatrix} \delta b_1 \\ \delta b_2 \\ \vdots \\ \delta b_n \end{pmatrix}.$$

Nun definieren wir

$$\delta a_{ij} = \begin{cases} \frac{r_i}{s_i} \Delta a_{ij} \text{sign}(\bar{x}_j) & \text{für } s_i > 0 \\ 0 & \text{für } s_i = 0 \end{cases}$$

und

$$\delta b_i = \begin{cases} -\frac{r_i}{s_i} \Delta b_i & \text{für } s_i > 0 \\ 0 & \text{für } s_i = 0 \end{cases}.$$

Wir haben zu zeigen, dass für die so definierten Größen

$$\bar{\mathbf{A}} = \mathbf{A} + \delta \mathbf{A}, \quad \bar{\mathbf{b}} = \mathbf{b} + \delta \mathbf{b}$$

einerseits $\bar{\mathbf{A}} \in \mathcal{A}$ und $\bar{\mathbf{b}} \in \mathcal{B}$ und andererseits $\bar{\mathbf{A}}\bar{\mathbf{x}} = \bar{\mathbf{b}}$ gilt. Die ersten beiden Aussagen folgen sofort aus der Voraussetzung $|\mathbf{r}(\bar{\mathbf{x}})| \leq \mathbf{s}$. Es folgt $|r_i| \leq s_i$ für $i = 1, \dots, n$ und damit

$$|\delta a_{ij}| \leq \Delta a_{ij}, \quad i, j = 1, \dots, n$$

und

$$|\delta b_i| \leq \Delta b_i, \quad i = 1, \dots, n.$$

Damit folgt aber $|\delta \mathbf{A}| \leq \Delta \mathbf{A}$ und $|\delta \mathbf{b}| \leq \Delta \mathbf{b}$, daher $\bar{\mathbf{A}} \in \mathcal{A}$ und $\bar{\mathbf{b}} \in \mathcal{B}$. Weiterhin folgt

$$\mathbf{r}(\bar{\mathbf{x}}) = \mathbf{b} - \mathbf{A}\bar{\mathbf{x}} = (\bar{\mathbf{b}} - \delta \mathbf{b}) - (\bar{\mathbf{A}} - \delta \mathbf{A})\bar{\mathbf{x}} = \bar{\mathbf{b}} - \bar{\mathbf{A}}\bar{\mathbf{x}} - \delta \mathbf{b} + \delta \mathbf{A}\bar{\mathbf{x}}.$$

Für $\delta \mathbf{b} - \delta \mathbf{A}\bar{\mathbf{x}}$ gilt komponentenweise im Falle $s_i > 0$

$$\begin{aligned} (\delta \mathbf{b} - \delta \mathbf{A}\bar{\mathbf{x}})_i &= \delta b_i - \sum_{j=1}^n \delta a_{ij} \bar{x}_j = -\frac{r_i}{s_i} \Delta b_i - \sum_{j=1}^n \frac{r_i}{s_i} \Delta a_{ij} \text{sign}(\bar{x}_j) \bar{x}_j \\ &= -\frac{r_i}{s_i} \left[\Delta b_i + \sum_{j=1}^n \Delta a_{ij} |\bar{x}_j| \right] = -\frac{r_i}{s_i} s_i = -r_i. \end{aligned}$$

Im Falle $s_i = 0$ gilt auch $r_i = 0$, $\delta b_i = 0$ und $\delta a_{ij} = 0$ für $j = 1, \dots, n$. Damit ist wieder $0 = (\delta \mathbf{b} - \delta \mathbf{A}\bar{\mathbf{x}})_i = -r_i$.

Insgesamt gilt dann $\delta \mathbf{b} - \delta \mathbf{A}\bar{\mathbf{x}} = -\mathbf{r}(\bar{\mathbf{x}})$ und weiter $\mathbf{r}(\bar{\mathbf{x}}) = \bar{\mathbf{b}} - \bar{\mathbf{A}}\bar{\mathbf{x}} + \mathbf{r}(\bar{\mathbf{x}})$, also $\bar{\mathbf{A}}\bar{\mathbf{x}} = \bar{\mathbf{b}}$. *

Mit Hilfe des Satzes von PRAGER und OETTLI wird aus der Größe des Residuums auf die Güte der Lösung geschlossen, falls Informationen über die Datenfehler in der Matrix und in der rechten Seite bekannt sind. Das Berechnen des Residuums sollte dabei mit großer Sorgfalt erfolgen (höhere Genauigkeit), da im allgemeinen Auslöschung auftritt. Genügt das Residuum $\mathbf{r}(\bar{\mathbf{x}})$ der Bedingung

$$|\mathbf{r}(\bar{\mathbf{x}})| \leq \Delta \mathbf{b} + \Delta \mathbf{A}|\bar{\mathbf{x}}|,$$

so ist \bar{x} als Lösung zu akzeptieren.

Oft hat man eine spezielle Fehlerkorrelation der Form $\Delta \mathbf{A} = \varepsilon |\mathbf{A}|$ und $\Delta \mathbf{b} = \varepsilon |\mathbf{b}|$ vorliegen. Hier sind alle Eingabedaten mit gleicher relativer Genauigkeit gegeben. In diesem Falle ergibt sich

$$|\mathbf{b} - \mathbf{A}\bar{x}| \leq \varepsilon (|\mathbf{b}| + |\mathbf{A}||\bar{x}|)$$

als Bedingung für die Brauchbarkeit der Lösung \bar{x} . Hat man andererseits ein \bar{x} berechnet, so lässt sich über die obige Ungleichung ein kleinstes Fehlerniveau $\underline{\varepsilon}$ berechnen, für das \bar{x} noch als Lösung akzeptierbar ist. Man erhält

$$\underline{\varepsilon} = \max_{i=1, \dots, n} \frac{|(\mathbf{b} - \mathbf{A}\bar{x})_i|}{(|\mathbf{b}| + |\mathbf{A}||\bar{x}|)_i}.$$

Falls für das wahre Fehlerniveau $\varepsilon < \underline{\varepsilon}$ gilt, ist \bar{x} nicht als Lösung akzeptierbar. Im Falle $\varepsilon \geq \underline{\varepsilon}$ ist \bar{x} zu akzeptieren.

8.16. Beispiel: Wir betrachten ein Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit

$$\mathbf{A} = \begin{pmatrix} 1.00 & 0.99 \\ 0.99 & 0.98 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 1.00 \\ 0.99 \end{pmatrix}.$$

Für die Näherungslösung

$$\bar{x} = \begin{pmatrix} 1.099 \\ 0.100 \end{pmatrix}$$

erhält man

$$\mathbf{r}(\bar{x}) = \begin{pmatrix} -0.19800 \\ -0.19601 \end{pmatrix}, \quad |\mathbf{A}||\bar{x}| + |\mathbf{b}| = \begin{pmatrix} 2.19800 \\ 2.17601 \end{pmatrix}.$$

Daraus berechnet man $\underline{\varepsilon} \approx 0.09$. Erst ab einem Fehler der Eingabedaten von 9% oder mehr ist die Näherungslösung akzeptierbar. Für die rechte Seite

$$\mathbf{b} = \begin{pmatrix} -1 \\ 1 \end{pmatrix}$$

und die Näherungslösung

$$\bar{x} = \begin{pmatrix} 19600 \\ -20000 \end{pmatrix}$$

erhält man

$$\mathbf{r}(\bar{x}) = \begin{pmatrix} 199 \\ 197 \end{pmatrix}, \quad |\mathbf{A}||\bar{x}| + |\mathbf{b}| = \begin{pmatrix} 39401 \\ 39005 \end{pmatrix}.$$

Hier ergibt sich $\underline{\varepsilon} \approx 0.005$. Schon für relative Fehler der Eingabedaten von 5% oder mehr ist die Lösung zu akzeptieren.



Dieses Beispiel zeigt auch, dass aus einem kleinen Residuum noch nicht auf eine genaue Lösung geschlossen werden darf.

8.2. Direkte Lösungsverfahren

Eine Möglichkeit, ein lineares Gleichungssystem zu lösen, besteht darin, das Lösen auf eine Folge von einfach zu lösenden Gleichungssystemen zurückzuführen. Wir lösen anstelle von

$$A\mathbf{x} = \mathbf{b}$$

Systeme

$$A^{(i)}\mathbf{x}^{(i)} = \mathbf{b}^{(i)}, \quad i = 1, \dots, k.$$

Dabei ist die Folge so zu konstruieren, dass einerseits $\mathbf{x}^{(k)}$ bei exakter Rechnung die Lösung des ursprünglichen Gleichungssystems ist und andererseits die Systeme mit einem Gesamtaufwand zu lösen sind, der den Aufwand beim Invertieren der Koeffizientenmatrix nicht übersteigt. Als Koeffizientenmatrizen bieten sich dabei Dreiecksmatrizen, Diagonalmatrizen und orthogonale Matrizen an. Ein solches Vorgehen bezeichnen wir als **direktes Lösungsverfahren**. Charakteristisch für ein direktes Lösungsverfahren ist, dass man bei exakter Rechnung nach endlich vielen Rechenschritten die exakte Lösung erhält. Die Anzahl der benötigten Rechenschritte ist dabei a priori abschätzbar.

8.2.1. Die LU-Zerlegung

Schon im alten China nutzte man das heute als GAUSS-Algorithmus bekannte Verfahren zum Lösen linearer Gleichungssysteme. Die Vorgehensweise ist bekannt: Man versucht, durch geeignete Linearkombination von Gleichungen, nach und nach die Variablen zu eliminieren. Bezogen auf die Koeffizientenmatrix A bedeutet das: Man versucht, durch Linearkombination der Zeilen unterhalb der Diagonalen Nullen zu erzeugen. Wir wollen uns einen Transformationsschritt ansehen und gleich überlegen, wie dieser Schritt mit Hilfe geeigneter Transformationsmatrizen darstellbar ist. Dabei setzen wir immer eine reguläre Matrix A voraus.

Es sei $A^{(0)} = A$ und $\mathbf{b}^{(0)} = \mathbf{b}$. Nach k Schritten wurden A und \mathbf{b} in $A^{(k)}$ und $\mathbf{b}^{(k)}$

transformiert. Die Matrix $\mathbf{A}^{(k)}$ besitzt dabei die Struktur

$$\mathbf{A}^{(k)} = \begin{pmatrix} a_{11}^{(k)} & \cdots & a_{1k}^{(k)} & a_{1,k+1}^{(k)} & \cdots & a_{1n}^{(k)} \\ 0 & \ddots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \ddots & a_{kk}^{(k)} & a_{k,k+1}^{(k)} & \cdots & a_{kn}^{(k)} \\ 0 & \cdots & 0 & a_{k+1,k+1}^{(k)} & \cdots & a_{k+1,n}^{(k)} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & a_{n,k+1}^{(k)} & \cdots & a_{nn}^{(k)} \end{pmatrix} = \left(\begin{array}{c|c} \mathbf{U}^{(k)} & \\ \hline \mathbf{O} & \mathbf{M}^{(k)} \end{array} \right)$$

mit $\mathbf{U}^{(k)} \in \mathbb{R}^{k \times n}$ und $\mathbf{M}^{(k)} \in \mathbb{R}^{(n-k) \times (n-k)}$. Nehmen wir nun an, dass $a_{k+1,k+1}^{(k)} \neq 0$ gilt, so sind in der $(k+1)$ -ten Spalte unterhalb der Diagonalen Nullen erzeugbar, indem zur i -ten Zeile ($i = k+2, \dots, n$) von $\mathbf{A}^{(k)}$ das $(-a_{i,k+1}^{(k)}/a_{k+1,k+1}^{(k)})$ -fache der $(k+1)$ -ten Zeile addiert wird.

Bezeichnen wir mit $\mathbf{z}_1^{(k)T}, \dots, \mathbf{z}_n^{(k)T}$ die Zeilen der Matrix $\mathbf{A}^{(k)}$, so gilt

$$\mathbf{z}_i^{(k+1)T} = \mathbf{z}_i^{(k)T} = \mathbf{e}_i^T \mathbf{A}^{(k)}, \quad i = 1, \dots, k+1$$

und für $i = k+2, \dots, n$

$$\begin{aligned} \mathbf{z}_i^{(k+1)T} &= \mathbf{z}_i^{(k)T} - \frac{a_{i,k+1}^{(k)}}{a_{k+1,k+1}^{(k)}} \mathbf{z}_{k+1}^{(k)T} \\ &= \mathbf{e}_i^T \mathbf{A}^{(k)} - \frac{a_{i,k+1}^{(k)}}{a_{k+1,k+1}^{(k)}} \mathbf{e}_{k+1}^T \mathbf{A}^{(k)} \\ &= \left(\mathbf{e}_i^T - \bar{l}_{i,k+1} \mathbf{e}_{k+1}^T \right) \mathbf{A}^{(k)} \end{aligned}$$

mit

$$\bar{l}_{i,k+1} = \frac{a_{i,k+1}^{(k)}}{a_{k+1,k+1}^{(k)}}.$$

Fasst man diese Gleichungen zusammen, so erhält man

$$\mathbf{A}^{(k+1)} = \begin{pmatrix} \mathbf{z}_1^{(k+1)T} \\ \vdots \\ \mathbf{z}_n^{(k+1)T} \end{pmatrix} = \begin{pmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_{k+1}^T \\ \mathbf{e}_{k+2}^T - \bar{l}_{k+2,k+1} \mathbf{e}_{k+1}^T \\ \vdots \\ \mathbf{e}_n^T - \bar{l}_{n,k+1} \mathbf{e}_{k+1}^T \end{pmatrix} \mathbf{A}^{(k)} = \mathbf{L}_{k+1}(-\bar{\mathbf{l}}_{k+1}) \mathbf{A}^{(k)}.$$

Die Transformationsmatrix $\mathbf{L}_{k+1}(-\bar{\mathbf{l}}_{k+1})$ ist eine LNT-Matrix. Sie hat die Struktur

$$\mathbf{L}_{k+1}(-\bar{\mathbf{l}}_{k+1}) = \begin{pmatrix} 1 & & & & & & & & \\ & \ddots & & & & & & & \\ & & 1 & & & & & & \\ & & & 1 & & & & & \\ & & & & 1 & & & & \\ & & & & & -\bar{l}_{k+2,k+1} & 1 & & \\ & & & & & & & \ddots & \\ & & & & & & & & -\bar{l}_{n,k+1} & \\ & & & & & & & & & 1 \end{pmatrix} = \mathbf{I} - \bar{\mathbf{l}}_{k+1} \mathbf{e}_{k+1}^T$$

mit

$$\bar{\mathbf{l}}_{k+1} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \frac{a_{k+2,k+1}^{(k)}}{a_{k+1,k+1}^{(k)}} \\ \vdots \\ \frac{a_{n,k+1}^{(k)}}{a_{k+1,k+1}^{(k)}} \end{pmatrix}.$$

Unter der Voraussetzung $a_{k+1,k+1}^{(k)} \neq 0$ lässt sich der $(k+1)$ -te Transformationsschritt in der Form

$$\mathbf{A}^{(k+1)} = \mathbf{L}_{k+1}(-\bar{\mathbf{l}}_{k+1}) \mathbf{A}^{(k)}$$

schreiben. Gilt $a_{k+1,k+1}^{(k)} = 0$, so ist mindestens ein Element $a_{i,k+1}^{(k)}$ mit $i = k+2, \dots, n$ von Null verschieden. Wäre das nicht der Fall, so wären alle Elemente $a_{i,k+1}^{(k)}$ mit

$i = k + 1, \dots, n$ gleich Null. Damit wäre aber die Matrix $\mathbf{A}^{(k)}$ singulär. Da aber alle Transformationsschritte mit regulären LNT-Matrizen durchgeführt wurden, müsste dann auch die Matrix \mathbf{A} im Widerspruch zur Voraussetzung singulär sein. Es sei nun $s(k+1) \in \{k+2, \dots, n\}$ ein Index mit $a_{s(k+1),k+1}^{(k)} \neq 0$. Wir tauschen die Zeilen $k+1$ und $s(k+1)$ in der Matrix $\mathbf{A}^{(k)}$ und erhalten eine Matrix⁴

$$\hat{\mathbf{A}}^{(k)} = \mathbf{T}_{s(k+1),k+1} \mathbf{A}^{(k)}.$$

In diesem Falle ergibt sich $\mathbf{A}^{(k+1)}$ gemäß

$$\mathbf{A}^{(k+1)} = \mathbf{L}_{k+1}(-\hat{\mathbf{l}}_{k+1}) \hat{\mathbf{A}}^{(k)} = \mathbf{L}_{k+1}(-\hat{\mathbf{l}}_{k+1}) \mathbf{T}_{s(k+1),k+1} \mathbf{A}^{(k)}.$$

Der Vektor $\hat{\mathbf{l}}_{k+1}$ ist durch

$$\hat{\mathbf{l}}_{k+1} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \frac{\hat{a}_{k+2,k+1}^{(k)}}{\hat{a}_{k+1,k+1}^{(k)}} \\ \vdots \\ \frac{\hat{a}_{n,k+1}^{(k)}}{\hat{a}_{k+1,k+1}^{(k)}} \end{pmatrix}$$

festgelegt, wobei die $\hat{a}_{ij}^{(k)}$ die Elemente der Matrix $\hat{\mathbf{A}}^{(k)}$ bezeichnen. Nach $n-1$ Schritten erhält man so eine obere Dreiecksmatrix $\mathbf{U} = \mathbf{A}^{(n-1)}$. Es gilt

$$\mathbf{U} = \mathbf{L}_{n-1}(-\hat{\mathbf{l}}_{n-1}) \mathbf{T}_{s(n-1),n-1} \cdots \mathbf{L}_2(-\hat{\mathbf{l}}_2) \mathbf{T}_{s(2),2} \mathbf{L}_1(-\hat{\mathbf{l}}_1) \mathbf{T}_{s(1),1} \mathbf{A}.$$

Im ursprünglichen GAUSS-Algorithmus werden diese Transformationen simultan auf die rechte Seite angewendet. Man erhält

$$\mathbf{c} = \mathbf{L}_{n-1}(-\hat{\mathbf{l}}_{n-1}) \mathbf{T}_{s(n-1),n-1} \cdots \mathbf{L}_2(-\hat{\mathbf{l}}_2) \mathbf{T}_{s(2),2} \mathbf{L}_1(-\hat{\mathbf{l}}_1) \mathbf{T}_{s(1),1} \mathbf{b}.$$

Damit hat man das Gleichungssystem

$$\mathbf{A}\mathbf{x} = \mathbf{b}$$

auf das äquivalente System

$$\mathbf{U}\mathbf{x} = \mathbf{c}$$

⁴Das Element $a_{s(k+1),k+1}$, das nach diesem Tausch auf der Position $(k+1, k+1)$ steht wird als **Pivotelement** bezeichnet.

transformiert. Dies lässt sich wegen der oberen Dreiecksgestalt von U einfach lösen. Diesen Prozess bezeichnet man als **Rücksubstitution**. Einen Algorithmus dazu geben wir am Ende dieses Abschnitts an.

In vielen Anwendungsfällen hat man mehrere Gleichungssysteme mit gleicher Koeffizientenmatrix A aber verschiedenen rechten Seiten b zu lösen. Oft sind die rechten Seiten dieser zu lösenden Gleichungssysteme Funktionen der schon berechneten Lösungen (z.B. bei vielen Iterationsverfahren zur Nullstellenberechnung nichtlinearer Gleichungssysteme). In diesen Fällen ist es günstig, sich die Transformationen $L_i(-\hat{l}_i)T_{s(i),i}$ in geeigneter Weise zu merken. Dazu trennen die eigentlichen Transformationsmatrizen von den Tauschmatrizen. Allgemein gilt für $1 \leq k < i \leq j \leq n$ (siehe Übungsaufgabe 13):

$$T_{ij}L_k(l) = L_k(T_{ij}l)T_{ij}.$$

Wir verwenden die Abkürzung

$$P_i = T_{s(n-1),n-1}T_{s(n-2),n-2} \cdots T_{s(i),i}, \quad i = 1, \dots, n-1,$$

also

$$P_{n-1} = T_{s(n-1),n-1}, \quad P_i = P_{i+1}T_{s(i),i}, \quad i = n-2, \dots, 1.$$

Dann gilt auch für $1 \leq k < i \leq n$

$$P_iL_k(l) = L_k(P_i l)P_i.$$

Damit folgt

$$\begin{aligned} & L_{n-1}(-\hat{l}_{n-1})T_{s(n-1),n-1}L_{n-2}(-\hat{l}_{n-2})T_{s(n-2),n-2} \cdots L_1(-\hat{l}_1)T_{s(1),1} \\ &= L_{n-1}(-\hat{l}_{n-1})P_{n-1}L_{n-2}(-\hat{l}_{n-2})T_{s(n-2),n-2} \cdots L_1(-\hat{l}_1)T_{s(1),1} \\ &= L_{n-1}(-\hat{l}_{n-1})L_{n-2}(-P_{n-1}\hat{l}_{n-2})P_{n-2} \cdots L_1(-\hat{l}_1)T_{s(1),1} \\ &\vdots \\ &= L_{n-1}(-\hat{l}_{n-1})L_{n-2}(-P_{n-1}\hat{l}_{n-2}) \cdots L_2(-P_3\hat{l}_2)L_1(-P_2\hat{l}_1)P_1. \end{aligned}$$

Insgesamt gilt dann

$$U = L_{n-1}(-\hat{l}_{n-1})L_{n-2}(-P_{n-1}\hat{l}_{n-2}) \cdots L_2(-P_3\hat{l}_2)L_1(-P_2\hat{l}_1)P_1A.$$

Setzt man noch

$$P = P_1, \quad l_{n-1} = \hat{l}_{n-1}, \quad l_i = P_{i+1}\hat{l}_i, \quad i = 1, \dots, n-1,$$

so folgt

$$U = L_{n-1}(-l_{n-1})L_{n-2}(-l_{n-2})L_{n-3}(-l_{n-3}) \cdots L_2(-l_2)L_1(-l_1)PA.$$

Multipliziert man diese Gleichung von links mit den Inversen der LNT-Matrizen, so ergibt sich

$$PA = L_1(l_1)L_2(l_2) \cdots L_{n-2}(l_{n-2})L_{n-1}(l_{n-1})U.$$

Nach Übungsaufgabe 14 ist das Produkt der LNT-Matrizen eine untere Einsdreiecksmatrix, und es gilt

$$L = L_1(l_1)L_2(l_2) \cdots L_{n-1}(l_{n-1})L_{n-2}(l_{n-2}) = I + (l_1, l_2, \dots, l_{n-1}, \mathbf{o}).$$

Man erhält somit die Matrix L ohne zusätzlichen Rechenaufwand aus den Vektoren l_1, \dots, l_{n-1} . Damit haben wir konstruktiv folgenden Satz bewiesen.

8.17. Satz: *Zu jeder regulären (n, n) -Matrix A existiert eine Permutationsmatrix P derart, dass sich PA eindeutig in das Produkt aus einer unteren Einsdreiecksmatrix L und einer regulären oberen Dreiecksmatrix U zerlegen lässt:*

$$PA = LU$$

$$= \begin{pmatrix} 1 & & & & \\ l_{21} & 1 & & & \mathbf{O} \\ \vdots & \vdots & \ddots & & \\ l_{n-1,1} & l_{n-1,2} & \cdots & 1 & \\ l_{n,1} & l_{n,2} & \cdots & l_{n,n-1} & 1 \end{pmatrix} \begin{pmatrix} u_{11} & u_{12} & \cdots & u_{1,n-1} & u_{1n} \\ & u_{22} & \cdots & u_{2,n-1} & u_{2n} \\ & & \ddots & \vdots & \vdots \\ \mathbf{O} & & & u_{n-1,n-1} & u_{n-1,n} \\ & & & & u_{nn} \end{pmatrix}$$

mit $u_{ii} \neq 0$ für $i = 1, \dots, n$.

Beweis: Es ist noch die Eindeutigkeit der Zerlegung zu zeigen. Nehmen wir dazu an, es gelte

$$PA = LU = \bar{L}\bar{U}$$

mit von L und U verschiedenen Matrizen \bar{L} und \bar{U} . Dann folgt

$$\bar{L}^{-1}L = \bar{U}U^{-1} = D.$$

Nach Übungsaufgabe 10 ist $\bar{L}^{-1}L$ untere Einsdreiecksmatrix. Analog ist $\bar{U}U^{-1}$ obere Dreiecksmatrix. Damit ist die Matrix D gleichzeitig untere Einsdreiecksmatrix und obere Dreiecksmatrix, daher gilt $D = I$, daher $\bar{L} = L$ und $\bar{U} = U$. \ast

Hat man einmal von einer Matrix A eine LU -Zerlegung $A = P^T LU$ berechnet, so lässt sich jedes Gleichungssystem $Ax = b$ lösen, indem man es auf das Lösen von zwei einfacheren Systemen zurückführt. Man hat dabei die beiden Systeme

$$Ly = Pb, \quad Ux = y$$

zu lösen. Dabei wird das Lösen des ersten Systems als Vorwärtssubstitution und das Lösen des zweiten Systems, wie schon erwähnt, als Rückwärtssubstitution bezeichnet.

Der bei der Zerlegung der Matrix A freiwerdende Speicherplatz (alle Elemente unterhalb der Hauptdiagonalen) darf genutzt werden, um die Elemente der unteren Dreiecksmatrix L zu speichern. Die Zeilenvertauschungen (Permutationsmatrix P) merkt man sich auf einem Feld p der Länge n ; der Index $p(i)$ gibt dann die wahre Nummer der Zeile an, die aktuell auf der Position i steht.

8.18. LU -Zerlegung mit Pivotisierung und explizitem Zeilentausch:

Es ist die reguläre (n, n) -Matrix A in ein Produkt $PA = LU$ mit einer Permutationsmatrix P , einer unteren Dreiecksmatrix L und einer oberen Dreiecksmatrix U zu zerlegen. Die Permutationsmatrix P sei dabei auf einem Feld p der Länge n gespeichert.

```

{Initialisierung}
for  $i = 1$  to  $n$  do
     $p(i) = i$ 
endfor
{LU-Zerlegung}
for  $j = 1$  to  $n - 1$  do
    {Pivotsuche}
    Wähle einen Index  $i^* \in \{i \mid a_{ij} \neq 0, \quad i = j, \dots, n\}$ .
    {Zeilentausch}
    if  $i^* \neq j$  then
         $pp = p(i^*); p(i^*) = p(j); p(j) = pp$ 
        for  $k = 1$  to  $n$  do
             $aa = a_{i^*k}; a_{i^*k} = a_{jk}; a_{jk} = aa$ 
        endfor
    endif
    {Transformation der Restmatrix}
    for  $i = j + 1$  to  $n$  do
         $a_{ij} = a_{ij} / a_{jj}$ 
        for  $k = j + 1$  to  $n$  do
             $a_{ik} = a_{ik} - a_{ij} \cdot a_{jk}$ 
        endfor
    endfor

```

endfor
endfor

Aufwand: $\sim n^3/3$ Additionen/Multiplikationen

Am Ende des Algorithmus stehen die wesentlichen Elemente der Matrix L auf dem unteren Dreieck der Matrix PA . Das obere Dreieck einschließlich der Diagonalen wird von der Matrix U eingenommen:

$$\begin{pmatrix} u_{11} & u_{12} & u_{13} & \cdots & u_{1,n-2} & u_{1,n-1} & u_{1n} \\ l_{21} & u_{22} & u_{23} & \cdots & u_{2,n-2} & u_{2,n-1} & u_{2n} \\ l_{31} & l_{32} & u_{33} & \cdots & u_{3,n-2} & u_{3,n-1} & u_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ l_{n-2,1} & l_{n-2,2} & l_{n-2,3} & \cdots & u_{n-2,n-2} & u_{n-2,n-1} & u_{n-2,n} \\ l_{n-1,1} & l_{n-1,2} & l_{n-1,3} & \cdots & l_{n-1,n-2} & u_{n-1,n-1} & u_{n-1,n} \\ l_{n1} & l_{n2} & l_{n3} & \cdots & l_{n,n-2} & l_{n,n-1} & u_{nn} \end{pmatrix}.$$

8.19. Vorwärtssubstitution bei explizitem Zeilentausch:

Es ist das lineare Gleichungssystem $Ly = Pb$ mit der regulären unteren Einsdreiecksmatrix L , der Permutationsmatrix P und einem beliebigen Vektor $b \in \mathbb{R}^n$ zu lösen. Die Matrix L ist auf dem unteren Dreieck der Matrix A , die Permutationsmatrix ist auf einem Feld p der Länge n gespeichert.

for $j = 1$ **to** n **do**

$$y_j = b_{p(j)}$$

endfor

for $j = 1$ **to** n **do**

for $i = 1$ **to** $j - 1$ **do**

$$y_i = y_i - a_{ij} \cdot y_j$$

endfor

endfor

Aufwand: $\sim n^2/2$ Additionen/Multiplikationen

8.20. Rücksubstitution bei explizitem Zeilentausch:

Es ist das lineare Gleichungssystem $Ux = y$ mit der regulären oberen Dreiecksmatrix U und einem beliebigen Vektor $y \in \mathbb{R}^n$ zu lösen. Die Matrix U ist auf dem oberen Dreieck der Matrix A gespeichert.

for $j = n$ **to** 1 **step** -1 **do**

$$x_j = y_j / a_{jj}$$

for $i = 1$ **to** $j - 1$ **do**

$$y_i = y_i - a_{ij} \cdot x_j$$

endfor
endfor

Aufwand: $\sim n^2/2$ Additionen/Multiplikationen

Es ist nicht notwendig, den Zeilentauch bei der Zerlegung der Matrix explizit durchzuführen. In der Permutationsmatrix, bzw. im Permutationsvektor ist alle wesentliche Information enthalten. Dieses Vorgehen bezeichnet man als **fiktiven Zeilentauch**. Man erhält folgende Algorithmen.

8.21. LU-Zerlegung, Pivotisierung und fiktiver Zeilentauch:

Es ist die reguläre Matrix A in ein Produkt $PA = LU$ mit einer Permutationsmatrix P , einer unteren Einsdreiecksmatrix L und einer oberen Dreiecksmatrix U zu zerlegen. Die Permutationsmatrix P sei dabei auf einem Feld p der Länge n gespeichert.

{Initialisierung}

for $j = 1$ **to** n **do**

$p(i) = i$

endfor

{LU-Zerlegung}

for $j = 1$ **to** $n - 1$ **do**

{Pivotsuche}

Wähle einen Index

$$i^* \in \left\{ i \mid a_{p(i),j} \neq 0, \quad i = j, \dots, n \right\}.$$

{Fiktiver Zeilentauch}

if $i^* \neq j$ **then**

$pp = p(i^*); p(i^*) = p(j); p(j) = pp$

endif

{Transformation der Restmatrix}

for $i = j + 1$ **to** n **do**

$a_{p(i),j} = a_{p(i),j} / a_{p(j),j}$

for $k = j + 1$ **to** n **do**

$a_{p(i),k} = a_{p(i),k} - a_{p(i),j} \cdot a_{p(j),k}$

endfor

endfor

endfor

Aufwand: $\sim n^3/3$ Additionen/Multiplikationen

Nach diesem Algorithmus liefert das Feld p die Reihenfolge, in der die Zeilen der Matrix verarbeitet wurden: Die Zeile $p(i)$ wurde im i -ten Schritt als Pivotzeile verwendet. Die gleiche Reihenfolge der Zeilen und damit auch der Komponenten der rechten Seite ist natürlich bei der Vorwärtssubstitution zu beachten.

