

# Numerische Lineare Algebra

## Eigenwertproblem: Shifts und Eigenvektorberechnung



$$Q^{-1}AQ = QTAQ = Q^2Q^{-1}AQ = A,$$

mit  $A$  wie in (1.3). Die Matrix  $A$  ist also über eine *orthogonale* Transformation *ähnlich zur oberen Hessenberg-Matrix*  $A$ . Also gilt:

Man kann eine Matrix  $A \in \mathbb{R}^n \times \mathbb{R}^n$  durch Householder-Transformation auf eine zu  $A$  ähnliche Matrix mit oberer Hessenberggestalt bringen.

Zur Berechnung der Eigenwerte der oberen Hessenberg-Matrix  $A$  benutzen wir den QR-Algorithmus. Man kann zeigen, dass die Hessenberggestalt der Matrix mehrere Vorteile bringt:

-wenn  $A$  eine nicht-reduzierbare obere Hessenbergmatrix ist, kann man die Identität als Anfangsmatrix beim QR-Algorithmus nehmen.

-dadurch, dass man beim QR-Algorithmus in einer Vorbereitungsphase die Matrix auf obere Hessenbergform bringt, braucht man nur die QR-Zerlegung einer Hessenberg-Matrix  $A_{k-1}$  zu berechnen. Falls man dazu Givens-Rotationen verwendet, ist der Aufwand für die Berechnung  $A_{k-1} = QR$ ,  $A_k = RQ$  nur  $O(n^2)$  Operationen. Falls  $A$  symmetrisch ist, ist dieser Aufwand nur  $O(n)$  Operationen.

## Der QR-Algorithmus

Basis für diesen Zusammenhang, der letztendlich im QR-Algorithmus mündet, ist der folgende Satz von Schur (1908), der hier in der Formulierung für reelle Matrizen wiedergegeben wird und schon in der Einführung zur Gesamtproblematik zitiert wurde.

**Satz (Satz von Schur).** Zu jeder reellen Matrix  $A$  der Ordnung  $N$  existiert eine orthogonale Matrix  $Q$  der Ordnung  $N$ , so dass die zu  $A$  ähnliche Matrix  $R = QTAQ$  die Quasidreiecksgestalt

$$R = QTAQ = \begin{pmatrix} R_{11} & \dots & R_{1m} \\ & \ddots & \\ & & R_{mm} \end{pmatrix} \begin{matrix} 0 \\ \vdots \\ 0 \end{matrix} \quad (1.4)$$

Besitzt.

(ohne Beweis).

Auf diesem Satz beruht nun *die Idee des QR-Algorithmus*: Man konstruiere obige Quasidreiecksmatrix  $R$  (deren Existenz aufgrund des Satzes gesichert ist) auf möglichst effektive Weise. Die Eigenwerte sind dann entweder direkt ablesbar oder aus den  $2 \times 2$  Untermatrizen einfach zu berechnen. Diese Konstruktion versucht man, über eine Folge von QR-Transformationen (bezüglich Hessenberg-Matrizen)  $H_k$ ,

$$H_k = Q_k R_k, \quad H_{k+1} = R_k Q_k, \quad k=1, 2, \dots$$

zu erreichen. Dann gilt folgende Konvergenzaussage

**Satz (Francis, 1961).**

- (i) Die Matrix  $H$  habe Eigenwerte  $\lambda_i$  mit der Eigenschaft, dass  

$$\lambda_1 > \lambda_2 > \dots > \lambda_N,$$

d.h. die Eigenwerte sind reell und nicht entartet. Die dazugehörigen Eigenvektoren  $x_i$  seien als Spalten in der regulären Matrix  $X \in \mathbb{R}^{N \times N}$  (der Reihe nach) zusammengefasst. Falls nun für  $X^{-1}$  die LU-Zerlegung existiert, dann konvergiert die Folge der QR-Transformationen  $H_k$  für  $k \rightarrow \infty$  gegen eine Rechtsdreiecksmatrix und es gilt

$$\lim_{k \rightarrow \infty} h_{ii} = \lambda_i, \quad i=1, \dots, N, \quad (1.5)$$

d.h. in der erzeugten Matrix sind die Eigenwerte betragsmäßig geordnet, mit dem absolut größten Eigenwert „links oben“

