

# Numerik von Eigenwertaufgaben

Steffen Börm

Stand 26. Juni 2020

Alle Rechte beim Autor.



# Inhaltsverzeichnis

<b>1</b>	<b>Einleitung</b>	<b>5</b>
<b>2</b>	<b>Beispiele</b>	<b>7</b>
2.1	Schwingende Saite . . . . .	7
2.2	Minipoly . . . . .	10
2.3	Lineares Anfangswertproblem . . . . .	12
<b>3</b>	<b>Theoretische Grundlagen</b>	<b>15</b>
3.1	Existenz von Eigenwerten . . . . .	15
3.2	Ähnlichkeitstransformationen . . . . .	19
3.3	Hilberträume . . . . .	22
3.4	Invariante Unterräume . . . . .	27
3.5	Selbstadjungierte und unitäre Matrizen . . . . .	29
3.6	Schur-Zerlegung . . . . .	31
3.7	Diagonalisierbarkeit durch unitäre Transformationen . . . . .	35
3.8	Nicht-unitäre Transformationen . . . . .	44
3.9	Eigenwerte nicht-negativer Matrizen* . . . . .	49
<b>4</b>	<b>Die Jacobi-Iteration</b>	<b>57</b>
4.1	Iterierte Ähnlichkeitstransformationen . . . . .	57
4.2	Zweidimensionaler Fall . . . . .	58
4.3	Höherdimensionaler Fall . . . . .	62
4.4	Algorithmus . . . . .	65
<b>5</b>	<b>Die Vektoriteration</b>	<b>69</b>
5.1	Grundidee . . . . .	69
5.2	Fehleranalyse . . . . .	81
5.3	Inverse Iteration mit und ohne Shift . . . . .	89
5.4	Inverse Iteration mit Rayleigh-Shift . . . . .	95
5.5	Orthogonale Iteration . . . . .	100
<b>6</b>	<b>Die QR-Iteration</b>	<b>117</b>
6.1	Grundidee . . . . .	117
6.2	Shift-Strategien und Deflation . . . . .	121
6.3	Hessenberg-Form . . . . .	123
6.4	Implizite Verfahren . . . . .	132
6.5	Singulärwertzerlegung* . . . . .	139

<b>7 Verfahren für Tridiagonalmatrizen</b>	<b>147</b>
7.1 Auswertung des charakteristischen Polynoms . . . . .	147
7.2 Sturmsche Ketten . . . . .	150
7.3 Trägheitssatz und Dreieckszerlegungen . . . . .	159
<b>8 Lanczos-Verfahren für schwachbesetzte Matrizen</b>	<b>165</b>
8.1 Zweidimensionales Modellproblem . . . . .	165
8.2 Krylow-Räume . . . . .	167
8.3 Arnoldi-Basis . . . . .	170
8.4 Konvergenz . . . . .	176
<b>9 Eigenwertverfahren für sehr große Matrizen</b>	<b>185</b>
9.1 Richardson-Iteration . . . . .	185
9.2 Optimale Dämpfung . . . . .	189
9.3 Vorkonditionierer . . . . .	192
9.4 Block-Verfahren . . . . .	200
9.5 Eigenwert-Mehrgitterverfahren . . . . .	202
<b>10 Verwandte Fragestellungen</b>	<b>209</b>
10.1 Verallgemeinerte Eigenwertprobleme . . . . .	209
10.2 Selbstdjungierte positiv definite verallgemeinerte Eigenwertprobleme	213
<b>Index</b>	<b>217</b>
<b>Literaturverzeichnis</b>	<b>219</b>

# 1 Einleitung

Ein *Eigenwert* eines linearen Operators  $L : V \rightarrow V$  auf einem  $K$ -Vektorraum  $V$  ist ein Element  $\lambda \in K$  des zugehörigen Körpers, für das ein von null verschiedener Vektor  $u \in V \setminus \{0\}$  mit

$$Lu = \lambda u \tag{1.1}$$

existiert. Diesen Vektor  $u$  nennt man einen *Eigenvektor* zu  $\lambda$ , das Paar  $(\lambda, u)$  nennt man ein *Eigenpaar* des Operators  $L$ .

Eigenwerte sind in naturwissenschaftlichen Anwendungen von Interesse, beispielsweise bei der Untersuchung des Resonanzverhaltens eines schwingenden Systems, aber auch bei der Klassifikation von Dokumenten, beispielsweise der Sortierung der Ergebnisse von Internet-Suchmaschinen, und bei der mathematischen Analyse von linearen Gleichungssystemen und linearen Anfangswertproblemen.

Obwohl die Gleichung (1.1) auf den ersten Blick einem gewöhnlichen linearen Gleichungssystem ähnelt, stellt sich bei genauerer Betrachtung heraus, dass durch den Zusammenhang zwischen  $\lambda$  und  $u$  ein nichtlineares System entsteht, das sich im Allgemeinen nicht mehr mit einem aus endlich vielen Rechenoperationen bestehenden Algorithmus lösen lässt.

Stattdessen kommen *iterative Verfahren* zum Einsatz, die beispielsweise eine Folge von Näherungen eines Eigenvektors oder Eigenwerts berechnen, die gegen exakte Lösungen konvergieren.

Dabei sind verschiedene Aufgabenstellung zu unterscheiden: In manchen Anwendungen ist man nur daran interessiert, einen bestimmten Eigenwert zu berechnen, beispielsweise um die niedrigste Resonanzfrequenz eines schwingungsfähigen Systems zu ermitteln. In anderen sind eine kleine Anzahl der größten oder kleinsten Eigenwerte von Interesse oder Eigenwerte, die in der Nähe eines gegebenen Punkts in der komplexen Ebene liegen. In wieder anderen ist eine vollständige Orthonormalbasis des gesamten Raums gesucht, mit deren Hilfe sich der Operator diagonalisieren lässt, um weitere Berechnungen zu vereinfachen.

Um die unterschiedlichen Anforderungen zu erfüllen werden eine Reihe von algorithmischen Ansätzen verwendet, beispielsweise lässt sich durch Potenzieren des Operators  $L$  eine Eigenvektor-Näherung bestimmen, Eigenwerte können als Minima oder Maxima geeigneter Funktionen beschrieben werden, alternativ im endlich-dimensionalen Fall auch als Nullstellen eines charakteristischen Polynoms.

Eine besondere Herausforderung stellt die Behandlung großer Matrizen dar, die beispielsweise bei der Analyse strukturmehchanischer oder elektromagnetischer Schwingungen von großer Bedeutung sind. In diesem Fall ist  $L$  ein Differentialoperator, der im Rahmen einer Diskretisierung durch eine Matrix approximiert wird, deren kleinste Eigenwerte und zugehörige Eigenvektoren gesucht werden. Die Matrix ist in der Regel

## *1 Einleitung*

schlecht konditioniert, so dass spezialisierte Iterationsverfahren zum Einsatz kommen müssen, die mit den bereits erwähnten verwandt sind, allerdings die besondere Form der Aufgabe berücksichtigen.

Ziel dieser Vorlesung ist es, einen Überblick über die grundlegende Theorie, die wichtigsten Verfahren und deren Analyse zu geben.

## **Danksagung**

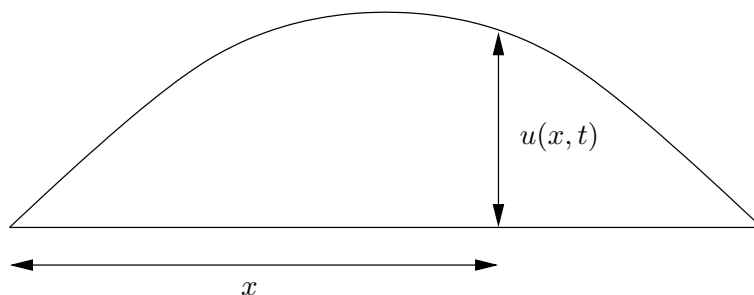
Ich bedanke mich bei Sabrina Reif und Robin-Thomas Léger für Hinweise auf Fehler in früheren Fassungen dieses Skripts und für Verbesserungsvorschläge.