8.22. Vorwärtssubstitution bei fiktivem Zeilentausch:

Es ist das lineare Gleichungssystem $\mathbf{L}\mathbf{y} = \mathbf{P}\mathbf{b}$ mit der regulären unteren Einsdreiecksmatrix \mathbf{L} , der Permutationsmatrix \mathbf{P} und einem beliebigen Vektor $\mathbf{b} \in \mathbb{R}^n$ zu lösen. Die Matrix \mathbf{L} ist auf dem unteren Dreieck der Matrix $\mathbf{P}\mathbf{A}$, die Permutationsmatrix ist auf einem Feld p der Länge n gespeichert.

```

for  $i = 1$  to  $n$  do
     $y_i = b_{p(i)}$ 
    for  $j = 1$  to  $i - 1$  do
         $y_i = y_i - a_{p(i),j} \cdot y_j$ 
    endfor
endfor

```

Aufwand: $\sim n^2/2$ Additionen/Multiplikationen

8.23. Rücksubstitution bei fiktivem Zeilentausch:

Es ist das lineare Gleichungssystem $\mathbf{U}\mathbf{x} = \mathbf{y}$ mit der regulären oberen Dreiecksmatrix \mathbf{U} und einem beliebigen Vektor $\mathbf{y} \in \mathbb{R}^n$ zu lösen. Die Matrix \mathbf{U} ist auf dem oberen Dreieck der Matrix \mathbf{A} gespeichert.

```

for  $j = n$  to  $1$  step  $-1$  do
     $x_j = y_j / a_{p(j),j}$ 
    for  $i = 1$  to  $j - 1$  do
         $y_i = y_i - a_{p(i),j} \cdot x_j$ 
    endfor
endfor

```

Aufwand: $\sim n^2/2$ Additionen/Multiplikationen

Bemerkung: In den Algorithmen zur LU -Zerlegung müsste in der Praxis noch der Fall berücksichtigt werden, dass kein geeignetes Pivotelement vorhanden ist. Der Schritt **Pivotsuche** ist dazu etwa in folgender Form abzuändern:

Pivotsuche: Bestimme einen Index $i^* \in \{j, \dots, n\}$ mit $|a_{i^*,j}| \geq \varepsilon$. Gilt für alle $i \in \{j, \dots, n\}$ $|a_{i,j}| < \varepsilon$, so ist die Matrix numerisch singulär. STOPP

8.2.2. Rundungsfehleranalyse der LU-Zerlegung

Führen wir die LU -Zerlegung auf einem realen Rechner durch, so erhalten wir i. a. nicht die exakten Dreiecksfaktoren \mathbf{L} und \mathbf{U} . Durch den Einfluss von Rundungsfehlern ergeben sich nur Näherungen \mathcal{L} und \mathcal{U} . Falls die LU -Zerlegung durchführbar war (kein Abbruch wegen numerischer Singularität), ist die berechnete Matrix \mathcal{U} regulär. Die Matrix \mathcal{L} ist als Einsdreiecksmatrix ohnehin regulär. Wir fassen in diesem

Falle \mathcal{L} und \mathcal{U} als exakte Zerlegung einer gestörten Matrix $A + \delta A$ auf. Es gilt dann

$$\mathcal{P}^T \mathcal{L} \mathcal{U} = A + \delta A.$$

Die Matrix \mathcal{P} ist wieder eine Permutationsmatrix. Sie wird sich von der exakten Permutationsmatrix unterscheiden, da der Rundungsfehlereinfluss dazu führen kann, dass die Zeilen der Matrix in anderer Reihenfolge verarbeitet werden. Die Größe der Störung δA wollen wir nun abschätzen. Wir führen eine Rückwärtsanalyse der LU -Zerlegung durch. Dazu nehmen wir der Einfachheit halber an, dass während des Algorithmus kein Zeilentausch notwendig ist. Die Zeilen der Matrix seien daher schon in der Reihenfolge geordnet, in der sie verarbeitet werden. Dann gilt $\mathcal{P} = P = I$. Wir erhalten damit

$$\delta A = \mathcal{L} \mathcal{U} - A.$$

Schauen wir uns nun den k -ten Transformationsschritt an. Es gilt

$$A^{(k)} = L_k(-l_k) A^{(k-1)}.$$

In Blockschreibweise erhält man

$$A^{(k)} = \left(\begin{array}{c|c} I^{(k-1)} & O \\ \hline O & L_1(-\bar{l}_k) \end{array} \right) \left(\begin{array}{c|c} U^{(k-1)} & \\ \hline O & M^{(k-1)} \end{array} \right)$$

mit $L_1(-\bar{l}_k) = I - \bar{l}_k e_1^T \in \mathbb{R}^{(n-k+1) \times (n-k+1)}$ und

$$\bar{l}_k = \begin{pmatrix} 0 \\ l_{k+1,k} \\ \vdots \\ l_{nk} \end{pmatrix} \in \mathbb{R}^{n-k+1} \quad \text{falls} \quad l_k = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ l_{k+1,k} \\ \vdots \\ l_{nk} \end{pmatrix} \in \mathbb{R}^n.$$

$I^{(k-1)}$ bezeichnet die Einheitsmatrix im $\mathbb{R}^{(k-1) \times (k-1)}$. Aus der obigen Blockdarstellung folgt

$$A^{(k)} = \left(\begin{array}{c|c} U^{(k-1)} & \\ \hline O & \bar{M}^{(k-1)} \end{array} \right) = \left(\begin{array}{c|c} U^{(k-1)} & \\ \hline O & L_1(-\bar{l}_k) M^{(k-1)} \end{array} \right).$$

Die Transformation wirkt nur auf die Restmatrix $M^{(k-1)}$. Wir haben daher nur Transformationen der Art

$$\bar{M} = L_1(-l)M$$

mit

$$M = \begin{pmatrix} m_{11} & m_{12} & \cdots & m_{1r} \\ m_{21} & m_{22} & \cdots & m_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ m_{r1} & m_{r2} & \cdots & m_{rr} \end{pmatrix} \in \mathbb{R}^{r \times r},$$

$$L_1(-l) = I - l e_1^T \in \mathbb{R}^{r \times r}, \quad l = \begin{pmatrix} 0 \\ \frac{m_{21}}{m_{11}} \\ \vdots \\ \frac{m_{r1}}{m_{11}} \end{pmatrix} \in \mathbb{R}^r$$

zu betrachten. Führen wir die Transformation auf einem realen Rechner aus, so erhalten wir statt der Matrix \bar{M} eine Matrix \tilde{M} mit

$$\tilde{M} = \text{gl}(L_1(-l)M).$$

Diese Matrix fassen wir nun als exaktes Transformationsergebnis einer gestörten Matrix $\hat{M} = M + \delta M$ auf. Es sei also

$$\tilde{M} = L_1(-\hat{l})\hat{M}$$

mit

$$\hat{l} = \begin{pmatrix} 0 \\ \frac{\hat{m}_{21}}{\hat{m}_{11}} \\ \vdots \\ \frac{\hat{m}_{r1}}{\hat{m}_{11}} \end{pmatrix}.$$

Betrachten wir zunächst nur die zweite Darstellung von \tilde{M} . Es gilt

$$\begin{aligned} \tilde{M} &= L_1(-\hat{l})\hat{M} = (I - \hat{l} e_1^T)\hat{M} = \hat{M} - \hat{l} e_1^T \hat{M} \\ &= \begin{pmatrix} \hat{m}_{11} & \hat{m}_{12} & \cdots & \hat{m}_{1r} \\ \hat{m}_{21} & \hat{m}_{22} & \cdots & \hat{m}_{2r} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{m}_{r1} & \hat{m}_{r2} & \cdots & \hat{m}_{rr} \end{pmatrix} - \begin{pmatrix} 0 \\ \frac{\hat{m}_{21}}{\hat{m}_{11}} \\ \vdots \\ \frac{\hat{m}_{r1}}{\hat{m}_{11}} \end{pmatrix} (\hat{m}_{11}, \hat{m}_{12}, \dots, \hat{m}_{1r}), \end{aligned}$$

also

$$\tilde{M} = \begin{pmatrix} \hat{m}_{11} & \hat{m}_{12} & \cdots & \hat{m}_{1r} \\ 0 & \hat{m}_{22} - \frac{\hat{m}_{21}\hat{m}_{12}}{\hat{m}_{11}} & \cdots & \hat{m}_{2r} - \frac{\hat{m}_{21}\hat{m}_{1r}}{\hat{m}_{11}} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \hat{m}_{r2} - \frac{\hat{m}_{r1}\hat{m}_{12}}{\hat{m}_{11}} & \cdots & \hat{m}_{rr} - \frac{\hat{m}_{r1}\hat{m}_{1r}}{\hat{m}_{11}} \end{pmatrix}. \quad (8.1)$$

Andererseits werden die Elemente von \tilde{M} aus den Elementen von M nach den Formeln

$$\begin{aligned} \tilde{m}_{1j} &= m_{1j} && \text{für } j = 1, \dots, r \\ \tilde{m}_{i1} &= 0 && \text{für } i = 2, \dots, r \\ \tilde{m}_{ij} &= \text{gl} \left(m_{ij} - \frac{m_{i1}m_{1j}}{m_{11}} \right) && \text{für } i, j = 2, \dots, r \end{aligned} \quad (8.2)$$

berechnet. Ein Vergleich von 8.1 und 8.2 liefert:

- Erste Zeile von \tilde{M} :

$$\left. \begin{aligned} \hat{m}_{1j} &= m_{1j} \\ m_{1j} + \delta m_{1j} &= m_{1j} \\ \delta m_{1j} &= 0 \end{aligned} \right\}, \quad j = 1, \dots, r.$$

- Rest der ersten Spalte von \tilde{M} :

$$0 = 0.$$

- Restmatrix von \tilde{M} :

$$\hat{m}_{ij} - \frac{\hat{m}_{i1}\hat{m}_{1j}}{\hat{m}_{11}} = \text{gl} \left(m_{ij} - \frac{m_{i1}m_{1j}}{m_{11}} \right), \quad i, j = 2, \dots, r.$$

Wegen $\delta m_{1j} = 0$ für $j = 1, \dots, r$ folgt daraus

$$m_{ij} + \delta m_{ij} - (m_{i1} + \delta m_{i1}) \frac{m_{1j}}{m_{11}} = \left[m_{ij} - \frac{m_{i1}}{m_{11}} (1 + \varepsilon_i) m_{1j} (1 + \varrho_{ij}) \right] (1 + \vartheta_{ij})$$

für $i, j = 2, \dots, r$. Von den Rundungsfehlern nehmen wir die Gültigkeit von

$$|\varepsilon_i|, |\varrho_{ij}|, |\vartheta_{ij}| \leq \text{eps}, \quad i, j = 2, \dots, r$$

an.

Wir erhalten so $(r-1)^2$ Gleichungen zur Festlegung der restlichen $r(r-1)$ Störungen δm_{ij} , $i = 1, \dots, r$, $j = 2, \dots, r$. Es bleiben daher noch $r-1$ Freiheitsgrade bei der Belastung der einzelnen Element von M mit Störungen. Wir setzen

$$\delta m_{i1} = m_{i1}\varepsilon_i, \quad i = 2, \dots, r.$$

Damit gilt

$$|\delta m_{i1}| \leq |m_{i1}|\text{eps}, \quad i = 2, \dots, r$$

und

$$\text{gl}(l_i) = \text{gl}\left(\frac{m_{i1}}{m_{11}}\right) = \frac{m_{i1}}{m_{11}}(1 + \varepsilon_i) = \frac{m_{i1} + \delta m_{i1}}{m_{11}} = \hat{l}_i, \quad i = 2, \dots, r.$$

Es folgt

$$m_{ij} + \delta m_{ij} - \frac{m_{i1}m_{1j}}{m_{11}}(1 + \varepsilon_i) = \left[m_{ij} - \frac{m_{i1}}{m_{11}}(1 + \varepsilon_i)m_{1j}(1 + \varrho_{ij}) \right] (1 + \vartheta_{ij})$$

für $i, j = 2, \dots, r$. Diese Gleichungen lösen wir nach den δm_{ij} auf und erhalten:

$$\delta m_{ij} = m_{ij}\vartheta_{ij} - \frac{m_{i1}m_{1j}}{m_{11}}(1 + \varepsilon_i)[(1 + \varrho_{ij})(1 + \vartheta_{ij}) - 1].$$

In erster Näherung folgt daraus

$$\delta m_{ij} \doteq m_{ij}\vartheta_{ij} - \frac{m_{i1}m_{1j}}{m_{11}}(1 + \varepsilon_i)(\varrho_{ij} + \vartheta_{ij}) \quad (8.3)$$

$$\doteq m_{ij}\vartheta_{ij} - l_i m_{1j}(\varrho_{ij} + \vartheta_{ij}) \quad (8.4)$$

$$|\delta m_{ij}| \leq (|m_{ij}| + 2|l_i||m_{1j}|)\text{eps}. \quad (8.5)$$

In dieser Abschätzung taucht noch l_i als Element der exakten Transformationsmatrix, also ein Element, das wir eigentlich nicht kennen, auf. Es lässt sich aber durch bekannte Elemente ersetzen. Aus der Darstellung

$$\begin{aligned} \tilde{m}_{ij} &= m_{ij}(1 + \vartheta_{ij}) - \frac{m_{i1}m_{1j}}{m_{11}}(1 + \varepsilon_i)(1 + \varrho_{ij})(1 + \vartheta_{ij}) \\ &= m_{ij}(1 + \vartheta_{ij}) - l_i m_{1j}(1 + \varepsilon_i)(1 + \varrho_{ij})(1 + \vartheta_{ij}) \end{aligned}$$

erhalten wir

$$l_i m_{1j} = \frac{m_{ij}}{(1 + \varepsilon_i)(1 + \varrho_{ij})} - \frac{\tilde{m}_{ij}}{(1 + \varepsilon_i)(1 + \varrho_{ij})(1 + \vartheta_{ij})}. \quad (8.6)$$

Setzen wir dies in 8.3 ein, so ergibt sich

$$\begin{aligned}
\delta m_{ij} &\doteq m_{ij}\vartheta_{ij} - m_{ij} \frac{(1+\varepsilon_i)(\varrho_{ij} + \vartheta_{ij})}{(1+\varepsilon_i)(1+\varrho_{ij})} + \tilde{m}_{ij} \frac{(1+\varepsilon_i)(\varrho_{ij} + \vartheta_{ij})}{(1+\varepsilon_i)(1+\varrho_{ij})(1+\vartheta_{ij})} \\
&\doteq m_{ij}\vartheta_{ij} - m_{ij} \frac{\varrho_{ij} + \vartheta_{ij}}{(1+\varrho_{ij})} + \tilde{m}_{ij} \frac{\varrho_{ij} + \vartheta_{ij}}{(1+\varrho_{ij})(1+\vartheta_{ij})} \\
&\doteq m_{ij}\vartheta_{ij} - m_{ij}(\varrho_{ij} + \vartheta_{ij})(1-\varrho_{ij}) + \tilde{m}_{ij}(\varrho_{ij} + \vartheta_{ij})(1-\varrho_{ij})(1-\vartheta_{ij}) \\
&\doteq m_{ij}(\vartheta_{ij} - \varrho_{ij} - \vartheta_{ij}) + \tilde{m}_{ij}(\varrho_{ij} + \vartheta_{ij}) \\
&\doteq -\varrho_{ij}m_{ij} + (\varrho_{ij} + \vartheta_{ij})\tilde{m}_{ij} \\
|\delta m_{ij}| &\stackrel{\cdot}{\leq} \text{eps}(|m_{ij}| + 2|\tilde{m}_{ij}|).
\end{aligned}$$

Für die gesamte Matrix δM erhalten wir so

$$\begin{aligned}
\delta m_{1j} &= 0 && \text{für } j = 1, \dots, r, \\
|\delta m_{i1}| &\stackrel{\cdot}{\leq} \text{eps} |m_{i1}| && \text{für } i = 2, \dots, r, \\
|\delta m_{ij}| &\stackrel{\cdot}{\leq} \text{eps} (|m_{ij}| + 2|\tilde{m}_{ij}|) && \text{für } i, j = 2, \dots, r.
\end{aligned}$$

In Matrixschreibweise lautet dieses Ergebnis für die Transformation

$$\begin{aligned}
M^{(k-1)} &\longrightarrow M^{(k)} \\
|\delta M^{(k-1)}| &\stackrel{\cdot}{\leq} \text{eps} \left[|M^{(k-1)}| + 2 \begin{pmatrix} 0 & \mathbf{o}^T \\ \mathbf{o} & |M^{(k)}| \end{pmatrix} \right].
\end{aligned}$$

Bezüglich einer absoluten und monotonen Matrixnorm folgt daraus

$$\|\delta M^{(k-1)}\| \stackrel{\cdot}{\leq} \text{eps} \left(\|M^{(k-1)}\| + 2\|M^{(k)}\| \right).$$

Betrachten wir nun wieder einen Transformationsschritt für die gesamte Matrix, so interpretieren wir die aus $A^{(k-1)}$ auf einem realen Rechner berechnete Matrix $\mathcal{A}^{(k)}$ als exaktes Transformationsergebnis einer gestörten Matrix $\hat{A}^{(k-1)} = A^{(k-1)} + \delta A^{(k-1)}$. Für die Störung $\delta A^{(k-1)}$ gilt

$$\delta A^{(k-1)} = \begin{pmatrix} \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \delta M^{(k-1)} \end{pmatrix}$$

und

$$\|\delta A^{(k-1)}\| = \|\delta M^{(k-1)}\|.$$

Insgesamt gilt dann für die berechnete Matrix $\mathcal{A}^{(k)}$

$$\mathcal{A}^{(k)} = \mathcal{L}^{(k-1)}(\mathcal{A}^{(k-1)} + \delta\mathcal{A}^{(k-1)})$$

mit der berechneten Transformationsmatrix $\mathcal{L}^{(k-1)} = \mathbf{L}_{k-1}(-\bar{\mathbf{l}}_{k-1})$. Für den Transformationsvektor $\bar{\mathbf{l}}_{k-1}$ gilt

$$\begin{aligned} \bar{\mathbf{l}}_{k-1} &= \text{gl}(\mathbf{l}_{k-1}) \\ &= \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \text{gl}\left(\frac{a_{k,k-1}^{(k-1)}}{a_{k-1,k-1}^{(k-1)}}\right) \\ \vdots \\ \text{gl}\left(\frac{a_{n,k-1}^{(k-1)}}{a_{k-1,k-1}^{(k-1)}}\right) \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \text{gl}\left(\frac{m_{k,k-1}^{(k-1)}}{m_{k-1,k-1}^{(k-1)}}\right) \\ \vdots \\ \text{gl}\left(\frac{m_{n,k-1}^{(k-1)}}{m_{k-1,k-1}^{(k-1)}}\right) \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \frac{m_{k,k-1}^{(k-1)} + \delta m_{n,k-1}^{(k-1)}}{m_{k-1,k-1}^{(k-1)}} \\ \vdots \\ \frac{m_{n,k-1}^{(k-1)} + \delta m_{n,k-1}^{(k-1)}}{m_{k-1,k-1}^{(k-1)}} \end{pmatrix}. \end{aligned}$$

Die aus $\mathcal{A}^{(k-1)}$ berechnete Transformationsmatrix $\mathcal{L}^{(k-1)}$ ist demnach als exakte Transformationsmatrix zur gestörten Matrix $\mathcal{A}^{(k-1)} + \delta\mathcal{A}^{(k-1)}$ interpretierbar. Für die berechnete Matrix \mathcal{U} erhalten wir dann

$$\begin{aligned} \mathcal{U} = \mathcal{A}^{(n-1)} &= \mathcal{L}^{(n-1)}(\mathcal{A}^{(n-2)} + \delta\mathcal{A}^{(n-2)}) \\ &= \mathcal{L}^{(n-1)}\left[\mathcal{L}^{(n-2)}(\mathcal{A}^{(n-3)} + \delta\mathcal{A}^{(n-3)}) + \delta\mathcal{A}^{(n-2)}\right] \\ &= \mathcal{L}^{(n-1)}\mathcal{L}^{(n-2)}\left(\mathcal{A}^{(n-3)} + \delta\mathcal{A}^{(n-3)} + \left(\mathcal{L}^{(n-2)}\right)^{-1}\delta\mathcal{A}^{(n-2)}\right). \end{aligned}$$

Die Matrizen $\left(\mathcal{L}^{(n-2)}\right)^{-1}$ und $\delta\mathcal{A}^{(n-2)}$ sind von folgender Struktur

$$\left(\mathcal{L}^{(n-2)}\right)^{-1} = \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & & \star & \\ \mathbf{O} & & & & 1 \end{pmatrix}, \quad \delta\mathcal{A}^{(n-2)} = \begin{pmatrix} \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \star \star \\ & \star \star \end{pmatrix}.$$

Man erkennt sofort, dass $(\mathcal{L}^{(n-2)})^{-1} \delta \mathbf{A}^{(n-2)} = \delta \mathbf{A}^{(n-2)}$ gilt. Damit folgt

$$\begin{aligned} \mathbf{u} &= \mathcal{L}^{(n-1)} \mathcal{L}^{(n-2)} \left(\mathbf{A}^{(n-3)} + \delta \mathbf{A}^{(n-3)} + \delta \mathbf{A}^{(n-2)} \right) \\ &= \mathcal{L}^{(n-1)} \mathcal{L}^{(n-2)} \left[\mathcal{L}^{(n-3)} \left(\mathbf{A}^{(n-4)} + \delta \mathbf{A}^{(n-4)} \right) + \delta \mathbf{A}^{(n-3)} + \delta \mathbf{A}^{(n-2)} \right] \\ &= \mathcal{L}^{(n-1)} \mathcal{L}^{(n-2)} \mathcal{L}^{(n-3)} \\ &\quad \left[\mathbf{A}^{(n-4)} + \delta \mathbf{A}^{(n-4)} + \left(\mathcal{L}^{(n-3)} \right)^{-1} \left(\delta \mathbf{A}^{(n-3)} + \delta \mathbf{A}^{(n-2)} \right) \right]. \end{aligned}$$

Wegen

$$\left(\mathcal{L}^{(n-3)} \right)^{-1} = \begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & 1 & & \\ & & * & 1 & \\ \mathbf{O} & & * & 0 & 1 \\ & & * & 0 & 0 & 1 \end{pmatrix}, \quad \delta \mathbf{A}^{(n-3)} + \delta \mathbf{A}^{(n-2)} = \begin{pmatrix} \mathbf{O} & \mathbf{O} \\ & * & * & * \\ \mathbf{O} & * & * & * \\ & * & * & * \end{pmatrix}$$

gilt wieder

$$\left(\mathcal{L}^{(n-3)} \right)^{-1} \left(\delta \mathbf{A}^{(n-3)} + \delta \mathbf{A}^{(n-2)} \right) = \delta \mathbf{A}^{(n-3)} + \delta \mathbf{A}^{(n-2)}.$$

Setzt man diesen Prozess fort, so erhält man

$$\mathbf{u} = \mathcal{L}^{(n-1)} \mathcal{L}^{(n-2)} \dots \mathcal{L}^{(1)} \left[\mathbf{A} + \delta \mathbf{A}^{(0)} + \dots + \delta \mathbf{A}^{(n-3)} + \delta \mathbf{A}^{(n-2)} \right].$$

Da sich

$$\mathcal{L} = \left(\mathcal{L}^{(1)} \right)^{-1} \dots \left(\mathcal{L}^{(n-1)} \right)^{-1}$$

ohne weitere Rundungsfehler aus den $\mathcal{L}^{(i)}$ ergibt, folgt

$$\mathcal{L} \mathbf{u} = \mathbf{A} + \delta \mathbf{A}, \quad \delta \mathbf{A} = \delta \mathbf{A}^{(0)} + \delta \mathbf{A}^{(1)} + \dots + \delta \mathbf{A}^{(n-2)}.$$

Bezüglich einer absoluten und monotonen Matrixnorm ist dann $\delta \mathbf{A}$ abschätzbar. Wir erhalten

$$\begin{aligned} \|\delta \mathbf{A}\| &\leq \sum_{i=0}^{n-2} \|\delta \mathbf{A}^{(i)}\| = \sum_{i=0}^{n-2} \|\delta \mathbf{M}^{(i)}\| \leq \text{eps} \sum_{i=0}^{n-2} \left(\|\mathbf{M}^{(i)}\| + 2\|\mathbf{M}^{(i+1)}\| \right) \\ &= \text{eps} \left[\|\mathbf{M}^{(0)}\| + \sum_{i=1}^{n-2} \|\mathbf{M}^{(i)}\| + 2 \sum_{i=1}^{n-1} \|\mathbf{M}^{(i)}\| \right] \\ &\leq \text{eps} \left[\|\mathbf{A}\| + 3 \sum_{i=1}^{n-1} \|\mathbf{M}^{(i)}\| \right]. \end{aligned}$$

Damit haben wir den folgenden Satz bewiesen.

8.24. Satz: *Der GAUSSsche Algorithmus sei für die (n, n) -Matrix \mathbf{A} durchführbar. Die berechneten Dreiecksfaktoren seien \mathcal{L} und \mathcal{U} . Dann existiert eine Störung $\delta\mathbf{A}$ mit*

$$\mathcal{L}\mathcal{U} = \mathbf{A} + \delta\mathbf{A}.$$

Die Störung genügt der Abschätzung

$$\|\delta\mathbf{A}\| \leq \text{eps } F(\mathbf{A})\|\mathbf{A}\|$$

mit der Fehlerkonstanten

$$F(\mathbf{A}) = 1 + 3 \sum_{i=1}^{n-1} \frac{\|\mathbf{M}^{(i)}\|}{\|\mathbf{A}\|}.$$

Der GAUSSsche Algorithmus ist nach diesem Satz gutartig, falls die Kondition der Matrizen $\mathbf{M}^{(i)}$ nicht beliebig groß wird. Ist \mathcal{A} eine Klasse von Matrizen, für die der GAUSSsche Algorithmus durchführbar ist und existiert eine Konstante $F < \infty$, so dass für alle Matrizen $\mathbf{A} \in \mathcal{A}$ die Ungleichung $F(\mathbf{A}) \leq F$ gilt, so ist die LU -Zerlegung auf \mathcal{A} numerisch gutartig.

Neben den Rundungsfehlern, die beim Berechnen der LU -Zerlegung auftreten, sind auch jene Rundungsfehler zu beachten, die beim Auflösen der Dreieckssysteme entstehen. Wir werden für die Vorwärts- und die Rücksubstitution wieder eine Rückwärtsanalyse durchführen. Es seien $\hat{\mathbf{y}}$ und $\hat{\mathbf{x}}$ die berechneten Lösungen der Gleichungssysteme

$$\mathcal{L}\mathbf{y} = \mathbf{b}, \quad \mathcal{U}\mathbf{x} = \mathbf{y}.$$

Wir nehmen nun an, dass $\hat{\mathbf{y}}$ und $\hat{\mathbf{x}}$ als exakte Lösungen gestörter Systeme interpretierbar sind. Es gelte also

$$(\mathcal{L} + \delta\mathcal{L})\hat{\mathbf{y}} = \mathbf{b}, \quad (\mathcal{U} + \delta\mathcal{U})\hat{\mathbf{x}} = \hat{\mathbf{y}}.$$

Wir beweisen zunächst den folgenden etwas allgemeineren Satz.

8.25. Satz: *Es sei $\hat{\mathbf{x}}$ die berechnete Lösung des Gleichungssystems*

$$\mathbf{L}\mathbf{x} = \mathbf{b}$$

mit der regulären unteren Dreiecksmatrix \mathbf{L} . Dann existiert eine Störung $\delta\mathbf{L}$ mit

$$(\mathbf{L} + \delta\mathbf{L})\hat{\mathbf{x}} = \mathbf{b}.$$

Die Störung δL ist ebenfalls untere Dreiecksmatrix und genügt der elementweisen Abschätzung

$$|\delta l_{ij}| \leq j \text{eps} |l_{ij}|, \quad i, j = 1, \dots, n.$$

Beweis: Das Lösen des Gleichungssystems $Lx = b$ erfolgt nach dem Algorithmus

```

for  $i = 1$  to  $n$  do
   $x_i = b_i$ 
  for  $j = 1$  to  $i - 1$  do
     $x_i = x_i - l_{ij} \cdot x_j$ 
  end
   $x_i = x_i / l_{ii}$ 
end

```

In Gleitpunktarithmetik gilt dann für das Berechnen einer Komponente \hat{x}_i :

```

 $\hat{x}_i = b_i$ 
for  $j = 1$  to  $i - 1$  do
   $\hat{x}_i = \text{gl}(\hat{x}_i - l_{ij} \hat{x}_j) = (\hat{x}_i - l_{ij} \hat{x}_j (1 + \varepsilon_j))(1 + \delta_j)$ 
end
 $\hat{x}_i = \text{gl}(\hat{x}_i / l_{ii}) = \hat{x}_i / (l_{ii} (1 + \varepsilon_i))$ 

```

Dabei nehmen wir wieder an, dass $|\varepsilon_j|, |\delta_j| \leq \text{eps}$ für $j = 1, \dots, i$ gilt. Damit folgt

$$\begin{aligned}
 \hat{x}_i &= \frac{\left(\dots \left((b_i - l_{i1} \hat{x}_1 (1 + \varepsilon_1))(1 + \delta_1) - l_{i2} \hat{x}_2 (1 + \varepsilon_2) \right) (1 + \delta_2) - \dots - l_{i,i-1} \hat{x}_{i-1} (1 + \varepsilon_{i-1}) \right) (1 + \delta_{i-1})}{l_{ii} (1 + \varepsilon_i)} \\
 &= \frac{b_i (1 + \delta_1) (1 + \delta_2) \dots (1 + \delta_{i-1}) - l_{i1} \hat{x}_1 (1 + \varepsilon_1) (1 + \delta_1) (1 + \delta_2) \dots (1 + \delta_{i-1})}{l_{ii} (1 + \varepsilon_i)} + \\
 &\quad + \frac{-l_{i2} \hat{x}_2 (1 + \varepsilon_2) (1 + \delta_2) \dots (1 + \delta_{i-1}) - \dots - l_{i,i-1} \hat{x}_{i-1} (1 + \varepsilon_{i-1}) (1 + \delta_{i-1})}{l_{ii} (1 + \varepsilon_i)} \\
 &= \frac{b_i - l_{i1} \hat{x}_1 (1 + \varepsilon_1) - l_{i2} \hat{x}_2 \frac{1 + \varepsilon_2}{1 + \delta_1} - \dots - l_{i,i-1} \hat{x}_{i-1} \frac{1 + \varepsilon_{i-1}}{(1 + \delta_1)(1 + \delta_2) \dots (1 + \delta_{i-2})}}{l_{ii} \frac{1 + \varepsilon_i}{(1 + \delta_1)(1 + \delta_2) \dots (1 + \delta_{i-1})}}.
 \end{aligned}$$

Nun gilt in erster Näherung

$$\frac{1 + \varepsilon_j}{(1 + \delta_1)(1 + \delta_2) \dots (1 + \delta_{j-1})} \doteq 1 + \varepsilon_j - \delta_1 - \delta_2 - \dots - \delta_{j-1}.$$

Damit folgt

$$\hat{x}_i = \frac{b_i - l_{i1} (1 + \beta_{i1}) \hat{x}_1 - l_{i2} (1 + \beta_{i2}) \hat{x}_2 - \dots - l_{i,i-1} (1 + \beta_{i1}) \hat{x}_{i-1}}{l_{ii} (1 + \beta_{ii})} \quad (8.7)$$

mit

$$\begin{aligned}\beta_{ij} &= \frac{1 + \varepsilon_j}{(1 + \delta_1)(1 + \delta_2) \cdots (1 + \delta_{j-1})} \\ &\doteq 1 + \varepsilon_j - \delta_1 - \delta_2 - \cdots - \delta_{j-1}\end{aligned}$$

und

$$|\beta_{ij}| \stackrel{\cdot}{\leq} j\text{eps}.$$

Interpretieren wir die berechnete Lösung \hat{x} als exakte Lösung eines gestörten Systems, so erhalten wir

$$(\mathbf{L} + \delta\mathbf{L})\hat{x} = \mathbf{b}.$$

Hieraus folgt

$$\hat{x}_i = \frac{b_i - (l_{i1} + \delta l_{i1})\hat{x}_1 - (l_{i2} + \delta l_{i2})\hat{x}_2 - \cdots - (l_{i,i-1} + \delta l_{i,i-1})\hat{x}_{i-1}}{l_{ii} + \delta l_{ii}}. \quad (8.8)$$

Ein Vergleich von 8.7 und 8.8 liefert

$$\delta l_{ij} = l_{ij}\beta_{ij}$$

und

$$|\delta l_{ij}| \stackrel{\cdot}{\leq} j\text{eps}|l_{ij}|$$

für $i, j = 1, \dots, n$. *

Analog gilt für die Lösung von Gleichungssystemen mit einer oberen Dreiecksmatrix der folgende Satz.

