

Numerik

Hinweis : Mit der Angabe „Buch“ ist das Buch „Numerische Mathematik I“ von Peter Deuffhard gemeint.

TODO :

Pivotstrategien : Beispiel aus dem Buch S.10/11 einfügen ?

Nachiteration Kapitel 1 : Beispiel rechnen für Genauigkeitsverbesserung

Grundlagen

Zeichen :

\doteq „bis auf Terme niedrigerer Ordnung“

Def. „hinreichende Bedingung“ , „notwendige Bedingung“ :

Bsp. : $A : n$ ist durch 4 teilbar $B : n$ ist durch 2 teilbar

$A \Rightarrow B \Leftrightarrow A$ ist hinreichende Bedingung für B

$\Leftrightarrow B$ ist notwendige Bedingung für A

Def. überbestimmtes / unterbestimmtes Gleichungssystem :

Ein Gleichungssystem mit m Gleichungen und n Unbekannten heißt unterbestimmt wenn $m < n$, quadratisch, wenn $m = n$ und überbestimmt, wenn $m > n$ ist.

Def. Matrix : $A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$

Bei Operationen auf den Zeilen einer Matrix (z.B. Elimination) Multiplikation von links.

Bei Operationen auf den Spalten einer Matrix Multiplikation von rechts.

Def. reguläre Matrix : $A \in R^{(n \times n)}$ ist regulär $\Leftrightarrow \text{Rang}(A) \neq 0 \Leftrightarrow \det(A) \neq 0$

Def. singuläre Matrix : $A \in R^{(n \times n)}$ ist singulär $\Leftrightarrow \text{Rang}(A) = 0 \Leftrightarrow \det(A) = 0$

\Leftrightarrow Gleichungssystem ist nicht eindeutig lösbar

(Eine nicht-singuläre Matrix nennt man reguläre Matrix.)

(Eine singuläre Matrix besitzt keine Inverse.)

Def. normale Matrix : Die Matrix A heißt normal $\Leftrightarrow A\bar{A}^T = \bar{A}^T A$

Def. unipotente Matrix : Diagonalelemente gleich eins

Def. positiv definite Matrix :

Eine symmetrische Matrix $A = A^T$ ist genau dann positiv definit, falls

$\langle x, Ax \rangle > 0$ für alle $x \neq 0$ (\Leftrightarrow alle Hauptunterdeterminanten sind positiv)

Def. strikt diagonaldominante Matrix :

$$|a_{ii}| > \sum_{\substack{j=1 \\ i \neq j}}^n |a_{ij}| \quad \text{für } i = 1, \dots, n$$

Def. konjugierte (ähnliche) Matrizen :

Zwei quadratische Matrizen A, B heißen konjugiert (ähnlich), wenn eine reguläre Matrix T existiert mit $B = TAT^{-1}$.

Def. Jacobi Matrix (bzw. Funktionalmatrix):

Für Funktionen $f : R^n \rightarrow R^m$ fasst man die partiellen Ableitungen in einer Matrix $J_f(x)$ zusammen.

Es gilt $J_f(x)[i, j] := \frac{\partial f_i(x)}{\partial x_j}$, d.h. $J_f(x) := \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{pmatrix}$

Für den Spezialfall $f : R^n \rightarrow R$ ergibt sich ein Zeilenvektor. Den zugehörigen transponierten (Spalten-)Vektor nennt man **Gradient** von f an der Stelle x ($grad f(x)$)

$$grad f(x) := \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix}$$

Def. Hessenberg-Matrix :

$$A = \begin{bmatrix} * & \dots & \dots & \dots & * \\ * & \ddots & & & \vdots \\ 0 & \ddots & \ddots & & \vdots \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & * & * \end{bmatrix}$$

Def. Vektornorm :

Eine Abbildung $\|\cdot\| : R^n \rightarrow R$ heißt Vektornorm auf R^n , falls :

- 1.) **Positivität**
 $\|x\| > 0$ für $x \in R^n$
 $\|x\| = 0 \Leftrightarrow x = 0$
- 2.) **Homogenität**
 $\|a \cdot x\| = |a| \cdot \|x\|$ ($a \in R, x \in R^n$)
- 3.) **Dreiecksungleichung**
 $\|x + y\| \leq \|x\| + \|y\|$

Vektornorm Beispiele :

$$\|x\|_2 := \sqrt{\sum_{i=1}^n x_i^2} \quad - \text{euklidische Norm}$$

$$\|x\|_1 := \sum_{i=1}^n |x_i| \quad - \text{1-Norm}$$

$$\|x\|_\infty := \max_{i=1, \dots, n} |x_i| \quad - \text{Maximumnorm}$$

Def. Matrixnorm :

Definition wie bei Vektornorm

Durch $\|A\|_p = \sup_{x \neq 0} \frac{\|Ax\|_p}{\|x\|_p} \quad (A \in \mathbb{R}^{n \times m})$ ist eine mit der Vektornorm $\|\cdot\|$ verträgliche Matrixnorm

definiert, die der Vektornorm $\|\cdot\|$ zugeordnete (induzierte) Matrixnorm.

Eine Matrixnorm heißt **submultiplikativ**, falls

$$\|AB\| \leq \|A\| \cdot \|B\| \quad (\text{z.B. Frobeniusnorm})$$

Matrixnorm Beispiele :

- Zeilensummennorm $\|A\|_\infty = \max_{i=1, \dots, m} \sum_{j=1}^n |a_{ij}|$
- Spaltensummennorm $\|A\|_1 = \max_{j=1, \dots, n} \sum_{i=1}^m |a_{ij}|$
- Spektralnorm $\|A\|_2 = \sqrt{\mathbf{r}(A^T A)}$
- Frobeniusnorm $\|A\|_F = \sqrt{\left(\sum_{i=1}^n \sum_{j=1}^n a_{ij}^2 \right)}$

Def. Spektralradius :

$$\mathbf{r}(A) := \max_{k=1, \dots, n} |\lambda_k(A)| \quad = \text{betragsmäßig größter Eigenwert der Matrix}$$

Def. Nullraum einer Matrix :

Für $A \in \mathbb{R}^{(m \times n)}$ heißt der Lösungsraum des homogenen linearen Gleichungssystems $Ax = 0$ der Nullraum oder der Kern von A .

Def. Spektrum einer Matrix :

Sei A eine (n, n) -Matrix. Dann heißt die Menge aller Eigenwerte von A das Spektrum von A .

1. Lineare Gleichungssysteme

- LGS $Ax = b$ ($x = A^{-1}b$) lässt sich für $\det A \neq 0$ mit der Cramerschen Regel berechnen :

$$x_j = \frac{1}{\det A} \cdot \begin{pmatrix} a_{11} & \dots & a_{1,j-1} & b_1 & a_{1,j+1} & \dots & a_{1,n} \\ a_{21} & \dots & a_{2,j-1} & b_2 & a_{2,j+1} & \dots & a_{2,n} \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ a_{n1} & \dots & a_{n,j-1} & b_n & a_{n,j+1} & \dots & a_{n,n} \end{pmatrix}, \text{ d.h. ersetze } j\text{-te Spalte von } A \text{ durch } b$$

Def. Determinante (nach Leibniz)

$$\det A = \sum_{p \in S_n} \text{sgn } p \cdot a_{1,p(1)} \cdot \dots \cdot a_{n,p(n)}$$

⇒ Aufwand zur Berechnung von $\det A$ beträgt $n \cdot n!$ Operationen.

Auch mit der rekursiven Bestimmung über Unterdeterminanten nach dem Laplaceschen Entwicklungssatz sind 2^n Operationen auszuführen.

gestaffeltes Gleichungssystem

$$\begin{aligned} r_{11}x_1 + r_{12}x_2 + \dots + r_{1n}x_n &= z_1 \\ r_{22}x_2 + \dots + r_{2n}x_n &= z_2 \\ \vdots & \\ r_{nn}x_n &= z_n \end{aligned} \quad , \text{d.h. } Rx = z$$

Auflösung durch Rückwärtssubstitution

$$x_n = z_n / r_{nn}$$

$$x_{n-1} = (z_{n-1} - r_{n-1,n}x_n) / r_{n-1,n-1}$$

usw.

Rechenaufwand :

für die i-te Zeile je n-i Additionen und Multiplikationen sowie eine Division

$$\text{D.h. } \sum_{i=1}^n (i-1) = \frac{n(n-1)}{2} \doteq \frac{n^2}{2} \text{ Multiplikationen und ebenso viele Additionen.}$$

1.1 Gaußsche Eliminationsmethode (LR-Zerlegung)

Ziel / Motivation : Umformung eines LGS in ein gestaffeltes LGS

$$A \rightarrow A^{(1)} \rightarrow A^{(2)} \dots \rightarrow A^{(n)} =: R, \text{ wobei}$$

$$A^{(k)} = \begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & \dots & \dots & a_{1n}^{(1)} \\ & a_{22}^{(2)} & \dots & \dots & \dots & a_{2n}^{(2)} \\ & & \ddots & & & \vdots \\ & & & a_{kk}^{(k)} & \dots & a_{kn}^{(k)} \\ & & & \vdots & & \vdots \\ & & & a_{nk}^{(k)} & \dots & a_{nn}^{(k)} \end{pmatrix}$$

mit einer $(n-k+1, n-k+1)$ Restmatrix rechts unten und den Pivotelementen $a_{kk}^{(k)}$

Algorithmus LR-Zerlegung :

$$\begin{aligned}
 l_{ik} &:= a_{ik}^{(k)} / a_{kk}^{(k)} && \text{für } i = k + 1, \dots, n \\
 a_{ij}^{(k+1)} &:= a_{ij}^{(k)} - l_{ik} \cdot a_{kj}^{(k)} \\
 b_i^{(k+1)} &:= b_i^{(k)} - l_{ik} \cdot b_k^{(k)}
 \end{aligned}$$

Da jeder Eliminationsschritt eine lineare Operation auf den Zeilen von A ist, lässt sich der Übergang von $A^{(k)}$ und $b^{(k)}$ zu $A^{(k+1)}$ und $b^{(k+1)}$ als Multiplikation mit deiner Matrix L_k (Frobenius-Matrix) von links darstellen.

Def. Frobenius-Matrix

$$\begin{pmatrix}
 1 & & & & 0 \\
 & 1 & & & \\
 & & \dots & & \\
 & & & 1 & \\
 & & & b_{i+1,i} & \dots \\
 & & & \dots & 1 \\
 0 & & & b_{n,i} & 1
 \end{pmatrix}$$

besondere Eigenschaft : L_k^{-1} entsteht aus L_k durch Vorzeichenwechsel der l_{ik}

$$Ax = b \quad \Rightarrow \quad Rx = z \quad \text{mit } R = L^{-1}A \quad \text{und } z = L^{-1}b$$

Algorithmus der Gauß-Elimination (mit Aufwandsangabe)

- | | |
|-------------------------------------|---|
| 1. $A = LR$ - Dreieckszerlegung | $\sum_{k=1}^{n-1} k^2 \doteq \frac{n^3}{3}$ |
| 2. $Lz = b$ - Vorwärtssubstitution | $\sum_{k=1}^{n-1} k \doteq \frac{n^2}{2}$ |
| 3. $Rx = z$ - Rückwärtssubstitution | $\sum_{k=1}^{n-1} k \doteq \frac{n^2}{2}$ |

Der Hauptaufwand besteht also in der LR-Zerlegung, die für verschiedene rechte Seiten b_1, \dots, b_j aber nur einmal durchgeführt werden muss.

Beispiel für Gaußsches Eliminationsverfahren / LR-Zerlegung :

Gegeben : $Ax = b$ hier : $\begin{pmatrix} 1 & 4 \\ 2 & 2 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$

1.) LR-Zerlegung

$$L = \begin{pmatrix} 1 & 0 \\ l_{21} & 1 \end{pmatrix}, \quad l_{ik} := \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} \quad \Rightarrow \quad l_{21} = \frac{a_{21}}{a_{11}} = 2 \quad \Rightarrow \quad L = \begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix}$$

$$R = L^{-1}A \quad \Rightarrow \quad R = \begin{pmatrix} 1 & 0 \\ -2 & 1 \end{pmatrix} \cdot \begin{pmatrix} 1 & 4 \\ 2 & 2 \end{pmatrix} = \begin{pmatrix} 1 & 4 \\ 0 & -6 \end{pmatrix}$$

2.) $Lz = b$ - Vorwärtssubstitution

$$\begin{array}{l|l} 1 & 0 \mid 2 & \rightarrow z_1 = 2 \\ 2 & 1 \mid 2 & \rightarrow z_2 = 2 - 4 = -2 \end{array}$$

3.) $Rx = z$ - Rückwärtssubstitution

$$\begin{array}{l|l} 1 & 4 \mid 2 & \rightarrow x_1 = 2 - 4 \cdot 1/3 = 2/3 \\ 0 & -6 \mid -2 & \rightarrow x_2 = 1/3 \end{array}$$

$$, \text{Lsg. : } \begin{pmatrix} 1 & 4 \\ 2 & 2 \end{pmatrix} \begin{pmatrix} 2/3 \\ 1/3 \end{pmatrix} = \begin{pmatrix} 2 \\ 2 \end{pmatrix}$$

1.2 Pivot-Strategien und Nachiteration

Motivation :

Wie man bereits an dem einfachen Beispiel

$$A = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \det A = -1, a_{11} = 0$$

sieht, gibt es Fälle, bei denen die Dreieckszerlegung versagt, obwohl $\det A \neq 0$. Eine Vertauschung der Zeilen führt jedoch zur denkbar einfachsten LR-Zerlegung

$$\bar{A} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = I = LR \quad \text{mit } L = R = I$$

- Bei der rechnerischen Realisierung des Gaußschen Eliminationsverfahrens können Schwierigkeiten nicht nur bei verschwindenden, sondern auch bei „zu kleinen“ Pivotelementen entstehen.

⇒ **Spaltenpivotstrategie**

Bei der Gauß-Elimination wählt man diejenige Zeile als Pivotzeile, die das betragsmäßig größte Element in der Pivotspalte besitzt.

Algorithmus Gauß-Elimination mit Spaltenpivotstrategie :

- a) Wähle im Eliminationsschritt $A^{(k)} \rightarrow A^{(k+1)}$ ein $p \in \{k, \dots, n\}$, so dass

$$\left| a_{pk}^{(k)} \right| \geq \left| a_{jk}^{(k)} \right| \quad \text{für } j = k, \dots, n$$

Die Zeile p soll die Pivotzeile werden.

- b) Vertausche die Zeilen p und k

$$A^{(k)} \rightarrow \tilde{A}^{(k)} \quad \text{mit } \tilde{a}_{ij}^{(k)} = \begin{cases} a_{kj}^{(k)}, & \text{falls } i = p \\ a_{pj}^{(k)}, & \text{falls } i = k \\ a_{ij}^{(k)}, & \text{sonst} \end{cases}$$

Nun gilt

$$\left| \tilde{l}_{ik} \right| = \left| \frac{\tilde{a}_{ik}^{(k)}}{\tilde{a}_{kk}^{(k)}} \right| = \left| \frac{\tilde{a}_{ik}^{(k)}}{a_{pk}^{(k)}} \right| \leq 1$$

- c) Führe den nächsten Eliminationsschritt angewandt auf $\tilde{A}^{(k)}$ aus,
 $\tilde{A}^{(k)} \rightarrow A^{(k+1)}$.

Bemerkung :

Anstelle der Spaltenpivotstrategie mit Zeilentausch kann man auch eine Zeilenpivotstrategie mit Spaltentausch durchführen. Beide Strategien erfordern im schlimmsten Fall $O(n^2)$ zusätzliche Operationen.

Kombiniert man beide Möglichkeiten und sucht die gesamte Restmatrix nach dem betragsgrößten Element ab, so benötigt man zusätzlich $O(n^3)$ Operationen. Diese teure vollständige Pivotsuche wird so gut wie nie angewandt.

Zur formalen Beschreibung der Dreieckszerlegung mit Spaltenpivotsuche werden im Folgenden Permutationsmatrizen $P \in \text{Mat}_n(\mathbb{R})$ verwendet. Jeder Permutation $\mathbf{p} \in S_n$ aus der symmetrischen Gruppe wird die Matrix

$$P_{\mathbf{p}} = [e_{p(1)} \dots e_{p(n)}] \quad \text{zugeordnet, wobei} \quad e_j = (\mathbf{d}_{1j}, \dots, \mathbf{d}_{nj})^T \quad \text{der } j\text{-te Einheitsvektor ist.}$$

Eine Permutation \mathbf{p} der Zeilen einer Matrix A lässt sich durch die Multiplikation von links mit $P_{\mathbf{p}}$ ausdrücken. (Analog eine Permutation der Spalten durch Multiplikation von rechts.)

Die Zuordnung $\mathbf{p} \rightarrow P_{\mathbf{p}}$ ist ein Gruppenhomomorphismus $S_n \rightarrow O(n)$ der symmetrischen Gruppe S_n in die orthogonale Gruppe $O(n)$.

Insbesondere gilt : $P^{-1} = P^T$.

Die Determinante einer Permutationsmatrix ist das Vorzeichen der zugehörigen Permutation, d.h.

$$\det P_{\mathbf{p}} = \text{sgn } \mathbf{p} \in \{\pm 1\},$$

also $+1$ falls \mathbf{p} durch eine gerade, und -1 , falls \mathbf{p} durch eine ungerade Anzahl von Transpositionen erzeugt wird.

Der folgende Satz zeigt, dass die Dreieckszerlegung mit Spaltenpivotsuche theoretisch nur versagen kann, falls die Matrix A singular ist.

Satz 1.1

Für jede invertierbare Matrix A existiert eine Permutationsmatrix P derart, dass eine Dreieckszerlegung

$$PA = LR$$

möglich ist. Dabei kann P so gewählt werden, dass alle Elemente von L betragsmäßig kleiner gleich 1 sind.

Beweis : Buch S.13/14

Bemerkung :

Es gilt : $\det A = \det(P) \cdot \det(LR) = \text{sgn}(\mathbf{p}_0) \cdot r_{11} \cdot \dots \cdot r_{nn}$

Allerdings : $\det(\mathbf{a} \cdot A) = \mathbf{a}^n \cdot \det A$

D.h. über die ganze Klasse dieser Transformationen betrachtet bleibt von der Determinante nur die invariante Boolesche Größe $\det A = 0$ oder $\det A \neq 0$ übrig.

Die Pivotstrategie lässt sich beliebig abändern, indem man die verschiedenen Zeilen mit unterschiedlichen statt gleichen Skalaren multipliziert. Dies führt zur praktisch enorm wichtigen Frage der Skalierung.

Def. Zeilenskalierung : Multiplikation von A mit einer Diagonalmatrix von links
 $A \rightarrow D_Z A$

Def. Spaltenskalierung : Multiplikation von A mit einer Diagonalmatrix von rechts
 $A \rightarrow A D_S$

Mathematisch gesprochen ändert man mit der Skalierung die Länge der Basisvektoren des Bildraumes (Zeilenskalierung) bzw. Urbildraumes (Spaltenskalierung) der durch die Matrix A beschriebenen linearen Abbildung. (physikalisch gesehen z.B. Änderung der Einheiten - z.B. von mm auf km)

Damit die Lösung des linearen Gleichungssystems $Ax = b$ nicht von dieser Wahl der Einheiten abhängt, muss das System in geeigneter Weise skaliert werden :

$$A \rightarrow \tilde{A} := D_Z A D_S, \quad ,$$

wobei

$$D_Z = \text{diag}(\mathbf{s}_1, \dots, \mathbf{s}_n) \quad \text{und} \quad D_S = \text{diag}(\mathbf{t}_1, \dots, \mathbf{t}_n) .$$

Auf den ersten Blick vernünftig scheinen die folgenden Strategien :

a) Äquilibration der Zeilen bezüglich einer Vektornorm

Sei A^i die i -te Zeile von A , und sei vorausgesetzt, dass keine Zeile eine Nullzeile ist.

Setzt man dann $D_S := I$ und

$$s_i := \|A^i\|^{-1} \quad \text{für } i = 1, \dots, n \quad ,$$

so haben alle Zeilen von \tilde{A} die Norm 1.

b) Äquilibration der Spalten

Seien alle Spalten A_j verschieden von Null.

Setzt man dann $D_Z := I$ und

$$t_j := \|A_j\|^{-1} \quad \text{für } j = 1, \dots, n \quad ,$$

so haben alle Spalten von \tilde{A} die Norm 1.

Problem :

Die berechnete Lösung \tilde{x} kann immer noch „ziemlich ungenau“ sein.

Man könnte die Lösung \tilde{x} verwerfen und versuchen, mit einer höheren Maschinengenauigkeit eine „bessere“ Lösung zu berechnen.

Dabei würden jedoch alle Informationen verloren gehen, die man bei der Berechnung von \tilde{x} gewonnen hat.

Dies wird bei der sogenannten **iterativen Verbesserung** oder **Nachiteration** vermieden, bei der das **Residuum**

$$r(y) := b - Ay = A(x - y)$$

explizit ausgewertet wird. Verschwände das Residuum, so wäre die exakte Lösung x gefunden.