- (ii) Hat die Matrix  $H$  Paare von konjugiert komplexen Eigenwerten derart, dass ihre Beträge und die Beträge der reellen Eigenwerte paarweise verschieden sind, und existiert die komplexe LU-Zerlegung von  $X^{-1}$ , dann konvergiert die Folge  $H_k$ ,  $k \rightarrow \infty$  gegen die Quasidreiecksmatrix  $A$

**Problem.** Obwohl aufgrund dieses Satzes die Konvergenz (und damit die Lösung des Eigenwertproblems) gesichert ist (zumindest unter der gegebenen Einschränkung, die allerdings fast immer erfüllt ist), kann diese Konvergenz sehr langsam sein. Es lässt sich nämlich zeigen, dass für große  $k$  die Subdiagonalelemente wie

$$h_{i+1,ik} \approx \lambda_{i+1} \lambda_{ik}, \quad i=1, \dots, N-1 \quad (1.6)$$

konvergieren. Damit gilt, dass

- (a) im reellen Fall,  $\lambda_{i+1}/\lambda_i < 1$  die Konvergenz nur *linear* ist und vom Verhältnis ( $\lambda_{i+1}/\lambda_i$ ) abhängt, wobei dieses Verhältnis nahe bei „1“ liegen kann und somit die Konvergenz „schleichend“ ist.

- (b) Für konjugiert komplexe Paare tritt keine Konvergenz ein, da  $\lambda_{i+1}/\lambda_i = 1$  (wie es auch der Satz behauptet; man beachte das Auftreten von entsprechenden  $2 \times 2$  Untermatrizen).

Aufgrund dieser Problematik ist der „reine“ QR-Algorithmus nur schlecht geeignet, das Eigenwertproblem zu lösen, da die erforderliche Rechenzeit teilweise unannehmbar wäre. Es gibt allerdings eine Möglichkeit, die Konvergenz drastisch zu beschleunigen, die auf der sog. „Spektralverschiebung“ beruht und auf zwei Arten in den Algorithmus „eingebaut“ werden kann. Im weiteren werden wir uns auf die sog. „explizite“ Spektralverschiebung konzentrieren und die entsprechenden Algorithmen (sowohl für rein reelle als auch für reelle und komplexe Eigenwerte) vorstellen. Bezüglich der zweiten Methode, der sog. „impliziten“ Spektralverschiebung, sei auf die Literatur verwiesen z.B. Schwarz, „Numerische Mathematik, Kap. 6.4.3 und „Numerical Recipes“, Kap. 11.3/4.

### Der QR-Algorithmus mit expliziter Spektralverschiebung für reelle Eigenwerte

Die oben erwähnte *Spektralverschiebung* beruht darauf, die Eigenwerte einer Matrix durch Subtraktion einer Diagonalmatrix mit identischen Elementen  $\sigma \in \mathbb{R}$  zu modifizieren. Die modifizierte Matrix

$$H = H - \sigma I$$

hat dann die Eigenwerte  $\lambda_i - \sigma$ ,  $i=1, \dots, N$ . Diese seien jetzt so indiziert, dass  

$$\lambda_1 - \sigma > \lambda_2 - \sigma > \dots > \lambda_N - \sigma$$

Damit konvergieren die Subdiagonalelemente der Folge

$$H_k \text{ mit } H_1 = H - \sigma I$$

mit

$$h_{i+1,i} \approx \lambda_{i+1} - \sigma \lambda_i - \sigma_k, \quad i=1, \dots, N-1 \quad (1.7)$$

Sei nun  $\sigma$  eine gute Näherung für  $\lambda_N$ , so dass

$$\lambda_N - \sigma \ll \lambda_i - \sigma, \quad i=1, \dots, N-1$$

gilt. Dann konvergiert das „letzte“ Subdiagonalelement  $h_{N,N-1}$  sehr schnell gegen Null (man beachte, dass der Quotient in Gl. (1.7) mit dieser Wahl sehr klein ist), und die Matrix zerfällt nach wenigen Iterationsschritten in die folgende Struktur (nachdem wir die Spektralverschiebung rückgängig gemacht haben, s.u.)

