



Numerik-Vorlesungen
Teil 3

Eigenwerte, Ausgleichsprobleme, Minimierung

Peter Szyler, Horst Hollatz

Letzte Änderung: 20. Januar 2006

Inhaltsverzeichnis

9. Eigenwerte symmetrischer Matrizen	1
9.1. Grundlegende Definitionen und Sätze	1
9.2. Störungstheorie	10
9.3. Das Jacobi-Verfahren	25
9.4. Die Vektor- und Teilraumiteration	34
9.5. Die Inverse Iteration nach Wielandt	47
9.6. Transformationen auf Tridiagonalform	51
9.7. Der Lanczos-Algorithmus	55
9.8. Der QR-Algorithmus	60
9.9. Der QR-Algorithmus für Tridiagonalmatrizen	74
10. Lineare Ausgleichsprobleme	87
10.1. Problemstellung und klassische Lösung	87
10.2. Störungstheorie	94
10.3. Normalgleichungsverfahren	107
10.4. Orthogonalisierungsverfahren	110
10.5. Aufgaben	115
11. Freie Minimierung	119
11.1. Einführung	119
11.1.1. Aufgabenstellung und grundlegende Begriffe	119
11.1.2. Differenzierbarkeit und Richtungsableitung	120
11.1.3. Optimalitätskriterien	126
11.2. Ein Modellalgorithmus	135
11.2.1. Schrittweiten bei glatter Zielfunktion	136
11.2.2. Konvergenz des Modellalgorithmus	142
11.3. Quasi-Newton-Verfahren	150
11.3.1. Gedämpftes und ungedämpftes Newton-Verfahren	150
11.3.2. Verfahren der Oren-Luenberger-Klasse	154
11.3.3. Algorithmen zur Aufdatierung von Zerlegungen	165
11.4. Trust-Region-Verfahren	167

Kapitel 9

Eigenwerte symmetrischer Matrizen

9.1. Grundlegende Definitionen und Sätze

Bekanntlich lassen sich für reelle und komplexe quadratische Matrizen Eigenwerte und Eigenvektoren definieren. Dabei sind die Eigenwerte und Eigenvektoren auch für eine reelle Matrix durchaus komplex.

Es sei $A \in \mathbb{C}^{n \times n}$. Eine Zahl $\lambda \in \mathbb{C}$ heißt **Eigenwert** der Matrix A , falls ein Vektor $x \in \mathbb{C}^n$ mit $x \neq \mathbf{o}$ existiert, so dass $Ax = \lambda x$ gilt. Jeder derartige Vektor heißt **Eigenvektor** zum Eigenwert λ .

Aus $Ax = \lambda x$ folgt $(A - \lambda I)x = \mathbf{o}$. Damit dieses Gleichungssystem nichttriviale Lösungen besitzt, muss der Rang der Matrix $A - \lambda I$ kleiner als n sein, d. h. es muss $\det(A - \lambda I) = 0$ gelten. Damit ergibt sich gleich die nächste Definition. Das Polynom $\varphi(\mu) = \det(A - \mu I)$ heißt **charakteristisches Polynom** der Matrix A . Offensichtlich gilt $\varphi \in \Pi_n$. Die Nullstellen des charakteristischen Polynoms sind gerade die Eigenwerte von A . Das charakteristische Polynom der Matrix A habe die Darstellung:

$$\varphi(\mu) = (-1)^n (\mu - \lambda_1)^{\sigma_1} (\mu - \lambda_2)^{\sigma_2} \dots (\mu - \lambda_k)^{\sigma_k}$$

mit paarweise verschiedenen $\lambda_1, \lambda_2, \dots, \lambda_k$. Dann heißen die Zahlen

$$\sigma_i = \sigma(\lambda_i), i = 1, \dots, k,$$

algebraische Vielfachheiten der Eigenwerte $\lambda_i, i = 1, \dots, k$. Wegen $\text{Grad}(\varphi) = n$ folgt $\sigma_1 + \sigma_2 + \dots + \sigma_k = n$. Die Eigenvektoren zu einem Eigenwert sind offensichtlich nicht eindeutig bestimmt. Sind x und y Eigenvektoren zum Eigenwert λ , so sind auch alle Linearkombinationen $\alpha x + \beta y \neq \mathbf{o}$ Eigenvektor zum Eigenwert λ . Die Menge aller Eigenvektoren zu einem festen Eigenwert bildet gemeinsam mit

dem Nullvektor einen Unterraum des \mathbb{C}^n .

Es sei $\mathbf{A} \in \mathbb{C}^{n \times n}$ und λ_i ein Eigenwert von \mathbf{A} .

$$\mathcal{L}_i = \mathcal{L}(\lambda_i) = \{ \mathbf{x} \in \mathbb{C}^n \mid (\mathbf{A} - \lambda_i \mathbf{I})\mathbf{x} = \mathbf{o} \}$$

heißt **Eigenraum** des Eigenwertes λ_i . Die Zahlen

$$\varrho_i = \varrho(\lambda_i) = \dim(\mathcal{L}_i) = n - \text{rg}(\mathbf{A} - \lambda_i \mathbf{I}), \quad i = 1, \dots, k$$

heißen **geometrische Vielfachheiten** der Eigenwerte $\lambda_i, i = 1, \dots, k$. Geometrische und algebraische Vielfachheit eines Eigenwertes brauchen nicht übereinzustimmen. Das zeigen schon die folgenden einfachen Beispiele.

9.1. Beispiel: Es sei $\mathbf{A} = \lambda \mathbf{I}$. Das charakteristische Polynom von \mathbf{A} lautet

$$\varphi(\mu) = \det(\mathbf{A} - \mu \mathbf{I}) = (\mu - \lambda)^n.$$

λ ist der einzige Eigenwert von \mathbf{A} . Jeder Vektor $\mathbf{x} \in \mathbb{R}^n$ mit $\mathbf{x} \neq \mathbf{o}$ ist Eigenvektor zum Eigenwert λ . Geometrische und algebraische Vielfachheit von λ sind also jeweils gleich n . ♡

9.2. Beispiel: Es sei

$$\mathbf{B} = \begin{pmatrix} \lambda & 1 & & & 0 \\ & \lambda & 1 & & \\ & & \cdot & \cdot & \\ & & & \cdot & \cdot \\ & & & & \lambda & 1 \\ 0 & & & & & \lambda \end{pmatrix} \in \mathbb{R}^{n \times n}.$$

Wir erhalten als charakteristisches Polynom $\varphi(\mu) = \det(\mathbf{B} - \mu \mathbf{I}) = (\mu - \lambda)^n$. Die Zahl λ ist einziger Eigenwert von \mathbf{A} mit der algebraischen Vielfachheit n . Für die geometrische Vielfachheit von λ ergibt sich aber

$$\varrho = n - \text{rg}(\mathbf{B} - \lambda \mathbf{I}) = n - (n - 1) = 1 \quad .$$

λ hat die geometrische Vielfachheit $1 < n$. ♡

Die algebraische und geometrische Vielfachheit eines Eigenwertes unterscheiden sich daher. Wir werden später zeigen, dass die geometrische Vielfachheit eines Eigenwertes nicht größer ist als seine algebraische Vielfachheit. Eine Matrix, die mindestens einen Eigenwert λ mit $\varrho(\lambda) < \sigma(\lambda)$ besitzt, heißt defektiv. Wollen wir Eigenwerte und Eigenvektoren einer Matrix berechnen, so ist einerseits zu untersuchen,

wie diese Größen auf Störungen in der Matrix reagieren, und andererseits, wie die Eigenwerte und Eigenvektoren verschiedener Matrizen miteinander zusammenhängen. Der erste Gesichtspunkt wird Gegenstand des nächsten Abschnittes sein. Zum zweiten Aspekt notieren wir folgenden Satz.

9.3. Satz: Ist $p(\mu) = \gamma_0 + \gamma_1\mu + \dots + \gamma_m\mu^m$ ein beliebiges Polynom, und definiert man für eine Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ die Matrix $p(\mathbf{A})$ durch

$$p(\mathbf{A}) = \gamma_0\mathbf{I} + \gamma_1\mathbf{A} + \dots + \gamma_m\mathbf{A}^m,$$

so besitzt die Matrix $p(\mathbf{A})$ den Eigenvektor \mathbf{x} zum Eigenwert $p(\lambda)$, falls \mathbf{x} Eigenvektor von \mathbf{A} zum Eigenwert λ ist. Insbesondere besitzt $\alpha\mathbf{A}$ den Eigenwert $\alpha\lambda$ und $\mathbf{A} + \tau\mathbf{I}$ den Eigenwert $\lambda + \tau$.

Beweis: Aus $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ folgt $\mathbf{A}^2\mathbf{x} = \mathbf{A}(\mathbf{A}\mathbf{x}) = \lambda\mathbf{A}\mathbf{x} = \lambda^2\mathbf{x}$ und allgemein

$$\mathbf{A}^i\mathbf{x} = \lambda^i\mathbf{x}.$$

Damit erhält man sofort $p(\mathbf{A})\mathbf{x} = p(\lambda)\mathbf{x}$. *

Weiterhin gilt

9.4. Satz: Ist λ Eigenwert der Matrix \mathbf{A} , so ist λ auch Eigenwert der Matrix \mathbf{A}^T und $\bar{\lambda}$ ist Eigenwert von \mathbf{A}^H .

Beweis: Es gilt:

$$\det(\mathbf{A}^T - \lambda\mathbf{I}) = \det((\mathbf{A} - \lambda\mathbf{I})^T) = \det(\mathbf{A} - \lambda\mathbf{I})$$

$$\det(\mathbf{A}^H - \bar{\lambda}\mathbf{I}) = \det((\mathbf{A} - \lambda\mathbf{I})^H) = \overline{\det(\mathbf{A} - \lambda\mathbf{I})}.$$

*

Wichtig sind Transformationen, die das Eigenwertspektrum einer Matrix unverändert lassen. Es sei $\mathbf{A} \in \mathbb{C}^{n \times n}$ eine beliebige Matrix und $\mathbf{T} \in \mathbb{C}^{n \times n}$ sei regulär. Dann bezeichnet man die Abbildung

$$\mathbf{A} \rightarrow \mathbf{B} = \mathbf{T}^{-1}\mathbf{A}\mathbf{T}$$

als **Ähnlichkeitstransformation**. Die Matrizen \mathbf{A} und \mathbf{B} heißen **ähnlich**.

Ist eine Matrix zu einer Diagonalmatrix ähnlich, so heißt sie **diagonalähnlich**. Die Eigenschaften einer Ähnlichkeitstransformation werden in folgendem Satz beschrieben.

9.5. Satz: Die Matrix $A \in \mathbb{C}^{n \times n}$ werde mittels einer regulären Matrix $T \in \mathbb{C}^{n \times n}$ in $B = T^{-1}AT$ transformiert. Dann haben A und B dieselben Eigenwerte mit denselben algebraischen und geometrischen Vielfachheiten.

Beweis: Aus

$$B - \lambda I = T^{-1}AT - \lambda T^{-1}T = T^{-1}(A - \lambda I)T$$

folgt

$$\det(B - \lambda I) = \det(T^{-1}) \det(A - \lambda I) \det(T) = \det(A - \lambda I).$$

Die charakteristischen Polynome von A und B stimmen daher überein. Damit haben A und B dieselben Eigenwerte mit denselben algebraischen Vielfachheiten. Für die geometrische Vielfachheit des Eigenwertes λ_i von B gilt wegen der Regularität von T :

$$\varrho_i^{(B)} = n - \operatorname{rg}(B - \lambda_i I) = n - \operatorname{rg}(T^{-1}(A - \lambda_i I)T) = n - \operatorname{rg}(A - \lambda_i I) = \varrho_i^{(A)}.$$

✱

Viele Verfahren zur Eigenwertberechnung beruhen darauf, dass man versucht, die Matrix durch Ähnlichkeitstransformationen so umzuformen, dass Eigenwerte und eventuell Eigenvektoren leicht ablesbar sind. Eine Möglichkeit zeigt der folgende Satz.

9.6. Komplexe SCHUR-Zerlegung: Zu jeder (n, n) -Matrix $A \in \mathbb{C}^{n \times n}$ existiert eine unitäre Matrix $U \in \mathbb{C}^{n \times n}$, so dass $U^H A U$ eine obere Dreiecksmatrix ist, in deren Diagonale die Eigenwerte von A stehen.

Beweis: Wir beweisen den Satz mittels Induktion über die Dimension n .

Für $n = 1$ ist die Aussage trivial.

Wir nehmen an, dass die Aussage für alle $(n - 1) \times (n - 1)$ -Matrizen gilt. Es sei nun $A \in \mathbb{C}^{n \times n}$ und λ_1 ein Eigenwert mit einem zugehörigen Eigenvektor $\mathbf{v}_1 \neq \mathbf{o}$. O.B.d.A. sei $\|\mathbf{v}_1\|_2 = 1$. Nun existieren weitere $n - 1$ Vektoren $\mathbf{v}_2, \mathbf{v}_3, \dots, \mathbf{v}_n$ finden, so dass $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ eine Orthonormalbasis des \mathbb{C}^n bilden. Die aus diesen Vektoren gebildete Matrix $V = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n)$ ist unitär, und es gilt

$$V^H A V = \begin{pmatrix} \lambda_1 & \mathbf{y}^T \\ \mathbf{o} & \bar{A} \end{pmatrix}, \quad \bar{A} \in \mathbb{C}^{(n-1) \times (n-1)}, \quad \mathbf{y} \in \mathbb{C}^{n-1}.$$

Nach Induktionsvoraussetzung existiert eine unitäre Matrix $W \in \mathbb{C}^{(n-1) \times (n-1)}$ derart, dass $W^H \bar{A} W$ eine obere Dreiecksmatrix ist. Die unitäre Matrix

$$U = V \begin{pmatrix} 1 & \mathbf{o}^T \\ \mathbf{o} & W \end{pmatrix} \in \mathbb{C}^{n \times n}$$

transformiert dann die Matrix A orthogonal ähnlich auf obere Dreiecksform, denn

$$U^H A U = \begin{pmatrix} 1 & \mathbf{o}^T \\ \mathbf{o} & \mathbf{W}^H \end{pmatrix} \begin{pmatrix} \lambda_1 & \mathbf{y}^T \\ \mathbf{o} & \bar{A} \end{pmatrix} \begin{pmatrix} 1 & \mathbf{o}^T \\ \mathbf{o} & \mathbf{W} \end{pmatrix} = \begin{pmatrix} \lambda_1 & \mathbf{y}^T \mathbf{W} \\ \mathbf{o} & \mathbf{W}^H \bar{A} \mathbf{W} \end{pmatrix}.$$

Damit ist der Satz bewiesen. *

Ist die Matrix A reell und besitzt nur reelle Eigenwerte, so ist die Transformationsmatrix U orthogonal wählbar. Man erhält in diesem Falle eine reelle obere Dreiecksmatrix als Ergebnis. Besitzt eine reelle Matrix komplexe Eigenwerte, so ist sie nur mittels unitärer Ähnlichkeitstransformationen auf obere Dreiecksform transformierbar. Strebt man nur eine obere Blockdreiecksmatrix als Ergebnis an, so ist die ganze Transformation im Reellen durchführbar.

9.7. Reelle SCHUR-Zerlegung: *Zu jeder reellen (n, n) -Matrix A existiert eine orthogonale Matrix Q mit*

$$Q^T A Q = \begin{pmatrix} \mathbf{R}_{11} & \mathbf{R}_{12} & \cdots & \mathbf{R}_{1m} \\ \mathbf{O} & \mathbf{R}_{22} & \cdots & \mathbf{R}_{2m} \\ \vdots & \ddots & \ddots & \vdots \\ \mathbf{O} & \cdots & \mathbf{O} & \mathbf{R}_{mm} \end{pmatrix},$$

wobei die Diagonalblöcke \mathbf{R}_{ii} entweder 1×1 -Matrizen oder 2×2 -Matrizen mit einem Paar konjugiert komplexer Eigenwerte sind.

Beweis: Wir führen den Beweis mittels Induktion über die Anzahl k der Paare konjugiert komplexer Eigenwerte von A . Für $k = 0$ besitzt die reelle Matrix A nur reelle Eigenwerte. Die Aussage folgt dann aus dem vorigen Satz. Wir nehmen an, dass eine derartige Zerlegung für alle Matrizen mit höchstens $k - 1$ konjugiert komplexen Eigenwertpaaren gilt. Es sei nun A eine Matrix mit k konjugiert komplexen Eigenwertpaaren. Ein komplexer Eigenwert sei $\lambda = \alpha + i\beta$ mit $\beta \neq 0$. Der Vektor $\mathbf{y} + iz$ mit $\mathbf{y}, z \in \mathbb{R}^n$ sei ein zugehöriger Eigenvektor. Dann gilt $A(\mathbf{y} + iz) = (\alpha + i\beta)(\mathbf{y} + iz)$ und damit $A\mathbf{y} = \alpha\mathbf{y} - \beta z$ und $Az = \beta\mathbf{y} + \alpha z$. In Matrixschreibweise lautet das

$$A(\mathbf{y}, z) = (\mathbf{y}, z) \begin{pmatrix} \alpha & \beta \\ -\beta & \alpha \end{pmatrix}.$$

Die Vektoren \mathbf{y} und z sind linear unabhängig, denn andernfalls wäre z ein reeller Eigenvektor zum komplexen Eigenwert λ , und das wäre ein Widerspruch zu $A \in \mathbb{R}^{n \times n}$. Der Rang der Matrix (\mathbf{y}, z) ist somit 2, und es existiert eine orthogonale Matrix V , so dass

$$(\mathbf{y}, z) = V \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{o} \end{pmatrix}$$

mit einer oberen Dreiecksmatrix $\mathbf{R}_1 \in \mathbb{R}^{2 \times 2}$. Es folgt weiter

$$\begin{aligned} \mathbf{V}^T \mathbf{A} \mathbf{V} \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{o} \end{pmatrix} &= \mathbf{V}^T \mathbf{A}(\mathbf{y}, \mathbf{z}) = \mathbf{V}^T(\mathbf{y}, \mathbf{z}) \begin{pmatrix} \alpha & \beta \\ -\beta & \alpha \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{R}_1 \\ \mathbf{o} \end{pmatrix} \begin{pmatrix} \alpha & \beta \\ -\beta & \alpha \end{pmatrix} = \begin{pmatrix} \bar{\mathbf{R}}_1 \\ \mathbf{o} \end{pmatrix}. \end{aligned}$$

Partitioniert man $\mathbf{V}^T \mathbf{A} \mathbf{V}$ entsprechend

$$\mathbf{V}^T \mathbf{A} \mathbf{V} = \begin{pmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{pmatrix}, \quad \mathbf{B}_{11} \in \mathbb{R}^{2 \times 2},$$

so folgt

$$\begin{pmatrix} \mathbf{B}_{11} \mathbf{R}_1 \\ \mathbf{B}_{21} \mathbf{R}_1 \end{pmatrix} = \begin{pmatrix} \bar{\mathbf{R}}_1 \\ \mathbf{o} \end{pmatrix}.$$

Da \mathbf{R}_1 regulär ist, gilt

$$\mathbf{R}_1^{-1} \mathbf{B}_{11} \mathbf{R}_1 = \begin{pmatrix} \alpha & \beta \\ -\beta & \alpha \end{pmatrix}, \quad \mathbf{B}_{21} = \mathbf{o}.$$

Die Matrix \mathbf{A} wird durch \mathbf{V} orthogonal ähnlich auf

$$\mathbf{V}^T \mathbf{A} \mathbf{V} = \begin{pmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{o} & \mathbf{B}_{22} \end{pmatrix},$$

transformiert. Das Paar konjugiert komplexer Eigenwerte (λ und $\bar{\lambda}$) von \mathbf{A} ist auch Eigenwertpaar von \mathbf{B}_{11} . Die restlichen Eigenwerte sind Eigenwerte von \mathbf{B}_{22} . Auf diese Matrix ist die Induktionsvoraussetzung anwendbar.

Damit folgt die Behauptung. *

Wir wollen nun noch den angekündigten Satz über den Zusammenhang zwischen den algebraischen und geometrischen Vielfachheiten von Eigenwerten einer Matrix beweisen.

9.8. Satz: *Es sei λ ein Eigenwert der Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ mit der algebraischen Vielfachheit $\sigma(\lambda)$ und der geometrischen Vielfachheit $\varrho(\lambda)$. Dann gilt*

$$1 \leq \varrho(\lambda) \leq \sigma(\lambda) \leq n.$$

Beweis: Die Gültigkeit der Ungleichungen $1 \leq \rho(\lambda)$ und $\sigma(\lambda) \leq n$ ist offensichtlich. Es bleibt nur $\rho(\lambda) \leq \sigma(\lambda)$ zu zeigen. Es sei dazu $\rho = \rho(\lambda)$. Dann existieren ρ linear unabhängige Eigenvektoren zum Eigenwert λ . Diese seien $\mathbf{v}_1, \dots, \mathbf{v}_\rho$. Wir ergänzen diese Vektoren durch Hinzunahme weiterer $n - \rho$ Vektoren zu einer Basis des \mathbb{C}^n und fassen alle Vektoren zu einer regulären Matrix $\mathbf{V} \in \mathbb{C}^{n \times n}$ zusammen. Wegen $\mathbf{V}\mathbf{e}_i = \mathbf{v}_i$ und $\mathbf{V}^{-1}\mathbf{v}_i = \mathbf{e}_i$ folgt für $i = 1, \dots, \rho$

$$\mathbf{V}^{-1}\mathbf{A}\mathbf{V}\mathbf{e}_i = \mathbf{V}^{-1}\mathbf{A}\mathbf{v}_i = \lambda\mathbf{V}^{-1}\mathbf{v}_i = \lambda\mathbf{e}_i.$$

Damit hat die Matrix $\mathbf{V}^{-1}\mathbf{A}\mathbf{V}$ die Gestalt

$$\mathbf{V}^{-1}\mathbf{A}\mathbf{V} = \begin{pmatrix} \lambda & 0 & * & \cdots & * \\ & \ddots & \vdots & & \vdots \\ 0 & \lambda & * & \cdots & * \\ 0 & \cdots & 0 & * & \cdots & * \\ \vdots & & \vdots & \vdots & & \vdots \\ 0 & \cdots & 0 & * & \cdots & * \end{pmatrix} = \begin{pmatrix} \lambda\mathbf{I} & \mathbf{B} \\ \mathbf{o} & \mathbf{C} \end{pmatrix} \quad \lambda\mathbf{I} \in \mathbb{C}^{\rho \times \rho}.$$

Das charakteristische Polynom von $\mathbf{V}^{-1}\mathbf{A}\mathbf{V}$ und damit auch das von \mathbf{A} lautet

$$\varphi(\mu) = \det(\mathbf{A} - \mu\mathbf{I}) = \det(\mathbf{V}^{-1}\mathbf{A}\mathbf{V} - \mu\mathbf{I}) = (\lambda - \mu)^\rho \det(\mathbf{C} - \mu\mathbf{I}).$$

Die Größe λ ist daher mindestens ρ -fache Nullstelle des charakteristischen Polynoms von \mathbf{A} . Die algebraische Vielfachheit $\sigma(\lambda)$ des Eigenwertes λ ist daher niemals kleiner als seine geometrische Vielfachheit $\rho(\lambda)$. *

Im Falle hermitescher Matrizen folgt aus der komplexen SCHUR-Zerlegung:

9.9. Satz: *Es sei $\mathbf{A} \in \mathbb{C}^{n \times n}$ eine hermitesche Matrix ($\mathbf{A}^H = \mathbf{A}$). Dann gelten folgende Aussagen:*

- (i) *\mathbf{A} besitzt nur reelle Eigenwerte.*
- (ii) *Eigenvektoren zu verschiedenen Eigenwerten sind orthogonal zueinander.*
- (iii) *Algebraische und geometrische Vielfachheiten der Eigenwerte stimmen überein.*

Beweis: Nach Satz 9.6 existiert eine unitäre Matrix $\mathbf{U} \in \mathbb{C}^{n \times n}$, so dass $\mathbf{\Lambda} = \mathbf{U}^H \mathbf{A} \mathbf{U}$ eine obere Dreiecksmatrix ist. Wegen $\mathbf{A}^H = \mathbf{A}$ folgt

$$\mathbf{\Lambda}^H = (\mathbf{U}^H \mathbf{A} \mathbf{U})^H = \mathbf{U}^H \mathbf{A}^H \mathbf{U} = \mathbf{\Lambda}.$$

Damit muss $\mathbf{\Lambda}$ reelle Diagonalmatrix sein, in deren Diagonale die reellen Eigenwerte von \mathbf{A} stehen. Die Spalten von \mathbf{U} sind die zugehörigen Eigenvektoren. Diese sind

paarweise orthogonal und bilden eine Basis des \mathbb{C}^n . Damit ist die Summe der geometrischen Vielfachheiten der Eigenwerte gleich n . Mit Satz 9.8 folgt dann, dass die geometrischen Vielfachheiten gleich den algebraischen Vielfachheiten sind. *

Im reellen Falle erhalten wir aus dem letzten Satz:

9.10. Satz: *Es sei A eine symmetrische Matrix ($A^T = A$). Dann gelten folgende Aussagen:*

- (i) *A besitzt nur reelle Eigenwerte und die Eigenvektoren sind reell wählbar.*
- (ii) *Eigenvektoren zu verschiedenen Eigenwerten sind orthogonal zueinander.*
- (iii) *Algebraische und geometrische Vielfachheiten der Eigenwerte stimmen überein.*

Hermitesche Matrizen sind daher durch unitäre Ähnlichkeitstransformation auf Diagonalform transformierbar. Es gibt noch weitere Matrizen, die ebenfalls diagonalähnlich sind. Das sind gerade alle nichtdefektiven Matrizen: Algebraische und geometrische Vielfachheiten der Eigenwerte stimmen überein. Bevor wir die Diagonalähnlichkeit der nichtdefektiven Matrizen zeigen, benötigen wir noch

9.11. Satz: *Eigenvektoren zu paarweise verschiedenen Eigenwerten sind linear unabhängig.*

Beweis: Es seien $\lambda_1, \dots, \lambda_s$ paarweise verschiedene Eigenwerte der Matrix $A \in \mathbb{C}^{n \times n}$ mit den zugehörigen Eigenvektoren $\mathbf{x}_1, \dots, \mathbf{x}_s \in \mathbb{C}^n$. Wir nehmen an, dass die Vektoren $\mathbf{x}_1, \dots, \mathbf{x}_s$ linear abhängig sind. O.B.d.A. seien die Vektoren

$$\mathbf{x}_1, \dots, \mathbf{x}_k \quad (k < s)$$

linear unabhängig. Die restlichen Vektoren lassen sich dann als Linearkombinationen dieser linear unabhängigen Vektoren darstellen. Es sei zum Beispiel

$$\mathbf{x}_s = \sum_{i=1}^k \alpha_i \mathbf{x}_i.$$

Es folgt

$$A\mathbf{x}_s = \sum_{i=1}^k \alpha_i A\mathbf{x}_i = \sum_{i=1}^k \alpha_i \lambda_i \mathbf{x}_i.$$

Andererseits gilt aber

$$A\mathbf{x}_s = \lambda_s \mathbf{x}_s = \sum_{i=1}^k \alpha_i \lambda_s \mathbf{x}_i.$$

Daraus folgt

$$\sum_{i=1}^k (\lambda_i - \lambda_s) \alpha_i \mathbf{x}_i = 0.$$

Wegen der linearen Unabhängigkeit der Vektoren $\mathbf{x}_1, \dots, \mathbf{x}_k$ verschwinden die Koeffizienten $(\lambda_i - \lambda_s) \alpha_i$ für $i = 1, \dots, k$ sämtlich. Da $\lambda_i \neq \lambda_s$, $i = 1, \dots, k$, folgt $\alpha_i = 0$, $i = 1, \dots, k$, und damit $\mathbf{x}_s = 0$ als Widerspruch. Damit ist die Annahme $k < s$ falsch. Die Vektoren $\mathbf{x}_1, \dots, \mathbf{x}_s$ sind linear unabhängig. *

Mit Hilfe dieses Satzes lässt sich nun die Diagonalähnlichkeit der nichtdefektiven Matrizen zeigen.

9.12. Satz: *Zu jeder nichtdefektiven komplexen (n, n) -Matrix \mathbf{A} gibt es eine komplexe, reguläre (n, n) -Matrix \mathbf{X} mit*

$$\mathbf{X}^{-1} \mathbf{A} \mathbf{X} = \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n) \in \mathbb{C}^{n \times n}.$$

Falls \mathbf{A} reell ist mit nur reellen Eigenwerten, ist \mathbf{X} reell wählbar.

Beweis: Da algebraische und geometrische Vielfachheiten der Eigenwerte übereinstimmen, existieren zu jedem Eigenwert λ_i , $i = 1, \dots, s$ genau σ_i linear unabhängige Eigenvektoren. Nach Satz 9.11 existieren damit $\sigma_1 + \sigma_2 + \dots + \sigma_s = n$ linear unabhängige Eigenvektoren der Matrix \mathbf{A} . Diese fassen wir zur Matrix $\mathbf{X} \in \mathbb{C}^{n \times n}$ zusammen. \mathbf{X} ist offensichtlich regulär und es gilt $\mathbf{A} \mathbf{X} = \mathbf{X} \Lambda$. *

Im Falle defektiver Matrizen ist nicht mehr zu erwarten, dass sie durch eine Ähnlichkeitstransformation auf Diagonalgestalt transformierbar sind. Hier gilt nur noch der folgende fundamentale Satz, den wir ohne Beweis angeben.

9.13. Satz: *Die Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ habe die paarweise verschiedenen Eigenwerte*

$$\lambda_1, \dots, \lambda_s$$

mit den algebraischen bzw. geometrischen Vielfachheiten

$$\sigma_i, \varrho_i, i = 1, \dots, s.$$

Zu jedem Eigenwert λ_i gibt es dann ϱ_i natürliche Zahlen $\nu_j^{(i)}$, $j = 1, \dots, \varrho_i$ mit

$$\sigma_i = \nu_1^{(i)} + \nu_2^{(i)} + \dots + \nu_{\varrho_i}^{(i)}$$

9.15. Satz: *Es sei*

$$p(z) = z^n + a_{n-1}z^{n-1} + \cdots + a_1z + a_0$$

ein Polynom mit den Nullstellen $\lambda_1, \dots, \lambda_n$. Dann gibt es zu jedem $\varepsilon > 0$ ein positives $\delta = \delta(\varepsilon)$ mit der Eigenschaft: Ist

$$q(z) = z^n + b_{n-1}z^{n-1} + \cdots + b_1z + b_0$$

ein Polynom mit $|b_j - a_j| \leq \delta$ für $j = 0, \dots, n-1$, so sind die Nullstellen μ_1, \dots, μ_n von q so nummerierbar, dass $|\mu_i - \lambda_i| \leq \varepsilon$ für $i = 1, \dots, n$ gilt.

Beweis: Es seien $\hat{\lambda}_1, \dots, \hat{\lambda}_m$ die paarweise verschiedenen Nullstellen von p . O.B.d.A. sei

$$\varepsilon < \frac{1}{2} \min \left\{ |\hat{\lambda}_i - \hat{\lambda}_j| \mid 1 \leq i < j \leq m \right\}.$$

Für jedes $i = 1, \dots, m$ definieren wir die offene Kreisscheibe

$$D_i = \{ z \in \mathbb{C} \mid |z - \hat{\lambda}_i| < \varepsilon \}$$

mit dem Rand γ_i . Weiterhin sei

$$m_i = \min \{ |p(z)| \mid z \in \gamma_i \}, \quad M_i = \max \left\{ \sum_{j=0}^{n-1} |z|^j \mid z \in \gamma_i \right\}.$$

Wegen

$$\varepsilon < \frac{1}{2} \min \{ |\hat{\lambda}_i - \hat{\lambda}_j| : 1 \leq i < j \leq m \}$$

besitzt p keine Nullstellen auf γ_i , so dass $m_i > 0$. Wählt man nun $\delta > 0$ so klein, dass $M_i \delta < m_i$ für $i = 1, \dots, m$ ist, so gilt für ein Polynom

$$q(z) = z^n + b_{n-1}z^{n-1} + \cdots + b_1z + b_0, \quad |b_j - a_j| \leq \delta, \quad j = 0, \dots, n-1$$

$$|q(z) - p(z)| \leq \sum_{j=0}^{n-1} |b_j - a_j| |z|^j \leq \delta \sum_{j=0}^{n-1} |z|^j \leq \delta M_i < m_i \leq |p(z)|, \quad z \in \gamma_i.$$

Aus dem Satz von ROUCHÉ folgt, dass p und q gleiche Anzahl von Nullstellen in D_i haben, nämlich gerade die Vielfachheit von $\hat{\lambda}_i$. Damit sind diese so nummerierbar, dass für $i = 1, \dots, n$ jeweils $|\mu_i - \lambda_i| \leq \varepsilon$ gilt. *

Unmittelbar aus diesem Satz folgt:

9.16. Satz: Ist $A \in \mathbb{C}^{n \times n}$ eine Matrix mit den Eigenwerten $\lambda_1, \dots, \lambda_n$ und $\|\circ\|$ eine Matrixgrenznorm, so gibt es zu jedem $\varepsilon > 0$ ein $\delta = \delta(\varepsilon) > 0$ mit der Eigenschaft: Ist $\delta A \in \mathbb{C}^{n \times n}$ eine Matrix mit $\|\delta A\| \leq \delta$, so sind die Eigenwerte μ_1, \dots, μ_n von $A + \delta A$ so nummerierbar, dass $|\mu_j - \lambda_j| \leq \varepsilon$, $j = 1, \dots, n$, gilt.

Aus der stetigen Abhängigkeit der Eigenwerte von den Matrixelementen folgt der nächste Satz, der eine Lokalisierung der Eigenwerte einer Matrix ermöglicht.

9.17. GERSCHGORIN: Es sei $A = (a_{ij}) \in \mathbb{C}^{n \times n}$. Weiterhin seien durch

$$G_i = \left\{ z \in \mathbb{C} \mid |z - a_{ii}| \leq r_i \right\}, \quad r_i = \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad i = 1, \dots, n$$

die GERSCHGORIN-Kreise definiert. Dann gilt:

- (i) Ist λ ein Eigenwert von A , so ist $\lambda \in \cup_{i=1}^n G_i$. Alle Eigenwerte von A sind also in der Vereinigung aller GERSCHGORIN-Kreise enthalten.
- (ii) Hat die Vereinigung \hat{G} von m Kreisen G_i einen leeren Durchschnitt mit den restlichen $n - m$ Kreisen, so enthält \hat{G} genau m Eigenwerte von A , wobei jeder Eigenwert entsprechend seiner algebraischen Vielfachheit zu zählen ist.

Beweis: Es sei λ ein Eigenwert von A . Ist $\lambda = a_{ii}$ für ein $i \in \{1, \dots, n\}$, so gilt die erste Behauptung offensichtlich.

Es sei $\lambda \neq a_{ii}$, $i = 1, \dots, n$, $D = \text{diag}(a_{11}, a_{22}, \dots, a_{nn})$; die Matrix $(\lambda I - D)$ ist regulär. Die Matrix $(\lambda I - D)^{-1}(A - D)$ hat den Eigenwert 1, denn aus $Ax = \lambda x$ folgt $(\lambda I - D)^{-1}(A - D)x = x$. Für eine beliebige Matrixgrenznorm gilt somit:

$$1 \leq \varrho((\lambda I - D)^{-1}(A - D)) \leq \|(\lambda I - D)^{-1}(A - D)\|.$$

Wählt man als Matrixnorm die ∞ -Norm, also die maximale Zeilenbetragssumme, so erhält man:

$$\begin{aligned} 1 \leq \|(\lambda I - D)^{-1}(A - D)\|_{\infty} &= \max \left\{ \sum_{\substack{j=1 \\ j \neq i}}^n \frac{|a_{ij}|}{|\lambda - a_{ii}|} \mid i = 1, \dots, n \right\} \\ &= \max \left\{ \frac{r_i}{|\lambda - a_{ii}|} \mid i = 1, \dots, n \right\}. \end{aligned}$$

Daher existiert ein $i \in \{1, \dots, n\}$ mit $\lambda \in G_i$. Zum Beweis des zweiten Teils des Satzes sei o.B.d.A.

$$\hat{G} = \bigcup_{i=1}^m G_i, \quad \tilde{G} = \bigcup_{i=m+1}^n G_i, \quad \hat{G} \cap \tilde{G} = \emptyset.$$

Für beliebiges $t \in [0, 1]$ definieren wir die Matrix $\mathbf{A}(t) = \mathbf{D} + t(\mathbf{A} - \mathbf{D})$. Offensichtlich gilt $\mathbf{A}(0) = \mathbf{D}$ und $\mathbf{A}(1) = \mathbf{A}$. Die GERSCHGORIN-Kreise für die Matrix $\mathbf{A}(t)$ sind dann durch

$$G_i(t) = \left\{ z \in \mathbb{C} \mid |z - a_{ii}| \leq t r_i \right\} \subset G_i, \quad i = 1, \dots, n$$

gegeben. Für $t \in [0, 1]$ liegen daher alle Eigenwerte von $\mathbf{A}(t)$ in

$$\hat{G} \cup \tilde{G} = \bigcup_{i=1}^n G_i.$$

Nun definieren wir

$$I = \left\{ t \in [0, 1] \mid \text{Genau } m \text{ Eigenwerte von } \mathbf{A}(t) \text{ liegen in } \hat{G} \right\}.$$

Weiterhin sei

$$t_0 = \sup \{ t \mid t \in I \}, \quad \varepsilon = \frac{1}{2} \min \{ |\hat{z} - \tilde{z}| \mid \hat{z} \in \hat{G}, \tilde{z} \in \tilde{G} \}.$$

Wegen $\mathbf{A}(0) = \mathbf{D}$ ist $0 \in I$; I ist daher nichtleer. Nun zeigen wir, dass auch $t_0 \in I$ gilt. Es seien dazu $\lambda_1(t_0), \dots, \lambda_n(t_0)$ die Eigenwerte von $\mathbf{A}(t_0)$. Es gilt

$$\|\mathbf{A}(t) - \mathbf{A}(t_0)\| = (t - t_0) \|\mathbf{A} - \mathbf{D}\|.$$

Nach Satz 9.16 gibt es zu ε ein $\delta = \delta(\varepsilon) > 0$, so dass für alle $t \in [0, 1]$ mit $|t - t_0| \leq \delta$ eine Numerierung der Eigenwerte $\lambda_1(t), \dots, \lambda_n(t)$ von $\mathbf{A}(t)$ mit

$$|\lambda_i(t) - \lambda_i(t_0)| \leq \varepsilon, \quad i = 1, \dots, n,$$

existiert. Nach Definition von t_0 existiert ein $t \in [t_0 - \delta, t_0] \cap I$. Für dieses t gilt: Genau m Eigenwerte $\lambda_i(t)$ von $\mathbf{A}(t)$ liegen in \hat{G} , die anderen in \tilde{G} . Wegen

$$|\lambda_i(t) - \lambda_i(t_0)| \leq \varepsilon = \frac{1}{2} \min \{ |\hat{z} - \tilde{z}| \mid \hat{z} \in \hat{G}, \tilde{z} \in \tilde{G} \}$$

liegen die ersten m Eigenwerte von $\mathbf{A}(t_0)$ ebenfalls in \hat{G} und die restlichen Eigenwerte in \tilde{G} . Damit gilt $t_0 \in I$. Nehmen wir nun an, dass $t_0 < 1$ gilt. Mit $\delta > 0$ und $t \in [t_0, t_0 + \delta] \cap I$ folgern wir genau wie eben, dass die Eigenwerte von $\mathbf{A}(t)$ so nummerierbar sind, dass $|\lambda_i(t) - \lambda_i(t_0)| \leq \varepsilon, \quad i = 1, \dots, n$, gilt. Wegen $t_0 \in I$ ist dann auch $t \in I$ im Widerspruch zur Definition von t_0 . Es muss daher $t_0 = 1$ gelten, und der Satz ist bewiesen. *

Eigentlich haben wir im ersten Teil des Beweises etwas mehr gezeigt, als zu beweisen war. Sind nämlich $A, D \in \mathbb{C}^{n \times n}$ und ist λ Eigenwert von A , so ist entweder λ Eigenwert von D oder

$$1 \leq \rho((\lambda I - D)^{-1}(A - D)) \leq \|(\lambda I - D)^{-1}(A - D)\|$$

und

$$1 \leq \rho((A - D)(\lambda I - D)^{-1}) \leq \|(A - D)(\lambda I - D)^{-1}\|$$

für jede Matrixgrenznorm $\|\circ\|$. Je nach Wahl von D und der Matrixnorm erhält man so verschiedene Abschätzungen für die Lage der Eigenwerte. Wählt man wieder

$$D = \text{diag}(a_{11}, a_{22}, \dots, a_{nn}),$$

aber diesmal $\|\circ\| = \|\circ\|_1$, also die Spaltensummennorm, so erhält man:

$$1 \leq \|(A - D)(\lambda I - D)^{-1}\|_1 = \max \left\{ \frac{1}{|\lambda - a_{jj}|} \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}| \mid j = 1, \dots, n \right\}.$$

Damit gelten die Aussagen des Kreisesatzes von GERSCHGORIN auch für die durch

$$G_j = \left\{ z \in \mathbb{C} \mid |z - a_{jj}| \leq q_j \right\}, \quad q_j = \sum_{\substack{i=1 \\ i \neq j}}^n |a_{ij}|, \quad j = 1, \dots, n$$

definierten Kreise.

Verwendet man die FROBENIUS-Norm, so existiert zu jedem Eigenwert ein $k \in \{1, \dots, n\}$ mit

$$|\lambda - a_{kk}| \leq \left(\sum_{\substack{i,j=1 \\ i \neq j}}^n |a_{ij}|^2 \right)^{1/2}.$$

Der Kreisesatz von GERSCHGORIN ist natürlich auch auf Matrizen anwendbar, die die gleichen Eigenwerte wie A besitzen, z. B. A^H oder jede durch Ähnlichkeitstransformation aus A hervorgegangenen Matrix.

9.18. Beispiel: Wir betrachten die Matrix

$$A = \begin{pmatrix} 0.9 & 0.0 & 0.0 \\ 0.0 & 0.4 & 0.0 \\ 0.0 & 0.0 & 0.2 \end{pmatrix} + 10^{-5} \begin{pmatrix} 0.1 & 0.4 & -0.2 \\ -0.1 & 0.5 & 0.1 \\ 0.2 & 0.1 & 0.3 \end{pmatrix}.$$

Nach Satz 9.17 ergeben sich folgende Aussagen über die Lokalisierung der Eigenwerte:

$$|\lambda_1 - (0.9 + 0.1 \cdot 10^{-5})| \leq 0.6 \cdot 10^{-5},$$

$$|\lambda_2 - (0.4 + 0.5 \cdot 10^{-5})| \leq 0.2 \cdot 10^{-5},$$

$$|\lambda_3 - (0.2 + 0.3 \cdot 10^{-5})| \leq 0.3 \cdot 10^{-5}.$$

Wendet man den Satz 9.17 auf die zu A ähnliche Matrix

$$P^{-1}AP, \quad P = \text{diag}(10^5, 1, 1)$$

an, so erhält man für den ersten Eigenwert folgende Abschätzung:

$$|\lambda_1 - (0.9 + 0.1 \cdot 10^{-5})| \leq 0.6 \cdot 10^{-10}.$$

Die beiden anderen GERSCHGORIN-Kreise sind nicht mehr disjunkt. Eine verbesserte Abschätzung für den zweiten Eigenwert erhält man, falls man $P = \text{diag}(1, 10^5, 1)$ verwendet:

$$|\lambda_2 - (0.4 + 0.5 \cdot 10^{-5})| \leq 0.2 \cdot 10^{-10}.$$

Analog ergibt sich mit $P = \text{diag}(1, 1, 10^5)$

$$|\lambda_3 - (0.2 + 0.3 \cdot 10^{-5})| \leq 0.3 \cdot 10^{-10}.$$



Der Kreisesatz von GERSCHGORIN diene dazu, die mögliche Lage der Eigenwerte einer gegebenen Matrix einzugrenzen. Wir wollen uns nun dem eigentlichen Thema dieses Abschnitts zuwenden, nämlich der Frage: Wie verhalten sich die Eigenwerte einer Matrix, falls die Matrixelemente gestört werden? Wir betrachten neben der Matrix $A \in \mathbb{C}^{n \times n}$ eine Störung $\delta A \in \mathbb{C}^{n \times n}$ und fragen nach dem Zusammenhang zwischen den Eigenwerten der Matrix A und der Matrix $A + \delta A$. Für beliebige Matrizen ist es wieder schwierig, genauere Aussagen zu treffen. Grenzt man die Klasse der zu betrachtenden Matrizen etwas ein (diagonalisierbare Matrizen, hermitesche Matrizen) so lässt sich das Verhalten der Eigenwerte dieser Matrizen gegenüber Störungen genauer beschreiben. Eine Matrix $A \in \mathbb{C}^{n \times n}$ heißt **diagonalisierbar**, falls eine reguläre Matrix $X \in \mathbb{C}^{n \times n}$ mit $X^{-1}AX = \text{diag}(\lambda_1, \dots, \lambda_n) = D$ existiert.

9.19. BAUER-FIKE: *Es sei $A \in \mathbb{C}^{n \times n}$ eine diagonalisierbare Matrix:*

$$X^{-1}AX = \text{diag}(\lambda_1, \dots, \lambda_n) = D.$$

Es sei weiterhin $\delta A \in \mathbb{C}^{n \times n}$ und λ ein Eigenwert von $A + \delta A$. Dann gilt bezüglich einer Matrixgrenznorm

$$\min \{ |\lambda - \lambda_j| \mid j = 1, \dots, n \} \leq \text{cond}(X) \|\delta A\|.$$

Beweis: Ist λ auch ein Eigenwert von \mathbf{A} , so ist die Behauptung offensichtlich wahr. Es sei nun $\lambda \neq \lambda_j$ für $j = 1, \dots, n$. Der Vektor \mathbf{x} sei ein zu λ gehörender Eigenvektor: $(\mathbf{A} + \delta\mathbf{A})\mathbf{x} = \lambda\mathbf{x}$. Dann folgt

$$\delta\mathbf{A}\mathbf{x} = (\lambda\mathbf{I} - \mathbf{A})\mathbf{x} = (\lambda\mathbf{I} - \mathbf{X}\mathbf{D}\mathbf{X}^{-1})\mathbf{x} = \mathbf{X}(\lambda\mathbf{I} - \mathbf{D})\mathbf{X}^{-1}\mathbf{x}$$

und daraus

$$\mathbf{X}^{-1}\mathbf{x} = (\lambda\mathbf{I} - \mathbf{D})^{-1}(\mathbf{X}^{-1}\delta\mathbf{A}\mathbf{X})\mathbf{X}^{-1}\mathbf{x}.$$

Wegen der Submultiplikativität der Matrixgrenznorm erhalten wir

$$\begin{aligned} \|\mathbf{X}^{-1}\mathbf{x}\| &\leq \|(\lambda\mathbf{I} - \mathbf{D})^{-1}\| \|\mathbf{X}^{-1}\delta\mathbf{A}\mathbf{X}\| \|\mathbf{X}^{-1}\mathbf{x}\| \\ &= \max \left\{ \frac{1}{|\lambda - \lambda_j|} \mid j = 1, \dots, n \right\} \|\mathbf{X}^{-1}\delta\mathbf{A}\mathbf{X}\| \|\mathbf{X}^{-1}\mathbf{x}\| \end{aligned}$$

und weiter

$$\min \{ |\lambda - \lambda_j| \mid j = 1, \dots, n \} \leq \text{cond}(\mathbf{X}) \|\delta\mathbf{A}\|.$$

*

Die Kondition der Matrix \mathbf{X} , die \mathbf{A} diagonalisiert, ist daher ein Maß für die Störanfälligkeit der Eigenwerte von \mathbf{A} . Für hermitesche Matrizen folgt sofort:

9.20. Satz: *Es sei $\mathbf{A} \in \mathbb{C}^{n \times n}$ eine hermitesche Matrix und $\delta\mathbf{A} \in \mathbb{C}^{n \times n}$. $\lambda_1, \dots, \lambda_n$ seien die (reellen) Eigenwerte von \mathbf{A} und λ sei ein Eigenwert der gestörten Matrix $\mathbf{A} + \delta\mathbf{A}$. Dann gilt bezüglich der Spektralnorm*

$$\min \{ |\lambda - \lambda_j| \mid j = 1, \dots, n \} \leq \|\delta\mathbf{A}\|_2.$$

Beweis: Für hermitesche Matrizen ist die Matrix \mathbf{X} , die \mathbf{A} diagonalisiert, unitär wählbar. Damit ist $\text{cond}_2(\mathbf{X}) = 1$ und mit Satz 9.19 folgt die Behauptung. *

Wie wir sehen, ist das Eigenwertproblem für hermitesche Matrizen (oder im Reellen für symmetrische Matrizen) besonders gut konditioniert. In diesem Falle lässt sich die Aussage des letzten Satzes sogar noch verschärfen. Vorher noch zwei wichtige Sätze über hermitesche Matrizen.

9.21. RAYLEIGH: *Für die hermitesche Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ seien die (reellen) Eigenwerte gemäß $\lambda_1 \geq \dots \geq \lambda_n$ geordnet. Durch $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ sei ein zugehöriges Orthornormalsystem von Eigenvektoren gegeben. Es gilt daher*

$$\mathbf{A}\mathbf{u}_i = \lambda_i\mathbf{u}_i, \quad \mathbf{u}_i^H \mathbf{u}_j = \delta_{ij}, \quad 1 \leq i, j \leq n.$$

Für $j = 1, \dots, n$ definieren wir den $(n + 1 - j)$ -dimensionalen Teilraum

$$\mathcal{M}_j = \text{span}(\mathbf{u}_j, \dots, \mathbf{u}_n).$$

Dann gilt

$$\lambda_j = \max \left\{ \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \mid \mathbf{x} \in \mathcal{M}_j \setminus \{\mathbf{o}\} \right\}, \quad j = 1, \dots, n.$$

Beweis: Für festes $j \in \{1, \dots, n\}$ sei $\mathbf{x} \in \mathcal{M}_j \setminus \{\mathbf{o}\}$ beliebig. Dann gilt

$$\mathbf{x} = \alpha_j \mathbf{u}_j + \dots + \alpha_n \mathbf{u}_n$$

und

$$\mathbf{A} \mathbf{x} = \lambda_j \alpha_j \mathbf{u}_j + \dots + \lambda_n \alpha_n \mathbf{u}_n.$$

Damit ergibt sich

$$\frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} = \frac{\sum_{i=j}^n \lambda_i |\alpha_i|^2}{\sum_{i=j}^n |\alpha_i|^2} \leq \lambda_j.$$

Es folgt daher

$$\sup \left\{ \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \mid \mathbf{x} \in \mathcal{M}_j \setminus \{\mathbf{o}\} \right\} \leq \lambda_j.$$

Andererseits ist $\mathbf{u}_j \in \mathcal{M}_j$ und $\mathbf{u}_j^H \mathbf{A} \mathbf{u}_j / \mathbf{u}_j^H \mathbf{u}_j = \lambda_j$, so dass

$$\max \left\{ \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \mid \mathbf{x} \in \mathcal{M}_j \setminus \{\mathbf{o}\} \right\} = \lambda_j.$$

✱

Bemerkungen: (i) Für eine hermitesche Matrix wird der Ausdruck

$$RQ(\mathbf{x}) = \mathbf{x}^H \mathbf{A} \mathbf{x} / \mathbf{x}^H \mathbf{x}$$

als RAYLEIGH-Quotient bezeichnet. Ist \mathbf{x} ein Eigenvektor von \mathbf{A} so ist $RQ(\mathbf{x})$ gleich dem zugehörigen Eigenwert. Ist \mathbf{x} nur eine Näherung für einen Eigenvektor von \mathbf{A} , so erwartet man, dass durch $RQ(\mathbf{x})$ eine Näherung für den zugehörigen Eigenwert gegeben ist.

(ii) Ist \mathbf{A} symmetrisch und reell, so gilt der Satz in analoger Weise im Reellen.

9.22. COURANT: Für die hermitesche Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ seien die (reellen) Eigenwerte gemäß $\lambda_1 \geq \dots \geq \lambda_n$ geordnet. Für $j = 1, \dots, n$ definieren wir die Mengen

$$\mathcal{N}_j = \left\{ \mathcal{N}_j \subseteq \mathbb{C}^{n \times n} \mid \mathcal{N}_j \text{ ist linearer Teilraum der Dimension } n + 1 - j \right\}.$$

Dann ist

$$\lambda_j = \min \left\{ \max \left\{ \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \mid \mathbf{x} \in \mathcal{N}_j \setminus \{\mathbf{o}\} \right\} \mid \mathcal{N}_j \in \mathcal{N}_j \right\}, \quad j = 1, \dots, n.$$

Beweis: Es sei

$$\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$$

ein Orthonormalsystem von Eigenvektoren zu den Eigenwerten $\lambda_1, \dots, \lambda_n$. Für festes $j \in \{1, \dots, n\}$ definieren wir den linearen Teilraum

$$\mathcal{L}_j = \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_j).$$

Der Durchschnitt von \mathcal{L}_j mit einem beliebigen Teilraum $\mathcal{N}_j \in \mathcal{N}_j$ enthält wegen

$$\dim(\mathcal{L}_j \cap \mathcal{N}_j) = \dim(\mathcal{L}_j) + \dim(\mathcal{N}_j) - \dim(\mathcal{L}_j \oplus \mathcal{N}_j) \geq 1$$

ein $\mathbf{x} \neq \mathbf{o}$. Aus $\mathbf{x} \in \mathcal{L}_j$ folgt $\mathbf{x}^H \mathbf{A} \mathbf{x} / \mathbf{x}^H \mathbf{x} \geq \lambda_j$ und damit

$$\min \left\{ \max \left\{ \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \mid \mathbf{x} \in \mathcal{N}_j \setminus \{\mathbf{o}\} \right\} \mid \mathcal{N}_j \in \mathcal{N}_j \right\} \geq \lambda_j.$$

Wählt man nun

$$\mathcal{N}_j = \mathcal{M}_j = \left\{ \mathbf{x} \in \mathbb{C}^n \mid \mathbf{u}_i^H \mathbf{x} = 0, \quad i = 1, \dots, j-1 \right\},$$

so folgt nach Satz 9.21

$$\max \left\{ \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \mid \mathbf{x} \in \mathcal{M}_j \setminus \{\mathbf{o}\} \right\} = \lambda_j.$$

Damit ist der Satz bewiesen. *

Nun sind wir in der Lage, das Verhalten der Eigenwerte einer hermiteschen Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ unter einer hermiteschen Störung $\delta \mathbf{A} \in \mathbb{C}^{n \times n}$ genauer anzugeben.

9.23. Satz: Die hermitesche Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ habe die (reellen) Eigenwerte

$$\lambda_1 \geq \dots \geq \lambda_n.$$

Für eine hermitesche Störung $\delta\mathbf{A} \in \mathbb{C}^{n \times n}$ betrachten wir die Matrix $\mathbf{A} + \delta\mathbf{A} \in \mathbb{C}^{n \times n}$ mit den (reellen) Eigenwerten $\mu_1 \geq \dots \geq \mu_n$. Bezüglich jeder Matrixgrenznorm gilt dann die Abschätzung

$$|\lambda_j - \mu_j| \leq \|\delta\mathbf{A}\|.$$

Beweis: Da \mathbf{A} und $\delta\mathbf{A}$ hermitesch sind, gilt

$$\frac{\mathbf{x}^H (\mathbf{A} - (\mathbf{A} + \delta\mathbf{A})) \mathbf{x}}{\mathbf{x}^H \mathbf{x}} = \frac{\mathbf{x}^H (-\delta\mathbf{A}) \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \leq \|\delta\mathbf{A}\|_2 = \varrho(\delta\mathbf{A}) \leq \|\delta\mathbf{A}\|.$$

Daraus folgt

$$\frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \leq \frac{\mathbf{x}^H (\mathbf{A} + \delta\mathbf{A}) \mathbf{x}}{\mathbf{x}^H \mathbf{x}} + \|\delta\mathbf{A}\|.$$

Für festes $j \in \{1, \dots, n\}$ folgt dann wie in Satz 9.22 für ein beliebiges $\mathcal{N}_j \in \mathcal{N}_j$

$$\begin{aligned} & \max \left\{ \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \mid \mathbf{x} \in \mathcal{N}_j \setminus \{\mathbf{o}\} \right\} \\ & \leq \max \left\{ \frac{\mathbf{x}^H (\mathbf{A} + \delta\mathbf{A}) \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \mid \mathbf{x} \in \mathcal{N}_j \setminus \{\mathbf{o}\} \right\} + \|\delta\mathbf{A}\| \end{aligned}$$

und weiter

$$\begin{aligned} & \min \left\{ \max \left\{ \frac{\mathbf{x}^H \mathbf{A} \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \mid \mathbf{x} \in \mathcal{N}_j \setminus \{\mathbf{o}\} \right\} \mid \mathcal{N}_j \in \mathcal{N}_j \right\} \\ & \leq \min \left\{ \max \left\{ \frac{\mathbf{x}^H (\mathbf{A} + \delta\mathbf{A}) \mathbf{x}}{\mathbf{x}^H \mathbf{x}} \mid \mathbf{x} \in \mathcal{N}_j \setminus \{\mathbf{o}\} \right\} \mid \mathcal{N}_j \in \mathcal{N}_j \right\} + \|\delta\mathbf{A}\|. \end{aligned}$$

Nach Satz 9.22 heißt das

$$\lambda_j \leq \mu_j + \|\delta\mathbf{A}\|, \quad j = 1, \dots, n.$$

Vertauscht man \mathbf{A} und $\mathbf{A} + \delta\mathbf{A}$, so erhält man

$$\mu_j \leq \lambda_j + \|\delta\mathbf{A}\|, \quad j = 1, \dots, n,$$

insgesamt also

$$|\lambda_j - \mu_j| \leq \|\delta\mathbf{A}\|, \quad j = 1, \dots, n.$$

Die letzten Sätze haben gezeigt, dass die Eigenwerte einer nichtdefektiven Matrix lokal lipschitzstetige Funktionen der Matrixelemente sind. Im Falle hermitescher Matrizen ist die Lipschitzkonstante gleich 1. Das Eigenwertproblem für hermitesche Matrizen (oder im Reellen für symmetrische Matrizen) ist besonders gut konditioniert.

Wir wollen nun noch untersuchen, wie sich die Eigenvektoren einer Matrix gegenüber Störungen verhalten.

9.24. Satz: Die hermitesche Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ habe die (reellen) Eigenwerte

$$\lambda_1, \dots, \lambda_n$$

und die zugehörigen orthonormierten Eigenvektoren

$$\mathbf{u}_1, \dots, \mathbf{u}_n.$$

Weiterhin sei $\delta\mathbf{A} \in \mathbb{C}^{n \times n}$ und (μ, \mathbf{u}) sei ein Eigenpaar von $\mathbf{A} + \delta\mathbf{A}$. Dann gelten folgende Aussagen:

(i) Falls für ein $j \in \{1, \dots, n\}$ $\mathbf{u}_j^H \mathbf{u} \neq 0$, so gilt

$$|\mu - \lambda_j| \leq \frac{\|\delta\mathbf{A}\mathbf{u}\|_2}{\mathbf{u}_j^H \mathbf{u}} \leq \frac{\|\delta\mathbf{A}\|_2}{\mathbf{u}_j^H \mathbf{u}}. \quad (9.1)$$

(ii) Falls für ein $j \in \{1, \dots, n\}$ $\mu \neq \lambda_j$, so gilt

$$|\cos(\angle(\mathbf{u}_j, \mathbf{u}))| \leq \frac{\|\delta\mathbf{A}\mathbf{u}\|_2}{|\mu - \lambda_j|} \leq \frac{\|\delta\mathbf{A}\|_2}{|\mu - \lambda_j|}. \quad (9.2)$$

(iii) Falls für ein $j \in \{1, \dots, n\}$

$$\gamma_j = \min \{ |\mu - \lambda_i| \mid i = 1, \dots, n, i \neq j \} > 0$$

so gilt

$$\sin(\angle(\mathbf{u}_j, \mathbf{u})) \leq \frac{\|\delta\mathbf{A}\mathbf{u}\|_2}{\gamma_j} \leq \frac{\|\delta\mathbf{A}\|_2}{\gamma_j}. \quad (9.3)$$

Beweis: Wir stellen den Vektor \mathbf{u} als Linearkombination der Eigenvektoren der Matrix \mathbf{A} dar:

$$\mathbf{u} = \sum_{i=1}^n \zeta_i \mathbf{u}_i = \mathbf{U}\mathbf{z}, \quad \mathbf{z} = (\zeta_1, \dots, \zeta_n)^T, \quad \|\mathbf{z}\|_2 = 1.$$

Dann gilt $\mathbf{z} = \mathbf{U}^H \mathbf{u}$ und $\zeta_i = \mathbf{u}_i^H \mathbf{u} = \cos(\angle(\mathbf{u}_i, \mathbf{u}))$ für $i = 1, \dots, n$. Es sei weiter

$$\begin{aligned} \mathbf{y} &= (\mu \mathbf{I} - \mathbf{A}) \mathbf{z} \\ &= (\mu \mathbf{U}^H \mathbf{U} - \mathbf{U}^H \mathbf{A} \mathbf{U}) \mathbf{z} \\ &= \mathbf{U}^H (\mu \mathbf{I} - \mathbf{A}) \mathbf{U} \mathbf{z} \\ &= \mathbf{U}^H (\mu \mathbf{I} - (\mathbf{A} + \delta \mathbf{A}) + \delta \mathbf{A}) \mathbf{u} \\ &= \mathbf{U}^H (\mu \mathbf{I} - (\mathbf{A} + \delta \mathbf{A})) \mathbf{u} + \mathbf{U}^H \delta \mathbf{A} \mathbf{u} \\ &= \mathbf{U}^H \delta \mathbf{A} \mathbf{u}. \end{aligned}$$

Für eine Komponente η_i des Vektors \mathbf{y} gilt dann $\eta_i = (\mu - \lambda_i) \zeta_i = \mathbf{u}_i^H \delta \mathbf{A} \mathbf{u}$. Ist nun j ein Index, für den $\zeta_j = \mathbf{u}_j^H \mathbf{u} \neq 0$, so folgt sofort

$$|\mu - \lambda_j| = \frac{|\mathbf{u}_j^H \delta \mathbf{A} \mathbf{u}|}{\mathbf{u}_j^H \mathbf{u}} \leq \frac{\|\mathbf{u}_j\|_2 \|\delta \mathbf{A} \mathbf{u}\|_2}{\mathbf{u}_j^H \mathbf{u}} = \frac{\|\delta \mathbf{A} \mathbf{u}\|_2}{\mathbf{u}_j^H \mathbf{u}} \leq \frac{\|\delta \mathbf{A}\|_2}{\mathbf{u}_j^H \mathbf{u}},$$

also der erste Teil des Satzes.