Der absolute Fehler $\Delta x_0 := x - x_0$ von $x_0 := \tilde{x}$ genügt der Gleichung

$$A \cdot \Delta x_0 = r(x_0) \quad (*)$$

Bei der Lösung dieser Korrekturgleichung (*) erhält man im Allgemeinen eine wiederum fehlerbehaftete

Korrektur $\tilde{\Delta}x_0 \neq \Delta x_0$. Trotzdem erwartet man, dass die Näherungslösung

$$x_1 := x_0 + \tilde{\Delta}x_0$$

„besser“ ist als x_0 .

Die Idee der Nachiteration besteht nun darin, diesen Prozess solange zu wiederholen, bis die Näherungslösung x_i „genau genug“ ist.

Dabei beachte man, das sich das Gleichungssystem (*) nur in der rechten Seite von dem Ursprünglichen unterscheidet, so das die Berechnung der Korrekturen Δx_i relativ wenig Aufwand erfordert.

(weiterführend in Kapitel 2.3.2 - „Beurteilung von Näherungslösungen“)

1.3 Cholesky Zerlegung (für spd-Matrizen)

Motivation : Bei Gleichungssystemen mit spd-Matrizen kann die Dreieckszerlegung stark vereinfacht werden.

Satz 1.2 :

Für jede spd-Matrix existiert eine eindeutig bestimmte Zerlegung der Form $A = LDL^T$, wobei L eine unipotente untere Dreiecksmatrix und D eine positive Diagonalmatrix (der Pivotelemente) ist.

Kor. 1.3 :

Da $D = \text{diag}(d_i)$ positiv ist, existiert $D^{\frac{1}{2}} = \text{diag}(\sqrt{d_i})$ und daher die Cholesky-Zerlegung

$$A = \bar{L}\bar{L}^T \quad , \text{wobei } \bar{L} \text{ die untere Dreiecksmatrix } \bar{L} := LD^{\frac{1}{2}} \text{ ist.}$$

Algorithmus der Cholesky Zerlegung

for $k := 1$ to n do

$$l_{kk} := \left(a_{kk} - \sum_{j=1}^{k-1} l_{kj}^2 \right)^{\frac{1}{2}}$$

for $i := k+1$ to n do

$$l_{ik} := \left(a_{ik} - \sum_{j=1}^{k-1} l_{ij} \cdot l_{kj} \right) / l_{kk}$$

endfor

endfor

Beispiel für Cholesky-Zerlegung :

Gegeben : $A = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$

1.) mit LR-Zerlegung ($A = LDL^T$)

$$A^{(1)} = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}, l_{ik} := \frac{a_{ik}^{(k)}}{a_{kk}^{(k)}} \Rightarrow l_{21} = \frac{1}{2} \Rightarrow L = \begin{pmatrix} 1 & 0 \\ 1/2 & 1 \end{pmatrix}$$

$$a_{ij}^{(k+1)} := a_{ij}^{(k)} - l_{ik} \cdot a_{kj}^{(k)} \Rightarrow a_{21}^{(2)} = 1 - \frac{1}{2} \cdot 2 = 0 \text{ und } a_{22}^{(2)} = 2 - \frac{1}{2} \cdot 1 = \frac{3}{2}$$

$$\Rightarrow A^{(2)} = \begin{pmatrix} 2 & 1 \\ 0 & 3/2 \end{pmatrix} \Rightarrow D = \begin{pmatrix} 2 & 0 \\ 0 & 3/2 \end{pmatrix}$$

Probe : $LDL^T = \begin{pmatrix} 1 & 0 \\ 1/2 & 1 \end{pmatrix} \cdot \begin{pmatrix} 2 & 0 \\ 0 & 3/2 \end{pmatrix} \cdot \begin{pmatrix} 1 & 1/2 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 2 & 0 \\ 1 & 3/2 \end{pmatrix} \cdot \begin{pmatrix} 1 & 1/2 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} = A$

2.) direkt ($A = \bar{L} \cdot \bar{L}^T$)

$$l_{kk} := \left(a_{kk} - \sum_{j=1}^{k-1} l_{kj}^2 \right)^{\frac{1}{2}} \Rightarrow l_{11} = \sqrt{\left(a_{11} - \sum_{j=1}^0 l_{1j}^2 \right)} = \sqrt{a_{11}} = \sqrt{2}$$

for $i := k+1$ to n do $l_{ik} := \left(a_{ik} - \sum_{j=1}^{k-1} l_{ij} \cdot l_{kj} \right) / l_{kk} \Rightarrow l_{21} = \frac{a_{21} - \sum_{j=1}^0 l_{1j} \cdot l_{2j}}{l_{11}} = \frac{1}{\sqrt{2}}$

$$l_{kk} := \left(a_{kk} - \sum_{j=1}^{k-1} l_{kj}^2 \right)^{\frac{1}{2}} \Rightarrow l_{22} = \sqrt{\left(a_{22} - \sum_{j=1}^1 l_{2j}^2 \right)} = \sqrt{a_{22} - l_{21}^2} = \sqrt{2 - \frac{1}{2}} = \sqrt{\frac{3}{2}}$$

Probe: $\bar{L} \cdot \bar{L}^T = \begin{pmatrix} \sqrt{2} & 0 \\ 1 & \sqrt{\frac{3}{2}} \end{pmatrix} \cdot \begin{pmatrix} \sqrt{2} & \frac{1}{\sqrt{2}} \\ 0 & \sqrt{\frac{3}{2}} \end{pmatrix} = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix} = A$

2. Fehleranalyse

Es gibt zwei Arten von Fehlern :

- **Eingabefehler** : Maschinengenauigkeit , Messfehler $- f(\tilde{x}) - f(x) \rightarrow$ Kondition eines Problems
 - relativer Fehler einer Gleitkommazahl d mit Mantissenlänge l : $\frac{|x - fl(x)|}{|x|} \leq eps := \frac{d^{1-l}}{2}$
 - Messfehler : absoluter Fehler/Toleranz $\mathbf{d} \cdot x \quad |x - \bar{x}| \leq \mathbf{d} \cdot x \quad (\bar{x} = \text{„wahrer“ Wert})$
- **Fehler im Algorithmus** : Rundungsfehler , Approximationsfehler $- \tilde{f}(x) - f(x) \rightarrow$ Stabilität eines Algorithmus
 - Rundungsfehler : $a \hat{=} b = (a \circ b)(1 + \mathbf{e})$ für ein $\mathbf{e} = \mathbf{e}(x, y)$ mit $|\mathbf{e}| \leq eps$
 - Approximationsfehler : Berechnung des Sinus durch Reihenentwicklung

Machtlos gegenüber Eingabefehlern.

Fehler im Algorithmus kann man durch Änderung des Verfahrens vermeiden bzw. verringern

Gleitpunktzahlen sind auf der reellen Achse nicht gleichverteilt !

Def. Maschinenepsilon $eps := \mathbf{b}^{1-l}$: $(\mathbf{b}$ - Basis l - Mantissenlänge)

= kleinste Maschinenzahl, die zu 1 addiert einen von 1 verschiedenen Wert ergibt.

Def. absoluter / relativer Fehler :

Der absolute Fehler einer Größe mit Sollwert \bar{x} und Istwert x ist $\mathbf{dx} := \|x - \bar{x}\|$.

Ist $\bar{x} \neq 0$, so ist der relative Fehler definiert als $\frac{\|x - \bar{x}\|}{\|\bar{x}\|}$.

2.1 Kondition eines Problems

= Charakterisierung des Verhältnisses von Eingabe- und Resultatmenge

Die Eingabe x ist logisch ununterscheidbar von allen Eingaben \tilde{x} , die sich im Rahmen der vorgegebenen Genauigkeit befinden. Statt der exakten Eingabe x betrachtet man die Eingabemenge E aller dieser gestörten Eingaben \tilde{x} .

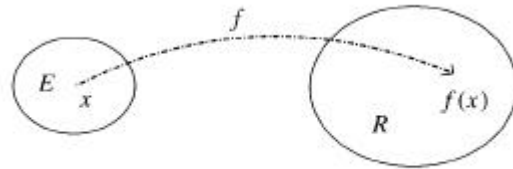


Abbildung 2.1. Eingabe- und Resultatmenge

Eine Maschinenzahl x repräsentiert demnach die Eingabemenge $E = \{\tilde{x} \in \mathbb{R} : |\tilde{x} - x| \leq \text{eps} \cdot |x|\}$.

Wüsste man hingegen, dass die Eingabe x mit einer absoluten Toleranz $d \cdot x$ vorliegt, so wäre die Eingabemenge $E = \{\tilde{x} \in \mathbb{R} : |\tilde{x} - x| \leq d \cdot x\}$.

Beispiel für die Kondition eines Problems : zeichnerische Bestimmung des Schnittpunktes zweier Geraden

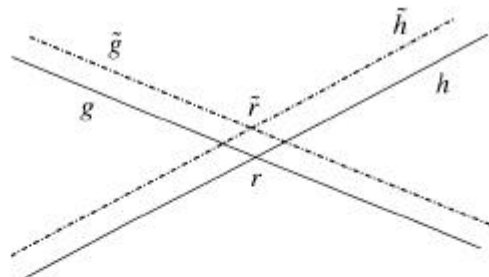


Abbildung 2.2. Schnittpunkt r zweier Geraden g, h (gut konditioniert)

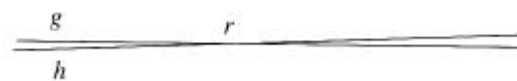


Abbildung 2.3. Schleifender Schnitt (schlecht konditioniert)

Beispiel für die Kondition eines Problems : zeichnerische Bestimmung der Nullstellen einer Funktion

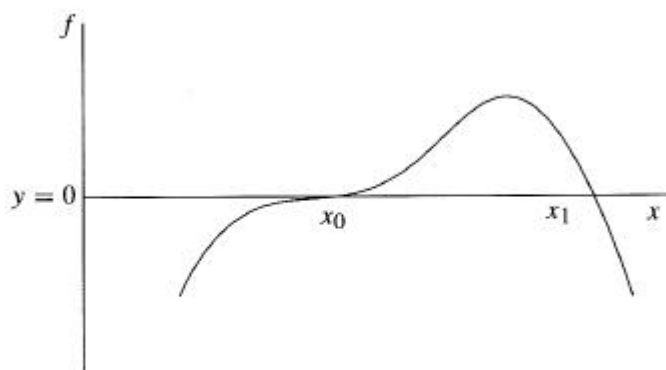


Abbildung 2.4. Schlecht konditionierte Nullstelle bei x_0 , gut konditionierte bei x_1

2.1.1 Normweise Konditionsanalyse

Die **absolute normweise Kondition** des Problems (f, x) ist die kleinste Zahl $k_{abs} \geq 0$,

so daß $\|f(\tilde{x}) - f(x)\| \leq k_{abs} \|\tilde{x} - x\|$ für $\tilde{x} \rightarrow x$

Das Problem ist schlecht bzw. unsachmäßig gestellt, wenn es keine solche Zahl gibt.

Die **relative normweise Kondition** des Problems (f, x) ist die kleinste Zahl $k_{rel} \geq 0$,

so daß $\frac{\|f(\tilde{x}) - f(x)\|}{\|f(x)\|} \leq k_{rel} \frac{\|\tilde{x} - x\|}{\|x\|}$ für $\tilde{x} \rightarrow x$

D.h. k_{abs} beschreibt die Verstärkung des absoluten Fehlers und k_{rel} die Verstärkung des relativen Fehlers.

Ein Problem ist gut konditioniert, falls seine Kondition klein ist.

Als Orientierung : $k_{rel} = 1$ entspricht der reinen Rundung des Resultats.

Für differenzierbare f gilt aufgrund des Mittelwertsatzes :

- $k_{abs} = \|f'(x)\|$ und $k_{rel} = \|x\| \cdot \frac{\|f'(x)\|}{\|f(x)\|}$

, wobei $\|f'(x)\|$ die Norm der Jacobi-Matrix $f'(x)$ in der zugeordneten Matrixnorm $\|A\| = \sup_{x \neq 0} \frac{\|Ax\|}{\|x\|}$ ist.

Beispiel : Kondition der Addition (bzw. Subtraktion)

$$f : \mathbb{R}^2 \rightarrow \mathbb{R} \quad \text{d.h.} \quad f(x, y) := x + y \quad \Rightarrow f'(x, y) = \begin{pmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \end{pmatrix} = (1 \quad 1)$$

Wählt man auf \mathbb{R}^2 die 1-Norm $\|(x \quad y)^T\| = |x| + |y|$ und auf \mathbb{R} den Betrag, so folgt für die Ableitung in der zugeordneten Matrixnorm, dass $\|f'(x, y)\| = \|(1 \quad 1)\| = 1$.

$$\Rightarrow \mathbf{k}_{abs} = 1 \quad \text{und} \quad \mathbf{k}_{rel} = \frac{|x| + |y|}{|x + y|}$$

Hieraus folgt, dass die Subtraktion zweier annähernd gleicher Zahlen bezüglich der relativen Kondition schlecht konditioniert ist, da in diesem Fall gilt : $|x + y| \ll |x| + |y| \Leftrightarrow \mathbf{k}_{rel} \gg 1$.

\Rightarrow **Vermeidbare Subtraktionen annähernd gleicher Zahlen vermeiden.**
Falls unvermeidbar, diese an den Anfang des Algorithmus stellen. (siehe Vorwärtsanalyse)

Beispiel für vermeidbare Auslöschung : Lösen einer quadratischen Gleichung

$$f(x) := x^2 + px + q \quad x_{1,2} = -\frac{p}{2} \pm \sqrt{\frac{p^2}{4} - q}$$

Liegt eine Nullstelle in der Nähe von Null, so tritt Auslöschung auf.

Nach dem Satz von Vieta jedoch ist q das Produkt der Nullstellen und somit lassen sich die beiden Lösungen

stabil berechnen :
$$x_1 = -\frac{p}{2} + \operatorname{sgn}(p) \sqrt{\frac{p^2}{4} - q} \quad , \quad x_2 = \frac{q}{x_1}$$

Beispiel : Kondition eines LGS $Ax=b$

Betrachtet man bei der Lösung des LGS nur den Vektor b als Eingabegröße, so wird das Problem von der Abbildung $f(b) := A^{-1}b$ beschrieben.

Diese ist bezüglich b linear, so daß $f'(b) = A^{-1}$.

$$\Rightarrow \mathbf{k}_{abs} = \|A^{-1}\| \quad \text{und} \quad \mathbf{k}_{rel} = \|b\| \cdot \frac{\|A^{-1}\|}{\|A^{-1}b\|} = \frac{\|Ax\|}{\|x\|} \cdot \|A^{-1}\|$$

Es lassen sich aber auch leicht Störungen in A berücksichtigen.

Die Abbildung $f(A) = A^{-1}b$ ist nichtlinear in A , aber differenzierbar.

Lemma 2.1 :

Die Abbildung $g : GL(n) \subset Mat_n(\mathbb{R}) \rightarrow GL(n)$ mit $g(A) = A^{-1}$ ist differenzierbar, und es gilt

$$g'(A) \cdot C = -A^{-1} \cdot C \cdot A^{-1} \quad \text{für alle } C \in Mat_n(\mathbb{R})$$

Beweis : Buch S.35

Aus dem Lemma folgt für die Ableitung der Lösung $f(A) = A^{-1}b$ des LGS nach A , dass

$$f'(A) \cdot C = -A^{-1} \cdot C \cdot A^{-1}b = -A^{-1} \cdot C \cdot x \quad \text{für } C \in \text{Mat}_n(R).$$

Hieraus folgen nun die Konditionen

$$\mathbf{k}_{abs} = \|f'(A)\| = \sup_{\|C\|=1} \|A^{-1} \cdot C \cdot x\| \leq \|A^{-1}\| \cdot \|x\|$$

$$\mathbf{k}_{rel} = \frac{\|A\|}{\|x\|} \|f'(A)\| \leq \|A\| \cdot \|A^{-1}\| \quad (\text{Abschätzung siehe oben})$$

Die Größe $\mathbf{k}(A) := \|A\| \cdot \|A^{-1}\|$ wird **Kondition der Matrix** A genannt.

Sie beschreibt insbesondere die relative Kondition eines linearen Gleichungssystems $Ax = b$ für alle möglichen rechten Seiten $b \in R^n$.

Eine andere Darstellung für $\mathbf{k}(A)$ ist
$$\mathbf{k}(A) := \frac{\max_{\|x\|=1} \|Ax\|}{\min_{\|x\|=1} \|Ax\|} \in [0, \infty]. \quad (\text{Buch Übung 2.12})$$

Vorteil : Auch für nicht invertierbare und rechteckige Matrizen wohldefiniert.

Aus dieser Darstellung ergeben sich sofort folgende drei Eigenschaften von $\mathbf{k}(A)$:

- 1) $\mathbf{k}(A) \geq 1$
- 2) $\mathbf{k}(a \cdot A) = \mathbf{k}(A)$ für alle $a \in R$, $a \neq 0$
- 3) $A \neq 0$ ist singulär $\Leftrightarrow \mathbf{k}(A) = \infty$

D.h. die Kondition $\mathbf{k}(A)$ einer Matrix ist im Gegensatz zur Determinante $\det A$ invariant unter den skalaren Transformationen $A \rightarrow a \cdot A$.

Zusammen mit den Eigenschaften 1) und 3) kann diese Konditionszahl daher eher zur Charakterisierung der Lösbarkeit eines LGS herangezogen werden als die Determinante.

2.1.2 Komponentenweise Konditionsanalyse

Günstiger, weil alle Eingaben in einen Rechner mit einem relativen Fehler in den einzelnen Komponenten behaftet sind und so einige Phänomene bei der normweisen Betrachtung nicht erklärt werden können.

Auch nimmt die normweise Betrachtung z.B. bei LGS keine Rücksicht auf eventuell vorliegende Spezialstrukturen einer Matrix A , sondern analysiert das Verhalten bezüglich beliebiger Störungen dA , d.h. auch solcher, die diese Spezialstruktur nicht erhalten.

Beispiel : Kondition einer Diagonalmatrix

Die Lösung eines Gleichungssystems $Ax = b$ mit einer Diagonalmatrix $A = \begin{pmatrix} 1 & 0 \\ 0 & e \end{pmatrix}$, $A^{-1} = \begin{pmatrix} 1 & 0 \\ 0 & e^{-1} \end{pmatrix}$

ist offensichtlich ein gut konditioniertes Problem, da die Gleichungen vollkommen unabhängig voneinander (entkoppelt) sind. (Dabei wird natürlich implizit vorausgesetzt, dass die Störungen die Diagonalgestalt erhalten.)

Die normweise Kondition $\mathbf{k}_\infty(A) = \|A\|_\infty \cdot \|A^{-1}\|_\infty = 1 \cdot e^{-1} = \frac{1}{e}$ wird für kleine $e \leq 1$ jedoch beliebig groß.

Erwartungsgemäß sollte die Kondition einer Diagonalmatrix, d.h. eines vollständig entkoppelten Gleichungssystems, stets 1 betragen – wie für eine skalare Gleichung.

\Rightarrow Einführung der komponentenweisen Konditionsanalyse

Def. komponentenweise Kondition :

Die (im Urbild) komponentenweise Kondition des Problems (f, x) ist die kleinste Zahl $k_{rel} \geq 0$, so dass

$$\frac{\|f(\tilde{x}) - f(x)\|_\infty}{\|f(x)\|_\infty} \leq k_{rel} \cdot \max_i \frac{|\tilde{x}_i - x_i|}{|x_i|} \quad \text{für } \tilde{x} \rightarrow x.$$

Alternativ kann man die relative komponentenweise Kondition auch wie folgt definieren :

$$\max_i \frac{|f_i(\tilde{x}) - f_i(x)|}{|f_i(x)|} \leq k_{rel} \cdot \max_i \frac{|\tilde{x}_i - x_i|}{|x_i|} \quad \text{für } \tilde{x} \rightarrow x.$$

Die Anwendung des Mittelwertsatzes $f(\tilde{x}) - f(x) = \int_{t=0}^1 f'(x + t \cdot (\tilde{x} - x)) \cdot (\tilde{x} - x) dt$ liefert komponentenweise

$$|f(\tilde{x}) - f(x)| \leq \int_{t=0}^1 |f'(x + t \cdot (\tilde{x} - x))| \cdot |(\tilde{x} - x)| dt \quad \text{und daher}$$

- $k_{rel} = \frac{\| |f'(x)| \cdot |x| \|_\infty}{\| f(x) \|_\infty}$

Beispiel : Kondition der Multiplikation

Die Multiplikation zweier reeller Zahlen wird durch die Abbildung $f : R^2 \rightarrow R$, $(x, y)^T \mapsto f(x, y) = xy$

beschrieben. Sie ist differenzierbar mit $f'(x, y) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right) = (y, x)$ und somit folgt

$$k_{rel} = \frac{\| |f'(x, y)| \cdot |(x, y)^T \|_\infty}{\| f(x, y) \|_\infty} = \frac{|yx| + |xy|}{|xy|} = 2.$$

→ Die Multiplikation ist daher noch gut konditioniert zu nennen.

Beispiel : Kondition des Skalarproduktes

$\langle x, y \rangle = \sum_{i=1}^n x_i y_i$ Es gilt, die Abbildung $f : R^{2n} \rightarrow R$, $(x, y) \mapsto \langle x, y \rangle$ an der Stelle (x, y) auszuwerten.

