

NUMERISCHE MATHEMATIK 1

Prof. Dr. Lars Diening

16. Januar 2012

Inhaltsverzeichnis

1	Einführung	1
1.1	Fehleranalyse	3
1.1.1	Gleitkommazahlen und Rundungsfehler	3
1.1.2	Rundung	4
1.1.3	Auslöschung	4
1.1.4	Kondition	5
1.1.5	Rundungsfehler	5
1.1.6	Konditionszahlen	6
1.1.7	Stabilität eines Verfahrens	8
1.1.8	Lineare Gleichungssysteme	9
1.1.9	Eigenwerte und Skalarprodukt	12
2	Interpolation, Approximation und Projektion	17
2.1	Polynominterpolation	17
2.1.1	Die Lagrange-Interpolationsaufgabe	18
2.1.2	Lagrangesche Basispolynome	18
2.1.3	Lagrangesche Darstellung	18
2.1.4	Interpolation von Funktionen	20
2.1.5	Spline-Interpolation	22
2.1.6	Trigonometrische Interpolation	24
3	Numerische Integration	27
3.1	Interpolatorische Quadraturformeln	27
3.1.1	Abgeschlossene Newton-Cotes-Formeln	28
3.1.2	Offene Newton-Cotes-Formeln	28
3.1.3	Gaußsche Quadraturformeln	30
4	Lineare Gleichungssysteme - direkte Verfahren	33
4.0.4	Eliminationsverfahren	33
5	Lineare Gleichungssysteme (Iterative Verfahren)	47

Kapitel 1

Einführung

Ein Hauptziel der Numerik ist die Simulation von komplexen Naturvorgängen auf Computern. Dadurch werden kostspielige Experimente wie beispielsweise Windkanalversuche oder Feuchtigkeitstests von Betonkonstruktionen durch günstige, beliebig wiederholbare ersetzt. Auch das Wetter wird auf ähnliche Weise voraus berechnet. Die verwendeten Verfahren beruhen auf einfachen Bausteinen (z.B. Integralberechnungen, Lösen linearer Gleichungssysteme, Berechnung von Nullstellen), die meist schon aus der Vor-Computer-Zeit kommen. Ziel dieses Skripts und dieser Vorlesung ist es, diese Bausteine vorzustellen. Vorallem aber gehört zur numerischen Lösung eines Problems unbedingt auch eine *Analyse des Fehlers* - hierbei unterscheiden wir die folgenden Fehlerarten:

- *Modellfehler*
 - *Idealisierungsfehler*: Zur Beschreibung des Problems wird ein mathematisches Modell gebildet. Oft müssen Vereinfachungen - etwa Linearisierungen - vorgenommen werden.
 - *Datenfehler*: Die Daten des mathematischen Modells sind oft nur ungenau bekannt. Das betrifft z.B. Koeffizienten oder Materialeigenschaften.
- *Numerische Fehler*
 - *Diskretisierungsfehler*: Kontinuierliche Prozesse werden durch endliche Prozesse ersetzt. Ein Beispiel hierfür sind die Treppenfunktionen beim Integral.
 - *Abbruchfehler*: Unendliche Algorithmen (z.B. Limesbildung) werden nach endlich vielen Schritten abgebrochen.
 - *Rundungsfehler*: Computer haben einen endlichen Zahlenbereich (etwa $\frac{1}{3} \approx 0.3333$)

Dies soll an einem Beispiel demonstriert werden: Ein Stahlseil der Länge $L = 1$ hängt an zwei Masten. Es soll die Auslenkung des Seils unter einer Last (Schwerkraft, Trapezkünstler) berechnet werden.

[Graphik-1: Stahlseil]

Die Auslenkung sei hierbei $y(x)$. Im *physikalischen Modell* gibt es zwei *Energien*:

- (a) die *potentielle Energie*.
- (b) die *Spannungsenergie*.

Zu (a): Die Belastung bei x sei $f(x)$. Damit ist die potentielle Energie $E_{pot} = - \int_0^1 f(x) y(x) dx$.

Zu (b): Nach dem *Hookeschen Gesetz* (lineare Elastizität) gilt: Spannungsenergie \propto Längenänderung. Die Änderung beträgt dann $\Delta x \cdot \left(\sqrt{1 + |y'(x)|^2} - 1 \right)$. Es folgt mit einer Materialkonstanten c :

$$\text{Spannungsenergie} = c \cdot \int_0^1 \left(\sqrt{1 + |y'(x)|^2} - 1 \right) dx = \int_0^1 \frac{|y'(x)|^2}{\sqrt{1 + |y'(x)|^2} + 1} dx$$

Für kleine Auslenkungen ($|y| \ll 1$) gilt

$$\frac{|y'(x)|^2}{\sqrt{1+|y'(x)|^2}+1} \approx \frac{1}{2}|y'(x)|$$

und mit der Taylorentwicklung, $z = |y'(x)|^2$: $\sqrt{1+z} - 1 \approx \frac{1}{2}z \pm \dots$. Damit erhalten wir das mathematische Modell:

$$J(y) = \int_0^1 \frac{1}{2}|y'(x)|^2 dx - \int_0^1 y(x) f(x) dx \rightarrow \min!$$

Nun soll die Funktion y bestimmt werden. Hierzu nützen wir die *1.Variation*. Sei $w \in C_0^\infty((0, \infty))$. Die 1. Variation berechnet sich zu

$$\begin{aligned} \delta J(y)(w) &= \frac{d}{d\varepsilon} J(y + \varepsilon \cdot w)|_{\varepsilon=0} = \int_0^1 c \cdot y'(x) w'(x) dx - \int_0^1 w(x) f(x) dx = \\ &= \int_0^1 (-cy''(x) - f(x)) w(x) dx \end{aligned}$$

Wir erhalten damit die Bedingung

$$\boxed{-cy'' = f}$$

auf $(0, 1)$ mit den Randwerten $y(0) = y(1) = 0$. Zur Lösung des Problems wird nun eine Diskretisierung vorgenommen:

$$t_i = i \cdot h, \quad i = 0, \dots, N+1, \quad h = \frac{1}{N+1} \quad (1.1)$$

$$f_i = f(t_i) \quad (1.2)$$

$$y''(t_i) = \frac{1}{h^2} (y(t_{i+1}) - 2y(t_i) + y(t_{i-1})) \quad (1.3)$$

Dies führt auf ein lineares Gleichungssystem:

$$\frac{c}{h^2} \begin{pmatrix} 2 & -1 & 0 & \dots & 0 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ \vdots & 0 & \ddots & \ddots & \vdots \end{pmatrix} \cdot \begin{pmatrix} \eta_1 \\ \vdots \\ \eta_n \end{pmatrix} = \begin{pmatrix} f_1 \\ \vdots \\ f_n \end{pmatrix}$$

bzw. in Kurzform $A\eta = b$. Dieses Gleichungssystem besitzt wegen $\det(A) \neq 0$ eine eindeutige Lösung. Diese kann wegen $D\eta = D\eta - A\eta + b$ für

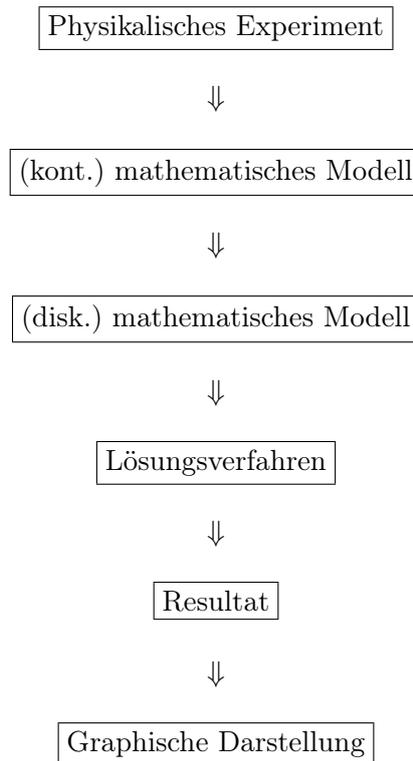
$$D = \begin{pmatrix} 2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & 2 \end{pmatrix}$$

ausgehend vom Startvektor $\eta^0 \in \mathbb{R}^N$ durch Iteration $\eta^{k+1} = \eta^k - D^{-1}[A\eta^k - b]$ angenähert werden. (Bemerkung: $Id - D^{-1}A$ ist Kontraktion.) Man kann zeigen, dass $\eta^k \rightarrow \eta$, $k \rightarrow \infty$. In der Praxis muss die Iteration jedoch abgerochen werden und es treten zusätzlich Rundungsfehler auf. Damit haben wir:

- Diskretisierungsfehler: $\max_i |\eta_i - y(t_i)|$
- Abbruchfehler: $\max_i |\eta_i^{(k)} - \eta_i|$

- Rundungsfehler: $\max_i |\tilde{\eta}_i^{(k)} - \eta_i^{(k)}|$

Zusammengefasst erhalten wir das folgende Diagramm:



1.1 Fehleranalyse

1.1.1 Gleitkommazahlen und Rundungsfehler

Bei der Abbildung der reellen Zahlen \mathbb{R} auf den endlichen Zahlbereich des Computers entstehen zwangsweise Fehler. Damit stellt sich die Frage: *Wie bildet man reelle Zahlen ab?*

Eine *normierte Gleitkommazahl* zur Basis $b \in \mathbb{N}, b \geq 2$ (meist $b = 2$ oder $b = 10$) ist eine Zahl der Form

$$x = \pm m \cdot b^{\pm e}$$

mit Mantisse $m = m_{-1}b^{-1} + m_{-2}b^{-2} + \dots + m_{-r}b^{-r}$ und Exponent $e = e_{s-1}b^{s-1} + \dots + e_1b^1 + e_0b^0$, wobei $m_i, e_i \in \{0, \dots, b-1\}$ ist. Für $x \neq 0$ ist die Darstellung durch $m_1 \neq 0$ eindeutig.

Aus technischen Gründen bietet sich auf dem Rechner das Binärsystem $b = 2$ an. Die in der obigen Form darstellbaren Zahlen heißen *Maschinenzahlen*. Sie bilden das *numerische Gleitkommagitter* $A(b, r, s)$. Da $A(b, r, s)$ endlich ist, gibt es eine größte und eine kleinste Zahl.

- (vom Betrag her) größte Zahl: $\pm (b-1)(b^{-1} + \dots + b^{-r}) \cdot b^{(b-1)(b^{s-1} + \dots + b^0)}$
- (vom Betrag her) kleinste Zahl: $\pm b^{-1} \cdot b^{-(b-1)(b^{s-1} + \dots + b^0)}$

BEISPIEL: $b = 10, r = 4$, Lichtgeschwindigkeit $= 0.2998 \cdot 10^9 \frac{m}{s}$.

IEEE-FORMAT: (double precision, REAL 8 in Fortran) Zur Darstellung haben wir 64 Bit (=8 Byte) zur Verfügung. Als Basis wird $b = 2$ gewählt. Für $x \neq 0$ beginnt die Darstellung wegen $m_1 \neq 0$ immer mit $m_1 = 1$. Daher muss man für m_1 kein Bit zur Darstellung opfern und man gewinnt ein Bit Genauigkeit. Die Zahlen mit $m = 11 \dots 1b$ und $m = 00 \dots 00b$ haben eine Sonderrolle im IEEE Format und dienen zur Darstellung von $\pm 0, \pm \infty$ (NaN - Not a Number).

1.1.2 Rundung

Bei der Darstellung von Maschinenzahlen soll so gerundet werden, dass

$$|x - \text{rd}(x)| = \min_{y \in A} |x - y|$$

wobei $x \mapsto \text{rd}(x)$ die Rundungsabbildung ist. Bei IEEE ($b = 2$) ist das die natürliche Rundung:

$$\text{rd}(\pm 0, m_1, \dots, m_{53} m_{54} \dots) = \left\{ \left(\begin{array}{l} \pm 0, m_1 \dots m_{53} \text{ falls } m_{54} = 0, \\ \text{rd}(x) = \pm 0, m_1 \dots m_{53} + 2^{-53} \text{ falls } m_{54} = 1 \end{array} \right) \right\}$$

für $x \in [x_{\text{negmin}}, x_{\text{negmax}}] \cup \{0\} \cup [x_{\text{posmin}}, x_{\text{posmax}}]$. Bei der Rundung entsteht der maximale Fehler

$$|x - \text{rd}(x)| \leq \underbrace{\frac{1}{2} b^{-r}}_{\text{letzte halbe Ziffer}} \cdot b^e$$

Das ergibt den relativen Rundungsfehler

$$\frac{|x - \text{rd}(x)|}{|x|} \leq \frac{1/2 b^{-r}}{|m| \cdot b^e} \leq \frac{1}{2} b^{1-r} = \text{Maschinengenauigkeit } \epsilon$$

Die arithmetischen Grundoperationen $\circ \in \{+, -, \cdot, /\}$ werden durch Maschinenoperationen ersetzt: $\otimes \in \{\oplus, \ominus, \odot, \oslash\}$. Es gilt

$$x \otimes y = \text{rd}(x \circ y) = (x \circ y) (1 + \epsilon) \text{ mit } |\epsilon| \leq \epsilon$$

Beachte, dass die Kommutativ- und Assoziativgesetze Im Allgemeinen *nicht* mehr gelten:

$$(x \oplus y) \oplus z \neq x \oplus (y \oplus z)$$

$$(x \oplus y) \odot z \neq (x \odot z) \oplus (y \odot z)$$

Insbesondere gilt

$$x \oplus y = x \text{ falls } |y| \leq \frac{|x|}{b} \epsilon$$

1.1.3 Auslöschung

Wir betrachten einführend ein *Beispiel* für $b = 10$ und $r = 3$: Für $x = 0.9995$ und $y = -0.9984$ ist $x + y = 0.0011$. Das Ergebnis sollte laso nach Rundung $\text{rd}(x + y) = 0.001$ sein. Addieren wir jedoch die gerundeten Werte $\text{rd}(x) = 1.000$ und $\text{rd}(y) = -0.998$ und runden anschließend wieder, so erhalten wir $\text{rd}(x) \oplus \text{rd}(y) = \text{rd}(0.002) = 0.002 = 0.200 \cdot 10^{-2}$. D.h. wir haben einen größeren relativen Fehler von $\frac{0.002}{0.0011} \approx 1.81 \dots$, obwohl wir mit 4 Stellen Genauigkeit gerechnet haben. Insofern ist Vorsicht geboten bei der Subtraktion ähnlich großer Zahlen. Durch geschickte Algorithmen kann man dieses Problem jedoch oft umgehen. Dies wollen wir am Beispiel der quadratischen Gleichung

$$x^2 + px + q = 0 \tag{1.4}$$

demonstrieren. Die Nullstellen sind gegeben durch

$$x_{1,2} = -\frac{p}{2} \pm \sqrt{\frac{p^2}{4} - q}$$

Konzentrieren wir uns auf den Fall $p > 0$ mit $|q| \ll \frac{p^2}{4}$. Dann kommt es nämlich bei x_1 zur Auslöschung, denn dann ist

$$\sqrt{\frac{p^2}{4} - q} \approx \frac{p}{2}$$

Dadurch erhält man für x_1 mit dieser Formel nur eine kleine relative Genauigkeit. Die Formel macht jedoch für x_2 keine Probleme. Aus dem Satz von *Vieta* folgt

$$x_1 x_2 = q$$

Somit können wir x_1 auch mittels der Beziehungen

$$x_2 = \frac{p}{2} + \sqrt{\frac{p^2}{4} - q}$$

$$x_1 = \frac{q}{x_2}$$

berechnen. Diese Berechnung für x_2 liefert deutlich bessere relative Fehler.

1.1.4 Kondition

Ein numerisches Problem heißt *gut konditioniert*, wenn kleine Strörungen der Daten nur kleine Änderungen der Ergebnisse zur Folge haben. (Beachte: Diese Definition wird meist auf den *relativen Fehler* angewandt, da dieser in der Numerik eine besonders wichtige Rolle spielt.)

1.1.5 Rundungsfehler

Wir betrachten nun das lineare Gleichungssystem

$$\begin{pmatrix} 1.2969 & 0.8648 \\ 0.2162 & 0.1441 \end{pmatrix} \cdot \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 0.86419999 \\ 0.14400001 \end{pmatrix} \quad (1.5)$$

mit der Kurzschreibweise $Ay = x$. Die Lösung ist gegeben durch

$$y = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} = \begin{pmatrix} 0.86419999 \\ 0.14400001 \end{pmatrix}$$

Betrachten wir jedoch das gerundete Gleichungssystem $A\hat{y} = \hat{x}$ (wobei wir x auf 4 Stellen runden):

$$\begin{pmatrix} 1.2969 & 0.8648 \\ 0.2161 & 0.1441 \end{pmatrix} \cdot \begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \end{pmatrix} = \begin{pmatrix} 0.8642 \\ 0.1440 \end{pmatrix} \quad (1.6)$$

so ergibt sich die Lösung

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_2 \end{pmatrix} = \begin{pmatrix} 2 \\ -2 \end{pmatrix}$$

Wir bemerken, dass *kleine* Änderungen von x *große* Änderungen von y zur Folge haben. Berechnen wir dies genauer, so erhalten wir

$$\Delta x = \hat{x} - x = \begin{pmatrix} 0.00000001 \\ -0.00000001 \end{pmatrix}$$

und

$$\Delta y = \hat{y} - y = \begin{pmatrix} 1.0089 \\ 1.5130 \end{pmatrix}$$

Dadurch ergeben sich die *relativen Fehler*

$$\frac{|\Delta x|}{|x|} \approx 1.1 \cdot 10^{-8} \text{ und } \frac{|\Delta y|}{|y|} \approx 0.56 \quad (1.7)$$

Der relative Fehler von x hat sich also beim Lösen des Gleichungssystems um den Faktor $5 \cdot 10^7$ verstärkt. Folglich ist das Problem $x \mapsto A^{-1}x$ *schlecht konditioniert*. Die Konditionierung linearer Gleichungssysteme werden wir später noch genauer untersuchen.

1.1.6 Konditionszahlen

Zu gegebenen Daten $x \in \mathbb{R}^m$ und Resultat $y \in \mathbb{R}^n$ seien $x + \Delta x$ und $y + \Delta y$ gestörte Daten. Dann heißen $|\Delta x|$ und $|\Delta y|$ (*absolute Fehler*) und $\frac{|\Delta x|}{|x|}$, $\frac{|\Delta y|}{|y|}$ (*relative Fehler*). Der relative Fehler ist in der Anwendung meist wichtiger als der Fehler selbst: So ist beispielsweise ein Fehler von $\pm 50\text{km}$ klein bei der Bestimmung des Abstands Mond-Erde, jedoch sehr groß beim Abstand Augsburg-München. Seien nun $f: \mathbb{R}^m \rightarrow \mathbb{R}^n$ eine glatte Funktion und sei $y = f(x)$. Für ein gestörtes Urbild $x + \Delta x$ definieren wir das gestörte Ergebnis $y + \Delta y = f(x + \Delta x)$. Wir wollen im Folgenden die Veränderung des relativen Fehlers komponentenweise untersuchen. Hierzu berechnen wir

$$\Delta y_i = f_i(x + \Delta x) - f_i(x) = \sum_{j=1}^m \partial_j f_i(x) \Delta x_j + R_i^f(x, \Delta x) \quad (1.8)$$

mit

$$\frac{|R_i^f(x, \Delta x)|}{|\Delta x|} \rightarrow 0, \quad |\Delta x| \rightarrow 0 \quad (1.9)$$

Es gilt sogar $|R_i^f(x, \Delta x)| \in O(|\Delta x|^2)$, falls $f \in C^2$ ist. Wir benutzen im Folgenden *Landausche Symbole*:

$$g(t) = O(h(t)), \quad t \rightarrow 0 \text{ falls } |g(t)| \leq c \cdot |h(t)| \text{ für ein } c > 0 \text{ und alle } |t| < t_0 \quad (1.10)$$

$$g(t) = o(h(t)), \quad t \rightarrow 0 \text{ falls } |g(t)| \leq c(t) |h(t)| \text{ mit } c(t) \rightarrow 0 \text{ für } t \rightarrow 0 \quad (1.11)$$

Sei nun $g \in C^2$. Dann gilt für einen Wert τ zwischen t und $t + \Delta t$

$$g(t + \Delta t) = g(t) + g'(t) \Delta t + g''(\tau) \frac{(\Delta t)^2}{2}$$

Dies können wir umschreiben zu

$$g(t + \Delta t) = g(t) + g'(t) \Delta t + O((\Delta t)^2)$$

und

$$\frac{g(t + \Delta t) - g(t)}{\Delta t} = g'(t) + O(\Delta t)$$

Für unsere obige Funktion f ergibt sich damit für $i = 1, \dots, n$:

$$\Delta y_i = \sum_{j=1}^m \delta_j f_i(x) \Delta x_j + O(|\Delta x|^2) \quad (1.12)$$

Dies ergibt den relativen Fehler in der i -ten Komponente:

$$\frac{\Delta y_i}{y_i} = \sum_{j=1}^m \underbrace{\left(\partial_j f_i(x) \frac{x_j}{f_i(x)} \right)}_{|\cdot|=k_{ij}} \frac{\Delta x_j}{x_j} + O\left(\frac{|\Delta x|^2}{|y_i|}\right) \quad (1.13)$$

Dabei ist der letzte Summand vernachlässigbar gegenüber dem ersten Teil, falls $|\Delta x| = o(|y_i|)$. Die k_{ij} heißen (*relative*) *Konditionszahlen* von f im Punkt x . Man nennt $f(x) = y$ *schlecht konditioniert*, falls $k_{ij} \gg 1$ und sonst *gut konditioniert*. Bei $|k_{ij}| < 1$ sprechen wir von *Fehlerdämpfung*, bei $|k_{ij}| > 1$ von *Fehlerverstärkung*. BEISPIEL. (Summe) Wie wir schon gesehen haben, ist $x \oplus y$ schlecht konditioniert, da es zur Auslöschung kommen kann. BEISPIEL. (Produkt) Sei $f(x_1, x_2) = x_1 \cdot x_2$. Dann ist

$$\partial_1 f(x_1, x_2) = x_2, \quad \partial_2 f(x_1, x_2) = x_1$$

Dies ergibt für die *relativen Konditionszahlen*

$$k_{11} = \partial_1 f(x) \frac{x_1}{f(x_1, x_2)} = 1$$

$$k_{12} = \partial_2 f(x) \frac{x_2}{f(x_1, x_2)} = 1$$

Damit ist das Problem gut konditioniert.