Gilt für ein $j \in \{1, \dots, n\}$ $\mu \neq \lambda_j$, so lässt sich abschätzen, wie stark \mathbf{u} und \mathbf{u}_j von der Orthogonalität abweichen:

$$|\cos(\angle(\mathbf{u}_j, \mathbf{u}))| = |\zeta_j| = \frac{|\mathbf{u}_j^H \delta \mathbf{A} \mathbf{u}|}{|\mu - \lambda_j|} \leq \frac{\|\delta \mathbf{A} \mathbf{u}\|_2}{|\mu - \lambda_j|}.$$

Damit wäre die zweite Aussage des Satzes bewiesen. Für den Beweis der dritten Aussage beachte man, dass für

$$\gamma_j = \min\{|\mu - \lambda_j| : i = 1, \dots, n, i \neq j\} > 0$$

$$(\sin(\angle(\mathbf{u}_j, \mathbf{u})))^2 = 1 - \zeta_j^2 = \sum_{\substack{i=1 \\ i \neq j}}^n \zeta_i^2 = \sum_{\substack{i=1 \\ i \neq j}}^n \frac{|\eta_i|^2}{|\mu - \lambda_j|^2} \leq \frac{\|\mathbf{y}\|_2^2}{\gamma_j^2} = \frac{\|\delta \mathbf{A} \mathbf{u}\|_2^2}{\gamma_j^2}$$

folgt. Nach Wurzelziehen erhält man die letzte Aussage des Satzes. *

Bemerkungen: (i) Es seien

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \text{ bzw. } \mu_1 \geq \mu_2 \geq \dots \geq \mu_n$$

die geordneten Eigenwerte der hermiteschen Matrizen \mathbf{A} bzw. $\mathbf{A} + \delta \mathbf{A}$. Ferner seien

$$\{\mathbf{u}_1, \dots, \mathbf{u}_n\} \text{ und } \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$$

zugehörige Orthonormalsysteme von Eigenvektoren. Setzen wir nun

$$\{\mu, \mathbf{u}\} = \{\mu_i, \mathbf{v}_i\},$$

so lautet (9.2):

$$|\cos(\angle(\mathbf{u}_j, \mathbf{v}_i))| \leq \frac{\|\delta\mathbf{A}\mathbf{v}_i\|_2}{|\mu_i - \lambda_j|} \leq \frac{\|\delta\mathbf{A}\|_2}{|\mu_i - \lambda_j|}.$$

Sind die Eigenwerte μ_i und λ_j genügend stark voneinander getrennt, und ist die Störung $\delta\mathbf{A}$ hinreichend klein, so sind die Eigenvektoren \mathbf{v}_i und \mathbf{u}_j fast orthogonal. Nach Satz 9.23 gilt aber

$$|\mu_i - \lambda_j| = |\mu_i - \lambda_i + \lambda_i - \lambda_j| \geq |\lambda_i - \lambda_j| - |\mu_i - \lambda_i| \geq |\lambda_i - \lambda_j| - \|\delta\mathbf{A}\|_2.$$

Damit ist μ_i hinreichend stark von λ_j getrennt, falls λ_i hinreichend stark von λ_j getrennt ist und die Störung $\delta\mathbf{A}$ hinreichend klein ist. Die Orthogonalität zweier Eigenvektoren einer hermiteschen Matrix wird durch eine kleine Störung wenig beeinflusst, falls die zugehörigen Eigenwerte weit genug voneinander entfernt sind.

(ii) Aus (9.3) erhalten wir für $\{\mu, \mathbf{u}\} = \{\mu_j, \mathbf{v}_j\}$

$$\sin(\angle(\mathbf{u}_j, \mathbf{v}_j)) \leq \frac{\|\delta\mathbf{A}\mathbf{v}_j\|_2}{\gamma_j} \leq \frac{\|\delta\mathbf{A}\|_2}{\gamma_j}, \quad \gamma_j = \min \{ |\mu_j - \lambda_i| \mid i = 1, \dots, n, i \neq j \}.$$

Sind γ_j hinreichend groß und $\|\delta\mathbf{A}\|_2$ hinreichend klein, so weichen die Richtungen der Eigenvektoren \mathbf{u}_j und \mathbf{v}_j nur wenig voneinander ab. Wegen

$$\gamma_j \geq \min \{ |\lambda_j - \lambda_i| \mid i = 1, \dots, n, i \neq j \} - \|\delta\mathbf{A}\|_2$$

ist dies der Fall, falls der Eigenwert λ_j zusätzlich genügend stark von allen anderen Eigenwerten λ_i getrennt ist. Für isolierte Eigenwerte einer hermiteschen Matrix haben daher kleine Störungen nur geringen Einfluss auf die Richtungen der zugehörigen Eigenvektoren.

(iii) In vielen Fällen gilt $\|\delta\mathbf{A}\mathbf{u}\|_2 \ll \|\delta\mathbf{A}\|_2\|\mathbf{u}\|_2$. Die Schranken mit $\|\delta\mathbf{A}\mathbf{u}\|_2$ sind dann wesentlich besser als die, die durch die Vergrößerung $\|\delta\mathbf{A}\mathbf{u}\|_2 \leq \|\delta\mathbf{A}\|_2$ entstehen. Da $\delta\mathbf{A}$ als hermitesch vorausgesetzt wurde, lässt sich wegen

$$|\mathbf{u}_i^H \delta\mathbf{A}\mathbf{u}| = |(\delta\mathbf{A}\mathbf{u})_i^H \mathbf{u}| \leq \|\delta\mathbf{A}\mathbf{u}_i\|_2$$

der Term $\|\delta\mathbf{A}\mathbf{u}\|_2$ in (9.1) bzw. (9.2) durch $\|\delta\mathbf{A}\mathbf{u}_j\|_2$ bzw. $\|\delta\mathbf{A}\mathbf{u}_i\|_2$ ersetzen.

Als letztes wollen wir in diesem Abschnitt noch sogenannte Residualkriterien angeben, die es uns ermöglichen, von einem Paar $\{\mu, \mathbf{u}\}$ zu entscheiden, ob es im Rahmen einer vorgegebenen Genauigkeit als Eigenpaar einer Matrix \mathbf{A} akzeptierbar ist.

9.25. Satz: Es seien $\mathbf{A} \in \mathbb{C}^{n \times n}$, $\mu \in \mathbb{C}$ und $\mathbf{u} \in \mathbb{C}^n$ mit $\|\mathbf{u}\|_2 = 1$ gegeben. Dann sind folgende Aussagen äquivalent:

1. Es existiert ein $\delta\mathbf{A} \in \mathbb{C}^{n \times n}$ mit

$$(\mathbf{A} + \delta\mathbf{A})\mathbf{u} = \mu\mathbf{u}, \quad \|\delta\mathbf{A}\|_2 \leq \varepsilon.$$

2. Es gilt

$$\|\mathbf{r}\|_2 \leq \varepsilon, \quad \mathbf{r} = \mathbf{A}\mathbf{u} - \mu\mathbf{u}.$$

Der Vektor $\mathbf{r} = \mathbf{A}\mathbf{u} - \mu\mathbf{u}$ heißt Residuum des Paares $\{\mu, \mathbf{u}\}$ in Bezug auf das Eigenwertproblem der Matrix \mathbf{A} .

Beweis: Aus $\mathbf{r} = \mathbf{A}\mathbf{u} - \mu\mathbf{u}$ und $(\mathbf{A} + \delta\mathbf{A})\mathbf{u} = \mu\mathbf{u}$ folgt

$$\mathbf{r} = -\delta\mathbf{A}\mathbf{u} \quad \|\mathbf{r}\|_2 \leq \|\delta\mathbf{A}\|_2 \leq \varepsilon.$$

Es sei nun $\eta = \|\mathbf{r}\|_2 \leq \varepsilon$. Dann existiert eine HOUSEHOLDER-Spiegelung $\mathbf{H} \in \mathbb{C}^{n \times n}$ mit $\|\mathbf{H}\|_2 = 1$, für die $\mathbf{r} = \eta\mathbf{H}\mathbf{u}$ gilt. Damit ist $\mathbf{A}\mathbf{u} - \mu\mathbf{u} = \eta\mathbf{H}\mathbf{u}$, also $\delta\mathbf{A} = -\eta\mathbf{H}$ und $\|\delta\mathbf{A}\|_2 = \|\eta\mathbf{H}\|_2 = \eta \leq \varepsilon$. *

Wir akzeptieren daher ein Paar $\{\mu, \mathbf{u}\}$ als Eigenpaar einer benachbarten Matrix $\mathbf{A} + \delta\mathbf{A}$, falls sein Residuum bezüglich des Eigenwertproblems der Matrix \mathbf{A} klein genug ist.

Die Algorithmen, die wir kennenlernen werden, liefern teilweise nur Näherungen für Eigenvektoren oder nur Näherungen für Eigenwerte. In diesen Fällen stellt sich die Aufgabe, zu einer gegebenen Näherung für einen Eigenvektor bzw. einen Eigenwert entsprechende gute Näherungen für einen zugehörigen Eigenwert bzw. Eigenvektor zu finden. Betrachten wir das erste Problem.

9.26. Satz: Es seien die hermitesche Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ und ein Vektor $\mathbf{u} \in \mathbb{C}^n$ mit $\|\mathbf{u}\|_2 = 1$ gegeben. Dann gilt:

1. Das lineare Ausgleichsproblem $\min \{ \|\mathbf{A}\mathbf{u} - \mu\mathbf{u}\|_2 \mid \mu \in \mathbb{R} \}$ hat die eindeutige Lösung

$$\mu^* = \mathbf{u}^H \mathbf{A}\mathbf{u} = RQ(\mathbf{u}).$$

Das zugehörige Residuum $\mathbf{r} = \mathbf{A}\mathbf{u} - \mu^*\mathbf{u}$ ist orthogonal zu \mathbf{u} , es gilt $\mathbf{u}^H \mathbf{r} = 0$.

2. Das Paar $\{\mu^*, \mathbf{u}\}$ ist Eigenpaar einer Matrix $\mathbf{A} + \delta\mathbf{A}$ mit

$$\delta\mathbf{A} \in \mathbb{C}^{n \times n}, \quad \|\delta\mathbf{A}\|_2 \leq \|\mathbf{r}\|_2.$$

3. Es gibt ein Eigenpaar $\{\lambda_j, \mathbf{u}_j\}$ von \mathbf{A} mit

$$|\mu^* - \lambda_j| = \min \{ |\mu^* - \lambda_i| \mid i = 1, \dots, n \} \leq \|\mathbf{r}\|_2.$$

Falls $\gamma_j = \min \{ |\mu^* - \lambda_i| \mid i = 1, \dots, n, i \neq j \} > \|\mathbf{r}\|_2$ gilt

$$\sin(\angle(\mathbf{u}, \mathbf{u}_j)) \leq \|\mathbf{r}\|_2 / \gamma_j.$$

Beweis: 1. Die Normalengleichungen zu dem gegebenen Ausgleichsproblem lauten $\mu \mathbf{u}^H \mathbf{u} = \mathbf{u}^H \mathbf{A} \mathbf{u}$. Daraus folgt wegen $\mathbf{u}^H \mathbf{u} = 1$ sofort die erste Aussage. Die Orthogonalität zwischen Residuum und dem Vektor \mathbf{u} erkennt man durch einfaches Nachrechnen.

2. Diese Aussage erhält man mit dem eben bewiesenen und Satz 9.25.

3. Die erste Aussage ergibt sich aus (ii) mit Hilfe von Satz 9.19. Für die Abschätzung des Winkels zwischen \mathbf{u} und \mathbf{u}_j wende man Satz 9.24 an. *

Die optimale Eigenwertnäherung zu einer Eigenvektornäherung \mathbf{v} besteht nach diesem Satz im schon erwähnten RAYLEIGH-Quotienten $RQ(\mathbf{v}) = \mathbf{v}^H \mathbf{A} \mathbf{v} / \mathbf{v}^H \mathbf{v}$. Zum Problem der Wahl einer Eigenvektornäherung bei gegebener Eigenwertnäherung erhalten wir den folgenden Satz.

9.27. Satz: Die hermitesche Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ besitze die (reellen) Eigenwerte

$$\lambda_1, \dots, \lambda_n.$$

Durch $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ sei ein zugehöriges Orthonormalsystem von Eigenvektoren gegeben. Für festes $\mu \in \mathbb{R}$ hat die Aufgabe

$$\min \{ \|\mathbf{A} \mathbf{u} - \mu \mathbf{u}\|_2 \mid \mathbf{u} \in \mathbb{C}^n, \|\mathbf{u}\|_2 = 1 \}$$

die Lösung $\mathbf{u} = \mathbf{u}_j$, wobei j durch

$$\min \{ |\lambda_i - \mu| \mid i = 1, \dots, n \} = |\lambda_j - \mu|$$

definiert ist. Für das minimale Residuum gilt

$$\|\mathbf{A} \mathbf{u}_j - \mu \mathbf{u}_j\|_2 = |\lambda_j - \mu|.$$

Beweis: Mit der unitären Matrix $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_n)$ und $\mathbf{u} = \mathbf{U} \mathbf{z}$ gilt

$$\|\mathbf{A} \mathbf{u} - \mu \mathbf{u}\|_2 = \|\mathbf{U}^H (\mathbf{A} \mathbf{u} - \mu \mathbf{u})\|_2 = \|\Lambda \mathbf{z} - \mu \mathbf{z}\|_2 = \sum_{i=1}^n (\lambda_i - \mu) z_i^2, \quad \sum_{i=1}^n z_i^2 = 1.$$

Damit erhält man für das äquivalente Minimierungsproblem

$$\min \{ \|\Lambda \mathbf{z} - \mu \mathbf{z}\|_2 \mid \mathbf{z} \in \mathbb{C}^n, \|\mathbf{z}\|_2 = 1 \}$$

sofort die Lösung $\mathbf{z} = \mathbf{e}_j$ und damit $\mathbf{u} = \mathbf{u}_j$. *

Die optimale Eigenvektornäherung zu einer vorgegebenen Eigenwertnäherung besteht also in einem Eigenvektor zum nächstgelegenen exakten Eigenwert. Ein effektives Verfahren zum näherungsweise Lösen dieser Aufgabe werden wir in Form der Inversen Iteration nach WIELANDT kennenlernen.

9.3. Das Jacobi-Verfahren

In diesem und den nächsten Abschnitten werden wir uns mit verschiedenen Verfahren zum Lösen von Eigenwertproblemen befassen. Dabei werden wir alle Verfahren im Reellen erläutern. Die Übertragung auf den komplexen Fall dürfte nicht schwierig sein. Prinzipiell könnte man zumindest im reellen symmetrischen Falle daran denken, die Eigenwerte als Nullstellen des charakteristischen Polynoms zu bestimmen. Dies führt im allgemeinen zu großen Fehlern. Denn während das symmetrische Eigenwertproblem gut konditioniert ist, reagieren die Nullstellen eines Polynoms i. a. äußerst empfindlich auf Änderungen in den Koeffizienten. Betrachten wir als Beispiel eine symmetrische Matrix mit den Eigenwerten $\lambda_1 = 1, \lambda_2 = 2, \dots, \lambda_{20} = 20$. Das charakteristische Polynom dieser Matrix lautet

$$\varphi(\lambda) = (\lambda - 1)(\lambda - 2) \cdots (\lambda - 20).$$

Ändern wir ein Matrixelement um den Wert $\varepsilon = 2^{-23} \approx 10^{-7}$ so ändern sich die Eigenwerte dieser Matrix nach Satz 9.23 höchstens um den Wert ε . Stört man den Koeffizienten von λ^{19} im charakteristischen Polynom um ε , so erhält man ein Polynom mit folgenden Nullstellen:

1.000	000	000	10.095	266	145 ± 0.643	500	$904i$
2.000	000	000	11.793	633	881 ± 1.652	329	$728i$
3.000	000	000	13.992	358	137 ± 2.518	830	$070i$
4.000	000	000	16.730	737	466 ± 2.812	624	$894i$
4.999	999	928	19.502	439	400 ± 1.940	330	$347i$
6.000	006	944					
6.999	697	234					
8.007	267	603					
8.917	250	249					
20.846	908	101					

Die Nullstellen des gestörten Polynoms weichen erheblich von den exakten Nullstellen ab. Es treten sogar fünf komplexe Nullstellenpaare auf. Dies Beispiel zeigt schon, dass der Weg der Eigenwertberechnung über das charakteristische Polynom nicht gangbar ist.

Es gibt mit dem QR -Algorithmus ein leistungsfähiges Werkzeug zum Lösen des

mit geeigneten Parametern p, q, c, s zu wählen. Die Matrizen $\mathbf{A}^{(k)}$, $k = 0, 1, \dots$, sind alle orthogonal ähnlich. Es gilt

$$\mathbf{A}^{(k)} = \mathbf{V}^{(k)T} \mathbf{A} \mathbf{V}^{(k)}, \quad \mathbf{V}^{(k)} = \begin{cases} \mathbf{I} & \text{für } k = 0 \\ \mathbf{G}^{(1)} \mathbf{G}^{(2)} \dots \mathbf{G}^{(k-1)} & \text{für } k \geq 1 \end{cases}.$$

Das Ziel der Transformationen ist zumindest näherungsweise Diagonalgestalt der Matrix. Als Maß für die Abweichung der Iterationsmatrizen von der Diagonalgestalt bietet sich die Summe der Quadrate der Nichtdiagonalelemente an. Wegen der Symmetrie von \mathbf{A} genügt es dabei, nur die Quadrate der Elemente oberhalb der Diagonalen aufzusummieren. Wir definieren darum:

$$\omega_k = \omega(\mathbf{A}^{(k)}) = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \left(a_{ij}^{(k)} \right)^2.$$

Im k -ten Schritt sind demnach die Parameter p, q, c, s so zu wählen, dass ω_{k+1} möglichst klein wird.

Es sei weiterhin

$$\mathbf{M}^{(k)} = \text{diag}(a_{11}^{(k)}, a_{22}^{(k)}, \dots, a_{nn}^{(k)}), \quad \mathbf{R}^{(k)} = \mathbf{A}^{(k)} - \mathbf{M}^{(k)}.$$

Dann ergibt sich

$$\mathbf{M}^{(k)} = \mathbf{A}^{(k)} - \mathbf{R}^{(k)} = \mathbf{V}^{(k)T} \left(\mathbf{A} + \delta \mathbf{A}^{(k)} \right) \mathbf{V}^{(k)}$$

mit der symmetrischen Störung

$$\delta \mathbf{A}^{(k)} = -\mathbf{V}^{(k)} \mathbf{R}^{(k)} \mathbf{V}^{(k)T}, \quad \text{für die} \quad \|\delta \mathbf{A}^{(k)}\|_F = \|\mathbf{R}^{(k)}\|_F = \sqrt{2\omega_k} \quad \text{gilt.}$$

Die letzten Gleichungen besagen, dass die Diagonalelemente der Matrix $\mathbf{A}^{(k)}$ die exakten Eigenwerte der gestörten Matrix $\mathbf{A} + \delta \mathbf{A}^{(k)}$ mit den Spalten von $\mathbf{V}^{(k)}$ als zugehörigen Eigenvektoren sind. Nach Satz 9.23 lassen sich dann die Eigenwerte λ_i von \mathbf{A} so ordnen, dass

$$|a_{ii}^{(k)} - \lambda_{l(i)}| \leq \|\delta \mathbf{A}^{(k)}\|_2 \leq \|\delta \mathbf{A}^{(k)}\|_F = \sqrt{2\omega_k} \quad \text{gilt.} \quad (9.4)$$

Für jedes k ist dabei $\{l(1), \dots, l(n)\}$ eine Permutation der Zahlen $\{1, \dots, n\}$ derart, dass die Eigenwerte von \mathbf{A} in der gleichen Weise wie die Diagonalelemente von $\mathbf{A}^{(k)}$ geordnet sind. Die Diagonalelemente von $\mathbf{A}^{(k)}$ approximieren damit die Eigenwerte von \mathbf{A} mit einem absoluten Fehler, der durch $\sqrt{2\omega_k}$ beschränkt ist.

Kommen wir nun zur Wahl der Parameter p, q, c, s . Ziel ist es in jedem Schritt, das Maß für die Nichtdiagonalität der Matrix so stark wie möglich zu verringern. Wir haben im k -ten Schritt die Matrix $\mathbf{G}^{(k)} = \mathbf{G}_{pq}(c, s)$ so zu wählen, dass

$$\Phi_k(p, q, c, s) = \omega_k - \omega_{k+1}$$

maximiert wird. Betrachten wir einen Transformationsschritt

$$\bar{\mathbf{A}} = \mathbf{G}_{pq}(c, s)^T \mathbf{A} \mathbf{G}_{pq}(c, s).$$

Bei dieser Transformation werden nur die Elemente der p -ten und q -ten Zeile und Spalte von \mathbf{A} geändert. Es gilt

$$\begin{aligned} \bar{a}_{pp} &= c^2 a_{pp} - 2csa_{pq} + s^2 a_{qq}, \\ \bar{a}_{qq} &= s^2 a_{pp} + 2csa_{pq} + c^2 a_{qq}, \\ \bar{a}_{pq} &= \bar{a}_{qp} = cs(a_{pp} - a_{qq}) + (c^2 - s^2)a_{pq}, \\ \bar{a}_{ip} &= \bar{a}_{pi} = ca_{ip} - sa_{iq}, \quad i \neq p, q, \\ \bar{a}_{iq} &= \bar{a}_{qi} = sa_{ip} + ca_{iq}, \quad i \neq p, q. \end{aligned}$$

Wegen $c^2 + s^2 = 1$ folgt

$$\bar{a}_{ip}^2 + \bar{a}_{iq}^2 = a_{ip}^2 + a_{iq}^2, \quad i \neq p, q.$$

Die Quadratsumme der Nichtdiagonalelemente wird nur durch die Elemente in den Positionen (p, q) und (q, p) verändert. Damit gilt

$$\Phi(p, q, c, s) = \omega - \bar{\omega} = \omega(\mathbf{A}) - \omega(\bar{\mathbf{A}}) = a_{pq}^2 - \bar{a}_{pq}^2.$$

Hält man p und q einmal fest, so wird $\Phi(p, q, c, s)$ offensichtlich maximal, falls c und s so gewählt werden, dass $\bar{a}_{pq} = 0$ gilt. Dann ist

$$\omega(\bar{\mathbf{A}}) = \omega(\mathbf{A}) - a_{pq}^2.$$

Als Bestimmungsgleichungen für c und s erhält man

$$\begin{aligned} cs(a_{pp} - a_{qq}) + (c^2 - s^2)a_{pq} &= 0, \\ c^2 + s^2 &= 1. \end{aligned}$$

Mit $t = s/c$ erhalten wir aus der ersten Gleichung

$$t^2 + 2\delta t - 1 = 0, \quad \delta = \frac{a_{qq} - a_{pp}}{2a_{pq}}.$$

Die betragskleinste Lösung dieser quadratischen Gleichung ist nach der Vorschrift

$$t = \begin{cases} 1 / \left(|\delta| + \sqrt{1 + \delta^2} \right) & \text{für } \delta \geq 0 \\ -1 / \left(|\delta| + \sqrt{1 + \delta^2} \right) & \text{für } \delta < 0 \end{cases} \quad (9.5)$$

numerisch stabil berechenbar. Damit erhält man die gesuchten Parameter zu

$$c = 1 / \sqrt{1 + t^2}, \quad s = tc. \quad (9.6)$$

Mit der Hilfsgröße

$$\tau = s / (1 + c) \quad (9.7)$$

ergeben sich jetzt folgende Transformationsformeln:

$$\bar{a}_{pp} = a_{pp} - ta_{pq}, \quad (9.8)$$

$$\bar{a}_{qq} = a_{qq} + ta_{pq}, \quad (9.9)$$

$$\bar{a}_{pq} = \bar{a}_{qp} = 0, \quad (9.10)$$

$$\bar{a}_{ip} = \bar{a}_{pi} = a_{ip} - s(a_{iq} + \tau a_{ip}), \quad i \neq p, q \quad (9.11)$$

$$\bar{a}_{iq} = \bar{a}_{qi} = a_{iq} + \tau(a_{ip} + \bar{a}_{ip}), \quad i \neq p, q. \quad (9.12)$$

lässt man die Wahl von p und q vorerst noch offen, so erhält man folgenden Basisalgorithmus:

9.28. JACOBI-Basisverfahren:

S0 (Initialisierung) Setze $\mathbf{A}^{(0)} = \mathbf{A}$, $k = 0$ und berechne $\omega_0 = \omega(\mathbf{A})$. Wähle ein $\varepsilon > 0$.

S1 (Abbruchtest) Falls $2\omega_k \leq \varepsilon^2$ STOPP.

S2 (Pivotwahl) Wähle Indizes $p = p(k)$ und $q = q(k)$ mit $a_{pq}^{(k)} \neq 0$.

S3 (Iterationsschritt)

S4 Berechne $t = t_k, c = c_k, s = s_k$ und $\tau = \tau_k$ gemäß (9.5)-(9.7).

S5 Berechne $\mathbf{A}^{(k+1)}$ aus $\mathbf{A}^{(k)}$ gemäß (9.8)-(9.12).

S6 Berechne $\omega_{k+1} = \omega_k - \left(a_{pq}^{(k)} \right)^2$.

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

Aufwand pro Schritt: $\sim 4n$ Add./Sub. + $\sim 3n$ Mult./Div. + 2 Quadratwurzeln.

Bemerkungen: (i) Das Verfahren darf auf dem Platz der Matrix, also auf dem oberen Dreieck, durchgeführt werden.

(ii) Ist der Abbruchtest in **S1** nach k Schritten erfüllt, so gilt nach (9.4)

$$|a_{ii}^{(k)} - \lambda_{l(i)}| \leq \varepsilon.$$

Die Diagonalelemente von $A^{(k)}$ approximieren die Eigenwerte von A mit der vorgegebenen Genauigkeit ε . Man sollte $\varepsilon \geq \text{eps} \|A\|_F$ wählen, da eine höhere Genauigkeit nicht erreichbar ist.

(iii) Falls ω_k klein wird (etwa $k \text{eps} \omega_1$), so sollte man statt der Aufdatierung in Schritt **S3** ω_k direkt aus $A^{(k)}$ berechnen.

(iv) Im KonvergenzFalle gilt $\omega_k \rightarrow 0$. Dann werden die Nichtdiagonalelemente von A_k klein. Das δ aus der zu lösenden quadratischen Gleichung kann dann groß werden. Wegen

$$|\tau| \leq \frac{|s|}{1 + \sqrt{2}/2} \leq |t| \leq \frac{1}{1 + |\delta|}$$

sind τ_k , s_k und t_k betragsmäßig klein, so dass die entsprechenden Korrekturformeln (9.8)-(9.12) gutartig sind.

Für eine optimale Wahl von p und q erkennt man sofort, dass a_{pq} das betragsgrößte Nichtdiagonalelement von A sein müßte. Mit dieser Strategie erhält man das klassische JACOBI-Verfahren. Zur Bestimmung dieses betragsgrößten Elementes ist in jedem Schritt ein großer Suchaufwand notwendig. Als einfachere Variante bietet es sich an, alle Nichtdiagonalelemente der Matrix zyklisch zu durchlaufen, so z.B. zeilenweise gemäß

$$(p, q) = (1, 2), \dots, (1, n), (2, 3), \dots, (2, n), \dots, (n-2, n-1), (n-2, n), (n-1, n).$$

Ein derartiges Verfahren bezeichnet man als zyklisches JACOBI-Verfahren. Beim zyklischen JACOBI-Verfahren kommt es vor, dass zu einem Element a_{pq} eine GIVENS-Rotation durchgeführt wird, obwohl es klein ist. Die Abnahme von ω_k , die dadurch erreicht wird, ist also unwesentlich. Darum modifizieren wir das zyklische JACOBI-Verfahren so, dass wir nur dann eine GIVENS-Rotation G_{pq} durchführen, falls das Element a_{pq} größer als ein vorgegebener Schwellwert ist. Dadurch werden ineffiziente Transformationen vermieden. Dieses Verfahren heißt Schwellwert-JACOBI-Verfahren. Als Schwellwert sollte man das quadratische Mittel der Nichtdiagonalelemente $\sqrt{\omega_k/N}$, $N = n(n-1)/2$ verwenden. Man beachte, dass jede GIVENS-Rotation neue Nichtnullelemente in der Matrix erzeugt, d.h. dass auch nach Durchlaufen aller Nichtdiagonalelemente die Matrix nicht auf Diagonalgestalt transformiert ist. Das JACOBI-Verfahren stellt im allgemeinen einen unendlichen Algorithmus zur Eigenwertberechnung dar. Die Konvergenz des Algorithmus wollen wir nun untersuchen.

9.29. Satz: Das JACOBI-Verfahren werde für eine symmetrische Matrix $\mathbf{A} \in \mathbb{C}^{n \times n}$ mit den Eigenwerten $\lambda_1, \dots, \lambda_n$ in exakter Arithmetik durchgeführt. Es existiere ein $\kappa < 1$, so dass die Pivotwahl in Schritt **S2** so erfolgt, dass jeweils

$$\omega_{k+1} \leq \kappa \omega_k$$

erfüllt ist. Dann gilt für die Folge $\{\mathbf{A}^{(k)}\}$

$$\lim_{k \rightarrow \infty} \mathbf{A}^{(k)} = \mathbf{M} = \text{diag}(\lambda_{m(1)}, \lambda_{m(2)}, \dots, \lambda_{m(n)}) \quad (9.13)$$

mit einer geeigneten Permutation $\{m(1), \dots, m(n)\}$ der Zahlen $\{1, \dots, n\}$. Weiterhin gilt:

(i) Für $i \neq j$ und $k \geq 0$ ist

$$|a_{ij}^{(k)}| \leq \|\mathbf{R}^{(k)}\|_F \leq \sqrt{2\omega_k} \leq (\sqrt{\kappa})^k \sqrt{2\omega_0}.$$

(ii) Im Falle $i = j$ existiert ein $k_0 > 0$, so dass für alle $k \geq k_0$

$$|a_{ii}^{(k)} - \lambda_{m(i)}| \leq \|\mathbf{M}^{(k)} - \mathbf{M}\| \leq \sqrt{2\omega_k} \leq (\sqrt{\kappa})^k \sqrt{2\omega_0}.$$

Beweis: (i) Diese Aussage folgt sofort aus $\omega_{k+1} \leq \kappa \omega_k$.

(ii) Für $i \neq p, q$ gilt

$$|a_{ii}^{(k+1)} - a_{ii}^{(k)}| = 0.$$

Für $i = p$ oder $i = q$ folgt aus den Transformationsformeln (9.8) und (9.12)

$$\begin{aligned} |a_{pp}^{(k+1)} - a_{pp}^{(k)}| &= |a_{pp}^{(k)} - t_k a_{pq}^{(k)} - a_{pp}^{(k)}| \\ &= |t_k| |a_{pq}^{(k)}| \\ |a_{qq}^{(k+1)} - a_{qq}^{(k)}| &= |a_{qq}^{(k)} + t_k a_{pq}^{(k)} - a_{qq}^{(k)}| \\ &= |t_k| |a_{pq}^{(k)}|. \end{aligned}$$

Damit erhält man

$$|a_{ii}^{(k+1)} - a_{ii}^{(k)}| \leq |t_k| |a_{pq}^{(k)}| \leq 1 \cdot (\sqrt{\kappa})^k \sqrt{2\omega_0},$$

und weiter für $r = 0, 1, \dots$

$$\begin{aligned} |a_{ii}^{(k+r+1)} - a_{ii}^{(k)}| &\leq |a_{ii}^{(k+r+1)} - a_{ii}^{(k+r)}| + \dots + |a_{ii}^{(k+1)} - a_{ii}^{(k)}| \\ &\leq \left[(\sqrt{\kappa})^{k+r} + \dots + (\sqrt{\kappa})^k \right] \sqrt{2\omega_0} \\ &= (\sqrt{\kappa})^k \left[(\sqrt{\kappa})^r + \dots + 1 \right] \sqrt{2\omega_0} \\ &\leq \frac{(\sqrt{\kappa})^k}{1 - \sqrt{\kappa}} \sqrt{2\omega_0}. \end{aligned}$$

Die Folge $\{a_{ii}^{(k)}\}$ ist damit CAUCHY-Folge und konvergiert gegen einen Grenzwert μ_i . Damit gilt (9.13) mit $M = \text{diag}(\mu_1, \mu_2, \dots, \mu_n)$. Da alle $A^{(k)}$ orthogonal ähnlich zu A sind, müssen diese μ_i Eigenwerte von A sein. Es existiert also eine Permutation $\{m(1), \dots, m(n)\}$ der Zahlen $\{1, \dots, n\}$, so dass $\mu_i = \lambda_{m(i)}$, $i = 1, \dots, n$. Die Abschätzung in (ii) folgt dann mit Satz 9.23. *

Bemerkungen: (i) Für das klassische JACOBI-Verfahren gilt

$$\omega_k \leq \frac{n(n-1)}{2} \left(a_{pq}^{(k)}\right)^2 = N \left(a_{pq}^{(k)}\right)^2.$$

Daraus folgt

$$\begin{aligned} \omega_{k+1} &= \omega_k - \left(a_{pq}^{(k)}\right)^2 \\ &= \left(1 - \frac{\left(a_{pq}^{(k)}\right)^2}{\omega_k}\right) \omega_k \\ &\leq \left(1 - \frac{1}{N}\right) \omega_k. \end{aligned}$$

Wir erhalten daher $\omega_{k+1} \leq \kappa \omega_k$ mit $\kappa = 1 - \frac{1}{N} < 1$.

(ii) Für das Schwellwert-JACOBI-Verfahren sollte

$$|a_{pq}^{(k)}| \geq \sqrt{\omega_k/N}$$

gelten. Damit folgt wieder

$$\begin{aligned} \omega_{k+1} &= \omega_k - \left(a_{pq}^{(k)}\right)^2 \\ &= \left(1 - \frac{\left(a_{pq}^{(k)}\right)^2}{\omega_k}\right) \omega_k \\ &\leq \left(1 - \frac{1}{N}\right) \omega_k, \end{aligned}$$

daher $\omega_{k+1} \leq \kappa \omega_k$ mit $\kappa = 1 - \frac{1}{N} < 1$.

(iii) Nach Satz 9.29 konvergieren die Elemente von $A^{(k)}$ mindestens linear mit dem Faktor κ gegen Null bzw. gegen die Eigenwerte von A . Da κ nahe bei 1 liegt, folgt

daraus langsame Konvergenz. Um eine vorgegebene Genauigkeit $\varepsilon = K\text{eps}\|\mathbf{A}\|_F$ zu erreichen würde dies eine Abschätzung

$$r = \frac{k}{N} \geq 2 \ln \left(\frac{1}{K\text{eps}} \right)$$

für die Anzahl der benötigten Zyklen ergeben. Praktische Erfahrungen zeigen aber, dass fast immer

$$r = \frac{k}{N} \leq 4 \dots 10$$

gilt. Grund für dieses Verhalten ist die asymptotisch quadratische Konvergenz des JACOBI-Verfahrens.

Als letztes wollen wir noch den Rundungsfehlereinfluss beim JACOBI-Verfahren abschätzen.

9.30. Rundungsfehleranalyse: Das JACOBI-Verfahren werde in Computerarithmetik mit einer vorgegebenen Genauigkeit $\varepsilon = K\text{eps}\|\mathbf{A}\|_F$ und einer Pivotwahl, für die $\omega_{k+1} \leq \kappa\omega_k$ mit $\kappa < 1$ gilt, durchgeführt. Dann bricht das Verfahren nach $k = Nr$ Schritten ab, wobei $r \leq 4 \dots 10$ gilt.

Die berechneten Diagonalelemente $\mu_i = a_{ii}^{(k+1)}$ sind exakte Eigenwerte einer Matrix $\mathbf{A} + \delta\mathbf{A}$, wobei für die symmetrische Störung $\delta\mathbf{A}$

$$\|\delta\mathbf{A}\| \leq \text{eps}(F + K)\|\mathbf{A}\|_F, \quad F \leq 18n^{3/2}r$$

gilt. (Praktisch gilt fast immer $F \leq k/n \approx nr/2$, oft sogar $F \leq 10$.)

Für die berechnete Matrix der Eigenvektoren $\tilde{\mathbf{V}}$ gilt

$$\|\tilde{\mathbf{V}} - \mathbf{V}\| \leq \text{eps}F_1, \quad F_1 \leq 6n^2$$

wobei \mathbf{V} die exakte Matrix der Eigenvektoren von $\mathbf{A} + \delta\mathbf{A}$ ist.

Der Aufwand des JACOBI-Verfahrens ist relativ hoch. Bricht man nach vier Zyklen ab, so benötigt man etwa $\sim 6n^3$ Additionen und Multiplikationen allein zum Berechnen der Eigenwerte und $\sim 12n^3$ zum Lösen des vollständigen Eigenwertproblems. Das QR -Verfahren, das wir später behandeln werden, löst dieselben Aufgaben mit einem Neuntel bzw. einem Drittel des Aufwands. Darum ist die Anwendung des JACOBI-Verfahrens auf Spezialfälle beschränkt, etwa das Berechnen der Eigenwerte mit geringer Genauigkeit oder das Berechnen der Eigenwerte mit voller Genauigkeit für Matrizen mit kleinen Nichtdiagonalelementen. In diesen Fällen kann man hoffen, mit ein bis zwei Zyklen ans Ziel zu kommen. Auch für Probleme kleiner Dimension ($n \leq 10$) ist das JACOBI-Verfahren konkurrenzfähig.

9.4. Die Vektor- und Teilraumiteration

In diesem Abschnitt werden wir iterative Verfahren betrachten, die für eine reelle symmetrische Matrix Eigenvektoren approximieren, die zu den betragsgrößten Eigenwerten gehören. Die Algorithmen selbst haben keine besondere Bedeutung, sie dienen als Grundlage für die inverse Iteration oder den QR -Algorithmus, die leistungsfähige Algorithmen zur Behandlung von Eigenwertproblemen darstellen.

Es seien $\lambda_1, \lambda_2, \dots, \lambda_n$ die gemäß

$$|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|$$

geordneten Eigenwerte der symmetrischen Matrix A . Ferner sei

$$\{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n\}$$

ein zugehöriges Orthonormalsystem von Eigenvektoren. Dann heißen die Eigenwerte

$$\{\lambda_1, \lambda_2, \dots, \lambda_p\}$$

dominant, falls $|\lambda_p| > |\lambda_{p+1}|$ gilt.

Der den dominanten Eigenwerten $\{\lambda_1, \lambda_2, \dots, \lambda_p\}$ zugeordnete Eigenraum

$$\mathcal{S}_p = \mathcal{S}_p(A) = \text{span}(\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p)$$

heißt p -ter dominanter Teilraum von A . Die Eigenvektoren $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_p$ bilden eine orthonormale Basis des \mathcal{S}_p und es gilt $\dim(\mathcal{S}_p) = p$. Weiterhin ist wegen $A\mathbf{u}_j = \lambda_j\mathbf{u}_j \in \mathcal{S}_p$ für $j = 1, 2, \dots, p$

$$A\mathcal{S}_p = \{\mathbf{y} = A\mathbf{x} : \mathbf{x} \in \mathcal{S}_p\} \subseteq \mathcal{S}_p. \quad (9.14)$$

Der Eigenraum \mathcal{S}_p ist somit invariant unter der durch A vermittelten linearen Abbildung. Ein Teilraum \mathcal{S}_p des \mathbb{R}^n , für den (9.14) gilt, heißt auch A -invarianter Teilraum des \mathbb{R}^n . Wir werden nun versuchen, \mathcal{S}_p zu approximieren. Dazu gehen wir von der Eigenwertzerlegung von A aus. Es sei

$$\begin{aligned} A &= U\Lambda U^T \\ &= \left[\lambda_1 \mathbf{u}_1 \mathbf{u}_1^T + \dots + \lambda_p \mathbf{u}_p \mathbf{u}_p^T \right] + \left[\lambda_{p+1} \mathbf{u}_{p+1} \mathbf{u}_{p+1}^T + \dots + \lambda_n \mathbf{u}_n \mathbf{u}_n^T \right]. \end{aligned}$$

Dann gilt

$$\begin{aligned} \mathbf{A}^k &= \mathbf{U} \mathbf{\Lambda}^k \mathbf{U}^T \\ &= \left[\lambda_1^k \mathbf{u}_1 \mathbf{u}_1^T + \cdots + \lambda_p^k \mathbf{u}_p \mathbf{u}_p^T \right] + \left[\lambda_{p+1}^k \mathbf{u}_{p+1} \mathbf{u}_{p+1}^T + \cdots + \lambda_n^k \mathbf{u}_n \mathbf{u}_n^T \right] \\ &= \lambda_p^k \left\{ \left[\left(\frac{\lambda_1}{\lambda_p} \right)^k \mathbf{u}_1 \mathbf{u}_1^T + \cdots + \mathbf{u}_p \mathbf{u}_p^T \right] \right. \\ &\quad \left. + \left[\left(\frac{\lambda_{p+1}}{\lambda_p} \right)^k \mathbf{u}_{p+1} \mathbf{u}_{p+1}^T + \cdots + \left(\frac{\lambda_n}{\lambda_p} \right)^k \mathbf{u}_n \mathbf{u}_n^T \right] \right\}. \end{aligned}$$

Wegen $|\lambda_p| > |\lambda_{p+1}|$ wird die Dominanz der ersten p Terme mit zunehmenden k -Potenzen immer größer. Für hinreichend großes k wird \mathbf{A}^k durch die ersten p Terme repräsentiert. Damit gilt für einen beliebigen Vektor $\mathbf{v} \in \mathbb{R}^n$

$$\mathbf{A}^k \mathbf{v} = \sum_{j=1}^n \lambda_j^k \left[\mathbf{u}_j^T \mathbf{v} \right] \mathbf{u}_j \approx \lambda_1^k \left[\mathbf{u}_1^T \mathbf{v} \right] \mathbf{u}_1 + \cdots + \lambda_p^k \left[\mathbf{u}_p^T \mathbf{v} \right] \mathbf{u}_p$$

falls nicht alle Skalarprodukte $\mathbf{u}_j^T \mathbf{v}, j = 1, \dots, p$, verschwinden. Für großes k liegt somit der Vektor $\mathbf{A}^k \mathbf{v}$ fast in \mathcal{S}_p .

Wenn wir geeignete p linear unabhängige Vektoren $\mathbf{v}_1, \dots, \mathbf{v}_p$ wählen, so werden die Vektoren $\mathbf{A}^k \mathbf{v}_1, \dots, \mathbf{A}^k \mathbf{v}_p$ für genügend großes k einen Teilraum des \mathbb{R}^n aufspannen, der \mathcal{S}_p gut approximiert. Das Berechnen von $\mathbf{A}^k \mathbf{v}$ erfolgt natürlich rekursiv gemäß $\mathbf{A}^{k+1} \mathbf{v} = \mathbf{A} \left(\mathbf{A}^k \mathbf{v} \right)$. Außerdem sind bei der praktischen Realisierung die Vektoren $\mathbf{A}^k \mathbf{v}$ zu normieren, da $\|\mathbf{A}^k \mathbf{v}\|$ im Falle $|\lambda_1| > 1$ unbeschränkt wächst und im Falle $|\lambda_1| < 1$ gegen Null geht.

Wir betrachten nun den einfachsten Fall $p = 1$. Das ist die sogenannte Vektoriteration.

9.31. Basisverfahren der Vektoriteration:

S0 (Initialisierung) Wähle Vektor $\mathbf{v}^{(0)} \in \mathbb{R}^n$ mit $\|\mathbf{v}^{(0)}\|_2 = 1$ und setze $k = 0$.

S1 (Iteration) Berechne $\mathbf{w}^{(k+1)} = \mathbf{A} \mathbf{v}^{(k)}$.

S2 (Normierung) Setze $\mathbf{v}^{(k+1)} = \mathbf{w}^{(k+1)} / \|\mathbf{w}^{(k+1)}\|_2$.

S3 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Aufwand pro Schritt:

- Eine Auswertung Matrix \times Vektor (**S1**)
- $\sim n$ Add./Sub. + $\sim 2n$ Mult./Div. + 1 Quadratwurzel (**S2**)

Bemerkungen: (i) Die Matrix \mathbf{A} wird bei der Iteration im Unterschied zu anderen Verfahren nicht verändert. Man benötigt eigentlich auch nur eine Prozedur zum

Berechnen von $\mathbf{A}\mathbf{x}$. Dies ist besonders bei schwach besetzten Matrizen hoher Dimension von Vorteil.

(ii) Die Vektoriteration wird auch als „VON-MISES-Iteration“ oder als Potenzmethode bezeichnet.

Der folgende Satz liefert uns Aussagen über das Konvergenzverhalten der Vektoriteration.

9.32. Satz: *Es sei λ_1 dominanter Eigenwert der symmetrischen Matrix \mathbf{A} , \mathbf{u}_1 bezeichne einen zugehörigen normierten Eigenvektor, der Startvektor $\mathbf{v}^{(0)}$ genüge der Bedingung*

$$\sigma = \mathbf{u}_1^T \mathbf{v}^{(0)} > 0. \quad (9.15)$$

Dann gilt für die Folge $\{\mathbf{v}^{(k)}\}$, die durch die Vektoriteration erzeugt wird,

$$0 \leq \tan \varphi_k \leq \kappa \tan \varphi_{k-1} \leq \kappa^k \tan \varphi_0 = \kappa^k \sqrt{1 - \sigma^2} / \sigma \quad (9.16)$$

mit $\kappa = |\lambda_1 / \lambda_2|$ und $\varphi_k = \angle(\mathbf{u}_1, \vartheta_k \mathbf{v}^{(k)}) \in [0, \pi/2)$, $\vartheta_k = (\operatorname{sgn}(\lambda_1))^k \in \{1, -1\}$.

Beweis: Es sei $\mathcal{S}_1 = \operatorname{span}(\mathbf{u}_1)$ und $\mathcal{Z}_2 = \operatorname{span}(\mathbf{u}_2, \dots, \mathbf{u}_n) = \mathcal{S}_1^\perp$. Dann lässt sich der Startvektor $\mathbf{v}^{(0)}$ in Komponenten bezüglich \mathcal{S}_1 und \mathcal{Z}_2 zerlegen:

$$\mathbf{v}^{(0)} = \vartheta_0 \mathbf{v}^{(0)} = \mathbf{u}_1 \cos \varphi_0 + \mathbf{z}_0 \sin \varphi_0, \quad \mathbf{z}_0 \in \mathcal{Z}_2, \quad \|\mathbf{z}_0\| = 1. \quad (9.17)$$

Wegen (9.15) gilt

$$\sigma = \cos \varphi_0 > 0,$$

und darum $0 \leq \varphi_0 < \pi/2$ und $\tan \varphi_0 = \sqrt{1 - \sigma^2} / \sigma \geq 0$. Multiplizieren wir 9.17 von links mit $\operatorname{sgn}(\lambda_1) \mathbf{A}$, so erhalten wir

$$\begin{aligned} \lambda_1 \mathbf{A} \vartheta_0 \mathbf{v}^{(0)} &= \vartheta_1 \mathbf{A} \mathbf{v}^{(0)} = \vartheta_1 \mathbf{w}^{(1)} \\ &= \operatorname{sgn}(\lambda_1) (\lambda_1 \mathbf{u}_1 \cos \varphi_0 + \mathbf{A} \mathbf{z}_0 \sin \varphi_0), \quad \mathbf{A} \mathbf{z}_0 \in \mathcal{Z}_2. \end{aligned}$$

Daraus folgt

$$\|\mathbf{w}^{(1)}\|^2 = \lambda_1^2 \cos^2 \varphi_0 + \|\mathbf{A} \mathbf{z}_0\|^2 \sin^2 \varphi_0 \geq \lambda_1^2 \cos^2 \varphi_0,$$

daher $\|\mathbf{w}^{(1)}\| \geq |\lambda_1| \cos \varphi_0 > 0$. Der Schritt **S2** im Algorithmus ist durchführbar und es gilt

$$\vartheta_1 \mathbf{v}^{(1)} = \vartheta_1 \mathbf{w}^{(1)} / \|\mathbf{w}^{(1)}\|_2 = \mathbf{u}_1 \cos \varphi_1 + \mathbf{z}_1 \sin \varphi_1, \quad \mathbf{z}_1 \in \mathcal{Z}_2, \quad \|\mathbf{z}_1\|_2 = 1,$$

wobei

$$\cos \varphi_1 = \frac{|\lambda_1| \cos \varphi_0}{\|\mathbf{w}^{(1)}\|}, \quad \sin \varphi_1 = \frac{\|\mathbf{A}z_0\| \sin \varphi_0}{\|\mathbf{w}^{(1)}\|}$$

und

$$z_1 = \begin{cases} \frac{\operatorname{sgn}(\lambda_1) \mathbf{A}z_0}{\|\mathbf{A}z_0\|} & \text{für } \mathbf{A}z_0 \neq \mathbf{o} \\ \mathbf{u}_2 & \text{für } \mathbf{A}z_0 = \mathbf{o}. \end{cases}$$

Daraus folgt $0 \leq \varphi_1 < \pi/2$ und

$$0 \leq \tan \varphi_1 = \frac{\|\mathbf{A}z_0\|}{|\lambda_1|} \tan \varphi_0 \leq \frac{|\lambda_2|}{|\lambda_1|} \tan \varphi_0 = \kappa \tan \varphi_0.$$

(9.16) gilt daher für $k = 1$. Der Rest folgt durch vollständige Induktion. *

Bemerkungen: (i) Der Vorzeichenfaktor ϑ_k hat nur beweistechnische Bedeutung.

(ii) $|\tan(\angle(\mathbf{u}, \mathbf{v}))|$ ist als „Entfernung“ zweier Richtungen interpretierbar. Orthogonale Richtungen sind unendlich weit voneinander entfernt, der Abstand der Richtungen paralleler Vektoren ist gleich Null.

(iii) Es gilt

$$\begin{aligned} \|\vartheta_k \mathbf{v}^{(k)} - \mathbf{u}_1\| &= \sqrt{(\vartheta_k \mathbf{v}^{(k)} - \mathbf{u}_1)^T (\vartheta_k \mathbf{v}^{(k)} - \mathbf{u}_1)} \\ &= \sqrt{\|\mathbf{v}^{(k)}\|^2 - 2\mathbf{u}_1^T \mathbf{v}^{(k)} + \|\mathbf{u}_1\|^2} \\ &= \sqrt{2(1 - \cos \varphi_k)} \\ &= 2 \sin(\varphi_k/2) \leq \tan \varphi_k. \end{aligned}$$

Wegen $\varphi_k \rightarrow 0$ folgt somit $\vartheta_k \mathbf{v}^{(k)} \rightarrow \mathbf{u}_1$.

(iv) Normalerweise möchte man neben einer Eigenvektornäherung auch eine entsprechende Approximation für den zugehörigen Eigenwert haben. Die beste Eigenwertnäherung zu einem gegebenen Eigenvektor ist nach Satz 9.26 der RAYLEIGH-Quotient:

$$\varrho_k = RQ(\mathbf{v}^{(k)}) = \mathbf{v}^{(k)T} \mathbf{A} \mathbf{v}^{(k)} = \mathbf{v}^{(k)T} \mathbf{w}^{(k+1)}.$$

Stellen wir $\mathbf{v}^{(k)}$ in der Form

$$\mathbf{v}^{(k)} = \mathbf{u}_1 \cos \varphi_k + z_k \sin \varphi_k, \quad z_k \in \mathcal{Z}_2, \quad \|z_k\| = 1$$

dar, so ergibt sich

$$\mathbf{A} \left(\vartheta_k \mathbf{v}^{(k)} \right) = \mathbf{u}_1 \lambda_1 \cos \varphi_k + \mathbf{A} \mathbf{z}_k \sin \varphi_k, \quad \mathbf{A} \mathbf{z}_k \in \mathcal{Z}_2$$

und weiter

$$\begin{aligned} \varrho_k = RQ(\vartheta_k \mathbf{v}^{(k)}) &= (\mathbf{u}_1 \cos \varphi_k + \mathbf{z}_k \sin \varphi_k)^T (\mathbf{u}_1 \lambda_1 \cos \varphi_k + \mathbf{A} \mathbf{z}_k \sin \varphi_k) \\ &= \lambda_1 \cos^2 \varphi_k + \mathbf{z}_k^T \mathbf{A} \mathbf{z}_k \sin^2 \varphi_k \\ &= \lambda_1 - (\lambda_1 - \mathbf{z}_k^T \mathbf{A} \mathbf{z}_k) \sin^2 \varphi_k. \end{aligned}$$

Nun gilt

$$\begin{aligned} |\lambda_1 - \mathbf{z}_k^T \mathbf{A} \mathbf{z}_k| &\leq |\lambda_1| + |\mathbf{z}_k^T \mathbf{A} \mathbf{z}_k| \\ &\leq |\lambda_1| + \|\mathbf{z}_k\| \|\mathbf{A} \mathbf{z}_k\| \\ &\leq |\lambda_1| + |\lambda_2| \\ &\leq 2|\lambda_1| = 2\|\mathbf{A}\|. \end{aligned}$$

Damit folgt endlich

$$|\varrho_k - \lambda_1| \leq 2\|\mathbf{A}\| \sin^2 \varphi_k \leq 2\|\mathbf{A}\| \tan^2 \varphi_k \leq 2\|\mathbf{A}\| \kappa^{2k} \tan^2 \varphi_0.$$

Die Folge der RAYLEIGH-Quotienten konvergiert linear mit dem Faktor κ^2 , also schneller als die Folge der Eigenvektorapproximationen.

Bei der praktischen Realisierung der Vektoriteration gilt wegen den auftretenden Rundungsfehlern nur noch der folgende Satz, den wir ohne Beweis angeben.

9.33. Satz: *Es sei λ_1 dominanter Eigenwert der symmetrischen Matrix \mathbf{A} :*

$$\kappa = |\lambda_2/\lambda_1| < 1,$$

\mathbf{u}_1 bezeichne einen zugehörigen normierten Eigenvektor. Die Vektoriteration werde in Computerarithmetik so durchgeführt, dass die berechneten Vektoren $\{\mathbf{w}^{(k+1)}\}$ der Beziehung

$$\mathbf{w}^{(k+1)} = (\mathbf{A} + \delta \mathbf{A}_k) \mathbf{v}^{(k)}, \quad \|\delta \mathbf{A}_k\| \leq \text{eps} F \|\mathbf{A}\|$$

genügen. Mit $\varepsilon = \text{eps}(F + n/2)$ gelte

$$\kappa + 10 < 1, \quad |\mathbf{u}_1^T \mathbf{v}^{(0)}| > \bar{\varepsilon} = \varepsilon / (1 - \kappa - \varepsilon). \quad (9.18)$$

Dann gilt für die den berechneten Vektoren $\{\mathbf{v}^{(k)}\}$ zugeordneten Winkel $\{\varphi_k\}$ die Abschätzung

$$0 \leq \tan \varphi_k \leq \tau_k, \quad \tau_k \geq \tau_{k+1}, \quad \lim_{k \rightarrow \infty} \tau_k = \bar{\varepsilon}.$$

Ist k_0 ein hinreichend großer Index mit $\tan \varphi_{k_0} < 1$, so gilt für alle $k \geq k_0$

$$\tan \varphi_{k+1} \leq \hat{\kappa} \tan \varphi_k + \hat{\varepsilon} \leq \hat{\kappa}^{k-k_0+1} \tan \varphi_{k_0} + \frac{\hat{\varepsilon}}{1-\hat{\kappa}}$$

mit

$$\hat{\kappa} = \frac{\kappa + \varepsilon}{1 - \sqrt{2\varepsilon}}, \quad \hat{\varepsilon} = \frac{\varepsilon}{1 - \sqrt{2\varepsilon}}.$$

Für alle $k \geq 1$ gilt weiterhin

$$\|\vartheta_k \mathbf{v}^{(k)} - \mathbf{u}_1\| \leq \tan \varphi_k + \text{eps} \frac{n}{2},$$

und die gemäß

$$\varrho_k = \mathbf{v}^{(k)T} \mathbf{w}^{(k+1)}$$

berechneten Eigenwertnäherungen genügen der Abschätzung

$$|\varrho_k - \lambda_1| \leq \|\mathbf{A}\| [2 \sin^2 \varphi_k + \text{eps}(F + 2n)].$$

Bemerkungen: (i) Durch die Bedingungen (9.18) wird einerseits gesichert, dass $|\lambda_1|$ hinreichend größer ist als $|\lambda_2|$, und andererseits wird verhindert, dass \mathbf{u}_1 und der Startvektor $\mathbf{v}^{(0)}$ fast orthogonal sind. Die zweite Bedingung ist nur von theoretischer Bedeutung. Praktisch konvergiert das Verfahren auch, falls \mathbf{u}_1 und $\mathbf{v}^{(0)}$ exakt orthogonal sind, da durch Rundungsfehler immer Komponenten in den Vektoren $\mathbf{v}^{(k)}$, $k > 1$, erzeugt werden, die parallel zu \mathbf{u}_1 sind.

(ii) Für vollbesetztes \mathbf{A} wäre die Konstante $F = n^{3/2}$. Ist \mathbf{A} jedoch wie in vielen Anwendungen schwach besetzt, so ergeben sich wesentlich kleinere Werte für F . Wir wollen nun den Fall betrachten, bei dem mehrere Vektoren gleichzeitig iteriert werden, um einen dominanten Teilraum zu approximieren. Das ist die sogenannte Teilraumiteration.

9.34. Basisverfahren der Teilraumiteration:

S0 Wähle einen Teilraum $\mathcal{Y}^{(0)} \subset \mathbb{R}^n$ mit $\dim(\mathcal{Y}^{(0)}) = p$ und setze $k = 0$.

S1 (Iteration) Definiere

$$\mathcal{Y}^{(k+1)} = \mathbf{A}\mathcal{Y}^{(k)} = \left\{ \mathbf{A}\mathbf{y} \mid \mathbf{y} \in \mathcal{Y}^{(k)} \right\}.$$

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

Um die Konvergenz des Algorithmus zu beschreiben, benötigen wir ein Maß für die Entfernung zweier Teilräume. Es sei $\mathcal{S} \subset \mathbb{R}^n$ ein Teilraum und $\mathbf{y} \in \mathbb{R}^n$ ein Vektor ungleich dem Nullvektor. Die Größe

$$\angle(\mathbf{y}, \mathcal{S}) = \min \{ \angle(\mathbf{y}, \mathbf{s}) \mid \mathbf{s} \in \mathcal{S}, \mathbf{s} \neq \mathbf{o} \}$$

heißt Winkel zwischen dem Vektor \mathbf{y} und dem Teilraum \mathcal{S} . Es sei $\mathcal{Y} \subset \mathbb{R}^n$ ein weiterer Teilraum. Dann heißt die Größe

$$\angle(\mathcal{Y}, \mathcal{S}) = \max \{ \angle(\mathbf{y}, \mathcal{S}) \mid \mathbf{y} \in \mathcal{Y}, \mathbf{y} \neq \mathbf{o} \}$$

Winkel zwischen den Teilräumen \mathcal{Y} und \mathcal{S} . Es gilt offensichtlich stets

$$0 \leq \angle(\mathcal{Y}, \mathcal{S}) \leq \pi/2.$$

Man beachte aber, dass für den Winkel zwischen zwei Teilräumen im allgemeinen **nicht** $\angle(\mathcal{Y}, \mathcal{S}) = \angle(\mathcal{S}, \mathcal{Y})$ gilt. Um dies einzusehen, betrachte man nur eine Gerade und eine Ebene im \mathbb{R}^3 . Im Falle $\dim(\mathcal{Y}) = \dim(\mathcal{S})$ gilt

$$\angle(\mathcal{Y}, \mathcal{S}) = \angle(\mathcal{S}, \mathcal{Y}) = \angle(\mathcal{Y}^\perp, \mathcal{S}^\perp).$$

Falls die Teilräume \mathcal{Y} und \mathcal{S} gleiche Dimension haben, ist $\tan(\angle(\mathcal{Y}, \mathcal{S}))$ zur Charakterisierung der Entfernung der beiden Teilräume verwendbar. Es gilt

$$\tan(\angle(\mathcal{Y}, \mathcal{S})) = 0$$

genau dann, wenn $\mathcal{Y} = \mathcal{S}$, und

$$\tan(\angle(\mathcal{Y}, \mathcal{S})) = \infty,$$

falls einer der Teilräume eine Richtung enthält, die orthogonal zum anderen ist. Nun formulieren wir einen Konvergenzsatz.

9.35. Satz: *Es seien $\{\lambda_1, \dots, \lambda_p\}$ dominante Eigenwerte der symmetrischen Matrix \mathbf{A} . $\mathcal{S}_p \subset \mathbb{R}^n$ sei der zugehörige p -te dominante Teilraum. Der Teilraum $\mathcal{Y}^{(0)} \subset \mathbb{R}^n$ genüge den Bedingungen*

$$\dim(\mathcal{Y}^{(0)}) = p, \quad \sigma = \cos(\angle(\mathcal{Y}^{(0)}, \mathcal{S}_p)) > 0.$$

Dann gilt für die durch die Teilraumiteration erzeugte Folge $\{\mathcal{Y}^{(k)}\}$

1. $\dim(\mathcal{Y}^{(k)}) = p, \quad k = 1, 2, \dots,$

2.

$$\tan \varphi_k \leq \kappa \tan \varphi_{k-1} \leq \kappa^k \tan \varphi_0 = \kappa^k \sqrt{1 - \sigma^2} / \sigma$$

mit

$$\kappa = |\lambda_{p+1} / \lambda_p| < 1, \quad \varphi_k = \sphericalangle(\mathbf{y}^{(k)}, \mathcal{S}_p).$$

Beweis: Es ist $\mathbf{y}^{(k)}$ in Komponenten bezüglich \mathcal{S}_p und

$$\mathcal{Z}_{p+1} = \text{span}(\mathbf{u}_{p+1}, \dots, \mathbf{u}_n) = \mathcal{S}_p^\perp$$

zu zerlegen. Der Rest verläuft analog zum Beweis des Satzes 9.18. *

Zur praktischen Realisierung der Teilraumiteration liegt es nahe, $\mathbf{y}^{(k)}$ durch p linear unabhängige Vektoren $\mathbf{v}_1^{(k)}, \dots, \mathbf{v}_p^{(k)}$ festzulegen. Faßt man diese Vektoren zu einer Matrix

$$\mathbf{V}^{(k)} = \left(\mathbf{v}_1^{(k)}, \dots, \mathbf{v}_p^{(k)} \right) \in \mathbb{R}^{n \times p}$$

zusammen, so besteht ein Iterationsschritt der Teilraumiteration in der Transformation

$$\mathbf{W}^{(k+1)} = \mathbf{A} \mathbf{V}^{(k)}.$$

Die Spalten der Matrix $\mathbf{W}^{(k+1)} \in \mathbb{R}^{n \times p}$ spannen dann den Teilraum $\mathcal{Y}^{(k+1)}$ auf. Zur Vermeidung von Über- und Unterlauf ist $\mathbf{W}^{(k+1)}$ zu normieren:

$$\mathbf{V}^{(k+1)} = \mathbf{W}^{(k+1)} \text{diag} \left(\frac{1}{\|\mathbf{w}_1^{(k+1)}\|}, \dots, \frac{1}{\|\mathbf{w}_p^{(k+1)}\|} \right).$$

Für großes k konvergieren alle Vektoren $\mathbf{v}_i^{(k)}$ gegen $\pm \mathbf{u}_1$ falls $\mathbf{u}_1^T \mathbf{v}_i^{(0)} \neq 0$ gilt. Die Vektoren $\mathbf{v}_1^{(k)}, \dots, \mathbf{v}_p^{(k)}$ werden mit wachsendem k immer stärker linear abhängig. Darum sollte man mit paarweise orthogonalen Vektoren starten und diese immer wieder orthonormieren. Weiterhin ist es günstig, mit einem Teilraum $\mathcal{Y}^{(0)}$ zu starten, dessen Dimension größer als p ist. Wir erhalten so den folgenden Algorithmus.

9.36. Teilraumiteration mit orthonormalen Basen:

S0 (Initialisierung) Wähle ein m mit $p \leq m \leq n$ und eine spaltenorthonormale Matrix $\mathbf{V}^{(0)} \in \mathbb{R}^{n \times m}$. Setze $k = 0$.

S1 (Iteration) Berechne $\mathbf{W}^{(k+1)} = \mathbf{A} \mathbf{V}^{(k)}$.

S2 (Orthogonalisierung) Bestimme eine spaltenorthonormale Matrix

$$\mathbf{V}^{(k+1)} \in \mathbb{R}^{n \times m}$$

und eine obere Dreiecksmatrix $\mathbf{R}^{(k+1)} \in \mathbb{R}^{m \times m}$ mit

$$\mathbf{W}^{(k+1)} = \mathbf{V}^{(k+1)} \mathbf{R}^{(k+1)}.$$

S3 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Aufwand pro Schritt:

- m Auswertungen Matrix \times Vektor (**S1**)
- $\sim K_1 nm^2$ Add./Mult. (**S2**)
($K_1 = 1$ für SCHMIDT'sches Orthogonalisierungsverfahren)

Bemerkungen: (i) Falls die Voraussetzungen von Satz 9.35 mit m statt p erfüllt sind, gelten natürlich alle entsprechenden Aussagen. In exakter Arithmetik wird dann durch $\mathbf{W}^{(k+1)}$ der m -te dominante Teilraum approximiert. Für

$$\mathcal{Y}^{(k+1)} = \text{span} \left(\mathbf{W}^{(k+1)} \right)$$

gilt

$$\dim \left(\mathcal{Y}^{(k+1)} \right) = \text{rg} \left(\mathbf{W}^{(k+1)} \right) = m.$$

Der Schritt **S2** im obigen Algorithmus ist daher durchführbar.