8.26. Satz: *Es sei \hat{x} die berechnete Lösung des Gleichungssystems*

$$\mathbf{U}\mathbf{x} = \mathbf{b}$$

mit der regulären oberen Dreiecksmatrix \mathbf{U} . Dann existiert eine Störung $\delta\mathbf{U}$ mit

$$(\mathbf{U} + \delta\mathbf{U})\hat{x} = \mathbf{b}.$$

Die Störung $\delta\mathbf{U}$ ist ebenfalls obere Dreiecksmatrix und genügt der elementweisen Abschätzung

$$|\delta u_{ij}| \leq (n - j + 1)\text{eps}|u_{ij}|, \quad i, j = 1, \dots, n.$$

Mit Hilfe der Sätze 8.24, 8.25 und 8.26 lässt sich der Gesamtrundungsfehler beim Lösen eines linearen Gleichungssystems mit Hilfe einer LU -Zerlegung abschätzen. Die berechnete Lösung \hat{x} ist exakte Lösung des Systems

$$(\mathcal{U} + \delta\mathcal{U})\hat{x} = \hat{y}.$$

\hat{y} ist exakte Lösung von

$$(\mathcal{L} + \delta\mathcal{L})\hat{y} = \mathbf{b}.$$

Insgesamt gilt

$$(\mathcal{L} + \delta\mathcal{L})(\mathcal{U} + \delta\mathcal{U})\hat{x} = \hat{y}, \quad (8.9)$$

wobei die Störungen $\delta\mathcal{L}$ und $\delta\mathcal{U}$ den Abschätzungen aus Satz 8.25 bzw. Satz 8.26 genügen. Multipliziert man in Gleichung 8.9 die linke Seite aus, so ergibt sich

$$\begin{aligned} (\mathcal{L}\mathcal{U} + \mathcal{L}\delta\mathcal{U} + \delta\mathcal{L}(\mathcal{U} + \delta\mathcal{U}))\hat{x} &= \mathbf{b} \\ (\mathbf{A} + \delta\mathbf{A} + \delta\mathbf{A}')\hat{x} &= \mathbf{b}. \end{aligned}$$

Die Störung $\delta\mathbf{A}$ beinhaltet den Rundungsfehlereinfluss der LU -Zerlegung. Für sie gilt nach Satz 8.24

$$|\delta\mathbf{A}| \stackrel{\cdot}{\leq} \text{eps} F(\mathbf{A}) \|\mathbf{A}\|.$$

In der Störung $\delta\mathbf{A}'$ steckt der Rundungsfehlereinfluss beim Lösen der Dreieckssysteme. Für sie ergibt sich elementweise

$$\begin{aligned} \delta a'_{ij} &= \sum_{k=1}^{\min\{i,j\}} [\ell_{ik} \delta u_{kj} + \delta \ell_{ik} (u_{kj} + \delta u_{kj})] \\ |\delta a'_{ij}| &= \left| \sum_{k=1}^{\min\{i,j\}} [\ell_{ik} \delta u_{kj} + \delta \ell_{ik} (u_{kj} + \delta u_{kj})] \right| \\ &\stackrel{\cdot}{\leq} \sum_{k=1}^{\min\{i,j\}} [|\ell_{ik}| |\delta u_{kj}| + |\delta \ell_{ik}| (|u_{kj}| + |\delta u_{kj}|)] \end{aligned}$$

Wendet man die Abschätzungen aus den Sätzen 8.25 und 8.26 an, so folgt

$$\begin{aligned} |\delta a'_{ij}| &\stackrel{\cdot}{\leq} \text{eps} \sum_{k=1}^{\min\{i,j\}} [(n-j+1)|\ell_{ik}| |u_{kj}| + k|\ell_{ik}| (|u_{kj}| + (n-j+1)\text{eps} |u_{kj}|)] \\ &\stackrel{\cdot}{\leq} \text{eps} \sum_{k=1}^{\min\{i,j\}} |\ell_{ik}| |u_{kj}| [n-j+1 + k + k(n-j+1)\text{eps}] \\ &\stackrel{\cdot}{\leq} \text{eps} \sum_{k=1}^{\min\{i,j\}} |\ell_{ik}| |u_{kj}| (n-j+1 + k). \end{aligned}$$

Wegen $k \leq \min\{i, j\} \leq j$ gilt $n - j + 1 + k \leq n + 1$. Damit ergibt sich

$$|\delta a'_{ij}| \leq \text{eps}(n+1) \sum_{k=1}^{\min\{i,j\}} |\ell_{ik}| |u_{kj}|.$$

Analog zu Gleichung 8.6 gilt für

$$\ell_{ik} u_{kj} = \ell_{ik} a_{kj}^{(k-1)}$$

die Abschätzung

$$|\ell_{ik}| |u_{kj}| \leq |a_{ij}^{(k)}| + |a_{ij}^{(k-1)}|.$$

Damit erhalten wir

$$\begin{aligned} |\delta a'_{ij}| &\leq \text{eps}(n+1) \left[\sum_{k=1}^{\min\{i,j\}} |a_{ij}^{(k)}| + \sum_{k=1}^{\min\{i,j\}} |a_{ij}^{(k-1)}| \right] \\ &\leq \text{eps}(n+1) \left[|a_{ij}^{(0)}| + \sum_{k=1}^{\min\{i,j\}} |a_{ij}^{(k)}| + \sum_{k=1}^{\min\{i,j\}-1} |a_{ij}^{(k)}| \right] \\ &\leq \text{eps}(n+1) \left[|a_{ij}^{(0)}| + 2 \sum_{k=1}^{\min\{i,j\}} |a_{ij}^{(k)}| \right]. \end{aligned}$$

In Matrixschreibweise gilt

$$|\delta \mathbf{A}'| \leq \text{eps}(n+1) \left[|\mathbf{A}| + 2 \sum_{k=1}^{n-1} \begin{pmatrix} \mathbf{O} & \mathbf{O} \\ \mathbf{O} & |\mathbf{M}^{(k)}| \end{pmatrix} \right],$$

und bezüglich einer absoluten und monotonen Matrixnorm

$$\|\delta \mathbf{A}'\| \leq \text{eps}(n+1) \left(\|\mathbf{A}\| + 2 \sum_{k=1}^{n-1} \|\mathbf{M}^{(k)}\| \right).$$

Damit haben wir folgenden Satz bewiesen.

8.27. Satz: Der GAUSSsche Algorithmus sei für eine (n, n) -Matrix \mathbf{A} durchführbar. $\bar{\mathbf{x}} \in \mathbb{R}^n$ sei die berechnete Lösung des Gleichungssystems $\mathbf{A}\mathbf{x} = \mathbf{b}$. Dann existiert eine Störung $\overline{\delta \mathbf{A}}$, so dass $\bar{\mathbf{x}}$ exakte Lösung des gestörten Systems

$$(\mathbf{A} + \overline{\delta \mathbf{A}})\bar{\mathbf{x}} = \mathbf{b}$$

ist. Für die Störung $\overline{\delta\mathbf{A}}$ gilt

$$\overline{\delta\mathbf{A}} = \delta\mathbf{A} + \delta\mathbf{A}',$$

wobei $\delta\mathbf{A}$ den Rundungsfehleranteil der LU-Zerlegung und $\delta\mathbf{A}'$ den Rundungsfehleranteil beim Lösen der Dreieckssysteme beschreibt. Es gelten die Abschätzungen

$$\|\delta\mathbf{A}\| \leq \text{eps } F(\mathbf{A}) \|\mathbf{A}\|$$

mit

$$F(\mathbf{A}) \doteq 1 + 3 \sum_{k=1}^{n-1} \frac{\|\mathbf{M}^{(k)}\|}{\|\mathbf{A}\|}$$

und

$$\|\delta\mathbf{A}'\| \leq \text{eps } F'(\mathbf{A}) \|\mathbf{A}\|$$

mit

$$F'(\mathbf{A}) \doteq (n+1) \left(1 + 2 \sum_{k=1}^{n-1} \frac{\|\mathbf{M}^{(k)}\|}{\|\mathbf{A}\|} \right) \leq (n+1)F(\mathbf{A}).$$

Insgesamt gilt

$$\|\overline{\delta\mathbf{A}}\| \leq \text{eps } \bar{F}(\mathbf{A}) \|\mathbf{A}\|$$

mit

$$\bar{F}(\mathbf{A}) \doteq (n+2)F(\mathbf{A}).$$

Für die Durchführbarkeit des Algorithmus ist

$$\text{eps cond}(\mathbf{A}) F(\mathbf{A}) < 1$$

hinreichend.

Betrachten wir nun noch den erzeugten Rundungsfehler $\delta\mathbf{x} = \bar{\mathbf{x}} - \mathbf{x}$. Für ihn gilt

$$\begin{aligned} \mathbf{A}\delta\mathbf{x} &= \mathbf{A}\bar{\mathbf{x}} - \mathbf{A}\mathbf{x} = (\mathbf{A} + \overline{\delta\mathbf{A}} - \overline{\delta\mathbf{A}})\bar{\mathbf{x}} - \mathbf{b} = (\mathbf{A} + \overline{\delta\mathbf{A}})\bar{\mathbf{x}} - \overline{\delta\mathbf{A}}\bar{\mathbf{x}} - \mathbf{b} \\ &= -\overline{\delta\mathbf{A}}\bar{\mathbf{x}} \quad \delta\mathbf{x} = -\mathbf{A}^{-1}\overline{\delta\mathbf{A}}\bar{\mathbf{x}}. \end{aligned}$$

Damit folgt

$$\|\delta\mathbf{x}\| \leq \|\mathbf{A}^{-1}\| \|\overline{\delta\mathbf{A}}\| \|\bar{\mathbf{x}}\| = \kappa \|\bar{\mathbf{x}}\|$$

mit $\kappa = \|A^{-1}\| \|\overline{\delta A}\|$. Für $\kappa < 1$ ist $A + \overline{\delta A}$ regulär. Dann ergibt sich

$$\begin{aligned} (A + \overline{\delta A})\delta x &= (A + \overline{\delta A})\bar{x} - (A + \overline{\delta A})x = b - Ax - \overline{\delta A}x \\ &= -\overline{\delta A}x, \quad \delta x = -(A + \overline{\delta A})^{-1}\overline{\delta A}x \end{aligned}$$

und weiter

$$\|\delta x\| = \|(A + \overline{\delta A})^{-1}\| \|\overline{\delta A}\| \|x\| \leq \frac{\|A^{-1}\| \|\overline{\delta A}\|}{1 - \kappa} \|x\| = \frac{\kappa}{1 - \kappa} \|x\|.$$

Bemerkung: Im allgemeinen sind die Schätzungen für den Rundungsfehlereinfluss beim Lösen der Dreieckssysteme zu pessimistisch. Der erzeugte Rundungsfehler ist oft wesentlich kleiner als die Schranke $\|A^{-1}\| \|\overline{\delta A}\| \|\bar{x}\|$. Praktische Erfahrungen zeigen, dass die Konstante in Satz 8.27 durch

$$\bar{F} \approx \begin{cases} F, & \text{falls } \text{cond}(A) \text{ groß} \\ nF, & \text{falls } \text{cond}(A) \text{ klein} \end{cases}$$

ersetzbar ist.

8.2.3. Pivotisierung und Skalierung

Im Schritt **Pivotsuche** der LU -Zerlegung wurde gesichert, dass das Pivotelement a_{i^*j} von Null verschieden ist. Um ein günstiges Fehlerverhalten zu realisieren, ist es nach Satz 8.24 notwendig, die Normen der Matrizen $M^{(i)}$ möglichst klein zu halten. Dies lässt sich durch geschickte Wahl des Pivotelements erreichen. Es gibt im wesentlichen drei Varianten. Dabei müsste der Schritt **Pivotsuche** im Algorithmus zur LU -Zerlegung in folgender Weise abgeändert werden.

- **Spaltenpivotisierung:**

Bestimme einen Index $i^* \in \{j, \dots, n\}$, so dass die Ungleichung $|a_{i^*j}| \geq |a_{ij}|$ für alle $i \in \{j, \dots, n\}$ gilt.

Als Pivotelement wird das betragsgrößte Element der Restspalte j genommen.

- **Zeilenpivotisierung:**

Bestimme einen Index $j^* \in \{j, \dots, n\}$, so dass die Ungleichung $|a_{jj^*}| \geq |a_{jk}|$ für alle $k \in \{j, \dots, n\}$ gilt.

Als Pivotelement wird das betragsgrößte Element der Restzeile j genommen. Danach erfolgt ein Tausch der Spalten j und j^* , der sich in einem zusätzlichen Permutationsvektor zu merken ist.

- **Totalpivotisierung:**

Bestimme Indizes $i^*, j^* \in \{j, \dots, n\}$, so dass die Ungleichung $|a_{i^*j^*}| \geq |a_{ik}|$ für

alle $i, k \in \{j, \dots, n\}$ gilt.

Als Pivotelement wird das betragsgrößte Element der Restmatrix genommen. Hier ist ein Zeilen- und ein Spaltentausch vorzunehmen.

Für einen Zerlegungsalgorithmus mit Spaltenpivotisierung wollen wir die Normen der Restmatrizen $\mathbf{M}^{(k)}$ abschätzen. Dazu betrachten wir wieder einen Transformationsschritt $\mathbf{A}^{(k-1)} \rightarrow \mathbf{A}^{(k)}$. Es gilt

$$a_{ij}^{(k)} = \hat{a}_{ij}^{(k-1)} - \hat{l}_{ik}^{(k-1)} \hat{a}_{kj}^{(k-1)}, \quad i, j = k+1, \dots, n$$

mit

$$\hat{l}_{ik}^{(k-1)} = \frac{\hat{a}_{ik}^{(k-1)}}{\hat{a}_{kk}^{(k-1)}}.$$

Mit $\hat{a}_{ij}^{(k-1)}$ bezeichnen wir die Elemente der Matrix $\hat{\mathbf{A}}^{(k)} = \mathbf{T}_{s(k)k} \mathbf{A}^{(k)}$. Da bei der Spaltenpivotisierung das Pivotelement $\hat{a}_{kk}^{(k-1)}$ als betragsgrößtes der Restspalte gewählt wird, gilt

$$|\hat{l}_{ik}^{(k-1)}| \leq 1, \quad i = k+1, \dots, n.$$

Damit folgt

$$\left| a_{ij}^{(k)} \right| \leq \left| \hat{a}_{ij}^{(k-1)} - \hat{l}_{ik}^{(k-1)} \hat{a}_{kj}^{(k-1)} \right| \leq \left| \hat{a}_{ij}^{(k-1)} \right| + \left| \hat{a}_{kj}^{(k-1)} \right|, \quad i, j = k+1, \dots, n.$$

Nun lässt sich die ∞ -Norm (Zeilensummennorm) von $\mathbf{M}^{(k)}$ abschätzen. Wir erhalten:

$$\begin{aligned} \|\mathbf{M}^{(k)}\|_{\infty} &= \max_{i=k+1, \dots, n} \left\{ \sum_{j=k+1}^n \left| a_{ij}^{(k)} \right| \right\} \\ &\leq \max_{i=k+1, \dots, n} \left\{ \sum_{j=k+1}^n \left(\left| \hat{a}_{ij}^{(k-1)} \right| + \left| \hat{a}_{kj}^{(k-1)} \right| \right) \right\} \\ &\leq \max_{i=k, \dots, n} \left\{ \sum_{j=k}^n \left| \hat{a}_{ij}^{(k-1)} \right| + \sum_{j=k}^n \left| \hat{a}_{kj}^{(k-1)} \right| \right\} \\ &= \max_{i=k, \dots, n} \left\{ \sum_{j=k}^n \left| a_{ij}^{(k-1)} \right| + \sum_{j=k}^n \left| a_{s(k)j}^{(k-1)} \right| \right\} \\ &= \|\mathbf{M}^{(k-1)}\|_{\infty} + \sum_{j=k}^n \left| a_{s(k)j}^{(k-1)} \right| \leq 2 \|\mathbf{M}^{(k-1)}\|_{\infty}. \end{aligned}$$

Durch fortlaufende Anwendung der letzten Ungleichung erhält man

$$\|\mathbf{M}^{(k)}\|_{\infty} \leq 2\|\mathbf{M}^{(k-1)}\|_{\infty} \leq 4\|\mathbf{M}^{(k-2)}\|_{\infty} \leq \dots \leq 2^k \|\mathbf{M}^{(0)}\|_{\infty} = 2^k \|\mathbf{A}\|_{\infty}.$$

Mit Satz 8.24 folgt dann

$$\|\delta\mathbf{A}\|_{\infty} \leq \text{eps} F(\mathbf{A}) \|\mathbf{A}\|_{\infty}$$

mit

$$\begin{aligned} F(\mathbf{A}) &= 1 + 3 \sum_{i=1}^{n-1} \frac{\|\mathbf{M}^{(i)}\|_{\infty}}{\|\mathbf{A}\|_{\infty}} \leq 1 + 3 \sum_{i=1}^{n-1} 2^i \\ &= 1 + 3(2^n - 2) = 3 \cdot 2^n - 5 \leq 3 \cdot 2^n. \end{aligned}$$

Analog erhält man für die Zeilenpivotisierung bezüglich der 1-Norm (Spaltensummennorm) die Abschätzung

$$\|\delta\mathbf{A}\|_1 \leq \text{eps} F(\mathbf{A}) \|\mathbf{A}\|_1$$

mit $F(\mathbf{A}) \leq 3 \cdot 2^n$.

Diese Fehlerschranken sind im allgemeinen zu pessimistisch. Für die meisten praktisch auftretenden Gleichungssysteme gilt im Falle der Spalten- oder Zeilenpivotisierung $F(\mathbf{A}) \leq 10n$. Damit reicht die Spaltenpivotisierung für die meisten Fälle zur Stabilisierung des GAUSSschen Algorithmus aus.

Im Falle der der Totalpivotisierung erhält man

$$F(\mathbf{A}) \leq 1.8 \cdot n^{2+\ln/4}.$$

Für praktisch auftretende Systeme liegt hier $F(\mathbf{A})$ in der Größenordnung von 1. Da aber der Aufwand an Vergleichsoperationen bei der Totalpivotisierung den entsprechenden Aufwand bei der Spalten- oder Zeilenpivotisierung bei weitem übersteigt ($n^3/3 \longleftrightarrow n^2/2$), verwendet man normalerweise nur die Spaltenpivotisierung.

8.28. Beispiel: Wir betrachten das Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit

$$\mathbf{A} = \begin{pmatrix} 0.005 & 1.000 \\ 1.000 & 1.000 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 0.5 \\ 1.0 \end{pmatrix}.$$

Die exakte Lösung lautet

$$\mathbf{x} = \begin{pmatrix} \frac{500}{995} \\ \frac{495}{995} \end{pmatrix} = \begin{pmatrix} 0.503\dots \\ 0.497\dots \end{pmatrix}.$$

Auf einem Rechner mit dem Maschinenzahlbereich $\mathbb{M}(10, 2, \dots)$ erhält man ohne Pivotisierung

$$\mathbf{A} \implies \begin{pmatrix} 1 & 0 \\ 200 & 1 \end{pmatrix} \begin{pmatrix} 0.005 & 1 \\ 0 & -200 \end{pmatrix} = \mathbf{LU}.$$

Vorwärtselemination und Rücksubstitution liefern

$$\mathbf{L}\mathbf{y} = \mathbf{b} \implies \mathbf{y} = \begin{pmatrix} 0.5 \\ -99 \end{pmatrix}$$

und

$$\mathbf{U}\hat{\mathbf{x}} = \mathbf{y} \implies \hat{\mathbf{x}} = \begin{pmatrix} 0 \\ 0.5 \end{pmatrix}.$$

Mit Spaltenpivotisierung würde das Element a_{21} als Pivotelement gewählt werden. Mit der Permutationsmatrix

$$\mathbf{P} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

ergibt sich

$$\mathbf{PA} \implies \begin{pmatrix} 1 & 0 \\ 0.005 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} = \bar{\mathbf{L}}\bar{\mathbf{U}}.$$

Als Lösung erhalten wir:

$$\bar{\mathbf{L}}\bar{\mathbf{y}} = \mathbf{P}^T\mathbf{b} \implies \bar{\mathbf{y}} = \begin{pmatrix} 1 \\ 0.5 \end{pmatrix}$$

und

$$\bar{\mathbf{U}}\bar{\mathbf{x}} = \bar{\mathbf{y}} \implies \bar{\mathbf{x}} = \begin{pmatrix} 0.5 \\ 0.5 \end{pmatrix}.$$

Die zweite Lösung ist im Rahmen der Maschinengenauigkeit exakt. Den beiden Zerlegungsvarianten entsprechen folgende Störungen in der Matrix:

$$\delta\mathbf{A} = \mathbf{A} - \mathbf{LU} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}, \quad \|\delta\mathbf{A}\| = 1$$

und

$$\delta\bar{\mathbf{A}} = \mathbf{A} - \mathbf{P}^T\bar{\mathbf{L}}\bar{\mathbf{U}} = \begin{pmatrix} 0 & 0.005 \\ 0 & 0 \end{pmatrix}, \quad \|\delta\bar{\mathbf{A}}\| = 0.005.$$

Multipliziert man im ursprünglichen Gleichungssystem die zweite Gleichung mit 200, so erhält man das äquivalente System

$$\tilde{\mathbf{A}}\mathbf{x} = \tilde{\mathbf{b}}$$

mit

$$\tilde{\mathbf{A}} = \begin{pmatrix} 1 & 200 \\ 1 & 1 \end{pmatrix}, \quad \tilde{\mathbf{b}} = \begin{pmatrix} 100 \\ 1 \end{pmatrix}.$$

Wendet man auf die Matrix $\tilde{\mathbf{A}}$ den Algorithmus zur LU -Zerlegung mit Spaltenpivotisierung an, so würde kein Spaltentausch erfolgen. Das Element $\tilde{a}_{11} = 1$ wird als Pivotelement akzeptiert. Im Maschinenzahlbereich $\mathbb{M}(10, 2, \dots)$ erhält man die Dreiecksfaktoren

$$\tilde{\mathcal{L}} = \begin{pmatrix} 1 & 0 \\ 1 & 1 \end{pmatrix}, \quad \tilde{\mathcal{U}} = \begin{pmatrix} 1 & 200 \\ 0 & -200 \end{pmatrix}.$$

Für die Dreieckssysteme werden dann folgende Lösungen berechnet

$$\tilde{\mathcal{L}}\tilde{\mathbf{y}} = \tilde{\mathbf{b}} \implies \tilde{\mathbf{y}} = \begin{pmatrix} 100 \\ -99 \end{pmatrix}$$

und

$$\tilde{\mathcal{U}}\tilde{\mathbf{x}} = \tilde{\mathbf{y}} \implies \tilde{\mathbf{x}} = \begin{pmatrix} 0 \\ 0.5 \end{pmatrix}.$$

Das ist dieselbe schlechte Lösung wie beim ursprünglichen System ohne Pivotisierung. Als entsprechende Störung der Matrix ergibt sich wieder

$$\delta\tilde{\mathbf{A}} = \tilde{\mathbf{A}} - \tilde{\mathcal{L}}\tilde{\mathcal{U}} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}.$$



Wie das Beispiel zeigte, lässt sich jedes Nichtnullelement der aktuellen Spalte durch Skalierung (Multiplikation der entsprechenden Zeile mit einer hinreichend großen Konstanten) zum betragsgrößten Element machen. Damit ist die einfache Pivotisierung ohne Berücksichtigung der Gesamtmatrix fraglich. Die Spaltenpivotisierung sichert, dass das Verhältnis der Normen $\|\mathbf{M}^{(k)}\|_\infty$ zu $\|\mathbf{M}^{(k-1)}\|_\infty$ nicht größer als 2 wird. Damit gilt die Abschätzung

$$\|\mathbf{M}^{(k)}\|_\infty \leq 2^k \|\mathbf{A}\|_\infty.$$

Die Rückwärtsanalyse der LU -Zerlegung lieferte die Aussage, dass die berechneten Faktoren \mathcal{L} und \mathcal{U} exakte Dreiecksfaktoren der gestörten Matrix $\mathbf{A} + \delta\mathbf{A}$ sind. Die Störung genüge dabei der Abschätzung

$$\|\delta\mathbf{A}\|_\infty \leq \text{eps} F(\mathbf{A}) \|\mathbf{A}\|_\infty$$

mit

$$F(\mathbf{A}) = 1 + 3 \sum_{k=1}^{n-1} \frac{\|\mathbf{M}^{(k)}\|_\infty}{\|\mathbf{A}\|_\infty}.$$

Für den gesamten erzeugten Rundungsfehler beim Lösen eines Gleichungssystems mittels LU -Zerlegung erhielten wir

$$\|\delta\mathbf{x}\|_\infty \leq \frac{\text{eps} \bar{F}(\mathbf{A}) \text{cond}_\infty(\mathbf{A})}{1 - \text{eps} \bar{F}(\mathbf{A}) \text{cond}_\infty(\mathbf{A})} \|\mathbf{x}\|_\infty$$

mit $\bar{F}(\mathbf{A}) = (n+2)F(\mathbf{A})$. Die Fehlerkonstante $F(\mathbf{A})$ lässt sich mittels Spaltenpivotisierung klein halten. Durch Skalierung, also Übergang zu einer Matrix $\tilde{\mathbf{A}} = \mathbf{D}\mathbf{A}$, wobei \mathbf{D} eine Diagonalmatrix ist, kann aber $\|\tilde{\mathbf{A}}\|_\infty$ so vergrößert werden, dass der Effekt der Pivotisierung zunichte gemacht wird. Wir wollen darum versuchen, durch geeignete Skalierung die Norm der Matrix zu verkleinern. Bezüglich der ∞ -Norm lässt sich eine optimale Skalierung angeben.

8.29. Satz: Für die reguläre (n, n) -Matrix \mathbf{A} sei die Diagonalmatrix \mathbf{D} wie folgt definiert.

$$\mathbf{D} = \text{diag}(d_1, \dots, d_n)$$

mit

$$d_i = \frac{\|\mathbf{A}\|_\infty}{\sum_{j=1}^n |a_{ij}|}, \quad i = 1, \dots, n.$$

Dann gilt für die Matrix $\tilde{\mathbf{A}} = \mathbf{D}\mathbf{A}$

1.

$$\|\tilde{\mathbf{A}}\|_\infty = \|\mathbf{A}\|_\infty,$$

2.

$$\frac{\|\mathbf{A}^{-1}\|_\infty}{\|\mathbf{D}\|_\infty} \leq \|\tilde{\mathbf{A}}^{-1}\|_\infty \leq \|\mathbf{A}^{-1}\|_\infty,$$

3.

$$\frac{\text{cond}_\infty(\mathbf{A})}{\|\mathbf{D}\|_\infty} \leq \text{cond}_\infty(\tilde{\mathbf{A}}) \leq \text{cond}_\infty(\mathbf{A}).$$

Beweis: Siehe Übungsaufgabe 19. *

Für die Skalierung aus Satz 8.29 gilt

$$\sum_{j=1}^n |\tilde{a}_{ij}| = \|\mathbf{A}\|_\infty, \quad i = 1, \dots, n.$$

Die Betragssummennorm aller Zeilen von $\tilde{\mathbf{A}}$ sind also gleich. Eine Matrix mit dieser Eigenschaft bezeichnen wir als **zeilenäquilibriert**. Der nächste Satz zeigt, dass gerade die zeilenäquilibrierten Matrizen die bezüglich der ∞ -Norm optimal skalierten sind.

8.30. Satz: *Unter allen durch Zeilenskalierung aus einer (n, n) -Matrix \mathbf{A} hervorgegangenen Matrizen $\tilde{\mathbf{A}} = \mathbf{D}\mathbf{A}$ hat jede zeilenäquilibrierte die kleinste Kondition bezüglich der ∞ -Norm.*

Beweis: Die Skalierungsmatrix \mathbf{D} aus Satz 8.29 erzeugt gerade aus der Matrix \mathbf{A} eine zeilenäquilibrierte Matrix $\tilde{\mathbf{A}}$. Es sei nun $\bar{\mathbf{D}}$ eine weitere Diagonalmatrix, für die ebenfalls $\|\bar{\mathbf{A}}\|_\infty = \|\bar{\mathbf{D}}\mathbf{A}\|_\infty = \|\mathbf{A}\|_\infty$ gilt. (Wegen $\text{cond}(\alpha\mathbf{A}) = \|\alpha\mathbf{A}\| \|(\alpha\mathbf{A})^{-1}\| = |\alpha| \|\mathbf{A}\| \frac{1}{|\alpha|} \|\mathbf{A}^{-1}\| = \text{cond}(\mathbf{A})$ genügt es, nur Skalierungen mit $\|\bar{\mathbf{D}}\mathbf{A}\| = \|\mathbf{A}\|$ zu betrachten.) Es sei weiterhin $\bar{\mathbf{D}} \geq \mathbf{O}$. Dann gilt

$$\|\bar{\mathbf{D}}\mathbf{A}\|_\infty = \max_{i=1, \dots, n} \left\{ \bar{d}_i \sum_{j=1}^n |a_{ij}| \right\} = \|\mathbf{A}\|_\infty.$$

Damit folgt

$$\bar{d}_i \leq \frac{\|\mathbf{A}\|_\infty}{\sum_{j=1}^n |a_{ij}|} = d_i, \quad i = 1, \dots, n,$$

und es gilt

$$\begin{aligned} \|\tilde{\mathbf{A}}^{-1}\|_\infty &= \|(\mathbf{D}\mathbf{A})^{-1}\|_\infty = \|\mathbf{A}^{-1}\mathbf{D}^{-1}\|_\infty = \|\mathbf{A}^{-1}\bar{\mathbf{D}}^{-1}\bar{\mathbf{D}}\mathbf{D}^{-1}\|_\infty \\ &\leq \|(\bar{\mathbf{D}}\mathbf{A})^{-1}\|_\infty \|\bar{\mathbf{D}}\mathbf{D}^{-1}\|_\infty = \|\bar{\mathbf{A}}^{-1}\|_\infty \max_{i=1, \dots, n} \left\{ \frac{\bar{d}_i}{d_i} \right\} \leq \|\bar{\mathbf{A}}^{-1}\|_\infty. \end{aligned}$$

Wegen $\|\tilde{\mathbf{A}}\|_\infty = \|\bar{\mathbf{A}}\|_\infty$ folgt daraus sofort $\text{cond}(\tilde{\mathbf{A}}) \leq \text{cond}(\bar{\mathbf{A}})$. *

Eine Zeilenäquilibrierung mit Spaltenpivotisierung führt daher zu einer Verkleinerung der Fehlerschranke

$$\frac{\text{eps}\bar{F}(\mathbf{A})\text{cond}_{\infty}(\mathbf{A})}{1 - \text{eps}\bar{F}(\mathbf{A})\text{cond}_{\infty}(\mathbf{A})} \|\mathbf{x}\|_{\infty}$$

und damit zu einer Verkleinerung des Rundungsfehlers.

8.31. LU-Zerlegung mit Spaltenpivotisierung, fiktivem Zeilentausch und fiktiver Skalierung:

Es ist die reguläre Matrix \mathbf{A} in ein Produkt $\mathbf{P}\mathbf{A} = \mathbf{L}\mathbf{U}$ mit einer Permutationsmatrix \mathbf{P} , einer unteren Einsdreiecksmatrix \mathbf{L} und einer oberen Dreiecksmatrix \mathbf{U} zu zerlegen. Die Permutationsmatrix \mathbf{P} sei dabei auf einem Feld p der Länge n gespeichert.