Die relativen Konditionszahlen sind komponentenweise definiert. D.h. k_{ij} hängt mit dem relativen Fehler von y_i bezüglich x_j zusammen. Man kann alternativ auch vektoriell rechnen und den relativen Fehler von dem gesamten Vektor y in Bezug auf den relativen Fehler von x analysieren. Dadurch erhält man

$$\frac{|\Delta y|}{|y|} = \underbrace{\left(|\nabla f| \cdot \frac{|x|}{|f(x)|} \right)}_{=k} \cdot \frac{|\Delta x|}{|x|} + O\left(\frac{|\Delta x|^2}{|y|}\right) \quad (1.14)$$

Im Falle von $f(x_1, x_2) = x_1 x_2$ ist dann

$$k = \frac{|x| \cdot |x|}{|x_1 x_2|} = \frac{x_1^2 + x_2^2}{x_1 x_2} \geq 2$$

sehr groß für $|x_1| \ll |x_2|$ und $|x_2| \ll |x_1|$. Dies führt zu falschen Einschätzungen. Dies ist z.B. auch der Fall bei der Funktion

$$g: \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mapsto \begin{pmatrix} 100 \cdot x_1 \\ x_2 \end{pmatrix}$$

Wir rechnen

$$\nabla g = \begin{pmatrix} 100 & 0 \\ 0 & 1 \end{pmatrix}, \quad |\nabla g| \approx 100, \quad k = \frac{|\nabla g| \cdot |x|}{|(100 \cdot x_1, x_2)^t|} \approx 100$$

Für die relativen Konditionszahlen haben wir jedoch $k_{11} = k_{22} = 1$ und $k_{12} = k_{21} = 0$. BEISPIEL. (Division) Sei $f(x_1, x_2) = x_1/x_2$. Dann ist

$$k_{11} = \left| \partial_1 f(x) \frac{x_1}{f(x)} \right| = 1$$

$$k_{22} = \left| \partial_2 f(x) \frac{x_2}{f(x)} \right| = 1$$

Aber wir wissen schon, dass f schlecht konditioniert ist für $x_2 \approx 0$. Wie kann das sein - mit anderen Worten - wo ist der Widerspruch? Das Problem liegt im vernachlässigten Term,

$$f''(x) O\left(\frac{|\Delta x|^2}{|y_i|}\right)$$

der für $x_2 \approx 0$ eben *nicht* vernachlässigt werden kann.

1.1.7 Stabilität eines Verfahrens

Im Folgenden wollen wir zeigen, dass verschiedene Algorithmen zur Lösung eines Problems sehr unterschiedliche *Stabilitätsverhalten* haben. Dies wollen wir am Beispiel der quadratischen Gleichung verdeutlichen. Betrachten wir hierzu

$$x^2 + px + q = 0 \quad (1.15)$$

Durch neue Variablen können wir dies in die für uns einfachere Form

$$x^2 - 2bx + c = 0$$

bringen. Dann sind die Lösungen gegeben durch

$$x_1 = b + \sqrt{b^2 - c}, \quad x_2 = b - \sqrt{b^2 - c}$$

Wir wollen im Folgenden speziell den Fall $b \gg c, b > 0$ betrachten, denn in diesem Fall ist $\sqrt{b^2 - c} \approx b$ und es wird bei x_2 möglicherweise zur Auslöschung kommen.

- 1. Algorithmus. Wir setzen

$$y_1 = b \cdot b$$

$$y_2 = y_1 - c$$

$$y_3 = \sqrt{y_2}$$

$$x_1 = y_4 = b + y_3$$

$$x_2 = y_5 = b - y_3$$

Die ersten 4 Schritte sind stabil. Beim 5. Schritt ergibt sich allerdings folgendes Problem:

$$\begin{aligned} k_{rel} &= \left| \frac{\partial y_5}{\partial y_4} \frac{y_4}{y_5} \right| = \left| -\frac{y_4}{y_5} \right| = \\ &= \left| \frac{-\sqrt{b^2 - c}}{b - \sqrt{b^2 - c}} \right| = \left| \frac{(b - \sqrt{b^2 - c}) \sqrt{b^2 - c}}{c} \right| \approx \frac{2|b|^2}{|c|} \gg 1 \end{aligned}$$

- 2. Algorithmus. Wir benutzen in diesem Algorithmus, dass $x_1 x_2 = c$ und setzen

$$y_1 = b \cdot b$$

$$y_2 = y_1 - c$$

$$y_3 = \sqrt{y_2}$$

$$x_1 = y_4 = b + y_3$$

$$x_2 = y_5 = \frac{c}{y_4}$$

$$\Rightarrow k_{rel} (3 \rightarrow 4) = \left| \frac{\partial y_4}{\partial y_3} \frac{y_3}{y_4} \right| = \left| \frac{y_3}{y_4} \right| = \left| \frac{\sqrt{b^2 - c}}{b + \sqrt{b^2 - c}} \right| \leq 1$$

$$k_{rel} (4 \rightarrow 5) = \left| \frac{\partial y_5}{\partial y_4} \frac{y_4}{y_5} \right| = \left| \frac{-c}{y_4^2} \cdot \frac{y_4}{y_5} \right| = 1$$

Dieses Problem ist also gut konditioniert.

1.1.8 Lineare Gleichungssysteme

Kommen wir nun zur Fehleranalyse des Problems $Ay = x$, d.h. $A^{-1}x = y$ mit $x \in \mathbb{R}^m, y \in \mathbb{R}^n, A \in \mathbb{R}^{m \times n}$ sowie $A^{-1} \in \mathbb{R}^{n \times n}$ (falls invertierbar). Aus der linearen Algebra benötigen wir die folgenden Begriffe:

- \mathbb{K} -Vektorraum
- Norm
- Hilbertraum
- Skalarprodukt
- Cauchy-Schwarz-Ungleichung ($|\langle x, y \rangle| \leq \|x\| \cdot \|y\|$)
- Eigenvektor
- Determinante
- charakteristisches Polynom

Die Eigenschaften einer Norm $\|\cdot\|$ auf einem \mathbb{K} -Vektorraum V sind:

$$(N1) \quad \|x\| \geq 0,$$

$$(N2) \quad \|x\| = 0 \Leftrightarrow x = 0$$

$$(N3) \quad \forall \alpha \in \mathbb{K}, x \in V: \|\alpha \cdot x\| = |\alpha| \cdot \|x\|$$

$$(N4) \quad \forall x, y \in V: \|x + y\| \leq \|x\| + \|y\|$$

BEISPIEL. Die folgenden Ausdrücke definieren Normen auf dem \mathbb{K}^n .

$$(a) \quad \|x\|_2 = \sqrt{\langle x, x \rangle} = \sqrt{\sum_{i=1}^n |x_i|^2}$$

$$(b) \quad \|x\|_1 = \sum_{i=1}^n |x_i|$$

$$(c) \quad \|x\|_\infty = \max_{i=1}^n |x_i|$$

$$(d) \quad \|x\|_p = (\sum_{i=1}^n |x_i|^p)^{1/p}$$

Ferner gilt $\|x\|_p \rightarrow \|x\|_\infty, p \rightarrow \infty$. Wir erinnern uns an die Normäquivalenz auf \mathbb{K}^n :

Theorem 1.16

Auf \mathbb{K}^n sind alle Normen äquivalent. D.h. zu zwei Normen $\|\cdot\|$ und $\|\cdot\|$ gibt es Zahlen $k, K > 0$ mit

$$k\|x\| \leq \|\cdot\| \leq K\|x\| \quad (1.17)$$

Beweis: Es genügt, die Aussage für $\|\cdot\| = \|\cdot\|_\infty$ zu zeigen. Sei nun $\|\cdot\|$ eine beliebige andere Norm und e_1, \dots, e_n eine Basis von \mathbb{K}^n . Dann gilt:

$$\|x\| = \left\| \sum_{i=1}^n x_i e_i \right\| \leq \sum_{i=1}^n |x_i| \cdot \|e_i\| \leq \|x\|_\infty \cdot \underbrace{\sum_{i=1}^n \|e_i\|}_{=C_0} = C_0 \|x\|_\infty$$

Das klärt $\|\cdot\| \leq K \|x\|_\infty$. Die Identität $Id: (\mathbb{K}^n, \|x\|_\infty) \rightarrow (\mathbb{K}^n, \|\cdot\|)$ ist stetig. Sei nun $S = \{x \in \mathbb{K}^n : \|x\|_\infty = 1\}$ die Einheitskugel bzgl. der $\|\cdot\|_\infty$ -Norm. Da S beschränkt und abgeschlossen in $(\mathbb{K}^n, \|\cdot\|_\infty)$, ist S kompakt in $(\mathbb{K}^n, \|\cdot\|_\infty)$. Daher nimmt $Id: (\mathbb{K}^n, \|x\|_\infty) \rightarrow (\mathbb{K}^n, \|\cdot\|)$ sein Minimum und Maximum auf S an. Sei daher

$$k = \min \{\|x\| : x \in S\} \quad (1.18)$$

$$K = \max \{ \|x\| : x \in S \} \quad (1.19)$$

Da $\|\cdot\|$ positiv definit ist, gilt $0 < a \leq b$ und damit

$$\forall x \in S: \quad k \leq \|x\| \leq K \quad (1.20)$$

$$\Leftrightarrow \forall x \in \mathbb{K}^n \setminus \{0\}: \quad k \leq \left\| \frac{x}{\|x\|_\infty} \right\| \leq K$$

$$\Leftrightarrow \forall x \in \mathbb{K}^n \setminus \{0\}: \quad k \|x\|_\infty \leq \|x\| \leq \|x\|_\infty$$

Da für $0 \in \mathbb{K}^n$ ohnehin $\|0\| = \|0\|_\infty = 0$ gilt, erhalten wir für alle $x \in \mathbb{K}^n$:

$$\Leftrightarrow \forall x \in \mathbb{K}^n: \quad k \|x\|_\infty \leq \|x\| \leq \|x\|_\infty \quad (1.21) \quad \blacksquare$$

Eine Norm *induziert* eine Topologie bzw. einen Konvergenzbegriff durch $x_n \rightarrow x$ genau dann, wenn $\|x_n - x\| \rightarrow 0, n \rightarrow \infty$. Die durch $\|\cdot\|_\infty$ erzeugte Konvergenz entspricht der komponentenweisen Konvergenz. Folglich gilt dies nach dem Normäquivalenzsatz für jede Norm. Da $\mathbb{K}^{n \times n} \cong \mathbb{K}^{n^2}$ gilt, kann man analog Normen für Matrizen betrachten. Besonders wichtig sind *Operatornormen* und *verträgliche Normen*.

Definition 1.22

Eine Norm $\|\cdot\|$ auf $\mathbb{K}^{n \times n}$ heißt *verträglich mit $\|\cdot\|$ auf \mathbb{K}^n , falls*

$$\forall x \in \mathbb{K}^{n \times n}, A \in \mathbb{K}^{n \times n}: \quad \|Ax\| \leq \|A\| \cdot \|x\| \quad (1.23)$$

Sie heißt Matrixnorm, falls sie zusätzlich submultiplikativ ist, d.h.

$$\forall A, B \in \mathbb{K}^{n \times n}: \quad \|AB\| \leq \|A\| \cdot \|B\| \quad (1.24)$$

BEMERKUNG: Eine Matrixnorm macht $\mathbb{K}^{n \times n}$ zu einer *Banachalgebra*. Die Abbildung $(x, y) \mapsto \langle x, y \rangle = \sum_{i=1}^n x_i \bar{y}_i$ definiert ein Skalarprodukt auf \mathbb{K}^n , wobei \bar{y}_i die komplex-konjugierte Zahl von y_i ist. Für eine Matrix $A \in \mathbb{K}^{n \times n}$ sei A^T die transponierte Matrix und $\bar{A}^T = A^*$ die adjungierte Matrix. Es gilt

$$\langle Ax, y \rangle = \langle x, A^* y \rangle$$

Theorem 1.25

Sei $\|\cdot\|$ eine Norm auf \mathbb{K}^n , so wird durch

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{\|x\| \leq 1} \|Ax\| \quad (1.26)$$

eine zu $\|\cdot\|$ *verträgliche Matrixnorm auf $\mathbb{K}^n \times n$ definiert. Diese heißt auch Operatornorm oder auch induzierte Matrixnorm. Sie ist die kleinste verträgliche Norm.*

BEISPIEL: Die *Frobeniusnorm*. Für $A, B \in \mathbb{K}^{n \times n}$ definiert

$$\langle A, B \rangle = \text{tr}(B^* A) \quad (1.27)$$

ein Skalarprodukt auf $\mathbb{K}^{n \times n}$, wobei $\text{tr}(B) = \sum_{i=1}^n B_{ii}$ die *Spur* der Matrix B ist. Daher wird durch

$$\|A\|_{\text{Frob}} = \sqrt{\langle A, A \rangle} = \sqrt{\text{tr}(A^* A)} = \left(\sum_{i,j} |a_{ij}|^2 \right)^{1/2} \quad (1.28)$$

eine Norm auf $\mathbb{K}^{n \times n}$ definiert. Damit ist $\|\cdot\|_{\text{Frob}}$ eine Hilbertraumnorm. Es gilt $\|Ax\|_2 \leq \|A\|_{\text{Frob}}\|x\|_2$, d.h. $\|\cdot\|_{\text{Frob}}$ ist verträglich mit $\|\cdot\|_2$. Weiterhin gilt

$$\|AB\|_{\text{Frob}} \leq \|A\|_{\text{Frob}}\|B\|_{\text{Frob}}$$

Damit ist $\|\cdot\|_{\text{Frob}}$ eine mit $\|\cdot\|_2$ verträgliche Matrixnorm. *Beachte aber:* $\|\cdot\|_{\text{Frob}}$ ist *nicht* die durch $\|\cdot\|_2$ induzierte Matrixnorm. BEMERKUNG: Für induzierte Matrixnormen gilt stets $\|E_n\| = 1$. Wegen $\|E_n\|_{\text{Frob}} = \sqrt{n}$ ist $\|\cdot\|_{\text{Frob}}$ keine induzierte Matrixnorm. Auch $\frac{1}{\sqrt{n}}\|\cdot\|_{\text{Frob}}$ kann dies nicht beheben, da dann im Allgemeinen $\|AB\| \leq \|A\| \cdot \|B\|$ verletzt ist, z.B. durch die Matrix

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$$

Wir bestimmen nun die durch $\|\cdot\|_1, \|\cdot\|_2, \|\cdot\|_\infty$ induzierten Matrixnormen, welche wir ebenso mit $\|\cdot\|_1, \|\cdot\|_2, \|\cdot\|_\infty$ bezeichnen.

Theorem 1.29

Es gilt

$$\|A\|_1 = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \quad (\text{maximale Spaltensumme}) \quad (1.30)$$

$$\|A\|_\infty = \max_{1 \leq k \leq n} \sum_{i=1}^n |a_{ik}| \quad (\text{maximale Zeilensumme}) \quad (1.31)$$

Beweis: Wir beschränken uns auf den Fall $\|\cdot\|_\infty$. der Fall $\|\cdot\|_1$ ist analog und wird in den Übungen behandelt. Sei also $\|z\|_\infty = \max_{1 \leq k \leq n} |z_k|$, so gilt

$$\|Ax\|_\infty = \left\| \sum_k a_{jk} x_k \right\|_\infty = \max_j \left| \sum_k a_{jk} x_k \right| \leq \max_j \sum_k |a_{jk}| \cdot \|x\|_\infty = \|A\|_\infty \|x\|_\infty$$

Hieraus folgt schon $\|A\|_\infty \leq \|A\|_\infty$. Es bleibt noch zu zeigen, dass $\|A\|_\infty \geq \|A\|_\infty$. Wir können ohne Einschränkung $A \neq 0$ annehmen. Dann gilt $\|A\|_\infty > 0$. Sei nun j_0 so, dass $\|A\|_\infty = \sum_{k=1}^n |a_{j_0,k}|$, d.h. das Maximum wird bei j_0 angenommen. Sei nun $z_k = \text{sgn}(a_{j_0,k})$, so ist $a_{j_0,k} z_k = |a_{j_0,k}|$ und $\|z\|_\infty = 1$. Es folgt

$$|(Az)_{j_0}| = \sum_k a_{j_0,k} z_k = \sum_k |a_{j_0,k}| = \|A\|_\infty$$

Dies ergibt

$$\|Az\|_\infty \geq \|A\|_\infty$$

Da $\|z\|_\infty = 1$, folgt hieraus wie gewünscht die Behauptung. ■

Zu der von $\|\cdot\|_2$ induzierten Matrixnorm kommen wir später. Da $\|\cdot\|_{\text{Frob}}$ mit $\|\cdot\|_2$ verträglich ist, wissen wir aber schon mal, dass

$$\|A\|_2 \leq \|A\|_{\text{Frob}}$$

Die Frobeniusnorm ist leicht zu berechnen, aber nur eine ungenaue obere Schranke für $\|A\|_2$.

1.1.9 Eigenwerte und Skalarprodukt

Die Eigenwerte einer Matrix A sind die Nullstellen des charakteristischen Polynoms $\chi_A(\lambda) = \det(\lambda I - A)$ (bis auf Vorzeichen). Sei $\sigma(A)$ die Menge der Eigenwerte von A (das sog. *Spektrum*). Zu jedem Eigenwert λ existieren Eigenvektoren $w \in \mathbb{K}^n \setminus \{0\}$ mit

$$Aw = \lambda w$$

Sei nun $\|\cdot\|$ eine Norm auf \mathbb{K}^n , so gilt

$$|\lambda| \cdot \|w\| = \|\lambda w\| = \|Aw\| \leq \|A\| \cdot \|w\|$$

Hieraus folgt $|\lambda| \leq \|A\|$ für jeden Eigenwert λ von A . Damit gilt

$$\rho(A) = \max_{\lambda \in \sigma(A)} |\lambda| \leq \|A\|$$

Der abgeschlossene Ball um Null in \mathbb{C} mit Radius $\rho(A)$ ist der kleinste in Null zentrierte Ball, der das Spektrum $\sigma(A)$ enthält. Aus diesem Grund heißt $\rho(A)$ der *Spektralradius* von A .

Definition 1.32

Eine Matrix $A \in \mathbb{K}^{n \times n}$ heißt hermitesch oder selbstadjungiert, falls

$$A = A^* \tag{1.33}$$

Reelle, hermitesche Matrizen heißen symmetrisch.

Theorem 1.34

Die von $\|\cdot\|_2$ induzierte Matrixnorm heißt Spektralnorm.

Theorem 1.35

Es gilt

$$\|A\|_2 = \sqrt{\max\{|\lambda| : \lambda \in \sigma(A^*A)\}} = \sqrt{\rho(A^*A)} \tag{1.36}$$

Ist A hermitesch, so gilt sogar

$$\|A\|_2 = \max\{|\lambda| : \lambda \in \sigma(A)\} = \rho(A) \tag{1.37}$$

Beweis: Sei A hermitesch. Wir wissen schon, dass $\rho(A) \leq \|A\|_2$. Es bleibt also zu zeigen, dass $\|A\|_2 \leq \rho(A)$. Aus der linearen Algebra wissen wir, dass dann A nur reelle Eigenwerte besitzt. Ferner existiert ein Orthonormalsystem aus Eigenvektoren w^1, \dots, w^n mit

$$Aw^i = \lambda_i w^i \text{ und } \langle w^i, w^j \rangle = \delta_{ij}$$

Wir zeigen nun $\|A\|_2 = \max_i |\lambda_i|$. Jedes $x \in \mathbb{K}^n$ besitzt eine Darstellung

$$x = \sum_i \alpha_i w^i$$

Hieraus folgt

$$\|x\|_2^2 = \langle x, x \rangle = \sum_{i,j} \alpha_i \bar{\alpha}_j \langle w^i, w^j \rangle = \sum_i |\alpha_i|^2$$

$$\|Ax\|_2^2 = \langle Ax, Ax \rangle = \sum_{i,j} \lambda_i \alpha_i \bar{\lambda}_j \bar{\alpha}_j \langle w^i, w^j \rangle = \sum_i |\lambda_i \alpha_i|^2$$

Damit erhalten wir die Abschätzung

$$\|Ax\|_2 \leq \left(\max_i |\lambda_i| \right) \|x\|_2 = \rho(A) \|x\|_2$$

Damit haben wir gezeigt, dass $\|A\|_2 \leq \rho(A)$. Dies ergibt die Behauptung für hermitesche A . Kommen wir nun zum allgemeinen Fall. Sei $B \in \mathbb{K}^{n \times n}$ beliebig. Dann ist $A = B^*B$ hermitesch und positiv semidefinit. Es gilt

$$\begin{aligned} \|B\|_2^2 &= \sup_{\|x\| \leq 1} \|Bx\|_2^2 = \sup_{\|x\| \leq 1} \langle Bx, Bx \rangle \leq \sup_{\|x\| \leq 1} \langle B^*Bx, x \rangle \\ &\leq \sup_{\|x\| \leq 1} \|B^*Bx\| \cdot \|x\| \leq \|B^*B\| = \rho(B^*B) \end{aligned}$$

nach dem ersten Teil. Es gilt also $\|B\|_2 \leq \rho(B^*B)$. Sei nun w_j ein Orthonormalsystem von A wie oben zu Eigenwerten λ_j . Dann gilt

$$\|Bw_j\|_2^2 = \langle Bw_j, Bw_j \rangle = \langle B^*Bw_j, w_j \rangle = \langle \lambda_j w_j, w_j \rangle = \lambda_j \|w_j\|_2^2 = |\lambda_j| \cdot \|w_j\|_2^2$$

Hierbei haben wir im letzten Schritt benutzt, dass die Eigenwerte nicht-negativ sind, da A positiv semidefinit ist. Also gilt

$$\|B\|^2 \geq \max_j |\lambda_j| = \rho(A) = \rho(B^*B)$$

Dies liefert die fehlende Abschätzung $\|B\|^2 \geq \sqrt{\rho(B^*B)}$. ■

Definition 1.38

Eine Matrix $A \in \mathbb{K}^{n \times n}$ heißt positiv semidefinit, falls

$$\forall x \in \mathbb{K}^n: \langle Ax, x \rangle \geq 0 \tag{1.39}$$

Sie heißt positiv definit, falls

$$\forall x \in \mathbb{K}^n \setminus \{0\}: \langle Ax, x \rangle > 0 \tag{1.40}$$

Theorem 1.41

Eine symmetrische Matrix $A \in \mathbb{K}^{n \times n}$ ist positiv semidefinit genau dann, wenn alle Eigenwerte nicht-negativ sind. Eine symmetrische Matrix ist positiv definit genau dann, wenn alle Eigenwerte positiv sind. In diesem Fall sind alle Diagonalelemente positiv und das betragsmäßig größte Element liegt auf der Hauptdiagonalen.

Beweis: In der linearen Algebra. ■

BEMERKUNG: Ist $A \in \mathbb{K}^{n \times n}$ positiv semidefinit, so kann man analog zu dem vorherigen Lemma zeigen, dass es unter allen betragsmäßig größten Elementen von A eines auf der Hauptdiagonalen gibt. Im Folgenden heißt positiv (semi)definit immer positiv (semi)definit und hermitesch (symmetrisch). Nun aber zurück zur Fehleranalyse von $Ax = x$ bzw. $y = A^{-1}x$ mit $A \in \mathbb{K}^{n \times n}$ invertierbar.