(ii) Partitioniert man die Matrizen $\mathbf{W}^{(k)}$, $\mathbf{V}^{(k)}$ und $\mathbf{R}^{(k)}$ gemäß

$$\begin{aligned} \mathbf{W}^{(k)} &= \left(\mathbf{W}_1^{(k)}, \mathbf{W}_2^{(k)} \right), & \mathbf{W}_1^{(k)} &\in \mathbb{R}^{n \times p}, & \mathbf{W}_2^{(k)} &\in \mathbb{R}^{n \times (m-p)}, \\ \mathbf{V}^{(k)} &= \left(\mathbf{V}_1^{(k)}, \mathbf{V}_2^{(k)} \right), & \mathbf{V}_1^{(k)} &\in \mathbb{R}^{n \times p}, & \mathbf{V}_2^{(k)} &\in \mathbb{R}^{n \times (m-p)} \end{aligned}$$

und

$$\mathbf{R}^{(k)} = \begin{pmatrix} \mathbf{R}_{11}^{(k)} & \mathbf{R}_{12}^{(k)} \\ \mathbf{o} & \mathbf{R}_{22}^{(k)} \end{pmatrix}$$

mit

$$\mathbf{R}_{11}^{(k)} \in \mathbb{R}^{p \times p}, \quad \mathbf{R}_{12}^{(k)} \in \mathbb{R}^{p \times (m-p)}, \quad \mathbf{R}_{22}^{(k)} \in \mathbb{R}^{(m-p) \times (m-p)},$$

dann gilt

$$\mathbf{W}_1^{(k+1)} = \mathbf{A} \mathbf{V}_1^{(k)}$$

und

$$\mathbf{W}_1^{(k+1)} = \mathbf{V}_1^{(k+1)} \mathbf{R}_{11}^{(k+1)}.$$

Die Spalten von $\mathbf{V}_1^{(k)}$ bilden eine orthonormale Basis von $\mathcal{Y}_1^{(k)} = \text{span}(\mathbf{W}_1^{(k)})$, also desjenigen Teilraums, der durch Iteration aus $\mathcal{Y}_1^{(0)} = \text{span}(\mathbf{V}_1^{(0)})$ entsteht. Sind die ersten p Eigenwerte dominant, so konvergiert damit die Folge $\{\mathcal{Y}_1^{(k)}\}$ gegen den p -ten dominanten Teilraum \mathcal{S}_p .

(iii) Für die Computerrealisierung lassen sich zu Satz 9.33 analoge Aussagen beweisen. Wesentlich für die Konvergenz des Verfahrens ist dabei, dass $\kappa = |\lambda_{p+1}/\lambda_p|$ hinreichend kleiner als 1 ist.

(iv) An das Orthogonalisierungsverfahren in Schritt **S2** brauchen keine besonders hohen Qualitätsanforderungen gestellt werden, da nicht die paarweise Orthogonalität der Spalten von $\mathbf{V}^{(k)}$ für die Konvergenz des Verfahrens notwendig ist, sondern nur ihre hinreichend starke lineare Unabhängigkeit.

Mit den Matrizen $\mathbf{V}^{(k)}$ hat man orthogonale Basen der Teilräume $\mathcal{Y}^{(k)}$ berechnet, die Approximationen des p -ten dominanten Teilraums \mathcal{S}_p darstellen, sofern

$$\dim(\mathbf{V}^{(k)}) = p \text{ und } \kappa = |\lambda_{p+1}/\lambda_p| < 1$$

gilt. Die Spalten von $\mathbf{V}^{(k)}$ stellen nicht die besten Näherungen für die ersten p Eigenvektoren dar. Außerdem stellt sich natürlich die Frage nach den zugehörigen Eigenwertapproximationen. Wir haben die Aufgabe, zu einer gegebenen Näherung des p -ten dominanten Teilraums entsprechende Näherungen $(\mu_i, \mathbf{z}_i), i = 1, \dots, p$, für die ersten p Eigenpaare der Matrix \mathbf{A} zu konstruieren. Aus den Approximationen μ_i und \mathbf{z}_i bilden wir die Matrizen

$$\mathbf{M} = \text{diag}(\mu_1, \dots, \mu_p) \in \mathbb{R}^{p \times p}, \quad \mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_p) \in \mathbb{R}^{n \times p}.$$

Wären nun die Paare $(\mu_i, \mathbf{z}_i), i = 1, \dots, p$, exakte Eigenpaare von \mathbf{A} , so würde $\mathbf{AZ} = \mathbf{ZM}$ gelten. Als Maß für die Güte der Approximation der Eigenpaare von \mathbf{A} durch \mathbf{M} und \mathbf{Z} verwenden wir darum die Norm der Residuumsmatrix $\mathbf{R} = \mathbf{AZ} - \mathbf{ZM}$. Eine in diesem Sinne beste Approximation der ersten p Eigenpaare einer symmetrischen Matrix \mathbf{A} , bei einer durch die spaltenorthonormale Matrix $\mathbf{V} \in \mathbb{R}^{n \times p}$ gegebenen Approximation des p -ten dominanten Teilraums \mathcal{S}_p , erhalten wir durch Lösen der Aufgabe

$$\min \{ \|\mathbf{AZ} - \mathbf{ZM}\|_F \mid \mathbf{M} \in \mathcal{D}_p, \mathbf{Z} \in \mathcal{OB}(\text{span}(\mathbf{V})) \}$$

wobei \mathcal{D}_p die Menge aller p -dimensionalen Diagonalmatrizen ist und $\mathcal{OB}(\mathcal{Y})$ für die Menge aller orthonormalen Matrizen \mathbf{Z} steht, deren Spalten eine Orthonormalbasis des Teilraums \mathcal{Y} bilden. Eine Lösung dieses Problems liefert der folgende Satz.

9.37. Satz: Gegeben seien eine symmetrische Matrix \mathbf{A} und eine spaltenorthonormale (n, p) -Matrix \mathbf{V} . Weiterhin sei $\mathcal{Y} = \text{span}(\mathbf{V})$ der Teilraum, der durch die Spalten von \mathbf{V} aufgespannt wird. Aus \mathbf{V} und \mathbf{A} werde die symmetrische (p, p) -Matrix

$$\mathbf{P} = \mathbf{V}^T \mathbf{A} \mathbf{V}$$

gebildet. Es sei

$$\mathbf{P} = \mathbf{X} \mathbf{M} \mathbf{X}^T$$

die zugehörige Eigenwertzerlegung von \mathbf{P} :

$$\mathbf{M} = \text{diag}(\mu_1, \dots, \mu_p),$$

$$\mathbf{X} \in \mathbb{R}^{p \times p}, \quad \text{orthogonal.}$$

Dann löst das Paar (\mathbf{M}, \mathbf{Z}) mit $\mathbf{Z} = \mathbf{V} \mathbf{X}$ die Aufgabe

$$\min \{ \|\mathbf{A} \mathbf{Z} - \mathbf{Z} \mathbf{M}\|_F \mid \mathbf{M} \in \mathcal{D}_p, \mathbf{Z} \in \mathcal{OB}(\text{span}(\mathbf{V})) \}.$$

Beweis: Aus $\dim(\mathbf{Z}) = p$, $\mathcal{Y} = \text{span}(\mathbf{Z}) = \text{span}(\mathbf{V})$ und der Spaltenorthogonalität von \mathbf{Z} folgt, dass \mathbf{Z} die Darstellung $\mathbf{Z} = \mathbf{V} \tilde{\mathbf{X}}$ mit $\tilde{\mathbf{X}} \in \mathbb{R}^{p \times p}$ besitzt. Weiter folgt $\mathbf{I} = \mathbf{Z}^T \mathbf{Z} = \tilde{\mathbf{X}}^T \mathbf{V}^T \mathbf{V} \tilde{\mathbf{X}} = \tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$, also die Orthogonalität von $\tilde{\mathbf{X}}$. Damit gilt

$$\begin{aligned} \|\mathbf{R}\|_F &= \|\mathbf{A} \mathbf{Z} - \mathbf{Z} \mathbf{M}\|_F \\ &= \|(\mathbf{A} \mathbf{Z} - \mathbf{Z} \mathbf{M}) \tilde{\mathbf{X}}^T\|_F \\ &= \|\mathbf{A} \mathbf{V} \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T - \mathbf{V} \tilde{\mathbf{X}} \mathbf{M} \tilde{\mathbf{X}}^T\|_F \\ &= \|\mathbf{A} \mathbf{V} - \mathbf{V} \tilde{\mathbf{P}}\|_F, \quad \tilde{\mathbf{P}} = \tilde{\mathbf{X}} \mathbf{M} \tilde{\mathbf{X}}^T. \end{aligned}$$

Wir ergänzen \mathbf{V} zu einer orthogonalen (n, n) -Matrix $\tilde{\mathbf{V}} = (\mathbf{V}, \hat{\mathbf{V}})$. Es ergibt sich weiter

$$\begin{aligned} \|\mathbf{R}\|_F &= \|\tilde{\mathbf{V}}^T (\mathbf{A} \mathbf{V} - \mathbf{V} \tilde{\mathbf{P}})\|_F \\ &= \left\| \begin{pmatrix} \mathbf{V}^T \\ \hat{\mathbf{V}}^T \end{pmatrix} (\mathbf{A} \mathbf{V} - \mathbf{V} \tilde{\mathbf{P}}) \right\|_F \\ &= \left\| \begin{pmatrix} \mathbf{V}^T \mathbf{A} \mathbf{V} - \tilde{\mathbf{P}} \\ \hat{\mathbf{V}}^T \mathbf{A} \mathbf{V} \end{pmatrix} \right\|_F \quad \text{wegen} \quad \mathbf{V}^T \mathbf{V} = \mathbf{I}, \quad \hat{\mathbf{V}}^T \mathbf{V} = \mathbf{o}. \end{aligned}$$

Damit erhalten wir

$$\|\mathbf{R}\|_F = \sqrt{\|\mathbf{V}^T \mathbf{A} \mathbf{V} - \tilde{\mathbf{P}}\|_F^2 + \|\hat{\mathbf{V}}^T \mathbf{A} \mathbf{V}\|_F^2}$$

und es gilt

$$\begin{aligned} \min \{ \|\mathbf{R}\|_F \} &= \min \{ \|\mathbf{AZ} - \mathbf{ZM}\|_F \mid \mathbf{M} \in \mathcal{D}_p, \mathbf{Z} \in \mathcal{OB}(\text{span}(\mathbf{V})) \} \\ &\geq \|\hat{\mathbf{V}}^T \mathbf{AV}\|_F \end{aligned}$$

und

$$\min \{ \|\mathbf{R}\|_F \} = \|\hat{\mathbf{V}}^T \mathbf{AV}\|_F$$

genau dann, wenn

$$\tilde{\mathbf{P}} = \mathbf{V}^T \mathbf{AV}.$$

Ist $\tilde{\mathbf{P}} = \mathbf{XMX}^T$ die Eigenwertzerlegung von $\tilde{\mathbf{P}}$, so löst das Paar $(\mathbf{M}, \mathbf{VX})$ die vorgegebene Minimierungsaufgabe. *

Die in diesem Satz beschriebene Zuordnung

$$(\mathbf{A}, \mathbf{V}) \rightarrow (\mathbf{M}, \mathbf{Z})$$

wird RAYLEIGH-RITZ-Algorithmus (kurz RR-Algorithmus) genannt. Die Größen μ_i und \mathbf{z}_i heißen RITZsche Eigenwerte bzw. Eigenvektoren der Matrix \mathbf{A} bezüglich des Teilraums $\mathcal{Y} = \text{span}(\mathbf{V})$. Die Matrix \mathbf{P} mit den Elementen $\pi_{ij} = \mathbf{v}^{(i)T} \mathbf{A} \mathbf{v}^{(j)}$ stellt eine Verallgemeinerung des RAYLEIGH-Quotienten dar und wird Projektion von \mathbf{A} auf den Teilraum $\mathcal{Y} = \text{span}(\mathbf{V})$ genannt.

9.38. RAYLEIGH-RITZ-Algorithmus:

Es seien die symmetrische Matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ und die spaltenorthonormale Matrix $\mathbf{V} \in \mathbb{R}^{n \times p}$ gegeben.

- S0**
- Berechne $\mathbf{B} = \mathbf{AV} \in \mathbb{R}^{n \times p}$.
 - Berechne $\mathbf{P} = \mathbf{V}^T \mathbf{B} \in \mathbb{R}^{p \times p}$.

S1 Berechne die Eigenwertzerlegung $\mathbf{P} = \mathbf{XMX}^T$.

S2 Berechne $\mathbf{Z} = \mathbf{VX}$.

Aufwand:

- m Auswertungen Matrix \times Vektor + $\sim np^2/2$ Add./Mult. (**S0**)
- $\sim K_2 p^3$ Add./Mult. (**S1**) (K_2 hängt vom angewendeten Verfahren ab)
- $\sim np^3$ Add./Mult. (**S2**)

Bemerkungen: (i) Ist \mathcal{Y} invarianter Teilraum von \mathbf{A} , so gilt

$$\text{span}(\mathbf{AV}) \subseteq \text{span}(\mathbf{V}) = \mathcal{Y}.$$

Es gibt daher eine Matrix $C \in \mathbb{R}^{p \times p}$ mit $AV = VC$. Damit folgt

$$P = V^T AV = V^T VC = C.$$

Es gilt also $AV = VP$. (Diese Beziehung folgt nicht einfach aus $P = V^T AV$, da im allgemeinen nicht $VV^T = I$ gilt.) Aus der letzten Beziehung folgt $AVX = VPX$, und weiter $AZ - ZM = o$. Für einen invarianten Teilraum verschwindet die Residuumsmatrix $R = AZ - ZM$. Die RITZschen Eigenpaare sind dann exakte Eigenpaare der Matrix A .

(ii) Sind die Spalten von V gute Näherungen für die gesuchten Eigenvektoren, so ist P fast diagonal. Zum Berechnen der Eigenwertzerlegung im Schritt **S2** bietet sich damit das JACOBI-Verfahren an.

Für den Fall, dass \mathcal{Y} kein invarianter Teilraum ist, geben wir noch den folgenden Satz ohne Beweis an, der Aussagen über die Güte der RITZschen Näherungen in Abhängigkeit von der Residuumsmatrix R macht.

9.39. Satz: Für eine symmetrische (n, n) -Matrix A und eine spaltenorthonormale (n, p) -Matrix V seien durch

$$M = \text{diag}(\mu_1, \dots, \mu_p), \quad Z = (z_1, \dots, z_p)$$

die in exakter Arithmetik berechneten RITZschen Eigenwerte und Eigenvektoren gegeben.

$$R = AZ - ZM = (r_1, \dots, r_p)$$

sei die zugehörige Residuumsmatrix. Dann gilt

(i) Es gibt symmetrische Störungen δA_j mit

$$(A + \delta A_j)z_j = \mu_j z_j, \quad \|\delta A_j\|_2 \leq \|r_j\|_2, \quad j = 1, \dots, p.$$

(ii) Es gibt Eigenwerte $\lambda_{l(1)}, \dots, \lambda_{l(p)}$ von A mit

$$|\mu_j - \lambda_{l(j)}| \leq \|R\|_2, \quad j = 1, \dots, p,$$

$$\sqrt{\sum_{j=1}^p (\mu_j - \lambda_{l(j)})^2} \leq \sqrt{2} \|R\|_F.$$

9.5. Die Inverse Iteration nach Wielandt

Die Vektor- und die Teilraumiteration lieferten nur Approximationen für Eigenvektoren bzw. für Eigenräume, die zu dominanten Eigenwerten gehören. Die Konvergenz der Algorithmen ist im allgemeinen langsam. In vielen Anwendungen werden andere Eigenvektoren gesucht, so zum Beispiel die zu einer vorgegebenen Eigenwertnäherung gehörende beste Eigenvektornäherung oder ein Eigenvektor zum betragskleinsten Eigenwert. Diese Aufgaben lassen sich mit der Vektoriteration in der ursprünglichen Form nicht lösen. Wenden wir die Vektoriteration auf die Inverse A^{-1} der Matrix A an (vorausgesetzt A ist regulär), so approximieren wir auf diese Weise einen Eigenvektor zum betragskleinsten Eigenwert von A , denn aus der Eigenwertzerlegung von A

$$A = U\Lambda U^T$$

ergibt sich die Eigenwertzerlegung der Inversen von A

$$A^{-1} = U\Lambda^{-1}U^T.$$

Sind $\lambda_1, \dots, \lambda_n$ die Eigenwerte der regulären Matrix A , so sind $\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_n}$ die Eigenwerte von A^{-1} . Durch Spektralverschiebung $A_\mu = A - \mu I = U(\Lambda - \mu I)U^T$ lässt sich jeder Eigenwert von A zum betragskleinsten Eigenwert von A_μ machen. Dies sind die grundlegenden Ideen der Inversen Iteration.

9.40. Basisverfahren der Inversen Iteration:

S0 (Initialisierung) Wähle ein $\mu \neq \lambda_j, j = 1, \dots, n$ und einen Startvektor $\mathbf{v}^{(0)} \in \mathbb{R}^n$ mit $\|\mathbf{v}^{(0)}\| = 1$. Setze $k = 0$.

S1 (Iteration) Berechne $\mathbf{w}^{(k+1)}$ aus

$$A_\mu \mathbf{w}^{(k+1)} = (A - \mu I) \mathbf{w}^{(k+1)} = \mathbf{v}^{(k)}.$$

S2 (Normierung) Setze $\mathbf{v}^{(k+1)} = \mathbf{w}^{(k+1)} / \|\mathbf{w}^{(k+1)}\|$.

S3 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Bemerkungen: (i) Die Matrix $(A - \mu I)^{-1}$ muss nicht berechnet werden. Statt dessen werden im Schritt **S1** jeweils lineare Gleichungssysteme mit der Matrix A_μ gelöst. Dazu bietet sich eine LDL^T -Zerlegung oder eine LU -Zerlegung an, die nur einmal berechnet wird. Das Lösen der Gleichungssysteme im Schritt **S1** benötigt dann jeweils $\sim n^2$ Add./Mult. Das ist derselbe Aufwand wie bei der einfachen Vektoriteration.

(ii) Wenn wir annehmen, dass die Eigenwerte von A so geordnet sind, dass

$$|\lambda_1 - \mu| < |\lambda_2 - \mu| \leq |\lambda_3 - \mu| \leq \dots \leq |\lambda_n - \mu|$$

gilt, so lässt sich der Konvergenzsatz 9.32 unmittelbar anwenden, falls

$$\kappa = \left| \frac{\lambda_1 - \mu}{\lambda_2 - \mu} \right|$$

gesetzt wird. Die Konvergenzgeschwindigkeit ist um so größer, je besser die Eigenwertnäherung μ ist. Mit besserer Eigenwertnäherung verschlechtert sich die Kondition der Matrix \mathbf{A}_μ , so dass beim Lösen der linearen Gleichungssysteme in Schritt **S1** Probleme auftreten.

(iii) Für Schritt **S1** gilt im allgemeinen

$$(\mathbf{A}_\mu + \delta\mathbf{A}_k)\mathbf{w}^{(k+1)} = \mathbf{v}^{(k)}, \quad \|\delta\mathbf{A}_k\| \leq \text{eps}F\|\mathbf{A}_\mu\|,$$

wobei F vom Verfahren zum Lösen des linearen Gleichungssystems abhängt. Ist \mathbf{A}_μ im Sinne von $\text{eps}F\text{cond}(\mathbf{A}_\mu) = \bar{\kappa} < 1$ genügend gut konditioniert (μ ist keine zu gute Näherung für λ_1), so ist $\mathbf{A}_\mu + \delta\mathbf{A}_k$ regulär, und es gilt

$$\mathbf{w}^{(k+1)} = (\mathbf{A}_\mu + \delta\mathbf{A}_k)^{-1}\mathbf{v}^{(k)} = (\hat{\mathbf{A}}_\mu + \delta\hat{\mathbf{A}}_k)\mathbf{v}^{(k)}$$

mit

$$\hat{\mathbf{A}} = \mathbf{A}_\mu^{-1}, \quad \delta\hat{\mathbf{A}}_k = (\mathbf{A}_\mu + \delta\mathbf{A}_k)^{-1} - \mathbf{A}_\mu^{-1}.$$

Für die Störung $\delta\hat{\mathbf{A}}_k$ ergibt sich dann die Abschätzung

$$\|\delta\hat{\mathbf{A}}_k\| \leq \frac{\|\hat{\mathbf{A}}\|^2}{1 - \bar{\kappa}} \|\delta\mathbf{A}_k\| \leq \hat{F}\|\hat{\mathbf{A}}\|, \quad \hat{F} = \frac{\text{cond}(\mathbf{A}_\mu)}{1 - \bar{\kappa}}F.$$

Die Voraussetzungen von Satz 9.33 sind mit $(\hat{\mathbf{A}}, \hat{F})$ statt (\mathbf{A}, F) erfüllt, falls

$$\kappa + 10\text{eps} \left(\hat{F} + \frac{n}{2} \right) = \left| \frac{\lambda_1 - \mu}{\lambda_2 - \mu} \right| + 10\text{eps} \left(\hat{F} + \frac{n}{2} \right) < 1$$

gilt.

(iv) Falls die Eigenvektoren zu den p betragskleinsten Eigenwerten gesucht sind, lässt sich das Verfahren der Inversen Iteration natürlich auch im Sinne der Teilraumiteration verallgemeinern. Man wendet die Teilraumiteration mit \mathbf{A}^{-1} statt \mathbf{A} an. Das Verfahren konvergiert gegen den p -ten dominanten Teilraum von \mathbf{A}^{-1} , falls $0 < |\lambda_1| \leq \dots \leq |\lambda_p| < |\lambda_{p+1}|$ gilt. Ist μ eine gute Näherung für λ_1 , so ist \mathbf{A}_μ schlecht konditioniert. Die Bedingung $\bar{\kappa} < 1$ ist dann nicht mehr erfüllt. Wir dürfen daher den Satz 9.33 nicht anwenden. Trotzdem ist natürlich die Inverse Iteration anwendbar, um zu μ eine akzeptable Eigenvektorapproximation zu berechnen. Man sollte folgenden speziellen Algorithmus anwenden.

9.41. Inverse Iteration bei guter Eigenwertnäherung:

Es sei μ ein im Sinne von

$$|\mu - \lambda| \leq \text{eps} F_0 \|A\| = \varepsilon_0$$

gute Näherung für den Eigenwert λ von A .

S0 (Initialisierung)

- Berechne Zerlegung $A_\mu = P_\mu^T L_\mu U_\mu$.
- Berechne $\mathbf{w}^{(0)}$ aus $U_\mu \mathbf{w}^{(0)} = \mathbf{e} = (1, \dots, 1)^T$.
- Wähle ein $\varepsilon > 0$ und setze $k = 0$.

S1 (Iteration) Berechne $\mathbf{w}^{(k+1)}$ aus

$$A_\mu \mathbf{w}^{(k+1)} = P_\mu^T L_\mu U_\mu \mathbf{w}^{(k+1)} = \mathbf{v}^{(k)}.$$

S2 (Normierung) Setze

$$\omega_{k+1} = \|\mathbf{w}^{(k+1)}\|, \quad \mathbf{v}^{(k+1)} = \mathbf{w}^{(k+1)} / \omega_{k+1}.$$

S3 (Abbruchtest) Falls $\varepsilon \cdot \omega_{k+1} < 1$ so setze $k = k + 1$ und gehe zu Schritt **S1**.

Aufwand:

- $\sim n^3/3$ Add./Mult. (**S0**)
- $\sim n^2$ Add./Mult. pro Schritt für vollbesetztes A (**S1** und **S2**)

Ohne Beweis geben wir folgenden Satz über die Konvergenz des Algorithmus an.

9.42. Satz: Das obige Verfahren werde in Computerarithmetik durchgeführt. Dann gilt

- (i) Falls das Verfahren nach k Schritten mit $\varepsilon \cdot \omega_{k+1} \geq 1$ abbricht, so gibt es eine symmetrische Störung δA , so dass $\{\mu, \mathbf{v}\}$, mit $\mathbf{v} = \mathbf{v}^{(k+1)}$, der Beziehung

$$(A + \delta A)\mathbf{v} = \mu\mathbf{v}, \quad \|\delta A\| \sim \varepsilon \|A\|$$

genügt; $\{\mu, \mathbf{v}\}$ ist daher ein akzeptables Eigenpaar von A .

(ii) Das Abbruchkriterium

$$\varepsilon \cdot \omega_{k+1} \geq 1$$

ist meist nach wenigen Schritten erfüllt ($k \leq 3$).

Bemerkungen: (i) Für eine gute Eigenwertnäherung ist $|\mu - \lambda_1| \leq \text{eps} F_0 \|\mathbf{A}\|$. Dann gilt für die Kondition der Matrix \mathbf{A}_μ

$$\text{cond}(\mathbf{A}_\mu) = \left| \frac{\mu - \lambda_n}{\mu - \lambda_1} \right| \geq \frac{|\mu - \lambda_n|}{\text{eps} F_0 \|\mathbf{A}\|}.$$

Für kleines F_0 liegt $\text{cond}(\mathbf{A}_\mu)$ in der Größenordnung $1/\text{eps}$. Die Durchführbarkeit der LU -Zerlegung mit regulärer oberer Dreiecksmatrix \mathbf{U}_μ ist nicht gesichert. Praktisch wird \mathbf{U}_μ fast immer regulär sein mit eventuell kleinen Diagonalelementen. Sollte im Algorithmus sogar $(\mathbf{U}_\mu)_{jj} = 0$ gelten, so setze man $(\mathbf{U}_\mu)_{jj} = \text{eps} \|\mathbf{A}\|$.

(ii) Man sollte in Schritt **S0** $\varepsilon = \text{eps} \sqrt{n} \|\mathbf{A}\|$ wählen.

(iii) Normalerweise ist im Falle $\text{cond}(\mathbf{A}_\mu) \approx 1/\text{eps}$ damit zu rechnen, dass die Lösung von $\mathbf{A}_\mu \mathbf{w} = \mathbf{v}$ einen großen Fehler aufweist. Für die Inverse Iteration ist nur die Richtung $\mathbf{w}/\|\mathbf{w}\|$ wesentlich, die auch bei schlechter Näherung für \mathbf{w} noch gut approximiert wird, falls ein numerisch gutartiges Verfahren angewendet wird, und falls der Eigenwert λ_1 genügend stark von den anderen Eigenwerten getrennt ist. Nur für den unwahrscheinlichen Fall $\delta \mathbf{w} \approx -\mathbf{w}^* = -\mathbf{A}_\mu^{-1} \mathbf{v}$, also für kleines $\|\mathbf{w}\|$, weicht die berechnete Richtung stark von der Richtung des Eigenvektors \mathbf{u}_1 ab.

Ist μ keine gute Eigenwertnäherung, so lässt sich μ während der Iteration verbessern. Dazu bietet sich wieder der RAYLEIGH-Quotient an. Wir erhalten folgenden Algorithmus.

9.43. RAYLEIGH-Quotienten-Iteration:

S0 (Initialisierung) Wähle $\mathbf{v}^{(0)}$ mit $\|\mathbf{v}^{(0)}\|_2 = 1$ und ein $\varepsilon > 0$. Setze $k = 0$.

S1 (Berechnen der Verschiebung) Berechne

$$\varrho_k = RQ(\mathbf{v}^{(k)}) = \mathbf{v}^{(k)T} \mathbf{A} \mathbf{v}^{(k)}.$$

S2 (Inverse Iteration) Berechne $\mathbf{w}^{(k+1)}$ aus

$$(\mathbf{A} - \varrho_k \mathbf{I}) \mathbf{w}^{(k+1)} = \mathbf{v}^{(k)}.$$

S3 (Normierung) Setze $\omega_{k+1} = \|\mathbf{w}^{(k+1)}\|_2$ und $\mathbf{v}^{(k+1)} = \mathbf{w}^{(k+1)}/\omega_{k+1}$.

S4 (Abbruchtest) Falls $\varepsilon \cdot \omega_{k+1} < 1$ so setze $k = k + 1$ und gehe zu Schritt **S1**.

Aufwand pro Schritt:

- $\sim n^2$ Add./Mult. (**S1**)
- $\sim K_1 n^3$ Add./Mult. für vollbesetztes \mathbf{A} (**S2**)
- $\sim K n$ Add./Mult. für tridiagonales \mathbf{A} (**S2**)
- $\sim n$ Add. + $\sim 2n$ Mult. + 1 Quadratwurzel (**S3**)

Es gilt der folgende Satz (ohne Beweis):

9.44. Satz: Die RAYLEIGH-Quotienten-Iteration werde mit $\varepsilon = 0$ in exakter Arithmetik durchgeführt. Dann gilt:

- (i) Die Normen der Residuen $\mathbf{r}^{(k)} = (\mathbf{A} - \varrho_k \mathbf{I})\mathbf{v}^{(k)}$ fallen monoton, und die Folge $\{\varrho_k, \vartheta_k \mathbf{v}^{(k)}\}$ konvergiert fast immer gegen ein Eigenpaar $\{\lambda_j, \mathbf{u}_j\}$.
(Die Größen $\vartheta_k \in \{1, -1\}$ werden dabei so gewählt, dass $\mathbf{u}_j^T (\vartheta_k \mathbf{v}^{(k)}) \geq 0$.)
- (ii) Falls ϱ_k gegen λ_j konvergiert, so ist die Konvergenz kubisch im Sinne von

$$\tan \varphi_{k+1} \leq \eta (\tan \varphi_k)^3.$$

Dabei ist $\eta > 0$ und $\varphi_k = \angle(\mathbf{u}_j, \vartheta_k \mathbf{v}^{(k)})$.

Bemerkungen: (i) Falls $\mathbf{A} - \varrho_k \mathbf{I}$ bei der praktischen Realisierung numerisch singular wird, so sollte man wie in Bemerkung (i) auf Seite 50 verfahren.

(ii) Das im Satz 9.44 beschriebene ideale Konvergenzverhalten wird bei der Computerrealisierung gestört. Praktisch tritt fast immer eine Verbesserung ein, so dass mit schneller Konvergenz zu rechnen ist.

(iii) Allgemein lässt sich nichts darüber ausgesagen, gegen welchen Eigenwert λ_j von \mathbf{A} die RAYLEIGH-Quotienten ϱ_k konvergieren. Es braucht insbesondere nicht derjenige Eigenwert zu sein, dessen Abstand zu ϱ_1 am kleinsten ist.

9.6. Transformationen auf Tridiagonalform

Der Aufwand der meisten Algorithmen ist beträchtlich kleiner, falls sie auf Matrizen einfacherer Gestalt angewendet werden. Darum ist es oft günstig, die zu behandelnde Matrix in einem ersten Schritt durch Ähnlichkeitstransformation auf so eine Form zu bringen. Für symmetrische Matrizen arbeitet man dabei mit orthogonalen Transformationen. Als Ziel bietet sich Tridiagonalform an (Diagonalform wäre schon die Lösung des Problems). Wir suchen eine orthogonale (n, n) -Matrix \mathbf{Q} , die die symmetrische (n, n) -Matrix \mathbf{A} gemäß

$$\mathbf{A} \rightarrow \mathbf{T} = \mathbf{Q}^T \mathbf{A} \mathbf{Q} = \begin{pmatrix} a_1 & b_2 & 0 & \cdots & 0 \\ b_2 & a_2 & b_3 & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & b_{n-1} & a_{n-1} & b_n \\ 0 & \cdots & 0 & b_n & a_n \end{pmatrix} = \text{trid}(a_1, \dots, a_n, b_2, \dots, b_n)$$

auf Tridiagonalgestalt transformiert. Wird \mathbf{A} noch durch $\sim n^2/2$ Daten repräsentiert, so sind es bei \mathbf{T} nur noch $\sim 2n$ Elemente. Durch diese Datenreduktion tritt

eine wesentliche Verringerung des Aufwands zum Lösen des Eigenwertproblems ein. Die orthogonale Matrix \mathbf{Q} ist mit Hilfe von HOUSEHOLDER-Spiegelungen oder GIVENS-Drehungen konstruierbar. Wir wollen hier nur die Vorgehensweise bei der Anwendung von HOUSEHOLDER-Matrizen erläutern.

$$\begin{aligned} \text{Es sei } \mathbf{A}^{(1)} &= \mathbf{A} \\ \mathbf{A}^{(k+1)} &= \mathbf{H}^{(k)} \mathbf{A}^{(k)} \mathbf{H}^{(k)}, \quad k = 1, \dots, n-2 \\ \mathbf{T} &= \mathbf{A}^{(n-1)} \end{aligned}$$

Die Matrizen $\mathbf{A}^{(k)}$ haben dabei folgende Struktur

$$\mathbf{A}^{(k)} = \left(\begin{array}{cc|c} & & \mathbf{o} \\ & \mathbf{T}^{(k)} & \\ \hline & & \mathbf{b}^{(k)T} \\ \hline \mathbf{o} & \mathbf{b}^{(k)} & \mathbf{B}^{(k)} \end{array} \right)$$

mit $\mathbf{T}^{(k)} = \text{trid}(a_1, \dots, a_k, b_2, \dots, b_k) \in \mathbb{R}^{k \times k}$, $\mathbf{b}^{(k)} \in \mathbb{R}^{n-k}$ und $\mathbf{B}^{(k)} \in \mathbb{R}^{(n-k) \times (n-k)}$. Die HOUSEHOLDER-Transformation $\mathbf{H}^{(k)}$ wird nun im k -ten Schritt so gewählt, dass durch sie $\mathbf{b}^{(k)}$ auf ein Vielfaches des entsprechenden Einheitsvektors transformiert wird:

$$\mathbf{H}^{(k)} = \mathbf{I} \quad \text{falls } \mathbf{b}^{(k)} = \mathbf{o} \quad \text{bzw.} \quad (9.19)$$

$$\mathbf{H}^{(k)} = \begin{pmatrix} \mathbf{I}_k & \mathbf{o} \\ \mathbf{o} & \bar{\mathbf{H}}^{(k)} \end{pmatrix}, \quad \mathbf{I}_k \in \mathbb{R}^{k \times k}, \quad \bar{\mathbf{H}}^{(k)} \in \mathbb{R}^{(n-k) \times (n-k)} \quad (9.20)$$

Die Matrix $\bar{\mathbf{H}}^{(k)}$ wird gemäß

$$\bar{\mathbf{H}}^{(k)} = \mathbf{I} - \bar{\mathbf{v}}_k \bar{\mathbf{v}}_k^T / \gamma_k \quad (9.21)$$

mit

$$\bar{\mathbf{v}}_k = \mathbf{e}_1 - \mathbf{b}^{(k)} / \varrho_k, \quad (9.22)$$

$$(9.23)$$

$$\varrho_k = \begin{cases} \|\mathbf{b}^{(k)}\| & \text{für } a_{k+1,k}^{(k)} \leq 0 \\ -\|\mathbf{b}^{(k)}\| & \text{für } a_{k+1,k}^{(k)} > 0 \end{cases} \quad \text{und} \quad (9.24)$$

$$(9.25)$$

$$\gamma_k = \bar{\mathbf{v}}_k^T \bar{\mathbf{v}}_k / 2 \quad (\text{Damit gilt } \gamma_k = v_{k+1,k} \cdot) \quad (9.26)$$

berechnet. Damit erhält man

$$\begin{aligned} \mathbf{A}^{(k+1)} &= \mathbf{H}^{(k)} \mathbf{A}^{(k)} \mathbf{H}^{(k)} & (9.27) \\ &= \left(\begin{array}{c|c} \mathbf{T}^{(k)} & \mathbf{o} \\ \hline \mathbf{o} & \bar{\mathbf{b}}^{(k)T} \\ \hline \mathbf{o} & \bar{\mathbf{b}}^{(k)} & \bar{\mathbf{B}}^{(k+1)} \end{array} \right) = \left(\begin{array}{c|c} \mathbf{T}^{(k+1)} & \mathbf{o} \\ \hline \mathbf{o} & \mathbf{b}^{(k+1)T} \\ \hline \mathbf{o} & \mathbf{b}^{(k+1)} & \mathbf{B}^{(k+1)} \end{array} \right) & (9.28) \end{aligned}$$

mit

$$\bar{\mathbf{b}}^{(k)} = \bar{\mathbf{H}}^{(k)} \mathbf{b}^{(k)}, \quad (9.29)$$

$$\bar{\mathbf{B}}^{(k)} = \bar{\mathbf{H}}^{(k)} \mathbf{B}^{(k)} \bar{\mathbf{H}}^{(k)} = \begin{pmatrix} a_{k+1,k+1} & \mathbf{b}^{(k+1)T} \\ \mathbf{b}^{(k+1)} & \mathbf{B}^{(k+1)} \end{pmatrix}. \quad (9.30)$$

Im Unterschied zum HOUSEHOLDER-Verfahren bei linearen Gleichungssystemen oder linearen Ausgleichsproblemen tritt hier die Spiegelung $\mathbf{H}^{(k)}$ auf beiden Seiten auf. Mit \mathbf{B} ist auch $\bar{\mathbf{B}} = \mathbf{H}\mathbf{B}\mathbf{H}$ symmetrisch. Es ist daher nicht zweckmäßig, $\bar{\mathbf{B}}$ nach den Vorschriften $\bar{\mathbf{B}} = (\mathbf{H}\mathbf{B})\mathbf{H}$ oder $\bar{\mathbf{B}} = \mathbf{H}(\mathbf{B}\mathbf{H})$ zu berechnen. Hierbei würde durch den Einfluss von Rundungsfehlern die Symmetrie verloren gehen. Eine symmetrische Berechnungsvorschrift lässt sich folgendermaßen herleiten:

Es gilt $\mathbf{H} = \mathbf{I} - \mathbf{v}\mathbf{v}^T/\gamma$. Damit ergibt sich

$$\bar{\mathbf{B}} = \mathbf{H}\mathbf{B}\mathbf{H} \quad (9.31)$$

$$= (\mathbf{I} - \mathbf{v}\mathbf{v}^T/\gamma)\mathbf{B}(\mathbf{I} - \mathbf{v}\mathbf{v}^T/\gamma) \quad (9.32)$$

$$= \mathbf{B} - \mathbf{v}\mathbf{v}^T\mathbf{B}/\gamma - \mathbf{B}\mathbf{v}\mathbf{v}^T/\gamma + \mathbf{v}\mathbf{v}^T\mathbf{v}\mathbf{v}^T/\gamma^2 \quad (9.33)$$

$$= \mathbf{B} - \mathbf{w}\mathbf{v}^T - \mathbf{v}\mathbf{w}^T + \frac{\mathbf{v}^T\mathbf{w}}{\gamma}\mathbf{v}\mathbf{v}^T \quad (9.34)$$

mit

$$\mathbf{w} = \mathbf{B}\mathbf{v}/\gamma. \quad (9.35)$$

Das lässt sich noch weiter umformen zu

$$\bar{\mathbf{B}} = \mathbf{B} - (\mathbf{p}\mathbf{v}^T + \mathbf{v}\mathbf{p}^T) \quad (9.36)$$

mit

$$\mathbf{p} = \mathbf{w} - \frac{\mathbf{v}^T\mathbf{w}}{2\gamma}\mathbf{v}. \quad (9.37)$$

Hier ergibt sich \bar{B} aus B durch eine symmetrische Rang-2-Modifikation. Nach $n - 2$ Schritten erhält man durch diesen Algorithmus mit

$$A^{(n-1)} = \left(\begin{array}{c|c} T^{(n-1)} & \mathbf{o} \\ \hline \mathbf{o} & \bar{B}^{(n-1)} \end{array} \right) = T = \text{trid}(a_1, \dots, a_n, b_2, \dots, b_n)$$

die gesuchte Tridiagonalmatrix. Es gilt

$$T = Q^T A Q$$

mit

$$Q = H^{(1)} H^{(2)} \dots H^{(n-2)}.$$

Insgesamt ergibt sich der folgende Algorithmus:

9.45. Tridiagonalisierung einer symmetrischen Matrix mittels HOUSEHOLDER-Spiegelungen:

S0 (Initialisierung) Setze $k = 1$ und $A^{(1)} = A$.

S1 (Berechnen der HOUSEHOLDER-Spiegelung)
Berechne $H^{(k)}$ nach (9.19)-(9.26).

(Iteration) Berechne $A^{(k+1)} = H^{(k)} A^{(k)} H^{(k)}$ nach (9.28)-(9.37).

S2 (Abbruch) Falls $k < n - 2$ so setze $k = k + 1$ und gehe zu Schritt **S1**.

Aufwand: $\sim 2n^3/3$ Add./Mult. und $\sim n$ Quadratwurzeln falls A voll besetzt ist.

Bemerkungen:

(i) Der Algorithmus ist auf dem Platz von A durchführbar. Die Elemente $a_{ij}, i \geq j + 1$ werden mit den Vektoren $\mathbf{v}^{(k)}$ überspeichert. Die Subdiagonalelemente b_i sind in einem extra Feld zu speichern.

(ii) Der Algorithmus ist in dem Sinne gutartig, dass zur berechneten Matrix T eine orthogonale Matrix P und eine Störung δA existiert, so dass

$$A + \delta A = P T P^T$$

mit $\|\delta A\|_F \leq \text{eps} 4.2 n^2 \|A\|_F$ bzw. $\|\delta A\|_2 \leq \text{eps} 0.4 n^{5/2} (1 + 9.5/\sqrt{n}) \|A\|_2$ gilt.

9.7. Der Lanczos-Algorithmus

In diesem Abschnitt werden wir einen anderen Algorithmus zur Tridiagonalisierung einer Matrix kennenlernen. Durch diesen sogenannten LANCZOS-Algorithmus wird eine Folge von Tridiagonalmatrizen wachsender Dimension erzeugt, deren größte und kleinste Eigenwerte schnell gegen die entsprechenden Eigenwerte der zugrundeliegenden Matrix A konvergieren.

Es sei die symmetrische (n, n) -Matrix A gegeben. Für einen Vektor $q \in \mathbb{R}^n$ heißt die Folge

$$q, Aq, A^2q, \dots$$

KRYLOV-Sequenz. Die ersten i Vektoren einer KRYLOV-Sequenz spannen einen Teilraum des \mathbb{R}^n , den KRYLOV-Raum, auf:

$$\begin{aligned} \mathcal{K}_i(q, A) &= \text{span}(q, Aq, A^2q, \dots, A^{i-1}q), \quad i \geq 1, \\ \mathcal{K}_0(q, A) &= \{o\}. \end{aligned}$$

Für einen festen Vektor $q \in \mathbb{R}^n$ sei m die maximale Dimension der zugehörigen KRYLOV-Räume:

$$m = \max \left\{ \dim(\mathcal{K}_i(q, A)) \mid 1 \leq i \leq n \right\}.$$

Für den KRYLOV-Raum $\mathcal{K}_m(q, A)$ sind die Vektoren

$$q, Aq, A^2q, \dots, A^{m-1}q$$

linear unabhängig, und für den Vektor $A^m q$ gilt

$$A^m q \in \mathcal{K}_m(q, A).$$

Daraus folgt

$$A\mathcal{K}_m(q, A) \subseteq \mathcal{K}_m(q, A).$$

Damit ist $\mathcal{K}_m(q, A)$ ein A -invarianter Teilraum des \mathbb{R}^n .

Die Idee des LANCZOS-Algorithmus besteht nun darin, die Abbildung

$$\Phi: \mathcal{K}_m(q, A) \rightarrow \mathcal{K}_m(q, A), \quad \Phi(x) = Ax$$

bezüglich einer speziellen orthonormalen Basis $\{q_1, \dots, q_m\}$ von $\mathcal{K}_m(q, A)$ zu beschreiben. Die Vektoren q_1, \dots, q_m werden so gewählt, dass $\{q_1, \dots, q_i\}$ jeweils eine orthonormale Basis des $\mathcal{K}_i(q, A)$ bildet. Diese Vektoren lassen sich leicht berechnen:

9.46. LANCZOS-Algorithmus:

S0 (Initialisierung) Wähle Vektor $\mathbf{q} \in \mathbb{R}^n$ mit $\|\mathbf{q}\|_2 = 1$. Setze $\mathbf{q}_1 = \mathbf{q}$ und $k = 1$, $b_1 \mathbf{q}_0 = \mathbf{o}$.

S1 (Iteration) Berechne

$$a_k = \mathbf{q}_k^T \mathbf{A} \mathbf{q}_k,$$

$$\mathbf{r}_k = \mathbf{A} \mathbf{q}_k - a_k \mathbf{q}_k - b_k \mathbf{q}_{k-1},$$

$$b_{k+1} = \|\mathbf{r}_k\|_2.$$

S2 (Abbruchtest)

- Falls $b_{k+1} \neq 0$, so berechne $\mathbf{q}_{k+1} = \mathbf{r}_k / b_{k+1}$, setze $k = k + 1$ und gehe zu Schritt **S1**.
- Falls $b_{k+1} = 0$, so setze $m = k$. STOPP

Aufwand pro Schritt:

- Eine Auswertung Matrix \times Vektor
- $\sim 4n$ Add./Sub. + $\sim 5n$ Mult./Div. + 1 Quadratwurzel.

Die durch den LANCZOS-Algorithmus erzeugten Vektoren bilden die gesuchte Orthonormalbasis des $\mathcal{K}_m(\mathbf{q}, \mathbf{A})$, wie der folgende Satz aussagt.

9.47. Satz: Die nach dem LANCZOS-Algorithmus berechneten Vektoren sind paarweise orthogonal und es gilt

$$\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_i) = \mathcal{K}_i(\mathbf{q}, \mathbf{A}), \quad i = 1, \dots, m.$$

Für den Abbruchindex gilt

$$m = \max \left\{ \dim(\mathcal{K}_i(\mathbf{q}, \mathbf{A})) \mid 1 \leq i \leq n \right\}.$$

Beweis: Wir beweisen den Satz mittels vollständiger Induktion über i . Zunächst stellen wir fest, dass für $i = 1, \dots, m - 1$

$$\mathbf{A} \mathbf{q}_i = b_i \mathbf{q}_{i-1} + a_i \mathbf{q}_i + b_{i+1} \mathbf{q}_{i+1}$$

gilt. Für $i = 1$ ist $\mathbf{q}_1 = \mathbf{q}$ offensichtlich eine orthonormale Basis des $\mathcal{K}_1(\mathbf{q}, \mathbf{A})$. Nehmen wir nun an, dass für alle $i \leq j$ die Vektoren $\mathbf{q}_1, \dots, \mathbf{q}_i$ paarweise orthogonal sind, und dass

$$\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_i) = \mathcal{K}_i(\mathbf{q}, \mathbf{A})$$

gilt. Ist nun der berechnete Vektor $\mathbf{r}_j \neq \mathbf{o}$, so ist $b_{j+1} \neq 0$. Der Vektor \mathbf{q}_{j+1} ist eindeutig definiert, und es gilt $\|\mathbf{q}_{j+1}\|_2 = 1$. Um die Orthogonalitätsbeziehungen zu beweisen, sind drei Fälle zu unterscheiden.

(i) Für $1 \leq i \leq j-2$ gilt

$$\begin{aligned}
 b_{j+1} \mathbf{q}_i^T \mathbf{q}_{j+1} &= \mathbf{q}_i^T \mathbf{r}_j \\
 &= \mathbf{q}_i^T (\mathbf{A} \mathbf{q}_j - a_j \mathbf{q}_j - b_j \mathbf{q}_{j-1}) \\
 &= \mathbf{q}_i^T \mathbf{A} \mathbf{q}_j \\
 &= (\mathbf{A} \mathbf{q}_i)^T \mathbf{q}_j \\
 &= (b_i \mathbf{q}_{i-1} + a_i \mathbf{q}_i + b_{i+1} \mathbf{q}_{i+1})^T \mathbf{q}_j \\
 &= 0.
 \end{aligned}$$

(ii) Für $i = j-1$ gilt

$$\begin{aligned}
 b_{j+1} \mathbf{q}_{j-1}^T \mathbf{q}_{j+1} &= \mathbf{q}_{j-1}^T \mathbf{r}_j \\
 &= \mathbf{q}_{j-1}^T (\mathbf{A} \mathbf{q}_j - a_j \mathbf{q}_j - b_j \mathbf{q}_{j-1}) \\
 &= \mathbf{q}_{j-1}^T \mathbf{A} \mathbf{q}_j - b_j \mathbf{q}_{j-1}^T \mathbf{q}_{j-1} \\
 &= (\mathbf{A} \mathbf{q}_{j-1})^T \mathbf{q}_j - b_j \\
 &= (b_{j-1} \mathbf{q}_{j-2} + a_{j-1} \mathbf{q}_{j-1} + b_j \mathbf{q}_j)^T \mathbf{q}_j - b_j \\
 &= b_j - b_j \\
 &= 0.
 \end{aligned}$$

(iii) Für $i = j$ gilt

$$\begin{aligned}
 b_{j+1} \mathbf{q}_j^T \mathbf{q}_{j+1} &= \mathbf{q}_j^T \mathbf{r}_j \\
 &= \mathbf{q}_j^T (\mathbf{A} \mathbf{q}_j - a_j \mathbf{q}_j - b_j \mathbf{q}_{j-1}) \\
 &= \mathbf{q}_j^T \mathbf{A} \mathbf{q}_j - a_j \mathbf{q}_j^T \mathbf{q}_j \\
 &= a_j - a_j \\
 &= 0.
 \end{aligned}$$

Wegen

$$\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_i) = \mathcal{K}_i(\mathbf{q}, \mathbf{A}) \subseteq \mathcal{K}_j(\mathbf{q}, \mathbf{A})$$

für $i \leq j$ gilt außerdem $\mathbf{A} \mathbf{q}_j \in \mathcal{K}_{j+1}(\mathbf{q}, \mathbf{A})$. Da

$$\mathbf{q}_{j+1} \in \text{span}(\mathbf{q}_{j-1}, \mathbf{q}_j, \mathbf{A} \mathbf{q}_j)$$

folgt $\mathbf{q}_{j+1} \in \mathcal{K}_{j+1}(\mathbf{q}, \mathbf{A})$ und weiter

$$\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_{j+1}) \subseteq \mathcal{K}_{j+1}(\mathbf{q}, \mathbf{A}).$$

Da die Vektoren $\mathbf{q}_1, \dots, \mathbf{q}_{j+1}$ wegen ihrer paarweisen Orthogonalität linear unabhängig sind, gilt

$$\text{span}(\mathbf{q}_1, \dots, \mathbf{q}_{j+1}) = \mathcal{K}_{j+1}(\mathbf{q}, \mathbf{A}).$$

Damit folgt sofort $j + 1 \leq \max \left\{ \dim(\mathcal{K}_i(\mathbf{q}, \mathbf{A})) \mid 1 \leq i \leq n \right\}$ und für den Abbruchindex m ebenso

$$m \leq \max \left\{ \dim(\mathcal{K}_i(\mathbf{q}, \mathbf{A})) \mid 1 \leq i \leq n \right\}.$$

Andererseits gilt wegen der Abbruchbedingung

$$\mathbf{r}_m = \mathbf{A}\mathbf{q}_m - a_m\mathbf{q}_m - b_m\mathbf{q}_{m-1} = \mathbf{o},$$

daher

$$\mathbf{A}\mathbf{q}_m \in \text{span}(\mathbf{q}_{m-1}, \mathbf{q}_m) \subset \text{span}(\mathbf{q}_1, \dots, \mathbf{q}_m) = \mathcal{K}_m(\mathbf{q}, \mathbf{A}).$$

Für $i < m$ gilt

$$\mathbf{A}\mathbf{q}_i \in \mathcal{K}_{i+1}(\mathbf{q}, \mathbf{A}) \subseteq \mathcal{K}_m(\mathbf{q}, \mathbf{A}).$$

Damit ist aber

$$\mathbf{A}\mathcal{K}_m(\mathbf{q}, \mathbf{A}) \subseteq \mathcal{K}_m(\mathbf{q}, \mathbf{A}),$$

$\mathcal{K}_m(\mathbf{q}, \mathbf{A})$ ist also \mathbf{A} -invarianter Teilraum; der kleinste unter ihnen ist jedoch durch die Bedingung

$$j = \max \left\{ \dim(\mathcal{K}_i(\mathbf{q}, \mathbf{A})) \mid 1 \leq i \leq n \right\}$$

festgelegt. Damit gilt für den Abbruchindex

$$m \geq \max \left\{ \dim(\mathcal{K}_i(\mathbf{q}, \mathbf{A})) \mid 1 \leq i \leq n \right\}.$$

Zusammengefasst ergibt sich

$$m = \max \left\{ \dim(\mathcal{K}_i(\mathbf{q}, \mathbf{A})) \mid 1 \leq i \leq n \right\}.$$

Bemerkungen: (i) Die durch den LANCZOS-Algorithmus berechneten Vektoren genügen dem folgenden Gleichungssystem:

$$\begin{aligned} \mathbf{A}\mathbf{q}_1 &= a_1\mathbf{q}_1 + b_2\mathbf{q}_2 \\ \mathbf{A}\mathbf{q}_2 &= b_2\mathbf{q}_1 + a_2\mathbf{q}_2 + b_3\mathbf{q}_3 \\ &\vdots \\ \mathbf{A}\mathbf{q}_{m-1} &= b_{m-1}\mathbf{q}_{m-2} + a_{m-1}\mathbf{q}_{m-1} + b_m\mathbf{q}_m \\ \mathbf{A}\mathbf{q}_m &= b_m\mathbf{q}_{m-1} + a_m\mathbf{q}_m. \end{aligned}$$

Mit der spaltenorthonormalen (n, m) -Matrix

$$\mathbf{Q} = (\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_{m-1}, \mathbf{q}_m)$$

und der (m, m) -Tridiagonalmatrix

$$\mathbf{J} = \text{trid}(a_1, \dots, a_m, b_2, \dots, b_m) = \begin{pmatrix} a_1 & b_2 & 0 & \dots & 0 \\ b_2 & a_2 & b_3 & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & b_{m-1} & a_{m-1} & b_m \\ 0 & \dots & 0 & b_m & a_m \end{pmatrix}$$

gilt dann

$$\mathbf{A}\mathbf{Q} = \mathbf{Q}\mathbf{J}.$$

Im Falle $m = n$ hat man mit diesem Verfahren die Matrix \mathbf{A} auf symmetrische Tridiagonalform gebracht. Die Eigenwerte von \mathbf{J} sind auch Eigenwerte von \mathbf{A} . Ist $m < n$, so lässt sich durch Wahl eines neuen Startvektors $\bar{\mathbf{q}}$ mit $\mathbf{Q}^T \bar{\mathbf{q}} = \mathbf{o}$ eine weitere spaltenorthonormale Matrix $\bar{\mathbf{Q}}$ und eine Tridiagonalmatrix $\bar{\mathbf{J}}$ berechnen, die dann weitere Eigenwerte von \mathbf{A} liefern.

(ii) Der Vorteil des LANCZOS-Algorithmus besteht darin, dass die größten bzw. kleinsten Eigenwerte der Matrizen

$$\mathbf{J}_i = \text{trid}(a_1, \dots, a_i, b_2, \dots, b_i) = \begin{pmatrix} a_1 & b_2 & 0 & \dots & 0 \\ b_2 & a_2 & b_3 & \ddots & \vdots \\ 0 & \ddots & \ddots & \ddots & 0 \\ \vdots & \ddots & b_{i-1} & a_{i-1} & b_i \\ 0 & \dots & 0 & b_i & a_i \end{pmatrix}$$

schnell gegen die größten bzw. kleinsten Eigenwerte von \mathbf{A} konvergieren. Diese lassen sich effektiv mit der Vektoriteration bzw. der Inversen Vektoriteration bestimmen. Benötigt man alle Eigenwerte, so wendet man auf die Tridiagonalmatrix den QR -Algorithmus an, der im nächsten Abschnitt beschrieben wird.

(iii) Dieser Algorithmus ist für große schwach besetzte Matrizen geeignet, da der wesentliche Aufwand in der Operation $\text{Matrix} \times \text{Vektor}$ liegt.

9.8. Der QR-Algorithmus

Der QR -Algorithmus, der 1961 unabhängig voneinander von FRANCIS und KUBLANOVSKAJA entwickelt wurde, ist das zur Zeit effektivste Verfahren zum Lösen des vollständigen Eigenwertproblems einer symmetrischen Tridiagonalmatrix. In seiner Grundform ist er auch auf beliebige symmetrische Matrizen anwendbar. Wir setzen darum vorerst nur die Symmetrie der Matrix A voraus.

9.48. QR -Algorithmus ohne Verschiebungen:

S0 (Initialisierung) Setze

$$A^{(0)} = A, \quad V^{(0)} = I, \quad k = 0.$$

S1 (QR -Zerlegung) Berechne eine orthogonale (n, n) -Matrix $Q^{(k)}$ und eine obere Dreiecksmatrix $R^{(k)}$ mit

$$A^{(k)} = Q^{(k)} R^{(k)}.$$

S2 (Iteration) Berechne

$$A^{(k+1)} = R^{(k)} Q^{(k)}, \quad V^{(k+1)} = V^{(k)} Q^{(k)}.$$

S3 Setze $k = k + 1$ und gehe zu Schritt **S1**.

9.49. Satz: Die vom obigen Algorithmus erzeugten Matrizen $A^{(k)}$, $Q^{(k)}$, $R^{(k)}$ und $V^{(k)}$ sowie die Matrix

$$U^{(k+1)} = R^{(k)} R^{(k-1)} \dots R^{(0)}$$

haben folgende Eigenschaften:

(i) $A^{(k+1)}$ ist orthogonal ähnlich zu $A^{(k)}$:

$$A^{(k+1)} = Q^{(k)T} A^{(k)} Q^{(k)}, \quad k = 0, 1, \dots$$

(ii) Es gilt $A^{(k)} = V^{(k)T} A V^{(k)}$ für $k = 0, 1, \dots$

(iii) Es gilt $A^k = V^{(k)} U^{(k)}$ für $k = 1, 2, \dots$

Beweis: (i) Aus $A^{(k)} = Q^{(k)} R^{(k)}$ folgt $R^{(k)} = Q^{(k)T} A^{(k)}$ und damit

$$A^{(k+1)} = R^{(k)} Q^{(k)} = Q^{(k)T} A^{(k)} Q^{(k)}.$$

(ii) Die Aussage ergibt sich durch wiederholte Anwendung von (i).

(iii) Wir zeigen die Aussage mittels vollständiger Induktion. Zunächst gilt wegen

$$\mathbf{A}^{(k)} = \mathbf{V}^{(k)T} \mathbf{A} \mathbf{V}^{(k)}$$

$$\mathbf{V}^{(k)} \mathbf{A}^{(k)} = \mathbf{A} \mathbf{V}^{(k)}, \quad k = 0, 1, \dots$$

Offensichtlich gilt $\mathbf{A}^1 = \mathbf{A} = \mathbf{A}^{(0)} = \mathbf{Q}^{(0)} \mathbf{R}^{(0)} = \mathbf{V}^{(1)} \mathbf{U}^{(1)}$. Damit ist der Induktionsanfang gesichert. Wir nehmen nun an, dass für alle $j \leq k$

$$\mathbf{A}^j = \mathbf{V}^{(j)} \mathbf{U}^{(j)}$$

gilt. Dann folgt

$$\begin{aligned} \mathbf{A}^{k+1} &= \mathbf{A} \mathbf{A}^k = \mathbf{A} \mathbf{V}^{(k)} \mathbf{U}^{(k)} = \mathbf{V}^{(k)} \mathbf{A}^{(k)} \mathbf{U}^{(k)} \\ &= \mathbf{V}^{(k)} \mathbf{Q}^{(k)} \mathbf{R}^{(k)} \mathbf{U}^{(k)} = \mathbf{V}^{(k+1)} \mathbf{U}^{(k+1)}. \end{aligned}$$

✱

Der QR -Algorithmus ist eng mit der Teilraumiteration verwandt. Um dies zu zeigen, benötigen wir noch folgenden Satz.

9.50. Satz: *Es sei \mathbf{A} eine reguläre (n, n) -Matrix. Durch $\mathbf{A} = \mathbf{R}$ und $\mathbf{A} = \bar{\mathbf{Q}} \bar{\mathbf{R}}$ seien zwei Zerlegungen mit orthogonalen Matrizen \mathbf{Q} und $\bar{\mathbf{Q}}$ und oberen Dreiecksmatrizen \mathbf{R} und $\bar{\mathbf{R}}$ gegeben. Dann gilt*

$$\mathbf{Q} = \bar{\mathbf{Q}} \mathbf{D}, \quad \mathbf{R} = \mathbf{D} \bar{\mathbf{R}}$$

mit einer Diagonalmatrix \mathbf{D} , für die

$$|\mathbf{D}| = \mathbf{I}$$

gilt.

Beweis: Aus $\mathbf{R} = \bar{\mathbf{Q}} \bar{\mathbf{R}}$ folgt $\bar{\mathbf{Q}}^T \mathbf{Q} = \bar{\mathbf{R}} \mathbf{R}^{-1} = \mathbf{B}$. Die Matrix \mathbf{B} ist damit gleichzeitig obere Dreiecksmatrix und Orthogonalmatrix; daher ist \mathbf{B} eine Diagonalmatrix, deren Diagonalelemente sämtlich $+1$ oder -1 sind. ✱

Damit lässt sich jetzt der Zusammenhang zwischen dem QR -Algorithmus und der Teilraumiteration zeigen. Starten wir die Teilraumiteration mit der (n, n) -Matrix $\bar{\mathbf{V}}^{(0)}$, so hat man im k -ten Schritt $\bar{\mathbf{W}}^{(k+1)} = \mathbf{A} \bar{\mathbf{V}}^{(k)}$ zu berechnen und anschließend die QR -Zerlegung $\bar{\mathbf{W}}^{(k+1)} = \bar{\mathbf{V}}^{(k+1)} \bar{\mathbf{R}}^{(k+1)}$ durchzuführen. (Um Verwechslungen mit den Matrizen aus dem QR -Algorithmus zu vermeiden, sind die hier auftretenden

Matrizen mit einem oberen Querstrich gekennzeichnet.) Führen wir nun die Matrizen $\bar{A}^{(k)}$ gemäß $\bar{A}^{(k)} = \bar{V}^{(k)T} A \bar{V}^{(k)}$ ein, so gilt für diese

$$\begin{aligned}\bar{A}^{(k)} &= \bar{V}^{(k)T} \bar{W}^{(k+1)} \\ &= \bar{V}^{(k)T} \bar{V}^{(k+1)} \bar{R}^{(k+1)}.\end{aligned}$$

Durch die orthogonale Matrix $\bar{Q}^{(k)} = \bar{V}^{(k)T} \bar{V}^{(k+1)}$ und die obere Dreiecksmatrix $\bar{R}^{(k+1)}$ ist eine QR -Zerlegung von $\bar{A}^{(k)}$ gegeben. Für die Matrix $\bar{A}^{(k+1)}$ erhält man

$$\begin{aligned}\bar{A}^{(k+1)} &= \bar{V}^{(k+1)T} A \bar{V}^{(k+1)} \\ &= \bar{V}^{(k+1)T} \bar{V}^{(k)} \bar{A}^{(k)} \bar{V}^{(k)T} \bar{V}^{(k+1)} \\ &= \bar{R}^{(k+1)} \bar{V}^{(k)T} \bar{V}^{(k+1)} \\ &= \bar{R}^{(k+1)} \bar{Q}^{(k)}.\end{aligned}$$

Die Matrix $\bar{A}^{(k+1)}$ entsteht, indem die Faktoren der QR -Zerlegung von $\bar{A}^{(k)}$ in umgekehrter Reihenfolge multipliziert werden. Das entspricht genau der Vorgehensweise des QR -Algorithmus. Die QR -Zerlegung einer Matrix ist bis auf die Vorzeichenverteilung auf die Spalten von Q und die Zeilen von R eindeutig. Es ist daher zu erwarten, dass die Folgen $\{\bar{V}^{(k)}\}$, $\{\bar{R}^{(k)}\}$, $\{\bar{A}^{(k)}\}$ und $\{\bar{Q}^{(k)}\}$, die von der Teilraumiteration erzeugt werden, eng mit den Folgen $\{V^{(k)}\}$, $\{R^{(k)}\}$, $\{A^{(k)}\}$ und $\{Q^{(k)}\}$, die der QR -Algorithmus liefert, zusammenhängen. Es gilt der folgende Satz.

9.51. Satz: Für die reguläre symmetrische (n, n) -Matrix A seien die Folgen $\{\bar{V}^{(k)}\}$ und $\{\bar{R}^{(k)}\}$ durch die Teilraumiteration mit $\bar{V}^{(0)} = I$ erzeugt. Die Matrizen $\bar{A}^{(k)}$ und $\bar{Q}^{(k)}$ seien wie oben definiert. Der QR -Algorithmus liefere die Folgen $\{V^{(k)}\}$, $\{R^{(k)}\}$, $\{A^{(k)}\}$ und $\{Q^{(k)}\}$.