{Initialisierung}

Wähle eine Genauigkeitsschranke $\varepsilon > 0$.

for $i = 1$ **to** n **do**

$p(i) = i; d_i = 0$

for $j = 1$ **to** n **do**

$d_i = d_i + |a_{ij}|$

endfor

endfor

{LU-Zerlegung}

for $j = 1$ **to** $n - 1$ **do**

{Pivotsuche}

$piv = 0$

for $k = j$ **to** n **do**

if $d_j|a_{kj}| > piv$ **then**

$piv = d_j|a_{kj}|$

$i^* = k$

endif

endfor

if $piv \leq \varepsilon$ **then**

STOPP

endif

{Fiktiver Zeilentausch}

if $i^* \neq j$ **then**

$pp = p(i^*); p(i^*) = p(j); p(j) = pp$

endif

{Transformation der Restmatrix}

for $i = j + 1$ **to** n **do**

$a_{p(i),j} = a_{p(i),j}/a_{p(j),j}$

```

for  $k = j + 1$  to  $n$  do
     $a_{p(i),k} = a_{p(i),k} - a_{p(i),j} \cdot a_{p(j),k}$ 
endfor
endfor
endfor

```

Aufwand: $\sim n^3/3$ Additionen/Multiplikationen

8.32. Beispiel: Wir betrachten wieder das Gleichungssystem $\mathbf{Ax} = \mathbf{b}$ mit

$$\mathbf{A} = \begin{pmatrix} 0.005 & 1.000 \\ 1.000 & 1.000 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 0.5 \\ 1.0 \end{pmatrix}.$$

Für die Matrix \mathbf{A} ergibt sich die Skalierungsmatrix

$$\mathbf{D} = \begin{pmatrix} \frac{2}{1.005} & 0 \\ 0 & 1 \end{pmatrix}$$

und damit das zeilenäquilibrierte System

$$\tilde{\mathbf{A}} = \mathbf{DAx} = \begin{pmatrix} \frac{0.01}{1.005} & \frac{2}{1.005} \\ 1 & 1 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \frac{1}{1.005} \\ 1 \end{pmatrix} = \mathbf{Db} = \tilde{\mathbf{b}}.$$

Spaltenpivotisierung führt für dieses System zur Wahl von $\tilde{a}_{21} = 1$ als Pivotelement genau wie beim Beispiel 8.28 mit Spaltenpivotisierung. ♡

Führt man die Skalierung explizit durch (n^2 Multiplikationen), so treten zusätzliche Rundungsfehler auf. Eigentlich benötigt man die Skalierungsfaktoren aber nur zur Festlegung des jeweiligen Pivotelements. Wir suchen dann bei der Spaltenpivotisierung nicht nach dem betragsgrößten Element der Restspalte j , sondern nach dem Element, für das $d_i|a_{ij}|$ am größten ist. Diese Vorgehensweise nennt man **fiktive Skalierung**. Im obigen Algorithmus ist dies dargestellt.

Meist sind die Eingabedaten (Matrix, rechte Seite) fehlerbehaftet. Der Einfluss dieser Datenfehler übersteigt dabei oft den erzeugten Rundungsfehler. In diesem Falle ist es nicht so wichtig, den Rundungsfehler zu minimieren. Vielmehr ist der Einfluss des Datenfehler genau abzuschätzen. Es sei also die Matrix \mathbf{A} mit dem Datenfehler $\delta_0\mathbf{A}$ und die rechte Seite \mathbf{b} mit dem Fehlervektor $\delta_0\mathbf{b}$ belastet. Dann gilt die Abschätzung

$$\|\delta\mathbf{x}\| \leq \frac{\|\mathbf{A}^{-1}\|}{1 - \|\mathbf{A}^{-1}\| \|\delta_0\mathbf{A}\|} (\|\delta_0\mathbf{A}\| \|\mathbf{x}\| + \|\delta_0\mathbf{b}\|).$$

8.33. Beispiel: Es sei

$$\mathbf{A} = \begin{pmatrix} 1.00 & 0.50 & 0.00 \\ 0.97 & 0.48 & 0.00 \\ 0.75 & 0.00 & -0.75 \end{pmatrix}, \quad \delta_0 \mathbf{A} = \mathbf{O},$$

$$\mathbf{b} = \begin{pmatrix} 3.50 \\ 3.37 \\ -6.75 \end{pmatrix}, \quad \delta_0 \mathbf{b} = 10^{-3} \begin{pmatrix} 1 \\ 1 \\ 75 \end{pmatrix}.$$

Es folgt

$$\mathbf{A}^{-1} = \begin{pmatrix} -96 & 100 & 0 \\ 194 & -200 & 0 \\ -96 & 100 & -\frac{4}{3} \end{pmatrix}, \quad \mathbf{x} = \begin{pmatrix} 1 \\ 5 \\ 10 \end{pmatrix}, \quad \delta \mathbf{x} = 10^{-3} \begin{pmatrix} 4 \\ -6 \\ -96 \end{pmatrix}$$

mit $\|\delta \mathbf{x}\|_\infty = 0.096$. Die obige Abschätzung liefert aber

$$\|\delta \mathbf{x}\|_\infty \leq \|\mathbf{A}^{-1}\|_\infty \|\delta_0 \mathbf{b}\|_\infty = 394 \cdot 75 \cdot 10^{-3} = 29.55.$$

Diese Abschätzung ist zu pessimistisch. Multiplizieren wir nun die letzte Gleichung in dem System mit $\frac{1}{75}$, so erhalten wir

$$\bar{\mathbf{A}} = \begin{pmatrix} 1.00 & 0.50 & 0.00 \\ 0.97 & 0.48 & 0.00 \\ 0.01 & 0.00 & -0.01 \end{pmatrix}, \quad \delta_0 \bar{\mathbf{A}} = \mathbf{O},$$

$$\bar{\mathbf{b}} = \begin{pmatrix} 3.50 \\ 3.37 \\ -0.09 \end{pmatrix}, \quad \delta_0 \bar{\mathbf{b}} = 10^{-3} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix},$$

$$\bar{\mathbf{A}}^{-1} = \begin{pmatrix} -96 & 100 & 0 \\ 194 & -200 & 0 \\ -96 & 100 & -100 \end{pmatrix}, \quad \bar{\mathbf{x}} = \mathbf{x}, \quad \delta \bar{\mathbf{x}} = \delta \mathbf{x}.$$

Hier liefert die Fehlerschätzung

$$\|\delta \bar{\mathbf{x}}\|_\infty \leq \|\delta \bar{\mathbf{A}}^{-1}\|_\infty \|\delta_0 \bar{\mathbf{b}}\|_\infty = 394 \cdot 10^{-3} = 0.394.$$

Diese Abschätzung liegt in der Größenordnung des wahren Fehlers. ♡

Das Verhalten, das im letzten Beispiel geschildert wurde, ist auch im allgemeinen gültig. Durch Äquilibrierung des Fehlervektors der rechten Seite erreicht man realistischere Fehlerabschätzungen. Gleiches gilt bezüglich einer Zeilenäquilibrierung des Fehlers der Koeffizientenmatrix. Damit ergeben sich folgende Skalierungsregeln:

R1 Sind die Eingabedaten nicht fehlerbehaftet ($\delta_0 \mathbf{A} = \mathbf{O}$ und $\delta_0 \mathbf{b} = \mathbf{o}$) skaliere so, dass die Matrix zeilenäquilibriert wird:

$$d_i = \frac{\|\mathbf{A}\|_\infty}{\sum_{j=1}^n |a_{ij}|}, \quad i = 1, \dots, n.$$

R2 Ist die rechte Seite fehlerbehaftet mit $|\delta_0 \mathbf{b}| \leq \Delta \mathbf{b}$ und die Matrix exakt bzw. $\|\delta_0 \mathbf{A}\| \ll \|\Delta \mathbf{b}\|$ skaliere so, dass der Fehlervektor $\Delta \mathbf{b}$ äquilibriert wird:

$$d_i = \frac{\max_{j=1, \dots, n} \{\Delta b_j\}}{\Delta b_i} = \frac{\|\Delta \mathbf{b}\|_\infty}{\Delta b_i}, \quad i = 1, \dots, n.$$

R3 Ist die Matrix fehlerbehaftet mit $|\delta_0 \mathbf{A}| \leq \Delta \mathbf{A}$ und die rechte Seite exakt bzw. $\|\delta_0 \mathbf{b}\| \ll \|\Delta \mathbf{A}\|$ skaliere so, dass die Fehlermatrix $\Delta \mathbf{A}$ zeilenäquilibriert wird:

$$d_i = \frac{\|\delta_0 \mathbf{A}\|_\infty}{\sum_{j=1}^n |\Delta a_{ij}|}, \quad i = 1, \dots, n.$$

R4 Sind sowohl die rechte Seite als auch die Matrix fehlerbehaftet mit $|\delta_0 \mathbf{b}| \leq \Delta \mathbf{b}$ und $|\delta_0 \mathbf{A}| \leq \Delta \mathbf{A}$ skaliere so, dass der Fehlervektor

$$\Delta \mathbf{z} = \Delta \mathbf{A} \mathbf{x} + \Delta \mathbf{b}$$

äquilibriert wird:

$$d_i = \frac{\max_{j=1, \dots, n} \{\Delta z_j\}}{\Delta z_i} = \frac{\|\Delta \mathbf{z}\|_\infty}{\Delta z_i}, \quad i = 1, \dots, n.$$

Bemerkung: Zum Berechnen der Skalierungsfaktoren in Regel **R4** ist die Kenntnis der Lösung \mathbf{x} notwendig. Darum wenden wir einen zweistufigen Prozess an.

S1 Skaliere mit **R1** und berechne Lösung \mathbf{x} .

S2 Berechne $\Delta \mathbf{z}$ und \mathbf{D} nach **R4** und schätze $\|\delta \mathbf{x}\|_\infty$ gemäß

$$\|\delta \mathbf{x}\|_\infty \leq \frac{\|\mathbf{A}^{-1} \mathbf{D}^{-1}\|_\infty}{1 - \|\mathbf{A}^{-1} \mathbf{D}^{-1}\|_\infty \|\mathbf{D} \Delta \mathbf{A}\|_\infty} \|\Delta \mathbf{z}\|_\infty.$$

8.2.4. Symmetrische Matrizen

Die gesamte Information einer symmetrischen Matrix ist auf dem oberen (unteren) Dreieck der Matrix einschließlich der Diagonalen enthalten. Man benötigt so gegenüber einer unsymmetrischen Matrix nur etwa den halben Speicherplatz. Beim Lösen von Gleichungssystemen mit symmetrischen Koeffizientenmatrizen lässt sich der Rechenaufwand ebenfalls halbieren. Eine hinreichende Bedingung liefert der folgende Satz.

8.34. Satz: *Der GAUSSsche Algorithmus sei für die symmetrische (n, n) -Matrix A ohne Pivotisierung durchführbar. Dann gilt*

$$A = LU = LDL^T$$

mit $D = \text{diag}(u_{11}, \dots, u_{nn})$.

Beweis: Es ist zu zeigen, dass $U = DL^T$ gilt. Es sei $R = D^{-1}U$. Dann ist R obere Dreiecksmatrix, und es gilt $A = LDR$. Aus der Symmetrie von A folgt

$$LDR = (LDR)^T = R^T DL^T.$$

Daraus erhält man

$$\left(R^T\right)^{-1} LD = DL^T R^{-1}.$$

Auf der linken Seite der Gleichung steht eine untere Dreiecksmatrix mit den Diagonalelementen d_1, \dots, d_n . Auf der rechten Seite steht eine obere Dreiecksmatrix mit den Diagonalelementen d_1, \dots, d_n . Das ist aber nur möglich, falls $DL^T R^{-1} = D$ ist. Damit folgt dann $L^T = R$. *

Erhält man eine Zerlegung der Form $A = LDL^T$, so war offensichtlich die LU -Zerlegung ohne Pivotisierung durchführbar. Damit ist die Bedingung auch notwendig. Eine Pivotisierung würde im allgemeinen die Symmetrie der Matrix zerstören. Denkbar ist nur eine Diagonalpivotisierung, also ein simultaner Zeilen- und Spaltentausch. Dabei werden nur Diagonalelemente als Pivotelemente ausgewählt.

Wir wollen nun zeigen, dass die LDL^T -Zerlegung einer symmetrischen Matrix etwa nur halb soviel Rechenoperationen erfordert wie die LU -Zerlegung. In Übungsaufgabe 15 wird gezeigt, dass alle Restmatrizen $M^{(k)}$ bei der LU -Zerlegung einer symmetrischen Matrix ebenfalls symmetrisch sind, falls ohne Pivotisierung gearbeitet wird. Wegen $U = DL^T$ dürfen die Elemente l_{ij} auf dem oberen Dreieck von

A gespeichert werden. Mit den Elementen von D werden die entsprechenden Diagonalelemente von A überschrieben. Somit ergibt sich nach k Schritten folgende Situation:

$$A^{(k)} = \begin{pmatrix} d_1 & l_{21} & \cdots & l_{k1} & l_{k+1,1} & \cdots & l_{n1} \\ & d_2 & \ddots & l_{k2} & l_{k+1,2} & & l_{n2} \\ & & \ddots & \ddots & \vdots & \ddots & \vdots \\ & & & d_k & l_{k+1,k} & \cdots & l_{nk} \\ & & & & a_{k+1,k+1}^{(k)} & \cdots & a_{k+1,n}^{(k)} \\ & & & & & \ddots & \vdots \\ & & & & & & a_{nn}^{(k)} \end{pmatrix}.$$

(Dabei ist in der Abbildung nur das obere Dreieck der aktuellen Iterationsmatrix angegeben.) Im nächsten Transformationsschritt wird folgendes berechnet:

$$\begin{aligned} d_{k+1} &= a_{k+1,k+1}^{(k)} \quad (\text{falls } a_{k+1,k+1}^{(k)} \neq 0), \\ l_{i,k+1} &= \frac{a_{k+1,i}^{(k)}}{d_{k+1}}, \quad i = k+2, \dots, n, \\ a_{ij}^{(k+1)} &= a_{ij}^{(k)} - l_{i,k+1} a_{k+1,j}^{(k)}, \quad i = k+2, \dots, n, \quad j = i, \dots, n. \end{aligned}$$

Wir erhalten den folgenden Algorithmus.

8.35. LDL^T -Zerlegung:

Es ist die reguläre symmetrische Matrix A in ein Produkt $A = LDL^T$ mit einer unteren Einsdreiecksmatrix L und einer Diagonalmatrix D zu zerlegen. Von der Matrix A ist nur das obere Dreieck einschließlich der Diagonalen gespeichert.

{Initialisierung}

Wähle Genauigkeitsschranke $\varepsilon > 0$.

{ LDL^T -Zerlegung}

for $k = 1$ **to** $n - 1$ **do**

if $|a_{kk}| \leq \varepsilon$ **then**

 STOPP

endif

 {Transformation der Restmatrix}

for $i = k + 1$ **to** n **do**

$l = a_{ki} / a_{kk}$

for $j = i$ **to** n **do**

$a_{ij} = a_{ij} - l \cdot a_{kj}$

endfor

```

    aki = l
  endfor
endfor

```

Aufwand: $\sim n^3/6$ Additionen/Multiplikationen

Ist von einer symmetrischen Matrix eine LDL^T -Zerlegung bekannt, so löst man ein lineares Gleichungssystem in drei Schritten:

$$L\mathbf{y} = \mathbf{b}, \quad D\mathbf{z} = \mathbf{y}, \quad L^T\mathbf{x} = \mathbf{z}.$$

Man erhält folgenden Algorithmus.

8.36. Lösen eines linearen Gleichungssystems bei bekannter LDL^T -Zerlegung:

Es ist das lineare Gleichungssystem $A\mathbf{x} = \mathbf{b}$ mit der regulären symmetrischen Matrix A zu lösen. Von der Matrix A sei eine Zerlegung der Form $A = LDL^T$ mit einer unteren Einsdreiecksmatrix L und einer Diagonalmatrix D bekannt. Dabei ist die Matrix L^T auf dem oberen Dreieck der Matrix A und die Diagonalmatrix D auf der Diagonale der Matrix A gespeichert.

```

{L $\mathbf{y} = \mathbf{b}$ }
{ $\mathbf{b}$  wird mit  $\mathbf{y}$  überschrieben.}
for k = 1 to n - 1 do
  for i = k + 1 to n do
    bi = bi - aki · bk
  endfor
endfor
{D $\mathbf{z} = \mathbf{y}$ }
{ $\mathbf{y}$  steht auf dem Speicherplatz von  $\mathbf{b}$  und wird mit  $\mathbf{z}$  überschrieben.}
for i = 1 to n do
  bi = bi/aii
endfor
{LT $\mathbf{x} = \mathbf{z}$ }
{ $\mathbf{z}$  steht auf dem Speicherplatz von  $\mathbf{b}$  und wird mit  $\mathbf{x}$  überschrieben.}
for k = n to 2 step -1 do
  for i = 1 to k - 1 do
    bi = bi - aik · bk
  endfor
endfor

```

Aufwand: $\sim n^2$ Additionen/Multiplikationen

Es bleibt die Frage zu klären, für welche Matrizen eine LDL^T -Zerlegung möglich ist. In Satz 8.34 hatten wir gesehen, dass das genau dann der Fall ist, wenn der

GAUSSsche Algorithmus ohne Pivotisierung durchführbar ist. Eine Klasse von Matrizen, für die dies zutrifft sind die sogenannten diagonaldominanten Matrizen. Eine symmetrische (n, n) -Matrix \mathbf{A} heißt **streng diagonaldominant**, falls

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|$$

für $i = 1, \dots, n$ gilt. **Bemerkung:** Im unsymmetrischen Falle gilt analog: **Strenge Zeilendiagonal- dominanz** bzw. **Strenge Spaltendiagonaldominanz**, falls

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n$$

bzw.

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ji}|, \quad i = 1, \dots, n$$

gilt.

8.37. Satz: Für streng diagonaldominante symmetrische Matrizen \mathbf{A} ist die LDL^T -Zerlegung durchführbar. Die berechneten Faktoren \mathcal{L} und \mathcal{D} sind exakte Faktoren der Zerlegung einer gestörten Matrix

$$\mathbf{A} + \delta\mathbf{A} = \mathcal{L}\mathcal{D}\mathcal{L}^T.$$

Die Störung $\delta\mathbf{A}$ genügt der Abschätzung

$$\|\delta\mathbf{A}\|_\infty \leq \text{eps } F(\mathbf{A}) \|\mathbf{A}\|_\infty$$

mit $F(\mathbf{A}) = 1 + 3(n - 1) \approx 3n$.

Beweis: In den Übungsaufgaben 15 und 16 wird gezeigt, dass die Symmetrie und die Diagonaldominanz bei einem Transformationsschritt erhalten bleiben. Die Diagonaldominanz wird sogar verstärkt. Damit ist der gesamte Algorithmus durchführbar. Es gilt nach Übungsaufgabe 16

$$\|\mathbf{M}^{(n-1)}\|_\infty \leq \|\mathbf{M}^{(n-2)}\|_\infty \leq \dots \leq \|\mathbf{M}^{(1)}\|_\infty \leq \|\mathbf{M}^{(0)}\|_\infty = \|\mathbf{A}\|_\infty.$$

Mit Satz 8.27 folgt dann sofort die Existenz einer Störung $\delta\mathbf{A}$ mit $\mathbf{A} + \delta\mathbf{A} = \mathcal{L}\mathcal{D}\mathcal{L}^T$ und $\|\delta\mathbf{A}\|_\infty \leq \text{eps } (1 + 3(n - 1)) \|\mathbf{A}\|_\infty$. *

Bemerkung: Die Abschätzung

$$\|\delta A\|_\infty \leq \text{eps}(1 + 3(n-1))\|A\|_\infty$$

gilt auch für die LU -Zerlegung einer zeilendiagonaldominanten unsymmetrischen Matrix.

Eine weitere Klasse von Matrizen, für die der GAUSSsche Algorithmus ohne Pivottisierung durchführbar ist, ist die Klasse der positiv definiten Matrizen.

Eine symmetrische Matrix A heißt **positiv definit**, falls für alle Vektoren $x \in \mathbb{R}^n$ mit $x \neq \mathbf{o}$ $x^T A x > 0$ gilt. Gilt für alle Vektoren $x \in \mathbb{R}^n$ mit $x \neq \mathbf{o}$ $x^T A x < 0$, so heißt A **negativ definit**, in den übrigen Fällen **indefinit**. Wir zeigen zunächst die Durchführbarkeit der LDL^T -Zerlegung für positiv definite Matrizen.

8.38. Satz: Für jede symmetrische, positiv definite Matrix A ist die LDL^T -Zerlegung durchführbar. Man erhält eine Diagonalmatrix

$$D = \text{diag}(d_1, \dots, d_n), \quad d_i > 0, \quad i = 1, \dots, n.$$

Beweis: Wir beweisen den Satz mittels Induktion über die Dimension n der Matrix.

Für $n = 1$ ist die Behauptung offensichtlich richtig.

Wir nehmen an, dass jede symmetrische, positiv definite Matrix $A \in \mathbb{R}^{k \times k}$ in der Form $A = LDL^T$ mit $D \geq \mathbf{O}$ zerlegbar ist. Nun betrachten wir eine Matrix $\bar{A} \in \mathbb{R}^{(k+1) \times (k+1)}$. Diese Matrix sei folgendermaßen partitioniert:

$$\bar{A} = \begin{pmatrix} A & b \\ b^T & \alpha \end{pmatrix}, \quad A \in \mathbb{R}^{k \times k}.$$

Mit einem beliebigen Vektor $x \in \mathbb{R}^n$, $x \neq \mathbf{o}$ folgt

$$0 < (x^T, 0) \bar{A} \begin{pmatrix} x \\ 0 \end{pmatrix} = (x^T, 0) \begin{pmatrix} A & b \\ b^T & \alpha \end{pmatrix} \begin{pmatrix} x \\ 0 \end{pmatrix} = x^T A x.$$

Damit ist die Matrix A ebenfalls positiv definit, und es existiert wegen der Voraussetzung eine Zerlegung $A = LDL^T$. Wir erhalten

$$\bar{A} = \begin{pmatrix} LDL^T & b \\ b^T & \alpha \end{pmatrix}.$$

Nun machen wir für die Faktoren der Zerlegung von \bar{A} den Ansatz

$$\bar{L} = \begin{pmatrix} L & \\ c^T & 1 \end{pmatrix}, \quad \bar{D} = \begin{pmatrix} D & \\ & d \end{pmatrix}.$$

Es folgt

$$\begin{aligned}\bar{L}\bar{D}\bar{L}^T &= \begin{pmatrix} \mathbf{L} & \\ \mathbf{c}^T & 1 \end{pmatrix} \begin{pmatrix} \mathbf{D} & \\ & d \end{pmatrix} \begin{pmatrix} \mathbf{L}^T & \mathbf{c} \\ & 1 \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{LDL}^T & \mathbf{LDc} \\ \mathbf{c}^T \mathbf{DL}^T & \mathbf{c}^T \mathbf{Dc} + d \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{LDL}^T & \mathbf{b} \\ \mathbf{b}^T & \alpha \end{pmatrix}.\end{aligned}$$

Daraus erhalten wir

$$\mathbf{c} = \mathbf{D}^{-1} \mathbf{L}^{-1} \mathbf{b}$$

und

$$d = \alpha - \mathbf{b}^T \left(\mathbf{L}^{-1} \right)^T \mathbf{D}^{-1} \mathbf{L}^{-1} \mathbf{b} = \alpha - \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}.$$

Aus der Regularität von \mathbf{D} und \mathbf{L} folgt sofort die Existenz von \mathbf{c} . Es bleibt dann nur noch die Positivität von d zu zeigen. Dazu sei $\mathbf{y} \in \mathbb{R}^{k+1}$ die Lösung des Gleichungssystems $\bar{\mathbf{L}}^T \mathbf{y} = \mathbf{e}_{k+1}$. Dann folgt aus der positiven Definitheit von $\bar{\mathbf{A}}$

$$0 < \mathbf{y}^T \bar{\mathbf{A}} \mathbf{y} = \mathbf{y}^T \bar{\mathbf{L}} \bar{\mathbf{D}} \bar{\mathbf{L}}^T \mathbf{y} = \mathbf{e}_{k+1}^T \bar{\mathbf{D}} \mathbf{e}_{k+1} = d.$$

*

Ohne Beweis seien die Ergebnisse der Rundungsfehleranalyse der LDL^T -Zerlegung einer symmetrischen, positiv definiten Matrix angegeben.

8.39. Satz: Die LDL^T -Zerlegung ist für die symmetrische, positiv definite Matrix \mathbf{A} auf einem realen Rechner mit $d_i > 0$, $i = 1, \dots, n$, durchführbar, falls

$$\text{eps}_F \|\mathbf{A}^{-1}\|_2 \|\mathbf{A}\|_F \leq \frac{1}{2}$$

mit $F = n + 2 \ln(n/2) \approx n$ gilt.

Zu den berechneten Faktoren \mathcal{L} und \mathcal{D} existiert eine symmetrische Störung $\delta \mathbf{A}$ mit

$$\mathbf{A} + \delta \mathbf{A} = \mathcal{L} \mathcal{D} \mathcal{L}^T.$$

Die Störung genügt der Abschätzung

$$\|\delta \mathbf{A}\|_F \leq \text{eps}_F \|\mathbf{A}\|_F.$$

Die zu $\mathbf{b} \in \mathbb{R}^n$ berechnete Lösung $\bar{\mathbf{x}}$ ist exakte Lösung des gestörten Problems

$$(\mathbf{A} + \overline{\delta\mathbf{A}})\mathbf{x} = \mathbf{b}$$

mit der symmetrischen Störung $\overline{\delta\mathbf{A}}$, die der Abschätzung

$$\|\overline{\delta\mathbf{A}}\|_p \leq \text{eps}\bar{F}\|\mathbf{A}\|_p$$

mit

$$\bar{F} = \begin{cases} 2n^{3/2} + F \approx 2n^{3/2} & \text{für } p = F \\ n^{3/2} + n + n^{1/2}F \approx 2n^{3/2} & \text{für } p = 2 \end{cases}$$

genügt.

Für positiv definite Matrizen lässt sich auch eine Zerlegung der Form $\mathbf{A} = \hat{\mathbf{L}}\hat{\mathbf{L}}^T$ mit einer unteren Dreiecksmatrix $\hat{\mathbf{L}}$ angeben. Ein Vergleich mit der LDL^T -Zerlegung zeigt, dass dann $\hat{\mathbf{L}} = \mathbf{L}\mathbf{D}^{1/2}$ mit

$$\mathbf{D}^{1/2} = \text{diag}(\sqrt{d_1}, \dots, \sqrt{d_n})$$

gilt. Diese Zerlegung wird als **CHOLESKY-Zerlegung** bezeichnet.

8.40. CHOLESKY-Zerlegung:

Es ist die symmetrische, positiv definite Matrix \mathbf{A} in ein Produkt $\mathbf{A} = \hat{\mathbf{L}}\hat{\mathbf{L}}^T$ mit einer unteren Dreiecksmatrix $\hat{\mathbf{L}}$ zu zerlegen. Von der Matrix \mathbf{A} sind nur das untere Dreieck und die Diagonale gespeichert.

{Initialisierung}

Wähle Genauigkeitsschranke $\varepsilon > 0$.

{Zerlegung}

for $k = 1$ to n **do**

if $a_{kk} \leq \varepsilon$ **then**

 STOPP

endif

$$a_{kk} = \sqrt{a_{kk}}$$

 {Transformation der Restmatrix}

for $i = k + 1$ to n **do**

$$a_{ik} = a_{ik}/a_{kk}$$

endfor

for $i = k + 1$ to n **do**

for $j = i$ to n **do**

```

       $a_{ij} = a_{ij} - l \cdot a_{jk}$ 
    endfor
  endfor
endfor

```

Aufwand: $\sim n^3/6$ Additionen/Multiplikationen + n Wurzeln

Das Lösen eines Gleichungssystems $Ax = b$ erfolgt dann in zwei Schritten:

$$\hat{L}y = b, \quad \hat{L}^T x = y.$$

Für die Rundungsfehleranalyse der CHOLESKY-Zerlegung gelten alle Aussagen von Satz 8.39 mit $F = n + 1 + \frac{1}{2} \ln n$. Für großes n gilt wieder $F \approx n$ und $\bar{F} \approx 2n^{3/2}$.

Die LDL^T -Zerlegung ist nicht für alle regulären symmetrischen Matrizen durchführbar, auch nicht falls man Diagonalphivotisierung zulässt.

8.41. Beispiel: Für die Matrix

$$A = \begin{pmatrix} 0 & 1 & 2 \\ 1 & 0 & 1 \\ 2 & 1 & 0 \end{pmatrix}$$

existiert keine LDL^T -Zerlegung. ♡

Für derartige Matrizen existieren aber sogenannte Blockfaktorisierungen der Form

$$PAP^T = LDL^T$$

mit

$$L = \begin{pmatrix} I_1 & & & \\ G_{21} & I_2 & & O \\ \vdots & \vdots & \ddots & \\ G_{l1} & G_{l2} & \cdots & I_l \end{pmatrix}, \quad D = \begin{pmatrix} D_1 & & & \\ & D_2 & & \\ & & \ddots & \\ & & & D_l \end{pmatrix}.$$

Die Dimension der Diagonalblöcke I_i bzw. D_i beträgt jeweils 1 oder 2. Für diese Blockzerlegungen existieren auch Pivotstrategien, die die Stabilität der Algorithmen verbessern.

8.2.5. Orthogonalisierungsverfahren

Lösen wir ein lineares Gleichungssystem $Ax = b$ mittels einer LU -Zerlegung von A , so gilt

$$Ly = P^T b, \quad Ux = y.$$

Für den Einfluss von Fehlern in \mathbf{A} , \mathbf{b} , \mathbf{L} , \mathbf{U} und \mathbf{y} gelten die Abschätzungen

$$\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} \stackrel{\cdot}{\leq} \text{cond}(\mathbf{A}) \left[\frac{\|\delta \mathbf{A}\|}{\|\mathbf{A}\|} + \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|} \right]$$

bzw.

$$\frac{\|\delta \mathbf{y}\|}{\|\mathbf{y}\|} \stackrel{\cdot}{\leq} \text{cond}(\mathbf{L}) \left[\frac{\|\delta \mathbf{L}\|}{\|\mathbf{L}\|} + \frac{\|\delta \mathbf{b}\|}{\|\mathbf{b}\|} \right],$$

$$\frac{\|\delta \mathbf{x}\|}{\|\mathbf{x}\|} \stackrel{\cdot}{\leq} \text{cond}(\mathbf{U}) \left[\frac{\|\delta \mathbf{U}\|}{\|\mathbf{U}\|} + \frac{\|\delta \mathbf{y}\|}{\|\mathbf{y}\|} \right].$$

Die maximale Fehlerverstärkung beim Lösen der Dreieckssysteme wird daher durch die Größe $\text{cond}(\mathbf{L})\text{cond}(\mathbf{U})$ beschrieben. Wegen der Submultiplikativität der zugrunde liegenden Matrixnormen gilt aber

$$\text{cond}(\mathbf{A}) \leq \text{cond}(\mathbf{L})\text{cond}(\mathbf{U}).$$

Dabei wird das Produkt der Konditionszahlen von \mathbf{L} und \mathbf{U} die Konditionszahl von \mathbf{A} bei weitem übertreffen.