Da f differenzierbar ist mit $f'(x, y) = (y^T, x^T)$, folgt

$$k_{rel} = \frac{\| |(y^T, x^T)| \cdot |(x, y) \|_\infty}{\| \langle x, y \rangle \|_\infty} = \frac{|y^T| \cdot |x| + |x^T| \cdot |y|}{|\langle x, y \rangle|} = 2 \frac{\langle |x|, |y| \rangle}{|\langle x, y \rangle|}$$

Beispiel : komponentenweise Kondition eines linearen Gleichungssystems (Skeelsche Kondition)

$Ax = b$

- Betrachtet man nur den Vektor b als Eingabegröße, so wird das Problem von der Abbildung $f(b) = A^{-1}b$ beschrieben.
 $\Rightarrow f'(b) = A^{-1}$, da f bezüglich b linear

$$\Rightarrow k_{rel} = \frac{\| |f'(b)| \cdot |b| \|_{\infty}}{\| f(b) \|_{\infty}} = \frac{\| |A^{-1}| \cdot |b| \|_{\infty}}{\| A^{-1}b \|_{\infty}} = \frac{\| |A^{-1}| \cdot |b| \|_{\infty}}{\| x \|_{\infty}}$$

Mit dieser, von Skeel eingeführten, Zahl lassen sich die Störungen $\tilde{x} - x$, $\tilde{x} = A^{-1}\tilde{b}$ wie folgt abschätzen :

$$\frac{\| \tilde{x} - x \|_{\infty}}{\| x \|_{\infty}} \leq k_{rel} \cdot e \quad \text{für } | \tilde{b} - b | \leq e \cdot |b|$$

- Betrachtet man Störungen in A , so gilt :

Die Abbildung $f(A) = A^{-1}b$ ist differenzierbar mit $f'(A) \cdot C = -A^{-1}CA^{-1}b = -A^{-1}Cx$ für $C \in Mat_n(R)$.

$$\Rightarrow k_{rel} = \frac{\| |f'(A)| \cdot |A| \|_{\infty}}{\| f(A) \|_{\infty}} = \frac{\| |A^{-1}| \cdot |A| \cdot |x| \|_{\infty}}{\| x \|_{\infty}}$$

Fasst man diese Ergebnisse zusammen und betrachtet Störungen in A und b , so addieren sich die beiden relativen Konditionen auf, und man erhält die **Kondition für das Gesamtproblem** :

$$k_{rel} = \frac{\| |A^{-1}| \cdot |A| \cdot |x| + |A^{-1}| \cdot |b| \|_{\infty}}{\| x \|_{\infty}} \leq 2 \frac{\| |A^{-1}| \cdot |A| \cdot |x| \|_{\infty}}{\| x \|_{\infty}}$$

Setzt man für x den Vektor $e = (1, \dots, 1)$ ein, so ergibt sich eine Charakterisierung der komponentenweisen Kondition von $Ax = b$ für alle möglichen rechten Seiten b :

$$\frac{1}{2} k_{rel} \leq \frac{\| |A^{-1}| \cdot |A| \cdot |e| \|_{\infty}}{\| e \|_{\infty}} = \| |A^{-1}| \cdot |A| \|_{\infty} \quad \text{mit der sogenannten **Skeelschen Kondition** } k_C(A) := \| |A^{-1}| \cdot |A| \|_{\infty}$$

.Diese Kondition $k_C(A)$ erfüllt die gleichen drei Eigenschaften wie $k(A)$ (siehe früher)

und zusätzlich die eingangs geforderte Eigenschaft $k_C(D) = 1$ für jede Diagonalmatrix D .

Sie ist sogar invariant unter Zeilenskalierung, d.h. $k_C(D \cdot A) = k_C(A)$,

da $| (DA)^{-1} | \cdot |DA| = |A^{-1}| \cdot |D^{-1}| \cdot |D| \cdot |A| = |A^{-1}| \cdot |A|$.

2.2 Stabilitätskonzepte

Definition : $\tilde{f}(E) := \bigcup_{\Phi \in \tilde{f}} \Phi(E)$.

2.2.1 Vorwärtsanalyse

Analyse der Menge $\tilde{R} := \tilde{f}(E)$ aller durch Eingabefehler und Fehler im Algorithmus gestörten Resultate.

Ist \tilde{R} von der gleichen Größenordnung wie R , so ist der Algorithmus stabil im Sinne der Vorwärtsanalyse.

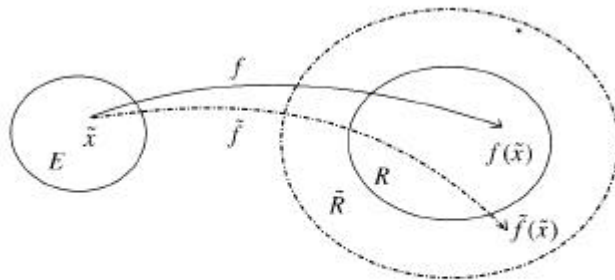


Abbildung 2.5. Eingabemenge und Resultatmengen bei der Vorwärtsanalyse

Der vom Algorithmus verursachte Fehler $\tilde{f}(x) - f(x)$ muss in Beziehung gesetzt werden zum unvermeidbaren Fehler $\mathbf{k}_{rel} \cdot eps$. (eps = Eingabefehler)

Def. Stabilitätsindikator der normweisen Vorwärtsanalyse :

= Faktor um den der Algorithmus den unvermeidbaren Fehler verstärkt.

= kleinstmögliche Zahl $\mathbf{s} \geq 0$, so dass für alle $\tilde{x} \in E$ gilt :

$$\frac{\|\tilde{f}(\tilde{x}) - f(\tilde{x})\|}{\|f(\tilde{x})\|} \leq \mathbf{s} \cdot \mathbf{k}_{rel} \cdot eps \quad , \text{ für } eps \rightarrow 0 \quad ,$$

wobei \tilde{f} im Folgenden die Gleitkommarealisierung eines Algorithmus zur Lösung des Problems (f, x) ist.

Analog:

Def. Stabilitätsindikator der komponentenweisen Vorwärtsanalyse :

= kleinstmögliche Zahl $\mathbf{s} \geq 0$, so dass

$$\max_i \frac{|\tilde{f}_i(\tilde{x}) - f_i(\tilde{x})|}{|f_i(\tilde{x})|} \leq \mathbf{s} \cdot \mathbf{K}_{rel} \cdot eps \quad , \text{ wobei } \mathbf{K}_{rel} = \text{komponentenweise relative Kondition von } (f, x)$$

- Der Algorithmus \tilde{f} heißt stabil im Sinne der Vorwärtsanalyse, falls \mathbf{s} kleiner als die Anzahl der hintereinander ausgeführten Elementaroperationen ist.

Lemma 2.2 :

Für die Elementaroperationen $\{+, -, *, /\}$ und ihre Gleitkommarealisierungen $\{\hat{+}, \hat{-}, \hat{*}, \hat{/}\}$ gilt :

$$\mathbf{s} \cdot \mathbf{k}_{rel} \leq 1$$

Beweis :

$$a \hat{\circ} b = (a \circ b)(1 + \mathbf{e}) \quad \text{für ein } \mathbf{e} \text{ mit } |\mathbf{e}| \leq eps. \quad \text{Hieraus folgt :}$$

$$\frac{|a \hat{\circ} b - a \circ b|}{|a \circ b|} = \frac{|(a \circ b)(1 + \mathbf{e}) - a \circ b|}{|a \circ b|} = |\mathbf{e}| \leq eps \quad \text{Q.e.d.}$$

Beispiel : Stabilität der Subtraktion

Im Fall der Auslöschung, $\mathbf{k} \gg 1$, ist der Stabilitätsindikator sehr klein, $\mathbf{s} \ll 1$.

Die Subtraktion ist in diesem Fall also hervorragend stabil und bei totaler Auslöschung sogar fehlerfrei, $a \hat{-} b = a - b$.

Lemma 2.3 :

Sowohl im norm- als auch im komponentenweisen Konzept gilt für den Stabilitätsindikator \mathbf{s}_f des

zusammengesetzten Algorithmus $\tilde{f} = \tilde{h} \circ \tilde{g}$, dass

$$\mathbf{s}_f \mathbf{k}_f \leq \mathbf{s}_h \mathbf{k}_h + \mathbf{s}_g \mathbf{k}_g \mathbf{k}_h.$$

Beweis : S.45

Schlußfolgerung :

- **Unvermeidbare Subtraktionen möglichst an den Anfang des Algorithmus stellen.**

Denn da für jede Elementaroperation $\mathbf{s} \cdot \mathbf{k}_{rel} \leq 1$ gilt, droht Gefahr für die Stabilität der zusammengesetzten Abbildung $f = h \circ g$ nur, falls es sich bei der zweiten Abbildung h um eine Subtraktion handelt, so dass $\mathbf{k}_h \gg 1$.

Lemma 2.4 :

Sind die Funktionen g und h aus Lemma 2.3 skalar und differenzierbar, so gilt für den Stabilitätsindikator \mathbf{s}_f des

zusammengesetzten Algorithmus $\tilde{f} = \tilde{h} \circ \tilde{g}$, dass

$$\mathbf{s}_f \leq \frac{\mathbf{s}_h}{\mathbf{k}_g} + \mathbf{s}_g.$$

Beweis : S.47

Schlußfolgerung :

- Ist die Kondition \mathbf{k}_g des ersten Teilproblems sehr klein, $\mathbf{k}_g \ll 1$, so wird der Algorithmus also instabil.

Eine kleine Kondition lässt sich auch als Informationsverlust interpretieren :

Eine Änderung der Eingabe hat so gut wie keinen Einfluss auf das Resultat.

Ein solcher Informationsverlust zu Anfang des Algorithmus hat daher Instabilität zur Folge.

Desweiteren schlägt sich eine Instabilität zu Beginn des Algorithmus (großes \mathbf{s}_g) voll auf den Gesamtalgorithmus durch.

TODO : Beispiel : $\cos mx$ S.47/48 ???

Beispiel : Summation von n reellen Zahlen (komponentenweise Vorwärtsanalyse)

Der einfachste Algorithmus für die Summe $s_n : R^n \rightarrow R$, $(x_1, \dots, x_n) \mapsto \sum_{i=1}^n x_i$ ist ihre rekursive Berechnung

gemäß $s_n = s_{n-1} \circ a_n$ für $n > 2$ und $s_2 = a_2$, wobei $a_n : R^n \rightarrow R^{n-1}$, $(x_1, \dots, x_n) \mapsto (x_1 + x_2, x_3, \dots, x_n)$ die Addition der ersten beiden Komponenten bezeichnet.

Konditionszahl und Stabilitätsindikator von a_n stimmen mit denen der Addition zweier Zahlen überein,

d.h. $k_{a_n} = k_+$ und $s_{a_n} = s_+$.

Mit den Bezeichnungen $k_j := k_{s_j}$ und $s_j := s_{s_j}$ gilt nach Lemma 2.3, dass

$$a_n k_n \leq s_{n-1} k_{n-1} + s_+ k_+ k_{n-1} = (s_{n-1} + s_+ k_+) \cdot k_{n-1} \leq (1 + s_{n-1}) \cdot k_{n-1} .$$

- Erinnerung : Konditionszahlen der Addition : $k_{abs} = 1$ und $k_{rel} = \frac{|x| + |y|}{|x + y|}$.

Für die Konditionen k_n gilt somit, dass $k_n = \frac{\sum_{i=1}^n |x_i|}{\left| \sum_{i=1}^n x_i \right|} \geq 1$ und $k_{n-1} = \frac{\sum_{i=3}^n |x_i| + |x_1 + x_2|}{\left| \sum_{i=1}^n x_i \right|} \leq k_n$,

und daher $s_n \leq 1 + s_{n-1}$. Da $s_2 = s_+ \leq \frac{1}{k_+} \leq 1$, folgt für den Stabilitätsindikator $s_n \leq n - 1$.

Bei den benötigten $n - 1$ Elementaroperationen ist der naive Algorithmus zur Summenbildung also numerisch stabil.

Beispiel : Ausführung des Skalarproduktes

Unterteilung der Berechnung des Skalarproduktes in die komponentenweise Multiplikation

$p : R^n \times R^n \rightarrow R^n$, $((x_i), (y_i)) \mapsto (x_i y_i)$ gefolgt von der im letzten Beispiel analysierten Summenbildung.

D.h. $f = s_n \circ p$.

Nach Lemma 2.3 gilt zusammen mit Lemma 2.2 und der Abschätzung des Stabilitätsindicators im letzten Beispiel, dass

$$s_f k_f \leq (s_n + s_p k_p) \cdot k_n \leq (1 + s_n) \cdot k_n \leq n \cdot k_n \quad \text{und daher}$$

$$s_f \leq n \frac{k_n}{k_f} = n \frac{\frac{\sum_{i=1}^n |x_i y_i|}{\left| \sum_{i=1}^n x_i y_i \right|}}{\frac{\sum_{i=1}^n |x_i y_i|}{2 \left| \sum_{i=1}^n x_i y_i \right|}} = \frac{n}{2}$$

Bei $2n - 1$ Elementaroperationen ist dieser Algorithmus für das Skalarprodukt also numerisch stabil.

2.2.2 Rückwärtsanalyse

Die durch den Algorithmus verursachten Fehler werden auf die Eingabegrößen zurückgespielt und so als zusätzliche Eingabefehler interpretiert.

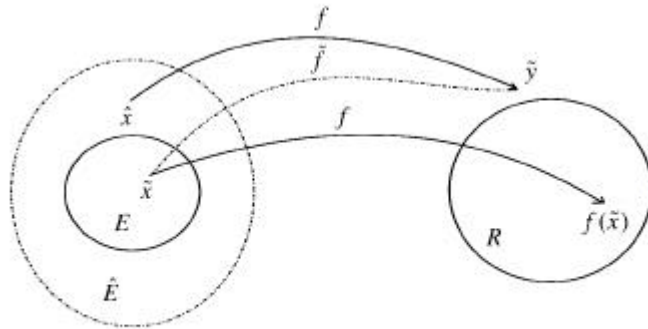


Abbildung 2.6. Eingabemengen und Resultatmenge bei der Rückwärtsanalyse

Die fehlerbehafteten Resultate $\tilde{y} = \tilde{f}(\tilde{x})$ werden als exakte Ergebnisse $\tilde{y} = f(\hat{x})$ zu gestörten Eingabegrößen \hat{x} aufgefaßt.

Dies gelingt nur, wenn $\tilde{f}(E)$ im Bild von f liegt.

Ist dies nicht der Fall, so ist eine Rückwärtsanalyse nicht möglich, und der Algorithmus instabil im Sinne der Rückwärtsanalyse.

Für nicht injektive Abbildungen f wird es mehr als ein $\hat{x} \in f^{-1}(\tilde{y})$ geben.

Hier wählt man das Element \hat{x} mit der geringsten Abweichung von der Eingabe \tilde{x} , d.h. $\|\hat{x} - \tilde{x}\| = \min$.

$$\Rightarrow \hat{E} := \left\{ \hat{x} : f(\hat{x}) = \tilde{f}(\tilde{x}) \text{ und } \|\hat{x} - \tilde{x}\| = \min \text{ für ein } \tilde{x} \in E \right\}$$

Das Verhältnis von \hat{E} zur Eingabemenge E ist ein Maß für die Stabilität im Sinn der Rückwärtsanalyse.

Def. normweiser Rückwärtsfehler :

Der normweise Rückwärtsfehler des Algorithmus \tilde{f} zur Lösung des Problems (f, x) , ist die kleinste Zahl, für die für alle $\tilde{x} \in E$ ein \hat{x} existiert, so dass

$$\frac{\|\hat{x} - \tilde{x}\|}{\|\tilde{x}\|} \leq h \quad \text{für } eps \rightarrow 0 .$$

Analog :

Def. komponentenweiser Rückwärtsfehler

= kleinste Zahl, für die für alle $\tilde{x} \in E$ ein \hat{x} existiert, so dass

$$\max_i \frac{|\hat{x}_i - \tilde{x}_i|}{|\tilde{x}_i|} \leq h \quad \text{für } eps \rightarrow 0 .$$

Der Algorithmus heißt **stabil bezüglich des relativen Eingabefehlers d** , falls

$$h < d .$$

Für den durch Rundung verursachten Eingabefehler $d = eps$ wird definiert :

Def. Stabilitätsindikator der Rückwärtsanalyse :

$$s_R := \frac{h}{eps}$$

In der Definition taucht die Kondition des Problems nicht auf.

Die Rückwärtsanalyse benötigt im Gegensatz zur Vorwärtsanalyse keine vorhergehende Konditionsanalyse des Problems. Auch sind die Ergebnisse leicht durch den Vergleich von Eingabefehler und Rückwärtsfehler zu interpretieren.

- **Aufgrund dieser Eigenschaften ist die Rückwärtsanalyse insbesondere bei komplexeren Algorithmen vorzuziehen.**

Lemma 2.5

$$S \leq S_R$$

Insbesondere folgt aus der Rückwärtsstabilität die Vorwärtsstabilität.

Beweis : S.50

Lemma 2.6

Gegeben sei die rekursive Vorschrift

$$\langle x, y \rangle := x_n y_n + \langle x^{(n-1)}, y^{(n-1)} \rangle$$

zur Auswertung des Skalarproduktes, wobei $x^{(n-1)} := (x_1, \dots, x_{n-1})^T$ und $y^{(n-1)} := (y_1, \dots, y_{n-1})^T$.

In dieser Form findet es sich meistens auf sequentiell arbeitenden Rechnern.

Hier nun die Frage nach der Stabilität bezüglich der Eingabe x , da diese später bei der Rückwärtsanalyse der Gauß-Elimination für die Analyse der Rückwärtssubstitution benötigt wird.

Die zugehörige Gleitpunktrealisierung des Skalarproduktes berechnet für $x, y \in R^n$ eine Lösung $\langle x, y \rangle_{fl}$, so dass

$$\langle x, y \rangle_{fl} = \langle \hat{x}, y \rangle \quad \text{für ein } \hat{x} \in R^n \text{ mit } |x - \hat{x}| \leq n \cdot eps \cdot |x|$$

, d.h. der relative komponentenweise Rückwärtsfehler beträgt $h \leq n \cdot eps$.

Das Skalarprodukt ist also (bei $2n - 1$ Elementaroperationen) stabil im Sinne der Rückwärtsanalyse.

TODO : Beweis : S.51 ???

2.3 Anwendung auf lineare Gleichungssysteme (Konzepte der Kondition und Stabilität)

Nochmals die Fragestellung : Wann ist ein lineares Gleichungssystem $Ax = b$ lösbar ?

Im Unterschied zu Kapitel 1, wo diese Frage auf der fiktiven Basis der reellen Zahlen beantwortet wurde, wird hier nun die oben hergeleitete Fehlertheorie herangezogen.

Entsprechend wird die charakteristische Größe $\det A$ durch Konditionszahlen zu ersetzen sein.

Bei der normweisen Konditionsanalyse folgte

$$\frac{\|\tilde{x} - x\|}{\|x\|} \leq k_{rel} \cdot d \quad , \text{ wobei } d \text{ der relative Eingabefehler von } A \text{ (oder } b \text{) ist.} \quad (*)$$

Als Resultat für die relative Kondition folgte bei normweiser Betrachtung

$$k_{rel} = \|A^{-1}\| \frac{\|Ax\|}{\|x\|} \leq \|A^{-1}\| \cdot \|A\| = k(A)$$

und bei komponentenweiser Betrachtung

$$k_{rel} = \frac{\| |A^{-1}| \cdot |A| \cdot |x| \|_{\infty}}{\|x\|_{\infty}} \leq \| |A^{-1}| \cdot |A| \|_{\infty} = k_C(A) .$$

Desweiteren folgte für eine Matrix $A \in Mat_n(R)$, welche nicht die Nullmatrix ist, dass

$$\det A = 0 \Leftrightarrow k(A) = \infty$$

Falls $k(A) < \infty$, wäre das lineare Gleichungssystem also für jede rechte Seite im Prinzip eindeutig lösbar.

Andererseits lässt (*) nur dann ein numerisch brauchbares Resultat erwarten, falls $k_{rel} \cdot d$ hinreichend klein ist.

Darüber hinaus muss anstelle einer einzelnen Matrix A die Menge aller von A nicht zu unterscheidenden Matrizen, also z.B.

$$E := \left\{ \tilde{A} : \|\tilde{A} - A\| \leq d \cdot \|A\| \right\}$$

betrachtet werden.

Es bietet sich deshalb an, eine Matrix A mit der relativen Genauigkeit d „fast singulär“ oder „numerisch singulär“ zu nennen, falls die zugehörige Eingabemenge mindestens eine (exakt) singuläre Matrix enthält.

Def. „fast singuläre“ Matrix :

Eine Matrix A heißt *fast singulär* oder *numerisch singulär* bezüglich der Kondition $k(A)$, falls

$$d \cdot k(A) \geq 1 \quad , \text{ wobei } d \text{ die relative Genauigkeit der Matrix } A \text{ ist.}$$

Für die Rundungsfehler bei der Eingabe von A in einen Rechner wird z.B. angenommen, dass $d = eps$.