Dann existieren (n, n) -Diagonalmatrizen $D^{(k)}$ mit $D^{(0)} = I$ und $|D^{(k)}| = I$ für $k = 1, 2, \dots$, so dass

$$\bar{V}^{(k)} = V^{(k)} D^{(k)}$$

und

$$\begin{aligned}\bar{A}^{(k)} &= D^{(k)} A^{(k)} D^{(k)}, \quad \bar{Q}^{(k)} = D^{(k)} Q^{(k)} D^{(k+1)} \quad \text{und} \quad \bar{R}^{(k+1)} \\ &= D^{(k+1)} R^{(k)} D^{(k)}\end{aligned}$$

gilt.

Beweis: Wir beweisen den Satz mittels vollständiger Induktion. Der Induktionsanfang für $k = 0$ ist offensichtlich. Aus der Annahme, dass die Aussagen bis zu einem Index k gelten, folgt dann

$$\begin{aligned}\bar{V}^{(k+1)} &= \bar{V}^{(k)} \bar{Q}^{(k)} \\ &= \left(V^{(k)} D^{(k)} \right) \left(D^{(k)} Q^{(k)} D^{(k+1)} \right) \\ &= V^{(k)} Q^{(k)} D^{(k+1)} \\ &= V^{(k+1)} D^{(k+1)}.\end{aligned}$$

Weiter ergibt sich

$$\begin{aligned}\bar{A}^{(k+1)} &= \bar{Q}^{(k)T} \bar{A}^{(k)} \bar{Q}^{(k)} \\ &= \left(D^{(k+1)} Q^{(k)T} D^{(k)} \right) \left(D^{(k)} A^{(k)} D^{(k)} \right) \left(D^{(k)} Q^{(k)} D^{(k+1)} \right) \\ &= D^{(k+1)} A^{(k+1)} D^{(k+1)}\end{aligned}$$

und

$$\begin{aligned}\bar{Q}^{(k+1)} \bar{R}^{(k+2)} &= \bar{A}^{(k+1)} \\ &= D^{(k+1)} A^{(k+1)} D^{(k+1)} \\ &= \left(D^{(k+1)} Q^{(k+1)} \right) \left(R^{(k+1)} D^{(k+1)} \right).\end{aligned}$$

Da die QR -Zerlegung einer regulären Matrix bis auf die Verteilung der Vorzeichen auf die Spalten von Q und die Zeilen von R eindeutig ist, existiert eine Diagonalmatrix $D^{(k+2)}$ mit $|D^{(k+2)}| = I$, so dass

$$\bar{Q}^{(k+1)} = \left(D^{(k+1)} Q^{(k+1)} \right) D^{(k+2)}, \quad \bar{R}^{(k+2)} = D^{(k+2)} \left(R^{(k+1)} D^{(k+1)} \right).$$

Damit ist der Satz bewiesen. *

Bemerkung: Wird in den Algorithmen bei der QR -Zerlegung jeweils $(R)_{jj} > 0$ für $j = 1, \dots, n$ gefordert, so gilt immer $D^{(k)} = I$. Dann stimmen die Matrizenfolgen beider Algorithmen exakt überein.

Da sich die Matrizen $V^{(k)}$ und $\bar{V}^{(k)}$ nur in den Vorzeichen der Spalten unterscheiden, gilt für eine Partitionierung

$$V^{(k)} = \left(V_1^{(k)}, V_2^{(k)} \right), \quad V_1^{(k)} \in \mathbb{R}^{n \times m}$$

und eine entsprechende Partitionierung

$$\bar{\mathbf{V}}^{(k)} = \left(\bar{\mathbf{V}}_1^{(k)}, \bar{\mathbf{V}}_2^{(k)} \right), \quad \bar{\mathbf{V}}_1^{(k)} \in \mathbb{R}^{n \times m}$$

$$\mathcal{Y}_1 = \text{span} \left(\mathbf{V}_1^{(k)} \right) = \bar{\mathcal{Y}}_1 = \text{span} \left(\bar{\mathbf{V}}_1^{(k)} \right).$$

Die Algorithmen erzeugen die gleichen Teilraumfolgen. Damit lassen sich alle Konvergenzaussagen der Teilraumiteration aus Satz 9.35 auf den QR -Algorithmus übertragen. Es gilt der folgende Satz.

9.52. Satz: Für die symmetrische (n, n) -Matrix \mathbf{A} mit den Eigenwerten

$$0 < |\lambda_n| \leq |\lambda_{n-1}| \leq \dots \leq |\lambda_{p+1}| < |\lambda_p| \leq \dots \leq |\lambda_2| \leq |\lambda_1|,$$

und zugehörigen orthonormierten Eigenvektoren

$$\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_n$$

werde der QR -Algorithmus in exakter Arithmetik durchgeführt. Mit

$$\mathcal{Y}^{(0)} = \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_p), \quad \mathcal{S}_p = \text{span}(\mathbf{u}_1, \dots, \mathbf{u}_p)$$

gelte

$$\sigma = \cos \left(\angle \left(\mathcal{Y}^{(0)}, \mathcal{S}_p \right) \right) > 0.$$

Die durch den QR -Algorithmus ohne Verschiebungen erzeugten Matrizen seien wie folgt partitioniert:

$$\mathbf{V}^{(k)} = \left(\mathbf{V}_1^{(k)}, \mathbf{V}_2^{(k)} \right), \quad \mathbf{V}_1^{(k)} \in \mathbb{R}^{n \times p}$$

und

$$\mathbf{A}^{(k)} = \begin{pmatrix} \mathbf{A}_{11}^{(k)} & \mathbf{A}_{12}^{(k)} \\ \mathbf{A}_{21}^{(k)} & \mathbf{A}_{22}^{(k)} \end{pmatrix}, \quad \mathbf{Q}^{(k)} = \begin{pmatrix} \mathbf{Q}_{11}^{(k)} & \mathbf{Q}_{12}^{(k)} \\ \mathbf{Q}_{21}^{(k)} & \mathbf{Q}_{22}^{(k)} \end{pmatrix}, \quad \mathbf{R}^{(k)} = \begin{pmatrix} \mathbf{R}_{11}^{(k)} & \mathbf{R}_{12}^{(k)} \\ \mathbf{o} & \mathbf{R}_{22}^{(k)} \end{pmatrix}$$

mit

$$\mathbf{A}_{11}^{(k)} \in \mathbb{R}^{p \times p}, \quad \mathbf{Q}_{11}^{(k)} \in \mathbb{R}^{p \times p}, \quad \mathbf{R}_{11}^{(k)} \in \mathbb{R}^{p \times p}.$$

$$\mathcal{Y}^{(k)} = \text{span} \left(\mathbf{V}_1^{(k)} \right).$$

Dann gilt

$$\tan \varphi_{k+1} \leq \kappa \tan \varphi_k \leq \kappa^{k+1} \tan \varphi_0 = \kappa^{k+1} \frac{\sqrt{1-\sigma^2}}{\sigma} = \varepsilon_{k+1} \quad (9.38)$$

mit

$$\kappa = \left| \frac{\lambda_{p+1}}{\lambda_p} \right|, \quad \varphi_k = \sphericalangle(\mathcal{Y}^{(k)}, \mathcal{S}_p) = \sphericalangle(\mathcal{Y}^{(k)\perp}, \mathcal{Z}_{p+1}).$$

Dabei ist

$$\mathcal{Y}^{(k)\perp} = \text{span}(\mathbf{V}_2^{(k)}), \quad \mathcal{Z}_{p+1} = \text{span}(\mathbf{u}_{p+1}, \dots, \mathbf{u}_n) = \mathcal{S}_p^\perp.$$

Weiterhin gelten die Abschätzungen

$$\|\mathbf{A}_{12}^{(k)}\|_2 = \|\mathbf{A}_{21}^{(k)}\|_2 \leq (1 + \kappa)\varepsilon_k \|\mathbf{A}\|_2 \leq 2\varepsilon_k \|\mathbf{A}\|_2, \quad (9.39)$$

$$\|\mathbf{Q}_{12}^{(k)}\|_2, \|\mathbf{Q}_{21}^{(k)}\|_2 \leq (1 + \kappa)\varepsilon_k \leq 2\varepsilon_k, \quad (9.40)$$

$$\|\mathbf{R}_{12}^{(k)}\|_2 \leq (1 + \kappa)^2 \varepsilon_k \|\mathbf{A}\|_2 \leq 4\varepsilon_k \|\mathbf{A}\|_2. \quad (9.41)$$

Beweis: Die Gültigkeit von (9.38) folgt sofort aus Satz 9.35 falls man berücksichtigt, dass im Falle $\dim(\mathcal{Y}) = \dim(\mathcal{S})$ die Beziehung

$$\sphericalangle(\mathcal{Y}, \mathcal{S}) = \sphericalangle(\mathcal{S}, \mathcal{Y}) = \sphericalangle(\mathcal{Y}^\perp, \mathcal{S}^\perp)$$

gilt. Um (9.40) zu beweisen gehen wir von

$$\mathbf{Q}^{(k)} = \mathbf{V}^{(k)T} \mathbf{V}^{(k+1)}$$

aus. Es gilt dann

$$\mathbf{Q}_{21}^{(k)} = \mathbf{V}_2^{(k)T} \mathbf{V}_1^{(k+1)} = \left[\mathbf{U}^T \mathbf{V}_2^{(k)} \right]^T \left[\mathbf{U}^T \mathbf{V}_1^{(k+1)} \right]$$

mit $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_p, \mathbf{u}_{p+1}, \dots, \mathbf{u}_n) = (\mathbf{U}_1, \mathbf{U}_2)$. Damit erhalten wir

$$\mathbf{U}^T \mathbf{V}_1^{(k+1)} = \begin{pmatrix} \mathbf{U}_1^T \mathbf{V}_1^{(k+1)} \\ \mathbf{U}_2^T \mathbf{V}_1^{(k+1)} \end{pmatrix} = \begin{pmatrix} \mathbf{F} \\ \mathbf{G} \end{pmatrix}.$$

Wegen der Orthogonalität der Matrizen U und $V^{(k+1)}$ gilt

$$\|U^T V_1^{(k+1)}\|_2 = \|V_1^{(k+1)}\|_2 = 1$$

und damit

$$\|F\|_2 \leq 1.$$

Es lässt sich zeigen, dass $\|G\|_2 = \sin \varphi_{k+1}$. Damit ergibt sich $\|G\|_2 \leq \varepsilon_{k+1}$. Für

$$U^T V_2^{(k+1)} = \begin{pmatrix} U_1^T V_2^{(k+1)} \\ U_2^T V_2^{(k+1)} \end{pmatrix} = \begin{pmatrix} G' \\ F' \end{pmatrix}$$

erhält man analog $\|F'\|_2 \leq 1$ und $\|G'\|_2 \leq \varepsilon_k$. Das liefert insgesamt

$$\begin{aligned} \|Q_{21}^{(k)}\|_2 &= \|G'^T F + F'^T G\|_2 \\ &\leq \|G'\|_2 \|F\|_2 + \|F'\|_2 \|G\|_2 \\ &\leq \varepsilon_k + \varepsilon_{k+1} \\ &= (1 + \kappa)\varepsilon_k. \end{aligned}$$

Dieselbe Abschätzung erhält man für Q_{12} . Weiterhin gilt $A^{(k)} = Q^{(k)} R^{(k)}$, und damit $A_{21}^{(k)} = Q_{21}^{(k)} R_{11}^{(k)}$. Mit

$$\|R_{11}^{(k)}\|_2 \leq \|R^{(k)}\|_2 = \|A^{(k)}\|_2 = \|A\|_2$$

folgt dann

$$\|A_{21}^{(k)}\|_2 \leq \|Q_{21}^{(k)}\|_2 \|R_{11}^{(k)}\|_2 \leq (1 + \kappa)\varepsilon_k \|A\|_2 \leq 2\varepsilon_k \|A\|_2.$$

Die letzte Abschätzung erhält man aus der Darstellung $R^{(k)} = A^{(k+1)} Q^{(k)T}$. Hieraus ergibt sich

$$R_{12}^{(k)} = A_{11}^{(k+1)} Q_{21}^{(k)T} + A_{12}^{(k+1)} Q_{22}^{(k)T}.$$

Wegen $\|A_{11}^{(k+1)}\|_2 \leq \|A^{(k+1)}\|_2 \leq \|A\|_2$ und $\|Q_{22}^{(k)}\|_2 \leq 1$ folgt

$$\begin{aligned} \|R_{12}^{(k)}\|_2 &\leq \|A_{11}^{(k+1)}\|_2 \|Q_{21}^{(k)}\|_2 + \|A_{12}^{(k+1)}\|_2 \|Q_{22}^{(k)}\|_2 \\ &\leq \|A\|_2 (1 + \kappa)\varepsilon_k + (1 + \kappa)\varepsilon_{k+1} \|A\|_2 \cdot 1 \\ &= (1 + \kappa)(\varepsilon_k + \varepsilon_{k+1}) \|A\|_2 \\ &= (1 + \kappa)^2 \varepsilon_k \|A\|_2 \\ &\leq 4\varepsilon_k \|A\|_2. \end{aligned}$$

Bemerkungen: (i) Nach Ungleichung (9.38) konvergieren die Teilräume $\mathcal{Y}^{(k)}$ bzw. $\mathcal{Y}^{(k)\perp}$ linear mit dem Konvergenzfaktor $\kappa < 1$ gegen \mathcal{S}_p bzw. \mathcal{Z}_{p+1} . Im Sinne von $\tan \varphi_{k+1} / \tan \varphi_k < \kappa$ ist diese Konvergenz monoton. Nach (9.39)-(9.41) konvergieren die Nichtdiagonalblöcke

$$\mathbf{A}_{12}^{(k)}, \mathbf{A}_{21}^{(k)}, \mathbf{Q}_{12}^{(k)}, \mathbf{Q}_{21}^{(k)}$$

und

$$\mathbf{R}_{12}^{(k)} \text{ für } k \rightarrow \infty \text{ wie } \varepsilon_k$$

gegen \mathbf{o} .

(ii) In den meisten Fällen werden mehrere dominante Eigenwerte existieren. Es gilt dann zum Beispiel für zwei Indizes p und p'

$$0 < |\lambda_n| \leq |\lambda_{n-1}| \leq \dots \leq |\lambda_{p'+1}| < |\lambda_{p'}| \leq \dots \leq |\lambda_{p+1}| < |\lambda_p| \leq \dots \leq |\lambda_2| \leq |\lambda_1|.$$

Die Aussagen des Satzes sind dann sowohl für p als auch für p' gültig falls

$$\cos \left(\angle \left(\mathcal{Y}^{(0)}, \mathcal{S}_p \right) \right) > 0$$

und

$$\cos \left(\angle \left(\mathcal{Y}^{(0)}, \mathcal{S}_{p'} \right) \right) > 0$$

gilt. Partitioniert man die Matrizen $\mathbf{A}^{(k)}$ gemäß

$$\mathbf{A}^{(k)} = \begin{pmatrix} \mathbf{A}_{11}^{(k)} & \mathbf{A}_{12}^{(k)} & \mathbf{A}_{13}^{(k)} \\ \mathbf{A}_{21}^{(k)} & \mathbf{A}_{22}^{(k)} & \mathbf{A}_{23}^{(k)} \\ \mathbf{A}_{31}^{(k)} & \mathbf{A}_{32}^{(k)} & \mathbf{A}_{33}^{(k)} \end{pmatrix} \text{ mit } \mathbf{A}_{11}^{(k)} \in \mathbb{R}^{p \times p}, \mathbf{A}_{22}^{(k)} \in \mathbb{R}^{(p'-p) \times (n-p')},$$

so folgt aus Satz 9.52

$$\|\mathbf{A}_{21}^{(k)}\|_2, \|\mathbf{A}_{31}^{(k)}\|_2 \leq \left\| \begin{pmatrix} \mathbf{A}_{21}^{(k)} \\ \mathbf{A}_{31}^{(k)} \end{pmatrix} \right\|_2 \leq 2\varepsilon_k \|\mathbf{A}\|_2 \quad \text{und}$$

$$\|\mathbf{A}_{31}^{(k)}\|_2, \|\mathbf{A}_{32}^{(k)}\|_2 \leq \left\| \left(\mathbf{A}_{31}^{(k)}, \mathbf{A}_{32}^{(k)} \right) \right\|_2 \leq 2\varepsilon'_k \|\mathbf{A}\|_2.$$

Insgesamt gilt dann $\|\mathbf{A}_{31}^{(k)}\|_2 \leq 2 \min\{\varepsilon_k, \varepsilon'_k\} \|\mathbf{A}\|_2$.

(iii) Sind alle Eigenwerte dem Betrag nach paarweise verschieden, also

$$0 < |\lambda_n| < |\lambda_{n-1}| < \cdots < |\lambda_2| < |\lambda_1|,$$

und gilt für alle $p \in \{1, \dots, n\}$ $\cos\left(\angle\left(\mathbf{y}^{(0)}, \mathbf{s}_p\right)\right) > 0$, so konvergieren alle Nichtdiagonalelemente gegen Null. Die Diagonalelemente konvergieren dann gegen die Eigenwerte von \mathbf{A} .

Satz 9.52 zeigt, dass nach hinreichend vielen Schritten die Nichtdiagonalblöcke $\mathbf{A}_{12}^{(k)}$ und $\mathbf{A}_{21}^{(k)}$ klein werden. Für genügend großes k gilt also

$$\|\mathbf{A}_{12}^{(k)}\|_2 = \|\mathbf{A}_{21}^{(k)}\|_2 \leq \varepsilon = \text{eps}F \|\mathbf{A}\|_2$$

mit einer Konstanten F , die den Rundungsfehlereinfluss im Algorithmus beschreibt. Vernachlässigt man unter diesen Voraussetzungen die Nichtdiagonalblöcke, so entspricht dies einem Übergang zur Matrix

$$\hat{\mathbf{A}}^{(k)} = \begin{pmatrix} \mathbf{A}_{11}^{(k)} & \mathbf{o} \\ \mathbf{o} & \mathbf{A}_{22}^{(k)} \end{pmatrix}.$$

Es gilt dann

$$\hat{\mathbf{A}}^{(k)} = \mathbf{A}^{(k)} + \delta\hat{\mathbf{A}}^{(k)}$$

mit der Störung

$$\delta\hat{\mathbf{A}}^{(k)} = - \begin{pmatrix} \mathbf{o} & \mathbf{A}_{12}^{(k)} \\ \mathbf{A}_{21}^{(k)} & \mathbf{o} \end{pmatrix}.$$

Für die Störung $\delta\hat{\mathbf{A}}^{(k)}$ erhalten wir die Abschätzung

$$\|\delta\hat{\mathbf{A}}^{(k)}\|_2 \leq \varepsilon.$$

Die Eigenwerte der Matrix $\hat{\mathbf{A}}^{(k)}$ weichen höchstens um den Betrag ε von den Eigenwerten der Matrix $\mathbf{A}^{(k)}$ und damit von den Eigenwerten von \mathbf{A} ab. Das Eigenwertproblem für $\hat{\mathbf{A}}^{(k)}$ zerfällt in zwei Eigenwertprobleme geringerer Dimension, nämlich die für $\mathbf{A}_{11}^{(k)} \in \mathbb{R}^{p \times p}$ und für $\mathbf{A}_{22}^{(k)} \in \mathbb{R}^{(n-p) \times (n-p)}$. Der Übergang von $\mathbf{A}^{(k)}$ zu $\hat{\mathbf{A}}^{(k)}$

ist damit mit einer wesentlichen Aufwandsverminderung verbunden. Dieses Vorgehen bezeichnet man als Deflation. Sind nach weiteren Schritten die Nichtdiagonalblöcke für die Matrizen $\mathbf{A}_{11}^{(k)}$ oder $\mathbf{A}_{22}^{(k)}$ wieder hinreichend klein, so zerfällt das Eigenwertproblem in noch kleinere Teilprobleme. So wird sich durch fortwährende Deflation der Gesamtaufwand erheblich verringern. Der Satz 9.52 zeigte, dass die Konvergenzgeschwindigkeit, mit der die Nichtdiagonalblöcke gegen Null streben, durch den Faktor $\kappa = |\lambda_{p+1}/\lambda_p|$ bestimmt wird. Wie bei der Inversen Iteration kann man nun wieder versuchen, diesen Konvergenzfaktor durch eine geeignete Spektralschiebung zu verkleinern. Für ein vorgegebenes $\mu \in \mathbb{R}$ definieren wir die Matrix

$$\bar{\mathbf{A}} = \mathbf{A} - \mu \mathbf{I}.$$

Weiter nehmen wir an, dass die Eigenwerte $\lambda_i - \mu$ von $\bar{\mathbf{A}}$ so geordnet sind, dass

$$0 < |\lambda_1 - \mu| \ll |\lambda_2 - \mu| \leq \dots \leq |\lambda_n - \mu|$$

gilt; μ ist eine gute Näherung für einen einfachen Eigenwert von \mathbf{A} , stimmt aber mit diesem nicht überein. Wir bezeichnen die Eigenwerte der Matrix $\bar{\mathbf{A}}$ mit

$$\lambda_n[\bar{\mathbf{A}}] = \bar{\lambda}_n = \lambda_1 - \mu, \quad \lambda_{n-1}[\bar{\mathbf{A}}] = \bar{\lambda}_{n-1} = \lambda_2 - \mu, \dots, \lambda_1[\bar{\mathbf{A}}] = \bar{\lambda}_1 = \lambda_n - \mu.$$

Führt man den QR-Algorithmus mit der Startmatrix $\bar{\mathbf{A}}^{(0)} = \bar{\mathbf{A}} = \mathbf{A} - \mu \mathbf{I}$ durch, so ergibt sich für die Approximation des zum Eigenwert $\bar{\lambda}_n = \lambda_1 - \mu$ gehörenden Teilraums der Konvergenzfaktor

$$\bar{\kappa} = \left| \frac{\bar{\lambda}_n}{\bar{\lambda}_{n-1}} \right| = \left| \frac{\lambda_1 - \mu}{\lambda_2 - \mu} \right| \ll 1.$$

Nach k Schritten erhalten wir die Matrix

$$\bar{\mathbf{A}}^{(k)} = \begin{pmatrix} \bar{\mathbf{A}}_{11}^{(k)} & \mathbf{a}^{(k)} \\ \mathbf{a}^{(k)T} & \bar{a}_{nn}^{(k)} \end{pmatrix}, \quad \bar{\mathbf{A}}_{11}^{(k)} \in \mathbb{R}^{(n-1) \times (n-1)}, \quad \mathbf{a}^{(k)} \in \mathbb{R}^{n-1}.$$

Mit Satz 9.52 folgt, dass unter der Voraussetzung $\cos \bar{\varphi}_1 > 0$ mit

$$\bar{\varphi}_1 = \sphericalangle(\text{span}(\mathbf{e}_1, \dots, \mathbf{e}_{n-1}), \mathbf{S}_{n-1})$$

$$\|\mathbf{a}^{(k)}\|_2 \leq (1 + \bar{\kappa}) \bar{\varepsilon}_k \|\bar{\mathbf{A}}\|_2, \quad \bar{\varepsilon}_{k+1} = \bar{\kappa} \bar{\varepsilon}_k = \bar{\kappa}^k \tan \bar{\varphi}_1$$

gilt. Folglich wird $\|\mathbf{a}^{(k)}\|_2$ schnell klein, falls $\tan \bar{\varphi}_1$ nicht zu groß ist. Nach wenigen Schritten ist daher $\mathbf{a}^{(k)}$ vernachlässigbar, und $\bar{a}_{nn}^{(k)}$ ist dann eine akzeptable Näherung

für $\bar{\lambda}_n = \lambda_1 - \mu$.

Die Matrizen $\bar{\mathbf{A}}^{(k)}$ sind nach Satz 9.49 orthogonal ähnlich zur Matrix $\bar{\mathbf{A}} = \mathbf{A} - \mu\mathbf{I}$. Definieren wir nun über

$$\mathbf{A}^{(k)} = \bar{\mathbf{A}}^{(k)} + \mu\mathbf{I}$$

eine weitere Folge von Matrizen, so gilt wegen

$$\bar{\mathbf{A}}^{(k)} = \bar{\mathbf{Q}}^{(k)} \bar{\mathbf{R}}^{(k)}, \quad \bar{\mathbf{A}}^{(k+1)} = \bar{\mathbf{R}}^{(k)} \bar{\mathbf{Q}}^{(k)}$$

$$\begin{aligned} \mathbf{A}^{(k+1)} &= \bar{\mathbf{A}}^{(k+1)} + \mu\mathbf{I} \\ &= \bar{\mathbf{R}}^{(k)} \bar{\mathbf{Q}}^{(k)} + \mu\mathbf{I} \\ &= \bar{\mathbf{Q}}^{(k)T} \bar{\mathbf{Q}}^{(k)} \bar{\mathbf{R}}^{(k)} \bar{\mathbf{Q}}^{(k)} + \mu\mathbf{I} \\ &= \bar{\mathbf{Q}}^{(k)T} \bar{\mathbf{A}}^{(k)} \bar{\mathbf{Q}}^{(k)} + \mu\mathbf{I} \\ &= \bar{\mathbf{Q}}^{(k)T} (\mathbf{A}^{(k)} - \mu\mathbf{I}) \bar{\mathbf{Q}}^{(k)} + \mu\mathbf{I} \\ &= \bar{\mathbf{Q}}^{(k)T} \mathbf{A}^{(k)} \bar{\mathbf{Q}}^{(k)}. \end{aligned}$$

Die so definierten Matrizen $\mathbf{A}^{(k)}$ sind daher orthogonal ähnlich zu $\mathbf{A}^{(0)} = \mathbf{A}$. Partitionieren wir die Matrix $\mathbf{A}^{(k)}$ genau wie $\bar{\mathbf{A}}^{(k)}$, so gilt

$$\mathbf{A}^{(k)} = \begin{pmatrix} \mathbf{A}_{11}^{(k)} & \mathbf{a}^{(k)} \\ \mathbf{a}^{(k)T} & a_{nn}^{(k)} \end{pmatrix} = \begin{pmatrix} \bar{\mathbf{A}}_{11}^{(k)} + \mu\mathbf{I} & \mathbf{a}^{(k)} \\ \mathbf{a}^{(k)T} & \bar{a}_{nn}^{(k)} + \mu \end{pmatrix}.$$

Die Nichtdiagonalblöcke von $\mathbf{A}^{(k)}$ sind dieselben wie die von $\bar{\mathbf{A}}^{(k)}$. Sie konvergieren gegen Null. Gilt für genügend großes k

$$\|\mathbf{a}^{(k)}\|_2 \leq \varepsilon,$$

so lässt sich zur Matrix

$$\hat{\mathbf{A}}^{(k)} = \begin{pmatrix} \mathbf{A}_{11}^{(k)} & \mathbf{o} \\ \mathbf{o}^T & a_{nn}^{(k)} \end{pmatrix}$$

übergehen. Das Diagonalelement $a_{nn}^{(k)}$ ist in diesem Falle eine gute Näherung für den Eigenwert λ_1 von \mathbf{A} . Der Algorithmus ist dann mit der symmetrischen Matrix

$A_{11}^{(k)} \in \mathbb{R}^{(n-1) \times (n-1)}$ fortsetzbar. Die Konvergenzgeschwindigkeit des Algorithmus erhöht sich, falls die Verschiebung $\mu = \mu_k$ in jedem Schritt so festgelegt wird, dass sie eine möglichst gute Eigenwertapproximation darstellt. Wir erhalten damit folgenden Algorithmus:

9.53. QR-Algorithmus mit Verschiebungen:

S0 (Initialisierung) Setze $A^{(0)} = A, V^{(0)} = I$ und $k = 0$.

S1 (Spektralverschiebung) Wähle Verschiebungsparameter μ_k .

S2 (QR-Zerlegung) Berechne eine orthogonale Matrix $Q^{(k)}$ und eine obere Dreiecksmatrix $R^{(k)}$ mit

$$A^{(k)} - \mu_k I = Q^{(k)} R^{(k)}.$$

S3 (Iteration) Berechne

$$A^{(k+1)} = R^{(k)} Q^{(k)} + \mu_k I, \quad V^{(k+1)} = V^{(k)} Q^{(k)}.$$

S3 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Bemerkungen: (i) Auch für die von diesem Algorithmus erzeugten Matrizen gilt

$$A^{(k)} = V^{(k)T} A V^{(k)}, \quad k = 0, 1, \dots$$

Es wird wieder eine Folge von orthogonal ähnlichen Matrizen erzeugt.

(ii) Falls μ_k kein Eigenwert der Matrix A ist, ist die Matrix $A^{(k)} - \mu_k I$ regulär. Bis auf die Vorzeichen der Spalten bzw. Zeilen sind dann die Faktoren $Q^{(k)}$ und $R^{(k)}$ eindeutig festgelegt. Der Algorithmus ist auch für den Fall, dass μ_k Eigenwert der Matrix A ist, durchführbar. In diesem Falle treten bei der QR-Zerlegung Nullen in der Diagonalen von $R^{(k)}$ auf, ein Eigenwert lässt sich abgespalten, und die Dimension des Problems verringert sich.

Es blieb im Algorithmus noch die Frage nach der Wahl der Verschiebungsparameter μ_k offen. Als einfachste Wahl bietet es sich an, $\mu_k = a_{nn}^{(k)}$ zu wählen.

9.54. QR-Algorithmus mit RAYLEIGH-Quotienten-Verschiebungen:

S0 (Initialisierung) Setze

$$A^{(0)} = A, \quad V^{(0)} = I, \quad k = 0.$$

S1 (Spektralverschiebung) Wähle Verschiebungsparameter $\mu_k = a_{nn}^{(k)}$.

(QR-Zerlegung) Berechne eine orthogonale Matrix $Q^{(k)}$ und eine obere Dreiecksmatrix $R^{(k)}$ mit

$$A^{(k)} - \mu_k I = Q^{(k)} R^{(k)}.$$

S2 (Iteration) Berechne

$$\mathbf{A}^{(k+1)} = \mathbf{R}^{(k)}\mathbf{Q}^{(k)} + \mu_k\mathbf{I}, \quad \mathbf{V}^{(k+1)} = \mathbf{V}^{(k)}\mathbf{Q}^{(k)}.$$

S3 Setze $k = k + 1$ und gehe zu Schritt **S1**.

Es gibt einen Zusammenhang mit der RAYLEIGH-Quotienten-Iteration, wie der folgende Satz zeigt.

9.55. Satz: Wenn der Algorithmus 9.54 exakt durchgeführt wird und in jedem Schritt μ_k kein Eigenwert der Matrix \mathbf{A} ist, stimmen die n -ten Spalten $\mathbf{v}_n^{(k)} = \mathbf{V}^{(k)}\mathbf{e}_n$ der Matrizen $\mathbf{V}^{(k)}$ bis auf das Vorzeichen mit den Iterationsvektoren $\mathbf{v}^{(k)}$, die sich nach der RAYLEIGH-Quotienten-Iteration mit dem Startvektor $\mathbf{v}^{(0)} = \mathbf{e}_n$ ergeben, überein.

Beweis: Es gilt

$$\mu_k = a_{nn}^{(k)} = \mathbf{e}_n^T \mathbf{A}^{(k)} \mathbf{e}_n = \mathbf{e}_n^T \mathbf{V}^{(k)T} \mathbf{A} \mathbf{V}^{(k)} \mathbf{e}_n = \mathbf{v}_n^{(k)T} \mathbf{A} \mathbf{v}_n^{(k)} = \varrho(\mathbf{v}_n^{(k)}).$$

Die Größe μ_k ist der zu $\mathbf{v}_n^{(k)}$ gehörende RAYLEIGH-Quotient. Es folgt weiter

$$\begin{aligned} (\mathbf{A} - \mu_k \mathbf{I}) \mathbf{v}_n^{(k+1)} &= \mathbf{V}^{(k)} \left(\mathbf{A}^{(k)} - \mu_k \mathbf{I} \right)^T \mathbf{V}^{(k)T} \mathbf{v}_n^{(k+1)} \\ &= \mathbf{V}^{(k)} \left(\mathbf{Q}^{(k)} \mathbf{R}^{(k)} \right)^T \mathbf{V}^{(k)T} \mathbf{v}_n^{(k+1)} \\ &= \mathbf{V}^{(k)} \mathbf{R}^{(k)T} \mathbf{Q}^{(k)T} \mathbf{V}^{(k)T} \mathbf{v}_n^{(k+1)} \\ &= \mathbf{V}^{(k)} \mathbf{R}^{(k)T} \mathbf{V}^{(k+1)T} \mathbf{v}_n^{(k+1)} \\ &= \mathbf{V}^{(k)} \mathbf{R}^{(k)T} \mathbf{e}_n \\ &= r_{nn}^{(k)} \mathbf{V}^{(k)} \mathbf{e}_n \\ &= r_{nn}^{(k)} \mathbf{v}_n^{(k)}. \end{aligned}$$

Da μ_k kein Eigenwert von \mathbf{A} ist, muss $r_{nn}^{(k)} \neq 0$ sein, so dass der Vektor

$$\mathbf{w}_n^{(k+1)} = \mathbf{v}_n^{(k+1)} / r_{nn}^{(k)}$$

der Gleichung

$$(\mathbf{A} - \mu_k \mathbf{I}) \mathbf{w}_n^{(k+1)} = \mathbf{v}_n^{(k)}$$

genügt. Das ist genau die Iterationsvorschrift der RAYLEIGH-Quotienten-Iteration. Damit folgt die Behauptung. *

Aus Satz 9.55 folgt sofort, dass die Konvergenzaussagen aus Satz 9.44 übernehmbar sind. Insbesondere ist die Konvergenz asymptotisch kubisch. Falls für hinreichend großes k $r_{nn}^{(k)} \leq \varepsilon$ gilt, so ist μ_k als Eigenwert von \mathbf{A} akzeptierbar. Es gilt dann

$$\mathbf{A}^{(k+1)} = \begin{pmatrix} \bar{\mathbf{A}}^{(k+1)} & \mathbf{o} \\ \mathbf{o}^T & \mu_k \end{pmatrix} + \delta\mathbf{A}^{(k+1)}, \quad \|\delta\mathbf{A}^{(k+1)}\|_2 \leq \varepsilon.$$

Der Algorithmus ist dann mit $\bar{\mathbf{A}}^{(k+1)}$ statt $\mathbf{A}^{(k+1)}$ fortsetzbar.

Eine etwas feinere Verschiebungsstrategie ist in folgendem Algorithmus enthalten.

9.56. QR-Algorithmus mit WILKINSON-Verschiebungen:

S0 (Initialisierung) Setze

$$\mathbf{A}^{(0)} = \mathbf{A}, \quad \mathbf{V}^{(0)} = \mathbf{I}, \quad k = 0.$$

S1 (Spektralverschiebung)

- Berechne die Eigenwerte μ' und μ'' der symmetrischen Matrix

$$\mathbf{P}^{(k)} = \begin{pmatrix} a_{n-1,n-1}^{(k)} & a_{n-1,n}^{(k)} \\ a_{n,n-1}^{(k)} & a_{n,n}^{(k)} \end{pmatrix} \in \mathbb{R}^{2 \times 2}.$$

- Ordne μ' und μ'' so, dass

$$|\mu' - a_{n,n}^{(k)}| < |\mu'' - a_{n,n}^{(k)}|$$

bzw.

$$|\mu'| < |\mu''|$$

im Falle

$$|\mu' - a_{n,n}^{(k)}| = |\mu'' - a_{n,n}^{(k)}|$$

gilt.

- Setze $\mu_k = \mu'$.

S2 (QR-Zerlegung) Berechne eine orthogonale (n, n) -Matrix $\mathbf{Q}^{(k)}$ und eine obere (n, n) -Dreiecksmatrix $\mathbf{R}^{(k)}$ mit

$$\mathbf{A}^{(k)} - \mu_k \mathbf{I} = \mathbf{Q}^{(k)} \mathbf{R}^{(k)}.$$

S3 (Iteration) Berechne

$$A^{(k+1)} = R^{(k)}Q^{(k)} + \mu_k I, \quad V^{(k+1)} = V^{(k)}Q^{(k)}.$$

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

Bemerkung: Die Matrix $P^{(k)}$ ist gerade die zur Matrix $Q = (\mathbf{v}_{n-1}^{(k)}, \mathbf{v}_n^{(k)})$ gehörende Matrix aus dem RAYLEIGH-RITZ-Algorithmus. Durch μ' und μ'' sind gerade die RITZschen Eigenwerte bezüglich des Unterraums $\text{span}(\mathbf{v}_{n-1}^{(k)}, \mathbf{v}_n^{(k)})$ gegeben. Damit ist der Algorithmus 9.56 eine Erweiterung des Algorithmus 9.54 auf den zweidimensionalen Fall.

9.9. Der QR-Algorithmus für Tridiagonalmatrizen

Obwohl die Konvergenz des im vorigen Abschnitt angegebenen QR -Algorithmus asymptotisch kubisch ist, ist er gegenüber dem JACOBI-Verfahren noch nicht konkurrenzfähig. Normalerweise sind immer mehrere QR -Schritte zur Abspaltung eines Eigenwertes notwendig. Selbst wenn nach jedem QR -Schritt ein Eigenwert abgespalten würde, und falls man außerdem die Dimensionsreduktion in jedem Schritt berücksichtigt, so würde das Berechnen aller Eigenwerte rund $n^4/3$ Additionen und Multiplikationen kosten. Der Aufwand beim JACOBI-Verfahren liegt dagegen bei $6n^3$ Operationen, also um eine n -Potenz niedriger. Die Situation ändert sich grundlegend, falls man die Matrix mit den Methoden aus Abschnitt 9.6. orthogonal ähnlich auf Tridiagonalform transformiert, und den QR -Algorithmus dann auf diese Tridiagonalmatrix anwendet.

Bei der QR -Zerlegung einer Tridiagonalmatrix T bietet es sich wegen der speziellen Struktur an, mit GIVENS-Drehungen zu arbeiten. Mit Hilfe von $n - 1$ GIVENS-Transformationen lassen sich so nacheinander die Subdiagonalelemente an den Positionen $(2, 1), (3, 2), (4, 3), \dots, (n, n - 1)$ zu Null machen. Man erhält

$$G_{n-1,n}G_{n-2,n-1} \cdots G_{2,3}G_{1,2}T = R.$$

Mit

$$Q = G_{1,2}^T G_{2,3}^T \cdots G_{n-2,n-1}^T G_{n-1,n}^T$$

gilt dann

$$T = QR.$$

Da die GIVENS-Drehungen jeweils nur auf zwei Zeilen wirken, hat die Matrix R die spezielle Struktur

$$R = \begin{pmatrix} \times & \times & \times & & & \\ & \times & \times & \times & & \\ & & \ddots & \ddots & \ddots & \\ & & & \times & \times & \times \\ & & & & \times & \times \\ & & & & & \times \end{pmatrix}.$$

Bildet man nun die Matrix

$$\bar{T} = RQ = RG_{1,2}^T G_{2,3}^T \cdots G_{n-2,n-1}^T G_{n-1,n}^T,$$

so wird durch jede GIVENS-Transformation $G_{j,j+1}^T$ ein Subdiagonalelement in der Spalte j erzeugt. Die Matrix \bar{T} ist damit eine obere HESSENBERG-Matrix:

$$\bar{T} = \begin{pmatrix} \times & \times & \times & & & \\ \times & \times & \times & \times & & \\ & \ddots & \ddots & \ddots & \ddots & \\ & & \times & \times & \times & \times \\ & & & \times & \times & \times \\ & & & & \times & \times \end{pmatrix}.$$

Andererseits wissen wir, dass im QR -Algorithmus die Symmetrie der Matrizen erhalten bleibt. Damit folgt aber, dass \bar{T} tridiagonal ist. Da die QR -Zerlegung einer Matrix bis auf die Vorzeichen der Spalten von Q bzw. der Zeilen von R eindeutig ist, gilt diese Aussage natürlich auch bei Anwendung eines anderen Verfahrens zur QR -Zerlegung. Starten wir den QR -Algorithmus mit einer Tridiagonalmatrix, so wird eine Folge von Tridiagonalmatrizen erzeugt. Der Aufwand für die Durchführung eines QR -Schritts verringert sich damit von $\sim 4n^3/3$ Additionen und Multiplikationen im Falle einer vollbesetzten Matrix auf $\sim K_1 n$ Additionen, $\sim K_2 n$ Multiplikationen und n Quadratwurzeln, falls mit Tridiagonalmatrizen gearbeitet wird. Die Konstanten K_1 und K_2 hängen von der konkreten Implementierung ab. Die einmalige Transformation der Ausgangsmatrix auf Tridiagonalform benötigt weitere $\sim 2n^3/3$ Operationen. Das ist etwa der halbe Aufwand eines QR -Schritts für eine vollbesetzte Matrix. Man erkennt, dass sich der Gesamtaufwand des QR -Algorithmus entscheidend verringert. Der Übergang zu Tridiagonalmatrizen hat noch einen weiteren Vorteil. Eine symme-

Eine wichtige Eigenschaft von Tridiagonalmatrizen wird in folgendem Satz beschrieben.

9.57. Satz: *Alle Eigenwerte einer nichtzerfallenden symmetrischen Tridiagonalmatrix sind einfach. Es existiert eine Nummerierung mit*

$$\lambda_1 < \lambda_2 < \cdots < \lambda_{n-1} < \lambda_n.$$

Beweis: Da T als nichtzerfallend vorausgesetzt wurde, sind für jedes λ die ersten $n - 1$ Spalten der Matrix $T - \lambda I$ linear unabhängig. Es gilt also $\text{rg}(T - \lambda I) \geq n - 1$. Ist λ ein Eigenwert von T , so ist $T - \lambda I$ singular. Dann gilt aber $\text{rg}(T - \lambda I) \leq n - 1$. Insgesamt folgt dann $\text{rg}(T - \lambda I) = n - 1$. Die geometrische Vielfachheit des Eigenwertes λ ist 1, λ ist einfacher Eigenwert. *

Für eine Matrix mit mehrfachen Eigenwerten wird bei exakter orthogonaler Transformation die zugehörige Tridiagonalmatrix zerfallend sein. Es ist dann wieder zu erwarten, dass der Gesamtaufwand sich deutlich verringert. Sind einige Eigenwerte nur dicht benachbart, so muss nach Transformation auf Tridiagonalform keine Matrix entstehen, die „fast zerfallend“ in dem Sinne ist, dass einige Nebendiagonalelemente klein werden. Ein bekanntes Beispiel dafür ist die WILKINSON-Matrix, die bei deutlich von Null verschiedenen Nichtdiagonalelementen eng benachbarte Eigenwerte besitzt.

9.58. Beispiel: Die WILKINSON-Matrix:

$$W = \text{trid}(10, 9, \dots, 1, 0, 1, \dots, 9, 10; 1, \dots, 1)$$

hat die Eigenwerte

$$\lambda_{20} = 10.74619418290339\dots,$$

$$\lambda_{21} = 10.74619418290332\dots$$



Wendet man den QR-Algorithmus mit Verschiebungen auf eine Tridiagonalmatrix an, so ist die Tridiagonalform unter der Transformation

$$A^{(k)} \rightarrow A^{(k+1)} = Q^{(k)T} A^{(k)} Q^{(k)} = R^{(k)} Q^{(k)} + \mu_k I$$

ebenfalls invariant solange μ_k kein Eigenwert von A ist. Es lässt sich nun zeigen, dass falls $A^{(k)}$ nicht zerfallend ist und μ_k Eigenwert von A ist, die letzte Zeile der

Matrix $\mathbf{A}^{(k+1)}$ bei exakter Rechnung die Gestalt $(0, \dots, 0, \mu_k)$ hat. Das Subdiagonalelement $a_{n,n-1}^{(k+1)}$ verschwindet und μ_k wird damit als Eigenwert erkannt und abgespalten. Wendet man die RAYLEIGH-Quotienten-Verschiebung an, so konvergiert $|a_{n,n-1}^{(k)}| = |b_n^{(k)}|$ in exakter Arithmetik fast immer monoton und asymptotisch kubisch gegen 0. Bei Verwendung von WILKINSON-Verschiebungen liegt sogar globale Konvergenz vor. Die Konvergenz ist dabei mindestens quadratisch und meistens besser als kubisch.

Arbeitet man bei der numerischen Realisierung mit GIVENS-Drehungen, so zeigt es sich, dass die QR -Zerlegung und das anschließende Berechnen von \mathbf{RQ} in gewisser Weise gleichzeitig durchführbar ist. Hat man nämlich mit j GIVENS-Drehungen die Subdiagonalelemente in den Spalten $1, \dots, j$ von $\mathbf{A}^{(k)} - \mu_k \mathbf{I}$ eliminiert, so bleiben diese Spalten von den weiteren Transformationen unberührt. Sie werden weder benötigt noch geändert. Diese Spalten sind schon die entsprechenden Spalten der Matrix $\mathbf{R}^{(k)}$, und auf sie sind schon die ersten $j - 1$ Transformationen von rechts anwendbar. Ein Schritt des QR -Algorithmus ist somit im wesentlichen auf dem Platz der Matrix \mathbf{A} durchführbar.

9.59. Expliziter QR -Schritt:

S0 (Initialisierung) Wähle Verschiebungsparameter μ_k und bilde

$$\bar{\mathbf{A}} = \mathbf{A}^{(k)} - \mu_k \mathbf{I}.$$

Setze $\mathbf{G}_{01} = \mathbf{I}$ und $\bar{\mathbf{V}} = \mathbf{V}^{(k)}$.

S1 (Transformation) Für $j = 1, \dots, n - 1$ führe aus:

- Lege $\mathbf{G}_{j,j+1}$ so fest, dass $(\mathbf{G}_{j,j+1} \bar{\mathbf{A}})_{j+1,j} = 0$ gilt.
- Bilde $\bar{\mathbf{A}} = \mathbf{G}_{j,j+1} \bar{\mathbf{A}}$, $\bar{\mathbf{A}} = \bar{\mathbf{A}} \mathbf{G}_{j-1,j}$ und $\bar{\mathbf{V}} = \bar{\mathbf{V}} \mathbf{G}_{j,j+1}^T$.

S2 Berechne

$$\bar{\mathbf{A}} = \bar{\mathbf{A}} \mathbf{G}_{n-1,n}^T, \quad \mathbf{A}^{(k+1)} = \bar{\mathbf{A}} + \mu_k \mathbf{I}$$

und setze

$$\mathbf{V}^{(k+1)} = \bar{\mathbf{V}}.$$

Aufwand:

- $\sim 4n$ Add./Sub., $\sim 13n$ Mult./Div. und n Quadratwurzeln zum Berechnen der Matrix $\mathbf{A}^{(k+1)}$.
- $\sim 3n^2$ Add./Sub. und $\sim 3n^2$ Mult./Div. zum Berechnen von $\mathbf{V}^{(k+1)}$.

Bemerkung: Beim expliziten QR -Schritt ist es günstig, sofort

$$\mathbf{A}^{(k+1)} = \mathbf{R}^{(k)} \mathbf{Q}^{(k)} = \mathbf{Q}^{(k)T} (\mathbf{A}^{(k)} - \mu_k \mathbf{I}) \mathbf{Q}^{(k)}$$

zu setzen. Es gilt dann

$$\mathbf{A}^{(k+1)} = \mathbf{V}^{(k)T} [\mathbf{A} - (\mu_0 + \dots + \mu_k) \mathbf{I}] \mathbf{V}^{(k)}.$$

Die Verschiebungen sind dann gemäß $\sigma_0 = 0$, $\sigma_{k+1} = \sigma_k + \mu_k$ für $k = 0, 1, \dots$ gesondert aufzusummieren.

Im obigen Expliziten QR -Schritt wirkt es noch etwas störend, dass die GIVENS-Rotationen nicht simultan von links und von rechts anwendbar sind. Mit dem folgenden Satzes lässt sich der Algorithmus so modifizieren, dass dieser Mangel beseitigt wird.

9.60. Satz: Für die symmetrische (n, n) -Matrix \mathbf{A} gelte

$$\mathbf{A} = \mathbf{Q} \mathbf{T} \mathbf{Q}^T = \bar{\mathbf{Q}} \bar{\mathbf{T}} \bar{\mathbf{Q}}^T$$

mit Tridiagonalmatrizen

$$\mathbf{T} = \text{trid}(a_1, \dots, a_n, b_2, \dots, b_n), \quad \bar{\mathbf{T}} = \text{trid}(\bar{a}_1, \dots, \bar{a}_n, \bar{b}_2, \dots, \bar{b}_n)$$

und orthogonalen Matrizen

$$\mathbf{Q} = (\mathbf{q}_1, \dots, \mathbf{q}_n), \quad \bar{\mathbf{Q}} = (\bar{\mathbf{q}}_1, \dots, \bar{\mathbf{q}}_n).$$

Ist die Matrix \mathbf{T} nichtzerfallend, also $b_i \neq 0$ für $i = 2, \dots, n$, und gilt

$$\mathbf{q}_1 = \mathbf{Q} \mathbf{e}_1 = \bar{\mathbf{Q}} \mathbf{e}_1 = \bar{\mathbf{q}}_1,$$

so ist

$$\bar{\mathbf{Q}} = \mathbf{Q} \mathbf{D}, \quad \bar{\mathbf{T}} = \mathbf{D} \mathbf{T} \mathbf{D}$$

mit

$$\mathbf{D} = \text{diag}(d_1, \dots, d_n), \quad d_1 = 1, \quad |d_i| = 1, \quad i = 2, \dots, n.$$

Beweis: Aus $\mathbf{Q} \mathbf{T} \mathbf{Q}^T = \bar{\mathbf{Q}} \bar{\mathbf{T}} \bar{\mathbf{Q}}^T$ folgt $\bar{\mathbf{T}} \mathbf{P} = \mathbf{P} \mathbf{T}$ mit der orthogonalen Matrix

$$\mathbf{P} = (\mathbf{p}_1, \dots, \mathbf{p}_n) = \bar{\mathbf{Q}}^T \mathbf{Q}.$$

Es gilt

$$\mathbf{p}_1 = \mathbf{P}\mathbf{e}_1 = \bar{\mathbf{Q}}^T \mathbf{Q}\mathbf{e}_1 = \bar{\mathbf{Q}}^T \mathbf{q}_1 = \bar{\mathbf{Q}}^T \bar{\mathbf{q}}_1 = \mathbf{e}_1.$$

Damit folgt

$$\bar{\mathbf{T}}\mathbf{P}\mathbf{e}_1 = \bar{\mathbf{T}}\mathbf{e}_1 = \bar{a}_1\mathbf{e}_1 + \bar{b}_2\mathbf{e}_2 = \mathbf{P}\mathbf{T}\mathbf{e}_1 = a_1\mathbf{p}_1 + b_2\mathbf{p}_2.$$

Linksmultiplikation mit \mathbf{e}_1^T liefert $\bar{a}_1 = a_1$ und damit $\bar{b}_2\mathbf{e}_2 = b_2\mathbf{p}_2$, also $|\bar{b}_2| = |b_2|$ und $\mathbf{p}_2 = d_2\mathbf{e}_2$ mit $d_2 = \bar{b}_2/b_2$ (da \mathbf{T} als nichtzerfallend vorausgesetzt wurde gilt $b_2 \neq 0$). Damit gilt $|d_2| = 1$. In analoger Weise folgt

$$\begin{aligned} \bar{\mathbf{T}}\mathbf{P}\mathbf{e}_2 &= \bar{\mathbf{T}}\mathbf{p}_2 = d_2\bar{\mathbf{T}}\mathbf{e}_2 = d_2(\bar{b}_2\mathbf{e}_1 + \bar{a}_2\mathbf{e}_2 + \bar{b}_3\mathbf{e}_3) \\ &= \mathbf{P}\mathbf{T}\mathbf{e}_2 = b_2\mathbf{p}_1 + a_2\mathbf{p}_2 + b_3\mathbf{p}_3 = b_2\mathbf{e}_1 + d_2a_2\mathbf{e}_2 + b_3\mathbf{p}_3. \end{aligned}$$

Wegen $d_2\bar{b}_2 = \bar{b}_2^2/b_2 = b_2^2/b_2 = b_2$ folgt

$$d_2(\bar{a}_2 - a_2\mathbf{e}_2) + \bar{b}_3\mathbf{e}_3 - b_3\mathbf{p}_3 = 0.$$

Multipliziert man diese Gleichung von links mit \mathbf{e}_2^T und beachtet dabei, dass aus der Orthogonalität von \mathbf{P} und $\mathbf{p}_2 = d_2\mathbf{e}_2$ $\mathbf{e}_2^T\mathbf{p}_3 = 0$ folgt, so ergibt sich $\bar{a}_2 = a_2$ und $\mathbf{p}_3 = d_3\mathbf{e}_3$ mit $d_3 = d_2\bar{b}_3/b_3$, daher $|d_3| = 1$. Durch Induktion folgt dann weiter $\mathbf{p}_i = d_i\mathbf{e}_i$ mit $d_i = d_{i-1}\bar{b}_i/b_i$ (also $|d_i| = 1$) für $i = 4, \dots, n$. *

Bemerkung: Die Bedingung $\mathbf{q}_1 = \bar{\mathbf{q}}_1$ ist durch $\mathbf{q}_n = \bar{\mathbf{q}}_n$ ersetzbar. Der Beweis ist dann in umgekehrter Reihenfolge, mit Spalte n beginnend, durchzuführen.

Satz 9.60 sagt aus, dass in einem QR -Schritt die Matrix $\mathbf{Q}^{(k)}$ und damit die Matrizen $\mathbf{A}^{(k+1)}$ und $\mathbf{V}^{(k+1)}$ im wesentlichen festgelegt sind, falls $\mathbf{A}^{(k)}$ nichtzerfallend ist und die erste Spalte von $\mathbf{Q}^{(k)}$ festgelegt ist. Die Matrix $\mathbf{Q}^{(k)}$ wird im Algorithmus durch ein Produkt von $n-1$ GIVENS-Transformationen dargestellt. Damit ergibt sich

$$\mathbf{Q}^{(k)}\mathbf{e}_1 = \mathbf{G}_{12}^T\mathbf{G}_{23}^T\cdots\mathbf{G}_{n-1,n}^T\mathbf{e}_1 = \mathbf{G}_{12}^T\mathbf{e}_1.$$

Die Matrix \mathbf{G}_{12} legt im wesentlichen den gesamten Transformationsschritt fest. Bestimmen wir nun weitere Matrizen $\bar{\mathbf{G}}_{23}, \dots, \bar{\mathbf{G}}_{n-1,n}$ und damit eine Matrix

$$\bar{\mathbf{Q}}^{(k)} = \mathbf{G}_{12}\bar{\mathbf{G}}_{23}\cdots\bar{\mathbf{G}}_{n-1,n}$$

so, dass

$$\bar{\mathbf{A}}^{(k+1)} = \bar{\mathbf{Q}}^{(k)T}\mathbf{A}^{(k)}\bar{\mathbf{Q}}^{(k)} = \bar{\mathbf{G}}_{n-1,n}\cdots\bar{\mathbf{G}}_{23}\mathbf{G}_{12}\mathbf{A}^{(k)}\mathbf{G}_{12}^T\bar{\mathbf{G}}_{23}^T\cdots\bar{\mathbf{G}}_{n-1,n}^T$$

S0 (Initialisierung) Wähle Verschiebungsparameter μ_k .
Setze $\bar{\mathbf{A}} = \mathbf{A}^{(k)}$ und $\bar{\mathbf{V}} = \mathbf{V}^{(k)}$.

S1 (Start-Transformation)

- Lege \mathbf{G}_{12} so fest, dass $[\mathbf{G}_{12}(\bar{\mathbf{A}} - \mu_k \mathbf{I})]_{21} = 0$ gilt.
- Bilde $\bar{\mathbf{A}} = \mathbf{G}_{12} \bar{\mathbf{A}} \mathbf{G}_{12}^T$ und $\bar{\mathbf{V}} = \bar{\mathbf{V}} \mathbf{G}_{12}^T$.

S2 (Rest-Transformation) Für $j = 2, \dots, n-1$ führe aus:

- Lege $\mathbf{G}_{j,j+1}$ so fest, dass $(\mathbf{G}_{j,j+1} \bar{\mathbf{A}})_{j+1,j-1} = 0$ gilt.
- Bilde $\bar{\mathbf{A}} = \mathbf{G}_{j,j+1} \bar{\mathbf{A}} \mathbf{G}_{j,j+1}^T$ und $\bar{\mathbf{V}} = \bar{\mathbf{V}} \mathbf{G}_{j,j+1}^T$.

S3 Setze $\mathbf{A}^{(k+1)} = \bar{\mathbf{A}}$ und $\mathbf{V}^{(k+1)} = \bar{\mathbf{V}}$.

Aufwand:

- $\sim 6n$ Add./Sub. + $\sim 11n$ Mult./Div. + n Quadratwurzeln zum Berechnen von $\mathbf{A}^{(k+1)}$.
- $\sim 3n^2$ Add./Sub. + $\sim 3n^2$ Mult./Div. zum Berechnen von $\mathbf{V}^{(k+1)}$.

Bemerkungen: (i) Der wesentliche Vorteil des Impliziten QR -Schritts liegt darin, dass $\mathbf{A}^{(k)} - \mu_k \mathbf{I}$ nicht explizit berechnet werden muss. Dadurch kann es nicht zu einem Informationsverlust im Falle $|\mu_k| \gg |a_{jj}^{(k)}|$ kommen. Die Verschiebung μ_k wird nur bei der Festlegung der ersten GIVENS-Drehung \mathbf{G}_{12} benötigt.

(ii) Nach jedem Schritt sollten alle Nichtdiagonalelemente auf Kleinheit überprüft werden. Hinreichend kleine Elemente sollten Null gesetzt werden. Das führt nicht nur zu einer deutlichen Reduktion des Aufwands, sondern auch zu einer Verbesserung des Konvergenzverhaltens, da durch die Aufspaltung in Teilprobleme geringerer Dimension für diese Probleme besser angepasste Verschiebungsparameter wählbar sind.

(iii) Ein mögliches Kriterium zum Nullsetzen von Subdiagonalelementen wäre

$$|b_i^{(k)}| \leq \text{eps}(|a_{i-1}^{(k)}| + |a_i^{(k)}|)$$

für $i \in \{2, \dots, n\}$, wobei

$$\mathbf{A}^{(k)} = \text{trid}(a_1^{(k)}, \dots, a_n^{(k)}, b_2^{(k)}, \dots, b_n^{(k)})$$

gilt. Ein aufwendigeres, aber besseres Kriterium ist

$$|b_i^{(k)}| \leq \text{eps}(\beta_i^{(k)} + \gamma_i^{(k)})(1 + \alpha_i^{(k)} / \beta_i^{(k)}), \quad i \in \{2, \dots, n\}$$

mit $\alpha_i^{(k)} = |a_{i-1}^{(k)}| + |a_i^{(k)}|$ und

$$\gamma_i^{(k)} = |a_{i-1}^{(k)} - a_i^{(k)}|, \quad i \in \{2, \dots, n\}$$

sowie

$$\beta_i^{(k)} = |b_{i-1}^{(k)}| + |b_i^{(k)}| + |b_{i+1}^{(k)}|, \quad i \in \{2, \dots, n-1\}$$

und

$$\beta_n^{(k)} = |b_{n-1}^{(k)}| + |b_n^{(k)}|.$$

Dieser Test sollte nur ausgeführt werden, falls $|b_i^{(k)}| \leq \sqrt{\text{eps}} \|\mathbf{A}\|_\infty$ gilt.

Nach hinreichend vielen Schritten sind die Nullsetzungskriterien für alle Nichtdiagonalelemente erfüllt. Dann ist $\mathbf{A}^{(k)}$ diagonal, und die Diagonalelemente approximieren die Eigenwerte von \mathbf{A} . Die Spalten von $\mathbf{V}^{(k)}$ sind die zugehörigen Eigenvektornäherungen zur Startmatrix \mathbf{A} . Wurde die Matrix \mathbf{A} durch orthogonale Tridiagonalisierung aus einer Matrix $\bar{\mathbf{A}}$ gemäß $\bar{\mathbf{A}} = \mathbf{Q}\mathbf{A}\mathbf{Q}^T$ gewonnen, so sind die Eigenvektornäherungen für $\bar{\mathbf{A}}$ durch $\mathbf{Q}\mathbf{V}^{(k)}$ gegeben. Es bietet sich daher an, in diesem Falle den Algorithmus mit $\mathbf{V}^{(0)} = \mathbf{Q}$ statt mit $\mathbf{V}^{(0)} = \mathbf{I}$ zu starten.

Insgesamt ergibt sich folgendes Grobschema zur Realisierung des QR-Algorithmus für eine Tridiagonalmatrix \mathbf{T} . Die Indizes $1 \leq p \leq q \leq n$ geben dabei die Position des gerade zu bearbeitenden untersten nichtzerfallenden Diagonalblockes von \mathbf{T} an.

9.62. QR-Algorithmus für Tridiagonalmatrizen:

Gegeben seien die (n, n) -Tridiagonalmatrix

$$\mathbf{T} = \text{trid}(a_1, \dots, a_n, b_2, \dots, b_n)$$

und die orthogonale (n, n) -Matrix \mathbf{V} .

S0 (Initialisierung) Setze $k=1$, $q=n$.

S1 (Abbruchtest) Falls $q = 1$: STOPP.

S2 (Nullsetzen kleiner Nichtdiagonalelemente und Festlegen des zu aktuellen Diagonalblocks

$$\mathbf{T}^{(p,q)} = \text{trid}(a_p, \dots, a_q, b_{p+1}, \dots, b_q)$$

der Dimension $m = q - p + 1$ von \mathbf{T})

Für $p=q, q-1, \dots, 2$ führe aus:

- Ist $|b_p|$ hinreichend klein, so
 - setze $b_p = 0$.
 - Ist $p = q$, so setze $q = q - 1$ und gehe zu Schritt **S1**.
 - Ist $p < q$, so gehe zu Schritt **S3**.
- Setze $p = 1$.

S3 (*k*-ter *QR*-Schritt) Führe *QR*-Schritt mit WILKINSON-Verschiebung für

$$\mathbf{A}^{(k)} = \mathbf{T}^{(p,q)} \in \mathbb{R}^{m \times m} \quad , \quad \mathbf{V}^{(k)} = \mathbf{V}^{(p,q)} \in \mathbb{R}^{n \times m}$$

auf dem Platz von \mathbf{T} bzw. \mathbf{V} durch. Dabei bezeichnet $\mathbf{V}^{(p,q)}$ die aus den Spalten p bis q von \mathbf{V} gebildete spaltenorthonormale Matrix.

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

Bemerkungen: (i) Praktische Erfahrungen zeigen, dass man im Mittel etwa $k \approx 1.7n$ *QR*-Schritte zur Diagonalisierung einer Tridiagonalmatrix \mathbf{T} benötigt. Das Berechnen der ersten Eigenwerte beansprucht dabei etwas mehr Schritte (etwa 4 bis 6), die letzten Eigenwerte werden nach höchstens ein bis zwei Schritten abgespalten. Damit beträgt der Gesamtaufwand zur Eigenwertberechnung einer Tridiagonalmatrix $\sim 10n^2$ Additionen und Multiplikationen und $\sim n^2$ Quadratwurzeln. Meistens ist der Aufwand sogar noch geringer. Will man auch alle Eigenvektoren berechnen, so kostet die Aufdatierung von $\mathbf{V} \sim 3n^3$ Additionen und Multiplikationen, ist daher wesentlich aufwendiger.

(ii) Für eine vollbesetzte Matrix ist natürlich noch der Aufwand für die Tridiagonalisierung zu berücksichtigen. Hierbei kostet das Berechnen der Transformation $\mathbf{T} = \mathbf{Q}^T \mathbf{A} \mathbf{Q}$ ungefähr $2n^3/3$ Additionen und Multiplikationen. Die Bildung der Matrix \mathbf{Q} , die man für die Ermittlung der Eigenvektoren benötigt, kostet weitere $\sim 2n^3/3$ Additionen und Multiplikationen. Damit beträgt der Gesamtaufwand für das Berechnen der Eigenwerte einer vollbesetzten Matrix mit Hilfe des *QR*-Algorithmus mehr als $2n^3/3$ Additionen und Multiplikationen und $\sim n^2$ Quadratwurzeln. Das Lösen des vollständigen Eigenwertproblems für eine vollbesetzte symmetrische Matrix benötigt dagegen mindestens $11n^3/3$ Additionen und Multiplikationen, ist also wesentlich teurer.