## 2 Beispiele

In diesem Kapitel werden drei Beispiele für Eigenwertprobleme in der Praxis gegeben. Gleichzeitig werden drei Typen von Eigenwertproblemen charakterisiert: Die Berechnung eines einzelnen Eigenwert-Eigenvektor-Paares, die Berechnung von wenigen solchen Paaren und die Berechnung aller solcher Paare, also die Schur-Zerlegung einer Matrix.

### 2.1 Schwingende Saite

Wir untersuchen eine horizontal gespannte Saite der Länge  $\ell$ . Ihre vertikale Auslenkung in Abhängigkeit von dem Ort  $x \in [0, \ell]$  und der Zeit  $t \in \mathbb{R}_{\geq 0}$  wird durch eine Funktion  $u \in C^2([0, \ell] \times \mathbb{R}_{\geq 0})$  modelliert:



Da die Saite links und rechts eingespannt ist, gelten die Randbedingungen

$$u(0, \cdot) = 0 \quad \text{und} \quad u(\ell, \cdot) = 0. \quad (2.1)$$

Das (vereinfachte) Verhalten der Saite wird durch die *Wellengleichung*

$$\frac{\partial^2 u}{\partial t^2}(x, t) = c \frac{\partial^2 u}{\partial x^2}(x, t) + f(x, t) \quad (2.2)$$

für alle  $(x, t) \in ]0, \ell[ \times \mathbb{R}_{\geq 0}$  bestimmt. Hierbei bezeichnet  $c \in \mathbb{R}$  einen Materialparameter, in den etwa Eigenschaften wie die Dicke der Saite, ihre Elastizität und das Maß der aufgewendeten Spannkraft eingehen. Die Funktion  $f \in C([0, \ell] \times \mathbb{R}_{\geq 0})$  beschreibt die von Außen auf die Saite ausgeübte Kraft.

Wir interessieren uns für den Fall, in dem die äußere Kraft lediglich der Anregung dient und die Saite anschließend ohne weitere Beeinflussung schwingt. Wir sind also an Lösungen der Gleichung

$$\frac{\partial^2 u}{\partial t^2}(x, t) = c \frac{\partial^2 u}{\partial x^2}(x, t) \quad (2.3)$$

## 2 Beispiele

interessiert. Da die auftretenden Differentialoperatoren linear sind und keine vom Ort abhängigen Koeffizienten auftreten, entscheiden wir uns für den folgenden Separationsansatz:

$$u(x, t) = u_0(x) \cos(\omega t) \quad (2.4)$$

für eine Funktion  $u_0 \in C^2([0, \ell])$ . Dann gilt

$$\frac{\partial^2 u}{\partial t^2}(x, t) = -\omega^2 u(x, t)$$

für alle  $(x, t) \in ]0, \ell[ \times \mathbb{R}_{\geq 0}$  und die Wellengleichung (2.3) nimmt die Form

$$c \frac{\partial^2 u}{\partial x^2}(x, t) = -\omega^2 u(x, t)$$

an. Mit den Abkürzungen

$$L := -c \frac{\partial^2}{\partial x^2} \quad \text{und} \quad \lambda := \omega^2$$

und durch Elimination des Cosinus aus  $u$  erhalten wir die Gleichung

$$Lu_0 = \lambda u_0. \quad (2.5)$$

Sie besitzt offenbar immer die triviale Lösung  $u_0 = 0$ , die einer ruhenden Saite entspricht. Für eine schwingende Saite fordern wir deshalb  $u_0 \neq 0$  und erhalten das folgende Problem:

Finde  $u_0 \in C^2([0, \ell]) \setminus \{0\}$  und  $\lambda \in \mathbb{R}$  derart, dass

$$Lu_0 = \lambda u_0$$

erfüllt ist.

In diesem Kontext bezeichnet man  $\lambda$  als *Eigenwert* und  $u_0$  als *Eigenvektor* (beziehungsweise in diesem Fall als *Eigenfunktion*).

Da der Separationsansatz zur Elimination der Zeit aus der Gleichung (2.3) hilfreich war, liegt es nahe, ihn auch auf das Eigenwertproblem anzuwenden. Wir nehmen an, dass

$$u_0(x) = \sin(\alpha x)$$

für ein  $\alpha \in \mathbb{R}$  gilt. Da diese Funktion nur dann gleich Null ist, falls  $\alpha x \in \pi\mathbb{Z}$  erfüllt ist, folgt aus der Randbedingung (2.1) die Gleichung

$$\alpha \ell \in \pi\mathbb{Z}, \quad \text{also} \quad \alpha = k\pi/\ell$$

für ein  $k \in \mathbb{Z}$ . Der Operator  $L$  lässt sich für ein  $u_0$  von dieser Gestalt als

$$Lu_0(x) = c\alpha^2 u_0(x)$$



schreiben, so dass aus der Gleichung (2.5) die Beziehung

$$c\alpha^2 = \lambda = \omega^2$$

folgt. Wir können also zu jedem  $\alpha$  den entsprechenden Eigenwert  $\lambda$  und die korrespondierende *Eigenfrequenz*  $\omega = \sqrt{\lambda}$  berechnen und erhalten das Ergebnis, dass für jedes  $k \in \mathbb{N}$

$$\omega = \sqrt{c} \frac{\pi}{\ell} k$$

eine Eigenfrequenz ist, die zu einem nicht-trivialen Eigenvektor

$$u_0(x) = \sin(k\pi/\ell)$$

gehört.

Falls der Materialparameter  $c$  in der Wellengleichung nicht konstant, sondern vom Ort abhängig ist (falls sich beispielsweise die Dicke der Saite ändert), lässt sich der Separationsansatz  $u_0(x) = \sin(\alpha x)$  nicht mehr anwenden.

In diesem Fall behilft man sich mit einer numerischen Approximation des Operators  $L$ , um aus der kontinuierlichen Gleichung (2.5) eine Matrixgleichung zu machen, die dann mit Standardverfahren behandelt werden kann.

Ein Ansatz für diese Approximation ist die Finite-Differenzen-Methode, die im Wesentlichen auf der Taylor-Entwicklung von  $u_0$  basiert: Für  $x \in ]0, \ell[$  und  $h \in \mathbb{R}$  mit  $x+h, x-h \in ]0, \ell[$  gibt es  $\eta_+ \in ]x, x+h[$  und  $\eta_- \in ]x-h, x[$  derart, dass

$$\begin{aligned} u_0(x+h) &= u_0(x) + hu'_0(x) + h^2u''_0(x)/2 + h^3u_0^{(3)}(x)/6 + h^4u_0^{(4)}(\eta_+)/24, \\ u_0(x-h) &= u_0(x) - hu'_0(x) + h^2u''_0(x)/2 - h^3u_0^{(3)}(x)/6 + h^4u_0^{(4)}(\eta_-)/24 \end{aligned}$$

gelten, wobei  $u_0^{(3)}$  und  $u_0^{(4)}$  die dritte und vierte Ableitung von  $u_0$  bezeichnen. Durch Addition dieser Gleichungen erhalten wir

$$u_0(x+h) + u_0(x-h) = 2u_0(x) + h^2u''_0(x) + h^4(u_0^{(4)}(\eta_+) + u_0^{(4)}(\eta_-))/24$$

und können folgern, dass

$$Du_0(x, h) := \frac{2u_0(x) - u_0(x+h) - u_0(x-h)}{h^2} \quad (2.6)$$

eine Approximation für  $-u_0''(x)$  ist, die die Fehlerabschätzung

$$|Du_0(x, h) + u_0''(x)| \leq \frac{h^2}{12} \|u_0^{(4)}\|_\infty \quad (2.7)$$

erfüllt.