$$\left( \begin{array}{cccc|c} \times & \times & \times & \times & \times \\ \times & \times & \times & \times & \times \\ \mathbf{0} & \times & \times & \times & \times \\ \mathbf{0} & \mathbf{0} & \times & \times & \times \\ \dots & \dots & \dots & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \underbrace{\tilde{h}_{N,N-1}}_{\rightarrow 0} & \rightarrow \lambda_N \end{array} \right) \Rightarrow \left( \begin{array}{c|c} \hat{H} & \hat{h} \\ \dots & \dots \\ \mathbf{0}^\top & \lambda_N \end{array} \right)$$

Nachdem die Matrix zerfallen ist und wir somit den ersten Eigenwert  $\lambda_N$  bestimmt haben, sind die restlichen Eigenwerte diejenigen der Matrix  $H \in \mathbb{R}^{(N-1) \times (N-1)}$ , die wir genauso wie eben beschrieben behandeln.

Die *Spektralverschiebung* wird nicht nur einmal, sondern vor jedem QR-Schritt angewandt, wobei  $\sigma$  den Veränderungen angepasst wird. Am Ende des Schrittes wird die jeweilige Verschiebung wieder rückgängig gemacht. Damit lautet der QR-Algorithmus mit expliziter Spektralverschiebung

$$\begin{aligned} H_k - \sigma_k I &= Q_k R_k' \quad (\text{Zerlegung}) \\ &\Rightarrow R_k' = Q_k^\top (H_k - \sigma_k I) \\ &\Rightarrow \\ H_{k+1} &= R_k' Q_k + \sigma_k I \quad (\text{Transformation, Spektralverschiebung rückgängig gemacht}) \\ &= Q_k^\top (H_k - \sigma_k I) Q_k + \sigma_k I \\ &= Q_k^\top H_k Q_k \end{aligned} \quad (1.8)$$

Unter Verwendung der Spektralverschiebung haben wir also wiederum eine Matrix  $H_{k+1}$  erzeugt, die orthogonal ähnlich zu  $H_k$  ist (wie im originalen Algorithmus, allerdings jetzt mit einer veränderten Transformationsmatrix  $Q'$ ), aber die Konvergenz beschleunigt. Man beachte allerdings, dass die Eigenwerte jetzt nicht mehr bezüglich ihrer (absoluten) Größe  $\lambda_i$  angeordnet sind. Nachzutragen bleibt nun noch, welchen Wert man für das jeweilige  $\sigma_k$  annimmt. Optimal wäre natürlich  $\sigma = \lambda_N$ , aber  $\lambda_N$  wird ja gerade gesucht. In Praxis werden **zwei Möglichkeiten** zur Wahl von  $\sigma$  betrachtet.

(a) Man verwendet einfach das aktuelle „letzte“ Diagonalelement

$$\sigma_k = h_{NN}(k), \quad k=1, 2, \dots, \quad (1.9)$$

(b) Eine bessere Konvergenz (und eine Verallgemeinerung auf den komplexen Fall, siehe nächstes Kapitel) ergibt sich, wenn man  $\sigma$  über die Eigenwerte der aktuellen, letzten  $2 \times 2$ -Untermatrix nähert. Diese berechnen sich über

$$C_k = \begin{vmatrix} h_{N-1,N-1} & h_{N-1,N} \\ h_{N,N-1} & h_{NN} \end{vmatrix} = abcd$$

$$a - \lambda d - \lambda - bc = 0$$



Die ersten beiden Zeilen wurden dabei durch die Transformationen  $U_{1T}=U_{T1,2}$  und  $U_{2T}=U_{T2,3}$  modifiziert. Die folgenden Rotationen  $U_{3T}, U_{4T}, \dots$  betreffen nur noch die dritte und folgende Zeilen, d.h. Zeile 1 und 2 und *insbesondere Spalte 1 und 2 dieser Zeilen* werden durch die weitergehende Zerlegung nicht mehr verändert. Deshalb kann nun die Operation  $\cdot U_1$  (Rechtsmultiplikation) durchgeführt werden, da sie „nur“ Linearkombinationen der Spalten 1 und 2 bildet. Den nächsten Schritt bildet dann die Linksmultiplikation von  $U_{3T}$ , danach die Rechtsmultiplikation  $\cdot U_2$  usw. Mit dieser Vorgehensweise müssen die jeweiligen Drehwinkel nicht abgespeichert werden, und das gesamte gestaffelte Verfahren lässt sich wie folgt beschreiben.