8.42. Beispiel: Für die Matrix

$$\mathbf{A} = \begin{pmatrix} 1 & 0 & 0 & 1 \\ -1 & 1 & 0 & 1 \\ -1 & -1 & 1 & 1 \\ -1 & -1 & -1 & 1 \end{pmatrix}$$

ergeben sich die Faktoren

$$\mathbf{L} = \begin{pmatrix} 1 & & & \\ -1 & 1 & & \\ -1 & -1 & 1 & \\ -1 & -1 & -1 & 1 \end{pmatrix}, \quad \mathbf{U} = \begin{pmatrix} 1 & 0 & 0 & 1 \\ & 1 & 0 & 2 \\ & & 1 & 4 \\ & & & 8 \end{pmatrix}.$$

Weiterhin ist

$$\mathbf{A}^{-1} = \frac{1}{8} \begin{pmatrix} 4 & -2 & -1 & -1 \\ 0 & 4 & -2 & -2 \\ 0 & 0 & 4 & -4 \\ 4 & 2 & 1 & 1 \end{pmatrix},$$

$$\mathbf{L}^{-1} = \begin{pmatrix} 1 & & & \\ 1 & 1 & & \\ 2 & 1 & 1 & \\ 4 & 2 & 1 & 1 \end{pmatrix}, \quad \mathbf{U}^{-1} = \frac{1}{8} \begin{pmatrix} 8 & 0 & 0 & -1 \\ & 8 & 0 & -2 \\ & & 8 & -4 \\ & & & 1 \end{pmatrix}.$$

Damit gilt

$$\begin{aligned}\operatorname{cond}_1(\mathbf{A}) &= \|\mathbf{A}\|_1 \|\mathbf{A}^{-1}\|_1 = 4 \cdot 1 = 4, \\ \operatorname{cond}_1(\mathbf{L}) &= \|\mathbf{L}\|_1 \|\mathbf{L}^{-1}\|_1 = 4 \cdot 8 = 32, \\ \operatorname{cond}_1(\mathbf{U}) &= \|\mathbf{U}\|_1 \|\mathbf{U}^{-1}\|_1 = 15 \cdot 1 = 15.\end{aligned}$$

Hier ist $4 = \operatorname{cond}(\mathbf{A}) \ll \operatorname{cond}_1(\mathbf{L})\operatorname{cond}_1(\mathbf{U}) = 480$. ♡

Eine Forderung an Zerlegungen einer Matrix \mathbf{A} der Form $\mathbf{A} = \mathbf{A}_1 \mathbf{A}_2 \cdots \mathbf{A}_k$ wäre damit die, dass das Produkt der Konditionszahlen der Faktoren $\mathbf{A}_1, \dots, \mathbf{A}_k$ nicht wesentlich größer sein sollte als die Konditionszahl von \mathbf{A} . Noch besser wäre es, falls

$$\operatorname{cond}(\mathbf{A}) = \operatorname{cond}(\mathbf{A}_1) \cdots \operatorname{cond}(\mathbf{A}_k)$$

gelten würde.

Bezüglich der Spektralnorm lassen sich solche Zerlegungen sofort angeben. Da die Spektralnorm invariant gegenüber orthogonalen Transformationen ist, gilt für eine Zerlegung der Form $\mathbf{A} = \mathbf{Q}\mathbf{R}$ mit einer orthogonalen Matrix \mathbf{Q} :

$$\|\mathbf{A}\|_2 = \|\mathbf{Q}\mathbf{R}\|_2 = \|\mathbf{R}\|_2, \quad \|\mathbf{A}^{-1}\|_2 = \|\mathbf{R}^{-1}\mathbf{Q}^T\|_2 = \|\mathbf{R}^{-1}\|_2.$$

Damit folgt $\operatorname{cond}_2(\mathbf{A}) = \operatorname{cond}_2(\mathbf{R})$. Für orthogonale Matrizen gilt $\operatorname{cond}_2(\mathbf{Q}) = 1$, so dass sich insgesamt $\operatorname{cond}_2(\mathbf{A}) = \operatorname{cond}_2(\mathbf{Q})\operatorname{cond}_2(\mathbf{R})$ ergibt. Ist nun noch \mathbf{R} eine obere Dreiecksmatrix, so verwendet man die Zerlegung $\mathbf{A} = \mathbf{Q}\mathbf{R}$ zum Lösen von linearen Gleichungssystemen $\mathbf{A}\mathbf{x} = \mathbf{b}$. Man erhält \mathbf{x} als Lösung von

$$\mathbf{R}\mathbf{x} = \mathbf{Q}^T \mathbf{b}.$$

Die Frage nach der Existenz einer solchen Zerlegung lässt sich sofort beantworten. Es gilt

8.43. Satz: *Für jede reguläre Matrix \mathbf{A} existieren eine orthogonale Matrix \mathbf{Q} und eine obere Dreiecksmatrix \mathbf{R} , so dass $\mathbf{A} = \mathbf{Q}\mathbf{R}$ gilt.*

Beweis: Man wende das SCHMIDTsche Orthogonalisierungsverfahren auf die Spalten von \mathbf{A} an. ✱

Mit dem SCHMIDTsche Orthogonalisierungsverfahren hätte man sofort einen Algorithmus zum Berechnen einer \mathbf{QR} -Zerlegung. Leider ist das SCHMIDTsche Orthogonalisierungsverfahren aber nicht numerisch gutartig, wie das folgende Beispiel zeigt.

8.44. Beispiel: Wir betrachten die Matrix

$$A = \begin{pmatrix} 1 & 1 & 1 & 1 \\ \varepsilon & 0 & 0 & 0 \\ 0 & \varepsilon & 0 & 0 \\ 0 & 0 & \varepsilon & 0 \end{pmatrix}.$$

Für ε gelte $\varepsilon > \text{eps}$, aber $\varepsilon^2 < \text{eps}$, so dass $\text{gl}(1 + \varepsilon^2) = 1$. Das SCHMIDT'sche Orthogonalisierungsverfahren liefert dann

$$Q = \begin{pmatrix} 1 & 0 & 0 & 0 \\ \varepsilon & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & -1 \\ 0 & \frac{1}{\sqrt{2}} & 0 & 0 \\ 0 & 0 & \frac{1}{\sqrt{2}} & 0 \end{pmatrix}, \quad R = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & \varepsilon\sqrt{2} & 0 & 0 \\ 0 & 0 & \varepsilon\sqrt{2} & 0 \\ 0 & 0 & 0 & \varepsilon \end{pmatrix}.$$

Damit ist zwar eine exakte Zerlegung gegeben: $A = QR$, aber die Matrix Q ist extrem nichtorthogonal:

$$Q^T Q = \begin{pmatrix} 1 + \varepsilon^2 & -\frac{\varepsilon}{\sqrt{2}} & -\frac{\varepsilon}{\sqrt{2}} & -\varepsilon \\ -\frac{\varepsilon}{\sqrt{2}} & 1 & \frac{1}{2} & \frac{1}{\sqrt{2}} \\ -\frac{\varepsilon}{\sqrt{2}} & \frac{1}{2} & 1 & \frac{1}{\sqrt{2}} \\ -\varepsilon & \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 1 \end{pmatrix}.$$

Würde man mit dieser Zerlegung ein Gleichungssystem $Ax = b$ gemäß $Rx = Q^T b$ lösen, so würde man einen großen Fehler erhalten. Zum Beispiel ergibt sich für die rechte Seite

$$b = \begin{pmatrix} 4 \\ \varepsilon \\ \varepsilon \\ \varepsilon \end{pmatrix} \quad \text{mit exaktem} \quad x = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \quad \text{die Lösung} \quad \bar{x} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ -1 \end{pmatrix}.$$



Das SCHMIDT'sche Orthogonalisierungsverfahren ist daher in dieser Form für Computerrechnung nicht geeignet. Es existiert aber eine Modifizierung des Verfahrens, die diese numerischen Schwierigkeiten überwindet. Wir wollen jedoch ein anderes Verfahren zum Berechnen einer QR -Zerlegung einer Matrix A behandeln.

Das Householder-Verfahren

Erinnern wir uns zunächst an den GAUSSschen Algorithmus. Dort wurde mit Hilfe von elementaren Transformationsmatrizen (LNT-Matrizen) die Matrix schrittweise auf obere Dreiecksform transformiert. Im k -ten Schritt wurden durch die Transformation mit $L_k(-l_k)$ in der k -ten Spalte unterhalb der Diagonalen Nullen erzeugt. Wir wollen nun orthogonale Transformationsmatrizen konstruieren, mit denen wir ähnliches erreichen. Wir suchen nach einer orthogonalen Matrix P , die einen Vektor $\mathbf{a} \in \mathbb{R}^k$ auf ein Vielfaches des ersten Einheitsvektors $\mathbf{e}_1 \in \mathbb{R}^k$ transformiert:

$$P\mathbf{a} = \varrho\mathbf{e}_1.$$

Wir werden sehen, dass dies mit HOUSEHOLDER-Spiegelungen möglich ist. Es sei daher P eine HOUSEHOLDER-Matrix

$$P = H = I - 2\mathbf{u}\mathbf{u}^T, \quad \mathbf{u}^T\mathbf{u} = 1.$$

Den Vektor \mathbf{u} bestimmen wir so, dass

$$H\mathbf{a} = (I - 2\mathbf{u}\mathbf{u}^T)\mathbf{a} = \varrho\mathbf{e}_1$$

erfüllt ist. Da H orthogonal ist, gilt

$$\|\mathbf{a}\|_2 = \|H\mathbf{a}\|_2 = \|\varrho\mathbf{e}_1\|_2 = |\varrho| \|\mathbf{e}_1\|_2 = |\varrho|.$$

Damit folgt $\varrho = \pm\|\mathbf{a}\|_2$ und

$$(I - 2\mathbf{u}\mathbf{u}^T)\mathbf{a} = \pm\|\mathbf{a}\|_2\mathbf{e}_1, \quad (8.10)$$

$$\mathbf{a} - 2\mathbf{u}(\mathbf{u}^T\mathbf{a}) = \pm\|\mathbf{a}\|_2\mathbf{e}_1, \quad (8.11)$$

$$\mathbf{u} = \frac{\mathbf{a} \mp \|\mathbf{a}\|_2\mathbf{e}_1}{2(\mathbf{u}^T\mathbf{a})}. \quad (8.12)$$

Aus Gleichung 8.11 ergibt sich

$$\begin{aligned} \mathbf{a}^T [\mathbf{a} - 2\mathbf{u}(\mathbf{u}^T\mathbf{a})] &= \pm\|\mathbf{a}\|_2\mathbf{a}^T\mathbf{e}_1, \\ \|\mathbf{a}\|_2^2 - 2(\mathbf{u}^T\mathbf{a})^2 &= \pm\|\mathbf{a}\|_2\mathbf{a}^T\mathbf{e}_1, \\ (\mathbf{u}^T\mathbf{a})^2 &= \frac{\mp\|\mathbf{a}\|_2\mathbf{a}^T\mathbf{e}_1 - \|\mathbf{a}\|_2^2}{2} = \frac{\varrho a_1 - \varrho^2}{2}, \\ \mathbf{u}^T\mathbf{a} &= \sqrt{\frac{\varrho a_1 - \varrho^2}{2}} \end{aligned}$$

wobei a_1 die erste Komponente des Vektors \mathbf{a} bezeichnet. Setzt man dies in die Gleichung 8.12 ein, so erhält man

$$\mathbf{u} = \frac{\mathbf{a} - \varrho \mathbf{e}_1}{\sqrt{2\varrho(\varrho - a_1)}}$$

und

$$\begin{aligned} \mathbf{H} &= \mathbf{I} - 2\mathbf{u}\mathbf{u}^T = \mathbf{I} - 2 \frac{(\mathbf{a} - \varrho \mathbf{e}_1)(\mathbf{a} - \varrho \mathbf{e}_1)^T}{2\varrho(\varrho - a_1)} \\ &= \mathbf{I} - \frac{(\mathbf{a} - \varrho \mathbf{e}_1)(\mathbf{a} - \varrho \mathbf{e}_1)^T}{\varrho(\varrho - a_1)} = \mathbf{I} - \frac{\mathbf{v}\mathbf{v}^T}{\gamma} \end{aligned}$$

mit

$$\mathbf{v} = \mathbf{a} - \varrho \mathbf{e}_1 = \begin{pmatrix} a_1 - \varrho \\ a_2 \\ \vdots \\ a_k \end{pmatrix}, \quad \gamma = \varrho(\varrho - a_1).$$

Damit haben wir die gesuchte Transformation gefunden. Für die Wahl von ϱ stehen uns zwei Möglichkeiten zur Verfügung:

$$\varrho = \|\mathbf{a}\|_2 \quad \text{oder} \quad \varrho = -\|\mathbf{a}\|_2.$$

Wir werden das Vorzeichen von ϱ so wählen, dass beim Berechnen von $\varrho - a_1$ keine Auslöschung auftritt: ϱ und a_1 müssen entgegengesetzte Vorzeichen haben. Wir erhalten den folgenden Algorithmus.

8.45. Festlegen der HOUSEHOLDER-Matrix zur Transformation eines Vektors auf ein Vielfaches des ersten Einheitsvektors:

Es sei ein Vektor $\mathbf{a} \in \mathbb{R}^k$ gegeben.

S0 Berechne $\varrho = \sqrt{\mathbf{a}^T \mathbf{a}} = \|\mathbf{a}\|_2$.

S1 if $\varrho = 0$ then $\mathbf{H} = \mathbf{I}$ STOPP.

S2 if $a_1 > 0$ then $\varrho = -\varrho$.

Berechne $\mathbf{v} = \mathbf{a} - \varrho \mathbf{e}_1$ und $\gamma = \varrho(\varrho - a_1)$.

Die Matrix $\mathbf{H} = \mathbf{I} - \mathbf{v}\mathbf{v}^T/\gamma$ ist dann orthogonal und es gilt

$$\mathbf{H}\mathbf{a} = \varrho \mathbf{e}_1.$$

Aufwand: $\sim n$ Additionen/Multiplikationen, eine Quadratwurzel.

Das Berechnen der Transformation eines Vektors \boldsymbol{x} erfolgt nach folgendem Algorithmus.

8.46. HOUSEHOLDER-Transformation:

Gegeben sei eine (n, n) -HOUSEHOLDER-Matrix

$$\boldsymbol{H} = \boldsymbol{I} - \boldsymbol{v}\boldsymbol{v}^T / \gamma.$$

Zu berechnen ist für einen Vektor $\boldsymbol{x} \in \mathbb{R}^n$ der Vektor

$$\boldsymbol{y} = \boldsymbol{H}\boldsymbol{x}.$$

S0 Berechne $\beta = \boldsymbol{v}^T \boldsymbol{x} / \gamma$.

S1 Berechne $\boldsymbol{y} = \boldsymbol{x} - \beta \boldsymbol{v}$

Aufwand: $\sim 2n$ Additionen/Multiplikationen.

Bemerkung: Die Matrix \boldsymbol{H} braucht explizit nicht berechnet zu werden. Die gesamten Informationen sind im Vektor \boldsymbol{v} und in γ enthalten. Das explizite Berechnen der HOUSEHOLDER-Matrix ist aufwendig ($\sim n^2/2$ Additionen/Multiplikationen). Das Berechnen des Vektors $\boldsymbol{y} = \boldsymbol{H}\boldsymbol{x}$ würde sogar $\sim n^2$ Operationen kosten, falls man ihn in der Form Matrix mal Vektor berechnet.

Mit Hilfe der letzten beiden Algorithmen entwickeln wir nun einen Algorithmus zum Transformieren einer Matrix auf obere Dreiecksform. Wir nehmen an, dass wir nach $k-1$ Schritten die Matrix auf die Form

$$\boldsymbol{A}^{(k-1)} = \begin{pmatrix} a_{11}^{(k-1)} & \cdots & a_{1,k-1}^{(k-1)} & a_{1k}^{(k-1)} & \cdots & a_{1n}^{(k-1)} \\ 0 & \ddots & \vdots & \vdots & \ddots & \vdots \\ \vdots & \ddots & a_{k-1,k-1}^{(k-1)} & a_{k-1,k}^{(k-1)} & \cdots & a_{k-1,n}^{(k-1)} \\ 0 & \cdots & 0 & a_{kk}^{(k-1)} & \cdots & a_{kn}^{(k-1)} \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & a_{nk}^{(k-1)} & \cdots & a_{nn}^{(k-1)} \end{pmatrix} \\ = \left(\begin{array}{c|c} \boldsymbol{R}^{(k-1)} & \\ \hline \boldsymbol{O} & \boldsymbol{M}^{(k-1)} \end{array} \right)$$

mit $\boldsymbol{R}^{(k-1)} \in \mathbb{R}^{(k-1) \times n}$, $\boldsymbol{M}^{(k-1)} \in \mathbb{R}^{(n-k+1) \times (n-k+1)}$ transformiert haben. Im nächsten Schritt wählen wir eine HOUSEHOLDER-Matrix $\boldsymbol{H}^{(k)}$ so, dass die Teilmatrix $\boldsymbol{R}^{(k-1)}$ unverändert bleibt und die erste Spalte von $\boldsymbol{M}^{(k-1)}$ auf ein Vielfaches des entsprechenden Einheitsvektors abgebildet wird. Damit hat $\boldsymbol{H}^{(k)}$ folgende Struktur

$$\boldsymbol{H}^{(k)} = \begin{pmatrix} \boldsymbol{I}^{(k-1)} & \boldsymbol{O} \\ \boldsymbol{O} & \tilde{\boldsymbol{H}}^{(k)} \end{pmatrix}$$

mit

$$\tilde{H}^{(k)} M^{(k-1)} = \begin{pmatrix} \star & \star & \cdots & \star \\ 0 & \star & \cdots & \star \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \star & \cdots & \star \end{pmatrix} = \begin{pmatrix} \star & \star & \cdots & \star \\ 0 & & & \\ \vdots & & & M^{(k)} \\ 0 & & & \end{pmatrix}.$$

Damit gilt

$$\begin{aligned} H^{(k)} A^{(k-1)} &= \begin{pmatrix} I^{(k-1)} & O \\ O & \tilde{H}^{(k)} \end{pmatrix} \begin{pmatrix} R^{(k-1)} \\ O \mid M^{(k-1)} \end{pmatrix} \\ &= \begin{pmatrix} R^{(k-1)} \\ O \mid \tilde{H}^{(k)} M^{(k-1)} \end{pmatrix} \\ &= \begin{pmatrix} R^{(k-1)} \\ O \mid \begin{matrix} \star & \star & \cdots & \star \\ 0 & & & \\ \vdots & & & M^{(k)} \\ 0 & & & \end{matrix} \end{pmatrix} = \begin{pmatrix} R^{(k)} \\ O \mid M^{(k)} \end{pmatrix} \end{aligned}$$

mit $R^{(k)} \in \mathbb{R}^{k \times n}$ und $M^{(k)} \in \mathbb{R}^{(n-k) \times (n-k)}$. Nach $n-1$ Schritten ist schließlich $M^{(n-1)} \in \mathbb{R}^{1 \times 1}$. Dann ist $\bar{A}^{(n-1)}$ obere Dreiecksmatrix, und es gilt

$$R = A^{(n-1)} = H^{(n-1)} H^{(n-2)} \dots H^{(2)} H^{(1)} A$$

bzw.

$$A = H^{(1)} H^{(2)} \dots H^{(n-2)} H^{(n-1)} R = QR$$

mit der orthogonalen Matrix

$$Q = H^{(1)} H^{(2)} \dots H^{(n-2)} H^{(n-1)}.$$

Die Matrix Q braucht nicht explizit berechnet zu werden. Es genügt, die Vektoren $v^{(k)}$ und die γ_k zu speichern. So sind die HOUSEHOLDER-Transformationen

$$H^{(k)} = I - v^{(k)} v^{(k)T} / \gamma_k$$

und damit die Matrix Q jederzeit kostengünstig rekonstruierbar. Die Vektoren $v^{(k)}$, $k = 1, \dots, n-1$, dürfen ähnlich wie bei der LU -Zerlegung auf dem freiwerdenden

Speicherplatz unterhalb der Diagonalen abgelegt werden. Man benötigt in der k -ten Spalte $n - k + 1$ Speicherplätze für $\mathbf{v}^{(k)}$. Das ist ein Speicherplatz mehr als unterhalb der Diagonalen frei wird. Darum benutzt man üblicherweise die Positionen k bis n in der k -ten Spalte zum Speichern von $\mathbf{v}^{(k)}$. Die $\varrho_k = r_{kk}$ und die γ_k werden in zwei zusätzlichen Feldern der Länge n gespeichert. Mit diesen Konventionen geben wir nun den Algorithmus an.

8.47. QR-Zerlegung einer Matrix nach HOUSEHOLDER:

Es ist die reguläre (n, n) -Matrix \mathbf{A} in ein Produkt $\mathbf{A} = \mathbf{Q}\mathbf{R}$ mit einer orthogonalen Matrix \mathbf{Q} und einer oberen Dreiecksmatrix \mathbf{R} zu zerlegen.

{Initialisierung}

Wähle eine Genauigkeitsschranke $\varepsilon > 0$.

{QR-Zerlegung}

for $k = 1$ to $n - 1$ **do**

{Berechnen der Transformationsmatrix $\mathbf{H}^{(k)}$ }

$$\varrho_k = \sqrt{a_{kk}^2 + a_{k+1,k}^2 + \cdots + a_{nk}^2}$$

if $\varrho_k \leq \varepsilon$ **then**

STOPP

endif

if $a_{kk} > 0$ **then**

$$\varrho_k = -\varrho_k$$

endif

$$a_{kk} = a_{kk} - \varrho_k$$

$$\gamma_k = -\varrho_k \cdot a_{kk}$$

{Transformation der Restmatrix}

for $j = k + 1$ to n **do**

$$\beta = 0$$

for $i = k$ to n **do**

$$\beta = \beta + a_{ik} \cdot a_{ij}$$

endfor

$$\beta = \beta / \gamma_k$$

for $i = k$ to n **do**

$$a_{ij} = a_{ij} - \beta \cdot a_{ik}$$

endfor

endfor

endfor

$$\varrho_n = a_{nn}$$

Aufwand: $\sim 2n^3/3$ Additionen/Multiplikationen, n Quadratwurzeln.

Bei gegebener QR -Zerlegung einer Matrix \bar{A} lässt sich dann jedes lineare System $Ax = b$ in zwei Schritten lösen:

$$c = Q^T b, \quad Rx = c.$$

Wurde die Matrix Q als Produkt von HOUSEHOLDER-Matrizen abgespeichert, so stellt sich Q^T ebenfalls als Produkt von HOUSEHOLDER-Matrizen dar. Aus

$$Q = H^{(1)} H^{(2)} \dots H^{(n-2)} H^{(n-1)}$$

folgt wegen der Symmetrie der HOUSEHOLDER-Matrizen

$$Q^T = H^{(n-1)} H^{(n-2)} \dots H^{(2)} H^{(1)}.$$

Damit lässt sich $c = Q^T b$ nach folgendem Algorithmus berechnen.

8.48. Lösen eines Gleichungssystems $Ax = b$ bei gegebener QR -Zerlegung der Matrix A :

Es ist das Gleichungssystem $Ax = b$ zu lösen.

Für die Matrix A wurde mit obigem Algorithmus eine QR -Zerlegung berechnet.

{Berechnen von $c = Q^T b$ }

for $k = 1$ **to** $n - 1$ **do**

$$\beta = 0$$

for $i = k$ **to** n **do**

$$\beta = \beta + b_i \cdot a_{ik}$$

endfor

$$\beta = \beta / \gamma_k$$

for $i = k$ **to** n **do**

$$b_i = b_i - \beta \cdot a_{ik}$$

endfor

endfor

{Lösen von $Rx = c$ }

for $k = n$ **to** 1 **step** -1 **do**

for $i = k + 1$ **to** n **do**

$$b_k = b_k - b_i \cdot a_{ki}$$

endfor

$$b_k = b_k / \varrho_k$$

endfor

Aufwand: $\sim n^2$ Additionen/Multiplikationen.

Ohne Beweis sei ein Satz über das Rundungsfehlerverhalten des HOUSEHOLDER-Verfahrens angegeben.

8.49. Satz: Die HOUSEHOLDER-Orthogonalisierung ist für eine reguläre Matrix A mit $\varrho_k = r_{kk} \neq 0$, $k = 1, \dots, n$, durchführbar, falls

$$\kappa = \text{eps}F \text{cond}(A) < 1$$

mit

$$F = 3.14 \left(1 + \frac{3}{n}\right) n^{5/2} \approx 4n^{5/2}$$

gilt. Es existiert eine orthogonale Matrix \hat{Q} , so dass für den berechneten Dreiecksfaktor \mathcal{R} gilt

$$\hat{Q}\mathcal{R} = A + \delta A$$

mit

$$\|\delta A\| \leq \text{eps}F \|A\|.$$

Für die durch die Vektoren $\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(n-1)}$ repräsentierte Matrix \mathcal{Q} gilt

$$\tilde{\mathbf{c}} = \text{gl} \left(\mathcal{Q}^T \mathbf{b} \right) = \hat{Q}^T (\mathbf{b} + \delta \mathbf{b})$$

mit

$$\|\delta \mathbf{b}\| \leq \text{eps} \frac{F}{\sqrt{n}} \|\mathbf{b}\|.$$

Bemerkung: Das zu lösende Dreieckssystem $\mathcal{R}\mathbf{x} = \tilde{\mathbf{c}}$ entsteht damit durch eine exakte orthogonale Transformation mit der Matrix \hat{Q} aus dem gestörten Gleichungssystem

$$(A + \delta A)\mathbf{x} = \mathbf{b} + \delta \mathbf{b}.$$

Das Lösen eines linearen Gleichungssystems mit dem HOUSEHOLDER-Verfahren ist demnach ein numerisch gutartiger Prozess.

8.50. Beispiel: Wir betrachten wieder die Matrix

$$A = \begin{pmatrix} 1 & 1 & 1 & 1 \\ \varepsilon & 0 & 0 & 0 \\ 0 & \varepsilon & 0 & 0 \\ 0 & 0 & \varepsilon & 0 \end{pmatrix}$$

mit $\varepsilon > \text{eps}$ und $\varepsilon^2 < \text{eps}$.

Das HOUSEHOLDER-Verfahren liefert im Maschinenzahlbereich $\mathbb{M}(10, 2, \dots)$ die Faktoren

$$\mathbf{Q} = \begin{pmatrix} -1 & \frac{\varepsilon}{\sqrt{2}} & \frac{\varepsilon}{\sqrt{6}} & -\frac{\varepsilon}{\sqrt{3}} \\ -\varepsilon & -\frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} \\ 0 & \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{6}} & \frac{1}{\sqrt{3}} \\ 0 & 0 & \frac{2}{\sqrt{6}} & \frac{1}{\sqrt{3}} \end{pmatrix}, \quad \mathbf{R} = \begin{pmatrix} -1 & -1 & -1 & -1 \\ 0 & \varepsilon\sqrt{2} & \frac{\varepsilon}{\sqrt{2}} & \frac{\varepsilon}{\sqrt{2}} \\ 0 & 0 & \frac{3\varepsilon}{\sqrt{6}} & \frac{\varepsilon}{\sqrt{6}} \\ 0 & 0 & 0 & -\frac{\varepsilon}{\sqrt{3}} \end{pmatrix}.$$

Damit gilt

$$\delta\mathbf{A} = \mathbf{Q}\mathbf{R} - \mathbf{A} = \begin{pmatrix} 0 & \varepsilon^2 & \varepsilon^2 & \varepsilon^2 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}$$

und

$$\mathbf{Q}^T \mathbf{Q} - \mathbf{I} = \begin{pmatrix} \varepsilon^2 & 0 & 0 & 0 \\ 0 & -\frac{\varepsilon^2}{2} & \frac{\sqrt{3}\varepsilon^2}{6} & -\frac{\sqrt{6}\varepsilon^2}{6} \\ 0 & \frac{\sqrt{3}\varepsilon^2}{6} & \frac{\varepsilon^2}{6} & -\frac{\sqrt{2}\varepsilon^2}{3} \\ 0 & -\frac{\sqrt{6}\varepsilon^2}{6} & -\frac{\sqrt{2}\varepsilon^2}{3} & \frac{\varepsilon^2}{3} \end{pmatrix}.$$

Das ist eine im Rahmen der Maschinengenauigkeit exakte Zerlegung mit einer im Rahmen der Maschinengenauigkeit exakten orthogonalen Matrix. \heartsuit

8.3. Iterative Verfahren

8.3.1. Iterationsverfahren und ihre Konvergenz

Neben den direkten Lösungsverfahren für lineare Gleichungssysteme, die nach einer wohlbestimmten Maximalzahl von Rechenoperationen eine Lösung liefern, existieren iterative Verfahren, die zu einem linearen Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ und einem Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^n$ eine Folge $\{\mathbf{x}^{(i)}\}_{i \in \mathbb{N}} \subset \mathbb{R}^n$ liefern, die gegen die Lösung des Gleichungssystems konvergiert. Charakteristisch für die direkten Verfahren war, dass die Koeffizientenmatrix des Systems in ein Produkt einfacherer Matrizen zerlegt wurde. Damit dürfen wir die direkten Verfahren auch als Zerlegungsmethoden oder Faktorisierungsverfahren bezeichnen. Charakteristisch für die iterativen Verfahren

wird sein, dass man die Matrix A in unveränderter Form im Algorithmus verwendet. Wie wir später sehen werden, benötigen wir eigentlich nicht die Matrix A selbst, sondern nur einen Algorithmus, der zu gegebenem $x \in \mathbb{R}^n$ den Vektor $y = Ax$ berechnet. Damit sind diese Verfahren besonders für große Systeme mit schwach besetzten Matrizen attraktiv.

Damit ein Iterationsverfahren sinnvoll und gegenüber einem direkten Verfahren konkurrenzfähig ist, sollte es folgende Bedingungen erfüllen:

- A1** Der Aufwand für den Übergang $x^{(i)} \rightarrow x^{(i+1)}$ sollte möglichst gering sein. Die Anzahl der benötigten Rechenoperationen für einen Iterationsschritt sollte in der Größenordnung n^2 oder kleiner liegen. Der Aufwand für einen Iterationsschritt sollte den Aufwand für das Berechnen von Matrix \times Vektor nicht wesentlich überschreiten.
- A2** Das Verfahren sollte für beliebige Startvektoren $x^{(0)}$ konvergieren. Die Konvergenzgeschwindigkeit sollte möglichst groß sein.

Wir betrachten folgenden allgemeinen Ansatz: Es sei $Ax = b$ ein lineares Gleichungssystem mit der regulären Matrix A ; der Vektor $x = A^{-1}b$ ist die exakte Lösung dieses Systems. Weiterhin sei B eine beliebige reguläre Matrix. Dann gilt

$$\begin{aligned} (A - B + B)x &= b, \\ Bx &= (B - A)x + b, \\ x &= B^{-1}(B - A)x + B^{-1}b, \\ x &= (I - B^{-1}A)x + B^{-1}b. \end{aligned}$$

Das lineare Gleichungssystem $Ax = b$ ist zur Fixpunktaufgabe

$$x = \Phi(x) = (I - B^{-1}A)x + B^{-1}b$$

äquivalent. Zum Lösen der Fixpunktaufgabe liegt es nahe, ein Iterationsverfahren anzuwenden:

Wähle $x^{(0)} \in \mathbb{R}^n$,

Für $k = 0, 1, \dots$ berechne $x^{(k+1)}$ durch Lösen von

$$Bx^{(k+1)} = (B - A)x^{(k)} + b.$$

Je nach Wahl von B erhält man verschiedene Iterationsverfahren. Um der Anforderung **A1** zu genügen, ist B möglichst einfach zu wählen. Es bieten sich wieder Diagonalmatrizen, Dreiecksmatrizen oder orthogonale Matrizen an. Die Forderung **A2** ist offensichtlich erfüllt, falls Φ auf dem gesamten \mathbb{R}^n ein kontrahierender Operator ist. Das ist eine Forderung an die Matrix $M = I - B^{-1}A$. Es gilt der

8.51. Satz:

1. Das Iterationsverfahren

$$\mathbf{x}^{(k+1)} = (\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})\mathbf{x}^{(k)} + \mathbf{B}^{-1}\mathbf{b}$$

ist konvergiert dann und nur dann, wenn der Spektralradius der Matrix $M = \mathbf{I} - \mathbf{B}^{-1}\mathbf{A}$ kleiner als 1 ist:

$$\rho(\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}) < 1.$$

2. Hinreichend für die Konvergenz ist die Bedingung

$$\|(\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})\| < 1$$

bezüglich einer beliebigen Matrixnorm, die mit einer Vektornorm verträglich ist.