(iii) Sind nur wenige Eigenvektoren zu berechnen, so erhält man diese billiger mittels Inverser Iteration. Diese benötigt pro Eigenvektor einer Tridiagonalmatrix $\sim Kn$ Additionen und Multiplikationen. Verwendet man für \mathbf{Q} die Darstellung als Produkt von GIVENS-Drehungen, so benötigt das Berechnen eines Eigenvektors der Ausgangsmatrix \mathbf{A} aus einem Eigenvektor der Tridiagonalmatrix $\mathbf{T} \sim n^2$ Additionen und Multiplikationen. Beim Berechnen von Eigenvektoren zu dicht benachbarten Eigenwerten ist dabei eventuell eine Reorthogonalisierung notwendig.

(iv) Für eine Bandmatrix \mathbf{A} wird der *QR*-Algorithmus direkt mit \mathbf{A} ohne vorherige Tridiagonalisierung durchgeführt, da die Bandgestalt invariant unter den *QR*-Transformationen ist.

(v) Man kann zeigen, dass die bei Abbruch nach k Schritten vorliegende Matrix $\mathbf{M} = \mathbf{T}^{(k)}$ die exakten Eigenwerte einer gestörten Matrix $\mathbf{A} + \delta\mathbf{A}$ enthält. Es gilt

$$\mathbf{A} + \delta\mathbf{A} = \mathbf{V} \mathbf{M} \mathbf{V}^T$$

mit der symmetrischen (n, n) -Störung $\delta\mathbf{A}$, die der Abschätzung

$$\|\delta\mathbf{A}\|_2 \leq \text{eps}F\|\mathbf{A}\|_2$$

genügt. Dabei ist \mathbf{V} eine exakt orthogonale Matrix. Das Berechnen der Eigenwerte der Matrix \mathbf{A} mit dem QR -Algorithmus ist demnach ein numerisch gutartiger Prozess. Die berechnete Matrix $\tilde{\mathbf{V}}$ der Eigenvektoren approximiert dabei die exakte Matrix \mathbf{V} im Sinne von

$$\|\mathbf{V} - \tilde{\mathbf{V}}\|_2 \leq \text{eps}F_1$$

hinreichend genau.

(vi) Es besteht auch die Möglichkeit, den Algorithmus in inverser Reihenfolge durchzuführen. Dabei werden die Verschiebungen jeweils durch die links oben stehende 2×2 -Matrix festgelegt. Statt der QR -Zerlegung der Matrix $\mathbf{A}^{(k)} - \mu_k\mathbf{I}$ wird dann eine QL -Faktorisierung $\mathbf{A}^{(k)} - \mu_k\mathbf{I} = \mathbf{Q}^{(k)}\mathbf{L}^{(k)}$ mit einer unteren Dreiecksmatrix $\mathbf{L}^{(k)}$ durchgeführt. Diese Variante wird auch QL -Algorithmus genannt und heute in Programmpaketen häufig angewendet.

Kapitel 10

Lineare Ausgleichsprobleme

10.1. Problemstellung und klassische Lösung

Viele Anwendungen in der Physik, Technik oder Ökonomie führen auf folgende Problemstellung:

Zwischen zwei Größen y und z wird ein funktionaler Zusammenhang der Art

$$y = f(z; x_1, \dots, x_n)$$

angenommen. Die Variablen x_1, \dots, x_n sind unbekannte Parameter, die zu bestimmen sind. Dazu führt man m „Messungen“ durch. Man erhält Wertepaare (z_i, y_i) ($i = 1, \dots, m$) und m Gleichungen

$$y_i = f(z_i; x_1, \dots, x_n), \quad i = 1, \dots, m,$$

um aus diesen die unbekannt Parameter zu ermitteln. Für ein exaktes Bestimmen der Parameter sind mindestens n Messungen nötig. Da aber die Messungen meist fehlerbehaftet sind, wird man mehr als n Messungen durchführen, um eine größere statistische Sicherheit zu bekommen. Damit ergibt sich ein überbestimmtes Gleichungssystem, das im allgemeinen keine Lösung besitzt. Man wird darum das Gleichungssystem „so gut wie möglich“ lösen wollen. Das führt auf einen neuen Lösungsbegriff.

Es seien $\mathbf{y} = (y_1, \dots, y_m)^T \in \mathbb{R}^m$, $\mathbf{z} = (z_1, \dots, z_m)^T \in \mathbb{R}^m$, $\mathbf{x} = (x_1, \dots, x_n)^T \in \mathbb{R}^n$ und

$$\mathbf{F}(\mathbf{z}; \mathbf{x}) = \begin{pmatrix} f(z_1; x_1, \dots, x_n) \\ f(z_2; x_1, \dots, x_n) \\ \vdots \\ f(z_m; x_1, \dots, x_n) \end{pmatrix} = \begin{pmatrix} f_1(\mathbf{x}) \\ f_2(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{pmatrix}.$$

Das überbestimmte Gleichungssystem lautet dann

$$\mathbf{y} = \mathbf{F}(\mathbf{z}; \mathbf{x}).$$

Es liegt nahe, den Vektor \mathbf{x} so zu bestimmen, dass das Residuum

$$\mathbf{r}(\mathbf{x}) = \mathbf{y} - \mathbf{F}(\mathbf{z}; \mathbf{x})$$

möglichst klein wird. Wir lösen daher folgende Ersatzaufgabe

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{r}(\mathbf{x})\|,$$

wobei $\|\circ\|$ eine beliebige Vektornorm bezeichnen soll. Dieser neue Lösungsbegriff für das Gleichungssystem hat den Vorteil, dass er den alten Lösungsbegriff einschließt. Im Falle exakter Lösbarkeit stimmt \mathbf{x}^* mit der exakten Lösung überein.

Je nach Wahl der Norm ergeben sich verschiedene Aufgabenstellungen.

Verwendet man die Maximumnorm, so ergibt sich

$$\|\mathbf{r}(\mathbf{x})\|_\infty = \max_{k=1, \dots, n} |y_k - f(z_k; x_1, \dots, x_n)| = \max_{k=1, \dots, n} |y_k - f_k(\mathbf{x})|$$

und

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \max_{k=1, \dots, n} |y_k - f_k(\mathbf{x})|.$$

Diese Aufgabe bezeichnet man als **Diskretes TSCHEBYSCHJEFF-Problem**.

Wir werden jedoch die euklidische Vektornorm

$$\|\mathbf{r}(\mathbf{x})\|_2^2 = \sum_{k=1}^n [y_k - f(z_k; x_1, \dots, x_n)]^2 = \sum_{k=1}^n [y_k - f_k(\mathbf{x})]^2$$

verwenden. Hier gilt

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \sum_{k=1}^n [y_k - f_k(\mathbf{x})]^2.$$

Dieses Problem bezeichnet man als **Ausgleichsproblem im engeren Sinne**. Es wurde schon von GAUSS als Methode der kleinsten Quadrate untersucht.

Besitzen die Funktionen $f_k(\mathbf{x}) = f(z_k; x_1, \dots, x_n)$ stetige partielle Ableitungen nach den x_i , so lassen sich sofort notwendige Bedingungen für die Existenz eines Minimums angeben:

$$\frac{\partial}{\partial x_i} \sum_{k=1}^n [y_k - f_k(x_1, \dots, x_n)]^2 = 2 \sum_{k=1}^n \frac{\partial}{\partial x_i} f_k(x_1, \dots, x_n) [y_k - f_k(x_1, \dots, x_n)] = 0$$

für $i = 1, \dots, n$. Diese Gleichungen heißen **Normalgleichungen**. Besonders einfach wird die Situation, falls die Funktion f linear von den Parametern x_i abhängt. Wir stellen dann f in der Form

$$f(\mathbf{z}; \mathbf{x}) = x_1 \varphi_1(z) + x_2 \varphi_2(z) + \dots + x_n \varphi_n(z)$$

und die f_k für $k = 1, \dots, m$ in der Form

$$f_k(\mathbf{x}) = x_1\varphi_1(z_k) + x_2\varphi_2(z_k) + \dots + x_n\varphi_n(z_k)$$

dar. Mit den Bezeichnungen $a_{ki} = \varphi_i(z_k)$ für $i = 1, \dots, n$ und $k = 1, \dots, m$ gilt dann

$$\mathbf{F}(\mathbf{z}; \mathbf{x}) = \mathbf{A}\mathbf{x}$$

mit

$$\mathbf{A} = \begin{pmatrix} a_{11} & \dots & a_{1n} \\ \vdots & \ddots & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix}.$$

Wir erhalten das **lineare Ausgleichsproblem**

$$\mathbf{x}^* = \arg \min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$$

mit den Normalgleichungen

$$\mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{x}) = \mathbf{o}$$

beziehungsweise

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{y}.$$

Dieses Gleichungssystem ist stets lösbar (Übungsaufgabe 1).

Im Falle $\text{rg}(\mathbf{A}) = n$ (\mathbf{A} ist **spaltenregulär**) ist die Matrix $\mathbf{A}^T \mathbf{A}$ regulär und darüber hinaus sogar positiv definit (Übungsaufgabe 2). In diesem Falle ist das Normalgleichungssystem eindeutig lösbar. Wir wenden zum Lösen das CHOLESKY-Verfahren an. Spätere Überlegungen werden jedoch zeigen, dass dies nicht der günstigste Weg ist.

Zunächst wollen wir jedoch den Zusammenhang zwischen dem eigentlichen Ausgleichsproblem und den Normalgleichungen näher beleuchten.

10.1. Satz: *Das lineare Ausgleichsproblem*

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

besitzt mindestens eine Lösung \mathbf{x}^* . Ist $\bar{\mathbf{x}}$ eine weitere Lösung, so gilt $\mathbf{A}\mathbf{x}^* = \mathbf{A}\bar{\mathbf{x}}$. Das Residuum $\mathbf{r} = \mathbf{y} - \mathbf{A}\mathbf{x}^*$ genügt der Gleichung $\mathbf{A}^T \mathbf{r} = \mathbf{o}$. Jede Lösung \mathbf{x}^* des Ausgleichsproblems ist auch Lösung der Normalgleichungen und umgekehrt.

Beweis: Es sei

$$\mathcal{V} = \text{span}(\mathbf{A}) = \{ \mathbf{A}\mathbf{x} \mid \mathbf{x} \in \mathbb{R}^n \} \subseteq \mathbb{R}^m$$

der Teilraum, der von den Spalten von \mathbf{A} aufgespannt wird. Weiterhin bezeichne

$$\mathcal{V}^\perp = \left\{ \mathbf{r} \in \mathbb{R}^m \mid \mathbf{r}^T \mathbf{z} = 0 \forall \mathbf{z} \in \mathcal{V} \right\} = \left\{ \mathbf{r} \in \mathbb{R}^m \mid \mathbf{r}^T \mathbf{A} = \mathbf{o} \right\}$$

den zu \mathcal{V} gehörenden Orthogonalraum. Wegen $\mathbb{R}^m = \mathcal{V} \oplus \mathcal{V}^\perp$ lässt sich der Vektor $\mathbf{y} \in \mathbb{R}^m$ eindeutig in der Form $\mathbf{y} = \mathbf{s} + \mathbf{r}$ mit $\mathbf{s} \in \mathcal{V}$ und $\mathbf{r} \in \mathcal{V}^\perp$ zerlegen. Es gibt mindestens ein $\mathbf{x}^* \in \mathbb{R}^n$ mit $\mathbf{s} = \mathbf{A}\mathbf{x}^*$. Wegen $\mathbf{A}^T \mathbf{r} = \mathbf{o}$ gilt

$$\mathbf{A}^T \mathbf{y} = \mathbf{A}^T \mathbf{s} = \mathbf{A}^T \mathbf{A}\mathbf{x}^*.$$

\mathbf{x}^* ist somit Lösung der Normalgleichungen. Umgekehrt entspricht jeder Lösung $\bar{\mathbf{x}}$ der Normalgleichungen die Zerlegung

$$\mathbf{y} = \bar{\mathbf{s}} + \bar{\mathbf{r}}, \bar{\mathbf{s}} = \mathbf{A}\bar{\mathbf{x}} \in \mathcal{V}, \quad \bar{\mathbf{r}} = \mathbf{y} - \mathbf{A}\bar{\mathbf{x}} \in \mathcal{V}^\perp.$$

Wegen der Eindeutigkeit der Zerlegung muss aber $\mathbf{s} = \bar{\mathbf{s}}$ und damit $\mathbf{A}\mathbf{x}^* = \mathbf{A}\bar{\mathbf{x}}$ gelten. Weiterhin ist jede Lösung \mathbf{x}^* der Normalgleichungen Optimallösung des Ausgleichsproblems. Ist nämlich \mathbf{x} beliebig und $\mathbf{z} = \mathbf{A}\mathbf{x} - \mathbf{A}\mathbf{x}^* \in \mathcal{V}$ und $\mathbf{r} = \mathbf{y} - \mathbf{A}\mathbf{x}^*$, so gilt wegen $\mathbf{z}^T \mathbf{r} = 0$

$$\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 = \|\mathbf{r} - \mathbf{z}\|_2^2 = \|\mathbf{r}\|_2^2 + \|\mathbf{z}\|_2^2 \geq \|\mathbf{r}\|_2^2 = \|\mathbf{y} - \mathbf{A}\mathbf{x}^*\|_2^2.$$

Damit ist die Existenz einer Optimallösung und die Eindeutigkeit des zugehörigen Residuums bewiesen. *

Im spaltenregulären Falle ist das Ausgleichsproblem mittels einer CHOLESKYZerlegung der Matrix $\mathbf{A}^T \mathbf{A}$ lösbar. Eine andere (mehr theoretische) Möglichkeit ergibt sich aus der Singulärwertzerlegung der Matrix \mathbf{A} . Es sei $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ mit $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_l)$, $l = \min\{m, n\}$, die Singulärwertzerlegung von \mathbf{A} . Für die Singulärwerte gelte

$$\sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_l = 0$$

mit $r = \text{rg}(\mathbf{A})$. Wegen der Invarianz der euklidischen Vektornorm gegenüber orthogonalen Transformationen gilt

$$\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 = \left\| \mathbf{y} - \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \mathbf{x} \right\|_2^2 = \left\| \mathbf{U}^T \left(\mathbf{y} - \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \mathbf{x} \right) \right\|_2^2 = \left\| \mathbf{U}^T \mathbf{y} - \mathbf{\Sigma}\mathbf{V}^T \mathbf{x} \right\|_2^2.$$

Mit $\boldsymbol{\eta} = \mathbf{U}^T \mathbf{y}$ und $\boldsymbol{\xi} = \mathbf{V}^T \mathbf{x}$ erhalten wir das äquivalente Ausgleichsproblem

$$\min_{\boldsymbol{\xi} \in \mathbb{R}^n} \|\boldsymbol{\eta} - \boldsymbol{\Sigma} \boldsymbol{\xi}\|_2.$$

Hier gilt aber

$$\|\boldsymbol{\eta} - \boldsymbol{\Sigma} \boldsymbol{\xi}\|_2^2 = \sum_{i=1}^r (\eta_i - \sigma_i \xi_i)^2 + \sum_{i=r+1}^m \eta_i^2.$$

Der zweite Summand hängt nicht von $\boldsymbol{\xi}$ ab; er spielt daher bei der Minimierung keine Rolle. Für den ersten Summanden gilt

$$\sum_{i=1}^r (\eta_i - \sigma_i \xi_i)^2 \geq 0$$

und

$$\sum_{i=1}^r (\eta_i - \sigma_i \xi_i)^2 = 0 \iff \xi_i = \frac{\eta_i}{\sigma_i}, \quad i = 1, \dots, r.$$

Wir erhalten die Lösung des Ausgleichsproblems in der Form

$$\mathbf{x}^* = \mathbf{V} \boldsymbol{\xi}^*$$

mit

$$\boldsymbol{\xi}^* = \left(\frac{\eta_1}{\sigma_1}, \dots, \frac{\eta_r}{\sigma_r}, \xi_{r+1}^*, \dots, \xi_n^* \right)^T.$$

Im spaltenregulären Falle ($r = n$) ist die Lösung eindeutig bestimmt. Für $r < n$ sind die $\xi_{r+1}^*, \dots, \xi_n^*$ frei wählbar. Setzt man $\xi_{r+1}^* = \dots = \xi_n^* = 0$, so erhält man die normkleinste Lösung. Für das Residuum gilt

$$\mathbf{r}^* = \mathbf{y} - \mathbf{A} \mathbf{x}^* = \mathbf{U} \boldsymbol{\eta} - \mathbf{U} \boldsymbol{\Sigma} \mathbf{V}^T \mathbf{V} \boldsymbol{\xi}^* = \mathbf{U} (\boldsymbol{\eta} - \boldsymbol{\Sigma} \boldsymbol{\xi}^*) = \mathbf{U} \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \eta_{r+1} \\ \vdots \\ \eta_m \end{pmatrix}.$$

Die normkleinste Lösung des Ausgleichsproblems lässt sich formal als

$$\mathbf{x}^* = \mathbf{A}^+ \mathbf{y}$$

schreiben. Die Matrix \mathbf{A}^+ wird als **Pseudoinverse** bezeichnet. Sie ist für eine beliebige (m, n) -Matrix \mathbf{A} mit der Singulärwertzerlegung $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ durch

$$\mathbf{A}^+ = \mathbf{V}\mathbf{\Sigma}^+\mathbf{U}^T$$

definiert. Dabei gilt für

$$\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$$

$$\mathbf{\Sigma}^+ = \text{diag}(1/\sigma_1, \dots, 1/\sigma_r, 0, \dots, 0) \in \mathbb{R}^{n \times m}$$

mit $r = \text{rg}(\mathbf{A})$. Somit für

$$\mathbf{\Sigma} = \left(\begin{array}{ccc|cc} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_r & & \\ \hline & & & \mathbf{O} & \\ \hline & & & \mathbf{O} & \end{array} \right) \in \mathbb{R}^{m \times n}$$

ergibt sich

$$\mathbf{\Sigma}^+ = \left(\begin{array}{ccc|cc} \frac{1}{\sigma_1} & & & & \\ & \ddots & & & \\ & & \frac{1}{\sigma_r} & & \\ \hline & & & \mathbf{O} & \\ \hline & & & \mathbf{O} & \end{array} \right) \in \mathbb{R}^{n \times m}.$$

Für $m = n = r$ gilt $\mathbf{A}^+ = \mathbf{A}^{-1}$, so dass die Pseudoinverse als Verallgemeinerung der gewöhnlichen Inversen einer Matrix anzusehen ist. Eine etwas andere Charakterisierung der Pseudoinversen ist im folgenden Satz angegeben.

10.2. Satz:

1. Zu jeder (m, n) -Matrix \mathbf{A} gibt es genau eine (n, m) -Matrix \mathbf{A}^+ , so dass

$$\mathbf{x}^* = \mathbf{A}^+ \mathbf{y}$$

die normkleinste Lösung (Normallösung) des linearen Ausgleichsproblems

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

ist. Ist die Singulärwertzerlegung von \mathbf{A} durch $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ mit

$$\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0)$$

und

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$$

gegeben, so gilt

$$\mathbf{A}^+ = \mathbf{V}\Sigma^+\mathbf{U}^T$$

mit

$$\Sigma^+ = \text{diag}(1/\sigma_1, \dots, 1/\sigma_r, 0, \dots, 0).$$

2. \mathbf{A}^+ ist durch die PENROSE-Bedingungen

- (a) $\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$,
- (b) $\mathbf{A}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+$,
- (c) $(\mathbf{A}\mathbf{A}^+)^T = \mathbf{A}\mathbf{A}^+$,
- (d) $(\mathbf{A}^+\mathbf{A})^T = \mathbf{A}^+\mathbf{A}$

eindeutig festgelegt.

Beweis:

1. Wurde bereits gezeigt.
2. Die Gültigkeit der PENROSE-Bedingungen für

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T, \quad \mathbf{A}^+ = \mathbf{V}\Sigma^+\mathbf{U}^T$$

zeigt man durch einfaches Nachrechnen. Es bleibt noch zu beweisen, dass die Pseudoinverse durch die PENROSE-Bedingungen eindeutig festgelegt ist. Dazu nehmen wir an, dass die Matrix \mathbf{B} ebenfalls die PENROSE-Bedingungen erfüllt. Dann gilt

$$\begin{aligned}
 \mathbf{B} &= \mathbf{B}\mathbf{A}\mathbf{B} && \text{(nach Bedingung } b) \\
 &= \mathbf{B}\mathbf{A}\mathbf{A}^+\mathbf{A}\mathbf{B} && \text{(nach Bedingung } a) \\
 &= \mathbf{A}^T\mathbf{B}^T\mathbf{A}^+\mathbf{A}\mathbf{B} && \text{(nach Bedingung } d) \\
 &= \mathbf{A}^T\mathbf{B}^T\mathbf{A}^T\mathbf{A}^{+T}\mathbf{B} = (\mathbf{A}\mathbf{B}\mathbf{A})^T\mathbf{A}^{+T}\mathbf{B} && \text{(nach Bedingung } d) \\
 &= \mathbf{A}^T\mathbf{A}^{+T}\mathbf{B} && \text{(nach Bedingung } a) \\
 &= \mathbf{A}^+\mathbf{A}\mathbf{B} && \text{(nach Bedingung } d) \\
 &= \mathbf{A}^+\mathbf{B}^T\mathbf{A}^T && \text{(nach Bedingung } c) \\
 &= \mathbf{A}^+\mathbf{A}\mathbf{A}^+\mathbf{B}^T\mathbf{A}^T && \text{(nach Bedingung } b) \\
 &= \mathbf{A}^+\mathbf{A}^{+T}\mathbf{A}^T\mathbf{B}^T\mathbf{A}^T = \mathbf{A}^+\mathbf{A}^{+T}(\mathbf{A}\mathbf{B}\mathbf{A})^T && \text{(nach Bedingung } c) \\
 &= \mathbf{A}^+\mathbf{A}^{+T}\mathbf{A}^T && \text{(nach Bedingung } a) \\
 &= \mathbf{A}^+\mathbf{A}\mathbf{A}^+ && \text{(nach Bedingung } c) \\
 &= \mathbf{A}^+ && \text{(nach Bedingung } b).
 \end{aligned}$$



Für einige spezielle Matrizen lässt sich die Pseudoinverse sofort angeben. Es gilt:

- $\mathbf{O}^+ = \mathbf{O}$,
- $(\mathbf{A}^T)^+ = (\mathbf{A}^+)^T$,
- $(\lambda\mathbf{A})^+ = \frac{1}{\lambda}\mathbf{A}^+$, $\lambda \neq 0$,
- $\mathbf{A}^+ = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T$, $\text{rg}(\mathbf{A}) = n \leq m$,
- $\mathbf{A}^+ = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T)^{-1}$, $\text{rg}(\mathbf{A}) = m \leq n$,
- $\mathbf{A}^+ = \mathbf{A}^{-1}$, $\text{rg}(\mathbf{A}) = m = n$,
- $(\mathbf{A}\mathbf{B})^+ = \mathbf{B}^+\mathbf{A}^+$, $\mathbf{A} \in \mathbb{R}^{m \times r}$, $\mathbf{B} \in \mathbb{R}^{r \times n}$, $\text{rg}(\mathbf{A}) = \text{rg}(\mathbf{B}) = r$,
- $(\mathbf{U}\mathbf{A}\mathbf{V})^+ = \mathbf{V}^T\mathbf{A}^+\mathbf{U}^T$ für orthogonale Matrizen \mathbf{U} , \mathbf{V} ,
- $\begin{pmatrix} \mathbf{A} \\ \mathbf{O} \end{pmatrix}^+ = (\mathbf{A}^+, \mathbf{O})$,
- $(\mathbf{a}\mathbf{b}^T)^+ = \frac{\mathbf{b}\mathbf{a}^T}{(\mathbf{a}^T\mathbf{a})(\mathbf{b}^T\mathbf{b})}$ für beliebige Vektoren $\mathbf{a} \in \mathbb{R}^m$ und $\mathbf{b} \in \mathbb{R}^n$ mit $\mathbf{a}, \mathbf{b} \neq \mathbf{o}$.

10.2. Störungstheorie

Wir wollen wieder untersuchen, welchen Einfluss Störungen in der Matrix \mathbf{A} und im Vektor \mathbf{y} auf die Lösung haben. Dazu betrachten wir neben dem Ausgleichsproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

mit der Normallösung $\mathbf{x} = \mathbf{A}^+\mathbf{y}$ ein gestörtes Problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|(\mathbf{y} + \delta\mathbf{y}) - (\mathbf{A} + \delta\mathbf{A})\mathbf{x}\|_2$$

mit der Normallösung $\mathbf{x} + \delta\mathbf{x} = (\mathbf{A} + \delta\mathbf{A})^+(\mathbf{y} + \delta\mathbf{y})$. Durch die Störungen $\delta\mathbf{A}$ und $\delta\mathbf{y}$ wird der Fehler

$$\delta\mathbf{x} = (\mathbf{A} + \delta\mathbf{A})^+(\mathbf{y} + \delta\mathbf{y}) - \mathbf{A}^+\mathbf{y} = [(\mathbf{A} + \delta\mathbf{A})^+ - \mathbf{A}^+] \mathbf{y} + (\mathbf{A} + \delta\mathbf{A})^+ \delta\mathbf{y}$$

erzeugt. Zunächst untersuchen wir, wie sich die Pseudoinverse einer Matrix gegenüber Störungen verhält. Dazu zwei Beispiele.

10.3. Beispiel: Es sei $\mathbf{A} \neq \mathbf{O}$ durch die Singulärwertzerlegung $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ mit

$$\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0), \quad \sigma_1 \geq \dots \geq \sigma_r > 0$$

gegeben. Es gilt $\text{rg}(\mathbf{A}) = r$. Wir betrachten die spezielle Störung $\delta\mathbf{A} = \mathbf{U}\delta\mathbf{\Sigma}\mathbf{V}^T$ mit

$$(\delta\mathbf{\Sigma})_{ij} = \begin{cases} \varepsilon & \text{für } i = j = r \\ 0 & \text{sonst} \end{cases}.$$

Dann folgt $\mathbf{A} + \delta\mathbf{A} = \mathbf{U}(\mathbf{\Sigma} + \delta\mathbf{\Sigma})\mathbf{V}^T$ mit

$$\mathbf{\Sigma} + \delta\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_{r-1}, \sigma_r + \varepsilon, 0, \dots, 0).$$

Weiterhin gilt $\|\mathbf{A}^+\|_2 = \|\mathbf{\Sigma}^+\|_2 = 1/\sigma_r$ und $\|\delta\mathbf{A}\|_2 = \|d\text{Sigma}\|_2 = |\varepsilon|$. Falls außerdem

$$\kappa = \|\mathbf{A}^+\|_2 \|\delta\mathbf{A}\|_2 = \frac{|\varepsilon|}{\sigma_r} < 1$$

gilt, folgt

$$\sigma_r + \varepsilon \geq \sigma_r - |\varepsilon| = \sigma_r(1 - \kappa) > 0.$$

Damit gilt

$$\text{rg}(\mathbf{A} + \delta\mathbf{A}) = \text{rg}(\mathbf{\Sigma} + \delta\mathbf{\Sigma}) = \text{rg}(\mathbf{\Sigma}) = \text{rg}(\mathbf{A}) = r.$$

Für die Pseudoinverse der gestörten Matrix gilt dann

$$(\mathbf{A} + \delta\mathbf{A})^+ = \mathbf{V}(\mathbf{\Sigma} + \delta\mathbf{\Sigma})^+\mathbf{U}^T$$

mit

$$(\mathbf{\Sigma} + \delta\mathbf{\Sigma})^+ = \text{diag}(1/\sigma_1, \dots, 1/\sigma_{r-1}, 1/(\sigma_r + \varepsilon), 0, \dots, 0).$$

Für die Norm der Pseudoinversen der gestörten Matrix erhält man

$$\begin{aligned} \|(\mathbf{A} + \delta\mathbf{A})^+\|_2 &= \|(\mathbf{\Sigma} + \delta\mathbf{\Sigma})^+\|_2 \\ &= \max \left\{ \frac{1}{\sigma_{r-1}}, \frac{1}{\sigma_r + \varepsilon} \right\} \\ &\leq \max \left\{ \frac{1}{\sigma_r}, \frac{1}{\sigma_r + \varepsilon} \right\} \\ &\leq \frac{1}{\sigma_r - |\varepsilon|} = \frac{1}{\sigma_r(1 - \kappa)}. \end{aligned}$$

Damit gilt

$$\|(\mathbf{A} + \delta\mathbf{A})^+\|_2 \leq \frac{\|\mathbf{A}^+\|_2}{1 - \kappa}.$$

Wir wollen noch abschätzen, wie stark sich die Pseudoinverse der gestörten Matrix $\mathbf{A} + \delta\mathbf{A}$ von der Pseudoinversen der Matrix \mathbf{A} unterscheidet. Es gilt

$$(\mathbf{A} + \delta\mathbf{A})^+ - \mathbf{A}^+ = \mathbf{V}(\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma})^+ \mathbf{U}^T - \mathbf{V}\boldsymbol{\Sigma}^+ \mathbf{U}^T = \mathbf{V} [(\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma})^+ - \boldsymbol{\Sigma}^+] \mathbf{U}^T$$

mit

$$((\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma})^+ - \boldsymbol{\Sigma}^+)_{ij} = \begin{cases} \frac{1}{\sigma_r + \varepsilon} - \frac{1}{\sigma_r} & \text{für } i = j = r \\ 0 & \text{sonst} \end{cases}.$$

Damit folgt

$$\|(\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma})^+ - \boldsymbol{\Sigma}^+\|_2 = \left| \frac{1}{\sigma_r + \varepsilon} - \frac{1}{\sigma_r} \right| = \frac{|\varepsilon|}{\sigma_r |\sigma_r + \varepsilon|} \leq \frac{|\varepsilon|}{\sigma_r^2 (1 - \kappa)}$$

und

$$\|(\mathbf{A} + \delta\mathbf{A})^+ - \mathbf{A}^+\|_2 = \|(\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma})^+ - \boldsymbol{\Sigma}^+\|_2 \leq \frac{\|\mathbf{A}^+\|_2}{1 - \kappa} \|\delta\mathbf{A}\|_2.$$

Diese Abschätzungen entsprechen genau den Abschätzungen, die wir bei der Behandlung der Störungen einer regulären Matrix für die Inverse erhalten hatten. \heartsuit

10.4. Beispiel: Es sei $\mathbf{A} \neq \mathbf{O}$ durch die Singulärwertzerlegung $\mathbf{A} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$ mit

$$\boldsymbol{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_r, 0, \dots, 0), \quad \sigma_1 \geq \dots \geq \sigma_r > 0$$

gegeben. Es gelte $\text{rg}(\mathbf{A}) = r < \min\{m, n\}$. Wir betrachten die spezielle Störung $\delta\mathbf{A} = \mathbf{U}\delta\boldsymbol{\Sigma}\mathbf{V}^T$ mit

$$(\delta\boldsymbol{\Sigma})_{ij} = \begin{cases} \varepsilon & \text{für } i = j = r + 1 \\ 0 & \text{sonst} \end{cases}.$$

Dann folgt $\mathbf{A} + \delta\mathbf{A} = \mathbf{U}(\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma})\mathbf{V}^T$ mit $\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_r, \varepsilon, 0, \dots, 0)$ und

$$(\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma})^+ = \text{diag}\left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_r}, \frac{1}{\varepsilon}, 0, \dots, 0\right).$$

Weiterhin gilt $\text{rg}(\mathbf{A} + \delta\mathbf{A}) = \text{rg}(\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma}) = r + 1 > r = \text{rg}(\boldsymbol{\Sigma}) = \text{rg}(\mathbf{A})$ und

$$\|(\mathbf{A} + \delta\mathbf{A})^+\|_2 = \|(\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma})^+\|_2 = \max\left\{\frac{1}{\sigma_r}, \frac{1}{|\varepsilon|}\right\} \geq \frac{1}{|\varepsilon|} = \frac{1}{\|\delta\mathbf{A}\|_2},$$

sowie

$$\|(\mathbf{A} + \delta\mathbf{A})^+ - \mathbf{A}^+\|_2 = \|(\boldsymbol{\Sigma} + \delta\boldsymbol{\Sigma})^+ - \boldsymbol{\Sigma}^+\|_2 = \frac{1}{|\varepsilon|} = \frac{1}{\|\delta\mathbf{A}\|_2}.$$

Diese rangerhöhende Störung hat einen katastrophalen Einfluss auf die Pseudoinverse. Je kleiner die Störung ist, desto mehr weicht die Pseudoinverse der gestörten Matrix von der Pseudoinversen der exakten Matrix ab. Die Pseudoinverse einer Matrix ist bezüglich dieser rangerhöhenden Störung unstetig und unbeschränkt. Es liegt ein extrem schlechtes Unstetigkeitsverhalten vor. ♡

Im allgemeinen Falle gibt es genau die zwei Klassen von Störungen, die in den beiden Beispielen zum Ausdruck kamen. Zum Beweis der entsprechenden Störungssätze benötigen wir den folgenden Hilfssatz.

10.5. Satz: *Es seien \mathbf{A} und $\delta\mathbf{A}$ (m, n) -Matrizen. Die Größen σ_i und $\delta\sigma_i$ seien die gemäß*

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_l$$

bzw.

$$\delta\sigma_1 \geq \delta\sigma_2 \geq \dots \geq \delta\sigma_l$$

geordneten Singulärwerte der Matrizen \mathbf{A} bzw. $\mathbf{A} + \delta\mathbf{A}$, wobei $l = \min\{m, n\}$. Dann gelten die Abschätzungen

$$|\delta\sigma_i| \leq \|\delta\mathbf{A}\|_2$$

für $i = 1, \dots, l$.

Beweis: Wir definieren die Matrizen

$$\mathbf{B} = \left(\begin{array}{c|c} \mathbf{O} & \mathbf{A} \\ \hline \mathbf{A}^T & \mathbf{O} \end{array} \right) \in \mathbb{R}^{(m+n) \times (m+n)}, \quad \delta\mathbf{B} = \left(\begin{array}{c|c} \mathbf{O} & \delta\mathbf{A} \\ \hline \delta\mathbf{A}^T & \mathbf{O} \end{array} \right) \in \mathbb{R}^{(m+n) \times (m+n)}.$$

Aus der Singulärwertzerlegung $\mathbf{A} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T$ folgt

$$\mathbf{B} = \left(\begin{array}{c|c} \mathbf{O} & \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T \\ \hline \mathbf{V}\boldsymbol{\Sigma}^T\mathbf{U}^T & \mathbf{O} \end{array} \right) = \mathbf{S}\mathbf{C}\mathbf{S}^T, \quad \mathbf{S} = \left(\begin{array}{c|c} \mathbf{U} & \mathbf{O} \\ \hline \mathbf{O} & \mathbf{V} \end{array} \right), \quad \mathbf{C} = \left(\begin{array}{c|c} \mathbf{O} & \boldsymbol{\Sigma} \\ \hline \boldsymbol{\Sigma}^T & \mathbf{O} \end{array} \right).$$

Wegen der Orthogonalität der Matrix \mathbf{S} sind die Matrizen \mathbf{B} und \mathbf{C} ähnlich, sie besitzen daher die gleichen Eigenwerte. Die Eigenwertgleichung $\mathbf{C}\mathbf{z} = \lambda\mathbf{z}$ lässt sich mit

$$\mathbf{z} = \begin{pmatrix} \mathbf{y} \\ \mathbf{x} \end{pmatrix} \quad \text{als} \quad \begin{pmatrix} \boldsymbol{\Sigma}\mathbf{x} \\ \boldsymbol{\Sigma}^T\mathbf{y} \end{pmatrix} = \lambda \begin{pmatrix} \mathbf{y} \\ \mathbf{x} \end{pmatrix} \quad (10.1)$$

schreiben. Für $m \geq n$ hat Σ die Form

$$\Sigma = \begin{pmatrix} D \\ O \end{pmatrix}, \quad D = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_n \end{pmatrix}.$$

Dann folgt aus Gleichung 10.1 mit $\mathbf{y} = (\bar{\mathbf{y}}^T, \tilde{\mathbf{y}}^T)^T$

$$\begin{aligned} D\mathbf{x} &= \lambda\bar{\mathbf{y}}, \\ \mathbf{o} &= \lambda\tilde{\mathbf{y}}, \\ D\bar{\mathbf{y}} &= \lambda\mathbf{x}. \end{aligned}$$

Hieraus ergeben sich die Eigenwerte

$$\begin{aligned} \lambda_i &= \sigma_i, \quad i = 1, \dots, n, \\ \lambda_i &= -\sigma_{i-n}, \quad i = n+1, \dots, 2n, \\ \lambda_i &= 0, \quad i = 2n+1, \dots, m+n. \end{aligned}$$

Für $m \leq n$ ergibt sich entsprechend

$$\begin{aligned} \lambda_i &= \sigma_i, \quad i = 1, \dots, m, \\ \lambda_i &= -\sigma_{i-m}, \quad i = m+1, \dots, 2m, \lambda_i = 0, \quad i = 2m+1, \dots, m+n. \end{aligned}$$

Die Matrizen \mathbf{C} und \mathbf{B} besitzen somit in jedem Falle die Eigenwerte $\pm\sigma_i$, $i = 1, \dots, l$. Genauso zeigt man, dass die Matrix $\mathbf{A} + \delta\mathbf{A}$ die Eigenwerte $\pm(\sigma_i + \delta\sigma_i)$, $i = 1, \dots, l$, besitzt. Mit Satz 9.23 folgt dann $|\delta\sigma_i| \leq \|\delta\mathbf{A}\|_2$ für $i = 1, \dots, l$. *

Nun untersuchen wir, wie sich rangerhöhende Störungen auswirken. Es gilt der folgende Satz.

10.6. Satz: *Es sei \mathbf{A} eine (m, n) -Matrix mit $\text{rg}(\mathbf{A}) = r < l = \min\{m, n\}$, und $\delta\mathbf{A}$ sei eine (m, n) -Störung mit $\text{rg}(\mathbf{A} + \delta\mathbf{A}) > \text{rg}(\mathbf{A})$. Dann gilt*

$$\|(\mathbf{A} + \delta\mathbf{A})^+\|_2 \geq \frac{1}{\|\delta\mathbf{A}\|_2}.$$

Beweis: Für die geordneten Singulärwerte σ_i bzw. $\tilde{\sigma}_i$, $i = 1, \dots, l$, der Matrizen \mathbf{A} bzw. $\tilde{\mathbf{A}} = \mathbf{A} + \delta\mathbf{A}$ gilt nach Satz 10.5

$$|\sigma_i - \tilde{\sigma}_i| \leq \|\delta\mathbf{A}\|_2, \quad i = 1, \dots, l.$$

Speziell gilt

$$|\sigma_{r+1} - \tilde{\sigma}_{r+1}| = \tilde{\sigma}_{r+1} \leq \|\delta\mathbf{A}\|_2.$$

Damit folgt

$$\|(\mathbf{A} + \delta\mathbf{A})^+\|_2 = \frac{1}{\tilde{\sigma}_{r+1}} \geq \frac{1}{\|\delta\mathbf{A}\|_2}.$$

*

Bemerkung: Wegen

$$\|(\mathbf{A} + \delta\mathbf{A})^+ - \mathbf{A}^+\|_2 \geq \|(\mathbf{A} + \delta\mathbf{A})^+\|_2 - \|\mathbf{A}^+\|_2 \geq \frac{1}{\|\delta\mathbf{A}\|_2} - \|\mathbf{A}^+\|_2$$

kann die Pseudoinverse der gestörten Matrix beliebig stark von der Pseudoinversen der exakten Matrix abweichen.

Für Störungen, die den Rang der Matrix nicht erhöhen gilt:

10.7. Satz: Für die (m, n) -Matrix \mathbf{A} und die (m, n) -Störung $\delta\mathbf{A}$ sei

$$\text{rg}(\mathbf{A} + \delta\mathbf{A}) \leq \text{rg}(\mathbf{A}), \quad \kappa = \|\mathbf{A}^+\|_2 \|\delta\mathbf{A}\|_2 < 1.$$

Dann gilt

1.

$$\text{rg}(\mathbf{A} + \delta\mathbf{A}) = \text{rg}(\mathbf{A}),$$

2.

$$\|(\mathbf{A} + \delta\mathbf{A})^+\|_2 \leq \frac{\|\mathbf{A}^+\|_2}{1 - \kappa},$$

3.

$$\begin{aligned} (\mathbf{A} + \delta\mathbf{A})^+ - \mathbf{A}^+ &= -\mathbf{A}^+ \delta\mathbf{A} \mathbf{A}^+ + \mathbf{A}^+ \mathbf{A}^{+T} \delta\mathbf{A}^T (\mathbf{I} - \mathbf{A} \mathbf{A}^+) + \\ &\quad + (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \delta\mathbf{A}^T \mathbf{A}^{+T} \mathbf{A}^+ + O(\|\delta\mathbf{A}\|_2^2), \end{aligned}$$

4.

$$\|(\mathbf{A} + \delta\mathbf{A})^+ - \mathbf{A}^+\|_2 \leq \mu \frac{\|\mathbf{A}^+\|_2^2}{1 - \kappa} \|\delta\mathbf{A}\|_2$$

mit

$$\mu = \begin{cases} \frac{1+\sqrt{5}}{2} & fr \quad \operatorname{rg}(\mathbf{A}) < \min\{m, n\} \\ \sqrt{2} & fr \quad \operatorname{rg}(\mathbf{A}) = \min\{m, n\} < \max\{m, n\} \\ 1 & fr \quad \operatorname{rg}(\mathbf{A}) = m = n \end{cases} .$$

Beweis: Es sei $r = \operatorname{rg}(\mathbf{A})$. Mit $\sigma_i, i = 1, \dots, r$, bezeichnen wir wieder die geordneten Singulärwerte von der matrix \mathbf{A} und mit $\tilde{\sigma}_i = \sigma_i + \delta\sigma_i, i = 1, \dots, l$ die geordneten Singulärwerte der matrix $\tilde{\mathbf{A}} = \mathbf{A} + \delta\mathbf{A}$. Dann gilt

$$\tilde{\sigma}_i = \sigma_i + \delta\sigma_i \geq \sigma_i - |\delta\sigma_i|,$$

und wegen $|\delta\sigma_i| \leq \|\delta\mathbf{A}\|_2$ folgt

$$\tilde{\sigma}_i \geq \sigma_i - \|\delta\mathbf{A}\|_2 \geq \sigma_r - \|\delta\mathbf{A}\|_2 = \frac{1}{\|\mathbf{A}^+\|_2} - \|\delta\mathbf{A}\|_2 = \frac{1 - \kappa}{\|\mathbf{A}^+\|_2} > 0$$

für $i = 1, \dots, r$. Folglich gilt $\operatorname{rg}(\mathbf{A} + \delta\mathbf{A}) \geq r$ und mit der Bedingung $\operatorname{rg}(\mathbf{A} + \delta\mathbf{A}) \leq r$ folgt $\operatorname{rg}(\mathbf{A} + \delta\mathbf{A}) = r$. Damit ist $\tilde{\sigma}_{r+1} = \dots = \tilde{\sigma}_l = 0$ und

$$\|(\mathbf{A} + \delta\mathbf{A})^+\|_2 = \frac{1}{\tilde{\sigma}_r} \leq \frac{\|\mathbf{A}^+\|_2}{1 - \kappa}.$$

Damit sind die ersten beiden Aussagen bewiesen.

Es sei nun

$$\tilde{\mathbf{A}} = \mathbf{A} + \delta\mathbf{A}, \quad \mathbf{G} = \tilde{\mathbf{A}}^+ - \mathbf{A}^+$$

sowie

$$\mathbf{P} = \mathbf{A}\mathbf{A}^+, \quad \mathbf{S} = \mathbf{A}^+\mathbf{A}, \quad \tilde{\mathbf{P}} = \tilde{\mathbf{A}}\tilde{\mathbf{A}}^+, \quad \tilde{\mathbf{S}} = \tilde{\mathbf{A}}^+\tilde{\mathbf{A}}.$$

Damit folgt

$$\begin{aligned} \mathbf{G} &= [\tilde{\mathbf{S}} + (\mathbf{I} - \tilde{\mathbf{S}})](\tilde{\mathbf{A}}^+ - \mathbf{A}^+)[\mathbf{P} + (\mathbf{I} - \mathbf{P})] \\ &= \tilde{\mathbf{S}}\tilde{\mathbf{A}}^+\mathbf{P} + \tilde{\mathbf{S}}\tilde{\mathbf{A}}^+(\mathbf{I} - \mathbf{P}) - \tilde{\mathbf{S}}\mathbf{A}^+\mathbf{P} - \tilde{\mathbf{S}}\mathbf{A}^+(\mathbf{I} - \mathbf{P}) + \\ &\quad + (\mathbf{I} - \tilde{\mathbf{S}})\tilde{\mathbf{A}}^+\mathbf{P} + (\mathbf{I} - \tilde{\mathbf{S}})\tilde{\mathbf{A}}^+(\mathbf{I} - \mathbf{P}) \\ &\quad - (\mathbf{I} - \tilde{\mathbf{S}})\mathbf{A}^+\mathbf{P} - (\mathbf{I} - \tilde{\mathbf{S}})\mathbf{A}^+(\mathbf{I} - \mathbf{P}). \end{aligned}$$

Nun gilt

$$\tilde{\mathbf{S}}\tilde{\mathbf{A}}^+ = \tilde{\mathbf{A}}^+\tilde{\mathbf{A}}\tilde{\mathbf{A}}^+ = \tilde{\mathbf{A}}^+, \quad \tilde{\mathbf{A}}^+\mathbf{P} = \tilde{\mathbf{A}}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+,$$

und damit

$$(I - \tilde{S})\tilde{A}^+ = O, \quad A^+(I - P) = O.$$

Somit ergibt sich

$$\begin{aligned} G &= \underbrace{(\tilde{A}^+ P - \tilde{S} A^+)}_{G_1} + \underbrace{\tilde{A}^+(I - P)}_{G_2} + \underbrace{(I - \tilde{S})(-A^+)}_{G_3} \\ &= G_1 + G_2 + G_3. \end{aligned}$$

Für die einzelnen Summanden erhalten wir weiter:

$$\begin{aligned} G_1 &= \tilde{A}^+ P - \tilde{S} A^+, \\ &= \tilde{A}^+ A A^+ - \tilde{A}^+ \tilde{A} A^+, \\ &= \tilde{A}^+ (A - \tilde{A}) A^+, \\ &= -\tilde{A}^+ \delta A A^+. \\ G_2 &= \tilde{A}^+ (I - P), \\ &= \tilde{A}^+ \tilde{A} \tilde{A}^+ (I - P), \\ &= \tilde{A}^+ \tilde{A}^{+T} \tilde{A}^T (I - P). \end{aligned}$$

Berücksichtigt man, dass

$$A^T (I - P) = A^T - A^T A A^+ = A^T - A^T A^{+T} A^T = A^T - A^T = O$$

gilt, so folgt

$$\begin{aligned} G_2 &= \tilde{A}^+ \tilde{A}^{+T} (\tilde{A}^T - A^T) (I - P), \\ &= \tilde{A}^+ \tilde{A}^{+T} \delta A^T (I - P). \end{aligned}$$

Für G_3 folgt mit

$$(I - \tilde{S})\tilde{A}^T = \tilde{A}^T - \tilde{A}^+ \tilde{A} \tilde{A}^T = \tilde{A}^T - \tilde{A}^T \tilde{A}^{+T} \tilde{A}^T = \tilde{A}^T - \tilde{A}^T = O$$

$$\begin{aligned} G_3 &= -(I - \tilde{S}) A^+, \\ &= -(I - \tilde{S}) A^+ A A^+, \\ &= -(I - \tilde{S}) A^T A^{+T} A^+, \\ &= -(I - \tilde{S}) (A^T - \tilde{A}^T) A^{+T} A^+, \\ &= -(I - \tilde{S}) \delta A^T A^{+T} A^+. \end{aligned}$$

Mit Hilfe der Abschätzung

$$\|\tilde{\mathbf{A}}^+\|_2 \leq \frac{\|\mathbf{A}^+\|_2}{1 - \kappa}$$

erhalten wir weiter

$$\|\mathbf{G}_1\|_2 \leq \|\tilde{\mathbf{A}}^+\|_2 \|\mathbf{A}^+\|_2 \|\delta\mathbf{A}\|_2 \leq \frac{\|\mathbf{A}^+\|_2}{1 - \kappa} \|\mathbf{A}^+\|_2 \|\delta\mathbf{A}\|_2 = \alpha$$

und

$$\|\mathbf{G}_3\|_2 \leq \|\mathbf{I} - \tilde{\mathbf{S}}\|_2 \|\mathbf{A}^+\|_2^2 \|\delta\mathbf{A}\|_2 = \|\mathbf{A}^+\|_2^2 \|\delta\mathbf{A}\|_2 \leq \alpha.$$

(Man beachte, dass wegen $(\mathbf{I} - \tilde{\mathbf{S}})^2 = \mathbf{I} - \tilde{\mathbf{S}}$ gilt: $\|\mathbf{I} - \tilde{\mathbf{S}}\|_2 = 1$.)

Bei der Abschätzung von \mathbf{G}_2 würde der Faktor $\|\tilde{\mathbf{A}}^+\|_2^2$ auftauchen. Um dies zu vermeiden, formen wir \mathbf{G}_2 weiter um. Es gilt

$$\begin{aligned} \|\mathbf{G}_2\|_2 &= \|\tilde{\mathbf{A}}^+(\mathbf{I} - \mathbf{P})\|_2, \\ &= \|\tilde{\mathbf{A}}^+ \tilde{\mathbf{P}}(\mathbf{I} - \mathbf{P})\|_2, \\ &\leq \|\tilde{\mathbf{A}}^+\|_2 \|\tilde{\mathbf{P}}(\mathbf{I} - \mathbf{P})\|_2. \end{aligned}$$

In Übungsaufgabe 4 ist zu zeigen, dass $\|\tilde{\mathbf{P}}(\mathbf{I} - \mathbf{P})\|_2 = \|\mathbf{P}(\mathbf{I} - \tilde{\mathbf{P}})\|_2$ gilt. Damit folgt

$$\begin{aligned} \|\mathbf{G}_2\|_2 &\leq \|\tilde{\mathbf{A}}^+\|_2 \|\mathbf{P}(\mathbf{I} - \tilde{\mathbf{P}})\|_2, \\ &= \|\tilde{\mathbf{A}}^+\|_2 \|\mathbf{A}\mathbf{A}^+(\mathbf{I} - \tilde{\mathbf{P}})\|_2, \\ &= \|\tilde{\mathbf{A}}^+\|_2 \|\mathbf{A}^{+T} \mathbf{A}^T(\mathbf{I} - \tilde{\mathbf{P}})\|_2, \\ &= \|\tilde{\mathbf{A}}^+\|_2 \|\mathbf{A}^{+T}(\mathbf{A}^T - \tilde{\mathbf{A}}^T)(\mathbf{I} - \tilde{\mathbf{P}})\|_2, \\ &= \|\tilde{\mathbf{A}}^+\|_2 \|\mathbf{A}^{+T} \delta\mathbf{A}^T(\mathbf{I} - \tilde{\mathbf{P}})\|_2, \\ &\leq \frac{\|\mathbf{A}^+\|_2}{1 - \kappa} \|\mathbf{A}^+\|_2 \|\delta\mathbf{A}\|_2 \|\mathbf{I} - \tilde{\mathbf{P}}\|_2, \\ &= \frac{\|\mathbf{A}^+\|_2}{1 - \kappa} \|\mathbf{A}^+\|_2 \|\delta\mathbf{A}\|_2 = \alpha. \end{aligned}$$

Dabei wurde wieder ausgenutzt, dass $\tilde{\mathbf{A}}^T(\mathbf{I} - \tilde{\mathbf{P}}) = \mathbf{O}$ und $\|\mathbf{I} - \tilde{\mathbf{P}}\|_2 = 1$ gilt. Insgesamt erhalten wir die Abschätzung

$$\|\mathbf{G}\|_2 \leq 3\alpha.$$

Die verbesserten Werte sind in Übungsaufgabe 5 zu zeigen. Aus

$$\|(\mathbf{A} + \delta\mathbf{A})^+ - \mathbf{A}^+\|_2 \leq 3 \frac{\|\mathbf{A}^+\|_2}{1 - \kappa} \|\mathbf{A}^+\|_2 \|\delta\mathbf{A}\|_2$$

folgt

$$(\mathbf{A} + \delta\mathbf{A})^+ = \mathbf{A}^+ + O(\|\delta\mathbf{A}\|_2).$$

Setzt man dies in die Darstellung

$$\mathbf{G} = (\mathbf{A} + \delta\mathbf{A})^+ - \mathbf{A}^+ = \mathbf{G}_1 + \mathbf{G}_2 + \mathbf{G}_3$$

ein, so erhält man die linearisierte Formel

$$\begin{aligned} (\mathbf{A} + \delta\mathbf{A})^+ - \mathbf{A}^+ &= -\mathbf{A}^+ \delta\mathbf{A} \mathbf{A}^+ + \mathbf{A}^+ \mathbf{A}^{+T} \delta\mathbf{A}^T (\mathbf{I} - \mathbf{A} \mathbf{A}^+) + \\ &\quad (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \delta\mathbf{A}^T \mathbf{A}^{+T} \mathbf{A}^+ + O(\|\delta\mathbf{A}\|_2^2). \end{aligned}$$

✱

Bemerkungen: (i) Die Bedingung $\text{rg}(\mathbf{A} + \delta\mathbf{A}) \leq \text{rg}(\mathbf{A})$ ist insbesondere erfüllt, falls die Matrix \mathbf{A} den maximalen Rang $r = l = \min\{m, n\}$ hat. In diesem Falle erhöht keine Störung den Rang der Matrix. Die Bestimmung der Pseudoinversen ist dann eine lokal lipschitzstetige Aufgabe, daher ein korrekt gestelltes Problem. Das folgt auch schon aus der expliziten Darstellung von \mathbf{A}^+ für diese Fälle:

$$\begin{aligned} \mathbf{A}^+ &= (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T, & r = n, \\ \mathbf{A}^+ &= \mathbf{A}^T (\mathbf{A} \mathbf{A}^T)^{-1}, & r = m. \end{aligned}$$

(ii) Bei rangdefizienten Problemen, also $\text{rg}(\mathbf{A}) < l = \min\{m, n\}$, hat man immer mit Störungen zu rechnen, die den Rang der Matrix erhöhen. In diesem Falle ist die Norm der Pseudoinversen unbeschränkt und $(\mathbf{A} + \delta\mathbf{A})^+$ weicht beliebig stark von \mathbf{A}^+ ab. Das Bestimmen der Pseudoinversen für rangdefiziente Probleme ist daher ein inkorrekt gestelltes Problem. Wir wollen nun noch das Störungsverhalten des gesamten linearen Ausgleichsproblems betrachten.

10.8. Satz: *Gegeben seien die linearen Ausgleichsprobleme*

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

und

$$\min_{\mathbf{x} + \delta\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} + \delta\mathbf{y} - (\mathbf{A} + \delta\mathbf{A})(\mathbf{x} + \delta\mathbf{x})\|_2$$

mit den Normallösungen

$$\mathbf{x} = \mathbf{A}^+ \mathbf{y}, \quad \mathbf{x} + \delta \mathbf{x} = (\mathbf{A} + \delta \mathbf{A})^+ (\mathbf{y} + \delta \mathbf{y}).$$

Die Störung $\delta \mathbf{A}$ genüge den Bedingungen

$$\text{rg}(\mathbf{A} + \delta \mathbf{A}) \leq \text{rg}(\mathbf{A}), \quad \kappa = \|\mathbf{A}^+\|_2 \|\delta \mathbf{A}\|_2 < 1.$$

Dann gilt

1. Der Fehler $\delta \mathbf{x}$ besitzt die Darstellung

$$\delta \mathbf{x} = \delta \mathbf{x}' + O(\|\delta \mathbf{A}\|_2 (\|\delta \mathbf{A}\|_2 + \|\delta \mathbf{y}\|_2))$$

mit dem bezüglich $\delta \mathbf{A}$ und $\delta \mathbf{y}$ linearen Teil

$$\delta \mathbf{x}' = \mathbf{A}^+ (-\delta \mathbf{A} \mathbf{x} + \delta \mathbf{y}) + \mathbf{A}^+ \mathbf{A}^{+T} \delta \mathbf{A}^T \mathbf{r} + (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \delta \mathbf{A}^T \mathbf{A}^{+T} \mathbf{x}.$$

Dabei bezeichnet $\mathbf{r} = \mathbf{y} - \mathbf{A} \mathbf{x} = (\mathbf{I} - \mathbf{A} \mathbf{A}^+) \mathbf{x}$ wie üblich das Residuum von \mathbf{x} .

2. Der Fehler $\delta \mathbf{x}$ genügt der Abschätzung

$$\|\delta \mathbf{x}\|_2 \leq \frac{\|\mathbf{A}^+\|_2}{1 - \kappa} [\|\delta \mathbf{A}\|_2 (\omega \|\mathbf{x}\|_2 + \|\mathbf{A}^+\|_2 \|\mathbf{r}\|_2) + \|\mathbf{y}\|_2]$$

und für $\mathbf{x} \neq \mathbf{o}$

$$\frac{\|\delta \mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \frac{1}{1 - \kappa} \left[\left(\omega \text{cond}(\mathbf{A}) + \frac{\|\mathbf{r}\|_2}{\|\mathbf{A}\|_2 \|\mathbf{x}\|_2} \text{cond}^2(\mathbf{A}) \right) \frac{\|\delta \mathbf{A}\|_2}{\|\mathbf{A}\|_2} + \frac{\|\mathbf{A}^+\|_2 \|\mathbf{y}\|_2}{\|\mathbf{x}\|_2} \frac{\|\delta \mathbf{y}\|_2}{\|\mathbf{y}\|_2} \right]$$

mit $\text{cond}(\mathbf{A}) = \|\mathbf{A}\|_2 \|\mathbf{A}^+\|_2$ und

$$\omega = \begin{cases} \sqrt{2} & \text{fr } \text{rg}(\mathbf{A}) < n \\ 1 & \text{fr } \text{rg}(\mathbf{A}) = n \end{cases}.$$

Beweis: Mit den Bezeichnungen aus dem Beweis von Satz 10.7 gilt

$$\begin{aligned} \delta \mathbf{x} &= (\mathbf{A} + \delta \mathbf{A})^+ (\mathbf{y} + \delta \mathbf{y}) - \mathbf{A}^+ \mathbf{y}, \\ &= [(\mathbf{A} + \delta \mathbf{A})^+ - \mathbf{A}^+] \mathbf{y} + (\mathbf{A} + \delta \mathbf{A})^+ \delta \mathbf{y}, \\ &= \mathbf{G} \mathbf{y} + \tilde{\mathbf{A}}^+ \delta \mathbf{y}, \\ &= \mathbf{G}_1 \mathbf{y} + \mathbf{G}_2 \mathbf{y} + \mathbf{G}_3 \mathbf{y} + \tilde{\mathbf{A}}^+ \delta \mathbf{y}. \end{aligned}$$

Nun gilt

$$\begin{aligned}
 \|G_1 \mathbf{y}\|_2 &= \|\tilde{\mathbf{A}}^+ \delta \mathbf{A} \mathbf{A}^+ \mathbf{y}\|_2, \\
 &= \|\tilde{\mathbf{A}}^+ \delta \mathbf{A} \mathbf{x}\|_2, \\
 &\leq \|\tilde{\mathbf{A}}^+\|_2 \|\delta \mathbf{A}\|_2 \|\mathbf{x}\|_2, \\
 &\leq \frac{\|\mathbf{A}^+\|_2}{1-\kappa} \|\delta \mathbf{A}\|_2 \|\mathbf{x}\|_2.
 \end{aligned}$$

und

$$\begin{aligned}
 \|G_3 \mathbf{y}\|_2 &= \|(\mathbf{I} - \tilde{\mathbf{S}}) \delta \mathbf{A}^T \mathbf{A}^{+T} \mathbf{A}^+ \mathbf{y}\|_2, \\
 &= \|\delta \mathbf{A}\|_2 \|\mathbf{A}^+\|_2 \|\mathbf{x}\|_2.
 \end{aligned}$$

Weiter gilt

$$\|G_1 \mathbf{y} + G_3 \mathbf{y}\|_2^2 = \|G_1 \mathbf{y}\|_2^2 + 2\mathbf{y}^T G_1^T G_3 \mathbf{y} + \|G_3 \mathbf{y}\|_2^2.$$

Wegen

$$G_1^T G_3 = \mathbf{A}^{+T} \delta \mathbf{A}^T \underbrace{\tilde{\mathbf{A}}^{+T} (\mathbf{I} - \tilde{\mathbf{S}})}_{=0} \delta \mathbf{A}^T \mathbf{A}^{+T} \mathbf{A}^+ = \mathbf{O}$$

folgt

$$\begin{aligned}
 \|G_1 \mathbf{y} + G_3 \mathbf{y}\|_2^2 &= \|G_1 \mathbf{y}\|_2^2 + \|G_3 \mathbf{y}\|_2^2, \\
 &\leq \left(\frac{\|\mathbf{A}^+\|_2 \|\delta \mathbf{A}\|_2 \|\mathbf{x}\|_2}{1-\kappa} \right)^2 + (\|\mathbf{A}^+\|_2 \|\delta \mathbf{A}\|_2 \|\mathbf{x}\|_2)^2, \\
 &\leq 2 \left(\frac{\|\mathbf{A}^+\|_2 \|\delta \mathbf{A}\|_2 \|\mathbf{x}\|_2}{1-\kappa} \right)^2, \\
 \|G_1 \mathbf{y} + G_3 \mathbf{y}\|_2 &\leq \sqrt{2} \frac{\|\mathbf{A}^+\|_2}{1-\kappa} \|\delta \mathbf{A}\|_2 \|\mathbf{x}\|_2.
 \end{aligned}$$

Weiter gilt

$$\begin{aligned}
 \|G_2 \mathbf{y}\|_2 &= \|\tilde{\mathbf{A}}^+ (\mathbf{I} - \mathbf{P}) \mathbf{y}\|_2, \\
 &= \|\tilde{\mathbf{A}}^+ (\mathbf{I} - \mathbf{A} \mathbf{A}^+) \mathbf{y}\|_2, \\
 &= \|\tilde{\mathbf{A}}^+ \mathbf{r}\|_2, \\
 &= \|\tilde{\mathbf{A}}^+ (\mathbf{I} - \mathbf{A} \mathbf{A}^+) \mathbf{r}\|_2, \\
 &= \|G_2 \mathbf{r}\|_2, \\
 &\leq \|G_2\|_2 \|\mathbf{r}\|_2, \\
 &\leq \frac{\|\mathbf{A}^+\|_2^2}{1-\kappa} \|\delta \mathbf{A}\|_2 \|\mathbf{r}\|_2
 \end{aligned}$$

und

$$\begin{aligned}\|\tilde{\mathbf{A}}^+ \delta \mathbf{y}\|_2 &\leq \|\tilde{\mathbf{A}}^+\|_2 \|\delta \mathbf{y}\|_2, \\ &\leq \frac{\|\mathbf{A}^+\|_2}{1-\kappa} \|\delta \mathbf{y}\|_2.\end{aligned}$$

Insgesamt erhalten wir

$$\begin{aligned}\|\delta \mathbf{x}\|_2 &\leq \|\mathbf{G}_1 \mathbf{y} + \mathbf{G}_3 \mathbf{y}\|_2 + \|\mathbf{G}_2 \mathbf{y}\|_2 + \|\tilde{\mathbf{A}}^+ \delta \mathbf{y}\|_2, \\ &\leq \sqrt{2} \frac{\|\mathbf{A}^+\|_2}{1-\kappa} \|\delta \mathbf{A}\|_2 \|\mathbf{x}\|_2 + \frac{\|\mathbf{A}^+\|_2^2}{1-\kappa} \|\delta \mathbf{A}\|_2 \|\mathbf{r}\|_2 + \frac{\|\mathbf{A}^+\|_2}{1-\kappa} \|\delta \mathbf{y}\|_2.\end{aligned}$$

Damit ist die Abschätzung für den absoluten Fehler bewiesen. Für $\text{rg}(\mathbf{A}) = n$ gilt

$$\tilde{\mathbf{A}}^+ = \left(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}} \right)^{-1} \tilde{\mathbf{A}}^T$$

und damit

$$\mathbf{I} - \tilde{\mathbf{S}} = \mathbf{I} - \left(\tilde{\mathbf{A}}^T \tilde{\mathbf{A}} \right)^{-1} \tilde{\mathbf{A}}^T \tilde{\mathbf{A}} = \mathbf{O}.$$

Dann ist $\mathbf{G}_3 = \mathbf{O}$ und der Faktor $\sqrt{2}$ in der Abschätzung ist durch 1 zu ersetzen. Die Abschätzung des relativen Fehlers folgt sofort aus der Abschätzung für den absoluten Fehler nach Division durch $\|\mathbf{x}\|_2$. Die linearisierte Darstellung von $\delta \mathbf{x}$ ergibt sich aus

$$\begin{aligned}\delta \mathbf{x} &= \mathbf{G}_1 \mathbf{y} + \mathbf{G}_2 \mathbf{y} + \mathbf{G}_3 \mathbf{y} + \tilde{\mathbf{A}}^+ \delta \mathbf{y}, \\ &= -\tilde{\mathbf{A}}^+ \delta \mathbf{A} \mathbf{A}^+ \mathbf{y} + \mathbf{G}_2 \mathbf{r} + (\mathbf{I} - \tilde{\mathbf{S}}) \delta \mathbf{A}^T \mathbf{A}^{+T} \mathbf{A}^+ \mathbf{y} + \tilde{\mathbf{A}}^+ \delta \mathbf{y}, \\ &= -\tilde{\mathbf{A}}^+ \delta \mathbf{A} \mathbf{x} + \tilde{\mathbf{A}}^+ \tilde{\mathbf{A}}^{+T} \delta \mathbf{A}^T \mathbf{r} + (\mathbf{I} - \tilde{\mathbf{S}}) \delta \mathbf{A}^T \mathbf{A}^{+T} \mathbf{x} + \delta \mathbf{y}.\end{aligned}$$

Mit $\tilde{\mathbf{A}}^+ = \mathbf{A}^+ + O(\|\delta \mathbf{A}\|_2)$ folgt

$$\begin{aligned}\delta \mathbf{x} &= \mathbf{A}^+ (-\delta \mathbf{A} \mathbf{x} + \delta \mathbf{y}) + \mathbf{A}^+ \mathbf{A}^{+T} \delta \mathbf{A}^T \mathbf{r} + (\mathbf{I} - \mathbf{A}^+ \mathbf{A}) \delta \mathbf{A}^T \mathbf{A}^{+T} \mathbf{x} \\ &\quad + O(\|\delta \mathbf{A}\|_2 (\|\delta \mathbf{A}\|_2 + \|\delta \mathbf{y}\|_2)).\end{aligned}$$

✱

Bemerkungen: (i) Für konsistente Probleme ($\mathbf{r} = \mathbf{o}$) erhält man die Abschätzung

$$\frac{\|\delta \mathbf{x}\|_2}{\|\mathbf{x}\|_2} \leq \frac{1}{1-\kappa} \left[\omega_{\text{cond}}(\mathbf{A}) \frac{\|\delta \mathbf{A}\|_2}{\|\mathbf{A}\|_2} + \frac{\|\mathbf{A}^+\|_2 \|\mathbf{y}\|_2}{\|\mathbf{x}\|_2} \frac{\|\delta \mathbf{y}\|_2}{\|\mathbf{y}\|_2} \right].$$

Bis auf den Faktor ω entspricht das genau der Fehlerabschätzung bei linearen Gleichungssystemen.