Theorem 1.42

Sei $B \in \mathbb{K}^{n \times n}$ mit $\|B\| < 1$, so ist $I + B$ invertierbar und

$$(I + B)^{-1} = \sum_{k \geq 0} (-1)^k B^k \tag{1.43}$$

$$\|(I + B)^{-1}\| \leq \frac{1}{1 - \|B\|} \tag{1.44}$$

Beweis: Sei $A = \sum_{k \geq 0} (-1)^k B^k$. Wegen $\|(-1)^k B^k\| \leq \|B\|^k$ konvergiert die Reihe in $\mathbb{K}^{n \times n}$ absolut. Es gilt

$$\begin{aligned} (I + B)A &= (I + B) \sum_{k \geq 0} (-1)^k B^k = \lim_{K \rightarrow \infty} \sum_{k=0}^K (I + B) (-1)^k B^k = \\ &= \lim_{K \rightarrow \infty} \sum_{k=0}^K (-1)^k B^k + \sum_{k=0}^{K+1} (-1)^k B^k = \\ &= \lim_{K \rightarrow \infty} \left((-1)^{K+1} B^{K+1} + B^0 \right) = E_n \end{aligned}$$

wobei $(-1)^{K+1} B^{K+1}$ für $K \rightarrow \infty$. Für $\|A\|$ erhalten wir die Abschätzung:

$$\|(I + B)^{-1}\| \leq \sum_{k \geq 0} \|B\|^k = \frac{1}{1 - \|B\|} \quad \blacksquare$$

Wir kommen nun zu einem Hauptresultat dieses Kapitels, dem *Störungssatz*:

Theorem 1.45

Sei $A \in \mathbb{K}^{n \times n}$ regulär und $\|\delta A\| \|A^{-1}\| = \frac{\|\delta A\|}{\|A\|} \text{cond}(A) < 1$, so ist $A + \delta A$ regulär und es gilt für $(A + \delta A)(y + \delta y) = x + \delta x$ die Abschätzung

$$\frac{\|\delta y\|}{\|y\|} \leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\|\delta A\|}{\|A\|}} \left(\frac{\|\delta x\|}{\|x\|} + \frac{\|\delta A\|}{\|A\|} \right) \quad (1.46)$$

mit $\text{cond}(A) = \|A\| \cdot \|A^{-1}\| \geq 1$.

Beweis: Es gilt

$$\|A^{-1} \delta A\| \leq \|A^{-1}\| \cdot \|\delta A\| < 1 \Rightarrow A + \delta A = A (E_n + A^{-1} \delta A) \text{ ist regulär}$$

$$\Rightarrow (A + \delta A)^{-1} = (E_n + A^{-1} \delta A)^{-1} A^{-1}$$

Aus $(A + \delta A)(y + \delta y) = x + \delta x$ folgt dann

$$(A + \delta A) \delta y = \delta x - (\delta A) y \Rightarrow \delta y = (A + \delta A)^{-1} (\delta x - \delta A y)$$

$$\Rightarrow \|\delta y\| \leq \|(A + \delta A)^{-1}\| \cdot \|\delta x - (\delta A) y\| \leq$$

$$\leq \frac{1}{1 - \|A^{-1}\| \|\delta A\|} \|A^{-1}\| (\|\delta x\| + \|\delta A\| \cdot \|y\|) \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|\delta A\|} (\|\delta x\| + \|\delta A\| \cdot \|y\|)$$

$$\Rightarrow \frac{\|\delta y\|}{\|y\|} \leq \frac{\|A^{-1}\| \cdot \|A\|}{1 - \|A^{-1}\| \cdot \|\delta A\|} \left(\frac{\|\delta x\|}{\|A\| \cdot \|y\|} + \frac{\|\delta A\|}{\|A\|} \right) \leq$$

$$\leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\|\delta A\|}{\|A\|}} \left(\frac{\|\delta x\|}{\|x\|} + \frac{\|\delta A\|}{\|A\|} \right) \quad \blacksquare$$

BEMERKUNG: $\text{cond}(A)$ hängt von $\|\cdot\|$ ab.

(a) Für $\|\cdot\|_\infty$ gilt $\text{cond}(A) = \|A\|_\infty \|A^{-1}\|_\infty$, wobei $\|A\|_\infty$ die maximale Zeilensumme von A ist.

(b) Für $\|\cdot\|_2$ und A hermitesch gilt

$$\text{cond}(A) = \frac{\lambda_{\max}}{\lambda_{\min}} \quad (1.47)$$

Ist

$$\text{cond}(A) \frac{\|\delta A\|}{\|A\|} \ll 1$$

so ist

$$\frac{\|\delta y\|}{\|y\|} \leq 2\text{cond}(A) \left(\frac{\|\delta x\|}{\|x\|} + \frac{\|\delta A\|}{\|A\|} \right)$$

mit Verstärkungsfaktor $\text{cond}(A)$. Wir haben die folgende *Regel*: sei $\text{cond}(A) \approx 10^s$ und

$$\left(\frac{\|\delta x\|}{\|x\|} + \frac{\|\delta A\|}{\|A\|} \right) \approx 10^{-k}, k > s$$

so gilt

$$\frac{\|\delta y\|}{\|y\|} \approx 10^{s-k}$$

D.h. man verliert s Stellen Genauigkeit. Hierzu ein bereits besprochenes Beispiel:

$$A = \begin{pmatrix} 1.2969 & 0.8648 \\ 0.2161 & 0.1441 \end{pmatrix}, A^{-1} = 10^8 \cdot \begin{pmatrix} 0.1441 & -0.8648 \\ 0.2161 & 12.969 \end{pmatrix} \Rightarrow \|A\|_\infty = 2.1617$$

$$\|A^{-1}\|_\infty = 1.530 \cdot 10^8 \Rightarrow \text{cond}(A) \approx 3.3 \cdot 10^8$$

Es gehen also schlimmstenfalls 8 Stellen verloren. Das Problem ist demnach sehr schlecht konditioniert.

BEMERKUNG: Die Abschätzung ist im Wesentlichen scharf. Sei hierzu λ_1 der kleinste und λ_n der größte Eigenwert. sei weiter w_1, \dots, w_n ein Orthonormalsystem von Eigenvektoren und $\delta A = 0$ sowie $x = w_n + \varepsilon w_1$. Dann folgt

$$\frac{\|\delta x\|}{\|x\|} = \frac{\varepsilon \|w_1\|}{\|w_n\|} = \varepsilon \text{ und } Ay = w_n \Rightarrow y = \frac{1}{\lambda_n} w_n$$

$$A(\delta y) = \varepsilon w_n \Rightarrow \delta y = \varepsilon \frac{w_1}{\lambda_1}$$

$$\Rightarrow \frac{\|\delta y\|}{\|y\|} = \varepsilon \frac{\lambda_n}{\lambda_1} = \frac{\|\delta x\|}{\|x\|} \frac{\lambda_n}{\lambda_1}$$

$$\Rightarrow \frac{\|\delta y\|}{\|y\|} = \text{cond}(A) \frac{\|\delta x\|}{\|x\|}$$

Kapitel 2

Interpolation, Approximation und Projektion

Die Ziele dieses Kapitels sind

- Rekonstruiere eine Funktion f , die nur an endlich vielen Stellen x_0, \dots, x_n bekannt ist.
- Finde eine einfach auszuwertende Darstellung einer analytischen Funktion.

Hierbei nützen wir die folgenden Funktionenklassen:

- Polynome P :

$$P = a_0 + \dots + a_n x^n$$

- rationale Funktionen:

$$\frac{a_0 + \dots + a_n x^n}{b_0 + \dots + b_m x^m}$$

- trigonometrische Polynome:

$$\frac{1}{2}a_0 + \sum_{k=1}^{\infty} (a_k \cos(kx) + b_k \sin(kx))$$

Wir benutzen dabei die folgende Terminologie:

- *Interpolation*: Fixiere Funktion durch Funktionswerte
- *Approximation*: Minimiere Abstand zu $g \in \mathcal{P}$, z.B.

(a) $\max_{a \leq x \leq b} |f(x) - g(x)|$

(b) $\|f - g\|_2$

(c) $\max_{i=0, \dots, n} |f(x_i) - g(x_i)|$

2.1 Polynominterpolation

Wir definieren

$$\mathcal{P}_n = \{p(x) = a_0 + \dots + a_n x^n : a_i \in \mathbb{R}, i = 0, \dots, n\}$$

Das sind die Polynome vom Grad n .

2.1.1 Die Lagrange-Interpolationsaufgabe

Wir wollen nun zu $n + 1$ verschiedenen Stützstellen (*Knoten*) $x_0, \dots, x_n \in \mathbb{R}$ und Knotenwerten $y_0, \dots, y_n \in \mathbb{R}$ ein Polynom $p \in \mathcal{P}_n$ mit $p(x_i) = y_i$, $i = 0, \dots, n$ finden.

Theorem 2.1

Die Lagrange-Interpolationsaufgabe ist eindeutig lösbar.

Beweis: Zur Eindeutigkeit: Es seien $p, \tilde{p} \in \mathcal{P}_n$ mit $p(x_i) = \tilde{p}(x_i) = y_i$. Dann gilt für $q = p - \tilde{p} \in \mathcal{P}_n$ für alle $i = 0, \dots, n$: $q(x_i) = 0$. Damit hat q mindestens $n + 1$ Nullstellen. Mit dem Satz von Rolle ergibt sich für die i -te Ableitung von q : $q^{(i)}$ hat $n - i$ Nullstellen. Insbesondere hat dann aber auch $q^{(n)}$ eine Nullstelle, also $q^{(n)} = 0$. Zur Existenz: Für festes x_0, \dots, x_n ist die Abbildung $A: \mathbb{R}^{n+1} \rightarrow \mathbb{R}^{n+1}$,

$$(a_0, \dots, a_n) \mapsto (y_0, \dots, y_n)$$

mit $y_j = a_0 + \dots + a_n x_j^n$ linear und wegen Eindeutigkeit injektiv. Somit ist A auch surjektiv und es existiert eine Lösung. ■

2.1.2 Lagrangesche Basispolynome

Wir suchen nun Polynome, die an vorgegebenen Knoten x_j den Wert 1 und sonst den Wert 0 haben. Hierzu sei

$$L_i^n(x) = \prod_{j=0, j \neq i}^n \frac{x - x_j}{x_i - x_j} \in \mathcal{P}_n, \quad i = 0, \dots, n$$

Dann gilt

$$L_i^n(x_k) = \delta_{ik}$$

mit dem Kronecker-Delta δ_{ik} . Diese sog. *Lagrangeschen Basispolynome* sind also tatsächlich linear unabhängig.

2.1.3 Lagrangesche Darstellung

Das Polynom

$$p(x) = \sum_{i=0}^n y_i L_i^n(x)$$

erfüllt $p(x_k) = y_k$ für alle $k = 0, \dots, n$. Dennoch haben wir ein *Problem*: Nehmen wir einen Punkt hinzu, so ändern sich alle Koeffizienten.

Die *Lösung* sind die *Newtonschen Basispolynome*:

$$N_0(x) = 1, \quad N_i(x) = \prod_{j=0}^{i-1} (x - x_j)$$

Sie sind linear unabhängig, da sie verschiedenen Grad haben. Wir machen den folgenden Ansatz:

$$p(x) = \sum_{i=0}^n a_i N_i(x)$$

$$\Rightarrow y_0 = p(x_0) = a_0$$

$$\Rightarrow y_1 = p(x_1) = a_0 + a_1(x_1 - x_0)$$

$$\vdots$$

$$y_n = p(x_n) = a_0 + a_1(x_1 - x_0) + \cdots + a_n(x_n - x_0) \cdots (x_n - x_{n+1})$$

Hierdurch lassen sich die a_n bestimmen. Besser (d.h. numerisch stabiler) sind allerdings die *dividierten Differenzen*:

Theorem 2.2

Das Lagrangesche Interpolationspolynom ist gegeben durch

$$p(x) = \sum_{i=1}^n y[x_0, \dots, x_n] N_i(x)$$

mit den rekursiv definierten dividierten Differenzen

$$y[x_i] = y_i, \quad y[x_i, \dots, x_{i+k}] = \frac{y[x_{i+1}, \dots, y_{i+k}] - y[x_i, \dots, x_{i+k-1}]}{x_{i+k} - x_i}$$

Beweis: Sei $p_{i,i+k}$ mit $i \leq k$ das Lagrangesche Interpolationspolynom zu $(x_i, y_i), \dots, (x_k, y_k)$. Damit gilt also $p = p_{0,n}$. Wir machen nun eine Induktion nach k :

Induktionsanfang: $k = 0$: $p_{i,i} = y_i = y[x_i]$.

Induktionsschritt: $(k-1) \rightarrow k$:

$$p_{i,i+k}(x) = p_{i,i+k-1}(x) + a(x - x_i) \cdots (x - x_{i+k-1})$$

und

$$p_{i,i+k}(x) = \frac{(x - x_i)p_{i+1,i+k-1}(x) - (x - x_{i+k})p_{i,i+k-1}(x)}{x_{i+k} - x_i}$$

Bestimme nun a durch Koeffizientenvergleich des Monoms x^n :

$$p_{i+1,i+k}(x) \underbrace{=}_{IV} y[x_{i+1}, \dots, x_{i+k}] x^{k-1} + O(x^{n-1})$$

$$p_{i,i+k-1}(x) \underbrace{=}_{IV} y[x_i, \dots, x_{i+k-1}] x^{k-1} + O(x^{n-1})$$

$$\Rightarrow a = \frac{y[x_{i+1}, \dots, y_{i+k}] - y[x_i, \dots, x_{i+k-1}]}{x_{i+k} - x_i} = y[x_i, \dots, x_{i+k}]$$

■

Folgerung 2.3

Für jede Permutation $\sigma \in S_n$ gilt:

$$y[x_0, \dots, x_n] = y[x_{\sigma(0)}, \dots, x_{\sigma(n)}]$$

Beweis: p ist unabhängig von σ . Daher sind die Koeffizienten von x^n unabhängig. Es folgt: $y[x_0, \dots, x_n]$ sind unabhängig. ■

2.1.4 Interpolation von Funktionen

Gegeben seien $f(x_i), i = 0, \dots, N$ mit x_0, \dots, x_N paarweise verschieden. Wir erhalten dazu das Interpolationspolynom

$$p(x) = \sum_{i=0}^N f[x_0, \dots, x_i](x - x_0) \dots (x - x_{i-1}) = N_i(x)$$

wobei wir jetzt $f[x_0, \dots, x_i]$ statt $y[x_0, \dots, x_i]$ schreiben.

Theorem 2.4

Sei $f \in C[a, b]$, so gilt

$$f(x) - p(x) = f[x_0, \dots, x_n, x] \underbrace{\prod_{j=0}^n (x - x_j)}_{=N_n(x)}$$

(zunächst) mit $x \in [a, b] \setminus \{x_0, \dots, x_n\}$.

Beweis: Die Polynome vom Grad $n + 1$

$$\alpha(t) = f(x), \beta(t) = p(t) + f[x_0, \dots, x_n, x]N_n(t)$$

stimmen auf den $N + 2$ Stellen x_0, \dots, x_n, x überein. Damit gilt $\alpha = \beta$, da $p(t) + f[x_0, \dots, x_n, x]N_n(t)$ gleich $f(x)$ bei $t = x$ ist. ■

Theorem 2.5

Sei $f \in C^{n+1}[a, b]$ und $x \in [a, b] \setminus \{x_0, \dots, x_n\}$, so gilt (zunächst für paarweise verschiedene x_0, \dots, x_n, x)

$$f[x_0, \dots, x_n, x] = \int_0^1 \int_0^{t_n} \dots \int_0^{t_{n-1}} f(t_0 x_0 + \dots + t_n x_n + (1 - t_0 - \dots - t_n)x) dt_n \dots dt_0$$

Beweis: Der Beweis ist eine Übungsaufgabe und kann etwa mit vollständiger Induktion geführt werden. ■

Die Formel aus Theorem 2.5 erlaubt es, $f[x_0, \dots, x_n, x]$ auch für *nicht paarweise verschiedene* x_0, \dots, x_n, x zu definieren. Insbesondere gilt

$$f[\underbrace{x_0, \dots, x_0}_{(k+1)\text{-mal}}] = \frac{f^{(k)}(x)}{k!}$$

Allgemein gilt

$$f[x_0, \dots, x_k] = \frac{f^{(k)}(\zeta)}{k!}$$

für ein $\zeta \in \text{conv}(\{x_0, \dots, x_n\})$.

Die Formel

$$f(x) = p(x) + f[x_0, \dots, x_n, x]N_n(x) = \sum_{i=0}^n f[x_0, \dots, x_i]N_i(x) + f[x_0, \dots, x_n, x]N_n(x)$$

bleibt aus Stetigkeitsgründen korrekt.

Folgerung 2.6

Sei $f \in C^{n+1}[a, b]$, so gilt

$$f(x) - p(x) = \frac{f^{(n+1)}(\zeta)}{(n+1)!} N_n(x)$$

für ein $\zeta \in \text{conv}(x_0, \dots, x_n)$.

Beweis: Die Behauptung folgt sofort aus dem Mittelwertsatz der Integralrechnung und

$$\int_0^1 \int_0^{t_n} \dots \int_0^{t_{n-1}} 1 dt_n \dots dt_0 = \frac{1}{(n+1)!}$$

Also gilt für den Interpolationsfehler auf $[a, b]$:

$$|f(x) - p(x)| \leq \frac{1}{(n+1)!} \max_{\zeta \in [a, b]} |f^{(n+1)}(\zeta)|$$

Der erste Faktor auf der rechten Seite deutet auf eine gute Fehlerkontrolle hin, doch haben wir keine Garantie, dass der zweite Faktor nicht zu groß wird. Hierzu ein

BEISPIEL: Sei

$$f(x) = \frac{1}{1+x^2} \in C^\infty(\mathbb{R})$$

Dann ist

$$f'(x) = \frac{2x}{(1+x^2)^2}$$

$$\Rightarrow |f^{(n)}(x)| \approx 2^n n! O(|x|^{-2-n})$$

für große x . Zwar kürzt sich der Faktor $n!$ bei der Abschätzung

$$|f(x) - p(x)| \leq \frac{1}{(n+1)!} \max_{\zeta \in [a, b]} |f^{(n+1)}(\zeta)|$$

heraus, doch der Faktor 2^n bleibt. Deshalb haben wir hier keine gleichmäßige Konvergenz (am Rand schlechter).

BEMERKUNG: Der *Weierstraßsche Approximationssatz* besagt, dass jedes $f \in C[a, b]$ gleichmäßig durch Polynome approximiert werden kann. Die Vermutung, dass dies Lagrangesche Interpolationspolynome sind, ist jedoch *falsch*.

BEISPIELE:

- $f(x) = |x|$, $x \in [-1, 1]$. Hier bleibt der Fehler an den Rändern groß.
- $\exp(x)$, $\cos(x)$, $\sin(x)$ sind in Ordnung, da die Ableitung durch eine von n unabhängige Konstante K abgeschätzt werden kann:

$$\max_{a \leq \zeta \leq b} \|f^{(n)}(\zeta)\| \leq K$$

BEMERKUNG: *Richardsonsche Extrapolation*. Hier wollen wir einen Grenzwert $a(h)$, $h \rightarrow 0$ berechnen, bei dem $a(0)$ nicht direkt berechenbar ist, aber die Entwicklung

$$a(x) = a_0 + \sum_{j=1}^n a_j x^j + a_{n+1}(h) h^{n+1} + O(h^{n+1})$$

Sei dann p_n das Interpolationspolynom für die Stützstellen $\rho^k, \dots, \rho^{k+n}$ mit $\rho < 1$. Dann gilt (ohne Beweis)

$$a(0) - p_n(0) = O(h^{k(n+1)})$$

2.1.5 Spline-Interpolation

Lagrangesche Interpolationspolynome sind insbesondere schlecht für nicht glatte Funktionen. Dies liegt am geforderten C^∞ -Übergang an den Knoten. Bei *Splines* fordert man weniger. Wir betrachten die Vektorräume

$$S_n^{(k,r)}([a,b]) = \{p \in C^r([a,b]) : p|_{I_i} \in \mathcal{P}_k(I_i)\}$$

wobei wir wieder von den Stützstellen $a = x_0, \dots, x_n$ ausgehen mit $I_i = [x_i - 10x_i]$, $i = 1, \dots, n$. Die Funktionen sind also lokal Polynome vom Grad k und global r -mal stetig differenzierbar.

BEISPIEL: Kubischer Spline.

Definition 2.7

$s_n \in S_n^{(3,2)}$ heißt kubischer Spline der Zerlegung $a = x_0 < \dots < x_n = b$. falls $s''(a) = s''(b) = 0$, so heißt s natürlich.

Theorem 2.8

Zu (x_i, y_i) mit paarweise verschiedenen x_0, \dots, x_n existieren Splines s_n mit

$$s_n(x_i) = y_i$$

Unter zusätzlicher Vorgabe von $s_n''(a), s_n''(b)$ ist s_n eindeutig.

Beweis: Die Abbildung

$$\text{Koeff}(s_n) \mapsto (y_0, \dots, y_n)$$

ist linear. Wir zeigen nun, dass die Abbildung Koeff auch injektiv ist. Seien hierzu $s_n^{(1)}, s_n^{(2)}$ Splines und $s = s_n^{(1)} - s_n^{(2)}$. Dann ist $s \in S_n^{(3,2)}$ und $s(x_i) = y_i$ sowie $s''(a) = s''(b) = 0$. Für $w \in C^2[a,b]$ mit $w(x_i) = 0$ gilt:

$$\begin{aligned} \int_a^b s'' w'' dx &= \sum_i \int_{x_i}^{x_{i+1}} s'' w'' dx = \\ &= \sum_i \left(- \int_{x_i}^{x_{i+1}} s^{(3)} w' dx + [s'' w']_{x_i}^{x_{i+1}} \right) = \\ &= \sum_i \left(\int_{x_i}^{x_{i+1}} s^{(3)} w' dx \right) + [s'' w']_{x_0}^{x_n} = \\ &= \sum_i \left(\int_{x_i}^{x_{i+1}} s^{(4)} w dx - [s^{(3)} w]_{x_i}^{x_{i+1}} \right) = 0 \end{aligned}$$

wobei der letzte Schritt analog zum Vorherigen durchgeführt wird. Sei nun $w = s$. Dann erhalten wir

$$\int_a^b |s''|^2 dx = 0 \Rightarrow s'' = 0$$

Somit ist s linear und da s an den Punkten a und b gleich 0 ist, folgt $s \equiv 0$. Damit haben wir Koeff injektiv und den Spline eindeutig. Ferner ist Koeff auch linear und endlichdimensional, damit auch surjektiv. Es folgt die Lösbarkeit. ■

Theorem 2.9

Die interpolierenden natürlichen Splines minimieren das Integral

$$\mathcal{J}(f) = \int_0^\infty |f''(x)|^2 dx$$

bezüglich $f \in C^2[a, b]$, $f(x_i) = y_i$, $i = 0, \dots, n$.