Um nun  $L$  zu approximieren, wählen wir eine äquidistante Partitionierung  $0 = x_0 < x_1 < \dots < x_n < x_{n+1} = \ell$  des Intervalls  $[0, \ell]$ : Wir setzen  $x_i = hi$  mit  $h = \ell/(n+1)$ . Für jedes  $i \in \{1, \dots, n\}$  ersetzen wir dann  $-u_0''(x_i)$  durch

$$Du_0(x_i, h) = \frac{2u_0(x_i) - u_0(x_{i+1}) - u_0(x_{i-1}))}{h^2}.$$



Die Spielregeln sind einfach:

- Zu Beginn steht eine Spielfigur auf einem beliebigen Feld.
- Falls die Figur auf einem der Felder 1, 2 oder 3 steht, wird ein üblicher sechsseitiger Würfel geworfen, dessen Augenzahl angibt, wieviele Schritte im Uhrzeigersinn durchzuführen sind.
- Falls die Figur auf dem „Gefängnisfeld“ 4 steht, wird wieder gewürfelt. Falls eine 6 herauskommt, kann die Figur auf Feld 1 fliehen, sonst bleibt sie im Gefängnis, also auf Feld 4.

Da die Position der Spielfigur nach einer Reihe von Schritten vom Zufall abhängt, können wir keine präzisen Vorhersagen treffen. Wir können allerdings nach der *Wahrscheinlichkeit* fragen, mit der eine Spielfigur auf einem bestimmten Feld stehen wird.

Eine ähnliche Technik wird beispielsweise bei der Bewertung von Internet-Seiten durch Suchmaschinen verwendet: Die ursprüngliche Idee der Suchmaschine Google beruhte auf dem Modell eines Anwenders, der ausgehend von einer Webseite zufällig einen der auf dieser Seite enthaltenen Verweise anklickt. Diejenigen Seiten, auf denen sich der simulierte Anwender mit hoher Wahrscheinlichkeit aufhielt, wurden als besonders attraktiv bewertet und in der Ergebnisliste der Suchmaschine als erste angegeben.

Die Folge der von einer Spielfigur eingenommenen Positionen bildet eine sogenannte *Markoff-Kette*, deren wahrscheinlichkeitstheoretische Eigenschaften durch die Übergangswahrscheinlichkeiten zwischen den einzelnen Feldern eindeutig bestimmt ist.

So beträgt die Wahrscheinlichkeit, von Feld 1 auf Feld 1 zu wechseln, gerade  $1/6$ , da dazu eine 4 gewürfelt werden muss. Um von Feld 1 auf Feld 2 zu wechseln, genügen dagegen eine 1 *oder* eine 5, so dass diese Wahrscheinlichkeit  $1/3$  beträgt. Wir können die Übergangswahrscheinlichkeiten in einer Matrix  $\mathbf{P} \in \mathbb{R}^{4 \times 4}$  zusammenfassen, indem wir in der  $j$ -ten Spalte und  $i$ -ten Zeile eintragen, wie hoch die Wahrscheinlichkeit ist, von Feld  $j$  auf Feld  $i$  zu wechseln:

$$\mathbf{P} := \begin{pmatrix} 1/6 & 1/6 & 1/3 & 1/6 \\ 1/3 & 1/6 & 1/6 & 0 \\ 1/3 & 1/3 & 1/6 & 0 \\ 1/6 & 1/3 & 1/3 & 5/6 \end{pmatrix} \quad (2.10)$$

Für den Vektor  $\mathbf{e}_2 = (0, 1, 0, 0)$  gibt  $\mathbf{P}\mathbf{e}_2$  an, mit welcher Wahrscheinlichkeit die Spielfigur sich auf den einzelnen Feldern befindet, wenn von Feld 2 ausgehend ein Zug durchgeführt wurde. Der Vektor  $\mathbf{P}^2\mathbf{e}_2$  beschreibt die Wahrscheinlichkeitsverteilung nach 2 Zügen, der Vektor  $\mathbf{P}^n\mathbf{e}_2$  die nach  $n$  Zügen. Die Matrix  $\mathbf{P}$  ordnet also der Wahrscheinlichkeitsverteilung in einem Zug die für den folgenden Zug zu.

Es stellt sich die Frage, ob ein Grenzwert existiert, ob sich also nach einer gewissen Anzahl von Spielzügen eine stabile Wahrscheinlichkeitsverteilung einstellt. Ein solches sogenanntes *invariantes Wahrscheinlichkeitsmaß*  $\mathbf{w}$  ist durch die Fixpunktgleichung

$$\mathbf{P}\mathbf{w} = \mathbf{w} \quad (2.11)$$

## 2 Beispiele

charakterisiert. Offensichtlich kann  $\mathbf{w}$  auch als Eigenvektor zum Eigenwert 1 aufgefasst werden, wir erhalten also wieder ein Eigenwertproblem, diesmal allerdings mit bekanntem Eigenwert und unbekanntem Eigenvektor.

Existenzsätze für Eigenwerte lassen sich auf die Matrix  $\mathbf{P}$  anwenden und führen zu dem Ergebnis, dass sie einen Eigenwert 1 besitzt, zu dem es einen Eigenvektor mit positiven Koeffizienten gibt. Mit geeigneter Skalierung kann er als das gesuchte invariante Maß interpretiert werden.

Es lassen sich außerdem Standardverfahren anwenden, um diesen Eigenvektor zu berechnen: Man erhält näherungsweise

$$\mathbf{w} \approx \begin{pmatrix} 0.185 \\ 0.097 \\ 0.113 \\ 0.605 \end{pmatrix},$$

die Spielfigur wird also im Durchschnitt über 60% der Spielzüge im Gefängnis verbringen und am seltensten auf dem Feld 2 anzutreffen sein.

Das wiederholte Anwenden von  $\mathbf{P}$  auf Wahrscheinlichkeitsverteilungen lässt sich als ein Verfahren zur Bestimmung des größten Eigenwerts, in diesem Fall 1, interpretieren, so dass sich die im Rahmen der Analyse dieses Verfahrens erzielten Konvergenzaussagen direkt auf das Beispiel übertragen.

Das in diesem Kontext auftretende Eigenwertproblem kann als Spezialfall des im Falle der schwingenden Saite behandelten gesehen werden: Es wird lediglich nach *einem* Eigenvektor zu einem bestimmten Eigenwert gesucht.

### 2.3 Lineares Anfangswertproblem

Zum Abschluss des Beispiel-Kapitels untersuchen wir ein lineares Anfangswertproblem: Wir möchten zu einer Matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  und einem Vektor  $\mathbf{y}_0 \in \mathbb{R}^n$  eine vektorwertige Funktion  $\mathbf{y} \in C^1(\mathbb{R}_{\geq 0}, \mathbb{R}^n)$  finden, die die Gleichungen

$$\mathbf{y}(0) = \mathbf{y}_0 \quad \text{und} \quad \mathbf{y}'(t) = \mathbf{A}\mathbf{y}(t) \tag{2.12}$$

für alle  $t \in \mathbb{R}_{>0}$  erfüllt. Aus der Analysis ist bekannt, dass sich  $\mathbf{y}$  durch

$$\mathbf{y}(t) = \exp(t\mathbf{A})\mathbf{y}_0$$

ausdrücken lässt. Um nun für ein konkretes  $t \in \mathbb{R}_{>0}$  den Wert  $\mathbf{y}(t)$  zu bestimmen, müssen wir die Exponentialfunktion der Matrix  $t\mathbf{A}$  auswerten. Per Taylor-Entwicklung um 0 erhalten wir die Beziehung

$$\exp(t\mathbf{A}) = \sum_{i=0}^{\infty} \frac{t^i}{i!} \mathbf{A}^i.$$

Für numerische Zwecke ist diese Formel unbrauchbar, da die Berechnung der Matrizen  $\mathbf{A}^i$  für  $i$  zwischen 0 und einer Obergrenze  $p$  immerhin  $\mathcal{O}(pn^3)$  Operationen erfordert und

$p$  infolge der langsamen Konvergenz der Exponentialreihe für große Werte von  $t$  ebenfalls groß sein muss.