- Schritt 1 :  $U_{1T}$  von links, merke Drehwinkel (alt)
  - Schritt  $i: i \leq 2 \leq N-1$
- a)  $U_{iT}$  von links, neuer Drehwinkel: QR-Zerlegung in  $i, (i+1)$  ter Zeile
  - b)  $U_{i-1}$  von rechts mit altem Drehwinkel: Aufbau der H-Matrix in  $i, (i-1)$  ter Spalte
  - c) alter Drehwinkel = neuer Drehwinkel
  - d) fortlaufend Subtraktion von  $\sigma_k$  und finale Addition für Diagonale
- Schritt  $N: N-1$  von rechts, Addition von  $\sigma_k$  zu Diagonalelement  $N-1, N-1$  und  $N; N$

### Komplexe Eigenwerte: Der QR-Doppelschritt

Wie wir im vorangegangenen Abschnitt (und im obigen Beispiel) gesehen haben, können bei der „Diagonalisierung“ von reellen Matrizen natürlich auch konjugiert komplexe Eigenwerte auftreten, die wir mit unserer bisherigen Vorgehensweise nicht erfassen konnten: Einerseits ergibt die Wahl von  $\sigma_k = \sqrt{\lambda_k} \in \mathbb{R}$  entsprechend Gl. (1.9) keine vernünftige Näherung für einen komplexen Eigenwert, andererseits ergibt ein komplexer Wert von  $\sigma_k$  entsprechend Gl. (2.0) eine Matrix mit komplexer Diagonalen, die wir bislang nicht bearbeiten könnten. Der sog. „QR-Doppelschritt“ löst diese Problematik. Als Vorbereitung auf diesen betrachten wir noch einmal die quadratische  $2 \times 2$  Untermatrix in der „linken unteren“ Ecke,

$$C_k = \begin{pmatrix} h_{N-1, N-1} & h_{N-1, N} \\ h_{N, N-1} & h_{N, N} \end{pmatrix} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

jetzt aber nicht nur mit den dazugehörigen reellen, sondern ggf. auch komplexen Eigenwerten und definieren

$$\lambda_{1,2} = \alpha \pm i\beta \quad \lambda_{1,2} = \alpha - i\beta \quad (2.4)$$

wobei

$$\alpha = \frac{a+d}{2} \quad \text{und} \quad \beta = \frac{\sqrt{(a-d)^2 + 4bc}}{2}$$

die Diskriminante des Problems ist.

Falls die Eigenwerte nun komplex sind, definieren wir den QR-Doppelschritt als zweifache UR-Transformation, einmal mit  $\sigma_k = \lambda_1$  und danach mit  $\sigma_{k+1} = \lambda_2$ , wobei die UR-Transformation das komplexe Analogon der QR-Transformation ist, und die zugehörige UR-Zerlegung durch folgenden Satz spezifiziert ist.

**Satz.** Für jede komplexe Matrix  $A \in \mathbb{C}^{N \times N}$  existiert die unitäre UR-Zerlegung

$$A = U \cdot R \quad (2.5)$$

wobei  $U$  eine unitäre Matrix ( $U^* = U^{-1}$ ) und  $R$  eine komplexe Rechtsdreiecksmatrix mit reeller Diagonalen ist (ohne Beweis).