Beweis: Sei $N = \{w \in C^2[a, b] : w(x_i) = 0, i = 0, \dots, n\}$. Dann gilt nach obigem Beweis:

$$\int_a^b s'' w'' dx = 0$$

d.h. $(\delta \mathcal{J})(s_n)(w) = 0$. Somit folgt

$$\mathcal{J}(f) = \mathcal{J}(s_n + \tilde{w}) \text{ für ein } \tilde{w} \in N$$

$$= \mathcal{J}(s_n) + 2 \int_a^b s'' \tilde{w}'' dx + \mathcal{J}(\tilde{w}) \geq \mathcal{J}(s_n)$$

EXPLIZITE BERECHNUNG: Auf den Intervallen $[x_i, x_{i+1}]$ gilt

$$s_n|_{[x_{i-1}, x_i]} = p_i(x) = a_1^{(i)}(x - x_i) + a_2^{(i)}(x - x_i)^2 + a_3^{(i)}(x - x_i)^3$$

Wir erhalten somit

(a)

$$p_j(x_j) = y_j, \quad p_i(x_{i-1}) = y_{i-1} \Rightarrow$$

$$a_0^{(i)} = y_i$$

$$y_{i-1} - y_i = -a_1^{(i)} h_i + a_2^{(i)} h_i^2 - a_3^{(i)} h_i^3$$

(b) Die Randbedingung $p_1''(x_0) = p_n''(x_n) = 0$ ergibt

$$a_2^{(1)} - 3a_3^{(1)} h_1 = 0, \quad a_2^{(1)} = 0$$

(c) Die C^1 -Forderung $p_i'(x_i) = p_{i+1}'(x_i)$ ergibt

$$a_1^{(i)} = a_1^{(i+1)} - 2a_2^{(i+1)} h_{i+1} + a_3^{(i+1)} h_{i+1}^2, \quad i = 1, \dots, n-1$$

(d) Die C^2 -Forderung $p_i''(x_i) = p_{i+1}''(x_i)$ ergibt

$$a_2^{(i)} = a_2^{(i+1)} - 3a_3^{(i+1)} h_{i+1}$$

Wir erhalten insgesamt $4n$ Gleichungen. Drücke nun $a_3^{(i)}$ durch $a_2^{(i)}$ aus sowie $a_1^{(i)}$ durch $a_2^{(i)}$:

$$a_3^{(i)} = \frac{a_2^{(i)} - a_2^{(i+1)}}{3h_i}, \quad i = 1, \dots, n$$

Theorem 2.11

Zu $y_0, \dots, y_n \in \mathbb{C}$ existiert genau ein $\varphi \in \tau_{\mathbb{C}}^n$ mit

$$\varphi(x_j) = y_j, \quad j = 0, \dots, n$$

wobei $x_j = j \frac{2\pi}{N+1}$. Es gilt $\varphi(x) = \sum_{j=0}^n e^{ijx}$ mit

$$c_k = \frac{1}{n+1} \sum_{j=0}^n e^{-ijx_k} y_j, \quad k = 0, \dots, n$$

Beweis: Wir müssen das Gleichungssystem

$$c_n \omega_j^n + \dots + c_1 \omega_j^1 + c_0 = y_j, \quad j = 0, \dots, n$$

mit $\omega_j = e^{ix_j} = e^{\frac{j2\pi \cdot i}{n+1}}$ lösen. In Matrixform heißt das

$$\begin{pmatrix} 1 & \omega_0 & \dots & \omega_0^n \\ \vdots & & \vdots & \\ 1 & \omega_n & \dots & \omega_n^n \end{pmatrix} \cdot \begin{pmatrix} c_0 \\ \vdots \\ c_n \end{pmatrix} = \begin{pmatrix} y_0 \\ \vdots \\ y_n \end{pmatrix}$$

Die Matrix ist eine Vandermonde-Matrix V und wir schreiben kurz $Vc = y$. Da die ω_j paarweise verschieden sind, ist V auch invertierbar. (Wir erinnern uns, dass für die Vandermonde-Matrix gilt: $\det(v) = \prod_{0 \leq i, j \leq n} (\omega_i - \omega_j)$) Die Existenz und Eindeutigkeit ist damit gezeigt. Für die ω_j gilt $|\omega_j| = 1$, $\omega_j = \omega_1^j$, $\overline{\omega_j} = \omega_j^{-1}$. Unter Verwendung von $\sum_{j=0}^n \omega_1^k = 0$ erhält man

$$V^*V = (n+1)I$$

Die Matrix $(n+1)^{-1}V^*$ ist also die Inverse zu V (Übung) und wir erhalten aus $Vc = y$ die Beziehung

$$c = \frac{1}{n+1} V^* y$$

mit $(V^*)_{jk} = \overline{V_{kj}} = \overline{\omega_k^j} = e^{-ijx_k}$, woraus sich die Koeffizienten c_k wie gewünscht ergeben. ■

Kapitel 3

Numerische Integration

In diesem Kapitel widmen wir uns der Berechnung von Integralen $\int_a^b f(x) dx$. Viele Funktionen f besitzen keine explizit darstellbare Stammfunktion. In der Praxis müssen die meisten Integrale daher durch endliche Summen approximiert werden. Wir kennen bereits die Approximation durch Rechtecke mit Obersummen bzw. Untersummen sowie die Rechteckregel:

$$\int_a^b f(x) dx \approx \sum_{i=0}^{n-1} (x_{i+1} - x_i) f(x_i) \quad (3.1)$$

Allerdings gibt es bessere Verfahren, Integrale zu approximieren. Diese wollen wir nun besprechen.

3.1 Interpolatorische Quadraturformeln

Wir approximieren die Funktion f durch ein Interpolationspolynom

$$p_n(x) = \sum_{i=0}^n f(x_i) L_i^n(x) \quad (3.2)$$

und erhalten

$$I^{(n)}(f) = \int_a^b p_n(x) dx = \sum_{i=0}^n f(x_i) \underbrace{\int_a^b L_i^n(x) dx}_{=\alpha_i} \quad (3.3)$$

Hierbei notiert $L_i^{(n)}$ das *Lagrangesche Interpolationspolynom*. Insbesondere hängen die α_i nur von den x_0, \dots, x_n ab. Wir präzisieren dies nun:

Theorem 3.4

Sei

$$I(f) - I^{(n)}(f) = \int_a^b f[x_0, \dots, x_n, x] \prod_{j=0}^n (x - x_j) dx \quad (3.5)$$

Dann gilt die folgende Abschätzung:

$$|I(f) - I^{(n)}(f)| \leq \frac{\|f^{(n+1)}\|_\infty}{(n+1)!} \int_a^b \left| \prod_{j=0}^n (x - x_j) \right| dx \quad (3.6)$$

BEMERKUNG: Für ein Polynom $f \in \mathcal{P}^n$ ist die Interpolation sogar exakt.

Definition 3.7

Eine Quadraturformel heißt (mindestens) von der Ordnung n , falls sie alle Polynome $f \in \mathcal{P}^n$ exakt integriert.

Die $I^{(n)}$ sind also von der Ordnung $n + 1$.

3.1.1 Abgeschlossene Newton-Cotes-Formeln

Wir wählen nun äquidistante Stützstellen $a = x_0, \dots, x_n = b$, d.h.

$$x_i = a + i \frac{b-a}{n}, \quad i = 0, \dots, n$$

Mit der Variablentransformation

$$t = \frac{x-a}{b-a} \cdot n$$

erhalten wir mit $t \in [0, n]$ und der Schrittweite $H = \frac{b-a}{n}$

$$x = a + t \cdot H$$

Die Polynome $L_i^{(n)}$ lauten damit

$$L_i^{(n)}(x) = \prod_{j=0, j \neq i}^n \frac{x - x_j}{x_i - x_j} = \prod_{j=0, j \neq i}^n \frac{a + tH - a - jH}{a + iH - a - jH} = \prod_{j=0, j \neq i}^n \frac{t - j}{i - j} \quad (3.8)$$

Damit erhalten wir

$$\alpha_i = \int_a^b L_i^{(n)}(x) dx = H \cdot \int_0^n \prod_{j=0, j \neq i}^n \frac{t - j}{i - j} dt \quad (3.9)$$

Für $n = 1, 2, 3$ erhalten wir so die wichtigen folgenden Regeln:

NAME	FORMEL
Mittelpunktregel	$I^{(1)}(f) = \frac{b-a}{2} (f(a) + f(b))$
Simpsonregel	$I^{(2)}(f) = \frac{b-a}{6} (f(a) + 4f(\frac{a+b}{2}) + f(b))$
$\frac{3}{8}$ -Regel	$I^{(3)}(f) = \frac{b-a}{8} (f(a) + 3f(a+H) + 3f(b-H) + f(b))$

3.1.2 Offene Newton-Cotes-Formeln

Im Folgenden sei $x_0 = a + H$, $x_i = x_0 + iH$ und $b = x_n + H$. Wir betrachten also $n + 2$ Intervalle der Länge H , wobei a und b keine Stützstellen sind.

$$\begin{aligned} \hat{I}^{(0)}(f) &= (b-a) f\left(\frac{a+b}{2}\right) \text{ (Mittelpunktsregel)} \\ \hat{I}^{(1)}(f) &= \frac{b-a}{2} (f(a+H) + f(b-H)) \\ \hat{I}^{(2)} &= \frac{b-a}{3} (2f(a+H) - f\left(\frac{a+b}{2}\right) + f(b-H)) \end{aligned}$$

Wir wollen im Folgenden die *Trapezregel*, *Simpsonregel* und die *Mittelpunktsregel* untersuchen.

Theorem 3.10

Für die Restglieder der Trapezregel, Simpsonregel und der Mittelpunktsregel gelten die folgenden Beziehungen:

- *Mittelpunktsregel:*

$$I(f) - (b-a) f\left(\frac{a+b}{2}\right) = \frac{(b-a)^3}{24} f''(\zeta) \quad \text{für } f \in C^2[a, b] \quad (3.11)$$

- *Trapezregel:*

$$I(f) - \frac{a+b}{2} (f(a) + f(b)) = -\frac{(b-a)^3}{12} f''(\zeta) \quad \text{für } f \in C^2[a, b] \quad (3.12)$$

- *Simpsonregel:*

$$I(f) - \frac{b-a}{6} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right) = -\frac{(b-a)^5}{2880} f^{(4)}(\zeta) \quad \text{für } f \in C^4[a, b] \quad (3.13)$$

für gewisse Zwischenstellen $\zeta \in [a, b]$.

BEWEIS: Für die *Trapezregel* erhalten wir beispielsweise durch Anwendung der Zwischenwertsätze

$$I(f) - I^{(1)}(f) = \int_a^b f[a, b, x] (x-a)(x-b) dx =$$

$$f[a, b\zeta_1] \int_a^b (x-a)(x-b) dx =$$

$$\frac{1}{2} f''(\zeta_2) \cdot \frac{-(b-a)^3}{12}$$

Dabei notieren $f[a, b, x]$ die dividierten Differenzen. In der letzten Zeilen verwendeten wir den Zwischensatz der Differentialrechnung, hingegen in der vorherigen den Zwischenwertsatz der Integralrechnung:

Sei f stetig, g integrierbar und $g \geq 0$ (bzw. $g \leq 0$). Dann existiert ein $\zeta \in [a, b]$ mit

$$\int_a^b f(x) g(x) dx = f(\zeta) \int_a^b g(x) dx$$

Man beachte, dass hier die Voraussetzungen erfüllt sind, da $(x-a)(x-b) \leq 0$ gilt. Bei der *Simpsonregel* können wir den Zwischenwertsatz nicht direkt anwenden, da g das Vorzeichen wechselt:

$$I(f) - I^{(2)}(f) = \int_a^b f\left[a, \frac{a+b}{2}, b, x\right] (x-a) \left(x - \frac{a+b}{2}\right) (x-b) dx$$

Um dennoch den Zwischenwertsatz anwenden zu können, schreiben wir

$$I(f) - I^{(2)}(f) = \int_a^b \frac{f\left[a, \frac{a+b}{2}, b, x\right] - f\left[a, \frac{a+b}{2}, b, \frac{a+b}{2}\right]}{x - \frac{a+b}{2}} (x-a) \left(x - \frac{a+b}{2}\right)^2 (x-b) dx -$$

$$f \left[a, \frac{a+b}{2}, b, \frac{a+b}{2} \right] \cdot \int_a^b f(x-a) \left(x - \frac{a+b}{2} \right) (x-b) dx$$

Beachte nun, dass $(x-a) \left(x - \frac{a+b}{2} \right)^2 (x-b) \geq 0$ und $\int_a^b f(x-a) \left(x - \frac{a+b}{2} \right) (x-b) dx$ punktsymmetrisch bzgl. $\frac{a+b}{2}$ ist. Damit gilt

$$I(f) - I^{(2)}(f) = \int_a^b \frac{f \left[a, \frac{a+b}{2}, b, x \right] - f \left[a, \frac{a+b}{2}, b, \frac{a+b}{2} \right]}{x - \frac{a+b}{2}} \underbrace{\left(x - a \right) \left(x - \frac{a+b}{2} \right)^2 (x-b)}_{\leq 0} dx \stackrel{\text{ZWS}}{=}$$

$$\begin{aligned} & f \left[a, \frac{a+b}{2}, \frac{a+b}{2}, b, \zeta_1 \right] \cdot \int_a^b \left(x - a \right) \left(x - \frac{a+b}{2} \right)^2 (x-b) dx = \\ & = \frac{1}{4!} f^{(4)}(\zeta_2) \frac{1}{120} (b-a)^5 \end{aligned}$$

Die Mittelpunktsregel wird in den Übungen behandelt. \square

BEMERKUNG:

- Bei $I^{(n)}$ mit $n \geq 7$ und $\hat{I}^{(n)}$ mit $n \geq 2$ treten negative Gewichte auf. Dadurch kann es möglicherweise zur Auslöschung kommen.
- Im Allgemeinen gilt *nicht* $I^{(n)}(f) \rightarrow I(f)$ für $n \rightarrow \infty$, da die Lagrange-Interpolation im Allgemeinen kein konvergenter Prozess ist. Man wendet daher $I^{(n)}$ nur auf Teilintervalle an. Das heißt

$$I_N^{(n)}(f) = \sum_{i=0}^{N-1} I_{[x_i, x_{i+1}]}^{(n)}(f) \quad \text{für } h = \frac{b-a}{N} \quad (3.14)$$

$$\Rightarrow I(f) - I_N^{(n)}(f) = \sum_{i=0}^{N-1} \omega_n h^{m+2} f^{(m+1)}(\zeta_i) =$$

$$= \omega_n h^{m+2} N f^{(m+1)}(\zeta) = \omega_n (a-b) h^{m+1} f^{(m+1)}(\zeta)$$

Dabei sei ω_n ein Faktor für ein $m \geq n$. Die letzte Zeile stellt dann eine Fehlerformel dar.

3.1.3 Gaußsche Quadraturformeln

$I^{(2)}$ integriert nicht nur $p \in \mathcal{P}^2$, sondern auch $p \in \mathcal{P}^4$ exakt. Kann die Wahl der x_0, \dots, x_n noch verbessert werden, dass sogar $p \in \mathcal{P}^k$ mit $k \geq 5$ exakt integriert werden?

BEMERKUNG: Der Fehler ist für n gerade $I - I^{(n)}(f) = O\left((b-a)^{n+1}\right)$ für n ungerade dagegen $O\left((b-a)^{n+1}\right)$.

Lemma 3.15

$I^{(n)}$ kann für $q \in \mathcal{P}^k$ mit $k \geq 2n+2$ nicht exakt sein.

BEWEIS: Sei

$$p(x) = \prod_{i=0}^n (x - x_i)^2 \in \mathcal{P}^{2n+2} \Rightarrow 0 < \int_a^b p(x) dx = I^{(n)}(p) = 0$$

Das ist ein Widerspruch und somit kann $I(f) - I^{(n)}(f)$ höchstens von der Ordnung $(a-b)^{n+2} c f^{(n+2)}(\zeta)$ sein. \square

Wir wollen nun eine geschickte Wahl der x_0, \dots, x_n herleiten. Hierzu ergänzen wir x_0, \dots, x_n durch x_{n+1}, \dots, x_{2n+1} . Dann gilt

$$I(f) - I^{(2n+1)}(f) = I(f) - \sum_{i=0}^{2n+1} f[x_0, \dots, x_n] \int_a^b \prod_{j=0}^{i-1} (x - x_j) dx =$$

$$I(f) - I^{(n)}(f) - \sum_{i=n+1}^{2n+1} f[x_0, \dots, x_i] \int_a^b \prod_{j=0}^{i-1} (x - x_j) dx$$

Ziel ist es nun, den letzten Faktor

$$\int_a^b \prod_{j=0}^{i-1} (x - x_j) dx = 0$$

zu setzen, denn dann wäre $I^{(n)}(f)$ so genau wie $I^{(2n+1)}(f)$. Für $i = n+1, \dots, 2n+1$ schreiben wir

$$\int_a^b \prod_{j=0}^{i-1} (x - x_j) dx = \int_a^b \prod_{j=0}^n (x - x_j) q_i(x) dx = 0 \quad (3.16)$$

für ein Polynom $q_i \in \mathcal{P}^{i-n-1} \subset \mathcal{P}^n$. Es soll also für alle Polynome $q \in \mathcal{P}^n$ gelten:

$$\int_a^b \prod_{j=0}^n (x - x_j) q(x) dx = 0 \quad (3.17)$$

oder äquivalent dazu

$$\left\langle \prod_{j=0}^n (x - x_j), q \right\rangle = 0 \Leftrightarrow \prod_{j=0}^n (x - x_j) \perp \mathcal{P}^n \quad (3.18)$$

Da $\prod_{j=0}^n (x - x_j) \in \mathcal{P}^{n+1}$, wäre für die Lösung der Gleichung z.B. das Legendrepolynom L_{n+1} ein Kandidat. Mit dem folgenden Lemma zeigen wir, dass es tatsächlich eine Lösung ist:

Lemma 3.19

$L_m(x)$ hat paarweise verschiedene Nullstellen.

BEWEIS: Sei $N_m = \{\lambda \in (a, b) : \lambda \text{ ist Nullstelle ungerader Ordnung von } L_m\}$, d.h. N_m ist die Menge der Nullstellen, bei denen L_m das Vorzeichen wechselt. Wir behaupten nun $|N_m| = m$. Hierzu ein Widerspruchsbeweis. Wir nehmen an, $|N_m| < m$. Dann liegt $q(x) = \prod_{\lambda \in N_m} (x - \lambda)$ in \mathcal{P}^m und $\langle q, L_m \rangle = 0$. Aber qL_m hat nur Nullstellen gerader Ordnung und damit *keinen* Vorzeichenwechsel, also wäre

$$\int_a^b q(x) L_m(x) dx \neq 0$$

Und dies ist ein Widerspruch. \square .

Wählen wir also x_0, \dots, x_n als Nullstellen des Legendre-Polynoms L_{n+1} , so gilt

$$I(f) - I^{(n)}(f) = I(f) - I^{(2n+1)}(f) = c_{2n+1} (b-a)^{n+2} f^{(2n+2)}(\zeta) \quad (3.20)$$

d.h. $I^{(n)}(f)$ auf \mathcal{P}^{2n+1} .

<!-- Local IspellDict: german-new8 -->

Kapitel 4

Lineare Gleichungssysteme - direkte Verfahren

Das Ziel dieses Kapitels ist es, das Gleichungssystem

$$Ax = b$$

für eine (zunächst) reguläre Matrix A zu lösen.

4.0.4 Eliminationsverfahren

Im einfachsten Fall ist A eine obere Dreiecksmatrix:

$$\begin{pmatrix} a_{11} & \cdots & a_{1n} \\ & \ddots & \vdots \\ 0 & & a_{nn} \end{pmatrix}$$

Dann gilt

$$Ax = b \Rightarrow x_n = \frac{1}{a_{nn}} b_n$$

und

$$x_j = \frac{1}{a_{jj}} \left(b_j - \sum_{k=j+1}^n a_{jk} x_k \right), \quad j = 1, \dots, n+1$$

Im allgemeinen Fall benutzt man z.B. das Gaußsche Eliminationsverfahren mit den elementaren Umformungen

- (i) Vertauschen zweier Gleichungen
- (ii) Addition eines Vielfachen einer Zeile (Gleichung) zu einer anderen Zeile
- (iii) Multiplizieren einer Zeile mit einer Zahl $\neq 0$.

In der Praxis wendet man das Verfahren auf die Matrix $(A|b)$ oder $(A|E_n)$ an.

Beschreibung des Verfahrens

1.Schritt:

- (*Pivotsuche*) Finde $a_{m1} \neq 0, m \in \{1, \dots, n\}$ So ein m existiert, da A regulär ist.

- (*Pivotisierung*) Vertausche 1. und m -te Zeile. Das Resultat sei $\tilde{A}^{(0)}, \tilde{b}^{(0)}$.
- Ziehe für jedes $r \in \{2, \dots, n\}$ das $\frac{\tilde{a}_{r1}^{(0)}}{\tilde{a}_{11}^{(0)}}$ -fache der 1. Zeile von der r -ten Zeile ab. Das Resultat sei $\tilde{A}^{(1)}, \tilde{b}^{(1)}$.

Dies entspricht der Multiplikation mit einer Permutationsmatrix P und einer Matrix G , also $[A^1|b^1] = GP[A|b]$. mit

$$P = \begin{pmatrix} 0 & 0 & 0 & 1 & & \\ 0 & 1 & 0 & & & \\ 0 & 0 & \ddots & & & \\ & & & 1 & & \\ 1 & & & 0 & & \\ & & & & \ddots & \\ & & & & & 1 \end{pmatrix} \quad \text{1. und } m\text{-te Spalte} \quad (4.1)$$

$$G = \begin{pmatrix} 1 & & & 0 \\ -q_{21} & 1 & & \\ \vdots & 0 & \ddots & \\ -q_{n1} & & & 1 \end{pmatrix} \quad (4.2)$$

2. bis n -ter Schritt: Wende den ersten Schritt auf die um eine Zeile und eine Spalte verkleinerte Matrix $[A^{(1)}|b^{(1)}]$ an. Am Ende erhält man

$$[R|c] = G_{n-1}P_{n-1} \cdots G_1P_1[A|b] \quad (4.3)$$

mit einer oberen Dreiecksmatrix R .