Beispiel :

$$Ax = b_1 \quad , \quad A = \begin{pmatrix} 1 & 1 \\ 0 & \epsilon \end{pmatrix} \quad , \quad b_1 = \begin{pmatrix} 2 \\ \epsilon \end{pmatrix} \quad , \text{ wobei } 0 < \epsilon \ll 1 \text{ die Eingabe darstellt.}$$

Die Matrix A und die rechte Seite b_1 haben den gemeinsamen Eingabewert ϵ ; sie sind also miteinander verkoppelt. Daher ist die Kondition der Matrix

$$k(A) = \|A^{-1}\|_{\infty} \|A\|_{\infty} = \frac{1}{\epsilon} \gg 1$$

nicht aussagekräftig für die Lösung des Problems (f_1, ϵ) , $f_1(\epsilon) = A^{-1}b_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$.

Die Lösung ist sogar unabhängig von ϵ , d.h. $f_1'(\epsilon) = 0$, und somit ist das Problem für alle ϵ gut konditioniert.

Betrachtet man jedoch eine von ϵ unabhängige rechte Seite $b_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$, so folgt für die Lösung

$$f_2(\epsilon) = A^{-1}b_2 = x_2 = \begin{pmatrix} 1/\epsilon \\ 1/\epsilon \end{pmatrix}$$

und für die komponentenweisen Konditionen

$$k_{abs} = \|f_2'(\epsilon)\| = \frac{1}{\epsilon^2} \quad \text{und} \quad k_{rel} = \frac{\| |f_2'(\epsilon)| \cdot |\epsilon| \|}{\| f_2(\epsilon) \|} = 1 .$$

In diesem Fall ist nur noch die Richtung der Lösung gut konditioniert, was sich in der relativen Kondition widerspiegelt.

2.3.1 Rückwärtsanalyse der Gauß-Elimination

Wie gesehen, ist die Berechnung des Skalarproduktes in Gleitkommaarithmetik stabil im Sinne der Rückwärtsanalyse. Die im ersten Kapitel besprochenen Algorithmen zur Lösung eines linearen Gleichungssystems erfordern im Einzelnen nur die Auswertung von Skalarprodukten bestimmter Zeilen und Spalten von Matrizen, so dass sich auf der Grundlage dieser Überlegungen eine Rückwärtsanalyse der Gaußschen Eliminationsmethode durchführen lässt.

Satz 2.7 (Rückwärtsanalyse der Gauß-Elimination)

Die Gleitkommarealisierung der Vorwärtssubstitution zur Lösung eines gestaffelten Gleichungssystems $Lx = b$ berechnet eine Lösung \hat{x} , so dass eine untere Dreiecksmatrix \hat{L} existiert mit

$$\hat{L}\hat{x} = b \quad \text{und} \quad |\hat{L} - L| \leq n \cdot \text{eps} \cdot |L|,$$

d.h. für den komponentenweisen relativen Rückwärtsfehler gilt $h \leq n \cdot \text{eps}$ und die Vorwärtssubstitution ist stabil im Sinne der Rückwärtsanalyse.

Beweis:

Der Algorithmus für die Vorwärtssubstitution lässt sich wie beim Skalarprodukt rekursiv formulieren:

$$l_{kk}x_k = b_k - \langle l^{(k-1)}, x^{(k-1)} \rangle \quad \text{für } k = 1, \dots, n, \quad \text{wobei wieder gilt}$$

$$l^{(k-1)} := (l_{k,1}, \dots, l_{k,k-1})^T \quad \text{und} \quad x^{(k-1)} := (x_1, \dots, x_{k-1})^T$$

Realisiert in Gleitpunktarithmetik ergibt sich die Rekursion

$$l_{kk} \cdot (1 + \mathbf{d}_k) \cdot (1 + \mathbf{e}_k) \cdot \hat{x}_k = b_k - \langle l^{(k-1)}, \hat{x}^{(k-1)} \rangle_{fl}$$

, wobei die \mathbf{d}_k und \mathbf{e}_k mit $|\mathbf{d}_k|, |\mathbf{e}_k| \leq \text{eps}$ die relativen Fehler der Multiplikation bzw. Addition beschreiben.

Über die Gleitkommarealisierung des Skalarproduktes ist nach Lemma 2.6 bekannt, dass

$$\langle l^{(k-1)}, \hat{x}^{(k-1)} \rangle_{fl} = \langle \hat{l}^{(k-1)}, \hat{x}^{(k-1)} \rangle$$

für einen Vektor $\hat{l}^{(k-1)} = (\hat{l}_{k,1}, \dots, \hat{l}_{k,k-1})^T$ mit $|\hat{l}^{(k-1)} - l^{(k-1)}| \leq (k-1) \cdot \text{eps} \cdot |l^{(k-1)}|$.

Setzt man nun noch $\hat{l}_{kk} := l_{kk} \cdot (1 + \mathbf{d}_k) \cdot (1 + \mathbf{e}_k)$, so gilt wie behauptet

$$\hat{L}\hat{x} = b \quad \text{und} \quad |\hat{L} - L| \leq n \cdot \text{eps} \cdot |L|. \quad \text{Q.e.d.}$$

Als erstes Resultat zur Rückwärtsanalyse der Gauß-Elimination wird nun die Qualität der LR-Zerlegung beurteilt:

Lemma 2.8 (Rückwärtsstabilität der LR-Zerlegung)

A besitze eine LR-Zerlegung. Dann berechnet die Gauß-Elimination \hat{L} und \hat{R} , so dass

$$\hat{L}\hat{R} = \hat{A} \quad \text{für eine Matrix } \hat{A} \text{ mit } |\hat{A} - A| \leq n \cdot |\hat{L}| \cdot |\hat{R}| \cdot \text{eps}.$$

ohne Beweis

Satz 2.9 (Rückwärtsstabilität für Gauß-Algorithmus ohne Pivotisierung)

A besitze eine LR-Zerlegung. Dann berechnet das Gaußsche Eliminationsverfahren für das Gleichungssystem $Ax = b$ eine Lösung \hat{x} mit

$$\hat{A}\hat{x} = b.$$

ohne Beweis

Daraus leitet sich nun die folgende Aussage über den normweisen Rückwärtsfehler der Gauß-Elimination mit Spaltenpivoting ab, welche auf *Wilkinson* zurückgeht :

Satz 2.10

Die Gauß-Elimination mit Spaltenpivoting für das Gleichungssystem $Ax = b$ berechnet ein \hat{x} , so dass

$$\hat{A}\hat{x} = b \quad \text{für eine Matrix } \hat{A} \text{ mit } \frac{\|\hat{A} - A\|_\infty}{\|A\|_\infty} \leq 2n^3 \cdot r_n(A) \cdot eps$$

, wobei $r_n(A) := \frac{\mathbf{a}_{\max}}{\max_{i,j} |a_{ij}|}$ und \mathbf{a}_{\max} der größte Betrag eines Elementes ist, welcher im Laufe der

Elimination in den Restmatrizen $A^{(1)} = A$ bis $A^{(n)} = R$ auftritt.

Beweis :

Im Folgenden werden die im Laufe der Gauß-Elimination mit Spaltenpivoting $PA = LR$ berechneten Größen mit $\hat{P}, \hat{L}, \hat{R}, \hat{x}$ bezeichnet.

Dann besitzt $\hat{P}\hat{A}$ eine LR-Zerlegung und nach Satz 2.9 existiert eine Matrix \hat{A} , so dass $\hat{A}\hat{x} = b$ und

$$\|\hat{A} - A\| \leq 2n \cdot \hat{P}^T \cdot \|\hat{L}\| \cdot \|\hat{R}\| \cdot eps .$$

Wendet man hierauf die Supremumsnorm an, folgt wegen $\|\hat{P}\|_\infty = 1$, dass

$$\|\hat{A} - A\|_\infty \leq 2n \cdot \|\hat{L}\|_\infty \cdot \|\hat{R}\|_\infty \cdot eps . \quad (*)$$

Die Spaltenpivotstrategie sorgt dafür, dass alle Komponenten von \hat{L} vom Betrag kleiner oder gleich eins sind, d.h.

$$\|\hat{L}\|_\infty \leq n .$$

Die Norm von \hat{R} lässt sich abschätzen durch

$$\|\hat{R}\|_\infty \leq n \cdot \max_{i,j} |\hat{r}_{ij}| \leq n \cdot \mathbf{a}_{\max} .$$

D.h. $\|\hat{A} - A\|_\infty \leq 2n^3 \cdot \mathbf{a}_{\max} \cdot eps$

Da $\max_{i,j} |a_{ij}| \leq \|A\|_\infty$, folgt die Behauptung $\frac{\|\hat{A} - A\|_\infty}{\|A\|_\infty} \leq 2n^3 \cdot \frac{\mathbf{a}_{\max}}{\max_{i,j} |a_{ij}|} \cdot eps$.

Q.e.d.

Die Frage nach der Stabilität lässt sich nach Satz 2.10 nicht eindeutig beantworten. Es hängt offenbar von der Zahl $r_n(A)$ ab, ob sich die Matrix für die Gauß-Elimination eignet oder nicht. Im Allgemeinen lässt sich diese Größe nur durch

$$r_n(A) \leq 2^{n-1}$$

abschätzen, wobei diese Abschätzung scharf ist, da die Grenze für die von *Wilkinson* angegebene Matrix

$$A_W = \begin{pmatrix} 1 & & & 1 \\ -1 & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ -1 & \cdots & -1 & 1 \end{pmatrix} \text{ angenommen wird.}$$

Die Gauß-Elimination mit Spaltenpivoting ist also über die ganze Menge der invertierbaren Matrizen betrachtet nicht stabil. Für spezielle Klassen von Matrizen sieht die Situation jedoch wesentlich besser aus :

Art der Matrix	Spaltenpivoting	$r_n \leq$
invertierbar	Ja	2^{n-1}
obere Hessenberg-Matrix	Ja	n
A oder A^T strikt diagonaldominant	Überflüssig	2
tridiagonal	Ja	2
symmetrisch positiv definit	Nein	1
Statistik	Ja	$n^{2/3}$ (im Mittel)

So ist die Gauß-Elimination für spd-Matrizen (Cholesky-Verfahren) ebenso wie für strikt diagonaldominante Matrizen stabil. Gestützt von statistischen Überlegungen kann man sagen, dass die Gauß-Elimination „in der Regel“, d.h. für üblicherweise in der Praxis auftretende Matrizen, stabil ist.

2.3.2 Beurteilung von Näherungslösungen

Zu einem linearen Gleichungssystem $Ax = b$ liege die Näherungslösung \tilde{x} vor.

Frage : Wie gut ist diese Lösung ?

⇒ Das Residuum $r(x) := b - Ax$ muss möglichst klein sein.

Problem : Die Norm $\|r(\tilde{x})\|$ kann durch eine Zeilenskalierung $Ax = b \rightarrow (D_z A)x = D_z b$ beliebig verändert werden, obwohl das Problem selbst dadurch nicht verändert wird.

Dies lässt sich auch an der Invarianz der Skeelschen Kondition $k_c(A)$ bezüglich Zeilenskalierung ablesen.

Das Residuum kann daher höchstens dann zur Beurteilung der Näherungslösung \tilde{x} herangezogen werden, wenn es eine problemspezifische Bedeutung besitzt. Besser geeignet ist das Konzept des Rückwärtsfehlers.

Für eine Näherungslösung \tilde{x} eines LGS lassen sich die normweisen und komponentenweisen Rückwärtsfehler direkt angeben.

Satz 2.11

Der normweise relative Rückwärtsfehler einer Näherungslösung \tilde{x} eines linearen Gleichungssystems

$f(A, b) = f^{-1}b$ bezüglich $\|(A, b)\| := \|A\| + \|b\|$ ist

$$h_N(\tilde{x}) = \frac{\|A\tilde{x} - b\|}{\|A\| \cdot \|\tilde{x}\| + \|b\|}.$$

Beweis : Fachliteratur

Satz 2.12 (Prager & Oettli)

Der komponentenweise relative Rückwärtsfehler einer Näherungslösung \tilde{x} von $Ax = b$ ist

$$h_c(\tilde{x}) = \max_i \frac{|A\tilde{x} - b|_i}{(|A| \cdot |\tilde{x}| + b)_i}.$$

Beweis : Buch S.58

3. Lineare Ausgleichsprobleme

Verallgemeinerung des Lösungs Begriffes auf rechteckige (überbestimmte) Gleichungssysteme.

3.1 Gaußsche Methode der kleinsten Fehlerquadrate

3.1.1 Problemstellung :

Gegeben seien m Messpunkte (t_i, b_i) , die z.B. die Zustände b_i eines Objektes zu den Zeiten t_i beschreiben. Es wird angenommen, dass diesen Messungen eine Gesetzmäßigkeit zugrunde liegt, so dass sich die Abhängigkeit von b von t durch eine Modellfunktion \mathbf{j} mit $b(t) = \mathbf{j}(t; x_1, \dots, x_n)$ ausdrücken lässt, wobei in die Modellfunktion n unbekannte Parameter x_i eingehen.

Beispiel : Ohmsches Gesetz

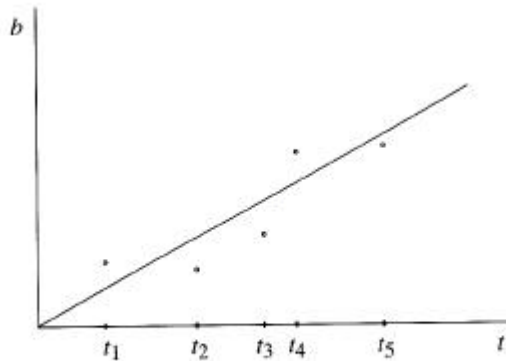


Abbildung 3.1. Lineare Ausgleichsrechnung beim Ohmschen Gesetz

Aufgabe : Eine Gerade durch den Nullpunkt legen, die dem Verlauf der Messungen „möglichst nahe“ kommt. Die Messungen sind in der Regel fehlerbehaftet, und auch die Modellfunktionen stellen stets nur eine ungefähre Beschreibung der Wirklichkeit dar. (So gilt das Ohmsche Gesetz annähernd nur in einem mittleren Temperaturbereich.) D.h. es wird gefordert, dass $b_i \approx \mathbf{j}(t_i; x_1, \dots, x_n)$ (für $i = 1, \dots, m$). Nun gibt es mehrere Möglichkeiten, die einzelnen Abweichungen $\Delta_i := b_i - \mathbf{j}(t_i; x_1, \dots, x_n)$ zu gewichten.

Aufgrund wahrscheinlichkeitstheoretischer Überlegungen wählte Gauß die Quadrate Δ_i^2 aus.

$$\text{D.h. } \Delta^2 := \sum_{i=1}^m \Delta_i^2 = \min . \quad (*)$$

Bemerkung : Äquivalenz zum Maximierungsproblem $e^{-\Delta^2} = \max$. Der Exponentialterm charakterisiert hierbei die Gaußsche Normalverteilung $f(x) = \frac{1}{s\sqrt{2\pi}} e^{-\frac{(x-m)^2}{2s^2}}$. \rightarrow Maximum-Likelihood-Methode.

In der Form (*) sind die Fehler der einzelnen Messungen alle gleich gewichtet. Normalerweise sind die Messungen (t_i, b_i) jedoch unterschiedlich genau. Zu jeder einzelnen Messung b_i gehört also eine absolute Messgenauigkeit bzw. Toleranz db_i .

$$\Rightarrow \text{Wichtung der einzelnen Fehler mit der Toleranz : } \sum_{i=1}^m \left(\frac{\Delta_i}{db_i} \right)^2 = \min .$$

3.1.2 Normalgleichungen

Geometrisch gesprochen sucht man bei der Lösung des linearen Ausgleichsproblems nach einem Punkt $z = Ax$ aus dem Bildraum $R(A)$ von A , der den kleinsten Abstand zu dem gegebenen Punkt b hat.

Für $m = 2$ und $n = 1$ ist $R(A) \subset \mathbb{R}^2$ entweder nur der Nullpunkt oder eine Gerade durch den Nullpunkt. (siehe Abbildung)

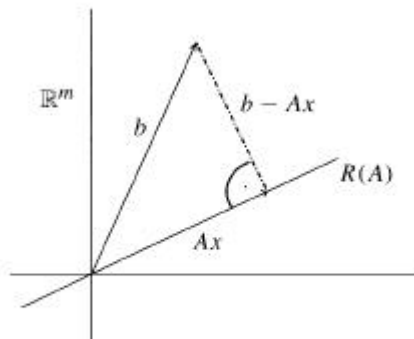


Abbildung 3.3. Projektion auf den Bildraum $R(A)$

Anschaulich ist klar, dass die Differenz $b - Ax$ gerade senkrecht auf dem Unterraum $R(A)$ stehen muss, damit der Abstand $\|b - Ax\|$ minimal ist.

Mit anderen Worten: Ax ist die **orthogonale** Projektion von b auf den Unterraum $R(A)$.

Satz 3.1

Sei V ein endlichdimensionaler euklidischer Vektorraum mit Skalarprodukt $\langle \cdot, \cdot \rangle$, $U \subset V$ ein Unterraum und

$$U^\perp = \{v \in V : \langle v, u \rangle = 0 \text{ für alle } u \in U\} \quad \text{sein orthogonales Komplement in } V.$$

Dann gilt für alle $v \in V$ bezüglich der von der Skalarprodukt induzierten Norm $\|v\| = \sqrt{\langle v, v \rangle}$, dass

$$\|v - u\| = \min_{u \in U} \|v - u\| \Leftrightarrow v - u \in U^\perp$$

Beweis: Buch S.70

Bemerkung 3.2

Damit ist die Lösung $u \in U$ von $\|v - u\| = \min$ eindeutig bestimmt und heißt **orthogonale Projektion** von v auf U .

Die Abbildung $P: V \rightarrow U$, $v \mapsto P \cdot v$ mit $\|v - Pv\| = \min_{u \in U} \|v - u\|$

ist linear und wird **orthogonale Projektion** von V auf U genannt.

Satz 3.3

Der Vektor $x \in \mathbb{R}^n$ ist genau dann Lösung des linearen Ausgleichsproblems $\|b - Ax\| = \min$, falls er die sogenannten **Normalgleichungen** $A^T Ax = A^T b$ erfüllt.
Insbesondere ist das lineare Ausgleichsproblem genau dann eindeutig lösbar, wenn der Rang von A maximal ist, d.h. $\text{Rang}(A) = n$.

Beweis :

Satz 3.1 angewandt auf $V = \mathbb{R}^m$ und $U = R(A)$ gilt

$$\begin{aligned} \|b - Ax\| = \min &\Leftrightarrow \langle b - Ax, Ax' \rangle = 0 && \text{für alle } x' \in \mathbb{R}^n \\ &\Leftrightarrow \langle A^T (b - Ax), x' \rangle = 0 && \text{für alle } x' \in \mathbb{R}^n \\ &\Leftrightarrow A^T (b - Ax) = 0 \\ &\Leftrightarrow A^T Ax = A^T b \end{aligned}$$

und daher die erste Aussage.

Der zweite Teil folgt nun aus der Tatsache, dass $A^T A$ genau dann invertierbar ist, wenn $\text{Rang}(A) = n$.

Q.e.d.

Bemerkung 3.4

Geometrisch besagen die Normalgleichungen gerade, dass $b - Ax$ eine Normale auf $R(A) \subset \mathbb{R}^m$ ist. Daher der Name.

3.1.3 Kondition des Projektionsproblems

Die relative Kondition des Projektionsproblems (P, b) bezüglich der Eingabe b hängt offenbar stark von dem Winkel ϑ ab, den b mit dem Unterraum V einschließt.

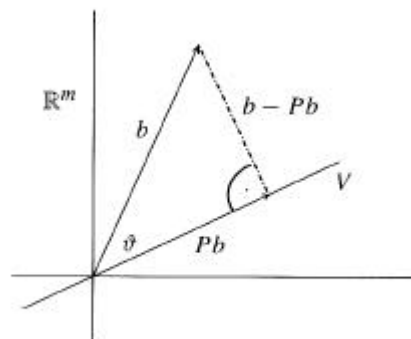


Abbildung 3.4. Projektion auf den Unterraum V

Ist der Winkel klein, d.h. $b \approx Pb$, so finden sich Störungen von b nahezu unverändert im Ergebnis Pb wieder.

Andererseits äußert sich eine kleine Störung von b in großen relativen Schwankungen von Pb , wenn b fast senkrecht auf V steht.

Lemma 3.5

Sei $P : R^m \rightarrow V$ die orthogonale Projektion auf einen Unterraum V des R^n .

Für die Eingabe b bezeichne J den Winkel zwischen b und V , d.h. $\sin J = \frac{\|b - Pb\|_2}{\|b\|_2}$.

Dann gilt für die relative Kondition des Problems (P, b) bezüglich der euklidischen Norm, dass

$$k = \frac{1}{\cos J} \|P\|_2.$$

Beweis :

Nach Pythagoras gilt $\|Pb\|^2 = \|b\|^2 - \|b - Pb\|^2$ und daher $\frac{\|Pb\|^2}{\|b\|^2} = 1 - \sin^2 J = \cos^2 J$.

Da P linear ist, ergibt sich daraus für die relative Kondition von (P, b) wie behauptet

$$k = \frac{\|b\|}{\|Pb\|} \|P'(b)\| = \frac{\|b\|}{\|Pb\|} \|P\| = \frac{1}{\sqrt{\frac{\|Pb\|^2}{\|b\|^2}}} = \frac{1}{\cos J} \|P\|.$$
 Q.e.d.