Das Problem lässt sich wesentlich vereinfachen, wenn man voraussetzt, dass  $\mathbf{y}_0$  ein Eigenvektor von  $\mathbf{A}$  zum Eigenwert  $\lambda$  ist, denn dann gilt

$$\mathbf{A}\mathbf{y}_0 = \lambda\mathbf{y}_0, \quad \mathbf{A}^2\mathbf{y}_0 = \lambda^2\mathbf{y}_0, \quad \mathbf{A}^i\mathbf{y}_0 = \lambda^i\mathbf{y}_0,$$

also lässt sich die Exponentialfunktion wesentlich einfacher darstellen:

$$\mathbf{y}(t) = \exp(t\mathbf{A})\mathbf{y}_0 = \sum_{i=0}^{\infty} \frac{t^i}{i!} \mathbf{A}^i \mathbf{y}_0 = \sum_{i=0}^{\infty} \frac{t^i}{i!} \lambda^i \mathbf{y}_0 = \exp(t\lambda)\mathbf{y}_0,$$

die Auswertung von  $\mathbf{y}(t)$  erfordert also nur noch die Berechnung des Werts der Exponentialfunktion für eine reelle Zahl. Diese Aufgabe lässt sich mit einem Aufwand von höchstens  $\mathcal{O}(p)$  bewerkstelligen, ist also *wesentlich* günstiger als der ursprüngliche Ansatz.

In der Praxis wird der Startvektor  $\mathbf{y}_0$  eher selten ein Eigenvektor von  $\mathbf{A}$  sein. Falls eine Basis  $(\mathbf{v}_\ell)_{\ell=1}^n$  aus Eigenvektoren zu den Eigenwerten  $(\lambda_\ell)_{\ell=1}^n$  von  $\mathbf{A}$  existiert, lässt sich  $\mathbf{y}_0$  durch

$$\mathbf{y}_0 = \sum_{\ell=1}^n \alpha_\ell \mathbf{v}_\ell$$

darstellen, so dass sich

$$\mathbf{y}(t) = \exp(t\mathbf{A})\mathbf{y}_0 = \sum_{\ell=1}^n \alpha_\ell \exp(t\mathbf{A})\mathbf{v}_\ell = \sum_{\ell=1}^n \alpha_\ell \exp(t\lambda_\ell)\mathbf{v}_\ell$$

ergibt und sich somit mit einem Aufwand von  $\mathcal{O}(np + n^2)$  berechnen lässt, falls die Koeffizienten  $\alpha_\ell$  vorliegen.

In der Regel werden die Koeffizienten nicht vorliegen, so dass sie erst berechnet werden müssen. Für eine allgemeine Basis  $(\mathbf{v}_\ell)_{\ell=1}^n$  erfordert das einen Aufwand von  $\mathcal{O}(n^3)$ , ist also unattraktiv.

Falls die Basis  $(\mathbf{v}_\ell)_{\ell=1}^n$  allerdings *orthonormal* ist, kann auch diese Hürde überwunden werden, weil dann

$$\alpha_\ell = \langle \mathbf{v}_\ell, \mathbf{y}_0 \rangle$$

gilt und sich die Koeffizienten also mit einem Aufwand von  $\mathcal{O}(n^2)$  bestimmen lassen.

Für die Berechnung von  $\mathbf{y}(t)$  sind wir also daran interessiert, *alle* Eigenwerte und Eigenvektoren einer Matrix  $\mathbf{A}$  zu berechnen und sicherzustellen, dass die Eigenvektoren eine orthonormale Basis bilden.

Schreibt man die Eigenvektoren  $(\mathbf{v}_\ell)_{\ell=1}^n$  als Spalten einer Matrix  $\mathbf{Q} \in \mathbb{R}^{n \times n}$ , so gilt  $\mathbf{Q}^* \mathbf{Q} = \mathbf{I}$ , da wir eine orthonormale Basis verwenden,  $\mathbf{Q}$  ist unitär, und aus

$$\mathbf{A}\mathbf{Q}\delta^{(i)} = \mathbf{A}\mathbf{v}_i = \lambda_i \mathbf{v}_i = \lambda_i \mathbf{Q}\delta^{(i)},$$

## 2 Beispiele

wobei  $\delta^{(i)}$  den  $i$ -ten kanonischen Einheitsvektor bezeichnet, folgt die Gleichung

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \underbrace{\begin{pmatrix} \lambda_1 & & & & \\ & \lambda_2 & & & \\ & & \ddots & & \\ & & & \lambda_{n-1} & \\ & & & & \lambda_n \end{pmatrix}}_{=:\mathbf{D}},$$

wir haben also eine orthonormale Matrix  $\mathbf{Q}$  gefunden, die  $\mathbf{A}$  diagonalisiert. Daraus lässt sich wiederum

$$\mathbf{A} = \mathbf{Q} \mathbf{D} \mathbf{Q}^* = \mathbf{Q}^* \mathbf{D}^* \mathbf{Q}^* = (\mathbf{Q} \mathbf{D} \mathbf{Q}^*)^* = \mathbf{A}^*$$

ableiten, also die negative Aussage, dass  $\mathbf{A}$  symmetrisch sein muss, damit eine orthonormale Basis existiert, die diese Matrix diagonalisiert. Wir werden später nachweisen, dass diese Bedingung bereits hinreichend ist.

## 3 Theoretische Grundlagen

Dieses Kapitel hat zwei Ziele: Einerseits sollen die elementaren Aussagen über die Eigenwerte und Eigenvektoren quadratischer Matrizen zur Verfügung gestellt werden, andererseits werden einige grundlegende Begriffe aus der linearen Algebra rekapituliert, die für die Untersuchung der in den späteren Kapiteln eingeführten Verfahren nützlich sein werden.

### 3.1 Existenz von Eigenwerten

Wir arbeiten im Folgenden im Körper  $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$  der reellen oder komplexen Zahlen. An einigen Stellen greifen wir auf den Fundamentalsatz der Algebra zurück, so dass wir uns auf den Fall  $\mathbb{K} = \mathbb{C}$  beschränken müssen.

Die Dimensionen der behandelten Matrizen werden wir mit  $n, m \in \mathbb{N}$  bezeichnen.

**Definition 3.1 (Eigenwerte und Eigenvektoren)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Eine Zahl  $\lambda \in \mathbb{K}$  heißt Eigenwert von  $\mathbf{A}$  genau dann, wenn es einen Vektor  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  gibt, der die Gleichung

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \tag{3.1}$$

erfüllt. Jeden derartigen Vektor bezeichnet man als Eigenvektor der Matrix  $\mathbf{A}$  zu dem Eigenwert  $\lambda$ .

Falls  $\mathbf{x}$  und  $\mathbf{y}$  Eigenvektoren einer Matrix  $\mathbf{A}$  zum gleichen Eigenwert  $\lambda$  sind, gilt dasselbe für alle von null verschiedenen Linearkombinationen der Vektoren:

$$\mathbf{A}(\mathbf{x} + \alpha\mathbf{y}) = \mathbf{A}\mathbf{x} + \alpha\mathbf{A}\mathbf{y} = \lambda\mathbf{x} + \alpha\lambda\mathbf{y} = \lambda(\mathbf{x} + \alpha\mathbf{y}).$$

Insbesondere sind Eigenvektoren zu einem Eigenwert niemals eindeutig bestimmt, stattdessen definieren sie einen Teilraum.

**Definition 3.2 (Eigenräume)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ , sei  $\lambda \in \mathbb{K}$  ein Eigenwert von  $\mathbf{A}$ . Dann bezeichnet  $\mathcal{E}_A(\lambda)$  den Raum, der von den Eigenvektoren von  $\mathbf{A}$  zu  $\lambda$  aufgespannt wird. Es gilt

$$\mathcal{E}_A(\lambda) := \text{Kern}(\lambda\mathbf{I} - \mathbf{A}).$$

Der Raum  $\mathcal{E}_A(\lambda)$  heißt Eigenraum von  $\mathbf{A}$  zu dem Eigenwert  $\lambda$ .

**Lemma 3.3 (Charakterisierung der Eigenwerte)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und sei  $\lambda \in \mathbb{K}$ . Es sind äquivalent:

### 3 Theoretische Grundlagen

1.  $\lambda$  ist ein Eigenwert von  $\mathbf{A}$ ,
2.  $\lambda\mathbf{I} - \mathbf{A}$  ist singulär,
3.  $\det(\lambda\mathbf{I} - \mathbf{A}) = 0$ .

*Beweis.* „1  $\Rightarrow$  2“: Sei zunächst  $\lambda$  ein Eigenwert. Dann gilt insbesondere

$$\text{Kern}(\lambda\mathbf{I} - \mathbf{A}) = \mathcal{E}_A(\lambda) \neq \{\mathbf{0}\},$$

also ist  $\lambda\mathbf{I} - \mathbf{A}$  nicht injektiv und damit insbesondere singulär.