BEMERKUNG: Man kann die wesentlichen Elemente der G_j in den frei werdenden Stellen von A speichern.

Definition 4.4

Eine Matrix A heißt nilpotent, falls es ein $k \in \mathbb{N}$ gibt, so dass $A^k = 0$ gilt. Eine Matrix A heißt unipotent, falls $A - I$ nilpotent ist.

BEMERKUNG: Obere Dreiecksmatrizen mit allen Diagonalelementen = 0 sind nilpotent, obere Dreiecksmatrizen mit Diagonalelementen = 1 sind unipotent.

Theorem 4.5

Mit Hilfe des Gaußschen Eliminationsalgorithmus erhält man die Zerlegung

$$PA = LR \quad (4.6)$$

mit der Permutationsmatrix $P = P_{n-1} \cdots P_1$, der unipotenten unteren Dreiecksmatrix L und der oberen Dreiecksmatrix R . Für $P = 1$ ist die Zerlegung eindeutig.

BEMERKUNG: Die G_j haben die Form

$$\begin{pmatrix} 1 & & 0 & & & \\ & \dots & 0 & & & \\ 0 & & 1 & & & \\ & & -q_{21} & 1 & & \\ & & \vdots & & \ddots & \\ 0 & & -q_{n1} & & 0 & 1 \end{pmatrix}$$

Theorem 4.7

Sei A regulär und diagonaldominant, d.h.

$$\sum_{k=1, k \neq j}^n |a_{jk}| \leq |a_{jj}| \quad (4.8)$$

Dann kann die Gaußelimination ohne Pivotisierung durchgeführt werden.

Für positiv definite Matrizen haben wir also die Zerlegung $A = LR$. Es gilt

$$D = \begin{pmatrix} r_{11} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & r_{nn} \end{pmatrix} \text{ und } \tilde{R} = D^{-1}R$$

Dann ist \tilde{R} unipotent. Da A symmetrisch ist, gilt

$$LR = A = A^T = (LR)^T = (LD\tilde{R})^T = \tilde{R}^T (DL^T)$$

Da die LR -Zerlegung eindeutig ist, gilt

$$L = \tilde{R}^T \text{ bzw. } R = DL^T \Rightarrow A = LR = LDL^T = \tilde{L}\tilde{L}^T \text{ mit}$$

$$\tilde{L} = LD^{1/2} = L \cdot \begin{pmatrix} d_{11}^{1/2} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & d_{nn}^{1/2} \end{pmatrix}$$

Theorem 4.9

Jede symmetrische, positiv definite Matrix A hat eine Cholesky-Zerlegung $A = \tilde{L}\tilde{L}^T$ mit einer unteren Dreiecksmatrix \tilde{L} .

BEMERKUNG: Wählt man die Diagonalelemente von \tilde{L} positiv, so ist die Cholesky-Zerlegung eindeutig.

Bandmatrizen

Definition 4.10

Eine Matrix $A = (a_{jk})_{j,k}$ heißt Bandmatrix vom Typ (m_l, m_r) mit $0 \leq m_l, m_r \leq n - 1$, falls

$$a_{jk} = 0 \text{ für } j + m_r < k < j - m_l \quad (4.11)$$

Das heißt, A ist auf der Hauptdiagonale und höchstens auf $m_r + m_l$ Nebendiagonalen ungleich Null. Die Größe $m = m_r + m_l$ heißt Bandbreite.

BEISPIEL: Bandmatrizen vom Typ

- $(n - 1, 0)$ heißen untere Dreiecksmatrizen
- $(0, n - 1)$ heißen obere Dreiecksmatrizen
- $(1, 1)$ heißen Tridiagonalmatrizen

BEISPIEL: Die $(4, 4)$ -Bandmatrix

$$A = \begin{pmatrix} B & -I & 0 & 0 \\ -I & B & -I & 0 \\ 0 & -I & B & -I \\ 0 & 0 & -I & B \end{pmatrix} \text{ und } B = \begin{pmatrix} 4 & -1 & 0 & 0 \\ -1 & 4 & -1 & 0 \\ 0 & -1 & 4 & -1 \\ 0 & 0 & -1 & 4 \end{pmatrix}$$

und der 4×4 -Einheitsmatrix ist diagonaldominant, aber nicht strikt.

Theorem 4.12

Ist die LR -Zerlegung von A ohne Zeilenvertauschung durchführbar und A vom Typ (m_l, m_r) , so sind L und R vom Typ $(m_l, 0)$ und $(0, m_r)$. Der Aufwand der Zerlegung ist

$$\frac{1}{3}nm_l m_r \quad (4.13)$$

BEMERKUNG: Es genügt das Band zu speichern. Dadurch wird das Problem reduziert und man kann im schnellen Speicher rechnen.

Nachiteration

Sei $LR = PA$ eine Zerlegung. Durch Rechenfehler hat man nur $\tilde{L}\tilde{R} = PA$. Dadurch ergibt sich die fehlerhafte Lösung

$$\tilde{L}\tilde{R}x_0 = b$$

Das heißt für den Defekt d^0 :

$$d^0 = b - Ax_0 \neq 0 \Rightarrow d^0 = A(x - x_0) \Rightarrow A^{-1}d^0 = x - x_0$$

Man müsste also x_0 um $A^{-1}d$ korrigieren. Stattdessen definiert man die Iteration

$$x^1 = x^0 + (\tilde{L}\tilde{R})^{-1} d^0 \text{ und } x^{j+1} = x^j + (\tilde{L}\tilde{R})^{-1} d^j$$

Es gilt

$$x^{j+1} = x^j + \underbrace{(\tilde{L}\tilde{R})^{-1} (b - Ax^j)}_{=F(x^j)}$$

Also $F(x) = x + (\tilde{L}\tilde{R})^{-1} (b - Ax)$. Damit erhalten wir

$$F(x) - F(y) = x - y + (\tilde{L}\tilde{R})^{-1} A(y - x) = \left(I - (\tilde{L}\tilde{R})^{-1} A \right) (x - y)$$

Die Iteration konvergiert, falls $\|I - (\tilde{L}\tilde{R})^{-1} A\| < 1$. Dies verbessert die Lösung. Es gilt

$$I - (\tilde{L}\tilde{R})^{-1} A = (\tilde{L}\tilde{R})^{-1} (\tilde{L}\tilde{R} - A) \text{ und } \tilde{L}\tilde{R} = A \left(I - A^{-1} (A - \tilde{L}\tilde{R}) \right)$$

Das heißt

$$(\tilde{L}\tilde{R})^{-1} = \left(I - A^{-1} (A - \tilde{L}\tilde{R}) \right)^{-1} A^{-1} \Rightarrow$$

$$\|(\tilde{L}\tilde{R})^{-1}\| \leq \frac{\|A^{-1}\|}{1 - \|A^{-1}\| \cdot \|A - \tilde{L}\tilde{R}\|} = \frac{\|A^{-1}\|}{1 - \text{cond}(A) \frac{\|A - \tilde{L}\tilde{R}\|}{\|A\|}}$$

$$\Rightarrow \|I - (\tilde{L}\tilde{R})^{-1} A\| \leq \frac{\text{cond}(A)}{1 - \text{cond}(A) \frac{\|A - \tilde{L}\tilde{R}\|}{\|A\|}} \cdot \frac{\|A - \tilde{L}\tilde{R}\|}{\|A\|}$$

BEMERKUNG: Man kann $\tilde{L}\tilde{R}$ mit geringerer Genauigkeit berechnen und dann mit hoher Genauigkeit nachiterieren.

Nichtreguläre Systeme

Sei $A \in \mathbb{R}^{m \times n}$ und $b \in \mathbb{R}^m$. Unser Ziel ist nun, $Ax = b$ zu lösen. Allerdings haben wir das folgende Problem: Falls $\text{rang}(A) < \text{rang}[A|b]$, so ist $Ax = b$ nicht im eigentlichen Sinne lösbar. Stattdessen suchen wir eine Lösung x , die das *Residuum* $\|Ax - b\|_2$ minimiert. Hierzu müssen wir

$$J(x) = \frac{1}{2} \|Ax - b\|_2^2 = \frac{1}{2} \langle Ax - b, Ax - b \rangle$$

minimieren. Es gilt

$$\delta J(x)(y) = \frac{d}{dt} J(x + ty) \Big|_{t=0} =$$

$$\frac{d}{dt} \left(\frac{1}{2} \langle Ax - b, Ax - b \rangle + t \langle Ax - b, Ay \rangle + \frac{1}{2} t^2 \langle Ay, Ay \rangle \right) \Big|_{t=0} = \langle Ax - b, Ay \rangle$$

Also muss das Minimum \bar{x}

$$\langle A\bar{x} - b, Ay \rangle = 0 \forall y$$

erfüllen. Dies bedeutet

$$\begin{aligned} \langle A^T (A\bar{x} - b), y \rangle &= 0 \Rightarrow \\ A^T A\bar{x} &= A^T b \end{aligned} \tag{4.14}$$

Dies ist die *Normalgleichung*.

Theorem 4.15

Es existiert stets ein $\bar{x} \in \mathbb{R}^n$ welches

$$x \mapsto \|Ax - b\|_2 \tag{4.16}$$

minimiert. Es gilt $A^T A\bar{x} = A^T b$ (Normalgleichung). Für $\text{rang}(A) = n$ ist \bar{x} eindeutig. Die allgemeinen Lösungen sind gegeben durch $\bar{x} + \text{Kern}(A)$.

BEWEIS: Zunächst zur Existenz einer Lösung der Normalgleichung. Es gilt $\mathbb{R}^n = \text{Bild}(A) \oplus \text{Kern}(A^T)$ (orthogonale Summe). Beachte: $\langle Ax, y \rangle = \langle x, A^T y \rangle$. Somit ist die Darstellung $b = \alpha + \beta$ mit $\alpha \in \text{Bild}(A)$ und $\beta \in \text{Kern}(A^T)$ eindeutig. Da $\alpha \in \text{Bild}(A)$, existiert ein $\bar{x} \in \mathbb{R}^n$ mit $A\bar{x} = \alpha$. Damit folgt

$$A^T A\bar{x} = A^T \alpha = A^T (\alpha + \beta) = A^T b$$

D.h., \bar{x} löst die Gleichung. Weiterhin gilt für alle x :

$$\|Ax - b\|_2^2 = \|A\bar{x} - b + A(x - \bar{x})\|_2^2 = \|A\bar{x} - b\|_2^2 + 2 \langle A\bar{x} - b, A(x - \bar{x}) \rangle + \|A(x - \bar{x})\|_2^2 \geq \|A\bar{x} - b\|_2^2$$

Das heißt, \bar{x} minimiert $x \mapsto \|Ax - b\|$. Umgekehrt wissen wir schon, dass jedes Minimum die Normalgleichung erfüllt. Seien nun \bar{x}, \hat{x} Lösungen, so gilt

$$b = \underbrace{A\bar{x}}_{\in \text{Bild}(A)} + \underbrace{(b - A\bar{x})}_{\in \text{Kern}(A^T)} = \underbrace{A\hat{x}}_{\in \text{Bild}(A)} + \underbrace{(b - A\hat{x})}_{\in \text{Kern}(A^T)}$$

Da aber $\mathbb{R}^n = \text{Bild}(A) \oplus \text{Kern}(A^T)$, gilt

$$A\bar{x} = A\hat{x} \Rightarrow \hat{x} - \bar{x} \in \text{Kern}(A) \Rightarrow \hat{x} =$$

$$\bar{x} + (\hat{x} - \bar{x}) \in \bar{x} + \text{Kern}(A)$$

Andererseits gilt $\|A(x + y) - b\| = \|Ax - b\|$ für jedes $y \in \text{Kern}(A)$. Damit folgt die Behauptung. \square

Gaußsche Ausgleichsrechnung

Gegeben seien Messpunkte (x_j, y_j) , $j = 1, \dots, n$. Unser Ziel ist nun: Finde zu gegebenen Funktionen u_1, \dots, u_n , $n < m$ eine Linearkombination

$$u = \sum_{k=1}^n c_k u_k$$

derart, dass

$$\Delta_2 = \left(\sum_{j=1}^m |u(x_j) - y_j|^2 \right)^{1/2} \quad (4.17)$$

minimiert wird. Setze hierzu

$$c = (c_1, \dots, c_n)^T, \quad y = (y_1, \dots, y_m)^T, \quad a_k = (u_k(x_1), \dots, u_k(x_m)), \quad k = 1, \dots, n, \quad A = [a_1 | \dots | a_n]$$

Dann müssen wir

$$F(x) = \frac{1}{2} \|Ac - y\|_2^2$$

minimieren. ($\Rightarrow A^T Ac = A^T y =$). Wir unterscheiden mehrere Fälle.

- Fall 1: $u_k(x) = x^{k-1}$, $k = 1, \dots, n$. Dann nennen wir u *Ausgleichsparabel*.
- Fall 2: $u_1(x) = 1$, $u_2(x) = x$ (linearer Fall).

x_i	-2	-1	0	1	2
y_i	0.5	0.5	2	3.5	3.5

Dies führt zu einem überbestimmten System

$$\begin{pmatrix} 1 & -2 \\ 1 & -1 \\ 1 & 0 \\ 1 & 1 \\ 1 & 2 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 0.5 \\ 0.5 \\ 2 \\ 3.5 \\ 3.5 \end{pmatrix}$$

Löse nun

$$\underbrace{\begin{pmatrix} 5 & 0 \\ 0 & 10 \end{pmatrix}}_{=A^T A} \cdot \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 10 \\ 9 \end{pmatrix}$$

Dies ergibt

$$\begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} 2 \\ 0.9 \end{pmatrix}$$

Somit ist $u(x) = 2 + 0.9x$ die *Ausgleichsgerade*.

BEMERKUNG: Man kann auch nichtlineare Größen durch Transformationen approximieren. Betrachte beispielsweise ein physikalisches Gesetz der Form

$$y(x) = \frac{a}{1 + bx}$$

Durch Umformen erhält man $\frac{1}{a} + \frac{b}{a}x = \frac{1}{y(x)}$ (dies ist linear in x mit den neuen Größen $\tilde{a} = \frac{1}{a}$, $\tilde{b} = \frac{b}{a}$).

Bei der Gaußschen Ausgleichsrechnung ist $\overline{A^T} A = A^* A$ wesentlich. Diese Matrix hat eine spezielle Struktur.

Lemma 4.18

Sei $A \in \mathbb{K}^{m \times n}$. Dann ist A^*A hermitesch (bzw. symmetrisch) und positiv semidefinit. Ist $\text{rang}(A) = n$, so ist A^*A positiv definit.

BEWEIS: Es gilt

$$\langle A^*A\zeta, \zeta \rangle = \langle A\zeta, A\zeta \rangle = \|A\zeta\|_2^2 \geq 0$$

und insbesondere auch $(A^*A)^* = A^*A$. Ist $\text{Rang}(A) = n$ (dies benötigt $m \geq n!$), so ist $A: \mathbb{R}^n \rightarrow \mathbb{R}^n$ injektiv. Damit folgt aus $Ax = 0$ schon $x = 0$. Also folgt aus $\langle A^*A\zeta, \zeta \rangle = 0 \Rightarrow \|A\zeta\|_2^2 = 0 \Rightarrow \zeta = 0$. Also ist A^*A positiv definit. \square

BEMERKUNG: In der Regel ist A^*A schlecht konditioniert.

BEMERKUNG: Bei symmetrischen, positiv definiten Matrizen kann man die LR -Zerlegung durch die Cholesky-Zerlegung vereinfachen:

$$\begin{pmatrix} \tilde{l}_{11} & 0 & 0 \\ \vdots & \ddots & 0 \\ \tilde{l}_{n1} & \cdots & \tilde{l}_{nn} \end{pmatrix} \cdot \begin{pmatrix} \tilde{l}_{11} & \cdots & \tilde{l}_{n1} \\ 0 & \ddots & \vdots \\ 0 & 0 & \tilde{l}_{nn} \end{pmatrix} = \begin{pmatrix} a_{11} & \cdots & a_{1n} \\ \vdots & & \vdots \\ a_{n1} & \cdots & a_{nn} \end{pmatrix}$$

Dann gilt

$$l_{11}^2 = a_{11} \Rightarrow l_{11} = \sqrt{a_{11}}$$

$$\tilde{l}_{j1} \cdot \tilde{l}_{11} = a_{j1}, j \geq 2 \Rightarrow \tilde{l}_{j1} = \frac{\tilde{l}_{11}}{\sqrt{a_{11}}}$$

Weiters folgt

$$\underbrace{\sum_{k=1}^i \tilde{l}_{ik}^2}_{\text{bekannt}} = a_{ii} \Rightarrow l_{ii} = \sqrt{a_{ii} - \tilde{l}_{i1}^2 - \cdots - \tilde{l}_{i,i-1}^2}$$

$$\tilde{l}_{ji} = \frac{1}{\tilde{l}_{ii}} a_{ji} - \tilde{l}_{j1} \tilde{l}_{i1} - \cdots - \tilde{l}_{j,i-1} \tilde{l}_{i,i-1}, i = i+1, \dots, n$$

BEMERKUNG: Rundung kann dieses Verfahren jedoch verhindern. Betrachte z.B. bei dreistelliger Rundung

$$A \begin{pmatrix} 1.03 & 1.10 \\ 1.03 & 1.11 \\ 1.07 & 1.15 \end{pmatrix} \Rightarrow A^T A = \begin{pmatrix} 3.43 & 3.60 \\ 3.60 & 3.76 \end{pmatrix}$$

aber

$$\left\langle A^T A \begin{pmatrix} -1 \\ 1 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \end{pmatrix} \right\rangle = -0.01 < 0$$

Wir wollen deshalb eine neue Methode entwickeln, die ohne die Normalengleichung auskommt.

Theorem 4.19

Sei $A \in \mathbb{K}^{m \times n}$ mit $m \geq n$ und $\text{Rang}(A) = n$. Dann existiert eine eindeutige Matrix $Q \in \mathbb{R}^{m \times n}$ mit $Q^*Q = E_n$ und eine eindeutige obere Dreiecksmatrix $R \in \mathbb{R}^{n \times n}$ mit positiven Diagonalelementen, so

dass

$$A = QR \tag{4.20}$$

BEWEIS: Zunächst zur Eindeutigkeit: $A = QR = Q_2 R_2 \Rightarrow$

$$\underbrace{Q_2^T Q}_\text{orthogonal} = R_2 R^{-1} = S$$

Zur Orthogonalität von $Q_2^T Q$:

$$(Q_2^T Q)^T Q_2^T Q = Q^T Q_2 Q_2^T Q = E_n$$

S ist dann eine obere Dreiecksmatrix mit positiven Diagonaleinträgen, denn

$$Q_2^T Q = E_n = R_2 R^{-1} \text{ und}$$

$$S = R_2 R^{-1} = (s_1 | \cdots | s_n), \langle s_1, s_1 \rangle = 1 \Rightarrow s_{11} = 1, \langle s_1, s_j \rangle = 0 \Rightarrow s_{1j} = 0 \text{ für } j \geq 2$$

Wir erhalten $R_2 = R$ und $Q = AR^{-1} = AR_2^{-1} = Q_2$. Jetzt zur Existenz: Wir wenden das *Gram-Schmidt-Verfahren* auf die Spaltenvektoren a_1, \dots, a_n von A an. Also

$$q_1 = \frac{a_1}{\|a_1\|}$$

$$\tilde{q}_k = a_k - \sum_{i=1}^{k-1} \langle a_k, q_i \rangle q_i$$

$$q_k = \frac{\tilde{q}_k}{\|\tilde{q}_k\|}$$

Da $\text{Rang}(A) = n$ gilt, sind a_1, \dots, a_n linear unabhängig und das Verfahren bricht somit nicht ab ($\|\tilde{q}_k\| \neq 0$). Ferner gilt $Q^* Q = E_n$ nach Konstruktion. Weiterhin ist

$$a_k = \tilde{q}_k + \sum_{i=1}^{k-1} \langle a_k, q_i \rangle q_i = \underbrace{\|\tilde{q}_k\|_2}_{=r_{kk}} q_k + \sum_{i=1}^{k-1} \underbrace{\langle a_k, q_i \rangle}_{=r_{ik}} q_i, \quad i < k$$

Setze nun noch $r_{ik} = 0$ für $i > k$, so folgt $A = QR$ und die Behauptung ist bewiesen. \square

Doch nun zurück zur Normalengleichung. Aus $A^T A x = A^T b$ und $A = QR$ wird $R^T R x = A^T b$, d.h. $R^T R x = A^T b$. Dies ist durch Rückwärtseinsetzen lösbar.

BEMERKUNG: Für den Aufwand der Verfahren gilt die Beziehung

$$\text{Aufwand}(QR) \approx 2 \cdot \text{Aufwand}(LR)$$

Dennoch haben wir das folgende Problem: Das Gram-Schmidt-Verfahren ist bei Rundungsfehlern ungeeignet, da die Orthonormalität der Spalten rasch verloren geht.

Householder-Verfahren

Für $v \in \mathbb{K}^m$ sei $v \otimes \bar{v} = v \bar{v}^T = (\bar{v}_i, v_j)_{i,j}$ das sog. *dyadische Produkt*.

Definition 4.21

Für $v \in \mathbb{K}^n$ mit $\|v\|_2 = 1$ heißt

$$S = I - 2v \otimes \bar{v} \tag{4.22}$$

Householder-Transformation.

Es gilt

$$S = S^* = S^{-1} \quad (4.23)$$

Denn es ist

$$S^* = I - 2(v \otimes \bar{v})^* = I - vv^* = S$$

und

$$SS = (I - 2vv^*)(I - 2vv^*) = I - 2vv^* - 2vv^* + 4vv^* = I, \text{ d.h. } S^{-1} = S$$

Damit ist S hermitesch und unitär.

Die Eigenwerte von S sind ± 1 .

- (i) Ist $w \in \text{orth}(v)$, so gilt $Pw = w \Rightarrow$ der Eigenwert ist 1. Nehme dann $n - 1$ linear unabhängige Vektoren aus $\text{orth}(v)$.
- (ii) $Pv = v - 2vv^T v = -v \Rightarrow$ der Eigenwert ist -1 und der Eigenvektor ist v . Für die Determinante von S gilt dann $\det(S) = 1$.