Lemma 3.6

Für eine Matrix $A \in Mat_{m,n}(R)$ von maximalem Rang $p = n$ gilt

$$k_2(A^T A) = k_2(A)^2.$$

Beweis :

Nach der Definition der Kondition einer rechteckigen Matrix gilt

$$k_2(A)^2 = \frac{\max_{\|x\|_2=1} \|Ax\|_2}{\min_{\|x\|_2=1} \|Ax\|_2} = \frac{\max_{\|x\|_2=1} \langle A^T Ax, x \rangle}{\min_{\|x\|_2=1} \langle A^T Ax, x \rangle} = \frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)} = k_2(A^T A)$$
 Q.e.d.

Satz 3.7 (Kondition des linearen Ausgleichsproblems)

Sei $A \in Mat_{m,n}(R)$, $m \geq n$ eine Matrix von vollem Spaltenrang, $b \in R^n$, und x die (eindeutige) Lösung des

Linearen Ausgleichsproblems $\|b - Ax\|_2 = \min$.

Es sei vorausgesetzt, dass $x \neq 0$ und J bezeichne den Winkel zwischen b und dem Bildraum $R(A)$ von A , d.h.

$\sin J = \frac{\|b - Ax\|_2}{\|b\|_2} = \frac{\|r\|_2}{\|b\|_2}$ mit dem Residuum $r = b - Ax$. Dann gilt für die relative Kondition von x in der

euklidischen Norm

- bezüglich Störungen in b

$$k \leq \frac{k_2(A)}{\cos J}$$

- bezüglich Störungen in A

$$k \leq k_2(A) + k_2(A)^2 \cdot \tan J.$$

Beweis : Buch S.73/74

3.1.4 Lösung der Normalgleichungen

Geht man davon aus, dass das lineare Ausgleichsproblem eindeutig lösbar ist, d.h. $\text{Rang}(A) = n$, so ist $A^T A$ eine spd-Matrix. Daher bietet sich das Cholesky-Verfahren an. Für den Aufwand zur Lösung des linearen Ausgleichsproblems mit Hilfe der Normalgleichungen ergibt sich dann (Anzahl der Multiplikationen) :

- a) Berechnung von $A^T A$: $\sim \frac{1}{2} n^2 m$
- b) Cholesky-Zerlegung von $A^T A$: $\sim \frac{1}{6} n^3$

Für $m \gg n$ überwiegt der Anteil von a), so dass der Aufwand insgesamt

$$\sim \frac{1}{2} n^2 m \quad \text{für } m \gg n \quad \text{und} \quad \sim \frac{2}{3} n^3 \quad \text{für } m \approx n \quad \text{beträgt.}$$

Die Lösung linearer Ausgleichsprobleme über die Normalgleichungen mit Hilfe der Cholesky-Zerlegung kann nur für große Residuen empfohlen werden, da bei kleinen Residuen der Übergang zu den Normalgleichungen eine Verschlechterung der Kondition bedeutet.

3.2 Orthogonalisierungsverfahren

Jeder Eliminationsprozess für lineare Gleichungssysteme, wie z.B. die Gauß-Elimination, lässt sich formal in der Form

$$A \xrightarrow{f_1} B_1 A \xrightarrow{f_2} B_2 B_1 A \xrightarrow{f_3} \dots \xrightarrow{f_k} B_k \dots B_1 A = R$$

darstellen, wobei die Matrizen B_j die Operationen auf der Matrix A beschreiben.

Bei der rekursiven Stabilitätsanalyse in Kapitel 2 folgte, dass die Stabilitätsindikatoren der Teilschritte eines Algorithmus verstärkt werden durch die Konditionen aller folgenden Teilschritte.

Beim oben beschriebenen Eliminationsprozess sind z.B. die Konditionen der Eliminationsmatrizen $B_j = L_j$ der Gaußschen Dreieckszerlegung nicht nach oben beschränkt, so dass es dabei zu Instabilitäten kommen kann.

Wählt man hingegen statt der L_j orthogonale Transformationen Q_j für die Elimination, so gilt bzgl. der euklidischen Norm, dass

$$\kappa(Q_j) = \|Q_j\|_2 \cdot \|Q_j^{-1}\|_2 = \|Q_j\|_2 \cdot \|Q_j^T\|_2 = 1$$

- + Diese sogenannten **Orthogonalisierungsverfahren** sind also auf jeden Fall stabil.
- Leider ist diese Stabilität mit einem etwas höheren Aufwand verbunden als z.B. bei der Gauß-Elimination.
- + Aufgrund der Invarianz der euklidischen Norm bezüglich orthogonaler Transformationen eignen sich Orthogonalisierungsverfahren zur Lösung linearer Ausgleichsprobleme.

Angenommen, man hätte die Matrix $A \in \text{Mat}_{m,n}(R)$ mit $m \geq n$ mit einer orthogonalen Matrix $Q \in O(m)$ auf obere

Dreiecksgestalt $Q^T A = \begin{bmatrix} * & \dots & * \\ & \ddots & \vdots \\ & & * \end{bmatrix} = \begin{bmatrix} R \\ 0 \end{bmatrix}$ gebracht. ($R =$ obere Dreiecksmatrix)

Dann lässt sich die Lösung des linearen Ausgleichsproblems $\|b - Ax\| = \min$ als Alternative zur Lösung der Normalgleichungen auch wie folgt bestimmen :

Satz 3.8

Sei $A \in \text{Mat}_{m,n}(\mathbb{R})$ mit $m \geq n$ und von maximalem Rang, $b \in \mathbb{R}^m$, und $Q \in O(m)$ eine orthogonale Matrix mit

$$Q^T A = \begin{pmatrix} R \\ 0 \end{pmatrix} \quad \text{und} \quad Q^T b = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix}, \quad \text{wobei } b_1 \in \mathbb{R}^n, b_2 \in \mathbb{R}^{m-n} \text{ und } R \in \text{Mat}_n(\mathbb{R}) \text{ eine (invertierbare)}$$

obere Dreiecksmatrix ist. Dann ist $x = R^{-1}b_1$ die Lösung des linearen Ausgleichsproblems $\|b - Ax\| = \min$.

Beweis :

Da $Q \in O(m)$, gilt für alle $x \in \mathbb{R}^m$

$$\|b - Ax\|^2 = \|Q^T(b - Ax)\|^2 = \left\| \begin{pmatrix} b_1 - Rx \\ b_2 \end{pmatrix} \right\|^2 = \|b_1 - Rx\|^2 + \|b_2\|^2 \geq \|b_2\|^2.$$

Wegen $\text{Rang}(A) = \text{Rang}(R) = n$ ist R invertierbar. Der erste Summand $\|b_1 - Rx\|^2$ verschwindet daher genau für $x = R^{-1}b_1$. Das Residuum $r := b - Ax$ verschwindet im Allgemeinen nicht, und es ist $\|r\| = \|b_2\|$. **Q.e.d.**

Die in Frage kommenden orthogonalen Transformationen lassen sich für $m = 2$ leicht geometrisch ableiten, nämlich als Drehungen (Rotationen) bzw. Spiegelungen (Reflexionen).

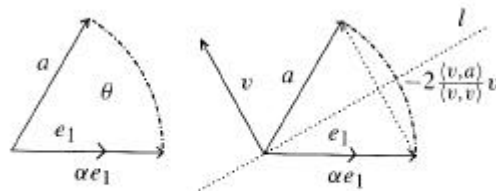


Abbildung 3.5. Drehung bzw. Spiegelung von a auf αe_1

Will man den Vektor $a \in \mathbb{R}^2$ mit einer orthogonalen Transformation auf ein Vielfaches $\mathbf{a} \cdot e_1$ des ersten Einheitsvektors abbilden, so folgt $\mathbf{a} = \|a\|$.

- 1. Möglichkeit : a um den Winkel q auf $\mathbf{a} \cdot e_1$ drehen, d.h.

$$a \mapsto \mathbf{a} \cdot e_1 = Q \cdot a \quad \text{mit } Q := \begin{pmatrix} \cos q & \sin q \\ -\sin q & \cos q \end{pmatrix}$$

- 2. Möglichkeit : a an der auf dem Vektor v senkrecht stehenden Geraden l spiegeln, d.h.

$$a \mapsto \mathbf{a} \cdot e_1 = a - 2 \frac{\langle v, a \rangle}{\langle v, v \rangle} v, \quad \text{wobei } v \text{ kollinear zur Differenz } a - \mathbf{a} \cdot e_1 \text{ ist.}$$

3.2.1 Givens-Rotationen

Definition : Als Givens-Rotationen bezeichnet man Matrizen der Form

$$\Omega_{kl} := \begin{bmatrix} I & & & \\ & c & s & \\ & & I & \\ & -s & c & \\ & & & I \end{bmatrix} \begin{matrix} \leftarrow k \\ \\ \leftarrow l \end{matrix} \in \text{Mat}_m(\mathbb{R})$$

, wobei I jeweils die Einheitsmatrix der passenden Dimension ist und $c^2 + s^2 = 1$. (Dabei sollen c und s an $\cos q$ und $\sin q$ erinnern.)

Geometrisch beschreibt die Matrix eine Drehung um den Winkel q in der $(k-l)$ -Ebene.

- Wendet man Ω_{kl} auf einen Vektor $x \in R^m$ an, so folgt

$$x \mapsto y = \Omega_{kl}x \quad \text{mit} \quad y_i = (\Omega_{kl}x)_i = \begin{cases} cx_k + sx_l & , \text{ falls } i = k \\ -sx_k + cx_l & , \text{ falls } i = l \\ x_i & , \text{ falls } i \neq k, l \end{cases} \quad (*)$$

- Multipliziert man eine Matrix $A = [A_1, \dots, A_n] \in Mat_{m,n}(R)$ von links mit Ω_{kl} , so operiert die Givens-Rotation auf den Spalten, d.h.

$$\Omega_{kl}A = [\Omega_{kl}A_1, \dots, \Omega_{kl}A_n].$$

Es werden daher wegen (*) nur die beiden Zeilen k und l der Matrix A verändert.

Dies ist insbesondere wichtig, wenn man bei der Transformation möglichst Besetzungsstrukturen der Matrix erhalten möchte.

Wie sind nun die Koeffizienten c und s zu bestimmen, um eine Komponente x_l des Vektors x zu eliminieren?

Da Ω_{kl} nur auf der (k, l) -Ebene operiert, genügt es, das Prinzip an dem Fall $m = 2$ zu erläutern.

Mit $x_k^2 + x_l^2 \neq 0$ und $c^2 + s^2 = 1$ gilt:

$$\begin{pmatrix} c & s \\ -s & c \end{pmatrix} \begin{pmatrix} x_k \\ x_l \end{pmatrix} = \begin{pmatrix} cx_k + sx_l \\ -sx_k + cx_l \end{pmatrix} = \begin{pmatrix} r \\ 0 \end{pmatrix}$$

$$\Rightarrow \text{II.} : x_k = \frac{cx_l}{s}$$

$$\Rightarrow \text{in I. einsetzen} : c^2 \frac{x_l}{s} + sx_l = r \quad \Rightarrow x_l(c^2 + s^2) = rs \quad \Rightarrow s = \frac{x_l}{r}$$

$$\Rightarrow \text{in II. einsetzen} : -\frac{x_l x_k}{r} + cx_l = 0 \quad \Rightarrow c = \frac{x_k}{r}$$

- für $|x_l| > |x_k|$: $t := \frac{x_k}{x_l}$, $s := \frac{1}{\sqrt{1+t^2}}$, $c := s \cdot t$
- für $|x_l| \leq |x_k|$: $t := \frac{x_l}{x_k}$, $c := \frac{1}{\sqrt{1+t^2}}$, $s := c \cdot t$

Damit vermeidet man zugleich Exponentenüberlauf.

Beispiel :

Spaltenweise werden nun die von Null verschiedenen Matrixkomponenten unterhalb der Diagonalen eliminiert.

(Die Indexpaare über den Pfeilen geben die Indizes der ausgeführten Givens-Rotation an.)

$$A = \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \end{bmatrix} \xrightarrow{(5,4)} \begin{bmatrix} * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ * & * & * & * \\ 0 & * & * & * \end{bmatrix} \xrightarrow{(4,3)} \dots \xrightarrow{(2,1)} \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \end{bmatrix} \xrightarrow{(5,4)} \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \\ 0 & * & * & * \\ 0 & 0 & * & * \end{bmatrix} \xrightarrow{(4,3)} \dots \xrightarrow{(5,4)} \begin{bmatrix} * & * & * & * \\ 0 & * & * & * \\ 0 & 0 & * & * \\ 0 & 0 & 0 & * \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

Der Aufwand für die QR-Zerlegung einer vollbesetzten Ausgangsmatrix $A \in \text{Mat}_{m,n}$ beträgt :

- $m \approx n$: $\sim \frac{n^2}{2}$ Quadratwurzeln und $\sim \frac{4n^3}{3}$ Multiplikationen
- $m \gg n$: $\sim mn$ Quadratwurzeln und $\sim 2mn^2$ Multiplikationen

Für $m = n$ erhält man somit eine Alternative zur Gaußschen Dreieckszerlegung aus Kapitel 1.

Die größere Stabilität muss jedoch mit einem erheblich höheren Aufwand, $\sim \frac{4n^3}{3}$ Multiplikationen gegenüber $\sim \frac{n^3}{3}$ bei Gauß, erkauft werden. Zu beachten ist aber, dass der Vergleich für dünnbesetzte Matrizen wesentlich günstiger ausfällt. Praktisch verwendet man sogenannte schnelle Givens-Rotationen, welche die Auswertung der Quadratwurzeln vermeiden und eine QR-Zerlegung einer zeilenskalierten Matrix DA berechnen.

3.2.2 Householder-Reflexionen

Definition : Als Householder-Reflexionen bezeichnet man Matrizen $Q \in \text{Mat}_n(\mathbb{R})$ der Form

$$Q = I - 2 \frac{vv^T}{v^T v} \quad \text{mit } v \in \mathbb{R}^n.$$

Diese Matrizen beschreiben die Reflexion an der auf v senkrecht stehenden Ebene. (siehe vorige Abbildung)
Insbesondere hängt Q nur von der Richtung von v ab.

Eigenschaften :

- Q ist symmetrisch, d.h. $Q^T = Q$
- Q ist orthogonal, d.h. $QQ^T = Q^T Q = I$
- Q ist involutorisch, d.h. $Q^2 = I$

- Wendet man Q auf einen Vektor $y \in \mathbb{R}^n$ an, so gilt

$$y \mapsto Q \cdot y = \left(I - 2 \frac{vv^T}{v^T v} \right) y = y - 2 \frac{\langle v, y \rangle}{\langle v, v \rangle} v.$$

Soll y auf ein Vielfaches $\mathbf{a} \cdot e_1$ des ersten Einheitsvektors e_1 abgebildet werden, d.h.

$$\mathbf{a} \cdot e_1 = y - 2 \frac{\langle v, y \rangle}{\langle v, v \rangle} v \in \text{span}(e_1), \text{ so folgt}$$

$$|\mathbf{a}| = \|y\|_2 \quad \text{und} \quad v \in \text{span}(y - \mathbf{a} \cdot e_1)$$

Daher lässt sich Q bestimmen durch $v := y - \mathbf{a} \cdot e_1$ mit $\mathbf{a} = \pm \|y\|_2$.

Um Auslöschung bei der Berechnung von $v = (y_1 - \mathbf{a}, y_2, \dots, y_n)^T$ zu vermeiden, wählt man $\mathbf{a} := -\text{sgn}(y_1) \cdot \|y\|_2$.

Wegen $\langle v, v \rangle = \langle y - \mathbf{a} \cdot e_1, y - \mathbf{a} \cdot e_1 \rangle = \|y\|_2^2 - 2\mathbf{a}\langle y, e_1 \rangle + \mathbf{a}^2 = -2\mathbf{a}(y_1 - \mathbf{a})$ lässt sich Qx für beliebige $x \in \mathbb{R}^n$ am einfachsten berechnen durch

$$Qx = x - 2 \frac{\langle v, x \rangle}{\langle v, v \rangle} v = x + \frac{\langle v, x \rangle}{\mathbf{a}(x_1 - \mathbf{a})} v$$

- Transformiert man mit Hilfe der Householder-Reflexionen eine Matrix $A = [A_1, \dots, A_n] \in \text{Mat}_{m,n}(\mathbb{R})$ in eine obere Dreiecksmatrix, so werden sukzessive die Elemente unterhalb der Diagonalen eliminiert.

Im ersten Schritt : $A \rightarrow A' := Q_1 A = \begin{bmatrix} \mathbf{a}_1 & & & \\ 0 & & & \\ \vdots & A_2' & \dots & A_n' \\ 0 & & & \end{bmatrix},$

wobei $Q_1 = I - 2 \frac{v_1 v_1^T}{v_1^T v_1}$ mit $v_1 := A_1 - \mathbf{a}_1 e_1$ und $\mathbf{a}_1 := -\text{sgn}(a_{11}) \cdot \|A_1\|_2$.

Nach dem k-ten Schritt wurde somit die Ausgangsmatrix A bis auf eine Restmatrix $T^{(k+1)} \in \text{Mat}_{m-k, n-k}(R)$ auf obere Dreiecksgestalt gebracht.

$$A^{(k)} = \begin{bmatrix} * & \dots & \dots & \dots & * \\ & \ddots & & & \vdots \\ & & * & \dots & * \\ & & 0 & & \\ & & \vdots & T^{(k+1)} & \\ & & 0 & & \end{bmatrix}$$

Bildet man nun die orthogonale Matrix $Q_{k+1} = \left[\begin{array}{c|c} I_k & 0 \\ \hline 0 & \bar{Q}_{k+1} \end{array} \right],$

wobei $\bar{Q}_{k+1} \in O(m-k)$ wie im ersten Schritt mit $T^{(k+1)}$ anstelle von A konstruiert wird, so lässt sich die nächste Spalte unterhalb der Diagonalen eliminieren.

Insgesamt erhält man so nach $p = \min(m-1, n)$ Schritten die obere Dreiecksmatrix

$$R = Q_p \dots Q_1 A$$

und somit wegen $Q_i^2 = I$ die Zerlegung

$$A = QR \quad \text{mit} \quad Q = Q_1 \dots Q_p .$$

Berechnet man nun also die Lösung des linearen Ausgleichsproblems $\|b - Ax\| = \min$ nach Satz 3.8, indem man mit Hilfe der Householder-Reflexionen $Q_j \in O(m)$ die QR-Zerlegung der Matrix $A \in \text{Mat}_{m,n}(R)$ mit $m \geq n$ berechnet, so gelangt man zu folgendem Verfahren :

- 1) $A = QR$, QR-Zerlegung mit Householder-Reflexionen
- 2) $(b_1, b_2)^T = Q^T b$ mit $b_1 \in R^n$ und $b_2 \in R^{m-n}$, Transformation von b
- 3) $Rx = b_1$, Auflösung des gestaffelten Systems

Speicherschema :

Speichert man die Diagonalelemente $r_{ii} = \mathbf{a}_i$ für $i = 1, \dots, p$ in einem separaten Vektor, so finden die Householder-Vektoren v_1, \dots, v_p in der unteren Hälfte von A Platz.

Eine andere Möglichkeit besteht darin, die Householder-Vektoren so zu normieren, dass die erste Komponente $\langle v_i, e_i \rangle$ jeweils 1 ist und nicht abgespeichert werden muss.

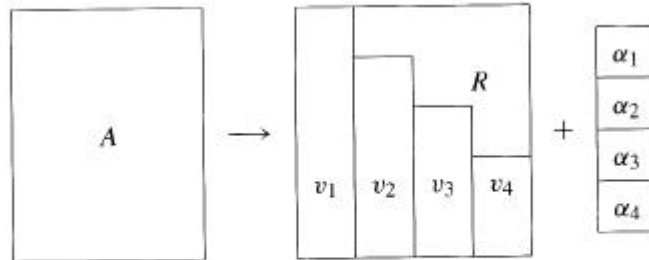


Abbildung 3.6. Speicheraufteilung bei der QR-Zerlegung mit Householder-Reflexionen für $m = 5$ und $n = 4$

Aufwand :

- a) $m \gg n$: $\sim 2n^2m$ Multiplikationen
- b) $m \approx n$: $\frac{2}{3}n^3$ Multiplikationen

Für $m \approx n$ benötigt man also in etwa den gleichen Aufwand wie bei dem Cholesky-Verfahren für Normalgleichungen. Für $m \gg n$ schneidet die QR-Zerlegung um Faktor 2 schlechter ab, hat aber die oben diskutierten Stabilitätsvorteile.

3.3 Verallgemeinerte Inverse

Die nach Satz 3.3 eindeutig bestimmte Lösung x des linearen Ausgleichsproblems $\|b - Ax\| = \min$ für

$A \in \text{Mat}_{m,n}(R)$, $m \geq n$ und $\text{Rang}(A) = n$ wird formal mit $x = A^+b$ bezeichnet.

Aus den Normalgleichungen folgt unter obigen Voraussetzungen, dass

$$\begin{aligned} A^T Ax = A^T b &\Rightarrow A^T AA^+b = A^T b \Rightarrow A^+b = (A^T A)^{-1} A^T b \\ &\Rightarrow A^+ = (A^T A)^{-1} A^T \end{aligned}$$

Da $A^+A = (A^T A)^{-1} A^T A = I$ gerade die Identität ergibt, nennt man A^+ auch die **Pseudoinverse** von A .