(ii) Für inkonsistente Probleme ($\mathbf{r} \neq \mathbf{o}$) tritt zusätzlich der Term

$$\frac{\|\mathbf{r}\|_2}{\|\mathbf{A}\|_2 \|\mathbf{x}\|_2} \text{cond}^2(\mathbf{A})$$

auf. Für hinreichend kleines $\|\mathbf{r}\|_2$ spielt er keine Rolle. Das Fehlerverhalten ist proportional zu $\text{cond}(\mathbf{A})$ wie bei konsistenten Problemen. Für stark inkonsistente Probleme wird dieser Term aber bestimmend. Das Fehlerverhalten ist dann proportional zu $\text{cond}^2(\mathbf{A})$ und somit bedeutend schlechter.

(iii) Der Fehlereinfluss von Störungen der Matrix hängt bei Ausgleichsproblemen über das Residuum auch von der rechten Seite \mathbf{y} ab.

10.3. Normalgleichungsverfahren

Wie wir im Abschnitt 10.1. gesehen haben, ist das Lösen des linearen Ausgleichsproblems

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

dem Lösen des Normalgleichungssystems

$$\mathbf{A}^T \mathbf{A} \mathbf{x} = \mathbf{A}^T \mathbf{y}$$

äquivalent. Im spaltenregulären Falle ($\text{rg}(\mathbf{A}) = n$) ist $\mathbf{A}^T \mathbf{A}$ regulär und positiv definit. Das Normalgleichungssystem ist dann eindeutig lösbar. Für $\text{rg}(\mathbf{A}) < n$ ist die Matrix $\mathbf{A}^T \mathbf{A}$ singulär. Das Normalgleichungssystem ist dann aber immer noch konsistent, hat aber unendlich viele Lösungen. Wir werden uns darum auf den spaltenregulären Fall beschränken. Hier bietet sich zum Lösen das CHOLESKY-Verfahren an.

10.9. Lösen des Normalgleichungssystems mit Hilfe einer CHOLESKY-Zerlegung:

Es ist das lineare Ausgleichsproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

mit der spaltenregulären Matrix (m, n) - \mathbf{A} zu lösen.

S0 *Berechne die (n, n) -Matrix $\mathbf{M} = \mathbf{A}^T \mathbf{A}$.*

S1 Bestimme eine CHOLESKY-Zerlegung $M = LL^T$ mit einer unteren Dreiecksmatrix L .

S2 Berechne $z = A^T y$.

S3 Löse die Dreieckssysteme $Lw = z$ und $L^T x = w$.

Aufwand:

S0 $\sim mn^2/2$ Additionen/Multiplikationen,

S1 $\sim n^3/6$ Additionen/Multiplikationen + n Quadratwurzeln,

S2 $\sim mn$ Additionen/Multiplikationen,

S3 $\sim n^2$ Additionen/Multiplikationen

Bemerkungen:

(i) In praktischen Fällen ist meist m größer als n . Dann liegt der Hauptaufwand im Berechnen der Matrix M .

(ii) Die Matrix M sollte als unteres Dreieck gesondert gespeichert werden, damit die Matrix A für eine Nachiteration weiter zur Verfügung steht.

Für das Rundungsfehlerverhalten dieses Verfahrens geben wir den folgenden Satz ohne Beweis an.

10.10. Satz: Für eine spaltenreguläre (m, n) -Matrix A ist der Algorithmus durchführbar, falls A im Sinne von

$$\hat{\kappa} = \text{eps} N_1 \text{cond}^2(A) \leq \frac{1}{2}$$

mit

$$N_1 = mn + F_1 \approx mn,$$

$$F_1 = n^{3/2} + n + n^{1/2}F \approx 2n^{3/2},$$

$$F = n + 1 + \frac{1}{2} \ln n \approx n$$

nicht zu schlecht konditioniert ist. Die berechneten Größen \hat{M} , \hat{L} , \hat{z} und \hat{x} genügen den Beziehungen

$$\begin{aligned} \hat{M} &= A^T A + \delta_0 M, & \|\delta_0 M\|_2 &\leq \text{eps} mn \|A\|_2^2, \\ \hat{z} &= (A + \delta_0 A)^T y, & \|\delta_0 A\|_2 &\leq \text{eps} m \sqrt{n} \|A\|_2, \\ \hat{L} \hat{L}^T &= M + \delta_1 M, & \|\delta_1 M\|_2 &\leq \text{eps} n F \|M\|_2, \\ (M + \delta_2 M) \hat{x} &= \hat{x}, & \|\delta_2 M\|_2 &\leq \text{eps} F_1 \|M\|_2. \end{aligned}$$

Bemerkung: Die Bedingung $\hat{\kappa} \leq 0.5$ schränkt die Klasse der mit diesem Algorithmus lösbaren Probleme wesentlich ein. Die Klasse der Probleme, für die das lineare Ausgleichsproblem korrekt gestellt ist, ist nach Satz 10.8 durch

$$\kappa = \|\mathbf{A}^+\|_2 \|\delta\mathbf{A}\|_2 < 1$$

charakterisiert. Nimmt man für $\delta\mathbf{A}$ nur den Darstellungsfehler auf dem Rechner an, so gilt $\|\delta\mathbf{A}\|_2 \leq \text{eps}\sqrt{n}\|\mathbf{A}\|_2$ und daher

$$\kappa = \text{eps}\sqrt{n}\text{cond}(\mathbf{A}) < 1.$$

Damit ist für alle Matrizen mit

$$\text{cond}(\mathbf{A}) < \frac{1}{\text{eps}\sqrt{n}}$$

das lineare Ausgleichsproblem korrekt gestellt. Aber nur für Matrizen mit

$$\text{cond}(\mathbf{A}) < \frac{1}{\sqrt{2\text{eps}mn}}$$

ist das Normalgleichungsverfahren anwendbar.

Mit Hilfe von Satz 10.10 lässt sich auch der erzeugte Rundungsfehler abschätzen.

10.11. Satz: *Unter den Voraussetzungen von Satz 10.10 gilt für den erzeugten Rundungsfehler $\delta\mathbf{x} = \hat{\mathbf{x}} - \mathbf{x}$, $\mathbf{x} = \mathbf{A}^+\mathbf{y} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{y}$, die Abschätzung*

$$\|\delta\mathbf{x}\|_2 \leq \text{eps} \frac{\text{cond}^2(\mathbf{A})}{1 - \hat{\kappa}} \left[N_1 \|\mathbf{x}\|_2 + m\sqrt{n} \frac{\|\mathbf{y}\|_2}{\|\mathbf{A}\|_2} \right].$$

Eine weitere Analyse zeigt, dass der erzeugte Rundungsfehler den unvermeidbaren Fehler beliebig stark übersteigen kann. Dieser Fall liegt besonders vor, falls die Norm des Residuums klein ist und $\|\mathbf{A}^+\mathbf{y}\|_2$ in der Größenordnung von $\|\mathbf{A}^+\|_2\|\mathbf{y}\|_2$ liegt. Andererseits ist das Normalgleichungsverfahren für stark inkonsistente Probleme ($\|\mathbf{r}\|_2 \gg 1$) mit kleinen Lösungen ($\|\mathbf{x}\|_2 \ll 1$) stabil und gutartig. Das sind gerade die Ausgleichsprobleme, bei denen die Fehlerverstärkung durch $\text{cond}^2(\mathbf{A})$ bestimmt wird.

Die numerische Instabilität des Normalgleichungsverfahrens lässt sich in ihrer Wirkung durch eine Nachiteration mildern.

10.12. Normalgleichungsverfahren mit Nachiteration:

Es ist das lineare Ausgleichsproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

mit der spaltenregulären (m, n) -Matrix \mathbf{A} zu lösen.

S0 Berechne die (n, n) -Matrix $M = A^T A$.

Bestimme eine CHOLESKY-Zerlegung $M = LL^T$ mit einer unteren Dreiecksmatrix L .

Setze $x^{(0)} = \mathbf{o}$ und $k = 0$.

S1 Berechne $r^{(k)} = \mathbf{y} - A\mathbf{x}^{(k)}$ und $\mathbf{g}^{(k)} = A^T r^{(k)}$.

S2 Löse die Dreieckssysteme

$$L\mathbf{w} = \mathbf{g}^{(k)}, \quad L^T \mathbf{h}^{(k)} = \mathbf{w}.$$

S3 Setze $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{h}^{(k)}$, $k = k + 1$ und gehe zu Schritt **S1**

Im ersten Schritt ergibt sich mit $\mathbf{x}^{(1)} = \mathbf{h}^{(0)}$ wieder die Lösung des üblichen Normalgleichungsverfahrens.

Nach hinreichend vielen Iterationen liefert dieses Verfahren eine Lösung $\mathbf{x}^{(k)}$, die als exakte Lösung eines benachbarten Ausgleichsproblems

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - (A + \delta A)\mathbf{x}\|_2$$

interpretierbar ist, falls die spaltenreguläre Matrix A der Bedingung

$$\bar{\kappa} = \text{eps} \bar{N}_1 \text{cond}^2(A) \leq 0.23$$

mit $\bar{N}_1 = (m + 1)(n + 1) + 2n^2$ genügt. Unter diesen Voraussetzungen liegt numerische Gutartigkeit vor.

Man beachte:

Hier spielt die Nachiteration eine ganz andere Rolle als bei linearen Gleichungssystemen. Bei linearen Gleichungssystemen diene sie zur Genauigkeitsverbesserung der Lösung eines ohnehin gutartigen Verfahrens. Hier wird dagegen erst durch die Nachiteration die Gutartigkeit des Verfahrens erreicht.

Die Konvergenzgeschwindigkeit des Verfahrens hängt entscheidend von $\bar{\kappa}$ ab. Die Klasse der lösbaren Probleme ist durch die Bedingung $\bar{\kappa} \leq 0.23$ wegen $\bar{\kappa} \sim \text{cond}^2(A)$ immer noch wesentlich eingeschränkt!

10.4. Orthogonalisierungsverfahren

Das in Abschnitt 8.2.5 beschriebene HOUSEHOLDER-Verfahren zum Lösen linearer Gleichungssysteme lässt sich unmittelbar auf das Lösen linearer Ausgleichsprobleme übertragen. Wegen der Invarianz der euklidischen Vektornorm gegenüber orthogonalen Transformationen sind nämlich die Probleme

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - A\mathbf{x}\|_2 \tag{10.2}$$

und

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{h} - \mathbf{S}\mathbf{x}\|_2, \quad (10.3)$$

wobei $\mathbf{h} = \mathbf{Q}^T \mathbf{y}$ und $\mathbf{S} = \mathbf{Q}^T \mathbf{A}$ mit einer orthogonalen (m, m) -Matrix \mathbf{Q} gilt, äquivalent. Es kommt nun darauf an, die orthogonale Matrix \mathbf{Q} so zu wählen, dass das Problem 10.3 möglichst einfach wird. Das ist zum Beispiel der Fall, wenn die Matrix \mathbf{S} die Struktur $\mathbf{S} = \begin{pmatrix} \mathbf{R} \\ \mathbf{O} \end{pmatrix}$ mit einer oberen (n, n) -Dreiecksmatrix \mathbf{R} besitzt. Mit der Partitionierung $\mathbf{h} = \begin{pmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \end{pmatrix}$, $\mathbf{h}_1 \in \mathbb{R}^n$ und $\mathbf{h}_2 \in \mathbb{R}^{m-n}$, gilt dann

$$\begin{aligned} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 &= \|\mathbf{h} - \mathbf{S}\mathbf{x}\|_2^2 \\ &= \left\| \begin{pmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \end{pmatrix} - \begin{pmatrix} \mathbf{R} \\ \mathbf{O} \end{pmatrix} \mathbf{x} \right\|_2^2 \\ &= \left\| \begin{pmatrix} \mathbf{h}_1 - \mathbf{R}\mathbf{x} \\ \mathbf{h}_2 \end{pmatrix} \right\|_2^2 \\ &= \|\mathbf{h}_1 - \mathbf{R}\mathbf{x}\|_2^2 + \|\mathbf{h}_2\|_2^2. \end{aligned}$$

Der zweite Summand in der letzten Gleichung ist konstant. Er spielt daher bei der Minimierung keine Rolle. Für den ersten Summanden gilt

$$\|\mathbf{h}_1 - \mathbf{R}\mathbf{x}\|_2^2 \geq 0, \quad \|\mathbf{h}_1 - \mathbf{R}\mathbf{x}\|_2^2 = 0 \iff \mathbf{R}\mathbf{x} = \mathbf{h}_1.$$

Die Lösung des linearen Ausgleichsproblems 10.2 ergibt sich somit nach der Transformation auf das Problem 10.3 aus dem linearen Gleichungssystem $\mathbf{R}\mathbf{x} = \mathbf{h}_1$. Der Algorithmus zur QR -Zerlegung aus Abschnitt 8.2.5 lässt sich direkt auf das Problem übertragen.

10.13. Lösen eines linearen Ausgleichsproblems mit HOUSEHOLDER-Transformationen:

Es ist das lineare Ausgleichsproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

mit der spaltenregulären (m, n) -Matrix \mathbf{A} zu lösen.

{Initialisierung}

Wähle eine Genauigkeitsschranke $\varepsilon > 0$.

{QR-Zerlegung}

for $k = 1$ **to** n **do**

```

{Berechnen der Transformationsmatrix  $\mathbf{H}^{(k)}$ }
 $\varrho_k = \sqrt{a_{kk}^2 + a_{k+1,k}^2 + \cdots + a_{mk}^2}$ 
if  $\varrho_k \leq \varepsilon$  then
  STOPP
endif
if  $a_{kk} > 0$  then
   $\varrho_k = -\varrho_k$ 
endif
 $a_{kk} = a_{kk} - \varrho_k$ 
 $\gamma_k = -\varrho_k \cdot a_{kk}$ 
{Transformation der Restmatrix}
for  $j = k + 1$  to  $n$  do
   $\beta = 0$ 
  for  $i = k$  to  $m$  do
     $\beta = \beta + a_{ik} \cdot a_{ij}$ 
  endfor
   $\beta = \beta / \gamma_k$ 
  for  $i = k$  to  $m$  do
     $a_{ij} = a_{ij} - \beta \cdot a_{ik}$ 
  endfor
  {Transformation des Vektors  $\mathbf{y}$ }
   $\beta = 0$ 
  for  $i = k$  to  $m$  do
     $\beta = \beta + a_{ik} \cdot y_{ij}$ 
  endfor
   $\beta = \beta / \gamma_k$ 
  for  $i = k$  to  $m$  do
     $y_i = y_i - \beta \cdot a_{ik}$ 
  endfor
endfor
endfor
 $\varrho_n = a_{nn}$ 
{Lösen von  $\mathbf{R}\mathbf{x} = \mathbf{c}$ }
for  $k = n$  to  $1$  step  $-1$  do
  for  $i = k + 1$  to  $n$  do
     $y_k = y_k - y_i \cdot a_{ki}$ 
  endfor
   $y_k = y_k / \varrho_k$ 
endfor

```

Aufwand:

Transf.: $\sim(m - n/3)n^2$ Additionen/Multiplikationen, n Quadratwurzeln,

Lösen: $\sim n^2/2$ Additionen/Multiplikationen.

Für $m \gg n$ ist der Aufwand damit etwa doppelt so groß wie der des Normalgleichungsverfahrens. Dafür zeigt aber eine Rundungsfehleranalyse, dass der Algorithmus gutartig ist und für eine große Klasse von Matrizen anwendbar ist. Wir geben den entsprechenden Satz ohne Beweis an.

10.14. Satz: Für eine spaltenreguläre (m, n) -Matrix \mathbf{A} ist der Algorithmus mit $r_{kk} \neq 0$ für $k = 1, \dots, n$ durchführbar, falls \mathbf{A} der Bedingung

$$\kappa = \text{eps}F \text{cond}(\mathbf{A}) < 1$$

mit

$$F = Km n^{3/2}, \quad K = 3.14(1 + 3/m) \approx 4$$

genügt. Für einen beliebigen Vektor $\mathbf{y} \in \mathbb{R}^m$ liefert der Algorithmus eine Lösung $\hat{\mathbf{x}} \in \mathbb{R}^n$, die als exakte Lösung eines benachbarten Problems

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|(\mathbf{y} + \delta\mathbf{y}) - (\mathbf{A} + \delta\mathbf{A})\mathbf{x}\|_2$$

interpretierbar ist. Die Störungen genügen den Abschätzungen

$$\begin{aligned} \|\delta\mathbf{A}\|_2 &\leq \text{eps}F_1 \|\mathbf{A}\|_2, \quad F_1 = F + n^{3/2} \approx F, \\ \|\delta\mathbf{y}\|_2 &\leq \text{eps} \frac{F}{\sqrt{n}} \|\mathbf{y}\|_2. \end{aligned}$$

Für das berechnete Residuum $\hat{\mathbf{r}} = \mathbf{r} + \delta\mathbf{r}$ gilt

$$\mathbf{r} + \delta\mathbf{r} = \mathbf{y} + \delta\mathbf{y} - (\mathbf{A} + \delta\mathbf{A})\hat{\mathbf{x}}$$

mit

$$\|\delta\mathbf{r}\|_2 \leq \text{eps} \frac{F}{\sqrt{n}} \|\mathbf{r}\|_2.$$

Bemerkungen (i) Das HOUSEHOLDER-Verfahren zum Lösen von linearen Ausgleichsproblemen ist ein gutartiger Prozess. Die Bedingung $\kappa = \text{eps}F \text{cond}(\mathbf{A}) < 1$ für die Anwendbarkeit des Verfahrens ist wesentlich schwächer als die Bedingung $\hat{\kappa} = \text{eps}N_1 \text{cond}^2(\mathbf{A}) \leq 0.5$ für die Anwendbarkeit des Normalgleichungsverfahrens.

Das HOUSEHOLDER-Verfahren ist trotz des etwa doppelt so großen Aufwands dem Normalgleichungsverfahren vorzuziehen.

(ii) Der Algorithmus zur QR -Zerlegung der rechteckigen (m, n) -Matrix A stimmt fast mit dem entsprechenden Algorithmus zur QR -Zerlegung einer quadratischen Matrix A überein. Es ist ein zusätzlicher n -ter Schritt auszuführen. Insbesondere darf die gleiche Speicherungsform gewählt werden.

(iii) In speziellen Anwendungen ist statt der Lösung x der Residuumsvektor $r = y - Ax$ gesucht; dieser lässt sich ohne Kenntnis von x berechnen. Es gilt

$$\begin{aligned} r &= y - Ax, \\ &= Qh - QSx, \\ &= Q \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} - Q \begin{pmatrix} R \\ O \end{pmatrix} R^{-1} h_1, \\ &= Q \begin{pmatrix} h_1 \\ h_2 \end{pmatrix} - Q \begin{pmatrix} h_1 \\ o \end{pmatrix}, \\ &= Q \begin{pmatrix} o \\ h_2 \end{pmatrix}. \end{aligned}$$

Auch bei Orthogonalisierungsverfahren ist eine Nachiteration möglich. Hier dient sie aber wieder nur einer Genauigkeitsverbesserung. Für spaltenreguläres A sei x^* die eindeutige Lösung des linearen Ausgleichsproblems. r^* sei das zugehörige minimale Residuum. Dann gilt $r^* + Ax^* = y$ und $A^T r^* = o$. Damit ist der Vektor $\begin{pmatrix} r^* \\ x^* \end{pmatrix}$

Lösung des linearen Gleichungssystems

$$\begin{pmatrix} I & A \\ A^T & O \end{pmatrix} \begin{pmatrix} r \\ x \end{pmatrix} = \begin{pmatrix} y \\ o \end{pmatrix}.$$

Es sei nun

$$A + \delta A = Q \begin{pmatrix} \mathcal{R} \\ O \end{pmatrix}$$

die näherungsweise QR -Zerlegung von A . Wendet man auf das obige Gleichungssystem ein übliches Iterationsverfahren mit

$$B = \begin{pmatrix} I & Q \begin{pmatrix} \mathcal{R} \\ O \end{pmatrix} \\ (\mathcal{R}^T, O) Q^T & O \end{pmatrix}$$

an, so erhält man:

10.15. Nachiteration für QR-Zerlegung bei linearen Ausgleichsproblemen:

Es ist das lineare Ausgleichsproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2$$

mit der spaltenregulären (m, n) -Matrix \mathbf{A} zu lösen.

S0 Berechne die QR-Zerlegung $\mathbf{A} = \mathbf{Q} \begin{pmatrix} \mathbf{R} \\ \mathbf{O} \end{pmatrix}$ mit einer orthogonalen (m, m) -

Matrix \mathbf{Q} und einer oberen (n, n) -Dreiecksmatrix \mathbf{R} .

Setze $\mathbf{x}^{(0)} = \mathbf{o}$, $\mathbf{r}^{(0)} = \mathbf{y}$ und $k = 0$.

S1 Berechne $\mathbf{e} = \mathbf{y} - \mathbf{r}^{(k)} - \mathbf{A}\mathbf{x}^{(k)}$ und $\mathbf{f} = -\mathbf{A}^T \mathbf{r}^{(k)}$.

Berechne

$$\mathbf{Q}^T \mathbf{e} = \begin{pmatrix} \mu \\ \mathbf{v} \end{pmatrix}, \quad \mu \in \mathbb{R}^n, \quad \mathbf{v} \in \mathbb{R}^{m-n}.$$

S2 Löse das Dreieckssystem $\mathbf{R}^T \mathbf{w} = \mathbf{f}$.

Löse das Dreieckssystem $\mathbf{R}\mathbf{h} = \mu - \mathbf{w}$.

S3 Berechne

$$\mathbf{s} = \mathbf{Q} \begin{pmatrix} \mathbf{w} \\ \mathbf{v} \end{pmatrix}.$$

Setze $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{h}$, $\mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} + \mathbf{s}$, $k = k + 1$ und gehe zu Schritt **S1**

10.5. Aufgaben

1. Man zeige: Für eine beliebige (m, n) -Matrix \mathbf{A} und einen beliebigen Vektor $\mathbf{y} \in \mathbb{R}^m$ ist das Gleichungssystem

$$\mathbf{A}^T \mathbf{A}\mathbf{x} = \mathbf{A}^T \mathbf{y}$$

stets lösbar.

2. Man zeige: Gilt für die (m, n) -Matrix \mathbf{A} $\text{rg}(\mathbf{A}) = n$, so ist die Matrix $\mathbf{M} = \mathbf{A}^T \mathbf{A}$ regulär und positiv definit.
3. Es liege das folgende mathematische Gesetz mit zwei unbekanntem Parametern α und β vor:

$$z = \alpha y + \beta$$

Weiterhin sei ein Satz von Messdaten $(y_i, z_i) \quad i = 1, \dots, m$ mit $y_i = i$ gegeben. Man bestimme mit Hilfe der linearen Ausgleichsrechnung die Parameter α und β !

- (a) Wie lauten die Normalgleichungen?
 (b) Wie lautet die CHOLESKY-Zerlegung der Matrix der Normalgleichungen? ($\mathbf{A}^T \mathbf{A} = \mathbf{L}\mathbf{L}^T$)
 (c) Man gebe eine Abschätzung von $\text{cond}_2(\mathbf{L})$ an.
 (d) Wie ändert sich die Kondition mit der Anzahl der Messdaten?

4. Für $\mathbf{A}, \tilde{\mathbf{A}}$ gelte $\text{rg}(\mathbf{A}) = \text{rg}(\tilde{\mathbf{A}}) = r$. Man zeige:

$$\|\tilde{\mathbf{P}}(\mathbf{I} - \mathbf{P})\|_2 = \|\mathbf{P}(\mathbf{I} - \tilde{\mathbf{P}})\|_2.$$

Hinweis: Unter Verwendung der Singulärwertzerlegung stelle man \mathbf{P} und $\tilde{\mathbf{P}}$ gemäß

$\mathbf{P} = \mathbf{U}\mathbf{\Pi}\mathbf{U}^T$ und $\tilde{\mathbf{P}} = \tilde{\mathbf{U}}\mathbf{\Pi}\tilde{\mathbf{U}}^T$ mit

$$\mathbf{\Pi} = \begin{pmatrix} \mathbf{I}^{(r)} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} \end{pmatrix}$$

dar und partitioniere $\mathbf{W} = \tilde{\mathbf{U}}^T \mathbf{U}$ gemäß

$$\mathbf{W} = \begin{pmatrix} \mathbf{W}_{11} & \mathbf{W}_{12} \\ \mathbf{W}_{21} & \mathbf{W}_{22} \end{pmatrix}.$$

Man beweise:

$$\|\tilde{\mathbf{P}}(\mathbf{I} - \mathbf{P})\| = \|\mathbf{W}_{12}\| \quad \text{und} \quad \|\mathbf{P}(\mathbf{I} - \tilde{\mathbf{P}})\| = \|\mathbf{W}_{21}\|$$

und folgere

$$\|\mathbf{W}_{12}\| = \|\mathbf{W}_{21}\|$$

aus der Orthogonalität von \mathbf{W} .

5. Für

$$\mathbf{y} \in \mathbb{R}^m \quad \text{und} \quad \mathbf{G}_1 = \tilde{\mathbf{A}}^+ \mathbf{P} - \tilde{\mathbf{S}} \mathbf{A}^+, \quad \mathbf{G}_2 = \tilde{\mathbf{A}}^+ (\mathbf{I} - \mathbf{P}), \quad \mathbf{G}_3 = -(\mathbf{I} - \tilde{\mathbf{S}}) \mathbf{A}^+$$

sowie

$$\mathbf{P} = \mathbf{A}\mathbf{A}^+, \quad \tilde{\mathbf{P}} = \tilde{\mathbf{A}}\tilde{\mathbf{A}}^+, \quad \tilde{\mathbf{S}} = \tilde{\mathbf{A}}^+ \tilde{\mathbf{A}}, \quad \mathbf{S} = \mathbf{A}^+ \mathbf{A} \quad \text{mit} \quad \|\mathbf{G}_i\| \leq \alpha$$

werde

$$\mathbf{z} = \mathbf{G}_1 \mathbf{y} + \mathbf{G}_3 \mathbf{y} = \mathbf{z}_1 + \mathbf{z}_2 + \mathbf{z}_3$$

gebildet. Man zeige:

- (a) $\mathbf{z}_1 \perp \mathbf{z}_3, \mathbf{z}_2 \perp \mathbf{z}_3$ folglich $\|\mathbf{z}\|^2 = \|\mathbf{z}_1 + \mathbf{z}_2\|^2 + \|\mathbf{z}_3\|^2$
 (b) mit $\mathbf{u} = \mathbf{P}\mathbf{y}, \mathbf{v} = (\mathbf{I} - \mathbf{P})\mathbf{y}$ gilt:

$$\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2 = \|\mathbf{y}\|^2 \text{ und } \mathbf{G}_1\mathbf{y} = \mathbf{G}_1\mathbf{u}, \mathbf{G}_2\mathbf{y} = \mathbf{G}_2\mathbf{v}, \mathbf{G}_3\mathbf{y} = \mathbf{G}_3\mathbf{u}.$$

Daraus folgere man weiter

$$\|\mathbf{z}\|^2 \leq \alpha^2 [(\|\mathbf{u}\| + \|\mathbf{v}\|)^2 + \|\mathbf{u}\|^2] = \alpha^2 \|\mathbf{y}\|^2 N$$

mit

$$N = (\cos \varphi + \sin \varphi)^2 + \cos^2 \varphi \leq \left[\frac{1 + \sqrt{5}}{2} \right]^2$$

(Dabei wurde $\|\mathbf{u}\| = \|\mathbf{y}\| \cos \varphi$ und $\|\mathbf{v}\| = \|\mathbf{y}\| \sin \varphi$ gesetzt.)

- (c) Im Falle $\text{rg}(\mathbf{A}) = n$ ist $\mathbf{z}_3 = \mathbf{o}$, also

$$\|\mathbf{z}\|^2 = \|\mathbf{z}_1 + \mathbf{z}_2\|^2 \leq 2\alpha^2 \|\mathbf{y}\|^2.$$

Im Falle $\text{rg}(\mathbf{A}) = m$ ist $\mathbf{z}_2 = \mathbf{o}$, also

$$\|\mathbf{z}\|^2 = \|\mathbf{z}_1\|^2 + \|\mathbf{z}_3\|^2 \leq 2\alpha^2 \|\mathbf{y}\|^2.$$

Im Falle $\text{rg}(\mathbf{A}) = m = n$ ist $\|\mathbf{z}\| = \|\mathbf{z}_1\| \leq \alpha \|\mathbf{y}\|$.

6. Es sei eine beliebige (m, n) -Matrix \mathbf{A} gegeben. Man zeige

$$\lim_{\rho \rightarrow 0} [\rho \mathbf{I} + \mathbf{A}^T \mathbf{A}]^{-1} \mathbf{A}^T = \mathbf{A}^+.$$

Hinweis: Singulärwertzerlegung von \mathbf{A} anwenden.

7. Für eine (m, n) -Matrix \mathbf{A} mit $m \geq n$ sei

$$\mathbf{M}(\mathbf{A}) = \begin{pmatrix} \mathbf{I} & \mathbf{A} \\ \mathbf{A}^T & \mathbf{O} \end{pmatrix}.$$

- (a) Man zeige, dass $\mathbf{M}(\mathbf{A})$ genau dann regulär ist, wenn $\text{rg}(\mathbf{A}) = n$ gilt.
 (b)

$$\mathbf{M}(\mathbf{A})^+ = \begin{pmatrix} \mathbf{I} - \mathbf{A}\mathbf{A}^+ & \mathbf{A}^{+T} \\ \mathbf{A}^+ & -\mathbf{A}^+ \mathbf{A}^{+T} \end{pmatrix}.$$

Kapitel 11

Freie Minimierung

11.1. Einführung

11.1.1. Aufgabenstellung und grundlegende Begriffe

In diesem Kapitel wollen wir folgendes Problem behandeln:

Es sei eine Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ gegeben. Zu bestimmen ist ein Punkt $\mathbf{x}^* \in \mathbb{R}^n$, der die Funktion f minimiert. Probleme dieser Art heißen freie Minimumprobleme bzw. Minimumprobleme ohne Nebenbedingungen (Restriktionen).

Wir unterscheiden zwischen globalen und lokalen Lösungen.

Ein Punkt $\mathbf{x}^* \in \mathbb{R}^n$, für den es eine Umgebung $U(\mathbf{x}^*)$ gibt, so dass $f(\mathbf{x}^*) \leq f(\mathbf{x})$ für alle $\mathbf{x} \in U(\mathbf{x}^*)$ gilt, heißt **lokale Lösung** des Minimumproblems

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}).$$

Im Falle $U(\mathbf{x}^*) = \mathbb{R}^n$ heißt der Punkt **globale Lösung**. Alle globalen Lösungen eines beliebigen Minimumproblems zu finden, ist i. a. eine schwierige Aufgabe. Wir werden uns damit begnügen, nach lokalen Lösungen oder nach Punkten zu suchen, die als lokale Lösungen in Frage kommen.

Die meisten Verfahren zum Lösen von Minimumproblemen erzeugen ausgehend von einem Startpunkt $\mathbf{x}^{(0)} \in \mathbb{R}^n$ eine Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ deren zugehörige Folge von Funktionswerten $\{f(\mathbf{x}^{(k)})\}_{k \in \mathbb{N}}$ monoton fallend ist.

Als Beispiel stellen wir uns folgende Situation vor. Wir befinden uns in einem Gebirge und werden von einem starken Gewitter überrascht. Natürlich wollen wir schnell talwärts in eine Hütte. Da wir nur die nahe Umgebung überblicken, werden wir eine Richtung wählen, in der es von unserem Standpunkt aus abwärts geht. In dieser Richtung gehen wir ein Stück voran und bestimmen dann wieder eine neue Richtung, usw.

Mathematisch gesprochen:

Wir sind in einem Punkt $\mathbf{x} \in \mathbb{R}^n$ und suchen eine Richtung $\mathbf{p} \in \mathbb{R}^n$, für die es ein $\alpha_0 > 0$ gibt, so dass

$$f(\mathbf{x} + \alpha\mathbf{p}) < f(\mathbf{x})$$

für alle $\alpha \in (0, \alpha_0)$ gilt. Gibt es eine solche **Abstiegsrichtung** \mathbf{p} , so bestimmen wir anschließend eine **Schrittweite** α mit $f(\mathbf{x} + \alpha\mathbf{p}) < f(\mathbf{x})$. Mit dem Punkt $\bar{\mathbf{x}} = \mathbf{x} + \alpha\mathbf{p}$ wird neu gestartet. Sind wir in einem Punkte \mathbf{x}^* angelangt, zu dem es keine Abstiegsrichtung gibt, also

$$f(\mathbf{x}^* + \alpha\mathbf{p}) \geq f(\mathbf{x}^*)$$

für alle Richtungen $\mathbf{p} \in \mathbb{R}^n$ und alle $\alpha \in (0, \alpha_0)$ gilt, so sind wir in einem Tal angelangt. Es ist nicht klar, ob es der tiefste Punkt in unserem Gebirge ist.

Wissen wir dagegen, dass eine Funktion f konvex ist, das heißt

$$f(\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y})$$

für alle $\mathbf{x} \in \mathbb{R}^n$ und alle $\alpha \in [0, 1]$, so ist ein lokales Minimum auch globales Minimum. Das sieht man folgendermaßen ein. Für beliebige $\mathbf{x} \in \mathbb{R}^n$ gilt wegen der Konvexität

$$f(\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*)) \leq (1 - \alpha)f(\mathbf{x}^*) + \alpha f(\mathbf{x}) \quad \forall \alpha \in [0, 1].$$

Daraus folgt

$$f(\mathbf{x}^* + \alpha(\mathbf{x} - \mathbf{x}^*)) - f(\mathbf{x}^*) \leq \alpha(f(\mathbf{x}) - f(\mathbf{x}^*)) \quad \forall \alpha \in [0, 1].$$

Da \mathbf{x}^* lokales Minimum ist, gilt für $0 < \alpha < \min\{\alpha_0, 1\}$

$$0 \leq \alpha(f(\mathbf{x}) - f(\mathbf{x}^*)).$$

Nach Division durch α folgt

$$f(\mathbf{x}^*) \leq f(\mathbf{x}).$$

\mathbf{x}^* ist somit globales Minimum.

Damit haben wir einige grundlegende Begriffe definiert.

11.1.2. Differenzierbarkeit und Richtungsableitung

Ist die Funktion $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ im Punkt \mathbf{x} stetig partiell differenzierbar, das heißt, alle partiellen Ableitungen $\frac{\partial F_i}{\partial x_j}$, $i = 1, \dots, m$ und $j = 1, \dots, n$, existieren im Punkt \mathbf{x} und sind stetig, so heißt \mathbf{F} im Punkt \mathbf{x} **stetig differenzierbar**. Die (m, n) -Matrix

$$\mathbf{F}'(\mathbf{x}) = \left(\frac{\partial F_i(\mathbf{x})}{\partial x_j} \right)_{\substack{i=1, \dots, m \\ j=1, \dots, n}}$$

heißt **Funktionalmatrix** von F in \mathbf{x} . Ist $f : \mathbb{R}^n \rightarrow \mathbb{R}$, so heißt

$$f'(\mathbf{x}) = \nabla f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{pmatrix}$$

Gradient von f in \mathbf{x} .

F heißt stetig differenzierbar auf einer offenen Menge $D \subset \mathbb{R}^n$ ($F \in C^1(D)$), falls F für jedes $\mathbf{x} \in D$ stetig differenzierbar ist.

11.1. Beispiel: Wir betrachten das nichtlineare Ausgleichsproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad f(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^m (F_i(\mathbf{x}))^2.$$

Die Funktionen $F_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, \dots, m$, seien dabei stetig differenzierbar. Es gilt

$$\frac{\partial f(\mathbf{x})}{\partial x_j} = \sum_{i=1}^m \frac{\partial F_i}{\partial x_j}(\mathbf{x}) \cdot F_i(\mathbf{x}), \quad j = 1, \dots, n.$$

Fasst man die F_i als Komponenten einer Funktion

$$F : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

auf, so gilt

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{F}(\mathbf{x})^T \mathbf{F}(\mathbf{x})$$

und

$$\nabla f(\mathbf{x}) = \mathbf{F}'(\mathbf{x})^T \mathbf{F}(\mathbf{x}).$$



11.2. Beispiel: Wir betrachten das nichtlineare Ausgleichsproblem bezüglich einer anderen Norm:

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad f(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|_\infty = \max_{i=1, \dots, m} |F_i(\mathbf{x})|$$

oder

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad f(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|_1 = \sum_{i=1}^m |F_i(\mathbf{x})|.$$

Hier sollte man auch für stetig differenzierbares F nicht mehr erwarten, dass die Funktion f differenzierbar ist.



Für die Fälle aus dem letzten Beispiel benötigen wir einen schwächeren Ableitungsbegriff. Eine Funktion

$$F : \mathbb{R}^n \longrightarrow \mathbb{R}^m$$

heißt im Punkte $\mathbf{x} \in \mathbb{R}^n$ **in Richtung** $\mathbf{p} \in \mathbb{R}^n$ **richtungsdifferenzierbar**, falls der Grenzwert

$$F'(\mathbf{x}; \mathbf{p}) = \lim_{\alpha \rightarrow +0} \frac{F(\mathbf{x} + \alpha \mathbf{p}) - F(\mathbf{x})}{\alpha}$$

existiert. Die Funktion $F'(\mathbf{x}; \mathbf{p})$ heißt dann **Richtungsableitung von F im Punkt \mathbf{x} in Richtung \mathbf{p}** .

Die Funktion F heißt im Punkt $\mathbf{x} \in \mathbb{R}^n$ **richtungsdifferenzierbar**, falls F in \mathbf{x} in jede Richtung $\mathbf{p} \in \mathbb{R}^n$ richtungsdifferenzierbar ist. In diesem Falle nennt man die Abbildung $F'(\mathbf{x}; \circ) : \mathbb{R}^n \longrightarrow \mathbb{R}^m$ **GATEAUX-Ableitung** von F in \mathbf{x} . Für den Zusammenhang der verschiedenen Ableitungsbegriffe gilt der folgende Satz.

11.3. Satz: *Ist die Funktion $F : \mathbb{R}^n \longrightarrow \mathbb{R}^m$ im Punkt $\mathbf{x} \in \mathbb{R}^n$ stetig differenzierbar, so ist F in \mathbf{x} richtungsdifferenzierbar. Die GATEAUX-Ableitung ist durch*

$$F'(\mathbf{x}; \mathbf{p}) = F'(\mathbf{x})\mathbf{p}$$

gegeben.

Im Falle konvexer Funktionen ist jede stationäre Lösung auch globale Lösung des Minimumproblems. Konvexe Funktionen haben noch andere nützliche Eigenschaften.

11.4. Satz: *Es sei $f : \mathbb{R}^n \longrightarrow \mathbb{R}$ auf dem \mathbb{R}^n konvex. Dann ist f in jedem Punkt $\mathbf{x} \in \mathbb{R}^n$ richtungsdifferenzierbar und es gilt*

$$f'(\mathbf{x}; \mathbf{p}) \leq f(\mathbf{x} + \mathbf{p}) - f(\mathbf{x}).$$

Die GATEAUX-Ableitung $f'(\mathbf{x}; \circ) : \mathbb{R}^n \longrightarrow \mathbb{R}$ von f in \mathbf{x} hat folgende Eigenschaften:

1. $f'(\mathbf{x}; \circ)$ ist nichtnegativ homogen:

$$f'(\mathbf{x}; \alpha \mathbf{p}) = \alpha f'(\mathbf{x}; \mathbf{p}) \quad \forall \mathbf{p} \in \mathbb{R}^n \forall \alpha > 0,$$

2. $f'(\mathbf{x}; \circ)$ ist subadditiv:

$$f'(\mathbf{x}; \mathbf{p} + \mathbf{q}) \leq f'(\mathbf{x}; \mathbf{p}) + f'(\mathbf{x}; \mathbf{q}) \quad \forall \mathbf{p}, \mathbf{q} \in \mathbb{R}^n,$$

3. $f'(\mathbf{x}; \circ)$ ist konvex:

$$f'(\mathbf{x}; \alpha \mathbf{p} + (1 - \alpha) \mathbf{q}) \leq \alpha f'(\mathbf{x}; \mathbf{p}) + (1 - \alpha) f'(\mathbf{x}; \mathbf{q}) \quad \forall \mathbf{p}, \mathbf{q} \in \mathbb{R}^n \forall \alpha \in [0, 1].$$

Beweis: Für beliebige aber feste $\mathbf{x}, \mathbf{p} \in \mathbb{R}^n$ definieren wir die Funktion

$$\Phi : (0, 1] \longrightarrow \mathbb{R}$$

durch

$$\Phi(\alpha) = \frac{f(\mathbf{x} + \alpha \mathbf{p}) - f(\mathbf{x})}{\alpha}.$$

Wir zeigen als erstes die Existenz des Grenzwertes

$$\lim_{\alpha \rightarrow +0} \Phi(\alpha).$$

Wegen der Konvexität von f gilt

$$\begin{aligned} f(\mathbf{x}) &= f\left(\frac{1}{1+\alpha}(\mathbf{x} + \alpha \mathbf{p}) + \frac{\alpha}{1+\alpha}(\mathbf{x} - \mathbf{p})\right) \\ &\leq \frac{1}{1+\alpha} f(\mathbf{x} + \alpha \mathbf{p}) + \frac{\alpha}{1+\alpha} f(\mathbf{x} - \mathbf{p}) \end{aligned}$$

für alle $\alpha \in (0, 1]$. Daraus folgt weiter

$$f(\mathbf{x}) - \frac{\alpha}{1+\alpha} f(\mathbf{x} - \mathbf{p}) \leq \frac{1}{1+\alpha} f(\mathbf{x} + \alpha \mathbf{p})$$

und

$$\begin{aligned} (1 + \alpha) f(\mathbf{x}) - \alpha f(\mathbf{x} - \mathbf{p}) &\leq f(\mathbf{x} + \alpha \mathbf{p}), \\ f(\mathbf{x}) - f(\mathbf{x} - \mathbf{p}) &\leq \frac{f(\mathbf{x} + \alpha \mathbf{p}) - f(\mathbf{x})}{\alpha} = \Phi(\alpha). \end{aligned}$$

Damit ist $\Phi(\alpha)$ nach unten beschränkt. Für ein β mit $0 < \beta \leq \alpha \leq 1$ gilt

$$\begin{aligned} \Phi(\beta) &= \frac{f(\mathbf{x} + \beta \mathbf{p}) - f(\mathbf{x})}{\beta} \\ &= \frac{f\left(\frac{\beta}{\alpha}(\mathbf{x} + \alpha \mathbf{p}) + \frac{\alpha - \beta}{\alpha} \mathbf{x}\right) - f(\mathbf{x})}{\beta} \\ &\leq \frac{\frac{\beta}{\alpha} f(\mathbf{x} + \alpha \mathbf{p}) + \frac{\alpha - \beta}{\alpha} f(\mathbf{x}) - f(\mathbf{x})}{\beta} \\ &= \frac{\beta f(\mathbf{x} + \alpha \mathbf{p}) + (\alpha - \beta) f(\mathbf{x}) - \alpha f(\mathbf{x})}{\alpha \beta} \\ &= \Phi(\alpha). \end{aligned}$$

Die Funktion Φ ist somit im Intervall $(0, 1]$ monoton fallend. Mit der oben bewiesenen Beschränktheit folgt die Existenz des Grenzwertes.

Weiter gilt:

$$\begin{aligned} f'(\mathbf{x}; \alpha \mathbf{p}) &= \lim_{\beta \rightarrow +0} \frac{f(\mathbf{x} + \beta \alpha \mathbf{p}) - f(\mathbf{x})}{\beta} \\ &= \lim_{\beta \rightarrow +0} \alpha \frac{f(\mathbf{x} + \beta \alpha \mathbf{p}) - f(\mathbf{x})}{\alpha \beta} \\ &= \alpha \lim_{\gamma \rightarrow +0} \frac{f(\mathbf{x} + \gamma \mathbf{p}) - f(\mathbf{x})}{\gamma} \\ &= \alpha f'(\mathbf{x}; \mathbf{p}), \end{aligned}$$

$$\begin{aligned} f'(\mathbf{x}; \mathbf{p} + \mathbf{q}) &= \lim_{\alpha \rightarrow +0} \frac{f(\mathbf{x} + \alpha(\mathbf{p} + \mathbf{q})) - f(\mathbf{x})}{\alpha} \\ &= \lim_{\alpha \rightarrow +0} \frac{f\left(\frac{1}{2}(\mathbf{x} + 2\alpha \mathbf{p}) + \frac{1}{2}(\mathbf{x} + 2\alpha \mathbf{q})\right) - f(\mathbf{x})}{\alpha} \\ &\leq \lim_{\alpha \rightarrow +0} \frac{\frac{1}{2}f(\mathbf{x} + 2\alpha \mathbf{p}) + \frac{1}{2}f(\mathbf{x} + 2\alpha \mathbf{q}) - f(\mathbf{x})}{\alpha} \\ &= \lim_{\alpha \rightarrow +0} \frac{f(\mathbf{x} + 2\alpha \mathbf{p}) - f(\mathbf{x}) + f(\mathbf{x} + 2\alpha \mathbf{q}) - f(\mathbf{x})}{2\alpha} \\ &= f'(\mathbf{x}; \mathbf{p}) + f'(\mathbf{x}; \mathbf{q}). \end{aligned}$$

Die Konvexität der GATEAUX-Ableitung folgt sofort aus der nichtnegativen Homogenität und der Subadditivität. *

Typische Beispiele konvexer Funktionen sind Vektornormen. Sie besitzen nach Satz 11.4 stets eine GATEAUX-Ableitung. Im Einzelfall ist es oft schwierig, diese zu berechnen. Es gilt

- für $f(\mathbf{x}) = \|\mathbf{x}\|_\infty = \max\{|x_i| : i = 1, \dots, n\}$

$$f'(\mathbf{x}; \mathbf{p}) = \begin{cases} \|\mathbf{p}\|_\infty & \text{für } \mathbf{x} = \mathbf{o} \\ \max\{\text{sign}(x_j)p_j \mid j \in J(\mathbf{x})\} & \text{für } \mathbf{x} \neq \mathbf{o} \end{cases}$$

mit $J(\mathbf{x}) = \{j \in \{1, \dots, n\} \mid |x_j| = \|\mathbf{x}\|_\infty\}$,

- für $f(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$

$$f'(\mathbf{x}; \mathbf{p}) = \sum_{j \in J(\mathbf{x})} |p_j| + \sum_{j \notin J(\mathbf{x})} \text{sign}(x_j) p_j$$

mit $J(\mathbf{x}) = \{j \in \{1, \dots, n\} \mid x_j = 0\}$,

- für $f(\mathbf{x}) = \|\mathbf{x}\|_2^2 = \sum_{i=1}^n x_i^2$

$$f'(\mathbf{x}; \mathbf{p}) = 2 \sum_{i=1}^n x_i p_i.$$

Im letzten Falle ist f sogar stetig differenzierbar und es gilt $f'(\mathbf{x}; \mathbf{p}) = f'(\mathbf{x})^T \mathbf{p}$.
Typische freie Minimumprobleme sind Probleme der Form

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad f(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|,$$

wobei $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ eine im allgemeinen nichtlineare Funktion und $\|\circ\|$ eine Vektornorm auf dem \mathbb{R}^m ist. Aufgaben diesen Typs heißen diskrete, nichtlineare Approximationsaufgaben. Man erhält im einzelnen:

- $\|\circ\| = \|\circ\|_2$: nichtlineares Ausgleichsproblem,
- $\|\circ\| = \|\circ\|_\infty$: diskretes, nichtlineares TSCHEBYSCHJEFF-Problem,
- $\|\circ\| = \|\circ\|_1$: diskrete L_1 -Approximationsaufgabe.

Sind die Komponenten F_i der Funktion \mathbf{F} hinreichend glatt, so darf man annehmen, dass f noch richtungsdifferenzierbar ist. So gilt zum Beispiel für die diskrete TSCHEBYSCHJEFF-Approximation:

11.5. Satz: Mit der Abbildung $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ sei $f : \mathbb{R}^n \rightarrow \mathbb{R}$ durch

$$f(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|_\infty = \max_{i=1, \dots, m} |F_i(\mathbf{x})|$$

definiert. Dann gilt:

Ist \mathbf{F} in $\mathbf{x}^* \in \mathbb{R}^n$ stetig differenzierbar, so ist die Funktion f in \mathbf{x}^* richtungsdifferenzierbar mit der GATEAUX-Ableitung

$$f'(\mathbf{x}^*; \mathbf{p}) = \begin{cases} \max \{|F'_i(\mathbf{x}^*) \mathbf{p}|, i = 1, \dots, m\} & \text{für } \mathbf{F}(\mathbf{x}^*) = \mathbf{o}, \\ \max \{\text{sign}(F_i(\mathbf{x}^*)) F'_i(\mathbf{x}^*) \mathbf{p}, i \in I(\mathbf{x}^*)\} & \text{für } \mathbf{F}(\mathbf{x}^*) \neq \mathbf{o} \end{cases}$$

mit $I(\mathbf{x}^*) = \{i \in \{1, \dots, m\} \mid |F_i(\mathbf{x}^*)| = \|\mathbf{F}(\mathbf{x}^*)\|_\infty\}$.

11.1.3. Optimalitätskriterien

Nun sind wir in der Lage, für das Minimumproblem eine erste notwendige Bedingung zu formulieren.

11.6. Satz: *Es sei $\mathbf{x}^* \in \mathbb{R}^n$ eine lokale Lösung des Minimumproblems*

$$\min \{ f(\mathbf{x}) \mid \mathbf{x} \in \mathbb{R}^n \}.$$

Ist die Zielfunktion f in \mathbf{x}^ richtungsdifferenzierbar, so gilt*

$$f'(\mathbf{x}^*; \mathbf{p}) \geq 0$$

für jede Richtung $\mathbf{p} \in \mathbb{R}^n$. Ist f in \mathbf{x}^ stetig differenzierbar, so gilt $f'(\mathbf{x}^*) = 0$.*

Beweis: Es existiert eine Umgebung $U(\mathbf{x}^*)$ mit $f(\mathbf{x}^*) \leq f(\mathbf{x})$ für alle $\mathbf{x} \in U(\mathbf{x}^*)$. Für ein beliebiges $\mathbf{p} \in \mathbb{R}^n$ existiert ein $\alpha_0 > 0$, so dass $\mathbf{x}^* + \alpha\mathbf{p} \in U(\mathbf{x}^*)$ und damit $f(\mathbf{x}^*) \leq f(\mathbf{x}^* + \alpha\mathbf{p})$ für alle $\alpha \in (0, \alpha_0)$ gilt. Daraus folgt

$$f'(\mathbf{x}^*; \mathbf{p}) = \lim_{\alpha \rightarrow +0} \frac{f(\mathbf{x}^* + \alpha\mathbf{p}) - f(\mathbf{x}^*)}{\alpha} \geq 0.$$

Ist f in \mathbf{x}^* stetig differenzierbar, so gilt $f'(\mathbf{x}^*; \mathbf{p}) = f'(\mathbf{x}^*)^T \mathbf{p} \geq 0$ für alle $\mathbf{p} \in \mathbb{R}^n$. Wir wählen speziell $\mathbf{p} = -f'(\mathbf{x}^*)$ und erhalten

$$-f'(\mathbf{x}^*)^T f'(\mathbf{x}^*) = -\|f'(\mathbf{x}^*)\|^2 \geq 0.$$

Andererseits gilt aber $\|f'(\mathbf{x}^*)\|^2 \geq 0$ und daher $\|f'(\mathbf{x}^*)\|^2 = 0$. *

Punkte, die diese notwendige Bedingung erfüllen, haben eine besondere Bedeutung bei der Konstruktion von Minimierungsverfahren. Ist die Funktion $f: \mathbb{R}^n \rightarrow \mathbb{R}$ im Punkt $\mathbf{x}^* \in \mathbb{R}^n$ richtungsdifferenzierbar, so heißt ein Punkt \mathbf{x}^* **stationär** von f oder eine **stationäre Lösung** des Minimumproblems $\min \{ f(\mathbf{x}) \mid \mathbf{x} \in \mathbb{R}^n \}$, falls $f'(\mathbf{x}^*; \mathbf{p}) \geq 0$ für alle Richtungen $\mathbf{p} \in \mathbb{R}^n$ gilt. Mit den Sätzen 11.4 und 11.5 erhält man für das diskrete TSCHEBYSCHEFF-Problem ein Optimalitätskriterium erster Ordnung¹.

11.7. Satz: *Gegeben sei die diskrete TSCHEBYSCHEFFsche Approximationsaufgabe*

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad f(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|_\infty = \max_{i=1, \dots, m} |F_i(\mathbf{x})|.$$

Die Funktionen $F_i: \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, \dots, m$, seien in \mathbf{x}^ stetig differenzierbar und es sei $\mathbf{F}(\mathbf{x}^*) \neq \mathbf{o}$. Dann gilt:*

Ist \mathbf{x}^ eine lokale Lösung des Minimumproblems, so ist \mathbf{x}^* auch eine stationäre Lösung und daher*

¹Optimalitätskriterium erster Ordnung bedeutet, dass nur Ableitungen erster Ordnung verwendet werden

1. $f'(\mathbf{x}^*; \mathbf{p}) = \max\{\text{sign}(F_i(\mathbf{x}^*))F_i'(\mathbf{x}^*)\mathbf{p}, i \in I(\mathbf{x}^*)\} \geq 0 \forall \mathbf{p} \in \mathbb{R}^n$ mit
 $I(\mathbf{x}^*) = \{i \in \{1, \dots, m\} \mid |F_i(\mathbf{x}^*)| = \|\mathbf{F}(\mathbf{x}^*)\|_\infty\}$.
2. Die erste Bedingung ist äquivalent zur Existenz von nichtnegativen reellen Zahlen λ_i^* , $i \in I(\mathbf{x}^*)$, mit

$$\sum_{i \in I(\mathbf{x}^*)} \lambda_i^* = 1, \quad \sum_{i \in I(\mathbf{x}^*)} \lambda_i^* \text{sign}(F_i(\mathbf{x}^*))F_i'(\mathbf{x}^*) = 0.$$

Im stetig differenzierbaren Falle werden die notwendigen Optimalitätskriterien natürlich einfacher. Für das nichtlineare Ausgleichsproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad f(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|_2^2 = \sum_{i=1}^m (F_i(\mathbf{x}))^2$$

mit einer stetig differenzierbaren Funktion $\mathbf{F} : \mathbb{R}^n \rightarrow \mathbb{R}^m$ ergibt sich

$$f'(\mathbf{x}^*) = 2\mathbf{F}'(\mathbf{x}^*)\mathbf{F}(\mathbf{x}^*) = 0$$

als notwendige Bedingung für ein lokales Minimum.

Neben den notwendigen Optimalitätskriterien erster Ordnung interessieren auch notwendige und hinreichende Optimalitätskriterien zweiter Ordnung. Dazu definieren wir:

Ist die Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ in einem Punkt $\mathbf{x} \in \mathbb{R}^n$ zweimal stetig partiell differenzierbar (alle partiellen Ableitungen zweiter Ordnung existieren in einer Umgebung von \mathbf{x} und diese sind in \mathbf{x} stetig), so heißt f in \mathbf{x} zweimal stetig differenzierbar. Die symmetrische (n, n) -Matrix

$$f''(\mathbf{x}) = \left(\frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x}) \right)_{1 \leq i, j \leq n}$$

heißt **HESSE-Matrix** von f im Punkt \mathbf{x} .

Eine Funktion f heißt zweimal stetig differenzierbar auf einer offenen Menge $D \in \mathbb{R}^n$, falls f für jedes $\mathbf{x} \in D$ zweimal stetig differenzierbar ist. Für zweimal stetig differenzierbare Funktionen kennen wir aus der Analysis notwendige und hinreichende Optimalitätskriterien.

11.8. Satz: Die Funktion f sei auf einer offenen Umgebung des Punktes $\mathbf{x}^* \in \mathbb{R}^n$ zweimal stetig differenzierbar. Dann gilt:

1. Ist \mathbf{x}^* eine lokale Lösung des Minimumproblems

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}),$$

so ist $f'(\mathbf{x}^*) = \mathbf{0}$ und $f''(\mathbf{x}^*)$ ist positiv semidefinit.

2. Ist $f'(\mathbf{x}^*) = \mathbf{o}$ und $f''(\mathbf{x}^*)$ positiv definit, so ist \mathbf{x}^* eine isolierte lokale Lösung des Minimumproblems: Es gibt eine Umgebung $U(\mathbf{x}^*)$ mit

$$f(\mathbf{x}^*) < f(\mathbf{x}) \quad \forall \mathbf{x} \in U(\mathbf{x}^*).$$

Im nichtglatten Falle ist die Situation wieder komplizierter. Es lassen sich aber für spezielle Probleme auch hinreichende Bedingungen formulieren.

11.9. Satz: Gegeben sei die diskrete TSCHEBYSCHEFFSche Approximationsaufgabe

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad f(\mathbf{x}) = \|\mathbf{F}(\mathbf{x})\|_\infty = \max_{i=1, \dots, m} |F_i(\mathbf{x})|.$$

Die Funktionen

$$F_i : \mathbb{R}^n \longrightarrow \mathbb{R}, \quad i = 1, \dots, m,$$

seien auf einer offenen Umgebung von $\mathbf{x}^* \in \mathbb{R}^n$ zweimal stetig differenzierbar und es sei $\mathbf{F}(\mathbf{x}^*) \neq \mathbf{o}$. Wir definieren die Menge

$$I(\mathbf{x}^*) = \{i \in \{1, \dots, m\} \mid |F_i(\mathbf{x}^*)| = \|\mathbf{F}(\mathbf{x}^*)\|_\infty\}$$

und nehmen an, dass nichtnegative reelle Zahlen λ_i^* , $i \in I(\mathbf{x}^*)$, existieren mit

$$\sum_{i \in I(\mathbf{x}^*)} \lambda_i^* = 1, \quad \sum_{i \in I(\mathbf{x}^*)} \lambda_i^* \text{sign}(F_i(\mathbf{x}^*)) F_i'(\mathbf{x}^*) = \mathbf{o}.$$

Weiterhin sei

$$T^* = \{ \mathbf{p} \in \mathbb{R}^n \mid F_i'(\mathbf{x}^*) \mathbf{p} = 0 \} \quad \forall i \in I(\mathbf{x}^*) \wedge \lambda_i^* > 0.$$

Gilt nun

$$\mathbf{p}^T \left\{ \sum_{i \in I(\mathbf{x}^*)} \lambda_i^* \text{sign}(F_i(\mathbf{x}^*)) F_i''(\mathbf{x}^*) \right\} \mathbf{p} > 0 \quad \forall \mathbf{p} \in T^* \setminus \{ \mathbf{o} \},$$

so ist \mathbf{x}^* eine isolierte Lösung des Minimumproblems.