**Der QR-Doppelschritt als zweifache UR-Transformation lässt sich nun folgendermaßen darstellen:**

a) erste Zerlegung

$$H_k - \sigma_{k+1} = U_k R_k \quad (2.6)$$

wobei  $(H_k - \sigma_{k+1})$  eine komplexe Diagonale hat (s.o.).

b) und zugehörige Transformation

$$H_{k+1} = R_k U_k + \sigma_{k+1} = U_k^\dagger H_k U_k \quad (2.7)$$

c) zweite Zerlegung

$$H_{k+1} - \sigma_{k+1} = U_{k+1} R_{k+1} \quad (2.8)$$

d) und zugehörige Transformation

$$H_{k+2} = R_{k+1} U_{k+1} + \sigma_{k+1} = U_{k+1}^\dagger H_{k+1} U_{k+1} \quad (2.9)$$

Mit (2.7) finden wir also zunächst

$$H_{k+2} = U_k U_{k+1}^\dagger H_k U_k U_{k+1} = Q^\dagger H_k Q \quad (3.4)$$

und  $H_{k+2}$  ist orthogonal ähnlich zu  $H_k$ , und wiederum eine Hessenberg-Matrix (UR-Transformationen erhalten die Hessenberg-Form).

Damit ergibt sich folgender Algorithmus für die Lösung des Eigenwertproblems einer reellen Matrix mit komplexen Eigenwerten

- i) Man berechne  $X = H_{k+2} - s H_{k+1}$  mit  $s$  und  $t$  aus  $\sigma_k, \sigma_{k+1} = \lambda_{1,2}, \dots$ ,
- ii) zerlege dann die reelle Matrix  $X = Q \cdot R$  und
- iii) führe die Transformation mit diesem  $Q$  durch:  $H_{k+2} = Q^\dagger H_k Q$ .
- iv) Die explizite Subtraktion und (finale) Addition der  $\sigma$ , die im Falle von reellen Eigenwerten notwendig war, entfällt hier, da die Spektralverschiebung schon bei der Berechnung von  $X$  berücksichtigt wird und bei der Transformation (d.h. der Rechtsmultiplikation mit  $Q = U_k U_{k+1}$ ) wieder rückgängig gemacht wird.

## Bestimmung der Eigenvektoren nach Durchführung des QR-Algorithmus

Im Folgenden wollen wir kurz diskutieren, wie man nach einer Bestimmung der Eigenwerte mittels QR-Algorithmus die dazugehörigen Eigenvektoren ermittelt.

### Direkte Berechnung

Die direkte Berechnung ist sehr aufwendig, da sie eine konsequente Buchhalten über alle orthogonalen Transformationen erfordert, die letztendlich zu der die Eigenwerte darstellenden Quasidreiecksmatrix  $R$  (1.4) geführt haben. Dies sind die Transformationen

$$Q = Q_1 Q_2 \dots Q_M$$

wenn die Ausgangsmatrix  $A$  über Givens in eine Hessenberg-Matrix  $H$  und diese auf  $R$  QR-transformiert wurde, so dass final  $R = Q^\dagger A Q$

gilt. Mit der originalen Eigenwertaufgabe



$$A \cdot x = \lambda \cdot x$$

folgt dann

$$QTAQR \cdot QTxy = \lambda QTxy$$

d.h. die Eigenwertaufgabe reduziert sich zunächst auf die Lösung von

$$R \cdot y = \lambda y \text{ mit } y := QTx$$

was sich aufgrund der Quasidreieckstruktur von R bei bekanntem  $\lambda$  i einfach durchführen lässt. Zur Bestimmung der tatsächlich gesuchten Eigenvektoren ist dann allerdings die Kenntnis von Q unerlässlich,

$$x_i = Q \cdot y_i$$

### Inverse Iteration

Das zweite Verfahren ist allgemeinerer Natur und setzt nur voraus, dass man die Eigenwerte einer Matrix zuvor schon bestimmt hat (wie auch immer). Aufgrund der iterativen Natur dieses Verfahrens eignet es sich auch zu einer Verbesserung der Qualität von schon zuvor bestimmten Eigenvektoren, z.B. mittels des Jacobi-Verfahrens.