Bei beliebigen Matrizen $A \in \text{Mat}_{m,n}(R)$ ist die Lösung von $\|b - Ax\| = \min$ jedoch im Allgemeinen nicht mehr eindeutig bestimmt.

Es bezeichne $\bar{P} : R^m \rightarrow R(A) \subset R^m$ die orthogonale Projektion von R^m auf den Bildraum $R(A)$.

Dann bilden die Lösungen nach Satz 3.1 einen affinen Unterraum

$$L(b) := \{x \in R^n : \|b - Ax\| = \min\} = \{x \in R^n : Ax = \bar{P}b\}.$$

Um dennoch Eindeutigkeit zu erzwingen, wählt man die bezüglich der euklidischen Norm kleinste Lösung $x \in L(b)$, welche wieder mit $x = A^+b$ bezeichnet werde.

Nun ist x gerade die orthogonale Projektion des Ursprungs $0 \in \mathbb{R}^n$ auf den affinen Unterraum $L(b)$. (siehe Abbildung)

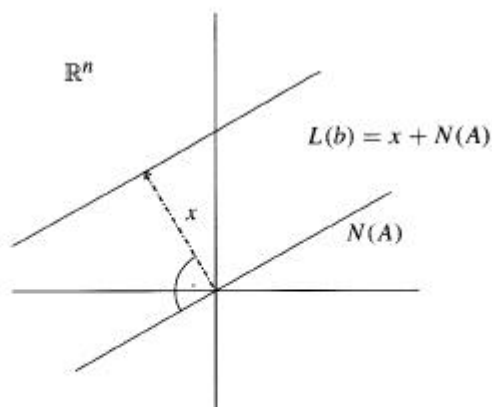


Abbildung 3.7. „Kleinste“ Lösung des linearen Ausgleichsproblems als Projektion von 0 auf $L(b)$

Ist $\bar{x} \in L(b)$ irgendeine Lösung von $\|b - Ax\| = \min$, so erhält man alle Lösungen, indem man den Nullraum $N(A)$ von A um \bar{x} verschiebt, d.h.

$$L(b) = \bar{x} + N(A) .$$

Daher muss die kleinste Lösung x senkrecht auf dem Nullraum $N(A)$ stehen, d.h. x ist der eindeutig bestimmte Vektor $x \in N(A)^\perp$ mit $\|b - Ax\| = \min$.

Def. Pseudoinverse :

Die Pseudoinverse einer Matrix $A \in \text{Mat}_{m,n}(\mathbb{R})$ ist die Matrix $A^+ \in \text{Mat}_{n,m}(\mathbb{R})$, so dass für alle $b \in \mathbb{R}^m$ der Vektor $x = A^+b$ die kleinste Lösung von $\|b - Ax\| = \min$ ist, d.h.

$$A^+b \in N(A)^\perp \quad \text{und} \quad \|b - AA^+b\| = \min .$$

Die Situation lässt sich durch das folgende kommutative Diagramm darstellen :

$$\begin{array}{ccc}
 \mathbb{R}^n & \xrightarrow{A} & \mathbb{R}^m \\
 & \xleftarrow{A^+} & \\
 P = A^+A & \downarrow \uparrow i & i \uparrow \downarrow \bar{P} = AA^+ \\
 R(A^+) = N(A)^\perp & \cong & R(A)
 \end{array}$$

Satz 3.9 (Penrose-Axiome)

Die Pseudoinverse $A^+ \in Mat_{n,m}(R)$ einer Matrix $A \in Mat_{m,n}(R)$ ist eindeutig charakterisiert durch folgende Eigenschaften :

- 1.) $(A^+A)^T = A^+A$
- 2.) $(AA^+)^T = AA^+$
- 3.) $A^+AA^+ = A^+$
- 4.) $AA^+A = A$

Beweis : Buch S.87

Bemerkung 3.10 :

Gilt nur ein Teil der Penrose-Axiome, so spricht man von einer **verallgemeinerten Inversen**.

Nun die Herleitung, wie sich die kleinste Lösung $x = A^+b$ für eine beliebige Matrix $A \in Mat_{n,m}(R)$ und $b \in R^m$ mit Hilfe der QR-Zerlegung berechnen lässt.

Sei $p := Rang(A) \leq \min(m, n)$ der Rang der Matrix A .

Zur Vereinfachung wird auf Permutationen verzichtet und A durch orthogonale Transformationen $Q \in O(m)$ auf obere

Dreiecksgestalt gebracht, d.h. $QA = \left[\begin{array}{c|c} R & S \\ \hline 0 & 0 \end{array} \right]$,

hier ist $R \in Mat_p(R)$ eine invertierbare obere Dreiecksmatrix und $S \in Mat_{p,n-p}(R)$.

Zerlegung der Vektoren x und Qb analog in

$$x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \quad \text{mit } x_1 \in R^p \quad \text{und} \quad x_2 \in R^{n-p}$$

$$Qb = \begin{pmatrix} b_1 \\ b_2 \end{pmatrix} \quad \text{mit } b_1 \in R^p \quad \text{und} \quad b_2 \in R^{m-p}$$

Lemma 3.11

Mit den obigen Bezeichnungen ist x genau dann eine Lösung von $\|b - Ax\| = \min$, falls

$$x_1 = R^{-1}b_1 - R^{-1}Sx_2 .$$

Beweis :

Aufgrund der Invarianz der euklidischen Norm unter orthogonalen Transformationen gilt

$$\|b - Ax\|^2 = \|Qb - QAx\|^2 = \|Rx_1 + Sx_2 - b_1\|^2 + \|b_2\|^2 .$$

Der Ausdruck wird genau dann minimal, wenn $Rx_1 + Sx_2 - b_1 = 0$.

Q.e.d.

Der Fall $p = Rang(A) = n$ entspricht dem bereits behandelten überbestimmten Gleichungssystem mit vollem Rang.

Die Matrix S verschwindet, und man erhält wie in Satz 3.8 die Lösung $x = x_1 = R^{-1}b_1$.

Für die Lösung im rangdefekten Fall $p < n$ gilt :

Lemma 3.12

Sei $p < n$, $V := R^{-1}S \in \text{Mat}_{p, n-p}(R)$ und $u := R^{-1}b_1 \in R^p$.

Dann ist die kleinste Lösung x von $\|b - Ax\| = \min$ gegeben durch $x = (x_1, x_2) \in R^p \times R^{n-p}$ mit

$$(I + V^T V) \cdot x_2 = V^T u \quad \text{und} \quad x_1 = u - Vx_2.$$

Beweis : Buch S.88

Für die Berechnung von x_2 kann man die Cholesky-Zerlegung benutzen, da $I + V^T V$ eine spd-Matrix ist.

Zusammengefasst erhält man folgendes Verfahren zur Berechnung der kleinsten Lösung $x = A^+ b$ von $\|b - Ax\| = \min$:

Algorithmus Pseudoinverse über QR-Zerlegung :

- 1.) QR-Zerlegung von A mit $p = \text{Rang}(A)$, wobei $Q \in O(m)$, $R \in \text{Mat}_p(R)$ obere Dreiecksmatrix und $S \in \text{Mat}_{p, n-p}(R)$
- 2.) Berechnung von $V \in \text{Mat}_{p, n-p}(R)$ aus $R \cdot V = S$
- 3.) Cholesky-Zerlegung von $I + V^T V$ ($I + V^T V = LL^T$), wobei $L \in \text{Mat}_{n-p}(R)$ untere Dreiecksmatrix
- 4.) $(b_1, b_2)^T := Qb$ mit $b_1 \in R^p$, $b_2 \in R^{n-p}$
- 5.) Berechnung von $u \in R^p$ aus $R \cdot u = b_1$
- 6.) Berechnung von $x_2 \in R^{n-p}$ aus $LL^T x_2 = V^T u$
- 7.) Setze $x_1 := u - Vx_2$

$$\text{Dann ist } x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = A^+ b.$$

Bemerkung : Die Schritte 1.) bis 3.) müssen für verschiedene rechte Seiten b nur einmal durchgeführt werden.

4. Iterative Lösung linearer Gleichungssysteme

Die bisher beschriebenen sogenannten **direkten Verfahren zur Lösung eines LGS** (Gauß-Elimination, Cholesky-Zerlegung, QR-Zerlegung mit Householder- oder Givens-Transformationen) haben folgende Eigenschaften :

- 1.) Die Verfahren gehen von beliebigen (bei der Cholesky-Zerlegung symmetrischen) vollbesetzten Matrizen aus.
- 2.) Der Aufwand zur Lösung des Gleichungssystems liegt bei $O(n^3)$ (Multiplikationen).

Dem gegenüber stehen jedoch in vielen Fällen große aber schwach besetzte Koeffizientenmatrizen (sparse matrix).

Auch haben große Gleichungssysteme ($n > 10^4$) häufig eine klare strukturierte Koeffizientenmatrix (z.B. Bandmatrizen, Blockmatrizen).

Die direkten Verfahren sind zur Behandlung derartiger Probleme ungeeignet; sie nutzen die spezielle Struktur nicht aus und dauern bei weitem zu lange.

Ansatz 1 :

Die Spezialstruktur der Matrix, insbesondere ihre Besetzungsstruktur (sparsity pattern) so weit wie möglich in direkten Verfahren ausnutzen.

Die Givens-Rotationen operieren jeweils nur auf zwei Zeilen (von links) oder Spalten (von rechts) einer Matrix und sind daher geeignet, eine Besetzungsstruktur weitgehend zu erhalten.

Die Householder-Transformationen dagegen sind dazu gänzlich ungeeignet. Sie zerstören bereits bei einem Schritt jedes Muster der Ausgangsmatrix.

Im Allgemeinen am schonendsten geht die Gauß-Elimination mit der Besetzungsstruktur von Matrizen um. Typischerweise wird dabei abwechselnd Spaltenpivotsuche mit evt. Zeilentausch und Zeilenpivotsuche mit evt. Spaltentausch ausgeführt, je nachdem, welche Strategie die meisten Nullelemente schont.

Ansatz 2 :

Iterative Verfahren zur Approximation der Lösung x .

Sinnvoll, da man in der Regel an der Lösung x nur bis auf eine vorgegebene Genauigkeit ϵ interessiert ist, die von der Genauigkeit der Eingabedaten abhängt.

Ziel ist die Konstruktion einer Iterationsvorschrift $x_{k+1} = \Phi(x_0, \dots, x_k)$, so dass

- a) die Folge $\{x_k\}$ der Iterierten möglichst schnell gegen die Lösung x konvergiert,
- b) x_{k+1} mit möglichst geringem Aufwand aus x_0, \dots, x_k berechnet werden kann.

Bei der zweiten Forderung verlangt man meist, dass die Auswertung von Φ nicht wesentlich mehr Aufwand erfordert als eine einfache Matrix-Vektor-Multiplikation.

Desweiteren beträgt der Aufwand für dünnbesetzte Matrizen $O(n)$ und nicht $O(n^2)$ (wie bei vollbesetzten Matrizen), da häufig die Anzahl der von Null verschiedenen Elemente in einer Zeile unabhängig von der Dimension n des Problems ist.

4.1 Klassische Iterationsverfahren

Den meisten klassischen Iterationsverfahren liegt die Idee der Fixpunktiteration zugrunde.

Einschub Fixpunktiteration

Äquivalente Umformung der Gleichung $f(x) = 0$ in eine Fixpunktgleichung

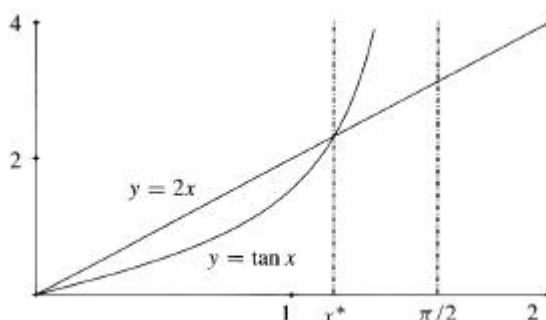
$$f(x) = x$$

und mit Hilfe der Iterationsvorschrift $x_{k+1} = f(x_k)$ mit $k = 0, 1, \dots$

für einen gegebenen Startwert x_0 eine Folge $\{x_0, x_1, \dots\}$ konstruieren, welche gegen einen Fixpunkt x^* mit $f(x^*) = x^*$ konvergiert, der Lösung der Gleichung ist, d.h. $f(x^*) = 0$.

Beispiel Fixpunktiteration :

$$f(x) := 2x - \tan x = 0$$



Anhand der Abbildung $(x^* \approx 1.2)$ wählt man den Startwert $x_0 = 1.2$.

Umformen der Gleichung in eine Fixpunktgleichung :

$$x = \frac{1}{2} \tan x =: f_1(x_k) \quad \text{oder} \quad x = \arctan(2x) =: f_2(x_k)$$

Es ergeben sich die Zahlenwerte

k	$x_{k+1} = \frac{1}{2} \tan x_k$	$x_{k+1} = \arctan(2x_k)$
0	1.2	1.2
1	1.2860	1.1760
2	1.70... > $\frac{p}{2}$	1.1687
3		1.1665
4		1.1658
5		1.1656
6		1.1655
7		1.1655

Die erste Folge divergiert ($\tan x$ hat einen Pol bei $\frac{p}{2}$ und $x_2 > \frac{p}{2}$), während die zweite konvergiert.

D.h. nicht jede naiv konstruierte Fixpunktiteration konvergiert.

Einschub Ende

Für ein Fixpunktverfahren $x_{k+1} = f(x_k)$ zur Lösung eines LGS $Ax = b$ muss eine Iterationsfunktion f derart konstruiert werden, dass sie genau einen Fixpunkt x^* besitzt und dieser gerade die exakte Lösung $x^* = x$ von $Ax = b$ ist.

⇒ Umformen der Gleichung $Ax = b$ in eine Fixpunktgleichung :

$$\begin{aligned} Ax = b & \Leftrightarrow Q^{-1}(b - Ax) = 0 \\ & \Leftrightarrow f(x) := \underbrace{(I - Q^{-1}A)}_{=:G} \cdot x + \underbrace{Q^{-1}b}_{=:c} = x \end{aligned} ,$$

wobei $Q \in GL(n)$ eine beliebige reguläre Matrix ist.

Natürlich muss dafür Sorge getragen werden, dass das Fixpunktverfahren $x_{k+1} = f(x_k) = Gx_k + c$ konvergiert

Def. : konsistentes Iterationsverfahren :

Ein Iterationsverfahren heißt konsistent, falls $x = A^{-1}b$ die Gleichung $x = Gx + c$ erfüllt.
Das GSV, das ESV und das SOR-Verfahren sind konsistent

Satz 4.1 (Notwendiges und hinreichendes Konvergenzkriterium)

Ein konsistentes Fixpunktverfahren / Iterationsverfahren $x_{k+1} = f(x_k) = Gx_k + c$ mit $G \in Mat_n(R)$ konvergiert genau dann für jeden Startwert $x_0 \in R^n$, wenn

$$r(G) < 1 \quad ,$$

wobei $r(G) = \max_j |I_j(G)|$ der **Spektralradius** der Iterationsmatrix G ist.

Beweis: Buch S.263

Da $r(G) \leq \|G\|$ für jede zugeordnete Matrixnorm, ist $\|G\| < 1$ hinreichend für $r(G) < 1$.

In diesem Fall kann man die Fehler $x_k - x = G^k(x_0 - x)$ durch

$$\|x_k - x\| \leq \|G\|^k \cdot \|x_0 - x\|$$

abschätzen.

Neben der Konvergenz wurde verlangt, dass sich $f(x) = Gx + c$ leicht berechnen lässt. Dazu muss die Matrix Q einfach zu invertieren sein.

Die am besten zu invertierende Matrix ist zweifellos $Q = I$.

⇒

Richardson-Verfahren

$$x_{k+1} = x_k - Ax_k + b \quad (G = I - A)$$

Geht man von einer spd-Matrix A aus, so ergibt sich für den Spektralradius von G

$$r(G) = r(I - A) = \max\left(|1 - \lambda_{\max}(A)|, |1 - \lambda_{\min}(A)|\right).$$

Eine notwendige Bedingung für die Konvergenz der Richardson-Iteration ist somit

$$\lambda_{\max}(A) < 2.$$

Für sich genommen ist diese Iteration demnach nur selten verwendbar.

Die nächst komplizierteren Matrizen sind die Diagonalmatrizen, so dass als zweite Möglichkeit für Q die Diagonale D von

$$A = L + D + R$$

in Frage kommt, wobei $D = (a_{11}, \dots, a_{nn})$ und

$$L := \begin{bmatrix} 0 & \cdots & \cdots & 0 \\ a_{21} & \ddots & & \vdots \\ \vdots & & \ddots & \vdots \\ a_{n1} & \cdots & a_{n,n-1} & 0 \end{bmatrix}, \quad R := \begin{bmatrix} 0 & a_{12} & \cdots & a_{1n} \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & a_{n-1,n} \\ 0 & \cdots & \cdots & 0 \end{bmatrix}$$

⇒

Jacobi-Verfahren (Gesamtschrittverfahren GSV)

$$x_{k+1} = (I - D^{-1}A) \cdot x_k + D^{-1}b = -D^{-1}(L + R) \cdot x_k + D^{-1}b$$

Eine hinreichende Bedingung für seine Konvergenz ist die strikte Diagonaldominanz von A .

Algorithmus – GSV :

1.) Wähle Anfangsnäherung $x^{(0)} \in R^n$ und setze $k := 0$

2.) Iterationsschritt $x^{(k)} \mapsto x^{(k+1)}$

$$x_i^{(k+1)} := \frac{b_i - \sum_{j=1, j \neq i}^n a_{ij} x_j^{(k)}}{a_{ii}} \quad i = 1, \dots, n$$

$$k := k + 1$$

3.) if „Konvergenz“ then stop else goto 2.)

Abbruchkriterium : $\|x^{(k)} - A^{-1}b\| \leq TOL$ (Genauigkeitsschranke)

Problem: $A^{-1}b$ nicht bekannt

Idee : Schätze $\|x^{(k)} - A^{-1}b\|$ auf Grundlage der Konvergenzanalyse des Iterationsverfahrens

Satz 4.2 (Hinreichende Konvergenz des GSV)

Die Jacobi-Iteration $x_{k+1} = -D^{-1}(L+R) \cdot x_k + D^{-1}b$ konvergiert für jeden Startwert x_0 gegen die Lösung $x = A^{-1}b$, falls A strikt diagonaldominant ist, d.h.

$$|a_{ii}| > \sum_{i \neq j} |a_{ij}| \quad \text{für alle } i = 1, \dots, n$$

Beweis :

Die Aussage folgt aus Satz 4.1, da

$$\mathbf{r}(D^{-1}(L+R)) \leq \|D^{-1}(L+R)\|_{\infty} = \max_i \sum_{i \neq j} \left| \frac{a_{ij}}{a_{ii}} \right|.$$

Q.e.d.

Nach den diagonalen folgen die gestaffelten Gleichungssysteme / Matrizen als nächstkomplexere. Setzt man für Q die untere Dreieckshälfte $Q := D + L$ an, so erhält man das

Gauß-Seidel-Verfahren (Einzelschrittverfahren ESV)

$$\begin{aligned} x_{k+1} &= (I - (D+L)^{-1}A) \cdot x_k + (D+L)^{-1}b \\ &= -(D+L)^{-1}R \cdot x_k + (D+L)^{-1}b. \end{aligned}$$

Satz 4.3 (Satz von Ostrowski und Reich)

Für symmetrische positiv definite Matrizen A gilt :

$$\mathbf{r}(G(\mathbf{w})) < 1 \quad \forall \mathbf{w} \in]0,2[\quad , \text{ wobei } G(\mathbf{w}) \text{ die Iterationsmatrix des SOR-Verfahrens ist.}$$

Inbesondere konvergiert das ESV / Gauß-Seidel-Verfahren für jede spd-Matrix.

Algorithmus – ESV :

1.) Wähle Anfangsnäherung $x^{(0)} \in R^n$ und setze $k := 0$

2.) Iterationsschritt $x^{(k)} \mapsto x^{(k+1)}$

$$x_i^{(k+1)} := \frac{b_i - \sum_{j=1}^n a_{ij}x_j^{(k+1)} - \sum_{j=i+1}^n a_{ij}x_j^{(k)}}{a_{ii}} \quad i = 1, \dots, n$$

$$k := k + 1$$

3.) if „Konvergenz“ then stop else goto 2.)

Zusammenfassung

Rechenaufwand : Beim GSV wie auch beim ESV pro Iterationsschritt $\leq n \cdot p$ Rechenoperationen, falls maximal p Nichtnullelemente pro Zeile von A

Speicherbedarf : GSV : 2 Vektoren der Länge n ($x_{\text{alt}}, x_{\text{neu}}$)
ESV : 1 Vektor der Länge n ($x_i^{(k+1)}$ wird auf den Platz von $x_i^{(k)}$ abgespeichert)

Achtung !!! : Durchführbarkeit nur für $a_{ii} \neq 0$, $i = 1, \dots, n$

Konvergenz : Konvergenz für beliebige Startvektoren $x^{(0)} \in \mathbb{R}^n$, falls $\|G^{k+1}\| \rightarrow 0$ für eine verträgliche Matrixnorm

$$\left(x^{(k+1)} - x \right) = G(x^{(k)} - x) = G \cdot G(x^{(k-1)} - x) = \dots = G^{k+1}(x^{(0)} - x)$$

hinreichende Bedingung : $\|G\| = q < 1 \Rightarrow \|G^{k+1}\| \leq \|G\|^{k+1} = q^{k+1} \rightarrow 0$

notwendige Bedingung : $r(G) < 1$ für den Spektralradius.