Beweis: Wir nehmen an, die Behauptung sei falsch. Es sei \mathbf{x}^* keine isolierte lokale Lösung. Dann gibt es eine gegen \mathbf{x}^* konvergierende Folge $\{\mathbf{x}_k\}_{k \in \mathbb{N}}$ mit $\mathbf{x}_k \neq \mathbf{x}^*$ für alle k , für die $f(\mathbf{x}_k) \leq f(\mathbf{x}^*)$ gilt. Wir stellen die \mathbf{x}_k in der Form $\mathbf{x}_k = \mathbf{x}^* + \alpha_k \mathbf{p}_k$ mit $\alpha_k > 0$ und $\|\mathbf{p}_k\| = 1$ dar. Wegen

$$\lim_{k \rightarrow \infty} \mathbf{x}_k = \mathbf{x}^* \quad \text{gilt} \quad \lim_{k \rightarrow \infty} \alpha_k = 0.$$

Aus $\{\mathbf{p}_k\}_{k \in \mathbb{N}}$ lässt sich eine konvergente Teilfolge auswählen. O.B.d.A. sei darum

$$\mathbf{x}_k = \mathbf{x}^* + \alpha_k \mathbf{p}_k, \quad \forall k : \alpha_k \geq 0, \quad \lim_{k \rightarrow \infty} \mathbf{p}_k = \mathbf{p} \neq \mathbf{o}.$$

Es sei weiterhin $\mathbf{r}_k = \alpha_k(\mathbf{p}_k - \mathbf{p})$. Dann gilt

$$\lim_{k \rightarrow \infty} \frac{\mathbf{r}_k}{\alpha_k} = \mathbf{o}.$$

Damit folgt

$$\begin{aligned} \lim_{k \rightarrow \infty} \frac{f(\mathbf{x}^* + \alpha_k \mathbf{p}_k) - f(\mathbf{x}^*)}{\alpha_k} &= \\ &= \lim_{k \rightarrow \infty} \frac{f(\mathbf{x}^* + \alpha_k \mathbf{p} + \mathbf{r}_k) - f(\mathbf{x}^*)}{\alpha_k} \\ &= \lim_{k \rightarrow \infty} \frac{f(\mathbf{x}^* + \alpha_k \mathbf{p}) - f(\mathbf{x}^*) + f(\mathbf{x}^* + \alpha_k \mathbf{p} + \mathbf{r}_k) - f(\mathbf{x}^* + \alpha_k \mathbf{p})}{\alpha_k} \\ &= \lim_{k \rightarrow \infty} \frac{f(\mathbf{x}^* + \alpha_k \mathbf{p}) - f(\mathbf{x}^*)}{\alpha_k} + \lim_{k \rightarrow \infty} \frac{f(\mathbf{x}^* + \alpha_k \mathbf{p} + \mathbf{r}_k) - f(\mathbf{x}^* + \alpha_k \mathbf{p})}{\alpha_k} \\ &= f'(\mathbf{x}^*, \mathbf{p}) + \lim_{k \rightarrow \infty} \frac{\|\mathbf{F}(\mathbf{x}^* + \alpha_k \mathbf{p} + \mathbf{r}_k)\|_\infty - \|\mathbf{F}(\mathbf{x}^* + \alpha_k \mathbf{p})\|_\infty}{\alpha_k}. \end{aligned}$$

Nun gilt

$$\begin{aligned} 0 &\leq \lim_{k \rightarrow \infty} \frac{\|\|\mathbf{F}(\mathbf{x}^* + \alpha_k \mathbf{p} + \mathbf{r}_k)\|_\infty - \|\mathbf{F}(\mathbf{x}^* + \alpha_k \mathbf{p})\|_\infty\|}{\alpha_k} \\ &\leq \lim_{k \rightarrow \infty} \frac{\|\mathbf{F}(\mathbf{x}^* + \alpha_k \mathbf{p} + \mathbf{r}_k) - \mathbf{F}(\mathbf{x}^* + \alpha_k \mathbf{p})\|_\infty}{\alpha_k} \\ &\leq \lim_{k \rightarrow \infty} \frac{\|\mathbf{F}'(\mathbf{x}^* + \alpha_k \mathbf{p}) + \mathbf{F}'(\mathbf{x}^* + \alpha_k \mathbf{p} + \vartheta \mathbf{r}_k) \mathbf{r}_k - \mathbf{F}(\mathbf{x}^* + \alpha_k \mathbf{p})\|_\infty}{\alpha_k} \\ &\leq \lim_{k \rightarrow \infty} \frac{\|\mathbf{F}'(\mathbf{x}^* + \alpha_k \mathbf{p} + \vartheta \mathbf{r}_k)\|_\infty \|\mathbf{r}_k\|_\infty}{\alpha_k} \quad (\vartheta_k \in (0, 1)) \\ &= \lim_{k \rightarrow \infty} \|\mathbf{F}'(\mathbf{x}^* + \alpha_k \mathbf{p} + \vartheta \mathbf{r}_k)\|_\infty \lim_{k \rightarrow \infty} \frac{\|\mathbf{r}_k\|_\infty}{\alpha_k} \\ &= 0. \end{aligned}$$

Insgesamt erhalten wir

$$\lim_{k \rightarrow \infty} \frac{f(\mathbf{x}^* + \alpha_k \mathbf{p}_k) - f(\mathbf{x}^*)}{\alpha_k} = f'(\mathbf{x}^*, \mathbf{p}).$$

Wegen $f(\mathbf{x}^* + \alpha_k \mathbf{p}_k) \leq f(\mathbf{x}^*)$ und $\alpha_k > 0$ folgt $f'(\mathbf{x}^*, \mathbf{p}) \leq \mathbf{o}$. Mit Satz 11.5 ergibt sich

$$f'(\mathbf{x}^*, \mathbf{p}) = \max_{i \in I(\mathbf{x}^*)} \text{sign}(F_i(\mathbf{x}^*)) F_i'(\mathbf{x}^*) \mathbf{p} \leq \mathbf{o}.$$

Nach der ersten Voraussetzung gilt

$$\sum_{i \in I(\mathbf{x}^*)} \lambda_i^* \text{sign}(F_i(\mathbf{x}^*)) F_i'(\mathbf{x}^*) \mathbf{p} = \mathbf{o}.$$

Wegen $\lambda_i^* \geq 0$ und $\text{sign}(F_i(\mathbf{x}^*)) F_i'(\mathbf{x}^*) \mathbf{p} \leq \mathbf{o}$ für $i \in I(\mathbf{x}^*)$ folgt daraus $F_i'(\mathbf{x}^*) \mathbf{p} = \mathbf{o}$ für $i \in I(\mathbf{x}^*)$. Es gilt daher $\mathbf{p} \in T^*$. Für $i \in I(\mathbf{x}^*)$ ist weiterhin

$$|F_i(\mathbf{x}_k)| \leq \|\mathbf{F}(\mathbf{x}_k)\|_\infty \leq \|\mathbf{F}(\mathbf{x}^*)\|_\infty = |F_i(\mathbf{x}^*)|,$$

und für hinreichend großes k ist $\text{sign}(F_i(\mathbf{x}_k)) = \text{sign}(F_i(\mathbf{x}^*))$. Damit folgt für alle $i \in I(\mathbf{x}^*)$ und für hinreichend großes k

$$\begin{aligned} 0 &\geq |F_i(\mathbf{x}_k)| - |F_i(\mathbf{x}^*)| \\ &= \text{sign}(F_i(\mathbf{x}^*)) [F_i(\mathbf{x}_k) - F_i(\mathbf{x}^*)] \\ &= \text{sign}(F_i(\mathbf{x}^*)) [F_i(\mathbf{x}^* + \alpha_k \mathbf{p}_k) - F_i(\mathbf{x}^*)] \\ &\geq \alpha_k \text{sign}(F_i(\mathbf{x}^*)) F_i'(\mathbf{x}^*) \mathbf{p}_k + \frac{1}{2} \alpha_k^2 \text{sign}(F_i(\mathbf{x}^*)) \mathbf{p}_k^T F_i''(\mathbf{x}^* + \vartheta_{ik} \alpha_k \mathbf{p}_k) \mathbf{p}_k \end{aligned}$$

mit $\vartheta_{ik} \in (0, 1)$. Multiplizieren wir die letzte Ungleichung mit $\lambda_i^* \neq 0$ und summieren über alle $i \in I(\mathbf{x}^*)$, so erhalten wir

$$\begin{aligned} 0 &\geq \alpha_k \sum_{i \in I(\mathbf{x}^*)} \lambda_i^* \text{sign}(F_i(\mathbf{x}^*)) F_i'(\mathbf{x}^*) \mathbf{p}_k \\ &\quad + \frac{\alpha_k^2}{2} \sum_{i \in I(\mathbf{x}^*)} \lambda_i^* \text{sign}(F_i(\mathbf{x}^*)) \mathbf{p}_k^T F_i''(\mathbf{x}^* + \vartheta_{ik} \alpha_k \mathbf{p}_k) \mathbf{p}_k \\ &= \frac{\alpha_k^2}{2} \mathbf{p}_k^T \left\{ \sum_{i \in I(\mathbf{x}^*)} \lambda_i^* \text{sign}(F_i(\mathbf{x}^*)) F_i''(\mathbf{x}^* + \vartheta_{ik} \alpha_k \mathbf{p}_k) \right\} \mathbf{p}_k. \end{aligned}$$

Der Grenzübergang $k \rightarrow \infty$ ($\mathbf{p}_k \rightarrow \mathbf{p}$ und $\vartheta_{ik} \rightarrow 0$) liefert

$$0 \geq \mathbf{p} \left\{ \sum_{i \in I(\mathbf{x}^*)} \lambda_i^* \text{sign}(F_i(\mathbf{x}^*)) F_i''(\mathbf{x}^*) \right\} \mathbf{p}$$

als Widerspruch zur Voraussetzung

$$0 < p \left\{ \sum_{i \in I(\mathbf{x}^*)} \lambda_i^* \operatorname{sign}(F_i(\mathbf{x}^*)) F_i''(\mathbf{x}^*) \right\} p.$$

Damit muss die Annahme falsch sein; \mathbf{x}^* ist isolierte lokale Lösung. *

Für konvexe Funktionen lassen sich allgemeine globale Konvergenzaussagen machen. Für hinreichend glatte Funktionen wollen wir die Konvexität durch geeignete Bedingungen charakterisieren.

Eine Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ heißt **gleichmäßig konvex** auf der konvexen Menge $D \subseteq \mathbb{R}^n$, falls eine Konstante $\gamma > 0$ existiert, so dass

$$(1 - \alpha)f(\mathbf{x}) + \alpha f(\mathbf{y}) - f((1 - \alpha)\mathbf{x} + \alpha\mathbf{y}) \geq \frac{\gamma}{2} \alpha(1 - \alpha) \|\mathbf{x} - \mathbf{y}\|_2^2$$

für alle $\mathbf{x}, \mathbf{y} \in D$ und alle $\alpha \in [0, 1]$ gilt.

Bemerkung: Eine Funktion f ist konvex auf D , falls die obige Bedingung mit $\gamma = 0$ gilt.

11.10. Beispiel: Wir betrachten für die (n, n) -symmetrische Matrix Q die quadratische Funktion

$$f(\mathbf{x}) = \mathbf{c}^T \mathbf{x} + \frac{1}{2} \mathbf{x}^T Q \mathbf{x}$$

für beliebige $\mathbf{x} \in \mathbb{R}^n$. Es gilt

$$\begin{aligned} (1 - \alpha)f(\mathbf{x}) + \alpha f(\mathbf{y}) - f((1 - \alpha)\mathbf{x} + \alpha\mathbf{y}) &= \\ &= (1 - \alpha)\mathbf{c}^T \mathbf{x} + \frac{1 - \alpha}{2} \mathbf{x}^T Q \mathbf{x} + \alpha \mathbf{c}^T \mathbf{y} + \frac{\alpha}{2} \mathbf{y}^T Q \mathbf{y} - \\ &\quad - \mathbf{c}^T [(1 - \alpha)\mathbf{x} + \alpha\mathbf{y}] - \frac{1}{2} [(1 - \alpha)\mathbf{x} + \alpha\mathbf{y}]^T Q [(1 - \alpha)\mathbf{x} + \alpha\mathbf{y}] \\ &= \frac{1 - \alpha}{2} \mathbf{x}^T Q \mathbf{x} + \frac{\alpha}{2} \mathbf{y}^T Q \mathbf{y} - \frac{(1 - \alpha)^2}{2} \mathbf{x}^T Q \mathbf{x} - \\ &\quad - \frac{\alpha(1 - \alpha)}{2} \mathbf{x}^T Q \mathbf{y} - \frac{\alpha(1 - \alpha)}{2} \mathbf{y}^T Q \mathbf{x} - \frac{\alpha^2}{2} \mathbf{y}^T Q \mathbf{y} \\ &= \frac{\alpha(1 - \alpha)}{2} \left[\mathbf{x}^T Q \mathbf{x} - \mathbf{x}^T Q \mathbf{y} - \mathbf{y}^T Q \mathbf{x} + \mathbf{y}^T Q \mathbf{y} \right] \\ &= \frac{\alpha(1 - \alpha)}{2} (\mathbf{x} - \mathbf{y})^T Q (\mathbf{x} - \mathbf{y}). \end{aligned}$$

Ist \mathbf{Q} positiv semidefinit, so gilt

$$(\mathbf{x} - \mathbf{y})^T \mathbf{Q}(\mathbf{x} - \mathbf{y}) \geq 0$$

für beliebige \mathbf{x} und \mathbf{y} . f ist in diesem Falle konvex. Ist \mathbf{Q} sogar positiv definit, so gilt

$$(\mathbf{x} - \mathbf{y})^T \mathbf{Q}(\mathbf{x} - \mathbf{y}) \geq \lambda_{\min} \|\mathbf{x} - \mathbf{y}\|_2^2$$

mit $\lambda_{\min} > 0$ als kleinstem Eigenwert von \mathbf{Q} . In diesem Falle ist f gleichmäßig konvex. ♡

Mit Hilfe von ersten und zweiten Ableitungen lassen sich Bedingungen für die Konvexität und gleichmäßige Konvexität angeben.

11.11. Satz: Die Menge $D \subseteq \mathbb{R}^n$ sei konvex und die Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ sei auf einer offenen Obermenge von D stetig differenzierbar. Dann gilt

1. f ist genau dann auf D konvex, wenn

$$f'(\mathbf{x})(\mathbf{y} - \mathbf{x}) \leq f(\mathbf{y}) - f(\mathbf{x})$$

für alle $\mathbf{x}, \mathbf{y} \in D$ gilt.

2. f ist genau dann auf D gleichmäßig konvex mit einer Konstanten $\gamma > 0$, wenn

$$\frac{\gamma}{2} \|\mathbf{x} - \mathbf{y}\|_2^2 + f'(\mathbf{x})(\mathbf{y} - \mathbf{x}) \leq f(\mathbf{y}) - f(\mathbf{x})$$

für alle $\mathbf{x}, \mathbf{y} \in D$ gilt.

Beweis: Wir nehmen zunächst an, die Funktion f sei konvex bzw. gleichmäßig konvex. Dann gilt

$$(1 - \alpha)f(\mathbf{x}) + \alpha f(\mathbf{y}) - f((1 - \alpha)\mathbf{x} + \alpha\mathbf{y}) \geq \frac{\gamma}{2} \alpha(1 - \alpha) \|\mathbf{x} - \mathbf{y}\|_2^2$$

für alle $\mathbf{x}, \mathbf{y} \in D$ und alle $\alpha \in [0, 1]$ mit $\gamma = 0$ für konvexes f bzw. $\gamma > 0$ für gleichmäßig konvexes f . Es folgt

$$\begin{aligned} \alpha(f(\mathbf{y}) - f(\mathbf{x})) &\geq f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})) - f(\mathbf{x}) + \frac{\gamma}{2} \alpha(1 - \alpha) \|\mathbf{x} - \mathbf{y}\|_2^2, \\ f(\mathbf{y}) - f(\mathbf{x}) &\geq \frac{f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})) - f(\mathbf{x})}{\alpha} + \frac{\gamma}{2} (1 - \alpha) \|\mathbf{x} - \mathbf{y}\|_2^2. \end{aligned}$$

Für $\alpha \rightarrow +0$ folgt daraus

$$f(\mathbf{y}) - f(\mathbf{x}) \geq f'(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{\gamma}{2} \|\mathbf{x} - \mathbf{y}\|_2^2.$$

Nehmen wir nun umgekehrt an, dass eine Konstante $\gamma \geq 0$ existiert, so dass für alle $\mathbf{x}, \mathbf{y} \in D$

$$f(\mathbf{y}) - f(\mathbf{x}) \geq f'(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{\gamma}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$

gilt. Für ein $\alpha \in [0, 1]$ ist $\mathbf{z} = (1 - \alpha)\mathbf{x} + \alpha\mathbf{y} \in D$. Dann gilt

$$\begin{aligned} f(\mathbf{x}) - f(\mathbf{z}) &\geq f'(\mathbf{z})(\mathbf{x} - \mathbf{z}) + \frac{\gamma}{2} \|\mathbf{x} - \mathbf{z}\|_2^2, \\ f(\mathbf{y}) - f(\mathbf{z}) &\geq f'(\mathbf{z})(\mathbf{y} - \mathbf{z}) + \frac{\gamma}{2} \|\mathbf{y} - \mathbf{z}\|_2^2. \end{aligned}$$

Wir multiplizieren die erste Ungleichung mit $1 - \alpha$ und die zweite Ungleichung mit α und addieren beide. Es ergibt sich

$$\begin{aligned} (1 - \alpha)f(\mathbf{x}) + \alpha f(\mathbf{y}) - f(\mathbf{z}) &\geq f'(\mathbf{z})[(1 - \alpha)\mathbf{x} + \alpha\mathbf{y} - \mathbf{z}] + \\ &\quad \frac{\gamma}{2} \left[(1 - \alpha) \|\mathbf{x} - \mathbf{z}\|_2^2 + \alpha \|\mathbf{y} - \mathbf{z}\|_2^2 \right]. \end{aligned}$$

Mit $\mathbf{z} = (1 - \alpha)\mathbf{x} + \alpha\mathbf{y}$ folgt daraus

$$\begin{aligned} (1 - \alpha)f(\mathbf{x}) + \alpha f(\mathbf{y}) - f((1 - \alpha)\mathbf{x} + \alpha\mathbf{y}) &\geq \\ &\geq \frac{\gamma}{2} \left\{ (1 - \alpha)(\mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{z} + \mathbf{z}^T \mathbf{z}) + \alpha(\mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T \mathbf{z} + \mathbf{z}^T \mathbf{z}) \right\} \\ &= \frac{\gamma}{2} \left\{ (1 - \alpha)\mathbf{x}^T \mathbf{x} - 2[(1 - \alpha)\mathbf{x} + \alpha\mathbf{y}]^T \mathbf{z} + (1 - \alpha)\mathbf{z}^T \mathbf{z} + \alpha\mathbf{y}^T \mathbf{y} + \alpha\mathbf{z}^T \mathbf{z} \right\} \\ &= \frac{\gamma}{2} \left\{ (1 - \alpha)\mathbf{x}^T \mathbf{x} + \alpha\mathbf{y}^T \mathbf{y} - \mathbf{z}^T \mathbf{z} \right\} \\ &= \frac{\gamma}{2} \left\{ (1 - \alpha)\mathbf{x}^T \mathbf{x} + \alpha\mathbf{y}^T \mathbf{y} - (1 - \alpha)^2 \mathbf{x}^T \mathbf{x} - 2(1 - \alpha)\alpha\mathbf{y}^T \mathbf{x} - \alpha^2 \mathbf{y}^T \mathbf{y} \right\} \\ &= \frac{\gamma}{2} \alpha(1 - \alpha) \|\mathbf{x} - \mathbf{y}\|_2^2. \end{aligned}$$

Für $\gamma = 0$ ist f damit konvex und für $\gamma > 0$ ist f gleichmäßig konvex. *

11.12. Satz: Es sei $D \subseteq \mathbb{R}^n$ konvex und die Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ sei auf einer offenen Obermenge von D zweimal stetig differenzierbar. Dann gilt:

1. Ist $f''(\mathbf{x})$ positiv definit für alle $\mathbf{x} \in D$, so ist f auf D konvex.

2. Existiert eine Konstante $\gamma > 0$, so dass für alle $\mathbf{p} \in \mathbb{R}^n$ und alle $\mathbf{x} \in D$

$$\gamma \|\mathbf{p}\|_2^2 \leq \mathbf{p}^T f''(\mathbf{x}) \mathbf{p}$$

gilt, so ist f auf D gleichmäßig konvex.

3. Ist D offen, so gelten in 1 und 2 die Umkehrungen.

Beweis: Für beliebige, aber feste $\mathbf{x}, \mathbf{y} \in D$ definieren wir die zweimal stetig differenzierbare Funktion

$$\Phi : [0, 1] \longrightarrow \mathbb{R}, \quad \Phi(\alpha) = f(\mathbf{x} + \alpha(\mathbf{y} - \mathbf{x})).$$

Es gilt

$$\Phi(1) - \Phi(0) - \Phi'(0) = \frac{1}{2} \Phi''(\alpha_0) \text{ mit } \alpha_0 \in (0, 1).$$

Damit folgt

$$f(\mathbf{y}) - f(\mathbf{x}) - f'(\mathbf{x})(\mathbf{y} - \mathbf{x}) = \frac{1}{2} (\mathbf{y} - \mathbf{x})^T f''(\mathbf{x} + \alpha_0(\mathbf{y} - \mathbf{x})) (\mathbf{y} - \mathbf{x}).$$

Wegen $\mathbf{x} + \alpha_0(\mathbf{y} - \mathbf{x}) \in D$ gilt $(\mathbf{y} - \mathbf{x})^T f''(\mathbf{x} + \alpha_0(\mathbf{y} - \mathbf{x})) (\mathbf{y} - \mathbf{x}) \geq \gamma \|\mathbf{y} - \mathbf{x}\|_2^2$.
Damit folgt

$$f(\mathbf{y}) - f(\mathbf{x}) - f'(\mathbf{x})(\mathbf{y} - \mathbf{x}) \geq \frac{\gamma}{2} \|\mathbf{y} - \mathbf{x}\|_2^2.$$

Nach Satz 11.11 ist f damit konvex und für positiv definites f'' ($\gamma > 0$) gleichmäßig konvex.

Es sei nun D offen und f auf D konvex ($\gamma = 0$) bzw. gleichmäßig konvex ($\gamma > 0$). Für beliebiges $\mathbf{x} \in D$ und eine beliebige Richtung $\mathbf{p} \in \mathbb{R}^n$ existiert ein $\alpha_0 > 0$, so dass $\mathbf{x} + \alpha \mathbf{p} \in D$ für alle $\alpha \in [0, \alpha_0]$. Nach Satz 11.11 gilt dann

$$f(\mathbf{x} + \alpha \mathbf{p}) - f(\mathbf{x}) \geq f'(\mathbf{x})(\alpha \mathbf{p}) + \frac{\gamma}{2} \alpha^2 \|\mathbf{p}\|_2^2$$

und

$$f(\mathbf{x}) - f(\mathbf{x} + \alpha \mathbf{p}) \geq -f'(\mathbf{x} + \alpha \mathbf{p})(\alpha \mathbf{p}) + \frac{\gamma}{2} \alpha^2 \|\mathbf{p}\|_2^2.$$

Addition der beiden Ungleichungen liefert

$$[f'(\mathbf{x} + \alpha \mathbf{p}) - f'(\mathbf{x})](\alpha \mathbf{p}) \geq \frac{\gamma}{2} \alpha^2 \|\mathbf{p}\|_2^2.$$

Daraus erhalten wir

$$\begin{aligned} \mathbf{p}^T f''(\mathbf{x})\mathbf{p} &= \lim_{\alpha \rightarrow +0} \frac{[f'(\mathbf{x} + \alpha\mathbf{p}) - f'(\mathbf{x})]\mathbf{p}}{\alpha} \\ &= \lim_{\alpha \rightarrow +0} \frac{[f'(\mathbf{x} + \alpha\mathbf{p}) - f'(\mathbf{x})]\alpha\mathbf{p}}{\alpha^2} \\ &\geq \frac{\gamma}{2}\alpha^2\|\mathbf{p}\|_2^2. \end{aligned}$$

✱

Für konvexe Funktionen lassen sich so globale Konvergenzaussagen machen. Ist andererseits für eine lokale Lösung \mathbf{x}^* die hinreichende Bedingung zweiter Ordnung erfüllt, also f in einer Umgebung $U(\mathbf{x}^*)$ zweimal stetig differenzierbar und $f''(\mathbf{x}^*)$ positiv definit, so existiert eine Umgebung von \mathbf{x}^* , in der f'' positiv definit und damit f gleichmäßig konvex ist. Die Konvexitätseigenschaften der Funktion f spielen damit auch für lokale Konvergenzaussagen und die Abschätzung von Konvergenzgeschwindigkeiten eine große Rolle.

11.2. Ein Modellalgorithmus

Die Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ sei in jedem Punkte $\mathbf{x} \in \mathbb{R}^n$ richtungsdifferenzierbar. Zum Lösen des Problems

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

wollen wir einen Algorithmus folgender Form anwenden.

11.13. Modellalgorithmus zum Lösen freier Minimumprobleme:

(Initialisierung) Wähle einen Startpunkt $\mathbf{x}^{(0)} \in \mathbb{R}^n$ und setze $k = 0$.

S0 (Abbruchkriterium)

Ist $f'(\mathbf{x}^{(k)}; \mathbf{p}) \geq 0$ für alle Richtungen $\mathbf{p} \in \mathbb{R}^n$, dann STOPP; $\mathbf{x}^{(k)}$ ist stationäre Lösung.

S1 (Richtungswahl) Wähle Abstiegsrichtung $\mathbf{p}^{(k)}$ mit

$$f'(\mathbf{x}^{(k)}; \mathbf{p}^{(k)}) < 0.$$

S2 (Schrittweitenbestimmung) Bestimme Schrittweite α_k , so dass

$$f(\mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}) < f(\mathbf{x}^{(k)}).$$

S3 (*Iterationsschritt*) Setze

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}, \quad k = k + 1$$

und gehe zu Schritt **S1**.

Legt man nun Strategien fest, nach denen Abstiegsrichtung und Schrittweite zu bestimmen sind, erhält man konkrete Verfahren.

11.2.1. Schrittweiten bei glatter Zielfunktion

Wir betrachten wieder das freie Minimumproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}).$$

Die Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ erfülle folgende Voraussetzungen:

V1 Für ein gegebenes $\mathbf{x}^{(0)} \in \mathbb{R}^n$ (i.a. der Startpunkt des Verfahrens) ist die Niveaumenge

$$L_0 = \left\{ \mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \leq f(\mathbf{x}^{(0)}) \right\} \text{ beschränkt.}$$

V2 Die Zielfunktion f ist auf einer offenen Menge $D \supseteq L_0$ stetig differenzierbar.

V3 Der Gradient $\nabla f(\mathbf{x}) = f'(\mathbf{x})^T$ ist auf L_0 lipschitzstetig: Es existiert eine Konstante $L > 0$ mit

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

für alle $\mathbf{x}, \mathbf{y} \in L_0$.

Bei vorgegebener Abstiegsrichtung $\mathbf{p} \in \mathbb{R}^n$ ist eine solche Schrittweite $\alpha > 0$ zu bestimmen, die einen hinreichend starken Abstieg sichert. Die Forderung $f(\mathbf{x}^{(k+1)}) < f(\mathbf{x}^{(k)})$ ist nicht hinreichend dafür, dass die Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ gegen eine lokale Lösung des Minimumproblems konvergiert.

11.14. Beispiel: Es sei $f(x) = x^2$ und $x_k = (-1)^k(\frac{1}{2} + 2^{-k})$. Dann gilt

$$f(x_{k+1}) = x_{k+1}^2 = \frac{1}{4} + 2^{-k-1} + 2^{-2k-2} < \frac{1}{4} + 2^{-k} + 2^{-2k} = x_k^2 = f(x_k).$$

Die Folge $\{x_k\}_{k \in \mathbb{N}}$ ist daher streng monoton fallend. Andererseits besitzt sie nur die beiden Häufungspunkte $-\frac{1}{2}$ und $\frac{1}{2}$, die beide keine lokalen Minima von f sind. \heartsuit

Das folgende Lemma benötigen wir, um für konkrete Verfahren den erzielten Abstieg $f(\mathbf{x}) - f(\mathbf{x} + \alpha\mathbf{p})$ nach unten abzuschätzen.

11.15. Satz: Die Zielfunktion f genüge den Voraussetzungen **V1**, **V2** und **V3**. Es seien $\mathbf{x} \in L_0$, $\mathbf{p} \in \mathbb{R}^n$ eine Abstiegsrichtung für f in \mathbf{x} und $\hat{\alpha} = \hat{\alpha}(\mathbf{x}; \mathbf{p})$ die erste positive Nullstelle der Funktion $\psi(\alpha) = f(\mathbf{x}) - f(\mathbf{x} + \alpha\mathbf{p})$. Dann gilt

$$f(\mathbf{x} + \alpha\mathbf{p}) \leq f(\mathbf{x}) + \alpha \nabla f(\mathbf{x})^T \mathbf{p} + \alpha^2 \frac{L}{2} \|\mathbf{p}\|_2^2$$

für alle $\alpha \in [0, \hat{\alpha}]$.

Beweis: Wegen Voraussetzung **V1** existiert $\hat{\alpha}$ und es ist

$$\mathbf{x} + \alpha\mathbf{p} \in L_0 \quad \forall \alpha \in [0, \hat{\alpha}].$$

Für ein $\alpha \in [0, \hat{\alpha}]$ erhält man

$$f(\mathbf{x} + \alpha\mathbf{p}) = f(\mathbf{x}) + \alpha \nabla f(\mathbf{x})^T \mathbf{p} + \int_0^\alpha [\underbrace{\nabla f(\mathbf{x} + \lambda\mathbf{p})}_{\in L_0} - \nabla f(\mathbf{x})]^T \mathbf{p} d\lambda.$$

Nach der CAUCHY-SCHWARZschen Ungleichung gilt

$$[\nabla f(\mathbf{x} + \lambda\mathbf{p}) - \nabla f(\mathbf{x})]^T \mathbf{p} \leq \|\nabla f(\mathbf{x} + \lambda\mathbf{p}) - \nabla f(\mathbf{x})\|_2 \|\mathbf{p}\|_2,$$

und mit Voraussetzung **V3** folgt

$$[\nabla f(\mathbf{x} + \lambda\mathbf{p}) - \nabla f(\mathbf{x})]^T \mathbf{p} \leq \lambda L \|\mathbf{p}\|_2^2.$$

Das liefert

$$\begin{aligned} f(\mathbf{x} + \alpha\mathbf{p}) &= f(\mathbf{x}) + \alpha \nabla f(\mathbf{x})^T \mathbf{p} + L \|\mathbf{p}\|_2^2 \int_0^\alpha \lambda d\lambda \\ &= f(\mathbf{x}) + \alpha \nabla f(\mathbf{x})^T \mathbf{p} + \alpha^2 \frac{L}{2} \|\mathbf{p}\|_2^2. \end{aligned}$$

*

Bemerkung: Für $\alpha = \hat{\alpha}$ erhält man wegen $f(\mathbf{x} + \hat{\alpha}\mathbf{p}) = f(\mathbf{x})$

$$-\frac{2\nabla f(\mathbf{x})^T \mathbf{p}}{L\|\mathbf{p}\|_2^2} \leq \hat{\alpha}.$$

Es liegt nun nahe, die Schrittweite α^* als globale oder erste lokale Lösung des eindimensionalen Minumumproblems

$$\min_{\alpha \in \mathbb{R}_+} f(\mathbf{x} + \alpha \mathbf{p})$$

zu bestimmen. Im ersten Falle gilt

$$f(\mathbf{x} + \alpha^* \mathbf{p}) = \min_{\alpha \in \mathbb{R}_+} f(\mathbf{x} + \alpha \mathbf{p})$$

und im zweiten Fall

$$\nabla f(\mathbf{x} + \alpha^* \mathbf{p})^T \mathbf{p} = 0, \quad \forall \alpha \in [0, \alpha^*) : \nabla f(\mathbf{x} + \alpha \mathbf{p})^T \mathbf{p} < 0.$$

In diesen Fällen spricht man von einer **exakten Schrittweite**. Im nächsten Satz wollen wir abschätzen, wie groß der erreichbare Abstieg für eine exakte Schrittweite ist.

11.16. Satz: Die Zielfunktion f genüge den Voraussetzungen **V1**, **V2** und **V3**. Der Punkt $\mathbf{x} \in L_0$ sei keine stationäre Lösung und $\mathbf{p} \in \mathbb{R}^n$ sei eine Abstiegsrichtung für f in \mathbf{x} . Weiterhin sei $\varphi(\alpha) = f(\mathbf{x} + \alpha \mathbf{p})$. Ist α^* die erste positive Nullstelle von φ' , so ist

1.

$$-\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{L \|\mathbf{p}\|_2^2} \leq \alpha^*,$$

2.

$$f(\mathbf{x}) - f(\mathbf{x} + \alpha^* \mathbf{p}) \geq \frac{1}{2L} \left(\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{\|\mathbf{p}\|_2} \right)^2.$$

Beweis: φ ist monoton fallend auf $[0, \alpha^*]$. Weiterhin ist $\alpha^* \leq \hat{\alpha}$, wobei $\hat{\alpha}$ die erste positive Nullstelle von $\varphi(0) - \varphi(\alpha)$ bezeichnet. Wegen der Lipschitzstetigkeit von ∇f auf L_0 gilt

$$\begin{aligned} 0 &= \nabla f(\mathbf{x} + \alpha^* \mathbf{p})^T \mathbf{p} + \nabla f(\mathbf{x})^T \mathbf{p} - \nabla f(\mathbf{x})^T \mathbf{p} \\ &= \nabla f(\mathbf{x})^T \mathbf{p} + [\nabla f(\mathbf{x} + \alpha^* \mathbf{p}) - \nabla f(\mathbf{x})]^T \mathbf{p} \\ &\leq \nabla f(\mathbf{x})^T \mathbf{p} + \alpha^* L \|\mathbf{p}\|_2^2. \end{aligned}$$

Daraus folgt

$$\alpha^* \geq -\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{L \|\mathbf{p}\|_2^2} = \tilde{\alpha}.$$

Wegen der Monotonie von φ gilt

$$\varphi(\alpha^*) \leq \varphi(\tilde{\alpha})$$

und damit

$$\begin{aligned} f(\mathbf{x} + \alpha^* \mathbf{p}) &\leq f(\mathbf{x} + \tilde{\alpha} \mathbf{p}) \\ &\leq f(\mathbf{x}) + \tilde{\alpha} \nabla f(\mathbf{x})^T \mathbf{p} + \tilde{\alpha}^2 \frac{L}{2} \|\mathbf{p}\|_2^2 \\ &= f(\mathbf{x}) + \left(-\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{L \|\mathbf{p}\|_2^2} \right) \nabla f(\mathbf{x})^T \mathbf{p} + \left(-\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{L \|\mathbf{p}\|_2^2} \right)^2 \frac{L}{2} \|\mathbf{p}\|_2^2 \\ &= f(\mathbf{x}) - \frac{1}{2L} \left(\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{\|\mathbf{p}\|_2} \right)^2. \end{aligned}$$

*

Leider lässt sich eine exakte Schrittweite i.a. nicht in endlich vielen Schritten berechnen. Darum verwendet man inexakte Schrittweiten.

Eine Variante ist die **POWELL-Schrittweite**:

Für ein $\mu \in (0, 1/2)$ und ein $\nu \in (\mu, 1)$ wird zu einem Punkt \mathbf{x} und einer Abstiegsrichtung \mathbf{p} eine Schrittweite α^* so bestimmt, dass

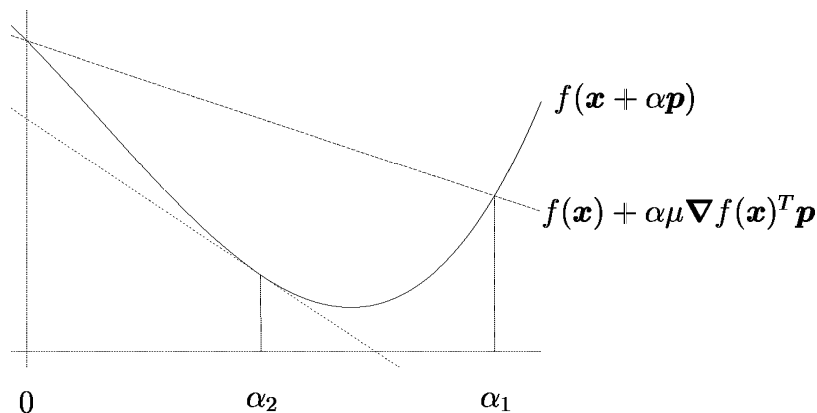
•

$$f(\mathbf{x} + \alpha^* \mathbf{p}) \leq f(\mathbf{x}) + \alpha^* \mu \nabla f(\mathbf{x})^T \mathbf{p} \quad \text{und}$$

•

$$\nabla f(\mathbf{x} + \alpha^* \mathbf{p})^T \mathbf{p} \geq \nu \nabla f(\mathbf{x})^T \mathbf{p}$$

gilt.



Für $0 \leq \alpha^* \leq \alpha_1$ ist die erste Bedingung erfüllt. Sie sichert einen hinreichend starken Abstieg.

Für $\alpha^* \geq \alpha_2$ ist die zweite Bedingung erfüllt. Sie verhindert, dass zu kleine Schrittweiten gewählt werden.

Wir wollen nun zeigen, dass stets eine POWELL-Schrittweite existiert; Außerdem ist der erreichbare Abstieg abzuschätzen.

11.17. Satz: Die Zielfunktion f genüge den Voraussetzungen **V1**, **V2** und **V3**. $\mathbf{x} \in L_0$ sei keine stationäre Lösung und $\mathbf{p} \in \mathbb{R}^n$ sei eine Abstiegsrichtung für f in \mathbf{x} . Für gegebenes $\mu \in (0, 1/2)$ und $\nu \in (\mu, 1)$ bezeichne

$$T_P(\mathbf{x}; \mathbf{p}) = \left\{ \alpha > 0 \mid \begin{array}{l} f(\mathbf{x} + \alpha\mathbf{p}) \leq f(\mathbf{x}) + \alpha\mu \nabla f(\mathbf{x})^T \mathbf{p}, \\ \nabla f(\mathbf{x} + \alpha\mathbf{p})^T \mathbf{p} \geq \nu \nabla f(\mathbf{x})^T \mathbf{p} \end{array} \right\}$$

die Menge der POWELL-Schrittweiten im Punkt \mathbf{x} in Richtung \mathbf{p} . Dann gilt

1. Die Menge $T_P(\mathbf{x}; \mathbf{p})$ ist nicht leer.
2. Es existiert eine nur von μ , ν und der LIPSCHITZ-Konstanten L abhängende Konstante $\Theta > 0$ mit

$$f(\mathbf{x}) - f(\mathbf{x} + \alpha\mathbf{p}) \geq \Theta \left(\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{\|\mathbf{p}\|_2} \right)^2$$

für alle $\alpha \in T_P(\mathbf{x}; \mathbf{p})$.

Beweis: Es sei

$$\Phi(\alpha) = f(\mathbf{x}) - f(\mathbf{x} + \alpha\mathbf{p}) + \alpha\mu \nabla f(\mathbf{x})^T \mathbf{p}$$

und α^* die erste positive Nullstelle von $\nabla f(\mathbf{x} + \alpha\mathbf{p})^T \mathbf{p}$. Dann gilt

$$\begin{aligned} \Phi'(\alpha) &= -\nabla f(\mathbf{x} + \alpha\mathbf{p})^T \mathbf{p} + \mu \nabla f(\mathbf{x})^T \mathbf{p}, \\ \Phi'(0) &= -\underbrace{(1-\mu)}_{>0} \underbrace{\nabla f(\mathbf{x})^T \mathbf{p}}_{<0} > 0, \\ \Phi'(\alpha^*) &= \mu \nabla f(\mathbf{x})^T \mathbf{p} < 0. \end{aligned}$$

Wegen $\Phi(0) = 0$ existiert dann ein $\bar{\alpha} \in (0, \alpha^*)$ mit $\Phi(\bar{\alpha}) > 0$ und $\Phi'(\bar{\alpha}) = 0$. Daraus folgt

$$\begin{aligned} f(\mathbf{x} + \bar{\alpha}\mathbf{p}) &> f(\mathbf{x}) + \bar{\alpha}\mu \nabla f(\mathbf{x})^T \mathbf{p}, \\ \nabla f(\mathbf{x} + \bar{\alpha}\mathbf{p})^T \mathbf{p} &= \mu \nabla f(\mathbf{x})^T \mathbf{p} > \nu \nabla f(\mathbf{x})^T \mathbf{p}, \end{aligned}$$

daher $\bar{\alpha} \in T_P(\mathbf{x}; \mathbf{p})$ und $T_P(\mathbf{x}; \mathbf{p}) \neq \emptyset$. Weiter gilt

$$\begin{aligned} -(1-\nu)\nabla f(\mathbf{x})^T \mathbf{p} &= \nu \nabla f(\mathbf{x})^T \mathbf{p} - \nabla f(\mathbf{x})^T \mathbf{p} \\ &\leq [\nabla f(\mathbf{x} + \alpha \mathbf{p}) - \nabla f(\mathbf{x})]^T \mathbf{p} \\ &\leq \|\nabla f(\mathbf{x} + \alpha \mathbf{p}) - \nabla f(\mathbf{x})\|_2 \|\mathbf{p}\|_2 \\ &\leq L \|\mathbf{x} + \alpha \mathbf{p} - \mathbf{x}\|_2 \|\mathbf{p}\|_2 \\ &= L\alpha \|\mathbf{p}\|_2^2. \end{aligned}$$

Daraus folgt

$$\alpha \geq -\frac{1-\nu}{L} \frac{\nabla f(\mathbf{x})^T \mathbf{p}}{\|\mathbf{p}\|_2^2}$$

und weiter

$$f(\mathbf{x} + \alpha \mathbf{p}) \leq f(\mathbf{x}) + \alpha \mu \nabla f(\mathbf{x})^T \mathbf{p} \leq f(\mathbf{x}) - \frac{\mu(1-\nu)}{L} \left(\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{\|\mathbf{p}\|_2} \right)^2.$$

Es gilt daher die Ungleichung mit $\Theta = \frac{\mu(1-\nu)}{L}$. *

Eine andere Variante der Schrittweitenbestimmung ist das ARMIJO-Prinzip:

Man wählt wieder ein $\mu \in (0, 1/2)$ und testet zunächst mit der Schrittweite $\alpha = 1$ die Gültigkeit der Ungleichung

$$f(\mathbf{x} + \alpha \mathbf{p}) \leq f(\mathbf{x}) + \alpha \mu \nabla f(\mathbf{x})^T \mathbf{p}.$$

Ist die Ungleichung erfüllt, so wird α als Schrittweite akzeptiert, Anderenfalls wird α kontrolliert verkleinert. Das könnte wie folgt geschehen:

S0. Wähle Zahlen $\mu \in (0, 1/2)$ und $0 < \varrho \leq \sigma < 1$. Setze $\alpha_0 = 1$ und $k = 0$.

S1. Falls $f(\mathbf{x} + \alpha_k \mathbf{p}) \leq f(\mathbf{x}) + \alpha_k \mu \nabla f(\mathbf{x})^T \mathbf{p}$, so STOPP. α_k ist die gesuchte Schrittweite.

S2. Wähle $\alpha_{k+1} \in [\varrho \alpha_k, \sigma \alpha_k]$.

S3. Setze $k = k + 1$ und gehe zu Schritt **S1**.

Für $\varrho = \sigma$ ist die ARMIJO-Schrittweite durch $\alpha = \varrho^j$ gegeben, wobei j die kleinste ganze Zahl ist, für die die Ungleichung

$$f(\mathbf{x} + \varrho^j \mathbf{p}) \leq f(\mathbf{x}) + \varrho^j \mu \nabla f(\mathbf{x})^T \mathbf{p}$$

erfüllt ist.

Eine andere Variante stammt von HAN (1981). Er ersetzt die Funktion

$$\varphi(\alpha) = f(\mathbf{x} + \alpha \mathbf{p})$$

durch eine quadratische Funktion $q(\alpha)$ mit

- $q(0) = \varphi(0) = f(\mathbf{x})$,
- $q(\alpha_k) = \varphi(\alpha_k) = f(\mathbf{x} + \alpha_k \mathbf{p})$,
- $q'(0) = \varphi'(0) = \nabla f(\mathbf{x})^T \mathbf{p}$.

Als Minimalpunkt dieser Ersatzfunktion ergibt sich

$$\alpha_k^* = -\frac{1}{2} \frac{\alpha_k^2 \nabla f(\mathbf{x})^T \mathbf{p}}{f(\mathbf{x} + \alpha_k \mathbf{p}) - f(\mathbf{x}) - \alpha_k \nabla f(\mathbf{x})^T \mathbf{p}}.$$

Nun wird die nächste Schrittweite durch

$$\alpha_{k+1} = \max\{0.1\alpha_k, \alpha_k^*\}$$

bestimmt.

Ist die Ungleichung $f(\mathbf{x} + \alpha_k \mathbf{p}) \leq f(\mathbf{x}) + \alpha_k \mu \nabla f(\mathbf{x})^T \mathbf{p}$ im k -ten Schritt nicht erfüllt, so folgt aus $\alpha_{k+1} \leq \alpha_k^*$

$$\alpha_{k+1} \leq \frac{1}{2(1-\mu)} \alpha_k.$$

Andererseits gilt immer $\alpha_{k+1} \geq 0.1\alpha_k$. Damit entsprechen die gemäß HAN berechneten Schrittweiten dem ARMIJO-Konzept mit $\varrho = 0.1$ und $\sigma = \frac{1}{2(1-\mu)}$.

Auch für ARMIJO-Schrittweiten lässt sich eine zu Satz 11.17 analoge Existenzaussage machen, die wir hier ohne Beweis angeben.

11.18. Satz: Die Zielfunktion f genüge den Voraussetzungen **V1**, **V2** und **V3**. $\mathbf{x} \in L_0$ sei keine stationäre Lösung und $\mathbf{p} \in \mathbb{R}^n$ sei eine Abstiegsrichtung für f in \mathbf{x} . Für gegebenes $\mu \in (0, 1/2)$ und $0 < \varrho \leq \sigma < 1$ sei α eine zugehörige ARMIJO-Schrittweite. Dann existiert eine Konstante $\Theta > 0$, die nur von μ , ϱ , σ und der LIPSCHITZ-Konstanten L abhängt, mit

$$f(\mathbf{x}) - f(\mathbf{x} + \alpha \mathbf{p}) \geq \Theta \min \left\{ \left(\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{\|\mathbf{p}\|_2} \right)^2 ; -\nabla f(\mathbf{x})^T \mathbf{p} \right\}.$$

11.2.2. Konvergenz des Modellalgorithmus

11.19. Satz: Die Zielfunktion f genüge den Voraussetzungen **V1**, **V2** und **V3**. Als Schrittweite im Modellalgorithmus werde

- die exakte Schrittweite $\alpha_k = \alpha^*(\mathbf{x}^{(k)}; \mathbf{p}^{(k)})$ oder

- die POWELL-Schrittweite $\alpha_k = \alpha_P(\mathbf{x}^{(k)}; \mathbf{p}^{(k)})$ oder
- die ARMIJO-Schrittweite $\alpha_k = \alpha_A(\mathbf{x}^{(k)}; \mathbf{p}^{(k)})$

verwendet. Weiterhin existiere eine Konstante $\gamma > 0$ mit

$$-\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)} \geq \gamma \|\nabla f(\mathbf{x}^{(k)})\| \|\mathbf{p}^{(k)}\|$$

und eine Konstante $\tau > 0$ mit

$$\|\mathbf{p}^{(k)}\| \geq \tau \|\nabla f(\mathbf{x}^{(k)})\|, \quad k = 0, 1, \dots$$

Dann gilt:

1. Jeder Häufungspunkt der erzeugten Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ ist eine stationäre Lösung des Minimumproblems

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}).$$

2. Besitzt das Minimumproblem genau eine stationäre Lösung \mathbf{x}^* in der Niveaumenge L_0 , so konvergiert die gesamte Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ gegen \mathbf{x}^* .

Beweis: Wegen Satz 11.16, Satz 11.17 und Satz 11.18 existiert eine Konstante $\Theta > 0$ mit

$$f(\mathbf{x}) - f(\mathbf{x} + \alpha \mathbf{p}) \geq \Theta \min \left\{ \left(\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{\|\mathbf{p}\|_2} \right)^2 ; -\nabla f(\mathbf{x})^T \mathbf{p} \right\}.$$

Mit den beiden Voraussetzungen folgt dann

$$f(\mathbf{x}) - f(\mathbf{x} + \alpha \mathbf{p}) \geq \Theta \min\{\gamma\tau, \gamma^2\} \|\nabla f(\mathbf{x}^{(k)})\|^2.$$

Da $\{f(\mathbf{x}^{(k)})\}_{k \in \mathbb{N}}$ eine monoton fallende, nach unten beschränkte Folge ist, muss die Folge $\{\nabla f(\mathbf{x}^{(k)})\}_{k \in \mathbb{N}}$ gegen den Nullvektor konvergieren. Ist \mathbf{x}^* ein Häufungspunkt von $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$, so existiert eine Teilfolge $\{\mathbf{x}^{(k_j)}\}_{j \in \mathbb{N}} \subseteq \{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ mit

$$\lim_{j \rightarrow \infty} \mathbf{x}^{(k_j)} = \mathbf{x}^*.$$

Wegen

$$\lim_{k \rightarrow \infty} \nabla f(\mathbf{x}^{(k)}) = \mathbf{o}$$

gilt dann aber $\nabla f(\mathbf{x}^*) = \mathbf{o}$. \mathbf{x}^* ist daher stationäre Lösung. Damit ist die erste Aussage bewiesen.

Wir nehmen nun an, das Minimumproblem besitzt genau eine stationäre Lösung \mathbf{x}^* .

Falls $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ nicht gegen \mathbf{x}^* konvergiert, existieren ein $\varepsilon > 0$ und eine Teilfolge $\{\mathbf{x}^{(k_j)}\}_{j \in \mathbb{N}} \subseteq \{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ mit

$$\|\mathbf{x}^{(k_j)} - \mathbf{x}^*\| \geq \varepsilon \quad j = 1, 2, \dots$$

Da L_0 als kompakt vorausgesetzt wurde, besitzt die Folge $\{\mathbf{x}^{(k_j)}\}_{j \in \mathbb{N}}$ einen Häufungspunkt $\hat{\mathbf{x}}$. Nach dem gerade bewiesenen ersten Teil gilt dann aber $\nabla f(\hat{\mathbf{x}}) = \mathbf{o}$, d. h. $\hat{\mathbf{x}}$ ist stationäre Lösung. Andererseits folgt aber aus $\|\hat{\mathbf{x}} - \mathbf{x}^*\| \geq \varepsilon > 0$, dass $\hat{\mathbf{x}} \neq \mathbf{x}^*$ sein muss, im Widerspruch zur Voraussetzung „ \mathbf{x}^* ist einzige stationäre Lösung“. Damit muss die Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ gegen \mathbf{x}^* konvergieren. *

Bemerkungen: (i) Die konkrete Art der Schrittweitenberechnung spielt in diesem Konvergenzsatz keine Rolle. Wichtig ist nur die Sicherung eines im Sinne von

$$f(\mathbf{x}) - f(\mathbf{x} + \alpha \mathbf{p}) \geq \Theta \min \left\{ \left(\frac{\nabla f(\mathbf{x})^T \mathbf{p}}{\|\mathbf{p}\|_2} \right)^2 ; -\nabla f(\mathbf{x})^T \mathbf{p} \right\}$$

hinreichend starken Abstiegs pro Schritt.

(ii) Für exakte und POWELL-Schrittweiten darf auf die Voraussetzung

$$\|\mathbf{p}^{(k)}\| \geq \tau \|\nabla f(\mathbf{x}^{(k)})\|$$

mit $\tau > 0$ verzichtet werden. Nun wollen wir wieder eine glatte, gleichmäßig konvexe Zielfunktion betrachten. Zur Vorbereitung beweisen wir den folgenden Hilfssatz.

11.20. Satz: Gegeben sei das freie Minimumproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}).$$

Die Zielfunktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ erfülle folgende Voraussetzungen:

K1 Für ein gegebenes $\mathbf{x}^{(0)} \in \mathbb{R}^n$ (i.a. der Startpunkt des Verfahrens) ist die Niveaumenge

$$L_0 = \left\{ \mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \leq f(\mathbf{x}^{(0)}) \right\} \text{ konvex.}$$

K2 Die Zielfunktion f ist auf einer offenen Menge $D \supseteq L_0$ stetig differenzierbar und auf L_0 gleichmäßig konvex: Es existiert eine Konstante $\gamma > 0$ mit

$$\frac{\gamma}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) \leq f(\mathbf{y}) - f(\mathbf{x})$$

für alle $\mathbf{x}, \mathbf{y} \in L_0$.

K3 Der Gradient $\nabla f(\mathbf{x}) = f'(\mathbf{x})^T$ ist auf L_0 lipschitzstetig: Es existiert eine Konstante $L > 0$ mit

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

für alle $\mathbf{x}, \mathbf{y} \in L_0$.

Dann ist die Niveaumenge L_0 kompakt und das Minimumproblem besitzt genau eine globale Lösung $\mathbf{x}^* \in L_0$. \mathbf{x}^* ist auch die einzige stationäre Lösung auf L_0 . Es gilt die Fehlerabschätzung

$$\frac{\gamma}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2 \leq f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{1}{2\gamma} \|\nabla f(\mathbf{x})\|_2^2 \quad (11.1)$$

für alle $\mathbf{x} \in L_0$.

Beweis: Die Niveaumenge L_0 ist offensichtlich abgeschlossen. Aus der gleichmäßigen Konvexität von f auf L_0 folgt

$$\frac{\gamma}{2} \|\mathbf{x} - \mathbf{x}^{(0)}\|_2^2 + \nabla f(\mathbf{x}^{(0)})^T (\mathbf{x} - \mathbf{x}^{(0)}) \leq f(\mathbf{x}) - f(\mathbf{x}^{(0)}) \leq 0.$$

Mit Hilfe der CAUCHY-SCHWARZschen Ungleichung erhält man daraus

$$\frac{\gamma}{2} \|\mathbf{x} - \mathbf{x}^{(0)}\|_2^2 \leq -\nabla f(\mathbf{x}^{(0)})^T (\mathbf{x} - \mathbf{x}^{(0)}) \leq \|\nabla f(\mathbf{x}^{(0)})\|_2 \|\mathbf{x} - \mathbf{x}^{(0)}\|_2,$$

und weiter

$$\|\mathbf{x} - \mathbf{x}^{(0)}\|_2 \leq \frac{2}{\gamma} \|\nabla f(\mathbf{x}^{(0)})\|_2.$$

Damit ist L_0 auch beschränkt und daher kompakt. f nimmt dann sein globales Minimum auf L_0 an. Dieses globale Minimum ist offensichtlich auch globale Lösung des Minimumproblems. Damit ist die Existenz einer globalen Lösung \mathbf{x}^* mit $\mathbf{x}^* \in L_0$ bewiesen.

Wir zeigen nun die Gültigkeit der Abschätzung 11.1, aus der auch die Eindeutigkeit

von \mathbf{x}^* folgt. Aus der gleichmäßigen Konvexität von f auf L_0 (**K2**) erhält man mit $\mathbf{y} = \mathbf{x}$, $\mathbf{x} = \mathbf{x}^*$ und $\nabla f(\mathbf{x}^*) = \mathbf{o}$

$$\frac{\gamma}{2} \|\mathbf{x} - \mathbf{x}^*\| \leq f(\mathbf{x}) - f(\mathbf{x}^*).$$

Das ist die erste zu beweisende Ungleichung. Zum Beweis der zweiten Ungleichung betrachten wir das Minimumproblem

$$\min_{\mathbf{p} \in \mathbb{R}^n} \frac{\gamma}{2} \|\mathbf{p}\|_2^2 + \nabla f(\mathbf{x})^T \mathbf{p}$$

für ein festes $\mathbf{x} \in L_0$. Die Lösung ist offensichtlich durch

$$\mathbf{p} = -\frac{1}{\gamma} \nabla f(\mathbf{x})$$

gegeben. Dann gilt für beliebige $\mathbf{p} \in \mathbb{R}^n$

$$\frac{\gamma}{2} \|\mathbf{p}\|_2^2 + \nabla f(\mathbf{x})^T \mathbf{p} \geq \frac{\gamma}{2} \left\| \frac{1}{\gamma} \nabla f(\mathbf{x}) \right\|_2^2 - \frac{1}{\gamma} \nabla f(\mathbf{x})^T \nabla f(\mathbf{x}) = -\frac{1}{2\gamma} \|\nabla f(\mathbf{x})\|_2^2.$$

Speziell für $\mathbf{p} = \mathbf{x}^* - \mathbf{x}$ gilt

$$\frac{\gamma}{2} \|\mathbf{x}^* - \mathbf{x}\|_2^2 + \nabla f(\mathbf{x})^T (\mathbf{x}^* - \mathbf{x}) \geq -\frac{1}{2\gamma} \|\nabla f(\mathbf{x})\|_2^2.$$

Wegen der gleichmäßigen Konvexität gilt

$$f(\mathbf{x}^*) - f(\mathbf{x}) \geq \frac{\gamma}{2} \|\mathbf{x}^* - \mathbf{x}\|_2^2 + \nabla f(\mathbf{x})^T (\mathbf{x}^* - \mathbf{x}),$$

also insgesamt

$$f(\mathbf{x}^*) - f(\mathbf{x}) \geq -\frac{1}{2\gamma} \|\nabla f(\mathbf{x})\|_2^2$$

beziehungsweise

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{1}{2\gamma} \|\nabla f(\mathbf{x})\|_2^2.$$

Wäre nun $\bar{\mathbf{x}}$ eine weitere stationäre Lösung in L_0 , so würde aus der Abschätzung 11.1

$$\|\bar{\mathbf{x}} - \mathbf{x}^*\| \leq \frac{1}{\gamma^2} \|\nabla f(\bar{\mathbf{x}})\|_2^2 = 0$$

folgen, d. h. $\bar{\mathbf{x}} = \mathbf{x}^*$.



Für gleichmäßig konvexe Zielfunktionen ist nun eine hinreichende Bedingung für die Konvergenz des Modellalgorithmus angebar.

11.21. Satz: Die Zielfunktion f genüge den Voraussetzungen **K1**, **K2** und **K3** aus Satz 11.20. Als Schrittweite im Modellalgorithmus werde

- die exakte Schrittweite $\alpha_k = \alpha^* \left(\mathbf{x}^{(k)}; \mathbf{p}^{(k)} \right)$ oder
- die POWELL-Schrittweite $\alpha_k = \alpha_P \left(\mathbf{x}^{(k)}; \mathbf{p}^{(k)} \right)$ oder
- die ARMIJO-Schrittweite $\alpha_k = \alpha_A \left(\mathbf{x}^{(k)}; \mathbf{p}^{(k)} \right)$

verwendet. Weiterhin sei

$$\delta_k = \min \left\{ -\frac{\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)}}{\|\nabla f(\mathbf{x}^{(k)})\|^2}; \left(\frac{\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)}}{\|\nabla f(\mathbf{x}^{(k)})\| \|\mathbf{p}^{(k)}\|} \right)^2 \right\}$$

für ARMIJO-Schrittweiten, bzw.

$$\delta_k = \left(\frac{\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)}}{\|\nabla f(\mathbf{x}^{(k)})\| \|\mathbf{p}^{(k)}\|} \right)^2$$

für exakte oder POWELL-Schrittweiten. Dann gilt:

1. Ist

$$\sum_{j=0}^{\infty} \delta_j = \infty,$$

so konvergiert die durch den Modellalgorithmus erzeugte Folge $\left\{ \mathbf{x}^{(k)} \right\}_{k \in \mathbb{N}}$ gegen die eindeutige globale Lösung \mathbf{x}^* .

2. Existiert ein $\delta > 0$ mit

$$\frac{1}{k+1} \sum_{j=0}^k \delta_j > \delta$$

für $k = 0, 1, \dots$, so existieren Konstanten $C > 0$ und $q \in (0, 1)$ mit

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq Cq^k$$

für $k = 0, 1, \dots$

Beweis: Nach Satz 11.16 und Satz 11.17 existiert sowohl für exakte als auch für und POWELL-Schrittweiten eine Konstante $\Theta > 0$ mit

$$f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k+1)}) = f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}) \geq \Theta \left(\frac{\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)}}{\|\mathbf{p}^{(k)}\|_2} \right)^2.$$

Nach Satz 11.18 existiert für ARMIJO-Schrittweiten eine Konstante $\Theta > 0$ mit

$$\begin{aligned} f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k+1)}) &= f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}) \\ &\geq \Theta \min \left\{ \left(\frac{\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)}}{\|\mathbf{p}^{(k)}\|_2} \right)^2 ; -\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)} \right\}. \end{aligned}$$

Mit der Definition von δ_k folgt dann

$$f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k+1)}) \geq \Theta \delta_k \|\nabla f(\mathbf{x}^{(k)})\|^2.$$

Wendet man die Fehlerabschätzung aus Satz 11.20 an, so erhält man

$$f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k+1)}) \geq 2\gamma\Theta\delta_k \left[f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*) \right].$$

Damit folgt

$$\begin{aligned} 0 &\leq f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^*) \\ &= f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^{(k)}) + f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*) \\ &\leq (1 - 2\gamma\Theta\delta_k) \left[f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*) \right] \\ &\leq \prod_{j=0}^k (1 - 2\gamma\Theta\delta_j) \cdot \left[f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*) \right] \\ &\leq \exp \left(-2\gamma\Theta \sum_{j=0}^k \delta_j \right) \cdot \left[f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*) \right]. \end{aligned}$$

Wegen

$$\lim_{k \rightarrow \infty} \sum_{j=0}^k \delta_j = \infty$$

gilt dann

$$\lim_{k \rightarrow \infty} \exp \left(-2\gamma\Theta \sum_{j=0}^k \delta_j \right) = 0$$

und

$$0 \leq \lim_{k \rightarrow \infty} f(\mathbf{x}^{(k+1)}) - f(\mathbf{x}^*) \leq 0,$$

also

$$\lim_{k \rightarrow \infty} f(\mathbf{x}^{(k)}) = f(\mathbf{x}^*).$$

Nach Satz 11.20 gilt weiterhin für $\mathbf{x}^{(k)} \in L_0$

$$\begin{aligned} \|\mathbf{x}^{(k)} - \mathbf{x}^*\| &\leq \sqrt{\frac{2}{\gamma} [f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*)]} \\ &\leq \sqrt{\frac{2}{\gamma} \exp\left(-2\gamma\Theta \sum_{j=0}^{k-1} \delta_j\right) [f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*)]}. \end{aligned}$$

Daraus ergibt sich sofort

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*.$$

Falls ein δ mit

$$\sum_{j=0}^{k-1} \delta_j \geq k\delta$$

für $k = 0, 1, \dots$ existiert, so folgt

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\| \leq e^{-k\gamma\Theta\delta} \sqrt{\frac{2}{\gamma} [f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*)]} = Cq^k$$

mit

$$C = \sqrt{\frac{2}{\gamma} [f(\mathbf{x}^{(0)}) - f(\mathbf{x}^*)]}, \quad q = e^{-\gamma\Theta\delta} \in (0, 1).$$

*

Bemerkungen: (i) Die konkrete Schrittweitenwahl spielt für die Konvergenzaussage wieder keine Rolle. Wesentlich ist nur die Sicherung eines hinreichend starken

Abstiegs.

(ii) Für exakte und POWELL-Schrittweiten ist

$$\delta_k = \left(\frac{\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)}}{\|\nabla f(\mathbf{x}^{(k)})\| \|\mathbf{p}^{(k)}\|} \right)^2$$

ein Maß für die Größe des Winkels zwischen der Richtung $\mathbf{p}^{(k)}$ und dem negativen Gradienten von f im Punkt $\mathbf{x}^{(k)}$. Die Bedingung

$$\sum_{j=0}^{\infty} \delta_j = \infty$$

besagt dann, dass diese Winkel nicht zu schnell gegen $\pi/2$ gehen dürfen.

11.3. Quasi-Newton-Verfahren

11.3.1. Gedämpftes und ungedämpftes Newton-Verfahren

Wir betrachten wieder das freie Minimumproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

mit der auf dem gesamten \mathbb{R}^n zweimal stetig differenzierbaren Zielfunktion

$$f : \mathbb{R}^n \longrightarrow \mathbb{R}.$$

Stationäre Lösungen dieses Problems sind Lösungen des im allgemeinen nichtlinearen Gleichungssystems

$$\nabla f(\mathbf{x}) = \mathbf{o}.$$

Zum Lösen dieses Systems könnte man das NEWTON-Verfahren anwenden. Hier wird ausgehend von einem Startpunkt $\mathbf{x}^{(0)} \in \mathbb{R}^n$ mit der Rekursionsformel

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left(\nabla^2 f(\mathbf{x}^{(k)}) \right)^{-1} \nabla f(\mathbf{x}^{(k)})$$

eine Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ berechnet. Der Konvergenzsatz für das NEWTON-Verfahren, angewendet auf dieses spezielle Problem, lautet:

11.22. Satz: Die Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ sei auf einer offenen Umgebung des Punktes $\mathbf{x}^* \in \mathbb{R}^n$ zweimal stetig differenzierbar. Es sei $\nabla f(\mathbf{x}^*) = \mathbf{o}$ und $\nabla^2 f(\mathbf{x}^*)$ sei regulär. Dann existiert ein $\delta > 0$, so dass für jedes

$$\mathbf{x}^{(0)} \in U_\delta(\mathbf{x}^*) = \{ \mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x} - \mathbf{x}^*\| < \delta \}$$

die durch das NEWTON-Verfahren

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \left(\nabla^2 f(\mathbf{x}^{(k)}) \right)^{-1} \nabla f(\mathbf{x}^{(k)})$$

erzeugte Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ definiert ist und superlinear gegen \mathbf{x}^* konvergiert:

$$\lim_{k \rightarrow \infty} \mathbf{x}^{(k)} = \mathbf{x}^*, \quad \lim_{k \rightarrow \infty} \frac{\|\mathbf{x}^{(k+1)} - \mathbf{x}^*\|}{\|\mathbf{x}^{(k)} - \mathbf{x}^*\|} = 0.$$

Ist außerdem $\nabla^2 f(\mathbf{x})$ auf einer hinreichend kleinen Kugel

$$U_\eta(\mathbf{x}^*) = \{ \mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x} - \mathbf{x}^*\| < \eta \}$$

lipschitzstetig mit der LIPSCHITZ-Konstanten L :

$$\forall \mathbf{x} \in U_\eta : \quad \|\nabla^2 f(\mathbf{x}) - \nabla^2 f(\mathbf{x}^*)\| \leq L \|\mathbf{x} - \mathbf{x}^*\|,$$

so konvergiert die Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ lokal quadratisch gegen \mathbf{x}^* .