„2  $\Rightarrow$  1“: Sei nun  $\lambda\mathbf{I} - \mathbf{A}$  singulär. Damit existiert insbesondere ein  $\mathbf{x} \in \text{Kern}(\lambda\mathbf{I} - \mathbf{A}) \setminus \{\mathbf{0}\}$ , und es folgt  $\lambda\mathbf{x} = \mathbf{A}\mathbf{x}$ , wir haben also einen Eigenvektor zu  $\lambda$  gefunden.

„2  $\Leftrightarrow$  3“: Die Determinante einer Matrix verschwindet genau dann, wenn die Matrix singulär ist. ■

Mit Hilfe der dritten Eigenschaft aus Lemma 3.3 können wir die Eigenwerte einer Matrix als Nullstellen eines Polynoms charakterisieren:

**Definition 3.4 (Charakteristisches Polynom)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Dann ist

$$p_A(\lambda) := \det(\lambda\mathbf{I} - \mathbf{A})$$

ein Polynom  $n$ -ten Grades. Es heißt das charakteristische Polynom der Matrix  $\mathbf{A}$ .

**Lemma 3.5 (Nullstellen von  $p_A$ )** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Dann ist  $\lambda \in \mathbb{K}$  genau dann eine Nullstelle von  $p_A$ , wenn  $\lambda$  ein Eigenwert von  $\mathbf{A}$  ist.

*Beweis.* Folgt direkt aus Lemma 3.3. ■

**Übungsaufgabe 3.6 (Begleitmatrix)** (vgl. [2, Abschnitt 7.4.6]) Seien  $n \in \mathbb{N}$  und  $c_0, \dots, c_{n-1} \in \mathbb{K}$  gegeben, und sei

$$\mathbf{A} := \begin{pmatrix} 0 & 0 & \cdots & 0 & -c_0 \\ 1 & 0 & \cdots & 0 & -c_1 \\ 0 & 1 & \cdots & 0 & -c_2 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & -c_{n-1} \end{pmatrix}.$$

*Beweisen Sie*

$$p_A(\lambda) = c_0 + c_1\lambda + \dots + c_{n-1}\lambda^{n-1} + \lambda^n \quad \text{für alle } \lambda \in \mathbb{K}.$$

Die Suche nach den Nullstellen eines beliebigen Polynoms lässt sich also immer auf ein Eigenwertproblem zurückführen.

Hinweis: Man könnte den Laplace'schen Entwicklungssatz auf die erste Spalte der Matrix  $\lambda\mathbf{I} - \mathbf{A}$  anwenden.



**Übungsaufgabe 3.7 (Gauß-Quadratur)** Die Stützstellen der Gauß-Quadraturformeln sind die Nullstellen der Legendre-Polynome, die durch  $L_0(x) = 1$ ,  $L_1(x) = x$  und die Rekurrenzgleichung

$$(n+1)L_{n+1}(x) = (2n+1)xL_n(x) - nL_{n-1}(x) \quad \text{für alle } n \in \mathbb{N}, x \in \mathbb{C} \quad (3.2)$$

gegeben sind. Die Berechnung dieser Nullstellen lässt sich auf Eigenwertprobleme zurückführen: Wenn wir eine Matrix finden, deren charakteristisches Polynom (bis auf Skalierung) ein Legendre-Polynom ist, sind die Eigenwerte gerade die Nullstellen.

(a) Seien Folgen  $(a_n)_{n=1}^{\infty}$  und  $(b_n)_{n=1}^{\infty}$  in  $\mathbb{C}$  und Matrizen

$$\mathbf{T}_n := \begin{pmatrix} a_1 & \bar{b}_1 & & \\ b_1 & a_2 & \ddots & \\ & \ddots & \ddots & \bar{b}_{n-1} \\ & & b_{n-1} & a_n \end{pmatrix} \quad \text{für alle } n \in \mathbb{N} \quad (3.3)$$

gegeben. Beweisen Sie, dass die durch  $p_0(x) = 1$ ,  $p_1(x) = x - a_1$  und

$$p_{n+1}(x) = (x - a_{n+1})p_n(x) - |b_n|^2 p_{n-1}(x) \quad \text{für alle } n \in \mathbb{N}, x \in \mathbb{C}$$

gegebenen Polynome gerade

$$\det(\lambda \mathbf{I} - \mathbf{T}_n) = p_n(\lambda) \quad \text{für alle } n \in \mathbb{N}, \lambda \in \mathbb{C}$$

erfüllen, dass also die charakteristischen Polynome der Tridiagonalmatrizen einer Rekurrenzgleichung genügen.

(b) Charakteristische Polynome sind immer normiert, die Koeffizienten ihrer höchsten Potenzen sind also immer gleich eins.

Beweisen Sie, dass die Koeffizienten der höchsten Potenzen der durch (3.2) gegebenen Legendre-Polynome die Gleichungen  $\alpha_0 = 1$ ,  $\alpha_1 = 1$ ,

$$\alpha_{n+1} = \frac{2n+1}{n+1} \alpha_n \quad \text{für alle } n \in \mathbb{N}$$

erfüllen.

(c) Geben Sie Matrizen  $\mathbf{T}_n$  der Form (3.3) an, die  $\alpha_n \det(\lambda \mathbf{I} - \mathbf{T}_n) = L_n(\lambda)$  für alle  $n \in \mathbb{N}$ ,  $\lambda \in \mathbb{C}$  erfüllen.

Bevor wir aus der Charakterisierung der Eigenwerte als Nullstellen des charakteristischen Polynoms einen ersten Existenzsatz gewinnen können, führen wir einige hilfreiche Begriffe ein:

### 3 Theoretische Grundlagen

**Definition 3.8 (Spektrum, Vielfachheiten)** Sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Die Menge

$$\sigma(\mathbf{A}) = \{\lambda \in \mathbb{C} : \lambda \mathbf{I} - \mathbf{A} \text{ ist singular}\}$$

heißt Spektrum der Matrix  $\mathbf{A}$ . Man beachte, dass dabei auch reellwertige Matrizen  $\mathbf{A}$  als komplexwertig aufgefasst werden.

Für jedes  $\lambda \in \sigma(\mathbf{A})$  bezeichnen wir mit  $\mu_A^a(\lambda)$  die Vielfachheit der Nullstelle  $\lambda$  des Polynoms  $p_A$  und mit  $\mu_A^g(\lambda)$  die Dimension des Kerns von  $\lambda \mathbf{I} - \mathbf{A}$ .  $\mu_A^a(\lambda)$  und  $\mu_A^g(\lambda)$  heißen algebraische und geometrische Vielfachheit des Eigenwerts  $\lambda$ .

Nun lässt sich für den komplexwertigen Fall eine erste Existenzaussage für Eigenwerte mit Hilfe des Fundamentalsatzes der Algebra gewinnen:

**Satz 3.9 (Existenz von Eigenwerten)** Sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Dann ist  $\sigma(\mathbf{A}) \neq \emptyset$  und es gilt

$$\sum_{\lambda \in \sigma(\mathbf{A})} \mu_A^a(\lambda) = n.$$

*Beweis.* Gemäß Fundamentalsatz der Algebra zerfällt  $p_A$  in Linearfaktoren, es gibt also  $\lambda_1, \dots, \lambda_n \in \mathbb{C}$  mit

$$p_A(\lambda) = \prod_{i=1}^n (\lambda - \lambda_i) \quad \text{für alle } \lambda \in \mathbb{C}.$$

Da jedes  $\lambda_i$  eine Nullstelle von  $p_A$  ist, gilt nach Lemma 3.5

$$\sigma(\mathbf{A}) = \{\lambda_i : i \in [1 : n]\}.$$

Wegen

$$\mu_A^a(\lambda_i) = \#\{j : \lambda_j = \lambda_i\}$$

folgt

$$\sum_{\lambda \in \sigma(\mathbf{A})} \mu_A^a(\lambda) = \sum_{\lambda \in \sigma(\mathbf{A})} \#\{j : \lambda_j = \lambda\} = n. \quad \blacksquare$$

Dieser Satz gilt in dieser Form nur für den komplexwertigen Fall, wie das folgende Beispiel illustriert:

**Beispiel 3.10** Sei  $\mathbf{A} \in \mathbb{R}^{2 \times 2}$  durch

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

gegeben. Es gilt

$$p_A(\lambda) = \det(\lambda \mathbf{I} - \mathbf{A}) = \lambda^2 + 1,$$

## 3.2 Ähnlichkeitstransformationen

also besitzt  $p_A$  keine reellen Nullstellen, die Matrix  $\mathbf{A}$  also keine reellen Eigenwerte.