Satz 4.4 (Zusammenhang GSV - ESV) – Satz von Stein und Rosenberg

Sei G_{GSV} die Iterationsmatrix des GSV und G_{ESV} die des ESV und alle Elemente von G_{GSV} größer Null.

Dann gilt :

$$0 < r(G_{GSV}) < 1 \quad \Rightarrow \quad 0 < r(G_{ESV}) < r(G_{GSV}) < 1$$

$$r(G_{GSV}) > 1 \quad \Rightarrow \quad 1 < r(G_{GSV}) < r(G_{ESV})$$

$$\text{oder} \quad r(G_{GSV}) = r(G_{ESV}) = 0 \quad \text{oder} \quad r(G_{GSV}) = r(G_{ESV}) = 1.$$

D.h. wenn das GSV konvergiert, so konvergiert auch das ESV , und wenn das GSV divergiert, so auch das ESV.
Desweiteren konvergiert das ESV schneller als das GSV.

Motivation : Konvergenzgeschwindigkeit verbessern

\Rightarrow

Relaxationsverfahren – SOR-Verfahren (successive over relaxation)

Das SOR-Verfahren unterscheidet sich nur wenig vom ESV. Zur Berechnung der jeweils folgenden Iterierten benutzt man die Idee des ESV, bildet aber anschließend noch einen gewichteten Mittelwert aus der alten und der neuen Komponente mit einem Parameter \mathbf{W} , den man ziemlich frei wählen kann. Für $\mathbf{W} = 1$ erhält man das alte ESV zurück.

$$\begin{aligned} x_{k+1} &= \mathbf{W} \cdot (Gx_k + c) + (1 - \mathbf{W}) \cdot x_k \\ &= G_{\mathbf{W}} x_k + \mathbf{W} \cdot c \quad \text{mit} \quad G_{\mathbf{W}} := \mathbf{W} \cdot G + (1 - \mathbf{W}) \cdot I , \end{aligned}$$

wobei $\mathbf{W} \in [0,1]$ ein sogenannter **Dämpfungsparameter** ist.

Auf diese Weise gewinnt man aus einer Fixpunktiteration $x_{k+1} = Gx_k + c$ eine ganze Schar von relaxierten Fixpunktiterationen mit der \mathbf{W} – abhängigen Iterationsfunktion

$$\mathbf{f}_{\mathbf{W}}(x) = \mathbf{W} \cdot \mathbf{f}(x) + (1 - \mathbf{W}) \cdot x = G_{\mathbf{W}} x + \mathbf{W} \cdot c .$$

Der Dämpfungsparameter \mathbf{W} ist nun so zu wählen, dass $r(G_{\mathbf{W}})$ möglichst klein wird.

Tatsächlich ist es für die sogenannten **symmetrisierbaren Iterationsverfahren** sogar möglich, durch geeignete Wahl von \mathbf{W} Konvergenz zu erzwingen, obwohl die Ausgangsiteration im Allgemeinen nicht konvergiert.

Def. symmetrisierbares Fixpunktverfahren :

Ein Fixpunktverfahren $x_{k+1} = Gx_k + c$ heißt symmetrisierbar, falls $I - G$ für jede spd-Matrix A ähnlich zu einer spd-Matrix ist, d.h. falls es eine reguläre Matrix $W \in GL(n)$ gibt, so dass

$$W \cdot (I - G) \cdot W^{-1}$$

eine spd-Matrix ist.

Relaxationsverfahren – Prof. Arnold

Idee : Berechne Näherung $\tilde{x}^{(k+1)}$ mit Iterationsverfahren und wähle

$$x^{(k+1)} := w \cdot \tilde{x}^{(k+1)} + (1 - w) \cdot x^{(k)}$$

mit einem Parameter $w \in R$ so, dass $\|x^{(k+1)} - A^{-1}b\|$ „möglichst klein“.

Wahl des Relaxationsparameters :

$w \in]0,1[$... under relaxation

$w > 1$... over relaxation

Hier Beschränkung auf Überrelaxation ($w > 1$):

JOR (Relaxation des GSV / Jacobi-Verfahrens)

$$x_i^{(k+1)} := (1 - w) \cdot x_i^{(k)} + w \frac{b_i - \sum_{j=1, i \neq j}^n a_{ij} x_j^{(k)}}{a_{ii}} = x_i^{(k)} + w \frac{b_i - \sum_{j=1}^n a_{ij} x_j^{(k)}}{a_{ii}} \quad i = 1, \dots, n$$

SOR (Relaxation des ESV / Gauß-Seidel-Verfahrens)

$$x_i^{(k+1)} := x_i^{(k)} + w \frac{b_i - \sum_{j=1}^n a_{ij} x_j^{(k+1)} - \sum_{j=i}^n a_{ij} x_j^{(k)}}{a_{ii}} \quad i = 1, \dots, n$$

Satz 4.5 (Satz von Kahan)

Für die Iterationsmatrix G_w des SOR-Verfahrens gilt :

$$r(G_w) \geq |w - 1|$$

Beweis : siehe Guido-Skript S.40

Def. zerlegbare Matrix (reduzibel) :

Eine Matrix $A \in R^{n \times n}$ heißt reduzibel / zerlegbar, wenn es eine Permutationsmatrix P gibt, so dass

$$P^T A P = \begin{pmatrix} C & D \\ 0 & F \end{pmatrix} ,$$

wobei die Blockmatrizen C und F quadratisch sind.

Falls kein Nullblock erzeugt werden kann, heißt die Matrix A unzerlegbar (irreduzibel) .

Ordnet man A einem gerichteten Graphen $G(A)$ so zu, dass A n Knoten K_1, \dots, K_n hat und das gilt :

$$G(A) \text{ enthält die Kante } K_i \rightarrow K_j \quad \Leftrightarrow \quad a_{ij} \neq 0 \quad ,$$

dann folgt :

A ist genau dann unzerlegbar, wenn $G(A)$ zusammenhängend ist, d.h. wenn es zu beliebigen $i, j \in \{1, \dots, n\}$ einen gerichteten Weg von K_i nach K_j gibt.

Beispiel zerlegbare Matrix :

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

$G(A)$ ist nicht zusammenhängend, da kein gerichteter Weg von K_1 nach K_3 existiert.

$$A \text{ ist zerlegbar, da } P^T A P = \begin{pmatrix} 1 & 0 & 3 \\ 0 & 1 & -1 \\ 0 & 2 & 1 \end{pmatrix} \quad \text{für } P = \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{pmatrix}$$

Def. konsistent geordnete Matrix :

Eine Matrix heißt konsistent geordnet, wenn die Eigenwerte von $-D^{-1}(\mathbf{a} \cdot L + \frac{1}{\mathbf{a}} R)$ für beliebiges $\mathbf{a} \neq 0$ von \mathbf{a} unabhängig sind.

Anriss Tschebyscheff-Beschleunigung :

Bei den bisherigen Fixpunktverfahren wird zur Berechnung von x_{k+1} nur die Information des letzten Iterationsschrittes herangezogen und die bereits berechneten Werte x_0, \dots, x_{k-1} nicht berücksichtigt.

Hier nun Verbesserung eines gegebenen Fixpunktverfahrens durch Konstruktion einer Linearkombination

$$y_k = \sum_{j=0}^k v_{kj} x_j$$

aus sämtlichen Werten x_0, \dots, x_k .

Bei geeigneter Wahl der Koeffizienten v_{kj} soll die Ersatzfolge $\{y_0, y_1, \dots\}$ schneller konvergieren als die Ausgangsfolge $\{x_0, x_1, \dots\}$.

Bestimmung der Koeffizienten :

Falls $x_0 = \dots = x_k = x$ bereits die Lösung ist, so muss auch $y_k = x$ gelten, woraus folgt, dass

$$\sum_{j=0}^k v_{kj} = 1 \quad .$$

....weiter behandelt im Buch S.270-274

4.2 Verfahren der konjugierten Gradienten (cg-Verfahren)

Motivation : Für LGS $Ax = b$ mit symmetrisch positiv definiten Koeffizientenmatrix A ein schnelles Verfahren finden.

Jede spd-Matrix A definiert in natürlicher Weise ein Skalarprodukt

$$(x, y) := \langle x, Ay \rangle \quad \text{mit der zugeordneten Norm} \quad \|y\|_A = \sqrt{(y, y)} \quad , \text{ der sogenannten Energienorm.}$$

Nutzung eines problemangepassten Orthogonalitätsbegriffs :

Def. konjugierte Vektoren :

Sei A eine spd-Matrix. Dann heißen zwei Vektoren x und y konjugiert, wenn gilt :

$$x^T \cdot A \cdot y = 0$$

Ist A die Einheitsmatrix, so bleibt das normale Skalarprodukt.

Die zusätzlich zwischengeschobene Matrix A produziert ein etwas allgemeineres Skalarprodukt.

Dadurch, dass die Matrix A symmetrisch positiv definit ist, bleiben die Eigenschaften des Skalarproduktes erhalten.

Verfahren:

Der Grundgedanke der Methode der konjugierten Gradienten besteht darin, statt das vorgegebene Gleichungssystem

$$Ax = b$$

zu lösen, das folgende Funktional zu minimieren :

$$F(x) = \frac{1}{2} x^T A x - b^T x \quad \left(= \frac{1}{2} \langle x, Ax \rangle - \langle x, b \rangle \right)$$

Der Gradient dieses Funktional lautet :

$$\text{grad } F(x) = Ax - b$$

Wurde für das Funktional ein Minimum bei x_0 errechnet, so verschwindet der Gradient an dieser Stelle :

$$\text{grad } F(x_0) = Ax_0 - b = 0 \quad ,$$

und damit ist x_0 die gesuchte Lösung des Gleichungssystems.

Beachte : Der Residuenvektor $r := Ax - b$ ist gerade gleich dem Gradienten des Funktional F .

Man startet mit einem beliebigen Vektor $x^{(0)}$.

Um möglichst schnell zum Minimum zu gelangen, legt man die Relaxationsrichtung $p^{(0)}$ in Richtung des negativen Gradienten, den man leicht als Residuenvektor berechnen kann. Denn dort ist der stärkste Abstieg zu erwarten.

$$p^{(0)} = -r^{(0)} = -\text{grad } F(x^{(0)}) = -(Ax^{(0)} - b)$$

Für den damit zu erzielenden ersten Näherungsvektor macht man den Ansatz

$$x^{(1)} = x^{(0)} + s_0 \cdot p^{(0)}$$

Die Bedingung, dass das Funktional F minimal werden möge, führt zu

$$s_0 = \frac{(r^{(0)})^T \cdot r^{(0)}}{(p^{(0)})^T \cdot A \cdot p^{(0)}}$$

Im allgemeinen Schritt wählt man als Relaxationsrichtung eine Linearkombination des aktuellen Residuenvektors, also des jeweiligen Gradienten, und der vorhergehenden Relaxationsrichtung :

$$p^{(k)} := -r^{(k)} + b_k \cdot p^{(k-1)}$$

Hier bestimmt sich b_k aus der Forderung, dass $p^{(k)}$ und $p^{(k-1)}$ konjugiert zueinander sein mögen :

$$b_k := \frac{(r^{(k)})^T \cdot r^{(k)}}{(r^{(k-1)})^T \cdot r^{(k-1)}}$$

Anschließend bestimmt man

$$s_k := \frac{(r^{(k)})^T \cdot r^{(k)}}{(p^{(k)})^T \cdot A \cdot p^{(k)}}$$

und berechnet den neuen Näherungsvektor für die Lösung

$$x^{(k+1)} := x^{(k)} + s_k \cdot p^{(k)}$$

Algorithmus cg-Verfahren :

$x^{(0)}$ beliebig
 $r^{(0)} := Ax^{(0)} - b$
 $p^{(0)} := -r^{(0)}$
 $k = 0$

if **accurate** then stop else
begin

$$\mathbf{s}_k := \frac{(r^{(k)})^T \cdot r^{(k)}}{(p^{(k)})^T \cdot A \cdot p^{(k)}}$$

$$\left(= \frac{\langle r^{(k)}, r^{(k)} \rangle}{\langle p^{(k)}, A \cdot p^{(k)} \rangle} \right)$$

$$x^{(k+1)} := x^{(k)} + \mathbf{s}_k \cdot p^{(k)}$$

$$k := k + 1$$

$$r^{(k)} := Ax^{(k)} - b$$

$$\mathbf{b}_k := \frac{(r^{(k)})^T \cdot r^{(k)}}{(r^{(k-1)})^T \cdot r^{(k-1)}}$$

$$\left(= \frac{\langle r^{(k)}, r^{(k)} \rangle}{\langle r^{(k-1)}, r^{(k-1)} \rangle} \right)$$

$$p^{(k)} := -r^{(k)} + \mathbf{b}_k \cdot p^{(k-1)}$$

end

(accurate = „ $r^{(k)}$ hinreichend klein“)

Algorithmus cg-Verfahren (Deuffhard S.277) :

$$p_1 := r_0 := b - Ax_0 ;$$

for $k := 1$ to k_{\max} do

$$\mathbf{s}_k := \frac{\langle r_{k-1}, r_{k-1} \rangle}{\langle p_k, p_k \rangle} = \frac{\langle r_{k-1}, Br_{k-1} \rangle}{\langle p_k, Ap_k \rangle} ;$$

$$x_k := x_{k-1} + \mathbf{s}_k \cdot p_k ;$$

if accurate then exit ;

$$r_k := r_{k-1} - \mathbf{s}_k \cdot Ap_k ;$$

$$\mathbf{b}_{k+1} := \frac{\langle r_k, r_k \rangle}{\langle r_{k-1}, r_{k-1} \rangle} ;$$

$$p_{k+1} := r_k + \mathbf{b}_{k+1} \cdot p_k ;$$

endfor

Bemerkungen :

- Pro Iterationsschritt wird tatsächlich im wesentlichen nur eine Matrix-Vektor-Multiplikation, nämlich $A \cdot p^{(k)}$ ausgeführt, womit das Verfahren den Forderungen bezüglich des Aufwands voll entspricht. (Zusätzlich werden noch Skalarprodukte und Additionen skalarer Vielfacher von Vektoren berechnet.)
 - Für das Abbruchkriterium „accurate“ wäre $\|x - x^{(k)}\|$ hinreichend klein, was jedoch so nicht ausführbar ist. Deswegen lautet es in der Praxis „ $\|r^{(k)}\|_2 = \|x - x^{(k)}\|_{A_2}$ hinreichend klein“.
- Wie sich in Kapitel 2.3.2 - „Beurteilung von Näherungslösungen“ zeigte, ist die Residuennorm jedoch kein geeignetes Maß für die Konvergenz : Gerade bei schlecht konditionierten Systemen, d.h. $\mathbf{k}(A) \gg 1$, können sich die Iterierten drastisch verbessern, obwohl die Normen der Residuen wachsen. (siehe weiter unten : „Vorkonditionierung“)
- Die Relaxationsrichtungen bestehen aus paarweise konjugierten Richtungen, demzufolge sind die Residuenvektoren paarweise orthogonal. In einem n-dimensionalen Raum gibt es aber nur n orthogonale Vektoren. Die Methode endet also theoretisch nach n Schritten. In der Praxis läuft aber durch das Auftreten von Rundungsfehlern doch ein iterativer Prozess ab.
 - Der aufwendigste Punkt ist die zweimalige Multiplikation von A mit $p^{(k)}$ bzw. $x^{(k)}$. Die zweite Multiplikation kann jedoch auf die erste zurückgeführt werden, denn es gilt :

$$\begin{aligned} r^{(k+1)} &= A \cdot x^{(k+1)} - b \\ &= A \cdot (x^{(k)} - s_k \cdot p^{(k)}) - b \\ &= r^{(k)} - s_k \cdot A \cdot p^{(k)} \end{aligned}$$

Betrachtungen zur Konvergenzgeschwindigkeit :

Satz 4.6

Der Approximationsfehler $x - x^{(k)}$ des cg-Verfahrens lässt sich in der Energienorm $\|y\|_A = \sqrt{\langle y, Ay \rangle}$ abschätzen durch

$$\|x - x^{(k)}\|_A \leq 2 \left(\frac{\sqrt{\mathbf{k}_2(A)} - 1}{\sqrt{\mathbf{k}_2(A)} + 1} \right)^k \|x - x^{(0)}\|_A ,$$

wobei $\mathbf{k}_2(A)$ die Kondition von A bezüglich der euklidischen Norm ist.

Beweis : Buch S.278/279

Kor. 4.7

Um den Fehler in der Energienorm um einen Faktor \mathbf{e} zu reduzieren, d.h.

$$\|x - x^{(k)}\|_A \leq \mathbf{e} \cdot \|x - x^{(0)}\|_A ,$$

benötigt man höchstens k cg-Iterationen, wobei k die kleinste ganze Zahl ist mit

$$k \geq \frac{1}{2} \sqrt{\mathbf{k}_2(A)} \cdot \ln \left(\frac{2}{\mathbf{e}} \right)$$

Beweis : Buch S.279/280

4.3 Vorkonditionierung

Die Abschätzungen der Konvergenzgeschwindigkeit sowohl für die Tschebyscheff-Beschleunigung als auch für das cg-Verfahren hängen monoton von der Kondition $k_2(A)$ bezüglich der euklidischen Norm ab.

Ziel : Transformation des Problems $Ax = b$ derart, so dass die entstehende Matrix von möglichst kleiner Kondition ist.

Idee : Anstelle von $Ax = b$ mit einer spd-Matrix $A \in Mat_n(R)$ löst man das für jede invertierbare Matrix $B \in GL(n)$ äquivalente Problem

$$\bar{A}\bar{x} = b \quad \text{mit} \quad \bar{A} := AB \quad \text{und} \quad \bar{x} := B^{-1}x .$$

Dabei ist darauf zu achten, dass die Symmetrie des Problems nicht zerstört wird, damit die Iterationsverfahren anwendbar bleiben.

Ist B ebenfalls symmetrisch positiv definit, so ist die Matrix $\bar{A} = AB$ zwar nicht mehr bezüglich des euklidischen Skalarproduktes $\langle \cdot, \cdot \rangle$ **selbstadjungiert**, wohl aber bezüglich des von B induzierten Produktes

$$(\cdot, \cdot)_B := \langle \cdot, B \cdot \rangle ,$$

da $(x, AB y)_B = \langle x, BAB y \rangle = \langle AB x, B y \rangle = (AB x, y)_B .$

Daher ist das cg-Verfahren wieder anwendbar, wenn man das Skalarprodukt geeignet überträgt :

$(\cdot, \cdot)_B$ übernimmt die Rolle des euklidischen Skalarproduktes $\langle \cdot, \cdot \rangle$ und das zugehörige „Energieprodukt“

$$(\cdot, \cdot)_{AB} = (AB \cdot, \cdot)_B = \langle AB \cdot, B \cdot \rangle$$

von $\bar{A} = AB$ die Rolle von (\cdot, \cdot) .

Daraus ergibt sich unmittelbar die folgende Iteration $\bar{x}_0, \bar{x}_1, \dots$ zur Lösung von $\bar{A}\bar{x} = b$:

$$p_1 := r_0 := b - AB\bar{x}_0 ;$$

for $k := 1$ to k_{\max} do

$$s_k := \frac{(r_{k-1}, r_{k-1})_B}{(p_k, p_k)_{AB}} = \frac{\langle r_{k-1}, Br_{k-1} \rangle}{\langle ABp_k, Bp_k \rangle} ;$$

$$\bar{x}_k := \bar{x}_{k-1} + s_k \cdot p_k ;$$

if accurate then exit ;

$$r_k := r_{k-1} - s_k \cdot ABp_k ;$$

$$b_{k+1} := \frac{(r_k, r_k)_B}{(r_{k-1}, r_{k-1})_B} = \frac{\langle r_k, Br_k \rangle}{\langle r_{k-1}, Br_{k-1} \rangle} ;$$

$$p_{k+1} := r_k + b_{k+1} \cdot p_k ;$$

endfor

Da man an einer Iteration für die eigentliche Lösung $x = B\bar{x}$ interessiert ist, ersetzt man daher die Zeile für die \bar{x}_k durch

$$x_k = x_{k-1} + s_k \cdot Bp_k .$$

Es fällt auf, dass nun die p_k nur in der letzten Zeile explizit auftauchen. Führt man aus diesem Grund die (A - orthogonalen)

Vektoren $q_k := Bp_k$ ein, so ergibt sich folgende sparsame Version des vorkonditionierten cg-Verfahrens :

Algorithmus pcg-Verfahren (Deuffhard S.283) :

(preconditioned conjugate gradient method)

$$r_0 := b - Ax_0 ;$$

$$q_1 := Br_0 ;$$

for $k := 1$ to k_{\max} do

$$s_k := \frac{(r_{k-1}, r_{k-1})_B}{(p_k, p_k)_{AB}} = \frac{\langle r_{k-1}, Br_{k-1} \rangle}{\langle Aq_k, q_k \rangle} ;$$

$$x_k = x_{k-1} + s_k \cdot q_k$$

if accurate then exit ;

$$r_k := r_{k-1} - s_k \cdot Aq_k ;$$

$$b_{k+1} := \frac{(r_k, r_k)_B}{(r_{k-1}, r_{k-1})_B} = \frac{\langle r_k, Br_k \rangle}{\langle r_{k-1}, Br_{k-1} \rangle} ;$$

$$q_{k+1} := Br_k + b_{k+1} \cdot q_k ;$$

endfor

Pro Iterationsschritt benötigt man jeweils eine Multiplikation mit der Matrix A (für Aq_k) bzw. mit B (für Br_k), also gegenüber dem ursprünglichen cg-Verfahren nur die Multiplikation mit B mehr.