Wir nehmen dazu an, dass für die Originalmatrix A ein numerischer Näherungswert  $\lambda$ , für den entsprechenden „exakten“ Eigenwert  $\lambda$  i bekannt ist. Definiert man nun die Folge

$$A - \lambda I z_k = z_{k-1} \tag{3.5}$$

mit einem geeignet definierten Startwert  $z_0$ , zum Beispiel  $(1, 1, 1, \dots, 1)^T$  oder einem Schätzwert für den gesuchten Eigenvektor, dann konvergiert die Folge  $z_k$ ,  $k=1, 2, \dots$  (unter einigen Einschränkungen) gegen den zu  $\lambda$  i gehörigen Eigenvektor.

Der Beweis zu dieser Aussage wird im Folgenden skizziert. Unter der (fast immer gerechtfertigten) Annahme, dass A einen vollständigen Satz von Eigenvektoren besitzt ( $x_j$ ,  $j=1, \dots, N$ ), lassen sich die iterierten Werte  $z$  nach diesen Eigenvektoren entwickeln, d.h.

$$z_{k-1} = \beta_j x_j$$

$$z_k = \alpha_j x_j$$

Einsetzen in (3.5) liefert

$$A - \lambda I \alpha_j x_j = \beta_j x_j$$

$$\alpha_j (\lambda_j - \lambda) x_j = \beta_j x_j$$

$$\Rightarrow \alpha_j = \beta_j (\lambda_j - \lambda)$$

Demzufolge lässt sich also  $z_k$  durch die Koeffizienten von  $z_{k-1}$  ausdrücken,

$$z_k = \beta_j (\lambda_j - \lambda) x_j$$

Sei nun der Startvektor durch die Entwicklung

$$z_0 = c_j x_j$$

gegeben, dann folgt sofort, dass

$$z_k = c_j (\lambda_j - \lambda)^k x_j \quad (3.6)$$

ist. (Man beachte die Potenz  $k$  im Nenner!). Falls nun der Eigenwert nicht entartet und  $\lambda_i$  ein guter Schätzwert für  $\lambda$  ist,

$$c_i \neq 0 \text{ und } 0 < \lambda_i - \lambda \ll \min_{j \neq i} |\lambda_j - \lambda| \quad (3.7)$$

dann gilt mit

$$z_k = \frac{1}{\lambda_i - \lambda} c_i x_i + \sum_{j \neq i} c_j \frac{(\lambda_j - \lambda)^k}{(\lambda_i - \lambda)^k} x_j \rightarrow 0 \quad (3.8)$$

dass  $z_k$  rasch gegen den gesuchten Eigenvektor  $x_i$  konvergiert, falls nach jeder Iteration  $z_k$  normiert wird!

Das Problem dieser Vorgehensweise liegt natürlich darin, dass wir de facto das lineare Gleichungssystem

$$(A - \lambda I)y = b$$

mit  $y = z_k$  und  $b = z_{k-1}$  lösen müssen, zum Beispiel über LU-Zerlegung (pro Eigenwert  $\lambda_i$  ist dann nur eine einmalige Zerlegung notwendig). Die (zu zerlegende) Matrix  $A - \lambda_i I$  ist aber nahezu singular, da  $\lambda \approx \lambda_i$ , d.h. das Gleichungssystem ist äußerst schlecht konditioniert, was zu großen Fehlern in der Lösung führt, egal wie „gut“ der Lösungsalgorithmus ist.

Dass die inverse Iteration überhaupt vernünftige Resultate liefert, liegt nun daran, dass im allgemeinen dieser Fehler, der bei der Lösung von linearen Gleichungssystemen gemacht wird, hier:

$$z_k - z_k^{\text{exakt}}$$

(mit  $z_k$  der numerische und  $z_k^{\text{exakt}}$  der „exakten“ Lösung), eine dominante Komponente in Richtung des Eigenvektors des betragskleinsten Eigenwertes der das Gleichungssystem darstellenden Matrix (hier:  $A - \lambda_i I$ ) hat.

Auf Grund der Voraussetzung (3.7) ist der betragskleinste Eigenwert aber  $\lambda_i - \lambda$ , d.h. der resultierende numerische Fehler liegt hauptsächlich in Richtung des gesuchten Eigenvektors  $x_i$ , d.h. wirkt sich kaum aus! *Das ist der eigentliche Grund, warum die inverse Iteration funktioniert.*