Fasst man  $\mathbf{A}$  als Matrix über  $\mathbb{C}$  auf, so gilt  $\sigma(\mathbf{A}) = \{i, -i\}$ , die Matrix besitzt also zwei rein imaginäre Eigenwerte, zu denen etwa

$$\mathbf{x}_1 := \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ i \end{pmatrix} \quad \text{und} \quad \mathbf{x}_2 := \frac{1}{\sqrt{2}} \begin{pmatrix} i \\ 1 \end{pmatrix}$$

Eigenvektoren sind, die sogar eine Orthonormalbasis von  $\mathbb{C}^2$  bilden (man beachte, dass im Komplexen das euklidische Skalarprodukt eine Sesquilinearform ist, keine Bilinearform).

Im dritten Beispiel aus Kapitel 2, dem linearen Anfangswertproblem, stellte sich die Frage nach der Existenz einer Basis aus Eigenvektoren einer gewissen Matrix. Da ein Eigenvektor nicht zu zwei Eigenwerten gehören kann, sind die Eigenräume zu unterschiedlichen Eigenwerten disjunkt (bis auf den in allen Eigenräumen enthaltenen Nullvektor), und es folgt, dass eine Basis aus Eigenvektoren genau dann existiert, wenn

$$\sum_{\lambda \in \sigma(A)} \mu_A^g(\lambda) = n$$

gilt. Das ist nicht immer der Fall, wie das folgende Beispiel demonstriert:

**Beispiel 3.11** *Wir untersuchen die Matrix*

$$\mathbf{A} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}.$$

Wegen

$$p_A(\lambda) = \det(\lambda \mathbf{I} - \mathbf{A}) = (\lambda - 1)^2$$

gilt  $\sigma(\mathbf{A}) = \{1\}$ . Um die geometrische Vielfachheit  $\mu_A^g(1)$  von  $\lambda = 1$  zu bestimmen, müssen wir die Dimension des Kerns von

$$\lambda \mathbf{I} - \mathbf{A} = \begin{pmatrix} 0 & -1 \\ 0 & 0 \end{pmatrix}$$

berechnen. Da das Bild dieser Matrix eindimensional ist, muss es ihr Kern ebenfalls sein. Demzufolge erhalten wir mit

$$\mathbf{x} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

zwar einen Eigenvektor, können aber keinen zweiten von ihm linear unabhängigen finden. Es gilt also  $\mu_A^g(1) = 2 > 1 = \mu_A^a(1)$ .

## 3.2 Ähnlichkeitstransformationen

Bei einer Diagonal- oder Dreiecksmatrix lassen sich die Eigenwerte und ihre algebraischen Vielfachheiten direkt an den Diagonalelementen ablesen, also wäre es nützlich, wenn wir allgemeine Matrizen auf diese spezielle Form bringen könnten.

### 3 Theoretische Grundlagen

Es stellt sich die Frage, wie eine Transformation aussehen muss, die mindestens die Eigenwerte unverändert lässt. Für die Klärung dieser Frage untersuchen wir eine Matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  mit einem Eigenwert  $\lambda \in \mathbb{K}$  und einem passenden Eigenvektor  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$ . Nach Definition gilt

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

Für eine beliebige Matrix  $\mathbf{B} \in \mathbb{K}^{n \times n}$  folgt daraus

$$\mathbf{B}\mathbf{A}\mathbf{x} = \lambda\mathbf{B}\mathbf{x}.$$

Das ist leider keine Eigenwert-Gleichung mehr. Falls  $\mathbf{B}$  regulär ist, können wir allerdings einen neuen Vektor

$$\hat{\mathbf{x}} := \mathbf{B}\mathbf{x}, \quad \mathbf{x} = \mathbf{B}^{-1}\hat{\mathbf{x}} \quad (3.4)$$

definieren und erhalten

$$\mathbf{B}\mathbf{A}\mathbf{B}^{-1}\hat{\mathbf{x}} = \lambda\hat{\mathbf{x}},$$

also wieder eine Eigenwert-Gleichung für die neue Matrix

$$\hat{\mathbf{A}} := \mathbf{B}\mathbf{A}\mathbf{B}^{-1}.$$

Der Wechsel von  $\mathbf{A}$  zu  $\hat{\mathbf{A}}$  lässt also Eigenwerte unverändert, während sich mit den Gleichungen (3.4) Eigenvektoren ineinander überführen lassen.

**Definition 3.12 (Ähnliche Matrizen)** Seien  $\mathbf{A}, \hat{\mathbf{A}} \in \mathbb{K}^{n \times n}$ . Die Matrizen  $\mathbf{A}$  und  $\hat{\mathbf{A}}$  heißen ähnlich, wenn eine reguläre Matrix  $\mathbf{B} \in \mathbb{K}^{n \times n}$  mit

$$\hat{\mathbf{A}} = \mathbf{B}\mathbf{A}\mathbf{B}^{-1} \quad (3.5)$$

existiert. Den Übergang von  $\mathbf{A}$  zu  $\hat{\mathbf{A}}$  bezeichnen wir als Ähnlichkeitstransformation.

Wir haben bereits gesehen, dass bei einer Ähnlichkeitstransformation Eigenwerte unverändert bleiben. Mit Hilfe der Determinanten-Multiplikationssatzes lässt sich sogar beweisen, dass das charakteristische Polynom ebenfalls unverändert bleibt, also insbesondere auch die algebraischen Vielfachheiten:

**Lemma 3.13 (Eigenwerte ähnlicher Matrizen)** Seien  $\mathbf{A}, \hat{\mathbf{A}} \in \mathbb{K}^{n \times n}$  ähnliche Matrizen. Dann gilt  $p_{\mathbf{A}} = p_{\hat{\mathbf{A}}}$ , also folgt insbesondere  $\sigma(\mathbf{A}) = \sigma(\hat{\mathbf{A}})$  und für alle  $\lambda \in \sigma(\mathbf{A})$  gilt  $\mu_{\mathbf{A}}^a(\lambda) = \mu_{\hat{\mathbf{A}}}^a(\lambda)$ .

*Beweis.* Sei  $\mathbf{B} \in \mathbb{K}^{n \times n}$  eine reguläre Matrix mit  $\hat{\mathbf{A}} = \mathbf{B}\mathbf{A}\mathbf{B}^{-1}$ . Aus dem Determinanten-Multiplikationssatz folgt

$$\begin{aligned} p_{\mathbf{A}}(\lambda) &= \det(\lambda\mathbf{I} - \mathbf{A}) = \det(\mathbf{I}) \det(\lambda\mathbf{I} - \mathbf{A}) \\ &= \det(\mathbf{B}\mathbf{B}^{-1}) \det(\lambda\mathbf{I} - \mathbf{A}) = \det(\mathbf{B}) \det(\lambda\mathbf{I} - \mathbf{A}) \det(\mathbf{B}^{-1}) \\ &= \det(\mathbf{B}(\lambda\mathbf{I} - \mathbf{A})\mathbf{B}^{-1}) = \det(\lambda\mathbf{B}\mathbf{B}^{-1} - \mathbf{B}\mathbf{A}\mathbf{B}^{-1}) \end{aligned}$$

$$= \det(\lambda \mathbf{I} - \widehat{\mathbf{A}}) = p_{\widehat{\mathbf{A}}}(\lambda)$$

für alle  $\lambda \in \mathbb{K}$ . ■

Eigenvektoren bleiben unter Ähnlichkeitstransformationen nicht unverändert, allerdings lässt sich explizit angeben, wie sie sich ändern. Insbesondere folgt, dass auch die geometrischen Vielfachheiten unverändert bleiben.