(Konvergenzanalyse siehe Buch S.284 - S.286)

Beispiel

Eine sehr einfache, aber häufig schon wirkungsvolle Vorkonditionierung ist die Inverse $B := D^{-1}$ der Diagonale D von A , die sogenannte diagonale Vorkonditionierung.

Beispiel

Wendet man die Cholesky-Zerlegung $A = LL^T$ aus Kapitel 1 auf eine symmetrische dünnbesetzte Matrix an, so ist zu beobachten, dass außerhalb der Besetzungsstruktur von A meist nur „verhältnismäßig kleine“ Elemente l_{ij} entstehen. Dies führt zur Idee der unvollständigen Cholesky-Zerlegung (incomplete Cholesky decomposition, IC). Sie besteht darin, diese Elemente einfach wegzulassen.

Bezeichnet man mit

$$P(A) := \{(i, j) : a_{ij} \neq 0\} \quad \text{die Indexmenge der nicht verschwindenden Elemente einer Matrix } A ,$$

so konstruiert man also anstelle von L eine Matrix \tilde{L} mit

$$P(\tilde{L}) \subset P(A) ,$$

indem man wie bei der Cholesky-Zerlegung vorgeht und $\tilde{l}_{ij} := 0$ setzt für alle $(i, j) \notin P(A)$.

Dabei erwartet man, dass $A \approx \tilde{A} := \tilde{L}\tilde{L}^T$.

5. Lineare Eigenwertprobleme

Eigenwertproblem : Gegeben : quadratische Matrix A

Gesucht : alle Eigenwerte $\lambda_i(A)$ mit zugehörigen Eigenvektoren x_i ($Ax_i = \lambda_i x_i$)
oder aber die m betragsmäßig größten / kleinsten Eigenwerte

Häufigstes Eigenwertproblem : Ist A symmetrisch ($A = A^T$), so gilt :

- $\lambda_i \in \mathbb{R}^n$
- Eigenvektoren zu verschiedenen Eigenwerten sind orthogonal
- Es gibt eine Orthonormalbasis des \mathbb{R}^n aus Eigenvektoren von A .

Vorgriff :

Das Eigenwertproblem ist nur für normale Matrizen gut konditioniert, deren wichtigste Klasse die reel-symmetrischen Matrizen sind. Für allgemeine Matrizen ist dagegen das Problem der Singulärwertzerlegung gut konditioniert und praktisch enorm relevant.

Satz 5.1 (Satz von Gerschgorin)

Die Eigenwerte einer Matrix $A \in \mathbb{R}^{n \times n}$ liegen in der Vereinigung $\bigcup_{i=1}^n K_i$ der Gerschgorin-Kreise

$$K_i := \left\{ z \in \mathbb{C} : |z - a_{ii}| \leq \sum_{j=1, j \neq i}^n |a_{ij}| \right\}$$

Beispiel : Gerschgorin-Kreise

$$A = \begin{pmatrix} 1 & 0.1 & -0.1 \\ 0 & 2 & 0.4 \\ -0.2 & 0 & 3 \end{pmatrix} \Rightarrow \begin{aligned} K_1 &= \{z : |z - 1| \leq 0.2\} \\ K_2 &= \{z : |z - 2| \leq 0.4\} \\ K_3 &= \{z : |z - 3| \leq 0.2\} \end{aligned}$$

5.1 Kondition des allgemeinen Eigenwertproblems

Lemma 5.2 (Stetigkeit des Eigenwertproblems)

Sei $\lambda_0 \in \mathbb{C}$ einfacher Eigenwert von $A \in \text{Mat}_n(\mathbb{C})$. Dann existiert eine stetig differenzierbare Abbildung

$$\mathbf{I} : V \subset \text{Mat}_n(\mathbb{C}) \rightarrow \mathbb{C}, \quad B \mapsto \mathbf{I}(B)$$

von einer Umgebung V von A in $\text{Mat}_n(\mathbb{C})$, so dass $\mathbf{I}(A) = \lambda_0$

und $\mathbf{I}(B)$ einfacher Eigenwert von B ist für alle $B \in V$.

Ist x_0 ein Eigenvektor von A zu λ_0 und y_0 ein (adjungierter) Eigenvektor von $A^* := \overline{A}^T$ zum Eigenwert $\overline{\lambda_0}$, d.h.

$$Ax_0 = \lambda_0 x_0 \quad \text{und} \quad A^* y_0 = \overline{\lambda_0} y_0,$$

so gilt für die Ableitung von \mathbf{I} an der Stelle A , dass

$$\mathbf{I}'(A) \cdot C = \frac{\langle Cx_0, y_0 \rangle}{\langle x_0, y_0 \rangle} \quad \text{für alle } C \in \text{Mat}_n(\mathbb{C}).$$

Beweis : Buch S.134/135

Satz 5.3 (Kondition des Eigenwertproblems)

Die absolute Kondition der Bestimmung eines einfachen Eigenwertes λ_0 einer Matrix $A \in Mat_n(C)$ bezüglich der 2-Norm ist

$$k_{abs} = \|I'(A)\| = \frac{\|x\| \cdot \|y\|}{|\langle x, y \rangle|} = \frac{1}{|\cos(\angle(x, y))|}$$

und die relative Kondition

$$k_{rel} = \frac{\|A\|}{|\lambda_0|} \|I'(A)\| = \frac{\|A\|}{|\lambda_0 \cdot \cos(\angle(x, y))|}$$

wobei x ein Eigenvektor von A zum Eigenwert λ_0 ist, d.h. $Ax = \lambda_0 x$, und y ein adjungierter Eigenvektor, d.h. $A^*y = \bar{\lambda}_0 y$.

Insbesondere ist das Eigenwertproblem für normale Matrizen gut konditioniert mit $k_{abs} = 1$.

5.2 Eigenwertberechnung

Motivation : Die Eigenwerte einer Matrix $A \in Mat_n(R)$ als Nullstellen des charakteristischen Polynoms $c_A(\lambda) = \det(A - \lambda \cdot I)$ zu berechnen ist allenfalls für $n = 2$ akzeptabel.
 ⇒ Entwicklung von Methoden, um die Eigenwerte und Eigenvektoren direkt zu bestimmen.

Einfachste Methode : direkte und inverse Vektoriteration

5.2.1 direkte Vektoriteration (power method)

Idee : Man iteriert die durch die Matrix $A \in Mat_n(R)$ gegebene Abbildung und definiert eine Folge $\{x_k\}_{k=0,1,\dots}$ für einen beliebigen Startvektor $x_0 \in R^n$ durch

$$x_{k+1} := Ax_k \quad \text{für } k = 0, 1, \dots$$

Ist ein einfacher Eigenwert λ von A betragsmäßig echt größer als alle anderen Eigenwerte von A , so wird vermutet, dass sich λ bei der Iteration gegenüber allen anderen Eigenwerten „durchsetzt“ und x_k gegen einen Eigenvektor von A zum Eigenwert λ konvergiert.

Der Einfachheit halber Beschränkung auf symmetrische Matrizen :

Satz 5.4

Sei λ_1 ein einfacher Eigenwert der symmetrischen Matrix $A \in Mat_n(R)$ und betragsmäßig echt größer als alle anderen Eigenwerte von A , d.h. $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$.

Sei ferner $x_0 \in R^n$ ein Vektor, der nicht senkrecht auf dem Eigenraum von λ_1 steht.

Dann konvergiert die Folge

$$y_k := \frac{x_k}{\|x_k\|} \quad \text{mit} \quad x_{k+1} := Ax_k$$

gegen einen normierten Eigenvektor von A zum Eigenwert λ_1 .

Beweis :

Sei h_1, \dots, h_n eine Orthonormalbasis von Eigenvektoren von A mit $Ah_i = \lambda_i h_i$.

Dann gilt $x_0 = \sum_{i=1}^n a_i h_i$ mit $a_1 = \langle x_0, h_1 \rangle \neq 0$.

Folglich ist

$$x_k = A^k x_0 = \sum_{i=1}^n a_i l_i^k h_i = a_1 l_1^k \underbrace{\left(h_1 + \sum_{i=2}^n \frac{a_i}{a_1} \left(\frac{l_i}{l_1} \right)^k h_i \right)}_{=: z_k}.$$

Da $|l_i| < |l_1|$ für alle $i = 2, \dots, n$, gilt $\lim_{k \rightarrow \infty} z_k = h_1$ und daher

$$y_k = \frac{x_k}{\|x_k\|} = \frac{z_k}{\|z_k\|} \rightarrow \pm h_1 \quad \text{für } k \rightarrow \infty. \quad \text{Q.e.d.}$$

Nachteile der direkten Vektoriteration :

- Man erhält nur den Eigenvektor zum betragsgrößten Eigenwert von A .
- Die Konvergenzgeschwindigkeit hängt vom Quotienten $\left| \frac{l_2}{l_1} \right|$ ab.

(D.h. liegen die Eigenwerte l_1 und l_2 betragsmäßig dicht zusammen, so konvergiert die direkte Vektoriteration nur sehr langsam.)

⇒

5.2.2 Inverse Iteration (inverse power method)

Angenommen, man hätte einen Schätzwert $\bar{l} \approx l_i$ eines beliebigen Eigenwertes l_i der Matrix A zur Verfügung, so dass

$$|\bar{l} - l_i| < |\bar{l} - l_j| \quad \text{für alle } j \neq i. \quad (*)$$

Dann ist $(\bar{l} - l_i)^{-1}$ der betragsgrößte Eigenwert der Matrix $(A - \bar{l} \cdot I)^{-1}$. Konsequenterweise konstruiert man deshalb die Vektoriteration für diese Matrix. Diese Idee liefert die Iterationsvorschrift

$$(A - \bar{l} \cdot I) \cdot x_{k+1} = x_k \quad \text{für } k = 0, 1, \dots$$

Man beachte, dass bei jedem Iterationsschritt ein lineares Gleichungssystem gelöst werden muss, jedoch nur für verschiedene rechte Seiten x_k . Die Matrix $A - \bar{l} \cdot I$ muss daher nur einmal zerlegt werden.

Nach Satz 5.4 konvergiert die Folge $y_k := \frac{x_k}{\|x_k\|}$ unter der Voraussetzung (*) für $k \rightarrow \infty$ gegen einen normierten Eigenvektor

von A zum Eigenwert l_i , falls nicht gerade der Startvektor x_0 senkrecht auf dem Eigenvektor h_i zum Eigenwert l_i steht.

Der **Konvergenzfaktor** ist dabei

$$\max_{j \neq i} \left| \frac{l_j - \bar{l}}{l_i - \bar{l}} \right| < 1.$$

Ist \bar{l} eine besonders gute Schätzung von l_i , so gilt

$$\left| \frac{l_j - \bar{l}}{l_i - \bar{l}} \right| \ll 1 \quad \text{für alle } j \neq i, \quad \text{das Verfahren konvergiert in diesem Fall sehr schnell.}$$

Durch geeignete Wahl von \bar{l} kann man also mit dieser Methode und einem nahezu beliebigen Startvektor x_0 einzelne Eigenwerte und Eigenvektoren herausgreifen.

Bemerkung :

Man beachte, dass die Matrix $A - \bar{I} \cdot I$ für „gut gewähltes“ $\bar{I} \approx I_i$ fast singularär ist.

Im vorliegenden Fall entstehen daraus jedoch keine numerischen Schwierigkeiten, da nur die Richtung des Eigenvektors gesucht ist, deren Berechnung gut konditioniert ist.

Beispiel Vektoriteration :

$$A = \begin{pmatrix} -1 & 3 \\ -2 & 4 \end{pmatrix}$$

Eigenwerte : $I_1 = 1$, $I_2 = 2$

Geht man von einer Approximation $\bar{I} = 1 - \mathbf{e}$ für I_1 mit $0 < \mathbf{e} \ll 1$ aus, so ist die Matrix

$$A - \bar{I} \cdot I = \begin{pmatrix} -2 + \mathbf{e} & 3 \\ -2 & 3 + \mathbf{e} \end{pmatrix} \quad \text{fast singularär und}$$

$$(A - \bar{I} \cdot I)^{-1} = \frac{1}{\mathbf{e}(\mathbf{e} + 1)} \begin{pmatrix} 3 + \mathbf{e} & -3 \\ 2 & -2 + \mathbf{e} \end{pmatrix}.$$

Da sich der Faktor $1/\mathbf{e}(\mathbf{e} + 1)$ bei der Normierung herauskürzt, ist die Berechnung der Richtung einer Lösung x von $(A - \bar{I} \cdot I)x = b$ gut konditioniert.

Dies lässt sich auch an der relativen komponentenweisen Kondition

$$k_{rel} = \frac{\| |(A - \bar{I} \cdot I)^{-1}| \cdot |b| \|_{\infty}}{\|x\|_{\infty}} \quad \text{bezüglich Störungen der rechten Seite ablesen.}$$

Für $b := (1,0)^T$ ergibt sich z.B.

$$x = (A - \bar{I} \cdot I)^{-1} b = |(A - \bar{I} \cdot I)^{-1}| \cdot |b| = \frac{1}{\mathbf{e}(\mathbf{e} + 1)} \begin{pmatrix} 3 + \mathbf{e} \\ 2 \end{pmatrix} \quad \text{und daher (mit Kapitel 2 - LGS),}$$

$$k_{rel} = k_c((A - \bar{I} \cdot I)^{-1}, b) = 1 .$$

Tatsächlich wird in Programmen bei einer (echt) singularären Matrix $A - \bar{I} \cdot I$ ein Pivotelement $\mathbf{e} = 0$ durch die relative Maschinengenauigkeit *eps* ersetzt und mit der nun fast singularären Matrix die inverse Vektoriteration durchgeführt.

5.2.3 QR-Algorithmus für symmetrische Eigenwertprobleme

Wie in Kapitel 5.1 gezeigt wurde, ist das Eigenwertproblem für symmetrische Matrizen gut konditioniert.

Motivation :

Effektive Methode entwickeln um sämtliche Eigenwerte einer reellen symmetrischen Matrix gleichzeitig zu berechnen.

Es ist bekannt, dass A nur reelle Eigenwerte $I_1, \dots, I_n \in R$ besitzt und eine Orthonormalbasis $h_1, \dots, h_n \in R^n$ aus

Eigenvektoren $Ah_i = I_i h_i$ existiert, d.h.

$$Q^T A Q = \Lambda = \text{diag}(I_1, \dots, I_n) \quad \text{mit } Q = [h_1, \dots, h_n] \in O(n) . \quad (*)$$

Die erste Idee wäre, Q direkt in endlich vielen Schritten zu bestimmen. Da die Eigenwerte die Wurzeln des charakteristischen Polynoms sind, hätte man damit auch ein endliches Verfahren zur Bestimmung der Nullstellen von Polynomen beliebigen Grades gefunden. Nach dem **Satz von Abel** jedoch existiert i.A. kein solches Verfahren. (basierend auf den Operationen +, -, *, /, $\sqrt{\quad}$)

Die zweite Idee, die von (*) nahegelegt wird, ist, A durch eine Ähnlichkeitstransformation (Konjugation), z.B. mit orthogonalen Matrizen, der Diagonalgestalt näherzubringen, da die Eigenwerte invariant unter Ähnlichkeitstransformationen sind.

Versucht man, eine symmetrische Matrix A durch Konjugation mit Householder-Matrizen auf Diagonalgestalt zu bringen, so erweist sich dies schnell als unmöglich :

$$\begin{bmatrix} * & \dots & \dots & * \\ \vdots & & & \vdots \\ \vdots & & & \vdots \\ * & \dots & \dots & * \end{bmatrix} \xrightarrow{Q_1} \begin{bmatrix} * & * & \dots & \dots & * \\ 0 & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ 0 & * & \dots & \dots & * \end{bmatrix} \xrightarrow{Q_1^T} \begin{bmatrix} * & 0 & \dots & \dots & 0 \\ \vdots & * & & & \vdots \\ \vdots & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ * & * & \dots & \dots & * \end{bmatrix}$$

Was die Multiplikation mit der Householder-Transformation von links geschafft hat, wird bei der Multiplikation von rechts wieder zerstört.

Anders sieht es aus, wenn man A nur auf **Tridiagonalgestalt** bringen will :

$$\begin{bmatrix} * & \dots & \dots & * \\ \vdots & & & \vdots \\ \vdots & & & \vdots \\ * & \dots & \dots & * \end{bmatrix} \xrightarrow{P_1} \begin{bmatrix} * & * & \dots & \dots & * \\ * & \vdots & & & \vdots \\ 0 & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ 0 & * & \dots & \dots & * \end{bmatrix} \xrightarrow{P_1^T} \begin{bmatrix} * & * & 0 & \dots & 0 \\ * & * & \dots & \dots & * \\ 0 & \vdots & & & \vdots \\ \vdots & \vdots & & & \vdots \\ 0 & * & \dots & \dots & * \end{bmatrix}$$

Lemma 5.5

Sei $A \in \text{Mat}_n(\mathbb{R})$ symmetrisch. Dann existiert eine orthogonale Matrix $P \in O(n)$, welche das Produkt von $n-2$ Householder-Reflexionen ist, so dass PAP^T Tridiagonalgestalt hat.

Beweis : Buch S.142

Somit wurde das Ausgangsproblem auf die Bestimmung der Eigenwerte einer symmetrischen Tridiagonalmatrix transformiert. Es existiert auch eine LR-Zerlegung, hier aber nur :

QR-Zerlegung : Diese existiert immer (keine Permutation nötig) und ist vor allen Dingen inhärent stabil.

Man definiert daher eine Folge $\{A_k\}_{k=1,2,\dots}$ von Matrizen durch

- a) $A_1 = A$
- b) $A_k = Q_k R_k$, QR-Zerlegung
- c) $A_{k+1} = R_k Q_k$.

Lemma 5.6

Die Matrizen A_k haben folgende Eigenschaften :

- 1.) Die Matrizen A_k sind alle konjugiert zu A .
- 2.) Ist A symmetrisch, so auch alle A_k .
- 3.) Ist A symmetrisch und tridiagonal, so auch alle A_k .

Beweis : Buch S.142/143

Satz 5.7 (Konvergenzeigenschaften für den einfachen Fall, dass die Beträge der Eigenwerte paarweise verschieden sind)

Sei $A \in \text{Mat}_n(\mathbb{R})$ symmetrisch mit den Eigenwerten $\lambda_1, \dots, \lambda_n$, so dass $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n| > 0$ und A_k, Q_k, R_k wie oben definiert. Dann gilt mit $A_k = (a_{ij}^{(k)})$:

- a) $\lim_{k \rightarrow \infty} Q_k = I$,
- b) $\lim_{k \rightarrow \infty} R_k = \Lambda$,
- c) $a_{ij}^{(k)} = O\left(\left|\frac{\lambda_i}{\lambda_j}\right|^k\right)$ für $i > j$.

Beweis : Buch S.144/145

Bemerkung :

Eine genauere Analyse zeigt, dass das Verfahren auch für mehrfache Eigenwerte $\lambda_i = \dots = \lambda_j$ konvergiert. Falls hingegen $\lambda_i = -\lambda_{i+1}$, so konvergiert das Verfahren nicht. Es bleiben 2×2 Blöcke stehen.

Liegen zwei Eigenwerte λ_i, λ_{i+1} betragsmäßig dicht beieinander, so konvergiert das Verfahren nur sehr langsam. Dies kann mit Hilfe der sogenannten **Shift-Strategien** verbessert werden. Im Prinzip versucht man, die beiden Eigenwerte dichter an den Nullpunkt zu schieben und so den Quotienten $|\lambda_{i+1} / \lambda_i|$ zu verkleinern.

Dazu verwendet man für jeden Iterationsschritt k einen **Shift-Parameter** s_k und definiert die Folge $\{A_k\}$ durch :

- a) $A_1 = A$
- b) $A_k - s_k I = Q_k R_k$ QR-Zerlegung
- c) $A_{k+1} = R_k Q_k + s_k I$

...detaillierter siehe Buch S. 145/146

Algorithmus QR-Algorithmus :

- a) Reduziere das Problem auf Tridiagonalgestalt, $A \rightarrow A_1 = PAP^T$, A_1 symmetrisch und tridiagonal, $P \in O(n)$.
- b) Approximiere die Eigenwerte mit dem QR-Algorithmus mit Givens-Rotationen angewandt auf A_1 , $\Omega \cdot A_1 \cdot \Omega^T \approx \Lambda$, Ω Produkt aller Givens-Rotationen $\Omega_{ij}^{(k)}$.
- c) Die Spalten von ΩP approximieren die Eigenvektoren von A : $\Omega P \approx [h_1, \dots, h_n]$

Aufwand :

- a) $\frac{4}{3}n^3$ Multiplikationen für die Transformation auf Tridiagonalgestalt,
- b) $O(n^2)$ Multiplikationen für den QR-Algorithmus.

Für große n überwiegt daher der Aufwand für die Reduktion auf Tridiagonalgestalt.