Algorithmus-Idee: Wähle S_1, \dots, S_n so, dass iterativ

$$A^{(i)} = S_i A^{(i-1)} \text{ mit } A^{(0)} = A$$

die Form

$$A^{(i)} = \begin{pmatrix} \star & \cdots & \star & \star & \star \\ 0 & \ddots & \star & & \star \\ 0 & 0 & \star & & \star \\ 0 & 0 & 0 & \star & \star \end{pmatrix}$$

hat. Nach n Schritten ist

$$\tilde{R} = A^{(n)} = \underbrace{S_n \cdots S_1}_{\text{unitär}} A = \tilde{Q}^* A \Rightarrow \tilde{Q} \tilde{R} = A$$

mit $R \in \mathbb{K}^{n \times n}$:

$$\tilde{R} = \begin{pmatrix} R \\ - \\ 0 \end{pmatrix}$$

$$\Rightarrow A = \tilde{Q} \tilde{R} = \begin{bmatrix} Q \\ \underbrace{\quad ? \quad} \\ \text{egal} \end{bmatrix} \cdot \begin{pmatrix} R \\ - \\ 0 \end{pmatrix} \Rightarrow A = QR$$

Bestimme nun die Matrizen S_1, \dots, S_n . Zu S_1 : Wir müssen a_1 an $\text{span}(e_1)$ spiegeln. Wir erhalten

$$v = \frac{a + \|a\|e_1}{\|a + \|a\|e_1\|}$$

Nach dem ersten Schritt ist dann

$$\begin{pmatrix} 1 & & & \\ 0 & \star & \cdots & \star \\ \vdots & \vdots & & \vdots \\ 0 & \star & \cdots & \star \end{pmatrix}$$

Im i -ten Schritt wende das Verfahren auf die reduzierte Matrix an:

$$S_2 = \begin{pmatrix} I & 0 \\ 0 & I - 2\tilde{v}\tilde{v}^* \end{pmatrix} = I - 2vv^*$$

mit

$$v = \begin{pmatrix} \underbrace{0 \cdots 0}_{i-1 \text{ Nullen}} & |\tilde{v} \end{pmatrix}$$

BEMERKUNG: Die durch die Householder-Transformation erhaltene QR -Zerlegung ist nicht die bekannte QR -Zerlegung. Dies liegt daran, dass die Diagonalelemente r_{ii} nicht positiv sein müssen.

BEISPIEL: Wir wollen nun die QR -Zerlegung der Matrix

$$\begin{pmatrix} 1 & 2 & 3 \\ 0 & 3 & 2 \\ 2 & 0 & 1 \end{pmatrix}$$

mittels Householder-Transformationen bestimmen. Im ersten Schritt ist

$$a_1 = \begin{pmatrix} 1 \\ 0 \\ 2 \end{pmatrix} \Rightarrow \|a_1\| = \sqrt{5}$$

$$\Rightarrow \tilde{v}_1 = a_1 \pm \|a_1\|e_1, \quad v_1 = \frac{\tilde{v}_1}{\|\tilde{v}_1\|}$$

Da $\text{sgn}(a_{11}) > 0$ ist, nehme $+$. Wir erhalten damit

$$\tilde{v}_1 = \begin{pmatrix} 1 + \sqrt{5} \\ 0 \\ 2 \end{pmatrix}, \quad \|\tilde{v}_1\|^2 = 10 + 2\sqrt{5}$$

Somit ist

$$S_1 = I - 2v_1v_1^* = I - 2\frac{\tilde{v}_1\tilde{v}_1^*}{\|\tilde{v}_1\|^2} = \begin{pmatrix} -1/\sqrt{5} & 0 & -2/\sqrt{5} \\ 0 & 1 & 0 \\ -2/\sqrt{5} & 0 & 1/\sqrt{5} \end{pmatrix}$$

$$\Rightarrow A^{(1)} = S_1A \approx \begin{pmatrix} -2.236 & -0.894 & -2.236 \\ 0 & 3 & -2 \\ 0 & -1.788 & -2.236 \end{pmatrix}$$

$$\tilde{a}_2 = \begin{pmatrix} 3 \\ -1.788 \end{pmatrix}, \quad \text{sgn}(\tilde{a}_{22}) > 0$$

$$\Rightarrow \tilde{v}_2 = \tilde{a}_2 + \|\tilde{a}_2\|e_2, \quad v_2 = \frac{\tilde{v}_2}{\|\tilde{v}_2\|} \Rightarrow \tilde{v}_2 = \begin{pmatrix} 6.492 \\ 1.788 \end{pmatrix}$$

Dann ist

$$\|\tilde{v}_2\|^2 \approx 45.346$$

und

$$\tilde{S}_2 = I_2 - \frac{\tilde{v}_2 \tilde{v}_2^*}{\|\tilde{v}_2\|^2} = \begin{pmatrix} -0.859 & 0.512 \\ 0.512 & 0.859 \end{pmatrix}$$

$$S_2 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \star & \star \\ 0 & \star & \tilde{S}_2 \end{pmatrix}$$

$$\Rightarrow A = S_2 A^{(1)} = \begin{pmatrix} -2.236 & -0.894 & -2.236 \\ 0 & -3.492 & -2.862 \\ 0 & 0 & -0.896 \end{pmatrix} = R$$

Da S_j unitär und hermitesch, folgt $S_2 S_1 A = R \Rightarrow A = S_1 S_2 R = QR$ und

$$Q = S_1 S_2 = \begin{pmatrix} -0.447 & -0.458 & -0.768 \\ 0 & -0.859 & 0.512 \\ -0.894 & 0.229 & 0.384 \end{pmatrix}$$

Möchte man $R_{ii} > 0$, so kann man QR durch $(QD)(DR)$ ersetzen mit

$$D = \begin{pmatrix} -1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & -1 \end{pmatrix}$$

Dann gilt $(QD)^*(QD) = D^*Q^*QD = D^*D = E_3$.

Singulärwertzerlegung

Die bisher vorgestellten Methoden zur Lösung linearer Gleichungssysteme oder Ausgleichsprobleme werden instabil, falls A schlecht konditioniert ist. Dadurch ist die Bestimmung der Ranges mittels LR oder QR -Zerlegung nicht mit genügender Sicherheit möglich. Die beste aktuelle Methode ist die *Singulärwertzerlegung* (englisch: *svd-singular value decomposition*). Hierbei wird A von links und rechts orthogonal transformiert. Sei nun $A \in \mathbb{R}^{m \times n}$ und $Q \in \mathbb{R}^{m \times m}$, $Z \in \mathbb{R}^{n \times n}$ orthonormal, so gilt

$$\|QAZ\|_2 = \|A\|_2 \text{ und } \|(QAZ)^{-1}\|_2 = \|Z^*A^{-1}Q^*\|_2 = \|A^{-1}\|_2$$

Somit folgt

$$\text{cond}_2(A) = \text{cond}_2(QAZ)$$

Die Kondition ändert sich also *nicht*.

Theorem 4.24

Sei $A \in \mathbb{R}^{m \times n}$. Dann existieren orthonormale Matrizen $U \in \mathbb{R}^{m \times m}$ und $V \in \mathbb{R}^{n \times n}$ so, dass

$$U^T A V = \text{diag}(\sigma_1, \dots, \sigma_p) = D \in \mathbb{R}^{m \times n} \quad (4.25)$$

mit $p = \min\{m, n\}$, $\sigma_1 \geq \dots \geq \sigma_p$ und

$$\text{diag}(\sigma_1, \dots, \sigma_p) = \begin{pmatrix} \sigma_1 & & & & & \\ & \ddots & & & & \\ & & \vdots & & & \\ & & & \sigma_p & & \\ & & & & \dots & \\ & & & & & 0 \end{pmatrix} \text{ bzw. } \begin{pmatrix} \sigma_1 & & & & & \\ & \ddots & & & & \\ & & \vdots & & & \\ & & & \sigma_p & & \\ & & & & \dots & \\ & & & & & 0 \end{pmatrix} \quad (4.26)$$

Man nennt die σ_i die Singulärwerte von A .

BEMERKUNG: Die so definierten σ_i sind eindeutig.

Aus dem Satz folgt

$$Av_i = \sigma_i u_i, \quad A^T u_i = \sigma_i v_i \text{ mit } \|u_i\| = \|v_i\| = 1$$

Man nennt die u_i die *linkssingulären* und die v_i die *rechtssingulären* Werte. Aus $U^T AV = D$ folgt

$$D^2 = (U^T AV)^T U^T AV = V^T A^T U U^T AV = V^T A^T AV \Rightarrow VD^2 = VV^T A^T AV = A^T AV$$

$$\Rightarrow \sigma_i^2 v_i = A^T Av_i$$

Somit sind die v_i Eigenvektoren von $A^T A$ zu Eigenwerten σ_i^2 . Analog sind die u_i die Eigenvektoren von AA^T zu Eigenwerten σ_i^2 . Hieraus lässt sich ein Existenzbeweis machen:

BEWEIS: Sei $B = A^T A \in \mathbb{R}^{n \times n}$. Dann ist B symmetrisch. Demnach existiert eine orthonormale Matrix $V = (v_1, \dots, v_n) \in \mathbb{R}^{n \times n}$ mit

$$V^T B V = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}$$

und $\lambda_1 \geq \dots \geq \lambda_n$ Eigenwerten von B . Da $\text{Rang}(A^T A) = \text{Rang}(A)$ und die λ_i positiv sind, können wir für $i = 1, \dots, r$

$$u_i = \frac{1}{\sqrt{\lambda_i}} A v_i \in \mathbb{R}^m$$

setzen. Es folgt

$$\langle u_i, u_j \rangle = \frac{1}{\sqrt{\lambda_i \lambda_j}} \langle v_i, A^T A v_j \rangle = \frac{\sqrt{\lambda_j}}{\sqrt{\lambda_i}} \delta_{ij}$$

Ergänze nun u_1, \dots, u_r zu einem vollständigen Orthonormalsystem u_1, \dots, u_m und setze

$$D = \begin{pmatrix} \sqrt{\lambda_1} & & & \\ & \ddots & & \\ & & \sqrt{\lambda_r} & \\ & & & 0 & \\ & & & & \ddots \end{pmatrix} \in \mathbb{R}^{m \times n}$$

Es gilt

$$(u_1, \dots, u_m) \begin{pmatrix} \sqrt{\lambda_1} & & \\ & \ddots & \\ & & \sqrt{\lambda_r} \end{pmatrix} = A(v_1, \dots, v_r) \Rightarrow UD = AV$$

Es kommen auf beiden Seiten nur Nullen dazu, und somit folgt $D = U^T AV$. \square

Betrachte nun $A = UDV^T$ bzw. $U^T AV = D$ mit

$$\text{diag}(\sigma_1, \dots, \sigma_p) = \begin{pmatrix} \sigma_1 & & 0 & \dots & 0 \\ & \ddots & \vdots & & \vdots \\ & & \sigma_p & 0 & \dots & 0 \end{pmatrix} \text{ bzw. } \begin{pmatrix} \sigma_1 & & & \\ & \ddots & & 0 \\ & & & \sigma_p \\ & 0 & \dots & 0 \\ \vdots & & & \vdots \\ 0 & \dots & & 0 \end{pmatrix}$$

Dann ist

- (i) $\text{Rang}(A) = r$
- (ii) $\text{Kern}(A) = \text{span}\{v_{r+1}, \dots, v_n\}$
- (iii) $\text{Bild}\{u_1, \dots, u_r\}$
- (iv) $A = UDV^T = \sum_{i=1}^r \sigma_i u_i v_i = \sum_{i=1}^r \sigma_i u_i v_i^T$

Betrachten wir nun das Problem zur Bestimmung des *numerischen Ranges*

$$\text{Rang}(A, \varepsilon) = \min_{\|A-B\|_2 \leq \varepsilon} \text{Rang}(B)$$

Es heißt A *numerisch rang-defizient*, falls

$$\text{Rang}(A, \varepsilon) < \min\{m, n\}$$

mit $\varepsilon = \text{eps}\|A\|_2$. Bei Messreihen ist ε an die Genauigkeit der Messwerte gekoppelt.

Theorem 4.27

(Fehlerabschätzung) Sei $A = UDV^T$ wie oben, so gilt für

$$A_k = \sum_{i=1}^r \sigma_i u_i \otimes v_i \quad k < r \tag{4.28}$$

die Abschätzung

$$\min_{\text{Rang}(B)=k} \|A - B\|_2 = \|A - A_k\|_2 = \sigma_{k+1} \tag{4.29}$$

Daraus ergibt sich für $r_\varepsilon = \text{Rang}(A, \varepsilon)$:

$$\sigma_1 \geq \dots \geq \sigma_{r_\varepsilon} > \varepsilon > \sigma_{r_\varepsilon-1} \geq \dots \geq \sigma_{\min\{m,n\}} \tag{4.30}$$

BEWEIS: Es sei $U^T A V = \text{diag}(\sigma_1, \dots, \sigma_r)$ und $U^T A_k V = \text{diag}(\sigma_1, \dots, \sigma_k, 0, \dots, 0)$. Damit folgt $U^T (A - A_k) V = \text{diag}(0, \dots, 0, \sigma_{k+1}, 0, \dots, 0)$ und daher

$$\|A - A_k\|_2 = \|\text{diag}(0, \dots, 0, \sigma_{k+1}, \dots, \sigma_r)\|_2 = \sigma_{k+1} \text{ bzw. } \|A - B\|_2 \geq \sigma_{k+1} \text{ für } \text{Rang}(B) = k$$

Dazu wähle eine Orthonormalbasis für $\text{Ker}(B) = \text{span}\{x_1, \dots, x_{n-k}\}$. Aus Dimensionsgründen folgt dann

$$\text{Ker}(B) \cap \text{span}\{v_1, \dots, v_{k+1}\} \neq \emptyset$$

Sei z aus dieser Menge mit $\|z\|_2 = 1$. Dann gilt

$$Bz = 0 \text{ und } Az = \sum_{i=1}^r \sigma_i (u_i \otimes v_i) z = \sum_{i=1}^{k+1} \sigma_i u_i v_i^T z$$

$$\Rightarrow \|A - B\|_2^2 \geq \|(A - B)z\|_2^2 = \|Az\|_2^2 = \sum_{i,j=1}^k |\sigma_i|^2 (v_i^T z)^2 \geq \sigma_{k+1}^2 \sum_{i=1}^{k+1} (v_i^T z)^2 = \sigma_{k+1}^2$$

□

Mit der Singulärwertzerlegung kann man auch das Ausgleichsproblem lösen. $\|Ax - b\|_2 = \min!$ führte zu $A^T A x = A^T b$. Ist $\text{Rang}(A) = n$, so gilt $x = (A^T A)^{-1} A^T b$. Ansonsten gibt es ∞ -viele Lösungen. Diejenige mit kleinster Norm $\|x\|_2$ heißt *Minimallösung* des Ausgleichsproblems.

Theorem 4.31

Sei $A = UDV^T$ die (eine) Singulärwertzerlegung von $A \in \mathbb{R}^{m \times n}$. Dann ist

$$\bar{x} = \sum_{i=1}^r \frac{u_i^T b}{\sigma_i} v_i \quad (4.32)$$

die eindeutige Minimallösung des Ausgleichsproblems. Der Fehler genügt der Beziehung

$$\|A\bar{x} - b\|_2^2 = \sum_{i=r+1}^m (u_i^T b)^2 \quad (4.33)$$

BEWEIS: Da $x \mapsto \|x\|$ strikt konvex ist, ist die Minimallösung eindeutig. Für jedes $x \in \mathbb{R}^n$ gilt

$$\|Ax - b\|_2^2 = \|AVV^T x - b\|_2^2 = \|(U^T AV) V^T x - U^T b\|_2^2 = \|DV^T x - U^T b\|_2^2$$

Sei $z = V^T x$ (d.h. $x = Vz$), so gilt

$$\|Ax - b\|_2^2 = \|Dz - U^T b\|_2^2 = \sum_{i=1}^r (\sigma_i z_i - u_i^T b)^2 + \sum_{i=r+1}^m (u_i^T b)^2$$

Die Minimallösung muss zusätzlich noch $\|x\|_2 = \|z\|_2$ minimieren. Das Minimum bezüglich z erfüllt

$$\sigma_i z_i = u_i^T b \Rightarrow z_i = \frac{1}{\sigma_i} u_i^T b \underset{x=Vz}{\Rightarrow} x = \sum_i v_i z_i = \sum_i \frac{u_i^T b}{\sigma_i} v_i$$

Der minimale Fehler ist dann

$$\|A\bar{x} - b\|_2^2 = \sum_{i=r+1}^m (u_i^T b)^2$$

□

BEMERKUNG: Sei $D^+ = \text{diag}\left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_r}, 0, \dots, 0\right) \in \mathbb{R}^{n \times m}$, so gilt $\bar{x} = VD^+U^T b = A^+$. Man nennt A auch die *Pseudoinverse* von A . Der Satz sagt dann aus, dass

$$\bar{x} = A^+ b \text{ und } \|A\bar{x} - b\|_2^2 = \|(I - AA^+) b\|_2^2$$

Die Pseudoinverse ist die eindeutige Lösung von

$$\min_{X \in \mathbb{R}^{n \times m}} \|AX - I\|_{\text{Frob}}$$

Es gilt

$$\text{Rang}(A) = n \Rightarrow A^+ = (A^T A)^{-1} A^T, \text{ Rang}(A) = m = n \Rightarrow A^+ = A^{-1}$$

BEMERKUNG: In der Praxis muss man $\text{Rang}(A, \varepsilon)$ bei der Definition von A^+ nehmen.

BEMERKUNG: Eine numerisch stabile Berechnung der Singulärwertzerlegung ist sehr aufwendig.

Kapitel 5

Lineare Gleichungssysteme (Iterative Verfahren)

Große Gleichungssysteme, $n \gg 1000 = 10^3$, benötigen für das Gaußverfahren viel Speicher $n^2 = 10^6$ und hohen Aufwand ($n^3 = 10^9$ Multiplikationen). Auch Bandssysteme mit $n = 10^6$ und $m = 10^2$ benötigen schon $n \cdot m = 10^8$ Speicherplatz. In der Praxis sind Matrizen jedoch dünn besetzt, d.h. nur 5 – 20 Einträge pro Zeile. Die folgenden iterativen Verfahren benötigen nur so viel Speicher, wie man für A (und *nicht* A^{-1}) braucht. Sie sind nicht exakt, sondern nähern sich iterativ der Lösung. Wir wollen im Folgenden das

- (a) *Jacobi-Verfahren* (Gesamtschrittverfahren)
- (b) *Gauß-Seidel-Verfahren* (Einzelschrittverfahren)

vorstellen.

Sei $A \in \mathbb{R}^{n \times n}$. Wir zerlegen A in $A = D + L + R$:

$$\begin{pmatrix} a_{11} & & & R \\ & \ddots & & \\ & & D & \\ L & & & \ddots \\ & & & & a_{nn} \end{pmatrix}$$

Dabei sind D die Diagonale, L die Einträge unter der Diagonalen und R die Einträge über der Diagonalen. Das Gleichungssystem $Ax = b$ können wir umschreiben als

$$a_{jj}x_j + \sum_{k=1, k \neq j}^n a_{jk}x_k = b_j, j = 1, \dots, n \quad \text{bzw.} \quad Dx + (L + R)x = b$$

Ist D diagonaldominant, so kann man die sukzessive Iteration mit $C = D$ versuchen. Das heißt:

$$f(x) = b - Ax = 0 \Rightarrow x^{k+1} = x^k + Cf(x^k) = x^k + D^{-1}(b - Ax^k) =$$

$$x^k + D^{-1}(b - Ax^k) = x^k + D^{-1}(b - (L + D + R)x^k) = D^{-1}(b - (L + R)x^k) = -D^{-1}(L + R)x^k + D^{-1}b$$

BEMERKUNG: $f'(x) = A$, $C = D^{-1} \approx A^{-1}$

Hieraus folgt das *Jacobi-Verfahren*:

$$x^{k+1} = - \underbrace{D^{-1}(L + R)}_{\text{Jacobimatrix}} x^k + D^{-1}b \tag{5.1}$$

Im *Gauß-Seidel-Verfahren* benutzt man, dass $D + L$ eine untere Dreiecksmatrix ist und mittels Vorwärtseinsetzen lösbar ist. Sei also $C = L + D$. Dann folgt

$$x^{k+1} = x^k + C(b - Ax^k) = (I - CA)x^k + Cb = (I - (L + D)^{-1}(L + D + R))x^k + (L + D)b =$$

$$- \underbrace{(L + D)^{-1}R}_{\text{Gauß-Seidel-Matrix}} x^k + (L + D)^{-1}b$$

Dies ist das *Gauß-Seidel-Verfahren*.

Man spricht allgemein von einem *Splittingverfahren*, wenn man A in $A = B + (A - B)$ zerlegt und statt A nun B invertiert. Das entspricht der sukzessiven Approximation mit $C = -B^{-1}$, d.h.

$$x^{k+1} = x^k + B^{-1}(b - Ax^k) = (I - B^{-1}A)x^k + B^{-1}b = B^{-1}(B - A)x^k + B^{-1}b$$

Dabei gilt dann $f'(x) = b - Ax$, $f'(x) = -A$, $B^{-1} \approx -A^{-1}$. Sei nun

$$g(x) = \underbrace{B^{-1}(B - A)}_{=M} x + \underbrace{B^{-1}b}_{=c} \Rightarrow g(x) - g(y) = M(x - y) \quad (5.2)$$

Dies ist eine Kontraktion, falls $\|M\| < 1$. Die Konvergenz folgt dann aus dem Banachschen Fixpunktsatz. Wenden wir zwei Schritte an, so erhalten wir

$$g^2(x) = g(g(x)) = Mg(x) + c = M(Mx + c) + c = M^2x + Mc + c$$

Dies konvergiert für $\|M^2\| < 1$. Für die Matrix

$$M = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}$$

ist

$$M^2 = \begin{pmatrix} cc0 & 0 \\ 0 & 0 \end{pmatrix}$$

sodass $\|M^2\| < 1$ ein schwächeres Kriterium ist. Dies führt schließlich zu der *Idee*:

$$x^k \mapsto g(x^k) \text{ konvergiert, falls } \lim_{n \rightarrow \infty} (\|M^n\|)^{1/n} < 1 \quad (5.3)$$

Bestimmen wir nun $\lim_{n \rightarrow \infty} (\|M^n\|)^{1/n}$.

Lemma 5.4

Es gilt $\|M^n\|^{1/n} \geq \rho(M)$ und

$$\lim_{n \rightarrow \infty} (\|M^n\|)^{1/n} = \rho(M) \quad (5.5)$$

für jede Matrixnorm, wobei $\rho(M)$ der Spektralradius ist.