Beweis:

1. Wir nehmen zunächst an, dass die Iteration konvergiert. Dann gilt

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}$$

Mit $\mathbf{A}\mathbf{x} = \mathbf{b}$. Für den Fehlervektor $\mathbf{f}^{(k)} = \mathbf{x}^{(k)} - \mathbf{x}$ folgt

$$\begin{aligned} \mathbf{f}^{(k+1)} &= \mathbf{x}^{(k+1)} - \mathbf{x} \\ &= (\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})\mathbf{x}^{(k)} + \mathbf{B}^{-1}\mathbf{b} - \mathbf{x} \\ &= (\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})(\mathbf{x}^{(k)} - \mathbf{x}) + (\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})\mathbf{x} + \mathbf{B}^{-1}\mathbf{b} - \mathbf{x} \\ &= (\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})\mathbf{f}^{(k)} + \mathbf{x} - \mathbf{B}^{-1}\mathbf{b} + \mathbf{B}^{-1}\mathbf{b} - \mathbf{x} \\ &= (\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})\mathbf{f}^{(k)}. \end{aligned}$$

Aus der Konvergenz folgt weiterhin

$$\lim_{k \rightarrow \infty} \mathbf{f}^{(k)} = \mathbf{o}.$$

Wählt man den Startvektor so, dass $\mathbf{f}^{(0)}$ ein Eigenvektor der Matrix $\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}$ zum Eigenwert λ ist, so gilt

$$\mathbf{f}^{(k)} = (\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})^k \mathbf{f}^{(0)} = \lambda^k \mathbf{f}^{(0)}.$$

Wegen $\mathbf{f}^{(0)} \neq \mathbf{o}$ muss dann aber

$$\lim_{k \rightarrow \infty} \lambda^k = 0$$

und damit $|\lambda| < 1$ gelten. Da λ ein beliebiger Eigenwert von $\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}$ war, folgt $\varrho(\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}) < 1$.

Nun nehmen wir an, dass $\varrho(\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}) < 1$ gilt. Es folgt

$$\lim_{k \rightarrow \infty} \varrho\left((\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})^k\right) = \lim_{k \rightarrow \infty} \left(\varrho(\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})\right)^k = 0,$$

also

$$\lim_{k \rightarrow \infty} (\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})^k = \mathbf{O}$$

und

$$\lim_{k \rightarrow \infty} \mathbf{f}^{(k)} = (\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})^k \mathbf{f}^{(0)} = \mathbf{o}.$$

2. Nach Satz 8.7 folgt aus $\|\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}\| < 1$ sofort

$$\varrho(\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}) \leq \text{lub}(\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}) \leq \|\mathbf{I} - \mathbf{B}^{-1}\mathbf{A}\| < 1$$

und mit der ersten Aussage des Satzes die Konvergenz des Verfahrens.

✱

Die Konvergenzgeschwindigkeit des Verfahrens hängt von $\varrho(\mathbf{I} - \mathbf{B}^{-1}\mathbf{A})$ ab. Je kleiner der Spektralradius, desto schneller die Konvergenz.

8.3.2. Das Jacobi- und das Gauß-Seidel-Verfahren

Es sei $\mathbf{A} \in \mathbb{R}^{n \times n}$ eine Matrix mit $a_{ii} \neq 0$ für $i = 1, \dots, n$. Wir wählen

$$\mathbf{B} = \text{diag}(a_{11}, a_{22}, \dots, a_{nn}).$$

Damit ergibt sich aus dem allgemeinen Ansatz für Iterationsverfahren das **JACOBI-Verfahren**:

8.52. JACOBI-Verfahren:

Es ist das Gleichungssystem $Ax = b$ zu lösen. Für die Matrix $A \in \mathbb{R}^{n \times n}$ gelte

$$a_{ii} \neq 0, \quad i = 1, \dots, n.$$

S0 Wähle einen Startvektor $x^{(0)} \in \mathbb{R}^{n \times n}$ und setze $i = 0$.

S1 Berechne

for $k = 1$ **to** n **do**

$$x_k^{(i+1)} = \frac{1}{a_{kk}} \left(b_k - \sum_{\substack{j=1 \\ j \neq k}}^n a_{kj} x_j^{(i)} \right)$$

end

S2 Setze $i = i + 1$ und gehe zu Schritt **S1**.

Das JACOBI-Verfahren wird auch als **Iteration in Gesamtschritten** bezeichnet. Als Iterationsmatrix $M = I - B^{-1}A$ erhält man

$$M = \begin{pmatrix} 0 & -\frac{a_{12}}{a_{11}} & \dots & -\frac{a_{1n}}{a_{11}} \\ -\frac{a_{21}}{a_{22}} & 0 & \dots & -\frac{a_{2n}}{a_{22}} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{a_{n1}}{a_{nn}} & -\frac{a_{n2}}{a_{nn}} & \dots & 0 \end{pmatrix}.$$

Nach Satz 8.51 konvergiert das Gesamtschrittverfahren genau dann, wenn $\rho(M) < 1$ gilt. Diese Bedingung läßt sich im allgemeinen schwer nachprüfen. Hinreichende Konvergenzbedingungen erhält man jedoch leicht.

8.53. Satz:**1. Starkes Zeilensummenkriterium**

Das Gesamtschrittverfahren konvergiert für alle Matrizen $A \in \mathbb{R}^{n \times n}$ mit

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \quad i = 1, \dots, n.$$

2. Starkes Spaltensummenkriterium

Das Gesamtschrittverfahren konvergiert für alle Matrizen $A \in \mathbb{R}^{n \times n}$ mit

$$|a_{jj}| > \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}| \quad \text{für } j = 1, \dots, n.$$

Beweis:

1. Aus

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \quad i = 1, \dots, n$$

folgt

$$\sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{|a_{ii}|} < 1 \quad i = 1, \dots, n$$

und

$$1 > \max_{1 \leq i \leq n} \left\{ \sum_{\substack{j=1 \\ j \neq i}}^n \left| \frac{a_{ij}}{a_{ii}} \right| \right\} = \|M\|_{\infty}.$$

Damit ist das hinreichende Konvergenzkriterium aus Satz 8.51 erfüllt.

2. Ist A streng spaltendiagonaldominant, so ist A^T streng zeilendiagonaldominant. Damit konvergiert das JACOBI-Verfahren für die Matrix A^T . Dann ist nach Satz 8.51 der Spektralradius der Iterationsmatrix $\bar{M} = I - B^{-1}A^T$ kleiner als 1. Da die Matrix

$$B\bar{M}B^{-1} = I - A^T B^{-1} = \left(I - B^{-1}A \right)^T$$

dieselben Eigenwerte wie \bar{M} hat, folgt

$$\varrho(M) = \varrho(M^T) = \varrho(B\bar{M}B^{-1}) = \varrho(\bar{M}) < 1.$$

Damit konvergiert das Gesamtschrittverfahren auch für die Matrix A .

*

Benutzt man in der Iterationsvorschrift des JACOBI-Verfahrens

$$x_i^{(k+1)} = \frac{1}{a_{ii}} \left(b_i - \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j^{(k)} \right)$$

beim Berechnen der i -ten Komponente des neuen Iterationsvektors statt der Komponenten

$$x_1^{(k)}, \dots, x_{i-1}^{(k)}$$

die schon berechneten neuen Komponenten

$$x_1^{(k+1)}, \dots, x_{i-1}^{(k+1)},$$

so ergibt sich das **GAUSS-SEIDEL-Verfahren**.

8.54. GAUSS-SEIDEL-Verfahren:

Es ist das Gleichungssystem $A\mathbf{x} = \mathbf{b}$ zu lösen. Für die Matrix $A \in \mathbb{R}^{n \times n}$ gelte

$$a_{ii} \neq 0, \quad i = 1, \dots, n.$$

S0 Wähle einen Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^{n \times n}$ und setze $i = 0$.

S1 Berechne

for $k = 1$ **to** n **do**

$$x_k^{(i+1)} = \frac{1}{a_{kk}} \left(b_k - \sum_{j=1}^{k-1} a_{kj} x_j^{(i+1)} - \sum_{j=k+1}^n a_{kj} x_j^{(i)} \right)$$

end

S2 Setze $i = i + 1$ und gehe zu Schritt **S1**.

Dieses Verfahren wird auch als **Einzelschrittverfahren** bezeichnet. Dem Verfahren entspricht die Wahl

$$B = \begin{pmatrix} a_{11} & & & \\ a_{21} & a_{22} & & \mathbf{O} \\ \vdots & \vdots & \ddots & \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{pmatrix}.$$

Als hinreichende Bedingungen für die Konvergenz des Einzelschrittverfahrens ergeben sich ebenfalls das starke Zeilen- oder Spaltensummenkriterium aus Satz 8.53.

8.55. Satz:

1. Starkes Zeilensummenkriterium

Das Einzelschrittverfahren konvergiert für alle Matrizen $A \in \mathbb{R}^{n \times n}$ mit

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \quad i = 1, \dots, n.$$

2. Starkes Spaltensummenkriterium

Das Einzelschrittverfahren konvergiert für alle Matrizen $A \in \mathbb{R}^{n \times n}$ mit

$$|a_{jj}| > \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}| \quad j = 1, \dots, n.$$

Beweis: Wir beweisen nur die erste Aussage. Es sei

$$A = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix} = -E + D - F$$

mit

$$E = - \begin{pmatrix} 0 & \cdots & \cdots & 0 \\ a_{21} & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ a_{n1} & \cdots & a_{n,n-1} & 0 \end{pmatrix}, \quad D = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$$

und

$$F = - \begin{pmatrix} 0 & a_{12} & \cdots & a_{1n} \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & a_{n-1,n} \\ 0 & \cdots & \cdots & 0 \end{pmatrix}.$$

Aus der Zeilendiagonaldominanz der Matrix A folgt $a_{ii} \neq 0$ für $i = 1, \dots, n$. Damit ist D regulär. Wir definieren weiterhin die Matrizen $L = D^{-1}E$ und $U = D^{-1}F$. Als Iterationsmatrizen für das JACOBI- bzw. GAUSS-SEIDEL-Verfahren erhalten wir

$$M_J = I - D^{-1}(-E + D - F) = L + U$$

und

$$M_{GS} = I - (-E + D)^{-1}(-E + D - F) = (-DL + D)^{-1}DU = (I - L)^{-1}U.$$

Für eine streng zeilendiagonaldominante Matrix A gilt dann $\|M_J\|_\infty < 1$. Wir zeigen, dass

$$\|M_{GS}\|_\infty \leq \|M_J\|_\infty < 1$$

gilt. Aus $\|M_J\|_\infty < 1$ folgt

$$\|M_J|e \leq \|M_J\|_\infty e < e \quad \text{mit} \quad e = (1, 1, \dots, 1)^T \in \mathbb{R}^n.$$

Weiterhin gilt $|M_J| = |L| + |U|$ wegen der speziellen Struktur der Matrizen L und U . Daraus folgt

$$|U|e = (|M_J| - |L|)e \leq (\|M_J\|_\infty I - |L|)e. \quad (8.13)$$

L und $|L|$ sind untere Dreiecksmatrizen mit verschwindender Diagonale. Man sieht leicht ein, dass für diese Matrizen

$$L^n = |L|^n = O$$

gilt. Damit existieren die Inversen von $I - L$ und $I - |L|$, und es gilt

$$O \leq \left| (I - L)^{-1} \right| = \left| I + L + \dots + L^{n-1} \right| \leq I + |L| + \dots + |L|^{n-1} = (I - |L|)^{-1}.$$

Mit dieser Ungleichung folgt aus 8.13

$$\begin{aligned} |M_{GS}|e &\leq \left| (I - L)^{-1} \right| |U|e \\ &\leq (I - |L|)^{-1} (\|M_J\|_\infty I - |L|)e \\ &= (I - |L|)^{-1} (I - |L| + (\|M_J\|_\infty - 1)I)e \\ &= \left(I + (\|M_J\|_\infty - 1)(I - |L|)^{-1} \right) e. \end{aligned}$$

Wegen $(I - |L|)^{-1} \geq I$ und $\|M_J\|_\infty < 1$ ergibt sich

$$|M_{GS}|e \leq [I + (\|M_J\|_\infty - 1)I]e = \|M_J\|_\infty e$$

und damit

$$\|M_{GS}\|_\infty \leq \|M_J\|_\infty.$$

*

Welches von beiden Verfahren für eine beliebige Matrix konvergiert und welches schneller konvergiert, läßt sich nicht ohne weiteres angeben.

8.56. Beispiel: Für die Matrix

$$A = \begin{pmatrix} 1 & -2 & 2 \\ -1 & 1 & -1 \\ -2 & -2 & 1 \end{pmatrix}$$

ergeben sich für das JACOBI- bzw. GAUSS-SEIDEL-Verfahren die Iterationsmatrizen

$$M_J = \begin{pmatrix} 0 & 2 & -2 \\ 1 & 0 & 1 \\ 2 & 2 & 0 \end{pmatrix} \quad \text{und} \quad M_{GS} = \begin{pmatrix} 0 & 2 & -2 \\ 0 & 2 & -1 \\ 0 & 8 & -6 \end{pmatrix}$$

mit $\varrho(M_J) = 0$ und $\varrho(M_{GS}) = 2(1 + \sqrt{2})$. Hier konvergiert das JACOBI-Verfahren, das GAUSS-SEIDEL-Verfahren dagegen nicht. Für die Matrix

$$A = \begin{pmatrix} 1 & \frac{1}{2} & \frac{1}{2} \\ -1 & 1 & -1 \\ -\frac{1}{2} & \frac{1}{2} & 1 \end{pmatrix}$$

ergeben sich

$$M_J = \begin{pmatrix} 0 & -\frac{1}{2} & -\frac{1}{2} \\ 1 & 0 & 1 \\ \frac{1}{2} & -\frac{1}{2} & 0 \end{pmatrix} \quad \text{und} \quad M_{GS} = \begin{pmatrix} 0 & -\frac{1}{2} & -\frac{1}{2} \\ 0 & -\frac{1}{2} & \frac{1}{2} \\ 0 & 0 & -\frac{1}{2} \end{pmatrix}$$

mit

$$\varrho(M_J) = \sqrt{5}/2, \quad \varrho(M_{GS}) = 1/2.$$

Hier konvergiert das GAUSS-SEIDEL-Verfahren; das JACOBI-Verfahren jedoch konvergiert nicht. ♡

Das starke Zeilen- bzw. Spaltensummenkriterium läßt sich abschwächen. Dazu führen wir den Begriff der unzerlegbaren Matrizen ein. Eine Matrix $A \in \mathbb{R}^{n \times n}$ heißt **unzerlegbar (irreduzibel)**, falls keine Permutationsmatrix $P \in \mathbb{R}^{n \times n}$ existiert, so dass $P^T A P$ die Gestalt

$$P^T A P = \begin{pmatrix} A_{11} & A_{12} \\ O & A_{22} \end{pmatrix}$$

mit $A_{11} \in \mathbb{R}^{k \times k}$, $A_{22} \in \mathbb{R}^{(n-k) \times (n-k)}$ und $1 \leq k \leq n-1$ besitzt. Die Unzerlegbarkeit einer Matrix läßt sich mit Hilfe der Graphentheorie überprüfen. Wir ordnen dazu einer Matrix $A \in \mathbb{R}^{n \times n}$ in folgender Weise einen gerichteten Graphen $G(A)$ zu. $G(A)$ besitzt die Knoten P_1, \dots, P_n . Vom Knoten P_i zum Knoten P_j existiert genau dann eine gerichtete Kante e_{ij} , wenn $a_{ij} \neq 0$ gilt. Die Matrix A ist genau dann unzerlegbar, wenn der Graph $G(A)$ zusammenhängend ist. Dabei heißt ein

Graph zusammenhängend, falls für jedes Paar (i, j) mit $i \neq j$ ein gerichteter Weg vom Knoten P_i zum Knoten P_j existiert.

Nun können wir ein schwächeres Konvergenzkriterium für das JACOBI-Verfahren formulieren.

8.57. Satz: Für eine unzerlegbare Matrix $A \in \mathbb{R}^{n \times n}$ gilt:

1. Schwaches Zeilensummenkriterium

Das JACOBI-Verfahren konvergiert, falls

$$|a_{ii}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| \quad \text{für } i = 1, \dots, n$$

und für mindestens ein i^*

$$|a_{i^*i^*}| > \sum_{\substack{j=1 \\ j \neq i^*}}^n |a_{i^*j}|$$

gilt.

2. Schwaches Spaltensummenkriterium

Das JACOBI-Verfahren konvergiert, falls

$$|a_{jj}| \geq \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}| \quad \text{für } j = 1, \dots, n$$

und für mindestens ein j^*

$$|a_{j^*j^*}| > \sum_{\substack{i=1 \\ i \neq j^*}}^n |a_{ij^*}|$$

gilt.

Beweis: Wir beweisen nur das Schwache Zeilensummenkriterium. Die Gültigkeit des Schwachen Spaltensummenkriteriums folgt dann analog zum Beweis des Starken Spaltensummenkriteriums. Aus dem Schwachen Zeilensummenkriterium folgt für das JACOBI-Verfahren

$$\|M_J\|_\infty = \|I - D^{-1}A\|_\infty \leq 1.$$

Als hinreichende Bedingung für die Konvergenz benötigen wir aber $\|M_J\|_\infty < 1$. Mit $e = (1, 1, \dots, 1)^T$ folgt

$$|M|e \leq e \quad \text{und} \quad |M|e \neq e.$$

Wir zeigen, dass $|M|^n e < e$ gilt. Daraus folgt dann

$$1 > \| |M|^n \|_\infty \geq (\| |M| \|_\infty)^n = (\|M\|)^n \geq (\rho(M))^n$$

und endlich $\rho(M) < 1$. Um die Beziehung $|M|^n e < e$ zu zeigen, betrachten wir die Ungleichungskette

$$e \stackrel{\geq}{\neq} |M|e \geq |M|^2 e \geq \dots \geq |M|^i e \geq |M|^{i+1} e \geq \dots$$

Daraus folgt

$$o \stackrel{\leq}{\neq} e - |M|e \leq e - |M|^2 e \leq \dots \leq e - |M|^i e \leq e - |M|^{i+1} e \leq \dots$$

Von den Vektoren $t^{(i)} = e - |M|^i e$ sind damit gewisse Komponenten gleich Null und die restlichen Komponenten größer als Null. Dann existiert eine Permutationsmatrix P , so dass

$$Pt^{(i)} = P(e - |M|^{i+1} e) = \begin{pmatrix} a \\ o \end{pmatrix} \quad \text{mit} \quad a \in \mathbb{R}^p \quad \text{und} \quad a > o.$$

Wegen $0 \neq t^{(i)}$ gilt $p < n$. Für den Vektor $t^{(i+1)} = e - |M|^{i+1} e$ gilt $t^{(i+1)} \geq t^{(i)}$, daher

$$Pt^{(i+1)} \geq Pt^{(i)} = \begin{pmatrix} a \\ o \end{pmatrix}.$$

Es gilt somit

$$Pt^{(i+1)} = \begin{pmatrix} b \\ o \end{pmatrix} \quad \text{mit} \quad b \in \mathbb{R}^q, \quad b > o$$

und $q \geq p$. *

8.3.3. Nachiteration

Gegeben sei eine näherungsweise LU -Zerlegung der Matrix A

$$P^T \mathcal{L}U = A + \delta A.$$

Wir wählen im allgemeinen Iterationsansatz

$$B = P^T \mathcal{L}U = A + \delta A$$

und erhalten

$$\begin{aligned} M &= I - B^{-1}A = I - B^{-1}(B - \delta A) \\ &= B^{-1}\delta A = (A + \delta A)^{-1}\delta A. \end{aligned}$$

Es gilt

$$\|M\| = \|(A + \delta A)^{-1}\delta A\| \leq \|(A + \delta A)^{-1}\| \|\delta A\|$$

und im Falle $\kappa = \|\delta A\| \|A^{-1}\| < 1$

$$\|M\| \leq \frac{\kappa}{1 - \kappa}.$$

Für die LU -Zerlegung mit Spaltenpivotisierung gilt

$$\kappa = \text{eps}F(A)\text{cond}(A)$$

mit $F(A) \approx 3 \cdot 2^n$. Die hinreichende Bedingung $\|M\| < 1$ ist erfüllt, falls $\kappa < 1/2$ gilt. Wir erhalten folgenden Algorithmus.

8.58. Nachiteration für die LU -Zerlegung:

Zu lösen sei das lineare Gleichungssystem $Ax = b$ mit $A \in \mathbb{R}^{n \times n}$ und $b \in \mathbb{R}^n$. Durch \mathcal{L} und U seien die berechneten Dreiecksfaktoren einer LU -Zerlegung der Matrix A gegeben. Es gelte

$$A + \delta A = P^T \mathcal{L}U.$$

S0 Wähle einen Startvektor $x^{(0)} \in \mathbb{R}^n$ und setze $i = 0$.

S1 Berechne das Residuum

$$r^{(i)} = b - Ax^{(i)}.$$

S2 Löse das Gleichungssystem $A\delta x^{(i)} = r^{(i)}$ näherungsweise gemäß

$$\mathcal{L}y = Pr^{(i)} \quad \text{und} \quad U\delta x^{(i)} = y.$$

S3 Setze $x^{(i+1)} = x^{(i)} + \delta x^{(i)}$, $i = i + 1$ und gehe zu Schritt **S1**.

Bemerkungen: (i) Wenn das Verfahren konvergiert, werden die $\mathbf{x}^{(i)}$ gute Näherungen für die Lösung des Gleichungssystems $\mathbf{Ax} = \mathbf{b}$ sein. Dann tritt aber beim Berechnen des Residuums in Schritt **S1** Auslöschung auf. Damit das ganze Verfahren noch sinnvoll bleiben soll, muss man an dieser Stelle besonders sorgfältig vorgehen. Am besten berechnet man das Residuum in einer höheren Genauigkeit.

(ii) Üblicherweise wählt man als Startvektor $\mathbf{x}^{(0)} = \mathbf{o}$. Damit entspricht der erste Schritt des Verfahrens der normalen Lösung des Gleichungssystems mittels LU -Zerlegung.

(iii) Die Konvergenz des Verfahrens ist im allgemeinen gut, falls \mathbf{A} nicht zu schlecht konditioniert ist. Es werden nur wenige Schritte ausgeführt.

Es gilt $\mathbf{Ax}^{(i)} = \mathbf{b} - \mathbf{r}^{(i)}$. Wir können $\mathbf{r}^{(i)}$ als eine Störung $\delta\mathbf{b}$ interpretieren. Nach Satz 8.14 gilt dann ($\delta\mathbf{A} = \mathbf{O}$)

$$\frac{1}{\text{cond}(\mathbf{A})} \frac{\|\mathbf{r}^{(i)}\|}{\|\mathbf{b}\|} \leq \frac{\|\delta\mathbf{x}^{(i)}\|}{\|\mathbf{x}^{(i)}\|} \leq \text{cond}(\mathbf{A}) \frac{\|\mathbf{r}^{(i)}\|}{\|\mathbf{b}\|}.$$

Diese Ungleichungen werden benutzt, um bei der Nachiteration die Konditionszahl von \mathbf{A} abzuschätzen. Man erhält

$$\text{cond}(\mathbf{A}) \geq \frac{\|\delta\mathbf{x}^{(i)}\|}{\|\mathbf{x}^{(i)}\|} \bigg/ \frac{\|\mathbf{r}^{(i)}\|}{\|\mathbf{b}\|}$$

beziehungsweise

$$\text{cond}(\mathbf{A}) \geq \frac{\|\mathbf{r}^{(i)}\|}{\|\mathbf{b}\|} \bigg/ \frac{\|\delta\mathbf{x}^{(i)}\|}{\|\mathbf{x}^{(i)}\|}.$$

Analoge Nachiterationsverfahren sind natürlich auch für andere direkte Verfahren möglich. Für die LDL^T -Zerlegung erhält man folgenden Algorithmus.

8.59. Nachiteration für die LDL^T -Zerlegung:

Zu lösen ist das lineare Gleichungssystem $\mathbf{Ax} = \mathbf{b}$ mit der symmetrischen Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ und $\mathbf{b} \in \mathbb{R}^n$. Durch \mathcal{L} und \mathcal{D} seien die berechneten Faktoren einer LDL^T -Zerlegung der Matrix \mathbf{A} gegeben. Es gelte

$$\mathbf{A} + \delta\mathbf{A} = \mathcal{L}\mathcal{D}\mathcal{L}^T.$$

S0 Wähle einen Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^n$ und setze $i = 0$.

Berechne das Residuum

$$\mathbf{r}^{(i)} = \mathbf{b} - \mathbf{Ax}^{(i)}.$$

S1 Löse das Gleichungssystem $\mathbf{A}\delta\mathbf{x}^{(i)} = \mathbf{r}^{(i)}$ näherungsweise gemäß

$$\mathcal{L}\mathbf{y} = \mathbf{r}^{(i)} \quad \text{und} \quad \mathcal{L}^T\delta\mathbf{x}^{(i)} = \mathcal{D}^{-1}\mathbf{y}.$$

S2 Setze $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} + \delta\mathbf{x}^{(i)}$, $i = i + 1$ und gehe zu Schritt **S1**.

Für ein Orthogonalisierungsverfahren ergibt sich das folgende Nachiterationsverfahren.

8.60. Nachiteration für die QR-Zerlegung:

Zu lösen ist das lineare Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit $\mathbf{A} \in \mathbb{R}^{n \times n}$ und $\mathbf{b} \in \mathbb{R}^n$. Durch \mathcal{Q} und \mathcal{R} seien die berechneten Faktoren einer QR-Zerlegung der Matrix \mathbf{A} gegeben. Es gelte

$$\mathbf{A} + \delta\mathbf{A} = \mathcal{Q}\mathcal{R}.$$

S0 Wähle einen Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^n$ und setze $i = 0$.

S1 Berechne das Residuum

$$\mathbf{r}^{(i)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(i)}.$$

S2 Löse das Gleichungssystem $\mathbf{A}\delta\mathbf{x}^{(i)} = \mathbf{r}^{(i)}$ näherungsweise gemäß

$$\mathbf{y} = \mathcal{Q}^T\mathbf{r}^{(i)} \quad \text{und} \quad \mathcal{R}\delta\mathbf{x}^{(i)} = \mathbf{y}.$$

S3 Setze $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} + \delta\mathbf{x}^{(i)}$, $i = i + 1$ und gehe zu Schritt **S1**.

Analoge Konditionszahlsschätzungen sind möglich. Das Residuum sollte wieder in höherer Genauigkeit berechnet werden.

8.3.4. Das cg-Verfahren von Hestenes und Stiefel

Das jetzt zu behandelnde Verfahren gehört nicht zu den Iterationsverfahren, wie sie sich aus dem allgemeinen Ansatz in Abschnitt 8.3.1. ergeben. Bei genauer Betrachtung müßte man es sogar zu den direkten Verfahren zählen, da man bei exakter Rechnung nach einer a priori bekannten Maximalzahl von Rechenoperationen die exakte Lösung des Gleichungssystems erhält. Da aber bei diesem Verfahren auch eine Folge von Vektoren konstruiert wird, die den exakten Lösungsvektor immer genauer approximieren, behandeln wir das Verfahren an dieser Stelle.

Wir betrachten zunächst ein lineares Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit einer symmetrischen, positiv definiten Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$. Diesem Gleichungssystem ordnen wir ein Optimierungsproblem zu. Es sei

$$F(\mathbf{z}) = \frac{1}{2}(\mathbf{A}\mathbf{z} - \mathbf{b})^T \mathbf{A}^{-1}(\mathbf{A}\mathbf{z} - \mathbf{b}) = \frac{1}{2}\mathbf{z}^T \mathbf{A}\mathbf{z} - \mathbf{b}^T \mathbf{z} + \frac{1}{2}\mathbf{b}^T \mathbf{A}^{-1}\mathbf{b}.$$

Wegen der positiven Definitheit von \mathbf{A} (und damit auch von \mathbf{A}^{-1}) gilt $F(\mathbf{z}) \geq 0$ für alle $\mathbf{z} \in \mathbb{R}^n$ und $F(\mathbf{z}) = 0$ genau dann, wenn $\mathbf{A}\mathbf{z} - \mathbf{b} = \mathbf{o}$. Das Minimierungsproblem

$$\mathbf{x} = \arg \min_{\mathbf{z} \in \mathbb{R}^n} F(\mathbf{z}) \quad (8.14)$$

ist dem Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ äquivalent. Minimierungsprobleme dieser Art löst man üblicherweise, indem man sie auf viele eindimensionale Minimierungsprobleme zurückführt.

S0 Wähle Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^n$ und setze $i = 0$.

S1 Wähle eine Richtung $\mathbf{h}^{(i)} \in \mathbb{R}^n$ und löse das eindimensionale Minimierungsproblem

$$\alpha_i = \arg \min_{\alpha \in \mathbb{R}} F(\mathbf{x}^{(i)} + \alpha \mathbf{h}^{(i)}).$$

S2 Setze $\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} + \alpha_i \mathbf{h}^{(i)}$.

S3 Setze $i = i + 1$ und gehe zu Schritt **S1**.

Für das gegebene Funktional $F(\mathbf{z})$ und beliebiges \mathbf{x} und eine beliebige Richtung \mathbf{h} läßt sich die Lösung des eindimensionalen Minimierungsproblems in Schritt **S1** leicht angeben. Es gilt der folgende

8.61. Satz: Für die symmetrische, positiv definite Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ und einen Vektor $\mathbf{b} \in \mathbb{R}^n$ sei das Funktional $F : \mathbb{R}^n \rightarrow \mathbb{R}_+$ folgendermaßen definiert:

$$F(\mathbf{z}) = \frac{1}{2} (\mathbf{A}\mathbf{z} - \mathbf{b})^T \mathbf{A}^{-1} (\mathbf{A}\mathbf{z} - \mathbf{b}).$$

Weiterhin seien ein Vektor $\mathbf{x} \in \mathbb{R}^n$ und ein Vektor $\mathbf{h} \in \mathbb{R}^n$ mit $\mathbf{h} \neq \mathbf{o}$ gegeben. Dann hat das Minimierungsproblem

$$\alpha^* = \arg \min_{\alpha \in \mathbb{R}} F(\mathbf{x} + \alpha \mathbf{h})$$

die eindeutige Lösung

$$\alpha^* = \frac{\mathbf{h}^T \mathbf{r}}{\mathbf{h}^T \mathbf{A} \mathbf{h}} \quad \text{mit} \quad \mathbf{r} = \mathbf{b} - \mathbf{A}\mathbf{x}.$$

Es gilt

$$\mathbf{h}^T [\mathbf{A}(\mathbf{x} + \alpha^* \mathbf{h}) - \mathbf{b}] = 0.$$

Beweis: Aus

$$F(\mathbf{z}) = \frac{1}{2} \mathbf{z}^T \mathbf{A} \mathbf{z} - \mathbf{b}^T \mathbf{z} + \frac{1}{2} \mathbf{b}^T \mathbf{A}^{-1} \mathbf{b}$$

folgt $F'(\mathbf{z}) = \mathbf{A} \mathbf{z} - \mathbf{b}$ und

$$\frac{d}{d\alpha} F(\mathbf{x} + \alpha \mathbf{h}) = \mathbf{h}^T [\mathbf{A}(\mathbf{x} + \alpha \mathbf{h}) - \mathbf{b}].$$

Eine notwendige Bedingung für ein Extremum ist

$$\left. \frac{d}{d\alpha} F(\mathbf{x} + \alpha \mathbf{h}) \right|_{\alpha^*} = 0.$$

Daraus ergibt sich

$$\mathbf{h}^T [\mathbf{A}(\mathbf{x} + \alpha^* \mathbf{h}) - \mathbf{b}] = 0$$

und

$$\alpha^* = \frac{\mathbf{h}^T (\mathbf{b} - \mathbf{A} \mathbf{x})}{\mathbf{h}^T \mathbf{A} \mathbf{h}} = \frac{\mathbf{h}^T \mathbf{r}}{\mathbf{h}^T \mathbf{A} \mathbf{h}}.$$

Es ist noch zu zeigen, dass α^* ein Minimum liefert. Es gilt

$$\begin{aligned} F(\mathbf{x} + \alpha^* \mathbf{h}) &= \frac{1}{2} (\mathbf{A} \mathbf{x} + \alpha^* \mathbf{A} \mathbf{h} - \mathbf{b})^T \mathbf{A}^{-1} (\mathbf{A} \mathbf{x} + \alpha^* \mathbf{A} \mathbf{h} - \mathbf{b}) \\ &= \frac{1}{2} (\mathbf{A} \mathbf{x} - \mathbf{b} + \alpha^* \mathbf{A} \mathbf{h})^T \mathbf{A}^{-1} (\mathbf{A} \mathbf{x} - \mathbf{b} + \alpha^* \mathbf{A} \mathbf{h}) \\ &= F(\mathbf{x}) + \frac{1}{2} \alpha^* \mathbf{h}^T (\mathbf{A} \mathbf{x} - \mathbf{b}) + \frac{1}{2} \alpha^* (\mathbf{A} \mathbf{x} - \mathbf{b})^T \mathbf{h} + \frac{1}{2} \alpha^{*2} \mathbf{h}^T \mathbf{A} \mathbf{h} \\ &= F(\mathbf{x}) - \alpha^* \mathbf{h}^T \mathbf{r} + \frac{1}{2} \alpha^{*2} \mathbf{h}^T \mathbf{A} \mathbf{h} \\ &= F(\mathbf{x}) - \frac{(\mathbf{h}^T \mathbf{r})^2}{\mathbf{h}^T \mathbf{A} \mathbf{h}} + \frac{1}{2} \frac{(\mathbf{h}^T \mathbf{r})^2}{(\mathbf{h}^T \mathbf{A} \mathbf{h})^2} \mathbf{h}^T \mathbf{A} \mathbf{h} \\ &= F(\mathbf{x}) - \frac{1}{2} \frac{(\mathbf{h}^T \mathbf{r})^2}{\mathbf{h}^T \mathbf{A} \mathbf{h}} \\ &\leq F(\mathbf{x}). \end{aligned}$$

Wir werden sehen, dass die Richtungen $\mathbf{h}^{(i)}$ so wählbar sind, dass nach höchstens n Schritten die Lösung des Minimierungsproblems 8.14 berechnet wird. Dazu definieren wir den Vektorraum

$$H_{i+1} = \text{span}(\mathbf{h}^{(0)}, \mathbf{h}^{(1)}, \dots, \mathbf{h}^{(i)}).$$

Das ist der Vektorraum, der von den Richtungen $\mathbf{h}^{(0)}, \mathbf{h}^{(1)}, \dots, \mathbf{h}^{(i)}$ aufgespannt wird. Nach dem Algorithmus ergibt sich $\mathbf{x}^{(i+1)}$ aus $\mathbf{x}^{(0)}$ nach

$$\mathbf{x}^{(i+1)} = \mathbf{x}^{(0)} + \alpha_0 \mathbf{h}^{(0)} + \alpha_1 \mathbf{h}^{(1)} + \dots + \alpha_i \mathbf{h}^{(i)}.$$

Es gilt also $\mathbf{x}^{(i+1)} - \mathbf{x}^{(0)} \in H_{i+1}$.