**Lemma 3.14 (Eigenvektoren)** *Seien  $\mathbf{A}, \widehat{\mathbf{A}} \in \mathbb{K}^{n \times n}$  ähnliche Matrizen und sei  $\mathbf{B} \in \mathbb{K}^{n \times n}$  eine reguläre Matrix mit  $\widehat{\mathbf{A}} = \mathbf{B}\mathbf{A}\mathbf{B}^{-1}$ . Für jeden Eigenwert  $\lambda \in \sigma(\mathbf{A}) = \sigma(\widehat{\mathbf{A}})$  gilt die Gleichung*

$$\mathcal{E}_{\widehat{\mathbf{A}}}(\lambda) = \mathbf{B}\mathcal{E}_{\mathbf{A}}(\lambda).$$

Insbesondere folgt  $\mu_{\widehat{\mathbf{A}}}^g(\lambda) = \mu_{\mathbf{A}}^g(\lambda)$ .

*Beweis.* Sei  $\lambda \in \sigma(\mathbf{A}) = \sigma(\widehat{\mathbf{A}})$ .

„ $\supseteq$ “: Sei  $\mathbf{x} \in \mathcal{E}_{\mathbf{A}}(\lambda)$ . Mit  $\widehat{\mathbf{x}} := \mathbf{B}\mathbf{x}$  gilt

$$\widehat{\mathbf{A}}\widehat{\mathbf{x}} = \mathbf{B}\mathbf{A}\mathbf{B}^{-1}\mathbf{B}\mathbf{x} = \mathbf{B}\mathbf{A}\mathbf{x} = \mathbf{B}\lambda\mathbf{x} = \lambda\widehat{\mathbf{x}},$$

also ist  $\widehat{\mathbf{x}} = \mathbf{B}\mathbf{x}$  ein Eigenvektor der Matrix  $\widehat{\mathbf{A}}$  zu dem Eigenwert  $\lambda$  und somit ein Element des Eigenraums  $\mathcal{E}_{\widehat{\mathbf{A}}}(\lambda)$ .

„ $\subseteq$ “: Sei  $\widehat{\mathbf{x}} \in \mathcal{E}_{\widehat{\mathbf{A}}}(\lambda)$ . Mit  $\mathbf{x} := \mathbf{B}^{-1}\widehat{\mathbf{x}}$  gilt

$$\mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{B}^{-1}\widehat{\mathbf{x}} = \mathbf{B}^{-1}\mathbf{B}\mathbf{A}\mathbf{B}^{-1}\widehat{\mathbf{x}} = \mathbf{B}^{-1}\widehat{\mathbf{A}}\widehat{\mathbf{x}} = \mathbf{B}^{-1}\lambda\widehat{\mathbf{x}} = \lambda\mathbf{x},$$

also ist  $\mathbf{x} = \mathbf{B}^{-1}\widehat{\mathbf{x}}$  ein Eigenvektor der Matrix  $\mathbf{A}$  zu dem Eigenwert  $\lambda$  und somit ein Element des Eigenraums  $\mathcal{E}_{\mathbf{A}}(\lambda)$ . ■

Von besonderem Interesse sind Ähnlichkeitstransformationen, die eine Matrix auf Diagonalgestalt bringen, denn bei einer Diagonalmatrix lassen sich nicht nur die Eigenwerte unmittelbar ablesen, sondern wir können die durch

$$\delta_i^{(j)} := \begin{cases} 1 & \text{falls } i = j, \\ 0 & \text{ansonsten} \end{cases} \quad \text{für alle } i, j \in [1 : n]$$

definierten *kanonischen Einheitsvektoren*  $\delta^{(1)}, \dots, \delta^{(n)} \in \mathbb{K}^n$  als Eigenvektoren verwenden, die aufgrund ihrer einfachen Gestalt häufig große Vorzüge bieten.

Beispielsweise können wir kanonische Einheitsvektoren in Kombination mit einer geeignet gewählten Ähnlichkeitstransformation verwenden, um eine Beziehung zwischen der geometrischen und der algebraischen Vielfachheit eines Eigenwert zu gewinnen.

**Lemma 3.15** *Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ , sei  $\lambda \in \sigma(\mathbf{A})$ . Dann gilt*

$$1 \leq \mu_{\mathbf{A}}^g(\lambda) \leq \mu_{\mathbf{A}}^a(\lambda).$$

### 3 Theoretische Grundlagen

*Beweis.* Die erste Ungleichung folgt direkt aus Lemma 3.3. Zum Nachweis der zweiten konstruieren wir eine geeignete Ähnlichkeitstransformation, indem wir  $p := \mu_A^g(\lambda)$  linear unabhängige Eigenvektoren  $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{K}^n$  von  $\mathbf{A}$  zu dem Eigenwert  $\lambda$  wählen und sie zu einer Basis  $(\mathbf{x}_i)_{i=1}^n$  ergänzen. Wir bezeichnen mit  $\mathbf{X} \in \mathbb{K}^{n \times n}$  die reguläre Matrix, deren Spalten die Vektoren  $(\mathbf{x}_i)_{i=1}^n$  sind, die also gerade

$$\mathbf{X}\delta^{(i)} = \mathbf{x}_i \quad \text{für alle } i \in [1 : n]$$

erfüllt. Daraus folgt

$$\mathbf{A}\mathbf{X}\delta^{(i)} = \mathbf{A}\mathbf{x}_i = \lambda_i\mathbf{x}_i = \lambda_i\mathbf{X}\delta^{(i)} \quad \text{für alle } i \in [1 : p],$$

so dass wir insbesondere

$$\mathbf{X}^{-1}\mathbf{A}\mathbf{X}\delta^{(i)} = \lambda_i\delta^{(i)} \quad \text{für alle } i \in [1 : p]$$

erhalten. Also verschwinden in den ersten  $p$  Spalten der Matrix  $\widehat{\mathbf{A}} := \mathbf{X}^{-1}\mathbf{A}\mathbf{X}$  jeweils die unteren  $n - p$  Zeilen, so dass die Matrix die Form

$$\widehat{\mathbf{A}} = \begin{pmatrix} \lambda\mathbf{I} & \mathbf{R} \\ 0 & \mathbf{C} \end{pmatrix}$$

für Matrizen  $\mathbf{R} \in \mathbb{K}^{p \times (n-p)}$  und  $\mathbf{C} \in \mathbb{K}^{(n-p) \times (n-p)}$  aufweist. Da  $\widehat{\mathbf{A}}$  und  $\mathbf{A}$  ähnlich sind, gilt nach Lemma 3.13  $p_A = p_{\widehat{\mathbf{A}}}$ , also folgt für alle  $\alpha \in \mathbb{K}$  die Gleichung

$$\begin{aligned} p_A(\alpha) &= p_{\widehat{\mathbf{A}}}(\alpha) = \det(\alpha\mathbf{I} - \widehat{\mathbf{A}}) \\ &= \det(\alpha\mathbf{I} - \lambda\mathbf{I}) \det(\alpha\mathbf{I} - \mathbf{C}) = (\alpha - \lambda)^p p_C(\alpha), \end{aligned}$$

und damit ist  $\lambda$  eine mindestens  $p$ -fache Nullstelle von  $p_A$ . ■

### 3.3 Hilberträume

Neben den bisher eingeführten algebraischen Techniken haben sich auch Konzepte der Analysis als sehr nützlich für die Untersuchung von Eigenwertproblemen erwiesen. Für uns sind dabei vor allem bestimmte Eigenschaften von Bedeutung, die in *Hilberträumen* gelten, also in Banach-Räumen, deren Norm von einem Skalarprodukt induziert ist.