Beweis: Falls $\rho(M) = 0$, so ist $M = 0$ - daraus folgt dann die Behauptung. Sei also $\rho(M) > 0$ und $\lambda_1, \dots, \lambda_m$ die Eigenwerte von M mit Eigenvektoren x_1, \dots, x_m . Dann gilt

$$\|M^n\| \cdot \|x_i\| \geq \|M^n x_i\| = \|\lambda_i^n x_i\| = |\lambda_i|^n \|x_i\|, \quad i = 1, \dots, m$$

$$\Rightarrow \|M^n\|^{1/n} \geq |\lambda_i|, \quad i = 1, \dots, m \Rightarrow \|M^n\|^{1/n} \geq \rho(M)$$

$$\Rightarrow \liminf_{n \rightarrow \infty} \|M^n\|^{1/n} \geq \rho(M)$$

beziehungsweise

$$\limsup_{n \rightarrow \infty} \|M^n\|^{1/n} \leq \rho(M)$$

Sei $M_\delta = \frac{1}{\rho(M)+\delta} \cdot M \Rightarrow \rho(M_\delta) < 1$. Wir behaupten, dass

$$\lim_{n \rightarrow \infty} \|M_\delta^n\|^{1/n} \leq 1 \quad (5.6)$$

Sei hierzu $M_\delta = U^{-1}J_\delta U$ die Jordanzerlegung von M_δ , also

$$J_\delta = \begin{pmatrix} J_1 & & 0 \\ & \ddots & \\ 0 & & J_r \end{pmatrix}, \quad J_i = \begin{pmatrix} \alpha_i & 1 & 0 \\ & \ddots & 1 \\ 0 & & \alpha_i \end{pmatrix}$$

mit $|\alpha_i| < 1$. Sei $q = \max_i |\alpha_i| = \frac{\rho(M)}{\rho(M)+\delta} < 1$. Dann gilt

$$\lim_{n \rightarrow \infty} J_i^n = 0$$

komponentenweise und auch bzgl. $\|\cdot\|$. Daher gilt $J_\delta^n \rightarrow 0$ bezüglich $\|\cdot\|$ und somit $\|M_\delta^n\| \leq \|U^{-1}\| \cdot \|J_\delta^n\| \cdot \|U\| \rightarrow 0$, $n \rightarrow \infty$. Damit folgt $\|M_\delta^n\| < 1/2$ für alle $n \geq N$. Somit

$$\limsup_{n \rightarrow \infty} \|M_\delta^n\|^{1/n} \leq \limsup_{n \rightarrow \infty} (1/2)^{1/n} = 1$$

Das ergibt die Zwischenbehauptung. Weiter folgt

$$\limsup_{n \rightarrow \infty} \|M^n\|^{1/n} = (\rho(M) + \delta) \limsup_{n \rightarrow \infty} \|M_\delta^n\|^{1/n} \leq \rho(M) + \delta$$

Da dies für alle δ gilt, folgt

$$\limsup_{n \rightarrow \infty} \|M^n\|^{1/n} \leq \rho(M)$$

Damit folgt die Behauptung. ■

BEMERKUNG: $\rho(B)$ ist *keine* Matrixnorm. Beispielsweise ist

$$\rho\left(\begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}\right) = 0 \quad \text{obwohl} \quad \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \neq \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$$

Für symmetrische $B \in \mathbb{R}^{n \times n}$ ist aber

$$\rho(B) = \|B\|_2 = \sup_{x \neq 0} \frac{\|Bx\|_2}{\|x\|_2} \quad (5.7)$$

Lemma 5.8

Zu jedem $M \in \mathbb{R}^{n \times n}$ und $\delta > 0$ existiert eine Norm $\|\cdot\|$, so dass die induzierte Matrixnorm

$$\rho(M) \leq \|M\| \leq \rho(M) + \delta \quad (5.9)$$

erfüllt.

Beweis: Die Beziehung $\rho(M) \leq \|M\|$ wurde bereits in Lemma 6.1. gezeigt. Für jede reguläre Matrix $C \in \mathbb{R}^{n \times n}$ ist $\|x\| = \|Cx\|_2$ eine Norm. Die induzierte Matrixnorm ist

$$\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|} = \sup_{x \neq 0} \frac{\|CAx\|_2}{\|Cx\|_2} = \sup_{x \neq 0} \frac{\|CAC^{-1}x\|_2}{\|x\|_2} = \|CAC^{-1}\|_2$$

Also gilt

$$\| \|M\| \| = \|CAC^{-1}\|_2$$

Sei nun $M = U^{-1}JU$ die Jordanzerlegung.

Rest des Beweises fehlt noch. ■

Theorem 5.10

Für $M \in \mathbb{R}^{n \times n}$, $x_0 \in \mathbb{R}^n$ betrachte die Iteration

$$x^{k+1} = Mx^k + c, \quad k \in \mathbb{N}_0 \quad (5.11)$$

Die Folge x^k konvergiert genau dann für jeden Startwert $x^0 \in \mathbb{R}^n$ gegen den Fixpunkt \bar{x} von $\bar{x} = M\bar{x} + c$, falls $\rho(M) < 1$. Im Falle der Konvergenz gilt

$$\sup_{x^0 \in \mathbb{R}^n} \limsup_{k \rightarrow \infty} \left(\frac{\|x^k - \bar{x}\|}{\|x^0 - \bar{x}\|} \right)^{1/k} = \rho(M) \quad (5.12)$$

Beweis: Sei $e^k = x^k - \bar{x}$. Dann gilt $e^{k+1} = x^{k+1} - \bar{x} = Mx^k - M\bar{x} = Me^k$. Nun unterscheiden wir mehrere Fälle.

- (a) Ist $\rho(M) < 1$, so existiert nach Lemma 6.2. eine Norm $\| \cdot \|$ mit $\|M\| < 1$. Dann gilt $\| \|e^k\| \| \leq \| \|M\| \| \cdot \| \|e^0\| \| \rightarrow 0$, $k \rightarrow \infty$. Die Folge $(e^k)_k$ konvergiert dann bezüglich $\| \cdot \|$, und da die Normen äquivalent sind, auch bezüglich $\| \cdot \|$.
- (b) Aus der Konvergenz für jeden Startwert folgt für $x^0 = w + x$ mit einem Eigenvektor x zum größten Eigenwert λ , dass $e^0 = w$ und $e^k = M^k w = \lambda^k w$. Dann gilt

$$\| \lambda^k w \| \rightarrow 0, \quad k \rightarrow \infty \Rightarrow |\lambda| < 1 \Rightarrow \rho(M) < 1$$

Außerdem erhalten wir

$$\left(\frac{\|e^k\|}{\|e^0\|} \right)^{1/k} = |\lambda|$$

Für den zweiten Teil des Beweises sei nun $\rho(M) < 1$. Dann konvergiert das Verfahren. Sei weiter $\| \cdot \|$ eine Norm mit

$$\rho(M) \leq \|M\| \leq \rho(M) + \delta$$

Aus der Äquivalenz der Normen erhalten wir dann $c_0 \|x\| \leq \|x\| \leq c_1 \|x\|$ und damit

$$\|e^k\| \leq c_1 \| \|e^k\| \| \leq c_1 \| \|M\| \| \cdot \| \|e^0\| \| \leq \frac{c_1}{c_0} (\rho(M) + \delta)^k \|e^0\|$$

Daraus erhalten wir

$$\begin{aligned} \left(\frac{\|e^k\|}{\|e^0\|} \right)^{1/k} &\leq \left(\frac{c_1}{c_0} \right)^{1/k} (\rho(M) + \delta) \\ \Rightarrow \limsup_{k \rightarrow \infty} \left(\frac{\|e^k\|}{\|e^0\|} \right)^{1/k} &\leq \rho(M) + \delta \end{aligned}$$

Da $\delta > 0$ beliebig war, folgt damit die Behauptung. ■

BEMERKUNG: Falls $\rho(M) \approx 0.99$, so braucht man immer noch ≈ 230 Schritte bis $(\rho(M))^k \leq 10^{-1}$. Wir diskutieren nun einige

Abbruchkriterien

(a) Nach dem Banachschen Fixpunktsatz gilt

$$\|x^k - \bar{x}\| \leq \frac{\|B\|}{1 - \|B\|} \underbrace{\|x^k - \bar{x}\|}_{=\delta x^k \text{ Update}}$$

$$\Rightarrow \frac{\|x^k - \bar{x}\|}{\|x^k\|} \leq \frac{\|B\|}{1 - \|B\|} \cdot \frac{\|\delta x^k\|}{\|x^k\|}$$

Daraus erhalten wir ein mögliches Abbruchkriterium:

$$\frac{\|B\|}{1 - \|B\|} \cdot \frac{\|\delta x^k\|}{\|x^k\|} \leq \varepsilon$$

Dennoch bleibt das Problem, dass $\rho(B)$ bzw. $\|B\|$ abgeschätzt werden muss.

(b) Betrachte das *Residuum* $Ax^k - b$. Zur Erinnerung: Die Gleichung $Ax = b$ führte zur Iteration

$$x^{k+1} = \underbrace{B^{-1}(B - A)}_{=M} x^k + \underbrace{B^{-1}b}_{=c}$$

$$\Rightarrow e^k = x^k - \bar{x} = A^{-1}(Ax^k - A\bar{x}) = A^{-1}(Ax^k - b)$$

$$\Rightarrow \|e^k\| \leq \|A^{-1}\| \cdot \|Ax^k - b\| \leq \text{cond}(A) \frac{1}{\|A\|} \|Ax^k - b\|$$

Daraus folgt

$$\frac{\|e^k\|}{\|x^k\|} \leq \text{cond}(A) \frac{1}{\|A\| \cdot \|x^k\|} \|Ax^k - b\| \leq \text{cond}(A) \frac{\|Ax^k - b\|}{\|b\|}$$

Allerdings haben wir hier zwei Nachteile:

- Ax^k muss berechnet werden.
- $\text{cond}(A)$ unbekannt.

BEMERKUNG: $Ax = b$ führte zum iterativen Verfahren

$$x^{k+1} = \underbrace{B^{-1}(B - A)}_{=M} x^k + \underbrace{B^{-1}b}_{=c}$$

Dabei haben wir zwei *Ziele*:

- (a) *gute Konvergenz*. Hierfür brauchen wir $\rho(M) \ll 1$.
- (b) $Mx^k = (I - B^{-1}A)x^k$ soll einfach berechenbar sein, d.h. $B^{-1}Ax^k$ soll einfach berechenbar sein.

Diese zwei Ziele sind konträr.

Zu (a): Optimal ist $B = A \Rightarrow M = 0 \Rightarrow \rho(M) = 0$.

Zu (b): Einfach ist $B = D$, jedoch ist oft $\rho(I - D^{-1}A) \approx 1$. Dann konvergiert das Verfahren aber langsam.

Beim Jacobi- und Gauß-Seidel-Verfahren ist (b) gut erfüllt. Dies geht aber auf Kosten von (a) - der Konvergenzgeschwindigkeit. Nun jedoch zurück zum Jacobi- und Gauß-Seidel-Verfahren:

Theorem 5.13

Sei $A \in \mathbb{R}^{n \times n}$ (bzgl. der Zeilen) strikt diagonaldominant, d.h.

$$\sum_{k=1, k \neq j}^n |a_{jk}| < |a_{jj}|, j = 1, \dots, n \quad (5.14)$$

so ist mit der Jacobimatrix $J = -D^{-1}(L + R)$ und der Gauß-Seidel-Matrix $H = -(D + R)^{-1}R$

$$\rho(J) < 1 \text{ und } \rho(H) < 1$$

Folglich konvergieren das Jacobi- und das Gauß-Seidel-Verfahren.

Beweis: Sei λ der größte Eigenwert von J mit dem Eigenvektor v , $\|v\|_\infty = 1$. Sei μ der größte Eigenwert von H mit Eigenvektor w mit $\|w\|_\infty = 1$. Dann gilt

$$|\lambda| \leq \|J\|_\infty = \|D^{-1}(L + R)\|_\infty \leq \max_j \frac{1}{|a_{jj}|} \sum_{k=1, k \neq j}^n |a_{jk}| < 1$$

nach Voraussetzung und

$$\mu w = Hw = -(D + L)^{-1}Rw \Rightarrow (\mu D + \mu L)w = -Rw \Rightarrow \mu w = -D^{-1}(\mu L + R)w$$

Wir erhalten

$$|\mu| \leq \|D^{-1}(\mu L + R)\|_\infty \leq \max_j \frac{1}{|a_{jj}|} \sum_{k=1, k \neq j}^n \max\{|\mu|, 1\} |a_{jk}|$$

$$< \max\{|\mu|, 1\} \text{ nach Voraussetzung} \Rightarrow |\mu| < 1$$

Damit ist der Satz bewiesen. ■

Sei nun A die oben betrachtete Bandmatrix. A ist *nicht* strikt diagonaldominant, aber diagonaldominant. Sie ist jedoch in den ersten und letzten 4 Zeilen strikt diagonaldominant. Dies werden wir nun nutzen:

Definition 5.15

Eine Matrix heißt irreduzibel, falls es keine Permutationsmatrix P gibt mit

$$PAP^T = \begin{pmatrix} \tilde{A}_{11} & 0 \\ \tilde{A}_{21} & \tilde{A}_{22} \end{pmatrix} \quad (5.16)$$

mit $\tilde{A}_{11} \in \mathbb{R}^{p \times p}$, $\tilde{A}_{22} \in \mathbb{R}^{(n-p) \times (n-p)}$ und $0 < p < n$.

Lemma 5.17

$A \in \mathbb{R}^{n \times n}$ ist irreduzibel genau dann, wenn der Graph

$$G(A) = \{\text{Knoten } K_1, \dots, K_n, \text{Kante } \overline{K_j K_k} \Leftrightarrow a_{jk} \neq 0, j, k = 1, \dots, n\} \quad (5.18)$$

Beweis: Die Reduzibilität ist äquivalent zur Existenz von $J, K \subset \{1, \dots, n\}$, $J, K \neq \emptyset$, $J \cap K = \emptyset$ mit $a_{jk} = 0$ für alle $j \in J, k \in K$. ■

Theorem 5.19

(Schwachtes Zeilensummenkriterium). Ist $A \in \mathbb{R}^{n \times n}$ irreduzibel, diagonaldominant und gilt

$$\sum_{k=1, k \neq r}^n |a_{rk}| < |a_{rr}| \quad (5.20)$$

für (mindestens) ein $r \in \{1, \dots, n\}$, so gilt $\rho(J), \rho(H) < 1$.

Beweis: Wir beschränken uns auf $\rho(J) < 1$. Sei λ der größte Eigenwert von J mit Eigenvektor v mit $\|v\|_\infty = 1$.

$$|\lambda| \leq \max_j \frac{1}{|a_{jj}|} \sum_{k=1, k \neq j}^n |a_{jk}| \leq 1, \quad |\mu| \leq 1 \text{ für } H$$

Genauer gilt

$$|v_j| = |\lambda| |v_j| = |[D^{-1}(L+R)v]_j| \leq \frac{1}{|a_{jj}|} \sum_{k \neq j} |a_{jk}| |v_k| \stackrel{j=r}{\Rightarrow} |v_r| < \|v\|_\infty = 1$$

Sei nun i_1, \dots, i_m Kette mit $i_1 = 1$ und $\{i_1, \dots, i_m\} = \{1, \dots, n\}$ (auch doppelte möglich).

Induktionsanfang: $|v_{i_1}| < 1$

Induktionsschritt:

$$|v_{i_m}| \leq \frac{1}{|a_{jj}|} \left(\sum_{k \neq j, i_{m-1}} |a_{jk}| |v_k| + |a_{j, i_{m-1}}| |v_{i_{m-1}}| \right) < 1$$

$$\Rightarrow |v_k| < 1 \text{ für alle } k = 1, \dots, n \Rightarrow \|v\|_\infty < 1$$

Damit erhalten wir einen Widerspruch und somit $\rho(J) < 1$. ■

SOR-Verfahren (successive overrelaxation method)

Da oft $\rho(J), \rho(H) \approx 1 - \varepsilon$, sodass das Jacobi- und Gauß-Seidel-Verfahren zu langsam konvergieren, versucht man, das Verfahren durch *Relaxation* zu verbessern. Bei Gauß-Seidel gilt:

$$Dx^{k+1} = -Lx^{k+1} - Rx^k + b$$

Hieraus wird durch Relaxation der neue Algorithmus des SOR-Verfahrens:

$$D\tilde{x}^{k+1} = -Lx^{k+1} - Rx^k + b \quad \text{und} \quad x^{k+1} = \omega\tilde{x}^{k+1} + (1-\omega)x^k \quad (5.21)$$

Beachte jedoch: *Dies ist eine Zirkeldefinition!* Durch Umformen erhalten wir

$$Dx^{k+1} = \omega D\tilde{x}^{k+1} + (1-\omega)Dx^k = -\omega Lx^{k+1} - \omega Rx^k + \omega b + (1-\omega)Dx^k$$

und somit

$$(D - \omega L)x^{k+1} = ((1-\omega)D - \omega R)x^k + \omega b \quad (5.22)$$

und durch Multiplikation mit $\frac{1}{\omega}$:

$$\left(\frac{1}{\omega}D + L\right)x^{k+1} = \left(\frac{1}{\omega}D - D - R\right)x^k + b$$

Damit ist das SOR-Verfahren ein *Splitting-Verfahren* mit

$$B_\omega := \frac{1}{\omega}D + L \quad \text{mit Iterationsmatrix} \quad M_\omega := B^{-1}(B - A) \quad (5.23)$$

BEMERKUNG: Dies ist auch für $\omega = 0$ wohldefiniert.

Lemma 5.24

Für $A \in \mathbb{R}^{n \times n}$ mit regulärem Diagonalanteil D gilt

$$\rho(M_\omega) \geq |\omega - 1|, \quad \omega \in \mathbb{R} \setminus \{0\} \quad (5.25)$$

Beweis: Wir berechnen die Determinante von M :

$$\begin{aligned} \det(M) &= \det(B^{-1}(B - A)) = \det\left(\left(\frac{1}{\omega}D + L\right)^{-1}\right) \det\left(\frac{1}{\omega}D - D - R\right) \\ &= \omega^n \left(1 - \frac{1}{\omega}\right)^n = (\omega - 1)^n \end{aligned}$$

Damit folgt

$$\rho(M) = \max_i |\lambda_i| \geq \left(\prod_{i=1}^n |\lambda_i|\right)^{1/n} = |\det(M)|^{1/n} = |\omega - 1| \quad \blacksquare$$

BEMERKUNG: Für $\rho(M) < 1$ muss also $\omega \in (0, 2)$ sein.

Theorem 5.26

Für positiv definite (symmetrische) Matrizen $A \in \mathbb{R}^{n \times n}$ gilt $\rho(M_\omega) < 1$ für $\omega \in (0, 2)$. Insbesondere ist das Gauß-Seidel-Verfahren konvergent.

Beweis: Es gilt $A = L + D + L^T$, da A symmetrisch ist. Sei im Folgenden $\omega \in (0, 2)$ sowie $\lambda \in \rho(M_\omega)$ mit Eigenvektor v , d.h. $M_\omega v = \lambda v$. Da $M_\omega = (D + \omega L)^{-1}((1 - \omega)D - \omega L^T)$ folgt

$$((1 - \omega)D - \omega L^T)v = \lambda(D + \omega L)v$$

und damit

$$\omega(D + L^T)v = (1 - \lambda)Dv - \lambda\omega Lv$$

Hieraus erhalten wir

$$\omega Av = \omega(D + L^T)v + \omega Lv = (1 - \lambda)Dv - \lambda\omega Lv + \omega Lv = (1 - \lambda)Dv + \omega(1 - \lambda)Lv$$

und

$$\begin{aligned} \lambda\omega Av &= \lambda\omega(D + L^T)v + \lambda\omega Lv = \lambda\omega(D + L^T)v + (1 - \lambda)Dv - \omega(D + L^T)v = \\ &= (1 - \lambda)(1 - \omega)Dv - \omega(1 - \lambda)L^T v \end{aligned}$$

Damit ergibt sich

$$wv^T Av = (1 - \lambda)v^T Dv + \omega(1 - \lambda)v^T Lv$$

und

$$\lambda wv^T Av = (1 - \lambda)(1 - \omega)v^T Dv - w(1 - \lambda)\underbrace{v^T L^T v}_{=v^T Lv}$$

Letztlich folgt

$$\omega(1 + \lambda)v^T Av = (1 - \lambda)(2 - \omega)v^T Dv$$

Dabei ist $v^T Av > 0$, da A positiv definit. Dann ist aber auch D positiv definit, also $v^T Dv > 0$. Es folgt, dass $\lambda \neq \pm 1$, denn sonst ist genau eine Seite gleich 0. Umstellen liefert

$$\frac{1 + \lambda}{1 - \lambda} = \frac{2 - \omega}{\omega} \frac{v^T Dv}{v^T Av} > 0 \Rightarrow \mu = \frac{1 + \lambda}{1 - \lambda} > 0$$

$$\Rightarrow \lambda = \frac{\mu - 1}{\mu + 1} \in (-1, +1) \Rightarrow |\lambda| < 1 \quad \blacksquare$$

Die qualitativen Aussagen lassen sich für spezielle Matrizen verschärfen:

Definition 5.27

Eine Matrix $A \in \mathbb{R}^{n \times n}$ mit $A = L + D + R$ heißt konsistent geordnet, falls für alle $\alpha \in \mathbb{C}$

$$\sigma(D^{-1}(\alpha L + \alpha^{-1}R)) = \sigma(D^{-1}(L + R)) \quad (5.28)$$

gilt (unabhängig von α mit der Jacobi-Matrix $J = D^{-1}(L + R)$).

BEMERKUNG:

- (a) Tridiagonalmatrizen sind stets konsistent geordnet.
- (b) Die Bandmatrix des obigen Beispiels ist ebenso konsistent geordnet.

Theorem 5.29

Sei $A \in \mathbb{R}^{n \times n}$ konsistent geordnet und $\omega \in (0, 2)$, so gilt

(a) $\lambda \in \sigma(J) \Rightarrow \lambda \in \sigma(J)$

(b) Gilt $\mu \in \sigma(J)$, so ist jedes λ mit

$$(\lambda + \omega - 1) = (\pm)\lambda^{1/2}\omega\mu$$

ein Eigenwert von M_ω .

(c) Ist $\lambda \neq 0$ mit $\lambda \in \sigma(M_\omega)$, so ist $\mu \in \sigma(J)$.

Beweis: Zu (a): Wir rechnen

$$\sigma(J) = \sigma(D^{-1}(L + R)) = \sigma(D^{-1}(-L - U)) = \sigma(-J) = -\sigma(J)$$

Zu (c): Ist $\lambda \in \sigma(H_\omega)$, $\lambda \neq 0$, so ist $H_\omega v = \lambda v$ äquivalent zu

$$((1 - \omega)I - \omega D^{-1}R)v = \lambda(I - \omega D^{-1}L)v$$

bzw.

$$(\lambda + \omega - 1)v = -\lambda^{-1/2}\omega J(\lambda^{1/2})v$$

Hier ist $\lambda^{1/2}$ eine Wurzel von λ . Dies ist allerdings äquivalent dazu, dass v ein Eigenvektor von $J(\lambda^{1/2})$ ist zum Eigenwert

$$\mu = \frac{\lambda + \omega - 1}{\lambda^{1/2}\omega}$$

(b) folgt völlig analog. ■

Folgerung 5.30

Sei A konsistent geordnet und positiv definit. Dann gilt

$$\rho(H) = \rho(J)^2 \quad (5.31)$$

d.h. das Gauß-Seidel-Verfahren ist doppelt so gut wie das Jacobi-Verfahren.