Wir wollen die Richtungen $\mathbf{h}^{(i)}$ und damit die Räume H_{i+1} so wählen, dass einerseits die Minimierung des Funktionals $F(\mathbf{z})$ auf dem gesamten \mathbb{R}^n auf eine schrittweise Minimierung des Funktionals auf den Teilräumen H_{i+1} zurückgeführt ist, und dass andererseits die Minimierung auf den Teilräumen möglichst einfach wird.

Wir nehmen an, wir hätten einen Teilraum H_i konstruiert und ein $\mathbf{x}^{(i)} \in \mathbf{x}^{(0)} + H_i$ bestimmt, das $F(\mathbf{z})$ auf $\mathbf{x}^{(0)} + H_i$ minimiert. $\mathbf{h}^{(i)}$ sei eine weitere Richtung, durch die H_{i+1} festgelegt ist. Jeder Vektor $\mathbf{y} \in \mathbf{x}^{(0)} + H_{i+1}$ ist dann in der Form

$$\mathbf{y} = \mathbf{x} + \alpha \mathbf{h}^{(i)}$$

mit $\mathbf{x} \in \mathbf{x}^{(0)} + H_i$ und $\alpha \in \mathbb{R}$ zerlegbar. Wenden wir das Funktional F auf \mathbf{y} an, so erhalten wir

$$\begin{aligned} F(\mathbf{y}) &= \frac{1}{2} \left[\mathbf{A} \left(\mathbf{x} + \alpha \mathbf{h}^{(i)} \right) - \mathbf{b} \right]^T \mathbf{A}^{-1} \left[\mathbf{A} \left(\mathbf{x} + \alpha \mathbf{h}^{(i)} \right) - \mathbf{b} \right] \\ &= F(\mathbf{x}) + \frac{1}{2} \alpha \mathbf{h}^{(i)T} (\mathbf{A}\mathbf{x} - \mathbf{b}) + \frac{1}{2} (\mathbf{A}\mathbf{x} - \mathbf{b})^T \mathbf{h}^{(i)} + \frac{1}{2} \alpha^2 \mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)} \\ &= F(\mathbf{x}) + \alpha \mathbf{h}^{(i)T} (\mathbf{A}\mathbf{x} - \mathbf{b}) + \frac{1}{2} \alpha^2 \mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)} \\ &= F(\mathbf{x}) + \alpha \mathbf{h}^{(i)T} \mathbf{A} \left(\mathbf{x} - \mathbf{x}^{(0)} \right) + \alpha \mathbf{h}^{(i)T} \mathbf{A} \mathbf{x}^{(0)} - \alpha \mathbf{h}^{(i)T} \mathbf{A} \mathbf{b} + \frac{1}{2} \alpha^2 \mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)}. \end{aligned}$$

Der zweite Summand in der letzten Gleichung hängt sowohl von \mathbf{x} als auch von $\mathbf{h}^{(i)}$ ab. Können wir die Richtung $\mathbf{h}^{(i)}$ so wählen, dass dieser Term für alle $\mathbf{x} \in \mathbf{x}^{(0)} + H_i$ verschwindet, so zerfällt das Minimierungsproblem von F auf $\mathbf{x}^{(0)} + H_{i+1}$ in zwei getrennte Minimierungsprobleme

$$\min_{\mathbf{x} \in \mathbf{x}^{(0)} + H_i} F(\mathbf{x})$$

und

$$\min_{\alpha \in \mathbb{R}} \left\{ \frac{1}{2} \alpha^2 \mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)} - \alpha \mathbf{h}^{(i)T} \mathbf{A} \mathbf{b} + \alpha \mathbf{h}^{(i)T} \mathbf{A} \mathbf{x}^{(0)} \right\}.$$

Die Forderung

$$\mathbf{h}^{(i)T} \mathbf{A} (\mathbf{x} - \mathbf{x}^{(0)}) = 0$$

für alle $\mathbf{x} \in \mathbf{x}^{(0)} + H_i$ ist genau dann erfüllt, wenn

$$\mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(j)} = 0 \quad \text{für } j = 0, 1, \dots, i-1$$

gilt. Zwei Vektoren $\mathbf{x} \in \mathbb{R}^n$ und $\mathbf{y} \in \mathbb{R}^n$, für die $\mathbf{x}^T \mathbf{A} \mathbf{y} = 0$ bezüglich einer symmetrischen, positiv definiten Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ gilt, heißen **A-konjugiert** oder **A-orthogonal**. Es gilt der folgende

8.62. Satz: Sind die Vektoren $\mathbf{h}^{(0)}, \dots, \mathbf{h}^{(i)}$ bezüglich der positiv definiten Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ paarweise konjugiert, so sind sie linear unabhängig.

Beweis: Es sei $\mathbf{h} = \alpha_0 \mathbf{h}^{(0)} + \dots + \alpha_i \mathbf{h}^{(i)} = \mathbf{o}$. Wir müssen zeigen, dass dann $\alpha_0 = \alpha_1 = \dots = \alpha_i = 0$ folgt. Es gilt

$$0 = \mathbf{h}^T \mathbf{A} \mathbf{h} = \left(\sum_{j=1}^i \alpha_j \mathbf{h}^{(j)} \right)^T \mathbf{A} \left(\sum_{j=1}^i \alpha_j \mathbf{h}^{(j)} \right) = \sum_{j=1}^i \sum_{k=1}^i \alpha_j \alpha_k \mathbf{h}^{(j)T} \mathbf{A} \mathbf{h}^{(k)}.$$

Wegen der A-Konjugiertheit der Vektoren entfallen alle Summanden mit $j \neq k$. Wir erhalten

$$0 = \sum_{j=1}^i \alpha_j^2 \mathbf{h}^{(j)T} \mathbf{A} \mathbf{h}^{(j)}.$$

Aus der positiven Definitheit von \mathbf{A} folgt $\mathbf{h}^{(j)T} \mathbf{A} \mathbf{h}^{(j)} > 0$ und daraus sofort

$$\alpha_0 = \alpha_1 = \dots = \alpha_i = 0.$$

✱

Betrachten wir nun wieder die Minimierung des Funktionals auf $\mathbf{x}^{(0)} + H_{i+1}$. Wurde die Richtung $\mathbf{h}^{(i)}$ so gewählt, dass $\mathbf{h}^{(j)T} \mathbf{A} \mathbf{h}^{(i)} = 0$ für $j = 0, \dots, i-1$ gilt, so folgt

$$\begin{aligned} \min_{\mathbf{y} \in \mathbf{x}^{(0)} + H_{i+1}} F(\mathbf{y}) &= \min_{\mathbf{x} \in \mathbf{x}^{(0)} + H_i, \alpha \in \mathbb{R}} F(\mathbf{x} + \alpha \mathbf{h}^{(i)}) \\ &= \min_{\alpha \in \mathbb{R}} \min_{\mathbf{x} \in \mathbf{x}^{(0)} + H_i} F(\mathbf{x} + \alpha \mathbf{h}^{(i)}) \\ &= \min_{\alpha \in \mathbb{R}} F(\mathbf{x}^{(i)} + \alpha \mathbf{h}^{(i)}) \\ &= F(\mathbf{x}^{(i)} + \alpha_i \mathbf{h}^{(i)}) \\ &= F(\mathbf{x}^{(i)}) - \frac{1}{2} \frac{\left(\mathbf{r}^{(i)T} \mathbf{h}^{(i)} \right)^2}{\mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)}} \end{aligned}$$

mit

$$\alpha_i = \frac{\mathbf{r}^{(i)T} \mathbf{h}^{(i)}}{\mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)}}.$$

Damit gilt

8.63. Satz: *Es sei $\mathbf{A} \in \mathbb{R}^{n \times n}$ eine symmetrische, positiv definite Matrix.*

$$\mathbf{h}^{(0)}, \mathbf{h}^{(1)}, \dots, \mathbf{h}^{(n-1)}, \quad \mathbf{h}^{(i)} \neq \mathbf{0}, \quad i = 0, 1, \dots, n-1$$

seien n paarweise \mathbf{A} -konjugierte Richtungen.

Für einen beliebigen Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^n$ seien die Vektoren $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}$ rekursiv durch

$$\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} + \alpha_i \mathbf{h}^{(i)}$$

mit

$$\alpha_i = \frac{\mathbf{r}^{(i)T} \mathbf{h}^{(i)}}{\mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)}} = \frac{(\mathbf{b} - \mathbf{A} \mathbf{x}^{(i)})^T \mathbf{h}^{(i)}}{\mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)}}$$

definiert. Dann gilt

1.

$$\mathbf{x}^{(i)} = \arg \min_{\mathbf{x} \in \mathbf{x}^{(0)} + H_i} F(\mathbf{x}) \quad \text{für } i = 1, \dots, n.$$

2.

$$\mathbf{A} \mathbf{x}^{(n)} = \mathbf{b}.$$

3.

$$\mathbf{h}^{(j)T} \mathbf{r}^{(k)} = 0 \quad \text{für alle } 0 \leq j < k \leq n.$$

Beweis:

1. Diese Aussage folgt sofort aus der Konstruktion der Vektoren $\mathbf{x}^{(i)}$.

2. Mit Satz 8.62 folgt $H_n = \mathbb{R}^n$. Dann gilt

$$\mathbf{x}^{(n)} = \arg \min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x})$$

$$\text{und } \mathbf{A} \mathbf{x}^{(n)} = \mathbf{b}.$$

3. Wir beweisen die Aussage mittels Induktion über k .

- Für $k = 1$ gilt

$$\begin{aligned} \mathbf{h}^{(0)T} \mathbf{r}^{(1)} &= \mathbf{h}^{(0)T} (\mathbf{b} - \mathbf{A}\mathbf{x}^{(1)}) \\ &= \mathbf{h}^{(0)T} (\mathbf{b} - \mathbf{A}\mathbf{x}^{(0)} - \alpha_0 \mathbf{A}\mathbf{h}^{(0)}) \\ &= \mathbf{h}^{(0)T} \mathbf{r}^{(0)} - \alpha_0 \mathbf{h}^{(0)T} \mathbf{A}\mathbf{h}^{(0)} \\ &= \mathbf{h}^{(0)T} \mathbf{r}^{(0)} - \frac{\mathbf{h}^{(0)T} \mathbf{r}^{(0)}}{\mathbf{h}^{(0)T} \mathbf{A}\mathbf{h}^{(0)}} \mathbf{h}^{(0)T} \mathbf{A}\mathbf{h}^{(0)} \\ &= 0. \end{aligned}$$

- Es sei für festes k $\mathbf{h}^{(j)T} \mathbf{r}^{(k)} = 0$ für $j = 0, \dots, k-1$. Wir zeigen, dass dann $\mathbf{h}^{(j)T} \mathbf{r}^{(k+1)} = 0$ für $j = 0, \dots, k$ gilt.

Wegen

$$\mathbf{r}^{(k+1)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k+1)} = \mathbf{b} - \mathbf{A}(\mathbf{x}^{(k)} + \alpha_k \mathbf{h}^{(k)}) = \mathbf{r}^{(k)} - \alpha_k \mathbf{A}\mathbf{h}^{(k)}$$

folgt

$$\mathbf{h}^{(j)T} \mathbf{r}^{(k+1)} = \mathbf{h}^{(j)T} \mathbf{r}^{(k)} - \alpha_k \mathbf{h}^{(j)T} \mathbf{A}\mathbf{h}^{(k)}.$$

Betrachten wir den Fall $j \leq k-1$. Wegen der Induktionsvoraussetzung gilt hier $\mathbf{h}^{(j)T} \mathbf{r}^{(k)} = 0$ und wegen der paarweisen \mathbf{A} -Konjugiertheit der Richtungen $\mathbf{h}^{(j)T} \mathbf{A}\mathbf{h}^{(k)} = 0$. Damit ist die Aussage bis $j = k-1$ bewiesen.

Für $j = k$ ergibt sich

$$\begin{aligned} \mathbf{h}^{(k)T} \mathbf{r}^{(k+1)} &= \mathbf{h}^{(k)T} \mathbf{r}^{(k)} - \alpha_k \mathbf{h}^{(k)T} \mathbf{A}\mathbf{h}^{(k)} \\ &= \mathbf{h}^{(k)T} \mathbf{r}^{(k)} - \frac{\mathbf{h}^{(k)T} \mathbf{r}^{(k)}}{\mathbf{h}^{(k)T} \mathbf{A}\mathbf{h}^{(k)}} \mathbf{h}^{(k)T} \mathbf{A}\mathbf{h}^{(k)} \\ &= 0. \end{aligned}$$

✱

Ist man nun in der Lage, effizient \mathbf{A} -konjugierte Richtungen $\mathbf{h}^{(0)}, \mathbf{h}^{(1)}, \dots, \mathbf{h}^{(n-1)}$ zu konstruieren, so erhält man mit dem Verfahren aus Satz 8.63 nach höchstens n Schritten die Lösung des Gleichungssystems $\mathbf{A}\mathbf{x} = \mathbf{b}$. Satz 8.63 zeigt auch, dass das Residuum $\mathbf{r}^{(i)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(i)}$ orthogonal zu allen $\mathbf{h} \in H_i$ ist. Es gilt daher entweder $\mathbf{r}^{(i)} = \mathbf{o}$ oder $\mathbf{r}^{(i)} \notin H_i$. Im ersten Falle ist $\mathbf{x}^{(i)}$ Lösung des Gleichungssystems. Im zweiten Falle bietet es sich an, für die neue Richtung $\mathbf{h}^{(i)}$ den Ansatz

$$\mathbf{h}^{(i)} = \mathbf{r}^{(i)} + \sum_{j=0}^{i-1} \beta_{ij} \mathbf{h}^{(j)} \tag{8.15}$$

zu verwenden. Dann gilt für die Residuen

$$\mathbf{r}^{(i)} = \mathbf{h}^{(i)} - \sum_{j=0}^{i-1} \beta_{ij} \mathbf{h}^{(j)} \in H_{i+1} \quad \text{für } i = 0, 1, \dots$$

Nach der dritten Aussage aus Satz 8.63 sind dann alle Residuen paarweise orthogonal. Außerdem folgt aus

$$\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} + \alpha_i \mathbf{h}^{(i)}$$

für die Residuen die Rekursionsformel

$$\mathbf{r}^{(i+1)} = \mathbf{r}^{(i)} - \alpha_i \mathbf{A} \mathbf{h}^{(i)}. \quad (8.16)$$

Dabei läßt sich für

$$\alpha_i = \frac{\mathbf{r}^{(i)T} \mathbf{h}^{(i)}}{\mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)}}$$

noch eine andere Darstellung finden. Multipliziert man Gleichung 8.15 von links mit $\mathbf{r}^{(i)T}$, so erhält man

$$\mathbf{r}^{(i)T} \mathbf{h}^{(i)} = \mathbf{r}^{(i)T} \mathbf{r}^{(i)}$$

und damit

$$\alpha_i = \frac{\mathbf{r}^{(i)T} \mathbf{r}^{(i)}}{\mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)}}.$$

Es sind zwei Fälle möglich:

1. $\alpha_i = 0$
Es folgt $\mathbf{r}^{(i)} = \mathbf{o}$. Dann ist $\mathbf{x}^{(i)}$ Lösung des Gleichungssystems $\mathbf{A} \mathbf{x} = \mathbf{b}$.
2. $\alpha_i \neq 0$
Es folgt aus Gleichung 8.16

$$-\mathbf{A} \mathbf{h}^{(i)} = \frac{\mathbf{r}^{(i+1)} - \mathbf{r}^{(i)}}{\alpha_i}.$$

Nun können wir darangehen, die β_{ij} so zu bestimmen, dass $\mathbf{h}^{(i)}$ zu allen $\mathbf{h}^{(k)}$, $k = 0, \dots, i-1$, \mathbf{A} -konjugiert ist. Wir erhalten

$$\begin{aligned} 0 &= \mathbf{h}^{(k)T} \mathbf{A} \mathbf{h}^{(i)} = \mathbf{h}^{(k)T} \mathbf{A} \mathbf{r}^{(i)} + \sum_{j=0}^{i-1} \beta_{ij} \mathbf{h}^{(k)T} \mathbf{A} \mathbf{h}^{(j)} \\ &= \mathbf{h}^{(k)T} \mathbf{A} \mathbf{r}^{(i)} + \beta_{ik} \mathbf{h}^{(k)T} \mathbf{A} \mathbf{h}^{(k)}. \end{aligned}$$

Daraus folgt

$$\beta_{ik} = -\frac{\mathbf{h}^{(k)T} \mathbf{A} \mathbf{r}^{(i)}}{\mathbf{h}^{(k)T} \mathbf{A} \mathbf{h}^{(k)}} = -\frac{\left(\mathbf{A} \mathbf{h}^{(k)}\right)^T \mathbf{r}^{(i)}}{\mathbf{h}^{(k)T} \mathbf{A} \mathbf{h}^{(k)}}.$$

Im Falle $\mathbf{r}^{(k)} \neq \mathbf{o}$ ($\alpha_k \neq 0$) ergibt sich

$$\beta_{ik} = \frac{\left(\mathbf{r}^{(k+1)} - \mathbf{r}^{(k)}\right)^T \mathbf{r}^{(i)}}{\alpha_k \mathbf{h}^{(k)T} \mathbf{A} \mathbf{h}^{(k)}} \quad \text{für } k = 0, \dots, i-1.$$

Für $k \leq i-2$ folgt dann wegen der paarweisen Orthogonalität der Residuen $\beta_{ik} = 0$.
Für $k = i-1$ erhalten wir

$$\begin{aligned} \beta_{i,i-1} &= \frac{1}{\alpha_{i-1}} \frac{\mathbf{r}^{(i)T} \mathbf{r}^{(i)}}{\mathbf{h}^{(i-1)T} \mathbf{A} \mathbf{h}^{(i-1)}} = \frac{\mathbf{h}^{(i-1)T} \mathbf{A} \mathbf{h}^{(i-1)}}{\mathbf{r}^{(i-1)T} \mathbf{r}^{(i-1)}} \frac{\mathbf{r}^{(i)T} \mathbf{r}^{(i)}}{\mathbf{h}^{(i-1)T} \mathbf{A} \mathbf{h}^{(i-1)}} \\ &= \frac{\mathbf{r}^{(i)T} \mathbf{r}^{(i)}}{\mathbf{r}^{(i-1)T} \mathbf{r}^{(i-1)}}. \end{aligned}$$

Damit haben wir das cg-Verfahren von HESTENES und STIEFEL hergeleitet.⁵

8.64. cg-Verfahren von HESTENES und STIEFEL:

Zu lösen ist das lineare Gleichungssystem $\mathbf{A} \mathbf{x} = \mathbf{b}$ mit einer symmetrischen, positiv definiten Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$.

S0 Wähle einen Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^n$.

S1 Berechne $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A} \mathbf{x}^{(0)}$.

Setze $\mathbf{h}^{(0)} = \mathbf{r}^{(0)}$ und $i = 0$.

S2 Falls $\mathbf{h}^{(i)} = \mathbf{o}$ STOPP. $\mathbf{x}^{(i)}$ ist die gesuchte Lösung.

⁵Der Name leitet sich aus dem Englischen ab. cg steht für conjugate gradients.

S3 Berechne

$$\begin{aligned}\alpha_i &= \frac{\mathbf{r}^{(i)T} \mathbf{r}^{(i)}}{\mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)}}, \\ \mathbf{x}^{(i+1)} &= \mathbf{x}^{(i)} + \alpha_i \mathbf{h}^{(i)}, \\ \mathbf{r}^{(i+1)} &= \mathbf{r}^{(i)} - \alpha_i \mathbf{A} \mathbf{h}^{(i)}, \\ \beta_i &= \frac{\mathbf{r}^{(i+1)T} \mathbf{r}^{(i+1)}}{\mathbf{r}^{(i)T} \mathbf{r}^{(i)}}, \\ \mathbf{h}^{(i+1)} &= \mathbf{r}^{(i+1)} + \beta_i \mathbf{h}^{(i)}.\end{aligned}$$

S4 Setze $i = i + 1$ und gehe zu Schritt **S2**.

Bemerkungen: (i) Im Algorithmus wird die Matrix (wie bei den anderen Iterationsverfahren) nicht verändert. Man benötigt wieder nur ein Unterprogramm, das zu einem gegebenen Vektor \mathbf{x} den Vektor $\mathbf{y} = \mathbf{A}\mathbf{x}$ berechnet.

(ii) Nach Satz 8.63 bricht das Verfahren bei exakter Rechnung nach höchstens n Schritten mit der Lösung des linearen Gleichungssystems ab.

(iii) Pro Schritt sind im Algorithmus ein Produkt Matrix \times Vektor, zwei Skalarprodukte und drei Operationen Vektor + Skalar \times Vektor zu berechnen. Für eine vollbesetzte Matrix beträgt der Aufwand maximal rund n^2 Additionen/Multiplikationen. Er wird durch die Operation Matrix \times Vektor bestimmt. Für eine schwach besetzte Matrix liegt der Aufwand für die Operation Matrix \times Vektor nur in der Größenordnung von n . Damit beträgt der Aufwand pro Schritt des Algorithmus nur Kn Additionen/Multiplikationen, wobei K wesentlich kleiner als n ist.

(iv) Der Gesamtaufwand beträgt für eine vollbesetzte Matrix rund n^3 Operationen, ist also größer als bei den bisher behandelten direkten Verfahren. Für schwach besetzte Matrizen erhalten wir jedoch einen Gesamtaufwand in der Größenordnung von n^2 . Hier ist das cg-Verfahren effektiv.

Satz 8.63 zeigt, dass für das Verhalten des cg-Verfahrens gerade die paarweise \mathbf{A} -Konjugiertheit der Richtungen $\mathbf{h}^{(i)}$ und damit die paarweise Orthogonalität der Residuen $\mathbf{r}^{(i)}$ wesentlich ist. In der Praxis werden diese Beziehungen durch den Einfluß von Rundungsfehlern oft verletzt sein. Man wird daher erwarten, dass man mehr als n Iterationen benötigt, um eine brauchbare Lösung zu berechnen. Praktische Erfahrungen zeigen aber, dass man gerade für großes n oft mit weniger Iterationsschritten auskommt. So wurde zum Beispiel 1971 von REID mit dem cg-Verfahren ein Gleichungssystem mit 4080 Gleichungen in 40 Schritten gelöst. Eine Begründung für dieses Verhalten liefert der folgende

8.65. Satz: *Es sei $\mathbf{A} \in \mathbb{R}^{n \times n}$ eine symmetrische, positiv definite Matrix mit genau k paarweise verschiedenen Eigenwerten $\lambda_1, \dots, \lambda_k$. Dann liefert das cg-Verfahren bei*

beliebigem Startvektor $\mathbf{x}^{(0)}$ nach höchstens k Schritten die Lösung des Gleichungssystems $\mathbf{A}\mathbf{x} = \mathbf{b}$.

Beweis: Zu den k Eigenwerten $\lambda_1, \dots, \lambda_k$ gehören k Eigenräume $\mathcal{U}_1, \dots, \mathcal{U}_k$. Da die Matrix \mathbf{A} symmetrisch ist, gilt

$$\dim(\mathcal{U}_1) + \dots + \dim(\mathcal{U}_k) = n$$

und

$$\mathcal{U}_i \perp \mathcal{U}_j \quad \text{für } i \neq j.$$

Damit lässt sich das Residuum $\mathbf{r}^{(0)}$ in folgender Weise zerlegen:

$$\mathbf{r}^{(0)} = p_1^{(0)} \mathbf{u}_1 + p_2^{(0)} \mathbf{u}_2 + \dots + p_k^{(0)} \mathbf{u}_k$$

mit $\mathbf{u}_i \in \mathcal{U}_i$ für $i = 1, \dots, k$. In Matrixschreibweise ergibt sich

$$\mathbf{r}^{(0)} = \mathbf{U}\mathbf{p}^{(0)}$$

mit der spaltenorthogonalen Matrix $\mathbf{U} = (\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k) \in \mathbb{R}^{n \times k}$ und $\mathbf{p}^{(0)} \in \mathbb{R}^k$. Für die Matrix \mathbf{U} gilt

$$\begin{aligned} \mathbf{A}\mathbf{U} &= \mathbf{A}(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_k) \\ &= (\mathbf{A}\mathbf{u}_1, \mathbf{A}\mathbf{u}_2, \dots, \mathbf{A}\mathbf{u}_k) \\ &= (\lambda_1 \mathbf{u}_1, \lambda_2 \mathbf{u}_2, \dots, \lambda_k \mathbf{u}_k) \\ &= \mathbf{U}\mathbf{\Lambda} \end{aligned}$$

mit $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_k) \in \mathbb{R}^{k \times k}$.

Mittels vollständiger Induktion zeigen wir nun, dass alle Vektoren $\mathbf{r}^{(i)}$ und $\mathbf{h}^{(i)}$, die vom cg-Verfahren erzeugt werden, Darstellungen der Form

$$\mathbf{r}^{(i)} = \mathbf{U}\mathbf{p}^{(i)} \quad \text{und} \quad \mathbf{h}^{(i)} = \mathbf{U}\mathbf{g}^{(i)}$$

mit $\mathbf{p}^{(i)} \in \mathbb{R}^k$ und $\mathbf{g}^{(i)} \in \mathbb{R}^k$ besitzen.

- Induktionsanfang: Für $i = 0$ gilt $\mathbf{h}^{(0)} = \mathbf{r}^{(0)} = \mathbf{U}\mathbf{p}^{(0)}$.
- Induktionsschritt: Es sei

$$\mathbf{r}^{(i)} = \mathbf{U}\mathbf{p}^{(i)} \quad \text{und} \quad \mathbf{h}^{(i)} = \mathbf{U}\mathbf{g}^{(i)}$$

mit $\mathbf{p}^{(i)} \in \mathbb{R}^k$ und $\mathbf{g}^{(i)} \in \mathbb{R}^k$. Ein Schritt des cg-Verfahrens liefert

$$\begin{aligned}\mathbf{r}^{(i+1)} &= \mathbf{r}^{(i)} + \alpha_i \mathbf{A} \mathbf{h}^{(i)} = \mathbf{U} \mathbf{p}^{(i)} + \alpha_i \mathbf{A} \mathbf{U} \mathbf{g}^{(i)} \\ &= \mathbf{U} \mathbf{p}^{(i)} + \alpha_i \mathbf{U} \Lambda \mathbf{g}^{(i)} = \mathbf{U} \left(\mathbf{p}^{(i)} + \alpha_i \Lambda \mathbf{g}^{(i)} \right) = \mathbf{U} \mathbf{p}^{(i+1)}\end{aligned}$$

mit $\mathbf{p}^{(i+1)} = \mathbf{p}^{(i)} + \alpha_i \Lambda \mathbf{g}^{(i)} \in \mathbb{R}^k$ und

$$\mathbf{h}^{(i+1)} = \mathbf{r}^{(i+1)} + \beta_i \mathbf{h}^{(i)} = \mathbf{U} \mathbf{p}^{(i)} + \beta_i \mathbf{U} \mathbf{g}^{(i)} = \mathbf{U} \left(\mathbf{p}^{(i+1)} + \beta_i \mathbf{g}^{(i)} \right) = \mathbf{U} \mathbf{g}^{(i+1)}$$

mit $\mathbf{g}^{(i+1)} = \mathbf{p}^{(i+1)} + \beta_i \mathbf{g}^{(i)} \in \mathbb{R}^k$.

Nun sind die $\mathbf{r}^{(i)}$ paarweise orthogonal, es gilt daher

$$\mathbf{r}^{(i)T} \mathbf{r}^{(j)} = 0 \quad \text{für } i \neq j.$$

Daraus folgt

$$0 = \mathbf{r}^{(i)T} \mathbf{r}^{(j)} = \left(\mathbf{U} \mathbf{p}^{(i)} \right)^T \left(\mathbf{U} \mathbf{p}^{(j)} \right) = \mathbf{p}^{(i)T} \mathbf{U}^T \mathbf{U} \mathbf{p}^{(j)} = \mathbf{p}^{(i)T} \mathbf{p}^{(j)}$$

für $i \neq j$. Die Vektoren $\mathbf{p}^{(i)}$ sind daher ebenfalls paarweise orthogonal. Wegen $\mathbf{p}^{(i)} \in \mathbb{R}^k$ gibt es nur k vom Nullvektor verschiedene Vektoren $\mathbf{p}^{(i)}$. Damit ist spätestens $\mathbf{p}^{(k)} = \mathbf{o}$ und damit $\mathbf{r}^{(k)} = \mathbf{o}$.

✱

Das modifizierte cg-Verfahren

In der ursprünglichen Form sind mit dem cg-Verfahren nur Gleichungssysteme lösbar, deren Koeffizientenmatrix symmetrisch und positiv definit ist. Das Verfahren läßt sich aber so modifizieren, dass es auf beliebige Gleichungssysteme anwendbar ist. Dazu betrachtet man das zur Lösung des linearen Gleichungssystems $\mathbf{A} \mathbf{x} = \mathbf{b}$ äquivalente Minimierungsproblem

$$x = \arg \min_{z \in \mathbb{R}^n} F(z),$$

wobei das Funktional F hier die Darstellung

$$F(z) = \frac{1}{2} (\mathbf{A}z - \mathbf{b})^T (\mathbf{A} \mathbf{A}^T)^{-1} (\mathbf{A}z - \mathbf{b})$$

besitzt. Die Matrix $\mathbf{A} \mathbf{A}^T$ ist genau dann positiv definit wenn \mathbf{A} regulär ist. Damit gilt

$$F(z) = 0 \iff \mathbf{A}z - \mathbf{b} = \mathbf{o}.$$

Eine analoge Vorgehensweise wie beim ursprünglichen cg-Verfahren liefert den folgenden Algorithmus.