Die Konvergenzeigenschaften dieses Verfahrens hängen stark von der Größe der Umgebungen $U_\delta(\mathbf{x}^*)$ und $U_\eta(\mathbf{x}^*)$ ab. Man sollte nun versuchen, durch Einführung von Schrittweiten globale Konvergenzeigenschaften zu erhalten. Damit kommt man zum gedämpften NEWTON-Verfahren:

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \alpha_k \left(\nabla^2 f(\mathbf{x}^{(k)}) \right)^{-1} \nabla f(\mathbf{x}^{(k)}).$$

Unter entsprechenden Konvexitätsvoraussetzungen an f lässt sich die Regularität und sogar die positive Definitheit von $\nabla^2 f(\mathbf{x}^{(k)})$ sichern. Dann ist die NEWTON-Richtung

$$\mathbf{p}^{(k)} = \left(\nabla^2 f(\mathbf{x}^{(k)}) \right)^{-1} \nabla f(\mathbf{x}^{(k)})$$

eine Abstiegsrichtung, und man wird die Konvergenz des gedämpften NEWTON-Verfahrens erwarten. Die Schrittweiten sollten so gewählt werden, dass das Verfahren bei hinreichender Annäherung an \mathbf{x}^* automatisch in das ungedämpfte NEWTON-Verfahren (Schrittweite 1) übergeht. Wir geben den entsprechenden Konvergenzsatz ohne Beweis an.

11.23. Satz: Die Zielfunktion f des freien Minimumproblems

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

genüge den Voraussetzungen:

1. Mit einem $\mathbf{x}^{(0)} \in \mathbb{R}^n$ ist die Niveaumenge

$$L_0 = \left\{ \mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \leq f(\mathbf{x}^{(0)}) \right\}$$

konvex.

2. Die Funktion f ist auf einer Umgebung von L_0 zweimal stetig differenzierbar und es existieren Konstanten $0 < m \leq M < \infty$, so dass

$$m \|\mathbf{p}\|_2^2 \leq \mathbf{p}^T \nabla^2 f(\mathbf{x}) \mathbf{p} \leq M \|\mathbf{p}\|_2^2$$

für alle $\mathbf{x} \in L_0$ und alle $\mathbf{p} \in \mathbb{R}^n$ gilt.

Für das gedämpfte NEWTON-Verfahren

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}, \quad \mathbf{p}^{(k)} = - \left(\nabla^2 f(\mathbf{x}^{(k)}) \right)^{-1} \nabla f(\mathbf{x}^{(k)}),$$

wobei zur Schrittweitenbestimmung das ARMIJO-Prinzip mit $\mu \in (0, 1/2)$ angewendet wird, gilt:

- Bricht das Verfahren nicht vorzeitig mit der Lösung \mathbf{x}^* ab, so konvergiert die vom Verfahren erzeugte Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ gegen \mathbf{x}^* .
- Für hinreichend großes k geht das Verfahren in das ungedämpfte NEWTON-Verfahren über.

Bemerkung: Falls f zweimal stetig differenzierbar und $\nabla^2 f(\mathbf{x}^*)$ positiv definit ist, ist die zweite Voraussetzung immer erfüllt. Die entsprechende Umgebung ist oft klein, wie das folgende Beispiel zeigt.

11.24. Beispiel: Wir betrachten die ROSENBROCK-Funktion

$$f(\mathbf{x}) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2.$$

Es gilt

$$\nabla f(\mathbf{x}) = \begin{pmatrix} 400x_1^3 - 400x_1x_2 + 2x_1 - 2 \\ -200x_1^2 + 200x_2 \end{pmatrix}$$

und

$$\nabla^2 f(\mathbf{x}) = \begin{pmatrix} 1200x_1^2 - 400x_2 + 2 & -400x_1 \\ -400x_1 & 200 \end{pmatrix}.$$

f besitzt das einzige globale Minimum $\mathbf{x}^* = (1, 1)^T$ mit

$$\nabla^2 f(\mathbf{x}^*) = \begin{pmatrix} 802 & -400 \\ -400 & 200 \end{pmatrix}.$$

Die HESSE-Matrix $\nabla^2 f$ besitzt im Punkt \mathbf{x}^* die Eigenwerte

$$\lambda_1 = 501 + \sqrt{250601} \approx 1001.6, \quad \lambda_2 = 501 - \sqrt{250601} \approx 0.4.$$

Sie ist daher positiv definit. Der nächste Punkt mit nur positiv semidefiniter Matrix (ein Eigenwert ist 0) ist

$$\tilde{\mathbf{x}} = (0.9979\dots, 1.0010\dots).$$

Der Abstand zwischen \mathbf{x}^* und $\tilde{\mathbf{x}}$ ist klein:

$$\|\mathbf{x}^* - \tilde{\mathbf{x}}\|_2 = 0.0022\dots$$

Für die ROSENBROCK-Funktion würde damit das gedämpfte NEWTON-Verfahren erst in einer Umgebung des globalen Minimums vom Radius $\delta \approx 0.0022$ konvergieren. ♡

Das gedämpfte NEWTON-Verfahren ist ein typisches lokales Verfahren. Weitab von einem lokalen Minimum wird die HESSE-Matrix im allgemeinen nicht positiv definit sein. Dann ist aber durch die NEWTON-Richtung

$$\mathbf{p}^{(k)} = - \left(\nabla^2 f(\mathbf{x}^{(k)}) \right)^{-1} \nabla f(\mathbf{x}^{(k)})$$

keine Abstiegsrichtung gegeben. Doch selbst dann, wenn man so nahe an einer lokalen Lösung ist, dass die positive Definitheit der HESSE-Matrix gesichert ist, hat das gedämpfte NEWTON-Verfahren noch einen großen Nachteil. In jedem Schritt ist nicht nur der Gradient von f zu berechnen, sondern auch noch die HESSE-Matrix. Oft ist aber die Funktion f nicht durch einen analytischen Ausdruck gegeben, so dass man schon bei der Gradientenberechnung Näherungsformeln anzuwenden hat. Dann ist die Approximation der HESSE-Matrix erst recht problematisch. Außerdem ist in jedem Schritt ein lineares Gleichungssystem der Dimension n zu lösen. Da es im allgemeinen nicht möglich ist, Informationen aus dem vorigen Schritt zu verwenden, liegt der Aufwand pro Schritt bestenfalls in der Größenordnung von $n^3/6$ Rechenoperationen.

11.3.2. Verfahren der Oren-Luenberger-Klasse

Das NEWTON-Verfahren und das gedämpfte NEWTON-Verfahren sind nur lokal konvergent und damit praktisch nicht von großem Nutzen. Sie geben aber Hinweise zur Konstruktion global konvergenter Verfahren. Wir haben gesehen, dass die NEWTON-Richtung

$$\mathbf{p}^{(k)} = - \left(\nabla^2 f(\mathbf{x}^{(k)}) \right)^{-1} \nabla f(\mathbf{x}^{(k)})$$

in der Nähe einer isolierten lokalen Lösung zu schnell konvergenten Verfahren führt. Wählt man nun positiv definite Matrizen \mathbf{H}_k und damit Richtungen der Form

$$\mathbf{p}^{(k)} = -\mathbf{H}_k \nabla f(\mathbf{x}^{(k)}),$$

so sind diese wegen

$$\mathbf{p}^{(k)T} \nabla f(\mathbf{x}^{(k)}) = -\nabla f(\mathbf{x}^{(k)})^T \mathbf{H}_k \nabla f(\mathbf{x}^{(k)}) < 0$$

Abstiegsrichtungen. Soll mit so definierten Abstiegsrichtungen ein im Sinne von Satz 11.19 konvergentes Verfahren erzeugt werden, so müssen Konstanten $\gamma > 0$ und $\tau > 0$ existieren, so dass für alle $k \geq 0$ die Ungleichungen

$$-\nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)} \geq \gamma \|\nabla f(\mathbf{x}^{(k)})\| \|\mathbf{p}^{(k)}\|$$

und

$$\|\mathbf{p}^{(k)}\| \geq \tau \|\nabla f(\mathbf{x}^{(k)})\|$$

erfüllt sind. Bezüglich der euklidischen Norm heißt das

$$\begin{aligned}
\frac{\nabla f(\mathbf{x}^{(k)})^T \mathbf{H}_k \nabla f(\mathbf{x}^{(k)})}{\|\nabla f(\mathbf{x}^{(k)})\|_2 \|\mathbf{H}_k \nabla f(\mathbf{x}^{(k)})\|_2} &= \frac{\left(\mathbf{H}_k^{-1/2} \mathbf{y}^{(k)}\right)^T \mathbf{H}_k \left(\mathbf{H}_k^{-1/2} \mathbf{y}^{(k)}\right)}{\left\|\mathbf{H}_k^{-1/2} \mathbf{y}^{(k)}\right\|_2 \left\|\mathbf{H}_k \mathbf{H}_k^{-1/2} \mathbf{y}^{(k)}\right\|_2} \\
&= \frac{\mathbf{y}^{(k)T} \mathbf{y}^{(k)}}{\left\|\mathbf{H}_k^{-1/2} \mathbf{y}^{(k)}\right\|_2 \left\|\mathbf{H}_k^{1/2} \mathbf{y}^{(k)}\right\|_2} \\
&= \frac{1}{\frac{\left\|\mathbf{H}_k^{-1/2} \mathbf{y}^{(k)}\right\|_2}{\|\mathbf{y}^{(k)}\|_2} \frac{\left\|\mathbf{H}_k^{1/2} \mathbf{y}^{(k)}\right\|_2}{\|\mathbf{y}^{(k)}\|_2}} \\
&\geq \frac{1}{\left\|\mathbf{H}_k^{-1/2}\right\|_2 \left\|\mathbf{H}_k^{1/2}\right\|_2} \\
&= \frac{1}{\text{cond}_2(\mathbf{H}_k^{1/2})} \\
&= \frac{1}{\sqrt{\text{cond}_2(\mathbf{H}_k)}}.
\end{aligned}$$

Die erste Bedingung ist erfüllt, falls für alle Matrizen \mathbf{H}_k

$$\text{cond}_2(\mathbf{H}_k) = \frac{\lambda_{max}^{(k)}}{\lambda_{min}^{(k)}} \leq \frac{1}{\gamma^2}$$

gilt. Weiterhin erhalten wir

$$\frac{\left\|\mathbf{H}_k \nabla f(\mathbf{x}^{(k)})\right\|_2}{\left\|\nabla f(\mathbf{x}^{(k)})\right\|_2} = \frac{\left\|\mathbf{H}_k \mathbf{H}_k^{-1} \mathbf{z}^{(k)}\right\|_2}{\left\|\mathbf{H}_k^{-1} \mathbf{z}^{(k)}\right\|_2} = \frac{\left\|\mathbf{z}^{(k)}\right\|_2}{\left\|\mathbf{H}_k^{-1} \mathbf{z}^{(k)}\right\|_2} \geq \frac{1}{\left\|\mathbf{H}_k^{-1}\right\|_2}.$$

Die zweite Bedingung ist erfüllt, falls für alle Matrizen \mathbf{H}_k

$$\frac{1}{\left\|\mathbf{H}_k^{-1}\right\|_2} = \lambda_{min}^{(k)} \geq \tau$$

gilt. (Dabei bezeichnen $\lambda_{max}^{(k)}$ und $\lambda_{min}^{(k)}$ den größten bzw. kleinsten Eigenwert der symmetrischen, positiv definiten Matrix \mathbf{H}_k .) Wir sehen, dass für Abstiegsrichtungen der Form $\mathbf{p}^{(k)} = -\mathbf{H}_k \nabla f(\mathbf{x}^{(k)})$ die Voraussetzungen von Satz 11.19 erfüllt sind, falls die Folge $\{\mathbf{H}_k\}_{k \in \mathbb{N}}$ eine Folge symmetrischer, gleichmäßig positiv definit und gleichmäßig beschränkter Matrizen ist. Unter diesen Bedingungen existieren Konstanten $0 < m \leq M < \infty$, so dass für alle $k = 0, 1, \dots$ und alle Vektoren $\mathbf{z} \in \mathbb{R}^n$

$$m \|\mathbf{z}\|_2^2 \leq \mathbf{z}^T \mathbf{H}_k \mathbf{z} \leq M \|\mathbf{z}\|_2^2$$

gilt.

Als einfachste Wahl bietet sich $\mathbf{H}_k = \mathbf{I}$ für alle $k = 0, 1, \dots$ an. Hier ist $m = M = 1$. Das zugehörige Verfahren heißt **Gradientenverfahren**. Es ist im Sinne von Satz 11.19 global konvergent. Bis zum Jahre 1959 war dieses Verfahren auch das gebräuchlichste zum Lösen von freien Minimumproblemen. Vom Gradientenverfahren ist lokal höchstens lineare Konvergenz zu erwarten. Darum liegt es nahe, andere Folgen $\{\mathbf{H}_k\}_{k \in \mathbb{N}}$ gleichmäßig positiv definit und gleichmäßig beschränkter Matrizen zu verwenden, die sich in der Nähe einer lokalen Lösung ähnlich wie die HESSE-Matrix verhalten. Entwickelt man den Gradienten an einer Stelle \mathbf{x} , so ergibt sich

$$\nabla f(\mathbf{x} + \mathbf{s}) = \nabla f(\mathbf{x}) + \nabla^2 f(\mathbf{x})\mathbf{s} + O(\|\mathbf{s}\|^2).$$

Wir lassen nun $O(\|\mathbf{s}\|^2)$ weg, ersetzen $\nabla^2 f(\mathbf{x})$ durch \mathbf{H}^{-1} und erhalten

$$\nabla f(\mathbf{x} + \mathbf{s}) = \nabla f(\mathbf{x}) + \mathbf{H}^{-1} \nabla f(\mathbf{x})\mathbf{s}$$

oder

$$\mathbf{H} (\nabla f(\mathbf{x} + \mathbf{s}) - \nabla f(\mathbf{x})) = \mathbf{s}.$$

Ist nun

$$\mathbf{x} = \mathbf{x}^{(k)}, \quad \mathbf{x} + \mathbf{s} = \mathbf{x}^{(k+1)}, \quad \mathbf{H} = \mathbf{H}_k,$$

so erhalten wir die **Quasi-NEWTON-Gleichung**

$$\mathbf{H}_k \left(\nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)}) \right) = \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}.$$

Bestimmt man nun die Matrizen \mathbf{H}_k so, dass sie nicht nur gleichmäßig positiv definit und gleichmäßig beschränkt sind, sondern dass sie auch in jedem Schritt die Quasi-NEWTON-Gleichung erfüllen, so ist neben der globalen Konvergenz auch lokal ein gutes Konvergenzverhalten zu erwarten. Diese Erwartung bestätigt sich mit Satz 11.27. Alle so erzeugten Verfahren bezeichnet man als **Quasi-NEWTON-Verfahren**. Die bekanntesten Vertreter sind das BFGS-Verfahren² und das DFP-Verfahren³. Aus verschiedenen Quasi-NEWTON-Verfahren wurde 1974 von OREN und LUENBERGER die folgende Klasse von Algorithmen definiert.

11.25. Verfahren der OREN-LUENBERGER-Klasse:

S0 (*Initialisierung*) Wähle einen Startvektor $\mathbf{x}^{(0)} \in \mathbb{R}^n$

und eine positiv definite Startmatrix $\mathbf{H}_0 \in \mathbb{R}^{n \times n}$ (z.B. $\mathbf{H}_0 = \mathbf{I}$). Setze $k = 0$.

²nach BROYDEN, FLETCHER, GOLDFARB und SHANNO, 1970

³nach DAVIDON, FLETCHER und POWELL, 1959 bzw. 1963

S1 (Abbruchbedingung) Falls $\nabla f(\mathbf{x}^{(k)}) = \mathbf{o}$, so STOPP; $\mathbf{x}^{(k)}$ ist stationäre Lösung.

S2 (Richtungswahl) Berechne eine Abstiegsrichtung

$$\mathbf{p}^{(k)} = -\mathbf{H}_k \nabla f(\mathbf{x}^{(k)}).$$

S3 (Schrittweitenbestimmung) Bestimme eine Schrittweite α_k als exakte, POWELL- oder ARMIJO-Schrittweite.

S4 (Iterationsschritt) Setze

$$\begin{aligned}\mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} + \alpha_k \mathbf{p}^{(k)}, \\ \mathbf{s}^{(k)} &= \alpha_k \mathbf{p}^{(k)} = \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}, \\ \mathbf{q}^{(k)} &= \nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)}).\end{aligned}$$

Wähle Konstanten $\gamma_k > 0$ und $\Theta_k \geq 0$ und berechne

$$\mathbf{H}_{k+1} = \Psi(\mathbf{H}_k, \mathbf{s}^{(k)}, \mathbf{q}^{(k)}; \gamma_k, \Theta_k)$$

wobei

$$\begin{aligned}\Psi(\mathbf{H}, \mathbf{s}, \mathbf{q}; \gamma, \Theta) &= \gamma \mathbf{H} + \left(1 + \gamma \Theta \frac{\mathbf{q}^T \mathbf{H} \mathbf{q}}{\mathbf{s}^T \mathbf{q}}\right) \frac{\mathbf{s} \mathbf{s}^T}{\mathbf{s}^T \mathbf{q}} \\ &\quad - \gamma \frac{1 - \Theta}{\mathbf{q}^T \mathbf{H} \mathbf{q}} \mathbf{H} \mathbf{q} \cdot \mathbf{q}^T \mathbf{H} \\ &\quad - \frac{\gamma \Theta}{\mathbf{s}^T \mathbf{q}} (\mathbf{s} \cdot \mathbf{q}^T \mathbf{H} + \mathbf{H} \mathbf{q} \cdot \mathbf{s}^T).\end{aligned}$$

Eine Formel, wie sie in Schritt **S4** verwendet wird, bezeichnet man als Update-Formel. Die neue Matrix $\Psi(\mathbf{H})$ entsteht hier aus der alten Matrix \mathbf{H} durch eine Rang-2-Modifikation. Es gilt $\text{rg}(\Psi(\mathbf{H}) - \gamma \mathbf{H}) \leq 2$.

Je nach Wahl der Parameter γ_k und Θ_k erhält man verschiedene Verfahren:

- $\gamma_k \equiv 1$ und $\Theta_k \equiv 0$. Wir erhalten das DFP-Verfahren

$$\Psi(\mathbf{H}, \mathbf{s}, \mathbf{q}) = \mathbf{H} + \frac{\mathbf{s} \mathbf{s}^T}{\mathbf{s}^T \mathbf{q}} - \frac{\mathbf{H} \mathbf{q} \cdot \mathbf{q}^T \mathbf{H}}{\mathbf{q}^T \mathbf{H} \mathbf{q}}.$$

- $\gamma_k \equiv 1$ und $\Theta_k \equiv 1$. Wir erhalten das BFGS-Verfahren

$$\Psi(\mathbf{H}, \mathbf{s}, \mathbf{q}) = \mathbf{H} + \left(1 + \frac{\mathbf{q}^T \mathbf{H} \mathbf{q}}{\mathbf{s}^T \mathbf{q}}\right) \frac{\mathbf{s} \mathbf{s}^T}{\mathbf{s}^T \mathbf{q}} - \frac{\mathbf{s} \cdot \mathbf{q}^T \mathbf{H} + \mathbf{H} \mathbf{q} \cdot \mathbf{s}^T}{\mathbf{s}^T \mathbf{q}}.$$

- $\gamma_k \equiv 1$ und

$$\Theta_k = \frac{\mathbf{s}^{(k)T} \mathbf{q}^{(k)}}{\mathbf{s}^{(k)T} \mathbf{q}^{(k)} - \mathbf{q}^{(k)T} \mathbf{H}_k \mathbf{q}^{(k)}}.$$

Wir erhalten das symmetrische Rang-1-Verfahren von BROYDEN

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{\left(\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{q}^{(k)}\right) \left(\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{q}^{(k)}\right)^T}{\mathbf{q}^{(k)T} \left(\mathbf{s}^{(k)} - \mathbf{H}_k \mathbf{q}^{(k)}\right)}.$$

(Hier ist $\Theta_k < 0$ möglich. \mathbf{H}_{k+1} kann dann indefinit werden.)

Wir wollen nun zeigen, unter welchen Bedingungen durch die Update-Formel der OREN-LUENBERGER-Klasse Folgen von positiv definiten Matrizen erzeugt werden.

11.26. Satz: *Es sei \mathbf{H} eine positiv definite Matrix. Für die Vektoren $\mathbf{q} \in \mathbb{R}^n$ und $\mathbf{s} \in \mathbb{R}^n$ gelte $\mathbf{s}^T \mathbf{q} > 0$. Dann ist die Matrix*

$$\begin{aligned} \mathbf{H}_+ &= \Psi(\mathbf{H}, \mathbf{s}, \mathbf{q}; \gamma, \Theta) \\ &= \gamma \mathbf{H} + \left(1 + \gamma \Theta \frac{\mathbf{q}^T \mathbf{H} \mathbf{q}}{\mathbf{s}^T \mathbf{q}}\right) \frac{\mathbf{s} \mathbf{s}^T}{\mathbf{s}^T \mathbf{q}} - \gamma \frac{1 - \Theta}{\mathbf{q}^T \mathbf{H} \mathbf{q}} \mathbf{H} \mathbf{q} \cdot \mathbf{q}^T \mathbf{H} \\ &\quad - \frac{\gamma \Theta}{\mathbf{s}^T \mathbf{q}} \left(\mathbf{s} \cdot \mathbf{q}^T \mathbf{H} + \mathbf{H} \mathbf{q} \cdot \mathbf{s}^T\right) \end{aligned}$$

für beliebige $\gamma > 0$ und $\Theta \geq 0$ ebenfalls positiv definit. Außerdem erfüllt \mathbf{H}_+ die Quasi-NEWTON-Gleichung

$$\mathbf{H}_+ \mathbf{q} = \mathbf{s}.$$

Beweis: Da \mathbf{H} positiv definit ist existiert eine Zerlegung $\mathbf{H} = \mathbf{L} \mathbf{L}^T$ mit einer unteren Dreiecksmatrix \mathbf{L} . Es sei $\mathbf{u} = \mathbf{L}^T \mathbf{y}$ mit einem beliebigen Vektor $\mathbf{y} \in \mathbb{R}^n$ und $\mathbf{y} \neq \mathbf{o}$.

Weiterhin sei $\mathbf{v} = \mathbf{L}^T \mathbf{q}$. Dann gilt

$$\begin{aligned}
\mathbf{y}^T \mathbf{H}_+ \mathbf{y} &= \gamma \mathbf{y}^T \mathbf{H} \mathbf{y} + \left(1 + \gamma \Theta \frac{\mathbf{q}^T \mathbf{H} \mathbf{q}}{\mathbf{s}^T \mathbf{q}}\right) \frac{\mathbf{y}^T \mathbf{s} \mathbf{s}^T \mathbf{y}}{\mathbf{s}^T \mathbf{q}} - \\
&\quad - \gamma \frac{1 - \Theta}{\mathbf{q}^T \mathbf{H} \mathbf{q}} \mathbf{y}^T \mathbf{H} \mathbf{q} \cdot \mathbf{q}^T \mathbf{H} \mathbf{y} - \frac{\gamma \Theta}{\mathbf{s}^T \mathbf{q}} \left(\mathbf{y}^T \mathbf{s} \cdot \mathbf{q}^T \mathbf{H} \mathbf{y} + \mathbf{y}^T \mathbf{H} \mathbf{q} \cdot \mathbf{s}^T \mathbf{y}\right) \\
&= \gamma \mathbf{u}^T \mathbf{u} + \left(1 + \gamma \Theta \frac{\mathbf{v}^T \mathbf{v}}{\mathbf{s}^T \mathbf{q}}\right) \frac{(\mathbf{s}^T \mathbf{y})^2}{\mathbf{s}^T \mathbf{q}} - \gamma \frac{1 - \Theta}{\mathbf{v}^T \mathbf{v}} (\mathbf{u}^T \mathbf{v})^2 \\
&\quad - \frac{2\gamma \Theta}{\mathbf{s}^T \mathbf{q}} (\mathbf{y}^T \mathbf{s})(\mathbf{u}^T \mathbf{v}) \\
&= \gamma \left[\mathbf{u}^T \mathbf{u} - \frac{(\mathbf{u}^T \mathbf{v})^2}{\mathbf{v}^T \mathbf{v}} \right] + \frac{(\mathbf{s}^T \mathbf{y})^2}{\mathbf{s}^T \mathbf{q}} \\
&\quad + \gamma \Theta \left[\frac{(\mathbf{v}^T \mathbf{v})(\mathbf{s}^T \mathbf{y})^2}{(\mathbf{s}^T \mathbf{q})^2} + \frac{(\mathbf{u}^T \mathbf{v})^2}{\mathbf{v}^T \mathbf{v}} - 2 \frac{(\mathbf{y}^T \mathbf{s})(\mathbf{u}^T \mathbf{v})}{\mathbf{s}^T \mathbf{q}} \right] \\
&= \gamma \frac{(\mathbf{u}^T \mathbf{u})(\mathbf{v}^T \mathbf{v}) - (\mathbf{u}^T \mathbf{v})^2}{\mathbf{v}^T \mathbf{v}} + \frac{(\mathbf{s}^T \mathbf{y})^2}{\mathbf{s}^T \mathbf{q}} + \gamma \Theta (\mathbf{v}^T \mathbf{v}) \left[\frac{\mathbf{s}^T \mathbf{y}}{\mathbf{s}^T \mathbf{q}} - \frac{\mathbf{u}^T \mathbf{v}}{\mathbf{v}^T \mathbf{v}} \right]^2.
\end{aligned}$$

Alle drei Terme sind nichtnegativ. Damit gilt zumindest

$$\mathbf{y}^T \mathbf{H}_+ \mathbf{y} \geq 0.$$

Der erste Term ist nach der CAUCHY-SCHWARZschen Ungleichung genau dann Null, wenn die Vektoren \mathbf{u} und \mathbf{v} linear abhängig sind: $\mathbf{u} = \beta \mathbf{v}$. Dann ist aber wegen der Regularität von \mathbf{L} auch $\mathbf{y} = \beta \mathbf{q}$. In diesem Falle folgt

$$\mathbf{y}^T \mathbf{H}_+ \mathbf{y} = \beta^2 (\mathbf{s}^T \mathbf{q}) + \gamma \Theta (\mathbf{v}^T \mathbf{v}) [\beta - \beta] = \beta^2 (\mathbf{s}^T \mathbf{q}) > 0.$$

Damit ist \mathbf{H}_+ positiv definit.

Die Gültigkeit der Quasi-NEWTON-Gleichung folgt durch einfaches Ausrechnen.

$$\begin{aligned}
\mathbf{H}_+ \mathbf{q} &= \gamma \mathbf{H} \mathbf{q} + \left(1 + \gamma \Theta \frac{\mathbf{q}^T \mathbf{H} \mathbf{q}}{\mathbf{s}^T \mathbf{q}}\right) \frac{\mathbf{s} \mathbf{s}^T \mathbf{q}}{\mathbf{s}^T \mathbf{q}} - \gamma \frac{1 - \Theta}{\mathbf{q}^T \mathbf{H} \mathbf{q}} \mathbf{H} \mathbf{q} \cdot \mathbf{q}^T \mathbf{H} \mathbf{q} - \\
&\quad - \frac{\gamma \Theta}{\mathbf{s}^T \mathbf{q}} \left(\mathbf{s} \cdot \mathbf{q}^T \mathbf{H} \mathbf{q} + \mathbf{H} \mathbf{q} \cdot \mathbf{s}^T \mathbf{q}\right) \\
&= \mathbf{H} \mathbf{q} [\gamma - \gamma(1 - \Theta) - \gamma \Theta] + \mathbf{s} \left[1 + \gamma \Theta \frac{\mathbf{q}^T \mathbf{H} \mathbf{q}}{\mathbf{s}^T \mathbf{q}} - \gamma \Theta \frac{\mathbf{q}^T \mathbf{H} \mathbf{q}}{\mathbf{s}^T \mathbf{q}}\right] \\
&= \mathbf{s}
\end{aligned}$$

Bemerkung: Die Bedingung $\mathbf{s}^T \mathbf{q} > 0$ ist eine Forderung an den Schrittweitenalgorithmus. Aus

$$\mathbf{s}^{(k)T} \mathbf{q}^{(k)} = \alpha_k \mathbf{p}^{(k)T} \left[\nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)}) \right] > 0$$

folgt

$$\mathbf{p}^{(k)T} \nabla f(\mathbf{x}^{(k+1)}) > \mathbf{p}^{(k)T} \nabla f(\mathbf{x}^{(k)}).$$

Für exakte Schrittweiten ist dann wegen

$$\mathbf{p}^{(k)T} \nabla f(\mathbf{x}^{(k+1)}) = 0, \quad \mathbf{p}^{(k)T} \nabla f(\mathbf{x}^{(k)}) < 0$$

diese Ungleichung immer erfüllt.

Für POWELL-Schrittweiten gilt wegen $\mathbf{p}^{(k)T} \nabla f(\mathbf{x}^{(k)}) < 0$ und $\nu < 1$ ebenfalls

$$\mathbf{p}^{(k)T} \nabla f(\mathbf{x}^{(k+1)}) \geq \nu \mathbf{p}^{(k)T} \nabla f(\mathbf{x}^{(k)}) > \mathbf{p}^{(k)T} \nabla f(\mathbf{x}^{(k)}).$$

Für ARMIJO-Schrittweiten kann die Bedingung $\mathbf{s}^{(k)T} \mathbf{q}^{(k)} > 0$ verletzt sein. Es können indefinite Matrizen entstehen.

Das Konvergenzverhalten von Verfahren der OREN-LUENBERGER-Klasse klärt der folgende Satz.

11.27. Satz: Mit der positiv definiten Matrix \mathbf{A} , dem Vektor $\mathbf{b} \in \mathbb{R}^n$ und einem $c \in \mathbb{R}$ definieren wir die quadratische Funktion $h : \mathbb{R}^n \rightarrow \mathbb{R}$ gemäß

$$h(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c.$$

Wendet man ein Verfahren der OREN-LUENBERGER-Klasse zum Lösen des Minimumproblems

$$\min_{\mathbf{x} \in \mathbb{R}^n} h(\mathbf{x})$$

an, so gilt für einen beliebigen Startpunkt $\mathbf{x}^{(0)} \in \mathbb{R}^n$ und eine beliebige positiv definite Startmatrix \mathbf{H}_0 bei Anwendung exakter Schrittweiten:

1. Es existiert ein kleinstes $m \leq n$, so dass $\mathbf{x}^{(m)} = \bar{\mathbf{x}} = -\mathbf{A}^{-1} \mathbf{b}$ die exakte Minimumstelle von h mit $\nabla h(\bar{\mathbf{x}}) = \mathbf{0}$ ist.
2. Für $0 \leq i \neq k \leq m-1$ ist

$$\mathbf{s}^{(i)T} \mathbf{q}^{(k)} = \mathbf{s}^{(i)T} \mathbf{A} \mathbf{s}^{(k)}$$

und für $0 \leq i \leq m-1$ ist

$$\mathbf{s}^{(i)T} \mathbf{q}^{(i)} > 0.$$

3. Für $0 \leq i < k \leq m - 1$ ist

$$\mathbf{s}^{(i)T} \nabla h(\mathbf{x}^{(k)}) = \mathbf{o}.$$

4. Für $0 \leq i < k \leq m - 1$ ist

$$\mathbf{H}_k \mathbf{q}^{(i)} = \gamma_{ik} \mathbf{s}^{(i)}$$

mit

$$\gamma_{ik} = \begin{cases} \gamma_{i+1} \gamma_{i+2} \cdots \gamma_{k-1} & \text{für } i < k - 1 \\ 1 & \text{für } i = k - 1 \end{cases}.$$

5. Ist $m = n$, so gilt zusätzlich

$$\mathbf{H}_m = \mathbf{H}_n = \mathbf{S} \mathbf{D} \mathbf{S}^{-1} \mathbf{A}^{-1}$$

mit

$$\mathbf{D} = \text{diag}(\gamma_{0n}, \dots, \gamma_{n-1,n}), \quad \mathbf{S} = \left(\mathbf{s}^{(0)}, \dots, \mathbf{s}^{(n-1)} \right).$$

Für $\gamma_i \equiv 1$ folgt $\mathbf{H}_n = \mathbf{A}^{-1}$.

Einen Beweis findet man in STOER „Numerische Mathematik“ Bd. 1.

Da sich jede hinreichend oft differenzierbare Funktion in der Nähe eines lokalen Minimums beliebig genau durch eine quadratische Funktion approximieren lässt, ist zu erwarten, dass die Quasi-NEWTON-Verfahren der OREN-LUENBERGER-Klasse auch bei der Anwendung auf beliebige Funktionen lokal ein gutes Konvergenzverhalten zeigen.

Es bleibt nun noch die Frage zu klären, wie die Parameter γ_k und Θ_k in jedem Schritt zu wählen sind. Praktische Erfahrungen zeigen, dass man mit dem BFGS-Verfahren ($\gamma_k \equiv 1$ und $\Theta_k \equiv 1$) gute Resultate erzielt. Andererseits folgte aus Satz 11.19, dass die Matrizen \mathbf{H}_k gleichmäßig positiv definit und gleichmäßig beschränkt sind, damit globale Konvergenz eintritt. Man kann nun versuchen, über die Parameter γ_k und Θ_k günstige Eigenschaften der Matrizen \mathbf{H}_k einzustellen.

Von OREN und SPEDICATO (1974) stammt das folgende Ergebnis. Es sei

$$\sigma_k = \mathbf{s}^{(k)T} \mathbf{H}_k^{-1} \mathbf{s}^{(k)},$$

$$\varrho_k = \mathbf{q}^{(k)T} \mathbf{H}_k \mathbf{q}^{(k)},$$

$$\tau_k = \mathbf{s}^{(k)T} \mathbf{q}^{(k)}.$$

Man setze

$$\gamma_k = \frac{\sigma_k}{\tau_k}, \quad \Theta_k = 0 \quad \text{falls} \quad \frac{\sigma_k}{\tau_k} \leq 1,$$

$$\gamma_k = \frac{\tau_k}{\varrho_k}, \quad \Theta_k = 1 \quad \text{falls} \quad \frac{\tau_k}{\varrho_k} \geq 1$$

beziehungsweise

$$\gamma_k = 1, \quad \Theta_k = \tau_k(\sigma_k - \tau_k)(\sigma_k \varrho_k - \tau_k^2) \quad \text{falls} \quad \frac{\tau_k}{\varrho_k} \leq 1 \leq \frac{\sigma_k}{\tau_k}.$$

Durch diese Wahl wird eine obere Schranke für $\text{cond}(\mathbf{H}_{k+1})/\text{cond}(\mathbf{H}_k)$ in jedem Schritt minimiert.

Eine andere Variante stammt von DAVIDON (1975). Er betrachtete Verfahren mit $\gamma_k \equiv 1$ und wählt

$$\Theta_k = \tau_k(\sigma_k - \tau_k)(\sigma_k \varrho_k - \tau_k^2) \quad \text{falls} \quad \tau_k \leq 2 \frac{\sigma_k \varrho_k}{\sigma_k + \varrho_k}$$

beziehungsweise

$$\Theta_k = \frac{\tau_k}{\tau_k - \varrho_k} \quad \text{sonst.}$$

Durch diese Wahl wird der Quotient $\lambda_{\max}/\lambda_{\min}$ bezüglich des allgemeinen Eigenwertproblems

$$\text{Bestimme ein } \lambda \in \mathbb{C} \text{ und ein } \mathbf{y} \in \mathbb{C}^n \text{ mit } \mathbf{y} \neq \mathbf{o} \text{ und } \mathbf{H}_{k+1}\mathbf{y} = \lambda\mathbf{H}_k\mathbf{y}.$$

minimiert.

Wie wir zu Satz 11.26 bemerkten, kann bei Verwendung von ARMIJO-Schrittweiten aus einer positiv definiten Matrix \mathbf{H}_k durch die update-Formel von OREN und LUNENBERGER eine indefinite Matrix \mathbf{H}_{k+1} entstehen. Die Sicherung der Bedingung

$$\nabla f(\mathbf{x}^{(k+1)})^T \mathbf{p}^{(k)} \geq \nu \nabla f(\mathbf{x}^{(k)})^T \mathbf{p}^{(k)},$$

durch die man hinreichend genaue Schrittweiten bekommen würde, ist in der Praxis aber nur mit einem unvertretbar großen Aufwand möglich. Wir wollen nun eine andere Variante betrachten, mit der man automatisch die gleichmäßige Beschränktheit und gleichmäßige positive Definitheit der Matrizen \mathbf{H}_k sichert. Die Abstiegsrichtungen werden in den bisher behandelten Verfahren in der Form

$$(\text{positiv definite Matrix } \mathbf{H}_k) \times (\text{negativer Gradient } -\nabla f(\mathbf{x}^{(k)}))$$

bestimmt. Hätte man statt der Matrix \mathbf{H}_k nur die Inverse $\mathbf{B}_k = \mathbf{H}_k^{-1}$, so ergäbe sich die Abstiegsrichtung $\mathbf{p}^{(k)}$ als Lösung des linearen Gleichungssystems

$$\mathbf{B}_k \mathbf{p}^{(k)} = -\nabla f(\mathbf{x}^{(k)}).$$

Die Matrix \mathbf{B}_k erfüllt die Quasi-NEWTON-Gleichung

$$\mathbf{B}_k \left(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \right) = \nabla f(\mathbf{x}^{(k+1)}) - \nabla f(\mathbf{x}^{(k)}).$$

Für die Matrizen \mathbf{B}_k lässt sich ebenfalls eine zweiparametrische Update-Formel angeben. Man braucht dazu nur die Update-Formel für die \mathbf{H}_k durch mehrmalige Anwendung der SHERMAN-MORRISON-Formel zu invertieren. Es ergibt sich

$$\mathbf{B}_{k+1} = \Psi(\mathbf{B}_k, \mathbf{q}_k, \mathbf{s}_k; \tilde{\gamma}_k, \tilde{\Theta}_k)$$

mit

$$\begin{aligned} \Psi(\mathbf{B}, \mathbf{q}, \mathbf{s}; \tilde{\gamma}, \tilde{\Theta}) &= \tilde{\gamma} \mathbf{B} + \left(1 + \tilde{\gamma} \tilde{\Theta} \frac{\mathbf{s}^T \mathbf{B} \mathbf{s}}{\mathbf{q}^T \mathbf{s}} \right) \frac{\mathbf{q} \mathbf{q}^T}{\mathbf{q}^T \mathbf{s}} \\ &\quad - \tilde{\gamma} \frac{1 - \tilde{\Theta}}{\mathbf{s}^T \mathbf{B} \mathbf{s}} \mathbf{B} \mathbf{s} \cdot \mathbf{s}^T \mathbf{B} \\ &\quad - \frac{\tilde{\gamma} \tilde{\Theta}}{\mathbf{q}^T \mathbf{s}} (\mathbf{q} \cdot \mathbf{s}^T \mathbf{B} + \mathbf{B} \mathbf{s} \cdot \mathbf{q}^T). \end{aligned}$$

Der Zusammenhang zu den entsprechenden Update-Formeln für \mathbf{H} ist durch

$$\tilde{\gamma} = \frac{1}{\gamma}, \quad \tilde{\Theta} = \frac{(1 - \Theta)(\mathbf{s}^T \mathbf{q})^2}{(1 - \Theta)(\mathbf{s}^T \mathbf{q})^2 + \Theta \mathbf{s}^T \mathbf{B} \mathbf{s} \cdot \mathbf{q}^T \mathbf{H} \mathbf{q}}$$

gegeben.

Wir erhalten

- aus dem DFP-Verfahren mit $\gamma_k \equiv 1$ und $\Theta_k \equiv 0$ $\tilde{\gamma}_k \equiv 1$ und $\tilde{\Theta}_k \equiv 1$, also

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \left(1 + \frac{\mathbf{s}^{(k)T} \mathbf{B}_k \mathbf{s}^{(k)}}{\mathbf{q}^{(k)T} \mathbf{s}^{(k)}} \right) \frac{\mathbf{q}^{(k)} \mathbf{q}^{(k)T}}{\mathbf{q}^{(k)T} \mathbf{s}^{(k)}} - \frac{\mathbf{q}^{(k)} \cdot \mathbf{s}^{(k)T} \mathbf{B}_k + \mathbf{B}_k \mathbf{s}^{(k)} \cdot \mathbf{q}^{(k)T}}{\mathbf{q}^{(k)T} \mathbf{s}^{(k)}},$$

- aus dem BFGS-Verfahren mit $\gamma_k \equiv 1$ und $\Theta_k \equiv 1$ $\tilde{\gamma}_k \equiv 1$ und $\tilde{\Theta}_k \equiv 0$, also

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{\mathbf{q}^{(k)} \mathbf{q}^{(k)T}}{\mathbf{q}^{(k)T} \mathbf{s}^{(k)}} - \frac{\mathbf{B}_k \mathbf{s}^{(k)} \cdot \mathbf{s}^{(k)T} \mathbf{B}_k}{\mathbf{s}^{(k)T} \mathbf{B}_k \mathbf{s}^{(k)}},$$

- aus dem symmetrischen Rang-1-Verfahren von BROYDEN mit $\gamma_k \equiv 1$ und

$$\Theta_k = \frac{\mathbf{s}^{(k)T} \mathbf{q}^{(k)}}{\mathbf{s}^{(k)T} \mathbf{q}^{(k)} - \mathbf{q}^{(k)T} \mathbf{H}_k \mathbf{q}^{(k)}}$$

$$\tilde{\gamma}_k \equiv 1 \text{ und}$$

$$\tilde{\Theta}_k = \frac{\mathbf{s}^{(k)T} \mathbf{q}^{(k)}}{\mathbf{s}^{(k)T} \mathbf{q}^{(k)} - \mathbf{q}^{(k)T} \mathbf{B}_k \mathbf{q}^{(k)}},$$

also

$$\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{(\mathbf{q}^{(k)} - \mathbf{B}_k \mathbf{s}^{(k)})(\mathbf{q}^{(k)} - \mathbf{B}_k \mathbf{s}^{(k)})^T}{\mathbf{s}^{(k)T} \mathbf{q}^{(k)} - \mathbf{q}^{(k)T} \mathbf{B}_k \mathbf{q}^{(k)}}.$$

Die Formeln ähneln den Update-Formeln für \mathbf{H}_k . Auf den ersten Blick scheint es aber so, als ob durch diese Vorgehensweise der Rechenaufwand stark ansteigt. Für das Berechnen von

$$\mathbf{p}^{(k)} = -\mathbf{H}_k \nabla f(\mathbf{x}^{(k)})$$

(Matrix \times Vektor) benötigt man rund n^2 Rechenoperationen, für das Lösen von

$$\mathbf{B}_k \mathbf{p}^{(k)} = -\nabla f(\mathbf{x}^{(k)})$$

(lineares Gleichungssystem) dagegen rund $n^3/6$ Rechenoperationen.

Die Matrix \mathbf{B}_{k+1} entsteht durch eine symmetrische Rang-2-Modifikation aus der Matrix $\tilde{\gamma}_k \mathbf{B}_k$, denn es gilt

$$\text{rg}(\mathbf{B}_{k+1} - \tilde{\gamma}_k \mathbf{B}_k) \leq 2.$$

Daher lässt sich der Aufwand zum Lösen der Gleichungssysteme beträchtlich verringern. Kennt man nämlich eine Zerlegung $\mathbf{B}_k = \mathbf{L}_k \mathbf{D}_k \mathbf{L}_k^T$, so lässt sich die entsprechende Zerlegung von \mathbf{B}_{k+1} mit einem Aufwand der Größenordnung n^2 berechnen. Dann ist der Aufwand zum Berechnen einer Abstiegsrichtung mittels \mathbf{B}_k von gleicher Größenordnung wie der Aufwand zum Berechnen einer Abstiegsrichtung mittels \mathbf{H}_k . Der große Vorteil der Verwendung einer LDL^T -Zerlegung der Matrix \mathbf{B}_k besteht darin, dass aus der Zerlegung immer die positive Definitheit von \mathbf{B}_k abzulesen ist. Denn $\mathbf{B} = \mathbf{L} \mathbf{D} \mathbf{L}^T$ ist genau dann positiv definit, wenn für die Diagonalmatrix $\mathbf{D} = \text{diag}(d_1, \dots, d_n)$ $d_i > 0$ für $i = 1, \dots, n$ gilt. Ist diese Bedingung irgendwann verletzt, so lässt sich die positive Definitheit der Matrix \mathbf{B}_{k+1} erzwingen, indem man die entsprechenden Elemente von \mathbf{D}_{k+1} durch positive Zahlen ersetzt. Dann ist zwar die Quasi-NEWTON-Gleichung nicht mehr erfüllt; dies lässt sich

aber als Neustart des Verfahrens mit dem aktuellen Iterationspunkt als $\boldsymbol{x}^{(0)}$ und der korrigierten Matrix \boldsymbol{B}_{k+1} als positiv definiten Startmatrix \boldsymbol{B}_0 interpretieren. Andererseits wird man von dem Verfahren verlangen, dass für hinreichend große k sich Schrittweiten $\alpha_k = 1$ einstellen. Dann lässt sich aber auch annehmen, dass die zweite Bedingung aus Satz 11.19 erfüllt und damit die positive Definitheit der Matrix \boldsymbol{B}_{k+1} gesichert ist. Die Eigenschaft $\alpha_k = 1$ für alle $k \geq k_0$ ist für jeden konkreten Algorithmus gesondert nachzuweisen.

11.3.3. Algorithmen zur Aufdatierung von Zerlegungen

Wir unterscheiden die Fälle $\bar{\boldsymbol{B}} = \boldsymbol{B} + \boldsymbol{v}\boldsymbol{v}^T$ oder $\bar{\boldsymbol{B}} = \boldsymbol{B} - \boldsymbol{v}\boldsymbol{v}^T$ und erhalten die folgenden beiden Algorithmen.

11.28. Aufdatierung einer LDL^T -Zerlegung $\bar{\boldsymbol{B}} = \boldsymbol{B} + \boldsymbol{v}\boldsymbol{v}^T$:

Von der positiv definiten (n, n) -Matrix \boldsymbol{B} sei eine Zerlegung

$$\boldsymbol{B} = \boldsymbol{L}\boldsymbol{D}\boldsymbol{L}^T$$

mit einer unteren (n, n) -Einsdreiecksmatrix \boldsymbol{L} und einer (n, n) -Diagonalmatrix \boldsymbol{D} bekannt.

Es sind eine untere Einsdreiecksmatrix $\bar{\boldsymbol{L}}$ und eine Diagonalmatrix $\bar{\boldsymbol{D}}$ derart zu berechnen, dass

$$\bar{\boldsymbol{B}} = \boldsymbol{B} + \boldsymbol{v}\boldsymbol{v}^T = \bar{\boldsymbol{L}}\bar{\boldsymbol{D}}\bar{\boldsymbol{L}}^T$$

gilt.

```

 $t_0 = 1; \boldsymbol{v}^{(1)} = \boldsymbol{v}$ 
for  $j = 1$  to  $n$  do
   $u_j = v_j^{(j)}$ 
   $t_j = t_{j-1} + u_j^2/d_j$ 
   $\bar{d}_j = d_j t_j / t_{j-1}$ 
   $\beta_j = u_j / (d_j t_j)$ 
  for  $k = j + 1$  to  $n$  do
     $v_k^{(j+1)} = v_k^{(j)} - u_j l_{kj}$ 
     $\bar{l}_{kj} = l_{kj} + \beta_j v_k^{(j+1)}$ 
  endfor
endfor

```

11.29. Aufdatierung einer LDL^T -Zerlegung $\bar{\boldsymbol{B}} = \boldsymbol{B} - \boldsymbol{v}\boldsymbol{v}^T$:

Von der positiv definiten (n, n) -Matrix \mathbf{B} sei eine Zerlegung $\mathbf{B} = \mathbf{L}\mathbf{D}\mathbf{L}^T$ mit einer unteren Einsdreiecksmatrix \mathbf{L} und einer Diagonalmatrix \mathbf{D} bekannt.

Es sind eine untere Einsdreiecksmatrix $\bar{\mathbf{L}}$ und eine Diagonalmatrix $\bar{\mathbf{D}}$ derart zu berechnen, dass

$$\bar{\mathbf{B}} = \mathbf{B} - \mathbf{v}\mathbf{v}^T = \bar{\mathbf{L}}\bar{\mathbf{D}}\bar{\mathbf{L}}^T$$

gilt.

Wähle ein $\delta > 0$

Löse das Gleichungssystem $\mathbf{L}\mathbf{u} = \mathbf{v}$

$$t_{n+1} = 1 - \mathbf{u}^T \mathbf{D}^{-1} \mathbf{u}$$

if $t_{n+1} < \delta$ **then**

$$t_{n+1} = \delta$$

endif

for $j = n$ **to** 1 **step** -1 **do**

$$t_j = t_{j+1} + u_j^2 / d_j$$

$$\bar{d}_j = d_j t_{j+1} / t_j$$

$$\beta_j = -u_j / (d_j t_{j+1})$$

$$v_j^{(j)} = u_j$$

for $k = j + 1$ **to** n **do**

$$v_k^{(j+1)} = v_k^{(j)} - u_j l_{kj}$$

$$\bar{l}_{kj} = l_{kj} + \beta_j v_k^{(j+1)}$$

endfor

endfor

Beim ersten Algorithmus gilt wegen $t_j \geq t_{j-1}$ $\bar{d}_j \geq d_j$. Die positive Definitheit der Matrix $\bar{\mathbf{B}}$ wird gegenüber der Matrix \mathbf{B} verbessert. Beim zweiten Algorithmus gilt wegen $t_j \leq t_{j-1}$ $\bar{d}_j \leq d_j$. Da aber immer $t_j \geq \delta > 0$ ist die positive Definitheit der Matrix $\bar{\mathbf{B}}$ gesichert.

Falls man im Algorithmus die positive Definitheit zu erzwingen hat, ist die berechnete Zerlegung nur noch Zerlegung einer benachbarten Matrix $\tilde{\mathbf{B}} = \bar{\mathbf{B}} + \delta \mathbf{B}$. Im Algorithmus erkennt man, dass in die Berechnung der Elemente von $\bar{\mathbf{D}}$ und $\bar{\mathbf{L}}$ nur die Elemente von \mathbf{D} , \mathbf{L} und \mathbf{v} eingehen. Ändert man irgendwelche Elemente, speziell die Elemente von $\bar{\mathbf{D}}$, ab, so hat das keinen Einfluss auf die Berechnung der restlichen Elemente. Damit ist es unproblematisch, Beziehungen der Form

$$0 < \delta \leq \min_{i=1, \dots, n} d_i \leq \max_{i=1, \dots, n} d_i \leq \Delta < \infty$$

und

$$\max_{i, j=1, \dots, n} |l_{ij}| \leq \Lambda < \infty$$

zu sichern. Hiermit sichert man aber auch die gleichmäßige Beschränktheit und gleichmäßige positive Definitheit der Matrizen B_k .

11.4. Trust-Region-Verfahren

Dieser Verfahrensgruppe liegt eine andere Idee zugrunde als jene, die zu den Quasi-NEWTON-Verfahren und NEWTON-Verfahren führte. Wir ersetzen das Minimumproblem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

lokal, d. h. in einer Kugelumgebung des aktuellen Iterationspunktes $\mathbf{x}^{(k)}$, durch ein einfacheres, leicht zu lösendes Modell. Wird mit der Lösung dieses vereinfachten Problems eine hinreichende Verminderung der Zielfunktion erreicht, so wird diese als neuer Iterationspunkt $\mathbf{x}^{(k+1)}$ akzeptiert. Erreicht man mit der Lösung des Ersatzproblems keine hinreichende Verminderung der Zielfunktion, so bedeutet das, dass man dem Modell auf einer zu großen Umgebung von $\mathbf{x}^{(k)}$ vertraut hat.⁴ Man wird die Umgebung verkleinern und das Ersatzproblem erneut lösen. Je nach Wahl des Hilfsproblems und der Norm zur Definition der Kugelumgebungen von $\mathbf{x}^{(k)}$ erhält man die verschiedensten Trust-Region-Verfahren. Als vereinfachte Probleme bieten sich die folgenden an.

- $f \in C^1(\mathbb{R})$:

$$f_{\mathbf{x}}(\mathbf{p}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{p}.$$

- $f \in C^2(\mathbb{R})$:

$$f_{\mathbf{x}}(\mathbf{p}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}) \mathbf{p}$$

oder

$$f_{\mathbf{x}}(\mathbf{p}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{B} \mathbf{p}$$

mit einer symmetrischen Matrix \mathbf{B} .

- $f(\mathbf{x}) = \|F(\mathbf{x})\|$ mit $F : \mathbb{R}^n \longrightarrow \mathbb{R}^m$:

$$f_{\mathbf{x}}(\mathbf{p}) = \|F(\mathbf{x}) + F'(\mathbf{x})\mathbf{p}\|.$$

⁴Daher stammt auch der Name der Verfahren. Den englischen Begriff „trust region“ würde man mit „Vertrauensbereich“ übersetzen.

Es sei nun $f_x(\mathbf{p})$ für jedes $\mathbf{x} \in \mathbb{R}^n$ gegeben. Eine einfache Version eines Trust-Region-Verfahrens könnte folgendermaßen aussehen.

11.30. Trust-region-Verfahren:

S0 Wähle Konstanten $0 < \varrho_1 < \varrho_2 < 1$, $0 < \sigma_1 < 1 < \sigma_2$ und $\Delta_0 > 0$.

Wähle einen Startpunkt $\mathbf{x}^{(0)} \in \mathbb{R}^n$ und setze $k = 0$.

S1 Berechne die Lösung \mathbf{p}^* der Aufgabe

$$\min_{\|\mathbf{p}\| \leq \Delta_k} f_{\mathbf{x}^{(k)}}(\mathbf{p}).$$

S2 Falls $f(\mathbf{x}^{(k)}) = f_{\mathbf{x}^{(k)}}(\mathbf{p}^*)$ STOPP

(Bei vernünftiger Wahl von f_x ist dann $\mathbf{x}^{(k)}$ zumindest stationäre Lösung es ursprünglichen Minimumproblems.)

S3 Berechne

$$r = \frac{f(\mathbf{x}^{(k)}) - f(\mathbf{x}^{(k)} + \mathbf{p}^*)}{f(\mathbf{x}^{(k)}) - f_{\mathbf{x}^{(k)}}(\mathbf{p}^*)}.$$

Falls $r \geq \varrho_1$, so $\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} + \mathbf{p}^*$.

Falls $r < \varrho_1$, wähle ein $\Delta_{k+1} \in (0, \sigma_1 \Delta_k]$.

Falls $\varrho_1 \leq r < \varrho_2$, wähle ein $\Delta_{k+1} \in [\sigma_1 \Delta_k, \Delta_k]$

Falls $\varrho_2 \leq r$, wähle ein $\Delta_{k+1} \in [\Delta_k, \sigma_2 \Delta_k]$

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

Bemerkung: Wird im Schritt **S2** nicht abgebrochen, so ist $f_{\mathbf{x}^{(k)}}(\mathbf{p}^*) < f(\mathbf{x}^{(k)})$. Die Größe r drückt damit aus, wie gut das Modell mit der tatsächlichen Funktion übereinstimmt. Je näher r an 1 liegt, desto genauer ist das Modell. Ist nun $r \geq \varrho_1$, so stimmen die tatsächliche Veränderung der Zielfunktion und die durch das Modell vorausgesagte Verminderung der Zielfunktion hinreichend gut überein. Wir können $\mathbf{x}^{(k)} + \mathbf{p}^*$ als neuen Iterationspunkt akzeptieren. Ist sogar $r \geq \varrho_2$, so stimmt das Modell so gut, dass wir im nächsten Schritt den Vertrauensbereich vergrößern dürfen. Ist dagegen $r < \varrho_1$, so war das Modell zu ungenau oder der Bereich, auf dem wir es angewendet haben, zu groß. Wir müssen den Vertrauensbereich verkleinern und das Modellproblem auf dem verkleinerten Bereich noch einmal lösen.

11.31. Beispiel: Für $f \in C^1(\mathbb{R})$ ergibt sich das einfachste Trust-region-Verfahren, falls man die Modellfunktion $f_x(\mathbf{p}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{p}$ und zur Festlegung der Kugelumgebung die euklidische Norm verwendet. Wir lösen damit in jedem Schritt das Hilfsproblem

$$\min_{\|\mathbf{p}\|_2 \leq \Delta} f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{p}.$$

Ist \mathbf{p}^* eine Lösung dieses Problems und gilt $f_x(\mathbf{p}^*) = f(\mathbf{x})$, so ist $\nabla f(\mathbf{x})^T \mathbf{p} \geq 0$ für alle \mathbf{p} mit $\|\mathbf{p}\|_2 \leq \Delta$. Wählt man speziell $\mathbf{p} = -\alpha \nabla f(\mathbf{x})$ mit einem $\alpha > 0$, so dass $\|\mathbf{p}\|_2 \leq \Delta$, so folgt $\|\nabla f(\mathbf{x})\|_2 \leq 0$ und weiter $\nabla f(\mathbf{x}) = \mathbf{o}$. Ist also für die Modellfunktion der Test im Schritt **S2** des Verfahrens erfüllt, so ist \mathbf{x} stationäre Lösung des Minimumproblems. Für $\nabla f(\mathbf{x}) \neq \mathbf{o}$ ergibt sich die Lösung des Hilfsproblems zu

$$\mathbf{p}^* = -\Delta \frac{\nabla f(\mathbf{x})}{\|\nabla f(\mathbf{x})\|_2}.$$



Betrachten wir nun glatte Zielfunktionen $f \in C^2(\mathbb{R})$. Für derartige Probleme hatten wir das NEWTON-Verfahren angewendet. Das Verfahren hat aber den Nachteil, dass die NEWTON-Richtung

$$\mathbf{p} = -\left(\nabla^2 f(\mathbf{x})\right)^{-1} \nabla f(\mathbf{x})$$

nicht immer eine Abstiegsrichtung ist. Das Trust-Region-Verfahren ist auch im Falle einer indefiniten HESSE-Matrix noch durchführbar. Wir verwenden dazu die Ersatzfunktion

$$f_x(\mathbf{p}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}) \mathbf{p}.$$

Das Hilfsproblem

$$(P_{x,\Delta}) \quad \min_{\|\mathbf{p}\|_2 \leq \Delta} f_x(\mathbf{p}) = \min_{\|\mathbf{p}\|_2 \leq \Delta} \left(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}) \mathbf{p} \right)$$

besitzt, da f_x stetig ist, stets eine globale Lösung \mathbf{p}^* . Falls $\nabla^2 f(\mathbf{x})$ indefinit ist, ist f_x nicht mehr gleichmäßig konvex. In diesem Falle können auch zusätzliche lokale Lösungen auftreten.

Der folgende Satz gibt notwendige und hinreichende Bedingungen für ein globales Minimum des Hilfsproblems $(P_{x,\Delta})$ an.

11.32. Satz: *Wir betrachten die Aufgabe*

$$(P_\Delta) \quad \min_{\|\mathbf{p}\|_2 \leq \Delta} \varphi(\mathbf{p}), \quad \varphi(\mathbf{p}) = f + \mathbf{g}^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \mathbf{B} \mathbf{p},$$

wobei $\Delta > 0$, $f \in \mathbb{R}$, $\mathbf{g} \in \mathbb{R}^n$ und $\mathbf{B} \in \mathbb{R}^{n \times n}$ gelte. Die Matrix \mathbf{B} sei darüber hinaus symmetrisch. Dann ist $\mathbf{p}^* \in \mathbb{R}^n$ mit $\|\mathbf{p}^*\|_2 \leq \Delta$ genau dann eine globale Lösung von (P_Δ) , wenn ein $\lambda^* \geq 0$ existiert, so dass

1. $(\mathbf{B} + \lambda^* \mathbf{I}) \mathbf{p}^* = -\mathbf{g}$,
2. $\lambda^* (\|\mathbf{p}^*\|_2 - \Delta) = 0$,
3. $\mathbf{B} + \lambda^* \mathbf{I}$ ist positiv semidefinit

gilt. Ist $\mathbf{B} + \lambda^ \mathbf{I}$ sogar positiv definit, so ist \mathbf{p}^* eindeutige globale Lösung von (P_Δ) . Weiterhin ist $\varphi(\mathbf{p}^*) = f$ genau dann, wenn $\mathbf{g} = \mathbf{o}$ und \mathbf{B} positiv semidefinit ist. Es gilt die Abschätzung*

$$f - \varphi(\mathbf{p}^*) \geq \frac{1}{2} \|\mathbf{g}\|_2 \min \left\{ \Delta, \frac{\|\mathbf{g}\|_2}{\|\mathbf{B}\|_2} \right\}.$$

Ein Trust-Region-Verfahren mit dem Hilfsproblem $(P_{\mathbf{x}, \Delta})$ bricht nach Satz 11.32 genau dann in Schritt **S2** ab, wenn $\nabla f(\mathbf{x}^{(k)}) = \mathbf{0}$ und $\nabla^2 f(\mathbf{x}^{(k)})$ positiv semidefinit ist. Damit sind aber im Punkt $\mathbf{x}^{(k)}$ die Optimalitätsbedingungen zweiter Ordnung aus Satz 11.8 erfüllt.

Zum Abschluss wollen wir noch einen Konvergenzsatz für die Trust-Region-Variante des NEWTON-Verfahrens angeben.

11.33. Satz: *Gegeben sei das freie Minimumproblem*

$$(P) \quad \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x})$$

und ein Startpunkt $\mathbf{x}^{(0)} \in \mathbb{R}^n$. Die Niveaumenge

$$L_0 = \left\{ \mathbf{x} \in \mathbb{R}^n \mid f(\mathbf{x}) \leq f(\mathbf{x}^{(0)}) \right\}$$

sei kompakt. Die Zielfunktion sei auf einer offenen Obermenge von L_0 zweimal stetig differenzierbar und der Gradient von f sei auf L_0 lipschitzstetig. Wir betrachten ein Trust-Region-Verfahren mit dem Hilfsproblem

$$(P_{\mathbf{x}, \Delta}) \quad \min_{\|\mathbf{p}\|_2 \leq \Delta} f_{\mathbf{x}}(\mathbf{p}) = \min_{\|\mathbf{p}\|_2 \leq \Delta} \left(f(\mathbf{x}) + \nabla f(\mathbf{x})^T \mathbf{p} + \frac{1}{2} \mathbf{p}^T \nabla^2 f(\mathbf{x}) \mathbf{p} \right).$$

Falls das Verfahren nicht vorzeitig abbricht, liefert es eine Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ mit

1. *Jeder Häufungspunkt von $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ ist stationäre Lösung von (P) .*
2. *Die Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ besitzt mindestens einen Häufungspunkt in dem die Optimalitätsbedingungen zweiter Ordnung erfüllt sind.*
3. *Ist der Vektor \mathbf{x}^* Häufungspunkt von $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ mit positiv definiten HESSE-Matrix $\nabla^2 f(\mathbf{x}^*)$, so konvergiert die Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ gegen \mathbf{x}^* . Ist darüber hinaus $\nabla^2 f(\mathbf{o})$ auf einer Kugelumgebung von \mathbf{x}^* lipschitzstetig, so konvergiert die Folge $\{\mathbf{x}^{(k)}\}_{k \in \mathbb{N}}$ mindestens quadratisch.*

Index

Ähnlichkeitstransformation, 3
GERSCHGORIN-Kreis, 12
GIVENS-Rotation, 26

Abstiegsrichtung, 120
Ausgleichsproblem, 88
 lineares, 89

Eigenraum, 2
Eigenvektor, 1
Eigenwert, 1
Eigenwerte
 dominante, 34

Funktion
 gleichmäßig konvexe, 131
 richtungsdifferenzierbare, 122
 stetig differenzierbare, 120
Funktionalmatrix, 121

GATEAUX-Ableitung, 122
Gradient, 121
Gradientenverfahren, 156

HESSE-Matrix, 127

Lösung
 globale, 119
 lokale, 119
 stationäre, 126

Matrix
 ähnliche, 3
 defektive, 2
 diagonalähnliche, 3
 diagonalisierbar, 15
 pseudoinverse, 92
 spaltenreguläre, 89

Normalform
 JORDANSche, 10
Normalgleichung, 88

Polynom
 charakteristisches, 1
POWELL-Schrittweite, 139
Punkt
 stationärer, 126
Quasi-NEWTON-Gleichung, 156
Richtungsableitung, 122
Schrittweite, 120
 exakte, 138
Teilraum
 dominanter, 34
TSCHEBYSCHJEFF-Problem
 diskretes, 88
Vielfachheit
 algebraische, 1
 geometrische, 2
Winkel, 40