Beweis: Er folgt einfach aus $\lambda = \lambda^{1/2}\mu \Rightarrow \lambda = \mu^2$. ■

Hieraus lässt sich der optimale Relaxationsparameter ablesen. Es gilt

$$\omega_{\text{opt}} = \operatorname{argmin}_{\omega \in (0,2)} \rho(H_\omega) \quad (5.32)$$

Und zwar gilt für konsistent geordnete, positiv definite A :

$$\rho(H_\omega) = \begin{cases} \omega - 1 & \text{für } \omega_{\text{opt}} \leq \omega \\ \frac{1}{2}(\rho(J)^2\omega \pm \sqrt{(\rho(J)^2\omega^2 - 4(\omega - 1))^2}) & \text{für } \omega \leq \omega_{\text{opt}} \end{cases}$$

[Graphik_para_gera]

Da die Parabel bei ω_{opt} die Steigung $-\infty$ hat, sollte man lieber ω_{opt} leicht *größer* wählen.

Theorem 5.33

Sei A konsistent geordnet und positiv definit. Seien ferner die Eigenwerte von J reell mit $\rho(J) < 1$. Dann gilt:

$$\omega_{\text{opt}} = \frac{2}{q + \sqrt{1 - (\rho(J))^2}}, \quad \rho(H_{\omega_{\text{opt}}}) = \omega_{\text{opt}} - 1 = \frac{1 - \sqrt{1 - (\rho(J))^2}}{1 + \sqrt{1 - (\rho(J))^2}} \quad (5.34)$$

BEISPIELE:

$$(a) \quad (\rho(J))^2 = \rho(H_1) = 0.99 \Rightarrow \sqrt{1 - (\rho(J))^2} = 0.1 \Rightarrow \rho(H_{\omega_{\text{opt}}}) = \frac{1-0.1}{1+0.1} \approx 0.81$$

$$(b) \quad (\rho(J))^2 = \rho(H_1) = 0.84 \Rightarrow \sqrt{1 - (\rho(J))^2} = 0.4 \Rightarrow \rho(H_{\omega_{\text{opt}}}) = \frac{0.6}{1.4} = 0.43$$

Abstiegsverfahren

Im Folgenden betrachten wir primär (symmetrische) positiv definite Matrizen.

Lemma 5.35

Ist $A \in \mathbb{R}^{n \times n}$ positiv definit und symmetrisch, so ist die Lösung von $Ax = b$ gleich dem eindeutigen Minimierer von

$$\mathcal{E}: \mathbb{R}^{n \times n} \rightarrow \mathbb{R}, \quad \mathcal{E}(y) = \frac{1}{2} \langle Ay, y \rangle - \langle b, y \rangle \quad (5.36)$$

Beweis: \mathcal{E} ist konvex. Es gilt

$$\mathcal{E}(y) \geq c_0 \|y\|_2^2 - \|b\|_2 \|y\|_2 \geq (c_0 \|y\|_2 - \|b\|_2) \|y\|_2 \geq c_1$$

\mathcal{E} ist somit nach unten beschränkt. Damit existiert ein Minimierer. Es ist $\mathcal{E} \in C^1$ und weiters

$$(\delta\mathcal{E})(y)(z) = \frac{1}{2} \langle Ay, z \rangle + \frac{1}{2} \langle Az, y \rangle - \langle b, z \rangle = \langle Ay, z \rangle - \langle b, z \rangle$$

da A symmetrisch ist. Beim Minimum \bar{x} gilt:

$$\langle A\bar{x}, z \rangle - \langle b, z \rangle = 0 \quad \forall z$$

Damit ergibt sich $A\bar{x} = b$, und da A regulär ist, $x = \bar{x}$. Das ist aber die Behauptung. ■

BEMERKUNG: Sei x eine Lösung von $Ax = b$ (bzw. Minimierer von \mathcal{E}). Dann gilt für alle y :

$$\mathcal{E}(y) - \mathcal{E}(x) = 1/2 \langle ay, y \rangle - 1/2 \langle Ax, x \rangle - \langle b, y \rangle + \langle b, x \rangle = 1/2 \langle Ay, y \rangle - 1/2 \langle Ax, x \rangle + \langle b, x - y \rangle \quad (5.37)$$

$$= \frac{1}{2} \left(\langle Ay, Ay \rangle + \langle Ax, x \rangle - 2 \langle Ax, y \rangle \right) = \frac{1}{2} \langle A(y - x), y - x \rangle.$$

Der Gradient von \mathcal{E} ist gegeben durch

$$\langle \nabla \mathcal{E}(y), z \rangle = \delta \mathcal{E}(y)(z) =$$

$$\langle Ay, z \rangle - \langle b, z \rangle \Rightarrow \nabla \mathcal{E}(y) = Ay - b$$

Dies ist die Richtung des *steilsten Abstiegs*. Die *Idee von Abstiegsverfahren* ist die Folgende:

- (a) Bestimme die Abstiegsrichtung r^k so, dass \mathcal{E} in Richtung r^k kleiner wird.
- (b) Setze $x^{k+1} = x^k + \alpha_r r^k$ mit einem Faktor $\alpha_r \in \mathbb{R}$ so, dass $\alpha \mapsto \mathcal{E}(x^k + \alpha r^k)$ minimiert wird.

Anhand von (b) bestimmen wir nun das *Abstiegsverfahren*. Sei $h(\alpha) = \mathcal{E}(x^k + \alpha r^k)$. Das Minimum erfüllt $h'(\alpha) = 0$, und daher gilt

$$\begin{aligned} 0 = h'(\alpha) &= \delta \mathcal{E}(x^k + \alpha r^k)(r^k) = \langle A(x^k + \alpha r^k) - b, r^k \rangle = \langle Ax^k - b, r^k \rangle + \alpha \langle Ar^k, r^k \rangle \\ &\Rightarrow \alpha = -\frac{\langle Ax^k - b, r^k \rangle}{\langle Ar^k, r^k \rangle} = -\frac{\langle g^k, r^k \rangle}{\langle Ar^k, r^k \rangle} \end{aligned}$$

mit Gradient $g^k = Ax^k - b$. Also erhalten wir für das allgemeine *Abstiegsverfahren*:

$$g^k = Ax^k - b \quad \text{und} \quad \alpha_k = -\frac{\langle g^k, r^k \rangle}{\langle Ar^k, r^k \rangle} \quad \text{und} \quad x^{k+1} = x^k - \alpha_k r^k \quad (5.38)$$

Wir diskutieren nun Gradientenverfahren.

Theorem 5.39

Sei $A \in \mathbb{R}^{n \times n}$ symmetrisch und positiv definit, so konvergiert das Gradientenverfahren, d.h. Startwert $x_0 \in \mathbb{R}^n$,

$$g^k = Ax^k - b \quad \text{und} \quad \alpha_k = -\frac{\langle g^k, r^k \rangle}{\langle Ar^k, r^k \rangle} \quad \text{und} \quad x^{k+1} = x^k - \alpha_k r^k \quad (5.40)$$

gegen die Lösung von $Ax = b$.

BEMERKUNG: Falls $\langle Ag^k, g^k \rangle = 0$, so gilt $g^k = 0$, also $Ax^k - b = 0$. Damit ist die Lösung schon gefunden.

Beweis Beweis des Satzes.: Sei x die Lösung von $Ax = b$. Dann gilt

$$\mathcal{E}(x^k) - \mathcal{E}(x) = 1/2 \langle A(x^k - x), x^k - x \rangle = 1/2 \langle e^k, e^k \rangle$$

$$\Rightarrow \mathcal{E}(x^k) - \mathcal{E}(x^{k+1}) = (\mathcal{E}(x^k) - \mathcal{E}(x)) - (\mathcal{E}(x^{k+1}) - \mathcal{E}(x)) = -1/2 \alpha_k^2 \langle Ag^k, g^k \rangle + \alpha_k \langle Ae^k, g^k \rangle$$

Nach der Definition von α_k gilt

$$\mathcal{E}(x^k) - \mathcal{E}(x^{k+1}) = 1/2 \frac{\langle g^k, g^k \rangle^2}{\langle Ag^k, g^k \rangle} \Rightarrow$$

$$\frac{\mathcal{E}(x^k) - \mathcal{E}(x^{k+1})}{\mathcal{E}(x^k) - \mathcal{E}(x)} = \frac{\langle g^k, g^k \rangle^2}{\langle Ag^k, g^k \rangle \cdot \langle Ae^k, e^k \rangle} = \frac{\langle g^k, g^k \rangle^2}{\langle Ag^k, g^k \rangle \cdot \langle A^{-1}g^k, g^k \rangle}$$

A ist positiv definit. Damit gilt $\lambda \langle \zeta, \zeta \rangle \leq \langle A\zeta, \zeta \rangle \leq \Lambda \langle \zeta, \zeta \rangle$. Somit haben wir

$$\frac{1}{\Lambda} \langle \zeta, \zeta \rangle \leq \langle A^{-1}\zeta, \zeta \rangle \leq \frac{1}{\lambda} \langle \zeta, \zeta \rangle$$

Es folgt hieraus

$$\begin{aligned} \frac{\mathcal{E}(x^k) - \mathcal{E}(x^{k+1})}{\mathcal{E}(x^k) - \mathcal{E}(x)} &\geq \frac{\lambda}{\Lambda} \\ \Rightarrow \mathcal{E}(x^{k+1}) - \mathcal{E}(x) &\leq -\frac{\lambda}{\Lambda}(\mathcal{E}(x^k) - \mathcal{E}(x)) \leq \underbrace{\left(1 - \frac{\lambda}{\Lambda}\right)}_{=q < 1} (\mathcal{E}(x^k) - \mathcal{E}(x)) \\ &\Rightarrow \mathcal{E}(x^k) - \mathcal{E}(x) \leq q^k (\mathcal{E}(x^0) - \mathcal{E}(x)) \end{aligned}$$

Dies bedeutet Konvergenz linearer Ordnung. Weiters gilt

$$\begin{aligned} \mathcal{E}(x^k) - \mathcal{E}(x) &\rightarrow 0, \quad k \rightarrow \infty, \quad 1/2 \langle Ae^k, e^k \rangle \geq \frac{\lambda}{2} \|e^k\|_2^2 \\ \Rightarrow \|e^k\|_2^2 &\leq \frac{2}{\lambda} q^k (\mathcal{E}(x^0) - \mathcal{E}(x)) \rightarrow 0 \end{aligned} \quad \blacksquare$$

BEMERKUNG: Sei $\langle x, y \rangle_A = \langle x, Ay \rangle$ mit symmetrischem, positiv definiten A . Dann ist $\langle \cdot, \cdot \rangle_A$ ein Skalarprodukt mit Norm $\|x\|_A = \sqrt{\langle x, x \rangle_A} = \sqrt{\langle x, Ax \rangle}$. Wir haben gezeigt, dass

$$\|e^k\|_A^2 = \langle e^k, Ae^k \rangle = 2(\mathcal{E}(x^k) - \mathcal{E}(x^0))$$

Hieraus folgt $\|e^k\|_A^2 \leq q^k \|e^0\|_A^2$.

Wir wollen im Folgenden die Konvergenzgeschwindigkeit verbessern.

Lemma 5.41

(Lemma von Kantorowitsch) Sei A symmetrisch und positiv definit mit $\lambda = \min_i \lambda_i$ und $\Lambda = \max_i \lambda_i$.

Dann gilt für alle $\zeta \neq 0$:

$$\frac{\langle \zeta, \zeta \rangle^2}{\langle A\zeta, \zeta \rangle \langle A^{-1}\zeta, \zeta \rangle} \geq \frac{4\lambda\Lambda}{(\Lambda + \lambda)^2} \quad (5.42)$$

BEMERKUNG: Dies ist besser als $\frac{\lambda}{\Lambda}$.

Beweis: Da A positiv definit ist, existiert eine orthogonale Matrix Q mit $A = Q^T D Q$ und $D = \text{diag}(\lambda_1, \dots, \lambda_n)$.

Somit gilt

$$I = \frac{\langle \zeta, \zeta \rangle^2}{\langle A\zeta, \zeta \rangle \langle A^{-1}\zeta, \zeta \rangle} \stackrel{z=Q\zeta, A^{-1}=Q^T D^{-1}Q}{=} \frac{\langle z, z \rangle^2}{\langle Dz, z \rangle \langle D^{-1}z, z \rangle} \text{ für } z \neq 0$$

Definiere nun $\alpha_i = \frac{z_i^2}{\|z\|_2^2}$, so gilt $\sum_i \alpha_i = 1$, $\alpha_i \geq 0$ und

$$F(z) = \frac{1}{\left(\sum_i \lambda_i \alpha_i\right) \left(\sum_i \frac{1}{\lambda_i} \alpha_i\right)}$$

Betrachte die Punkte $P_i = \left(\lambda_i, \frac{1}{\lambda_i}\right)$.

[Graphik Cucumber_Kurve]

Setze $\bar{\lambda} = \sum_i \lambda_i \alpha_i$. Da $\lambda \mapsto \frac{1}{\lambda}$, liegen alle P_1, \dots, P_n unterhalb der Geraden $\overline{P_1 P_n}$, die beschrieben wird durch

$$\lambda \mapsto \frac{1}{\lambda_1} + \frac{\frac{1}{\lambda_n} - \frac{1}{\lambda_1}}{\lambda_n - \lambda_1} (\lambda - \lambda_1) = \frac{\lambda_1 + \lambda_n - \lambda}{\lambda_1 \lambda_n} = g(\lambda)$$

Sei $Q = (\sum_i \lambda_i \alpha_i, \sum_i \frac{1}{\lambda_i} \alpha_i) = \sum \alpha_i P_i$. Dann gilt $Q \in \text{conv}(P_1, \dots, P_n)$ und damit

$$\begin{aligned} \sum_i \frac{1}{\lambda_i} \alpha_i &\leq g\left(\underbrace{\sum_i \lambda_i \alpha_i}_{=\bar{\lambda}}\right) \leq \frac{\lambda_1 + \lambda_n - \bar{\lambda}}{\lambda_1 \lambda_n} \\ \Rightarrow F(z) &= \frac{1}{\bar{\lambda} \sum \frac{1}{\lambda_i} \alpha_i} \geq \frac{\lambda_1 \lambda_n}{\bar{\lambda} (\lambda_1 + \lambda_n - \bar{\lambda})} \end{aligned}$$

Schließlich erhalten wir

$$F(z) \geq \frac{\lambda_1 \lambda_n}{\left(\frac{\lambda_1 + \lambda_n}{2}\right)^2} = \frac{4\lambda_1 \lambda_n}{(\lambda_1 + \lambda_n)^2}. \quad \blacksquare$$

Theorem 5.43

Für das Gradientenverfahren gilt die Fehlerabschätzung

$$\|e^k\|_A^2 = \|x^k - x\|_A \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^k \|e^0\|_A \quad (5.44)$$

für $k \in \mathbb{N}$. Dabei ist $\kappa = \frac{\Lambda}{\lambda} = \text{cond}_2(A)$.

Beweis: Im vorherigen Satz haben wir gezeigt, dass

$$\frac{\mathcal{E}(x^{k+1}) - \mathcal{E}(x)}{\mathcal{E}(x^k) - \mathcal{E}(x)} = \left(-\frac{\langle g^k, g^k \rangle^2}{\langle Ag^k, g^k \rangle \langle A^{-1}g^k, g^k \rangle} + 1\right) \frac{\mathcal{E}(x^k) - \mathcal{E}(x)}{\mathcal{E}(x^k) - \mathcal{E}(x)}$$

Mit

$$\mathcal{E}(x^k) - \mathcal{E}(x) = \frac{1}{2} \langle e^k, Ax^k \rangle = \frac{1}{2} \|e^k\|_A^2$$

folgt

$$\|e^{k+1}\|_A^2 \leq \left(1 - \frac{\langle g^k, g^k \rangle^2}{\langle Ag^k, g^k \rangle \langle A^{-1}g^k, g^k \rangle}\right) \|e^k\|_A^2 \stackrel{\leq}{\underbrace{\hspace{2cm}}} \text{Kantor.}$$

$$\left(1 - 4\frac{\lambda\Lambda}{(\lambda + \Lambda)^2}\right) \|e^k\|_A^2 = \left(\frac{\lambda - \Lambda}{\lambda + \Lambda}\right)^2 \|e^k\|_A^2$$

$$\Rightarrow \|e^{k+1}\|_A \leq \frac{\Lambda - \lambda}{\Lambda + \lambda} \|e^k\|_A \Rightarrow \|e^k\|_A \leq \left(\frac{\Lambda - \lambda}{\Lambda + \lambda}\right)^k \|e^0\|_A = \left(\frac{\kappa - 1}{\kappa + 1}\right)^k \|e^0\|_A \quad \blacksquare$$

Im Gradientenverfahren gilt

$$g^k = Ax^k - b \quad \text{und} \quad \alpha_k = -\frac{\langle g^k, r^k \rangle}{\langle Ar^k, r^k \rangle} \quad \text{und} \quad x^{k+1} = x^k - \alpha_k r^k$$

$$\Rightarrow g^{k+1} = Ax^{k+1} - b = Ax^k - \alpha_k Ag^k - b = g^k - \alpha_k Ag^k \Rightarrow \langle g^{k+1}, g^k \rangle = \langle g^k, g^k \rangle - \alpha_k \langle Ag^k, g^k \rangle = 0$$

Das bedeutet: *Aufeinanderfolgende Richtungen sind orthogonal.*

[Graphik energy_sketch_niveau_ellipse]

Hier wird als ein Zickzackkurs beschrieben.

CG- Verfahren (conjugate gradient method)

Das Gradientenverfahren *denkt* nur von Schritt zu Schritt. Die Grundidee der CG-Verfahren ist die Folgende:

Finde sukzessive Abstiegsrichtungen d^k , welche $\langle \cdot, \cdot \rangle_A$ -orthogonal sind.

Genauer heißt das: Sei $\{d^1, \dots, d^n\}$ Basis von \mathbb{R}^n , die $\langle \cdot, \cdot \rangle_A$ -orthogonal ist, d.h.

$$\langle d^j, d^k \rangle_A = \delta_{jk} \langle Ad^j, d^k \rangle, \quad j, k = 1, \dots, n$$

Dann gilt für alle

$$y = \sum_{k=1}^n t_k d^k$$

$$\mathcal{E}(z + y) = \frac{1}{2} \langle A(z + y), z + y \rangle - \langle z + y, b \rangle = \frac{1}{2} \langle Az, z \rangle - \langle z, b \rangle + \frac{1}{2} \langle Ay, y \rangle + \langle Az, y \rangle - \langle b, y \rangle =$$

$$= \mathcal{E}(z) + \frac{1}{2} \langle Ay, y \rangle + \langle Az - b, y \rangle =$$

$$= \mathcal{E}(z) + \frac{1}{2} \sum_{j,k=1}^n t_j t_k \langle Ad^j, d^k \rangle + \sum_{j=1}^n t_j \langle Az - b, d^j \rangle = \mathcal{E}(z) + \sum_{j=1}^n \left(\frac{1}{2} \langle Ad^j, d^j \rangle t_j^2 - \langle Az - b, d^j \rangle \right)$$

Sei

$$F_j(t_j) = \frac{1}{2} \langle Ad^j, d^j \rangle t_j^2 - \langle Az - b, d^j \rangle,$$

so ist

$$\mathcal{E}(z + y) = \mathcal{E}(z) + \sum_{j=1}^n F_j(t_j)$$

Damit ist das Problem entkoppelt bzgl. der t_j und man kann nacheinander $F_j(t_j)$ minimieren, um $\mathcal{E}(z + y)$ zu minimieren. Die d^1, \dots, d^n kann man sukzessive aus den Richtungen g^1, \dots, g^n bestimmen (nutze Gram-Schmidt). Wir formulieren nun den *Modellalgorithmus*:

- $x^1 \in \mathbb{R}^n$ Startwert
- $g^1 = \nabla \mathcal{E}(x^1)$
- $d^1 = -g^1$ (steilste Abstiegsrichtung)
- Für $k \geq 2$: $g^k = Ax^k - b = \nabla \mathcal{E}(x^k)$. Falls $g^k = 0 \Rightarrow$ Stopp (bzw. $\|g^k\| \leq \varepsilon$)
- Sonst für $k \geq 2$: Bestimme d^k aus $\text{span}\{g^1, \dots, g^k\}$ mit $\langle d^k, d^j \rangle_A = 0$ für $j = 1, \dots, k - 1$. Wir minimieren in Richtung d^k , also $t_k = \text{argmin} \mathcal{E}(x^k + t_k d^k)$. Somit $x^{k+1} = x^k + t_k d^k$.

Da $d^k \in \text{span}\{g^1, \dots, g^k\} = V_k$, gilt $x^{k+1} \in x^1 + V_k$. Da $t_k = \text{argmin}\mathcal{E}(x^k + t_k d^k)$ folgt

$$(\delta\mathcal{E})(x^k + t_k d^k)(d^k) = 0$$

gilt

$$\begin{aligned} \left\langle \underbrace{A(x^k + t_k d^k) - b}_{=g^{k+1}}, d^k \right\rangle &= 0 \\ \Rightarrow \langle g^k, d^k \rangle &= \langle Ax^k - b, d^k \rangle = -t_k \langle Ad^k, d^k \rangle = 0 \\ \Rightarrow t_k &= -\frac{\langle g^k, d^k \rangle}{\langle Ad^k, d^k \rangle} \end{aligned} \tag{5.45}$$

Ferner ist

$$g^{k+1} = A(x^k + t_k d^k) - b = g^k + t_k Ad^k \tag{5.46}$$

BEMERKUNG: Das CG-Verfahren bricht spätestens nach n Schritten ab, da dann die $\{d^1, \dots, d^n\}$ den ganzen Raum \mathbb{R}^n aufspannen. Dies gilt nur *ohne* Rundungsfehler. Dennoch haben wir schon nach wenigen Schritten $\ll n$ eine gute Approximationslösung.

Konvergenzgeschwindigkeit des CG-Verfahrens

Da $V_k = \text{span}\{d^1, \dots, d^k\} = \text{span}\{g^1, \dots, g^k\}$ und $g^{k+1} = g^k + \alpha_t Ad^k \in \text{span}\{g^k, \text{Aspan}\{g^1, \dots, g^k\}\}$, folgt leicht (z.B. mit Induktion)

$$V_k = \text{span}\{g^1, Ag^1, \dots, A^{k-1}g^1\} = \{p(A)g^1 : p \in \mathcal{P}\}$$

Dies folgt aber auch aus dem Wissen, dass $x^k + t_k d^k$ die Energie auf $x^0 + V_k$ minimiert.