8.66. Modifiziertes cg-Verfahren:

Lösen des linearen Gleichungssystem $\mathbf{Ax} = \mathbf{b}$ mit einer regulären Matrix.

S0 Wähle einen Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^n$.

S0 Berechne $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{Ax}^{(0)}$ und $\mathbf{h}^{(0)} = \mathbf{A}^T \mathbf{r}^{(0)}$

Setze $i = 0$.

S0 Falls $\mathbf{h}^{(i)} = \mathbf{0}$ STOPP. $\mathbf{x}^{(i)}$ ist die gesuchte Lösung.

S0 Berechne

$$\begin{aligned}\alpha_i &= \frac{\mathbf{r}^{(i)T} \mathbf{r}^{(i)}}{\mathbf{h}^{(i)T} \mathbf{h}^{(i)}}, \\ \mathbf{x}^{(i+1)} &= \mathbf{x}^{(i)} + \alpha_i \mathbf{h}^{(i)}, \\ \mathbf{r}^{(i+1)} &= \mathbf{r}^{(i)} - \alpha_i \mathbf{A} \mathbf{h}^{(i)}, \\ \beta_i &= \frac{\mathbf{r}^{(i+1)T} \mathbf{r}^{(i+1)}}{\mathbf{r}^{(i)T} \mathbf{r}^{(i)}}, \\ \mathbf{h}^{(i+1)} &= \mathbf{A}^T \mathbf{r}^{(i+1)} + \beta_i \mathbf{h}^{(i)}.\end{aligned}$$

S0 Setze $i = i + 1$ und gehe zu Schritt **S2**.

Beim modifizierten cg-verfahren sind die Richtungen und Residuen paarweise orthogonal:

$$\mathbf{h}^{(i)T} \mathbf{h}^{(j)} = 0 \quad \text{und} \quad \mathbf{r}^{(i)T} \mathbf{r}^{(j)} = 0 \quad \text{für} \quad i \neq j.$$

Außerdem gilt

$$\mathbf{h}^{(i)T} \mathbf{A}^{-1} \mathbf{r}^{(i)} = \mathbf{r}^{(i)T} \mathbf{r}^{(i)}.$$

Diese Aussagen folgen wie die entsprechenden Aussagen beim ursprünglichen cg-Verfahren aus der Herleitung der Rekursionsformeln. Aus der paarweisen Orthogonalität der Richtungen und Residuen folgt wieder, dass das Verfahren nach höchstens n Schritten die Lösung des linearen Gleichungssystems liefert. Analog zum Satz 8.65 ist die Anzahl der benötigten Iterationen genauer abschätzbar.

8.67. Satz: Gegeben sei ein lineares Gleichungssystem $\mathbf{Ax} = \mathbf{b}$ mit einer regulären (n, n) -Matrix \mathbf{A} , die genau k paarweise verschiedene Singulärwerte $\sigma_1, \dots, \sigma_k$ haben möge.

Dann liefert das modifizierte cg-Verfahren zu jedem Startvektor $\mathbf{x}^{(0)}$ nach höchstens k Schritten die Lösung des Gleichungssystems.

- Induktionsschritt: Wir setzen voraus, dass die Beziehungen

$$\mathbf{r}^{(i)} = \bar{\mathbf{U}}\mathbf{p}^{(i)} \quad \text{und} \quad \mathbf{h}^{(i)} = \bar{\mathbf{V}}\mathbf{g}^{(i)} \quad \text{mit} \quad \mathbf{p}^{(i)} \in \mathbb{R}^k \quad \text{und} \quad \mathbf{g}^{(i)} \in \mathbb{R}^k$$

gelten. Dann folgt

$$\begin{aligned} \mathbf{r}^{(i+1)} &= \mathbf{r}^{(i)} - \alpha_i \mathbf{A}\mathbf{h}^{(i)} = \bar{\mathbf{U}}\mathbf{p}^{(i)} - \alpha_i \mathbf{A}\bar{\mathbf{V}}\mathbf{g}^{(i)} \\ &= \bar{\mathbf{U}}\mathbf{p}^{(i)} - \alpha_i \bar{\mathbf{U}}\bar{\Sigma}\mathbf{g}^{(i)} = \bar{\mathbf{U}} \left(\mathbf{p}^{(i)} - \alpha_i \bar{\Sigma}\mathbf{g}^{(i)} \right) = \bar{\mathbf{U}}\mathbf{p}^{(i+1)} \end{aligned}$$

mit $\mathbf{p}^{(i+1)} = \mathbf{p}^{(i)} - \alpha_i \bar{\Sigma}\mathbf{g}^{(i)} \in \mathbb{R}^k$ und

$$\begin{aligned} \mathbf{h}^{(i+1)} &= \mathbf{A}^T \mathbf{r}^{(i+1)} + \beta_i \mathbf{h}^{(i)} = \mathbf{A}^T \bar{\mathbf{U}}\mathbf{p}^{(i+1)} + \beta_i \bar{\mathbf{V}}\mathbf{g}^{(i)} \\ &= \bar{\mathbf{V}}\bar{\Sigma}\mathbf{p}^{(i+1)} + \beta_i \bar{\mathbf{V}}\mathbf{g}^{(i)} = \bar{\mathbf{V}} \left(\bar{\Sigma}\mathbf{p}^{(i+1)} + \beta_i \mathbf{g}^{(i)} \right) = \bar{\mathbf{V}}\mathbf{g}^{(i+1)} \end{aligned}$$

mit $\mathbf{g}^{(i+1)} = \bar{\Sigma}\mathbf{p}^{(i+1)} + \beta_i \mathbf{g}^{(i)} \in \mathbb{R}^k$.

Aus der paarweisen Orthogonalität der Residuen folgt wieder die paarweise Orthogonalität der $\mathbf{p}^{(i)}$. Wegen $\mathbf{p}^{(i)} \in \mathbb{R}^k$ gibt es jedoch höchstens k vom Nullvektor verschiedene $\mathbf{p}^{(i)}$ die paarweise orthogonal sind. Damit ist spätestens $\mathbf{p}^{(k)} = \mathbf{o}$, $\mathbf{r}^{(k)} = \mathbf{o}$ und $\mathbf{A}\mathbf{x}^{(k)} = \mathbf{o}$. *

Vorkonditionierung der cg-Verfahren

Wir hatten schon bemerkt, dass für die Konvergenzaussagen des cg-Verfahrens die \mathbf{A} -Konjugiertheit der Richtungen wesentlich ist. Bei praktischen Rechnungen wird diese durch den Einfluß von Rundungsfehlern gestört. Andererseits wissen wir aus Satz 8.14, dass der Einfluß von Störungen um so geringer ist, je kleiner die Kondition der Matrix ist. Man kann nun versuchen, das Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ in ein solches System zu transformieren, bei dem die Koeffizientenmatrix besser konditioniert ist. Dieses Vorgehen nennt man **Vorkonditionierung**. Es sei dazu $\mathbf{L}\mathbf{L}^T = \mathbf{A} + \delta\mathbf{A}$ eine näherungsweise CHOLESKY-Zerlegung der Matrix \mathbf{A} . Wir wenden das cg-Verfahren auf das zu $\mathbf{A}\mathbf{x} = \mathbf{b}$ äquivalente System

$$\left(\mathbf{L}^{-1}\mathbf{A} \left(\mathbf{L}^{-1} \right)^T \right) \mathbf{L}^T \mathbf{x} = \mathbf{L}^{-1}\mathbf{b}$$

an. Die Matrix $\bar{\mathbf{A}} = \mathbf{L}^{-1}\mathbf{A} \left(\mathbf{L}^{-1} \right)^T = \mathbf{I} - \mathbf{L}^{-1}\delta\mathbf{A} \left(\mathbf{L}^{-1} \right)^T$ wird sich um so weniger von der Einheitsmatrix unterscheiden, je genauer die Zerlegung ist. Unter der Voraussetzung

$$\left\| \mathbf{L}^{-1}\delta\mathbf{A} \left(\mathbf{L}^{-1} \right)^T \right\| < 1$$

ergibt sich die Abschätzung

$$\text{cond}(\bar{\mathbf{A}}) \leq \frac{1 + \left\| \mathbf{L}^{-1} \delta \mathbf{A} (\mathbf{L}^{-1})^T \right\|}{1 - \left\| \mathbf{L}^{-1} \delta \mathbf{A} (\mathbf{L}^{-1})^T \right\|} = 1 + 2 \frac{\left\| \mathbf{L}^{-1} \delta \mathbf{A} (\mathbf{L}^{-1})^T \right\|}{1 - \left\| \mathbf{L}^{-1} \delta \mathbf{A} (\mathbf{L}^{-1})^T \right\|}.$$

Für $\|\delta \mathbf{A}\| \ll 1$ gilt dann $\text{cond}(\bar{\mathbf{A}}) \approx 1$. Es ist zu erwarten, dass das Rundungsfehlerverhalten und damit auch das Konvergenzverhalten des transformierten Systems günstiger sind als beim ursprünglichen Problem. Nach einigen Umformungen erhalten wir den folgenden Algorithmus.

8.68. Vorkonditioniertes cg-Verfahren:

Zu lösen ist das lineare Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit einer symmetrischen, positiv definiten Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$, von der eine näherungsweise CHOLESKY-Zerlegung $\mathbf{L}\mathbf{L}^T = \mathbf{A} + \delta \mathbf{A}$ bekannt ist.

S0 Wähle einen Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^n$.

S1 Berechne $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}$ und $\mathbf{h}^{(0)} = (\mathbf{L}\mathbf{L}^T)^{-1} \mathbf{r}^{(0)}$. Setze $i = 0$.

S2 Falls $\mathbf{h}^{(i)} = \mathbf{o}$ STOPP. $\mathbf{x}^{(i)}$ ist die gesuchte Lösung.

S3 Berechne

$$\begin{aligned} \alpha_i &= \frac{\mathbf{r}^{(i)T} (\mathbf{L}\mathbf{L}^T)^{-1} \mathbf{r}^{(i)}}{\mathbf{h}^{(i)T} \mathbf{A} \mathbf{h}^{(i)}}, \\ \mathbf{x}^{(i+1)} &= \mathbf{x}^{(i)} + \alpha_i \mathbf{h}^{(i)}, \\ \mathbf{r}^{(i+1)} &= \mathbf{r}^{(i)} - \alpha_i \mathbf{A} \mathbf{h}^{(i)}, \\ \beta_i &= \frac{\mathbf{r}^{(i+1)T} (\mathbf{L}\mathbf{L}^T)^{-1} \mathbf{r}^{(i+1)}}{\mathbf{r}^{(i)T} (\mathbf{L}\mathbf{L}^T)^{-1} \mathbf{r}^{(i)}}, \\ \mathbf{h}^{(i+1)} &= (\mathbf{L}\mathbf{L}^T)^{-1} \mathbf{r}^{(i+1)} + \beta_i \mathbf{h}^{(i)}. \end{aligned}$$

S4 Setze $i = i + 1$ und gehe zu Schritt **S2**.

Auch das modifizierte cg-Verfahren für beliebige reguläre Matrizen läßt sich vorkonditionieren. Dazu benutzt man eine näherungsweise LU-Zerlegung der Matrix.

8.69. Vorkonditioniertes, modifiziertes cg-Verfahren:

Zu lösen ist das lineare Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit einer regulären Matrix, von der eine näherungsweise LU-Zerlegung $\mathbf{P}^T \mathbf{L}\mathbf{U} = \mathbf{A} + \delta \mathbf{A}$ bekannt ist.

S0 Wähle einen Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^n$.

S1 Berechne $\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(0)}$ und $\mathbf{h}^{(0)} = (\mathbf{U}^T\mathbf{U})^{-1} \mathbf{A}^T \mathbf{P}^T (\mathbf{L}\mathbf{L}^T)^{-1} \mathbf{P}\mathbf{r}^{(0)}$.
Setze $i = 0$.

S2 Falls $\mathbf{h}^{(i)} = \mathbf{o}$ STOPP. $\mathbf{x}^{(i)}$ ist die gesuchte Lösung.

S3 Berechne

$$\alpha_i = \frac{(\mathbf{P}\mathbf{r}^{(i)})^T (\mathbf{L}\mathbf{L}^T)^{-1} (\mathbf{P}\mathbf{r}^{(i)})}{\mathbf{h}^{(i)T} (\mathbf{U}^T\mathbf{U}) \mathbf{h}^{(i)}},$$

$$\mathbf{x}^{(i+1)} = \mathbf{x}^{(i)} + \alpha_i \mathbf{h}^{(i)},$$

$$\mathbf{r}^{(i+1)} = \mathbf{r}^{(i)} - \alpha_i \mathbf{A}\mathbf{h}^{(i)},$$

$$\beta_i = \frac{(\mathbf{P}\mathbf{r}^{(i+1)})^T (\mathbf{L}\mathbf{L}^T)^{-1} (\mathbf{P}\mathbf{r}^{(i+1)})}{(\mathbf{P}\mathbf{r}^{(i)})^T (\mathbf{L}\mathbf{L}^T)^{-1} (\mathbf{P}\mathbf{r}^{(i)})},$$

$$\mathbf{h}^{(i+1)} = (\mathbf{U}^T\mathbf{U})^{-1} \mathbf{A}^T \mathbf{P}^T (\mathbf{L}\mathbf{L}^T)^{-1} \mathbf{P}\mathbf{r}^{(i+1)} + \beta_i \mathbf{h}^{(i)}.$$

S4 Setze $i = i + 1$ und gehe zu Schritt **S2**.

8.4. Aufgaben

1. Man zeige:

(a)

$$\text{lub}_1(\mathbf{A}) = \max_{1 \leq j \leq n} \left\{ \sum_{k=1}^n |a_{kj}| \right\}$$

(b)

$$\text{lub}_\infty(\mathbf{A}) = \max_{1 \leq i \leq n} \left\{ \sum_{k=1}^n |a_{ik}| \right\}.$$

2. Man beweise: Jede Grenznorm ist submultiplikativ.

3. Am Beispiel der Matrizen

$$\mathbf{A} = \begin{pmatrix} 0.6 & 0.6 \\ 0.6 & 0.6 \end{pmatrix}, \quad \mathbf{B} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix}$$

ist zu zeigen, dass $\|\mathbf{A}\|_2$ weder eine absolute noch eine monotone Norm ist.

4. Man zeige: Für reguläre Matrizen \mathbf{A} gilt:

$$\frac{1}{\text{lub}(\mathbf{A}^{-1})} = \min_{\mathbf{y} \neq \mathbf{o}} \frac{\|\mathbf{A}\mathbf{y}\|}{\|\mathbf{y}\|}.$$

5. Man zeige für eine reguläre Matrix \mathbf{A} und Vektoren $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$:

(a) Ist $\mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} \neq -1$, so gilt

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1} \mathbf{u}\mathbf{v}^T \mathbf{A}^{-1}}{1 + \mathbf{v}^T \mathbf{A}^{-1} \mathbf{u}}.$$

(b) Ist $\mathbf{v}^T \mathbf{A}^{-1} \mathbf{u} = -1$, so ist $(\mathbf{A} + \mathbf{u}\mathbf{v}^T)$ singulär.

(Hinweis: Finde einen Vektor $\mathbf{z} \neq \mathbf{o}$ mit $(\mathbf{A} + \mathbf{u}\mathbf{v}^T)\mathbf{z} = \mathbf{o}$!)

6. Es sei $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n)$ eine reguläre Matrix mit den Spalten $\mathbf{a}_i, i = 1, \dots, n$.

(a) $\hat{\mathbf{A}} = (\mathbf{a}_1, \dots, \mathbf{a}_{i-1}, \mathbf{b}, \mathbf{a}_{i+1}, \dots, \mathbf{a}_n)$ mit $\mathbf{b} \in \mathbb{R}^n$ sei die Matrix, in der gegenüber \mathbf{A} die i -te Spalte \mathbf{a}_i durch \mathbf{b} ersetzt wurde.

Man untersuche mit Hilfe der Formel aus Aufgabe 5 unter welchen Bedingungen $\hat{\mathbf{A}}^{-1}$ existiert und zeige, dass dann $\hat{\mathbf{A}}^{-1} = \mathbf{F}\mathbf{A}^{-1}$ gilt, wobei \mathbf{F} eine FROBENIUS-Matrix ist.

(b) Es sei $\mathbf{A} = (a_{ik})_{n,n}$ regulär. \mathbf{A}_α entstehe aus \mathbf{A} dadurch, dass ein einziges Element a_{ik} zu $a_{ik} + \alpha$ abgeändert wird. Für welche α existiert \mathbf{A}_α^{-1} ?

7. Es sei $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_n)$ mit $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_n > 0$ die Singulärwertzerlegung der (n, n) -Matrix \mathbf{A} .

(a) Wie lautet $\text{cond}_2(\mathbf{A})$ ausgedrückt in den σ_i ?

(b) Man gebe mit Hilfe von \mathbf{U} diejenigen Vektoren \mathbf{b} bzw. $\delta\mathbf{b}$ an, die in den Abschätzungen

i.

$$\|\delta\mathbf{x}\|_2 \leq \|\mathbf{A}^{-1}\|_2 \|\delta\mathbf{b}\|_2,$$

ii.

$$\frac{\|\delta\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \|\mathbf{A}\|_2 \|\mathbf{A}^{-1}\|_2 \frac{\|\delta\mathbf{b}\|_2}{\|\mathbf{b}\|_2} = \text{cond}_2(\mathbf{A}) \frac{\|\delta\mathbf{b}\|_2}{\|\mathbf{b}\|_2} \quad \text{und}$$

iii.

$$\|\mathbf{b}\|_2 \leq \|\mathbf{A}\|_2 \|\mathbf{x}\|_2$$

Gleichheit ergeben.

(c) Gibt es ein \mathbf{b} , so dass für alle $\delta\mathbf{b}$ in 7(b)ii. gilt:

$$\frac{\|\delta\mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \frac{\|\delta\mathbf{b}\|_2}{\|\mathbf{b}\|_2}?$$

Man bestimme solche Vektoren \mathbf{b} mit Hilfe von \mathbf{U} .

(Hinweis: Man betrachte Vektoren \mathbf{b} mit $\|\mathbf{A}^{-1}\|_2\|\mathbf{b}\|_2 = \|\mathbf{x}\|_2$.)

8. Man zeige, dass die FROBENIUS-Norm einer Matrix mit der euklidischen Vektornorm verträglich ist.

9. Es sei

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T, \quad \mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_l, 0, \dots, 0)$$

und

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0, \quad \sigma_{r+1} = \dots = \sigma_l = 0, \quad l = \min(m, n)$$

die Singulärwertzerlegung der (n, n) -Matrix \mathbf{A} . Man zeige:

(a) $\text{rg}(\mathbf{A}) = r$

(b) $\|\mathbf{A}\|_2 = \sigma_1$

(c) $\|\mathbf{A}\|_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2}$

10. Man zeige:

(a) Es seien \mathbf{L} und $\bar{\mathbf{L}}$ untere Einsdreiecksmatrizen. Dann ist das Produkt $\mathbf{L}\bar{\mathbf{L}}$ wieder untere Einsdreiecksmatrix.

(b) Mit \mathbf{L} ist auch \mathbf{L}^{-1} untere Einsdreiecksmatrix.

11. Man zeige, dass für beliebige Vektoren $\mathbf{x} \in \mathbb{R}^n$

(a)

$$\|\mathbf{x}\|_\infty \leq \|\mathbf{x}\|_2 \leq \sqrt{n}\|\mathbf{x}\|_\infty \quad \text{und}$$

(b)

$$\frac{1}{\sqrt{n}}\|\mathbf{x}\|_1 \leq \|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1.$$

gilt, und dass die auftretenden Konstanten nicht verbesserbar sind.

12. Wie groß ist der Aufwand an arithmetischen Operationen zum Lösen des Gleichungssystem

$$\mathbf{Ax} = \mathbf{b},$$

falls von der (n, n) -Matrix \mathbf{A} eine Zerlegung der Form $\mathbf{A} = \mathbf{LU}$ bekannt ist?

13. Man zeige, dass für $1 \leq k < i \leq s \leq n$

$$\mathbf{T}_{is}\mathbf{L}_k(\mathbf{l}) = \mathbf{L}_k(\mathbf{T}_{is}\mathbf{l})\mathbf{T}_{is}$$

gilt.

14. Man zeige:

$$\mathbf{L}_1(\mathbf{l}_1)\mathbf{L}_2(\mathbf{l}_2)\cdots\mathbf{L}_{n-1}(\mathbf{l}_{n-1}) = \mathbf{I} + (\mathbf{l}_1, \mathbf{l}_2, \dots, \mathbf{l}_{n-1}, 0).$$

15. Man zeige:

Wenn \mathbf{A} symmetrisch ist und $s(1) = 1$ wählbar ist, so ist auch $\mathbf{M}^{(1)}$ symmetrisch.

16. Man zeige, dass für zeilendiagonaldominante Matrizen:

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n$$

bei Durchführung des GAUSSschen Algorithmus ohne Pivotisierung gilt:

- (a) Alle Matrizen $\mathbf{M}^{(k)}$ sind zeilendiagonaldominant,
 (b) $\|\mathbf{M}^{(k)}\|_\infty \leq \|\mathbf{M}^{(k-1)}\|_\infty \leq \|\mathbf{A}\|_\infty$ für $k = 1, 2, \dots, n-1$.

(Hinweis: siehe KIELBASINSKI/SCHWETLICK Numerische lineare Algebra, S. 166.)

17. Man überlege sich, dass für die Matrix

$$\mathbf{A} = \begin{pmatrix} 1 & & & & & a \\ -1 & 1 & & & & a \\ -1 & -1 & 1 & & & a \\ \vdots & \vdots & \vdots & \ddots & & \vdots \\ -1 & -1 & -1 & \dots & 1 & a \\ -1 & -1 & -1 & \dots & -1 & a \end{pmatrix}$$

im Falle der Spaltenpivotisierung

$$\|\mathbf{A}\|_\infty = |a| + n - 1, \quad \|\mathbf{M}^{(k)}\|_\infty = 2^k |a| + n - k - 1$$

gilt, so dass der Quotient $\frac{\|\mathbf{M}^{(k)}\|_\infty}{\|\mathbf{A}\|_\infty}$ der Schranke 2^k für $|a| \rightarrow \infty$ beliebig nahe kommt.

18. Von der symmetrischen (n, n) -Matrix \mathbf{A} sei das obere Dreieck

- (a) zeilenweise,
- (b) spaltenweise

eindimensional gespeichert. Man schreibe ein Programm zum Berechnen von $\mathbf{y} = \mathbf{A}\mathbf{x}$ bei vorgegebenem $\mathbf{x} \in \mathbb{R}^n$.

19. Es sei \mathbf{A} eine reguläre Matrix mit den Zeilen $\mathbf{a}_i^T, i = 1, \dots, n$:

$$\mathbf{A} = \begin{pmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_n^T \end{pmatrix}.$$

Weiterhin sei \mathbf{D} eine Diagonalmatrix mit

$$\mathbf{D} = \text{diag} \left(\frac{\|\mathbf{A}\|_\infty}{\|\mathbf{a}_1^T\|_1}, \dots, \frac{\|\mathbf{A}\|_\infty}{\|\mathbf{a}_n^T\|_1} \right).$$

Man zeige, dass dann für die Matrix $\bar{\mathbf{A}} = \mathbf{D}\mathbf{A}$ folgendes gilt:

- (a) $\|\bar{\mathbf{A}}\|_\infty = \|\mathbf{A}\|_\infty$,
- (b) $\frac{\|\mathbf{A}^{-1}\|_\infty}{\max_i d_i} \leq \|\bar{\mathbf{A}}^{-1}\|_\infty \leq \|\mathbf{A}^{-1}\|_\infty$ und
- (c) $\frac{\text{cond}_\infty(\mathbf{A})}{\max_i d_i} \leq \text{cond}_\infty(\bar{\mathbf{A}}) \leq \text{cond}_\infty(\mathbf{A})$.

20. Es sei $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ die Singulärwertzerlegung der Matrix \mathbf{A} . Mit $\mathbf{c} = \mathbf{U}^T \mathbf{b}$ kann die Lösung von $\mathbf{A}\mathbf{x} = \mathbf{b}$ in der Form $\mathbf{x} = \mathbf{V}\mathbf{\Sigma}^{-1} \mathbf{c}$ geschrieben werden. Man zeige, dass unter der Voraussetzung

$$c_n^2 \geq \frac{0.01(c_1^2 + \dots + c_n^2)}{n}$$

in der 2-Norm die Abschätzung

$$K_b(\mathbf{A}, \mathbf{b}) = \frac{\|\mathbf{A}^{-1}\|_2 \|\mathbf{b}\|_2}{\|\mathbf{x}\|_2} \leq 10 \cdot \sqrt{n}$$

unabhängig von $\text{cond}_2(\mathbf{A})$ gilt.

21. Für das Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit

$$\mathbf{A} = \begin{pmatrix} 1.00 & 1.00 \\ 1.00 & 0.99 \end{pmatrix} \quad \mathbf{b} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

ist

$$\mathbf{A}^{-1} = \begin{pmatrix} -99 & 100 \\ 100 & -100 \end{pmatrix} \quad \mathbf{x} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

und $\text{cond}_\infty(\mathbf{A}) = 400$. Für die Störungen

(a)

$$\delta\mathbf{A} = 10^{-3} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} \quad \delta\mathbf{b} = 10^{-3} \begin{pmatrix} -1 \\ 1 \end{pmatrix}$$

(b)

$$\delta\mathbf{A} = 10^{-3} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \quad \delta\mathbf{b} = 10^{-3} \begin{pmatrix} -1 \\ -1 \end{pmatrix}$$

ist $\delta\mathbf{x}$ sowie $\|\delta\mathbf{x}\|_\infty$ zu berechnen und mit den Schranken gemäß der Vorlesung zu vergleichen.

22. Gegeben sei die positiv definite symmetrische Matrix $\bar{\mathbf{A}}$ mit der Partitionierung

$$\bar{\mathbf{A}} = \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{pmatrix}, \quad \mathbf{A} \in \mathbb{R}^{m \times m}.$$

Weiterhin sei $\bar{\mathbf{A}} = \mathbf{L}\mathbf{L}^T$ mit

$$\mathbf{L} = \begin{pmatrix} \mathbf{L}_{11} & \mathbf{O} \\ \mathbf{L}_{21} & \mathbf{L}_{22} \end{pmatrix}, \quad \mathbf{L}_{11} \in \mathbb{R}^{m \times m}.$$

Man zeige:

- (a) Die Matrix $M = C - B^T A^{-1} B$ ist positiv definit.
 (b) Es gilt $L_{22} L_{22}^T = M$.
 (c) Es gilt die Abschätzung

$$l_{ii}^2 \geq \min_{x \neq o} \left(\frac{x^T \bar{A} x}{x^T x} \right) = \frac{1}{\text{lub}_2(L^{-1})^2}, \quad i = 1, \dots, n.$$

- (d) Es gilt die Abschätzung

$$\text{lub}_2(L)^2 = \max_{x \neq o} \left(\frac{x^T \bar{A} x}{x^T x} \right) \geq l_{ii}^2, \quad i = 1, \dots, n.$$

- (e) Es gilt die Abschätzung

$$\text{cond}_2(L) \geq \max_{1 \leq i, k \leq n} \left| \frac{l_{ii}}{l_{kk}} \right|.$$

23. Gegeben seien $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ mit $\mathbf{a} \neq \mathbf{b}$ und $\|\mathbf{a}\|_2 = \|\mathbf{b}\|_2$.
 Man konstruiere eine HOUSEHOLDER-Spiegelung \mathbf{H} , für die $\mathbf{H}\mathbf{a} = \mathbf{b}$ gilt.

24. Man löse das Gleichungssystem $\mathbf{A}\mathbf{x} = \mathbf{b}$ mit

$$\mathbf{A} = \begin{pmatrix} \frac{1}{3} & -1 & \frac{5}{6} \\ \frac{2}{3} & 0 & \frac{1}{6} \\ \frac{2}{3} & \frac{1}{5} & \frac{1}{6} \end{pmatrix} \quad \mathbf{b} = \begin{pmatrix} \frac{1}{6} \\ \frac{5}{6} \\ \frac{31}{30} \end{pmatrix}$$

durch das HOUSEHOLDER-Verfahren!

25. Man zeige für eine (n, n) -Matrix \mathbf{P} , $\|\mathbf{P}\| < 1$:

$$(\mathbf{I} - \mathbf{P})^{-1} = \mathbf{I} + \sum_{i=1}^{\infty} \mathbf{P}^i$$

Hinweis: Man zeige für die Partialsummen

$$\mathbf{S}_n = \mathbf{I} + \sum_{i=1}^n \mathbf{P}^i \quad \text{gilt} \quad (\mathbf{I} - \mathbf{P})^{-1} - \mathbf{S}_n = (\mathbf{I} - \mathbf{P})^{-1} \mathbf{P}^{n+1}.$$

26. Wann konvergiert das JACOBI-Verfahren bzw. das GAUSS-SEIDEL-Verfahren für die Matrix

$$\mathbf{A} = \begin{pmatrix} \mathbf{I} & \mathbf{S} \\ \mathbf{S}^T & \mathbf{I} \end{pmatrix}, \quad \mathbf{I}, \mathbf{S} \in \mathbb{R}^{n \times n}?$$

27. Man zeige: Ist die (n, n) -Matrix A unzerlegbar, so ist die Matrix

$$M = I - B^{-1}A$$

für das JACOBI-Verfahren ebenfalls unzerlegbar.

Hinweis: B ist Diagonalmatrix.

28. Gegeben sei

$$A = \begin{pmatrix} 2 & 0 & -1 & -1 \\ 0 & 2 & -1 & -1 \\ -1 & -1 & 2 & 0 \\ -1 & -1 & 0 & 2 \end{pmatrix}.$$

Man zeige:

- (a) A ist unzerlegbar.
- (b) Das Gesamtschrittverfahren konvergiert nicht.

Index

- CHOLESKY-Zerlegung, 146
- Differentialgleichungssystem
 - schwach steifes, 32
 - steifes, 32
- Diskretisierungsfehler
 - globaler, 9, 50
 - lokaler, 8, 42
- Dreiecksmatrix
 - obere, 86
 - untere, 86
- Eigenvektor, 89
- Eigenwert, 89
- Einschrittverfahren
 - absolut stabiles, 34
- Einsdreiecksmatrix
 - obere, 86
 - untere, 86
- Einzel-schrittverfahren, 164
- EULER-Verfahren,
 - implizites, 31
- FROBENIUS-Matrix, 85
- GAUSS-SEIDEL-Verfahren, 164
- GIVENS-Drehung, 88
- GIVENS-Matrix, 88
- HOUSEHOLDER-Matrix, 89
- HOUSEHOLDER-Spiegelung, 89
- JACOBI-Verfahren, 161
- Kondition, 99
- Korrektor-Verfahren, 41
- Lösungsverfahren
 - direktes, 106
- LNT-Matrix, 85
- Matrix
 - indefinite, 144
 - irreduzible, 167
 - negativ definite, 144
 - numerisch reguläre, 93
 - numerisch singuläre, 93
 - positiv definite, 144
 - streng diagonaldominante, 143
 - unzerlegbare, 167
 - zeilenäquilibrierte, 135
- Matrixgrenznorm, 81
- Matrixnorm, 79
 - submultiplikative, 80
 - verträgliche, 80
- Mehrschrittverfahren
 - konsistentes, 42
 - konvergentes, 50
 - nullstabiles, 47
- Norm
 - absolute, 84
 - monotone, 84
- NT-Matrix, 85
- Permutationsmatrix, 88
- Pivotelement, 109
- Polynom
 - charakteristisches, 90
- Prediktor-Korrektor-Verfahren, 41
- Prediktor-Verfahren, 41
- Rücksubstitution, 110
- Skalierung
 - fiktive, 137
- Stabilitätsgebiet
 - absolutes, 34
- Tauschmatrix, 87
- Vektoren
 - konjugierte, 176
- Vektornorm, 77
- Verfahren

- diagonalimplizites, 35
- explizites, 40
- implizites, 40
- konsistentes, 8
- konvergentes, 9
- Vorkonditionierung, 186

- Zeilentausch
 - fiktiver, 114