Auf dem Raum  $\mathbb{K}^n$  verwendet man typischerweise das *euklidische Skalarprodukt*, das durch

$$\langle \mathbf{x}, \mathbf{y} \rangle := \sum_{i=1}^n \bar{x}_i y_i \quad \text{für alle } \mathbf{x}, \mathbf{y} \in \mathbb{K}^n$$

definiert ist. Es ist eine *Sesquilinearform*, erfüllt also die Gleichungen

$$\langle \mathbf{x} + \alpha\mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle + \bar{\alpha}\langle \mathbf{y}, \mathbf{z} \rangle, \quad (3.6a)$$

$$\langle \mathbf{x}, \mathbf{y} + \alpha \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{y} \rangle + \alpha \langle \mathbf{x}, \mathbf{z} \rangle, \quad (3.6b)$$

$$\langle \mathbf{x}, \mathbf{y} \rangle = \overline{\langle \mathbf{y}, \mathbf{x} \rangle} \quad \text{für alle } \mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{K}^n, \alpha \in \mathbb{K}. \quad (3.6c)$$

Das euklidische Skalarprodukt induziert die *euklidische Norm*

$$\|\mathbf{x}\| := \langle \mathbf{x}, \mathbf{x} \rangle^{1/2} = \left( \sum_{i=1}^n |x_i|^2 \right)^{1/2} \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n,$$

und aus dieser Beziehung folgt die *Cauchy-Schwarz-Ungleichung*

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\| \quad \text{für alle } \mathbf{x}, \mathbf{y} \in \mathbb{K}^n. \quad (3.7)$$

Beide Seiten dieser Ungleichung sind genau dann gleich, wenn die Vektoren  $\mathbf{x}$  und  $\mathbf{y}$  linear abhängig sind.

Eine nützliche Beziehung zwischen der Matrix-Vektor-Multiplikation und dem Skalarprodukt können wir mit Hilfe der *adjungierten Matrix* gewinnen:

**Definition 3.16 (Adjungierte Matrix)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times m}$ . Die durch

$$b_{ij} = \bar{a}_{ji} \quad \text{für alle } i \in [1 : m], j \in [1 : n]$$

definierte Matrix  $\mathbf{B} \in \mathbb{K}^{m \times n}$  heißt Adjungierte von  $\mathbf{A}$  und wird mit  $\mathbf{A}^* = \mathbf{B}$  bezeichnet. Im Falle  $\mathbb{K} = \mathbb{R}$  entspricht sie der transponierten Matrix  $\mathbf{A}^T$ , im Falle  $\mathbb{K} = \mathbb{C}$  der hermiteschen Matrix  $\mathbf{A}^H = \bar{\mathbf{A}}^T$ .

Bei unseren Untersuchungen spielt die Beziehung zwischen dem Skalarprodukt und der Adjungierten eine wichtige Rolle.

**Lemma 3.17 (Adjungierte und Skalarprodukt)** Seien eine Matrix  $\mathbf{A} \in \mathbb{K}^{n \times m}$  sowie Vektoren  $\mathbf{x} \in \mathbb{K}^m$  und  $\mathbf{y} \in \mathbb{K}^n$  gegeben. Dann gilt

$$\langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{A}^*\mathbf{y} \rangle.$$

*Beweis.* Nach Definition des Skalarprodukts und der Matrix-Vektor-Multiplikation gilt

$$\langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle = \sum_{i=1}^n \overline{\left( \sum_{j=1}^m a_{ij} x_j \right)} y_i = \sum_{i=1}^n \sum_{j=1}^m \bar{a}_{ij} \bar{x}_j y_i = \sum_{j=1}^m \bar{x}_j \left( \sum_{i=1}^n \bar{a}_{ij} y_i \right) = \langle \mathbf{x}, \mathbf{A}^*\mathbf{y} \rangle.$$

Das ist die gesuchte Identität. ■

Diese Gleichung lässt sich vielfältig einsetzen. Als Beispiel beweisen wir die folgende Aussage über die Adjungierte eines Produkts zweier Matrizen:

**Lemma 3.18 (Adjungierte eines Produkts)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times m}$  und  $\mathbf{B} \in \mathbb{K}^{m \times k}$  mit  $n, m, k \in \mathbb{N}$  gegeben. Dann gilt

$$(\mathbf{A}\mathbf{B})^* = \mathbf{B}^* \mathbf{A}^*.$$

Falls  $\mathbf{A} \in \mathbb{K}^{n \times n}$  invertierbar ist, gilt  $(\mathbf{A}^*)^{-1} = (\mathbf{A}^{-1})^*$ .

### 3 Theoretische Grundlagen

*Beweis.* Nach Lemma 3.17 gelten die Gleichungen

$$\langle \mathbf{x}, (\mathbf{AB})^* \mathbf{y} \rangle = \langle \mathbf{ABx}, \mathbf{y} \rangle = \langle \mathbf{Bx}, \mathbf{A}^* \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{B}^* \mathbf{A}^* \mathbf{y} \rangle \quad \text{für alle } \mathbf{x} \in \mathbb{K}^k, \mathbf{y} \in \mathbb{K}^n.$$

Daraus folgt

$$\langle \mathbf{x}, ((\mathbf{AB})^* - \mathbf{B}^* \mathbf{A}^*) \mathbf{y} \rangle = 0 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^k, \mathbf{y} \in \mathbb{K}^n,$$

und indem wir  $\mathbf{x} := ((\mathbf{AB})^* - \mathbf{B}^* \mathbf{A}^*) \mathbf{y}$  einsetzen erhalten wir

$$\|((\mathbf{AB})^* - \mathbf{B}^* \mathbf{A}^*) \mathbf{y}\|^2 = 0 \quad \text{für alle } \mathbf{y} \in \mathbb{K}^n,$$

so dass sich unmittelbar die erste Gleichung ergibt.

Sei nun  $\mathbf{A} \in \mathbb{K}^{n \times n}$  invertierbar. Dann gelten mit der ersten Gleichung

$$\mathbf{I} = \mathbf{I}^* = (\mathbf{A}^{-1} \mathbf{A})^* = \mathbf{A}^* (\mathbf{A}^{-1})^*, \quad \mathbf{I} = \mathbf{I}^* = (\mathbf{A} \mathbf{A}^{-1})^* = (\mathbf{A}^{-1})^* \mathbf{A}^*,$$

also ist  $(\mathbf{A}^{-1})^*$  eine Rechts- und Linksinverse der Matrix  $\mathbf{A}^*$ , also deren Inverse. ■

Bei der Analyse numerischer Näherungsverfahren stellt sich häufig die Frage danach, wie sich in den einzelnen Stufen des Algorithmus' eingeführte Approximationsfehler auf den Gesamtfehler auswirken. Deshalb müssen wir in der Lage sein, zu messen, wie sehr sich Matrizen unterscheiden. Zu diesem Zweck definieren wir die *Spektralnorm*:

**Definition 3.19 (Spektralnorm)** Seien  $n, m \in \mathbb{N}$ . Die Spektralnorm auf  $\mathbb{K}^{n \times m}$  ist durch

$$\|\mathbf{A}\| := \max \left\{ \frac{\|\mathbf{Az}\|}{\|\mathbf{z}\|} : \mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\} \right\} \quad \text{für alle } \mathbf{A} \in \mathbb{K}^{n \times m}$$

definiert.

Einige wesentliche Eigenschaften der Spektralnorm fasst das folgende Lemma zusammen.

**Lemma 3.20 (Spektralnorm)** Seien  $n, m, k \in \mathbb{N}$ . Die Spektralnorm ist verträglich mit der euklidischen Norm, es gilt nämlich

$$\|\mathbf{Az}\| \leq \|\mathbf{A}\| \|\mathbf{z}\| \quad \text{für alle } \mathbf{A} \in \mathbb{K}^{n \times m}, \mathbf{z} \in \mathbb{K}^m. \quad (3.8a)$$

Die Spektralnorm ist auch submultiplikativ, erfüllt also

$$\|\mathbf{AB}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \quad \text{für alle } \mathbf{A} \in \mathbb{K}^{n \times m}, \mathbf{B} \in \mathbb{K}^{m \times k}. \quad (3.8b)$$

Die Spektralnorm lässt sich alternativ durch

$$\|\mathbf{A}\| = \max \left\{ \frac{|\langle \mathbf{y}, \mathbf{Az} \rangle|}{\|\mathbf{y}\| \|\mathbf{z}\|} : \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}, \mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\} \right\} \quad \text{für alle } \mathbf{A} \in \mathbb{K}^{n \times m} \quad (3.8c)$$

darstellen. Aus dieser Gleichung folgen die Beziehungen

$$\|\mathbf{A}\| = \|\mathbf{A}^*\| = \|\mathbf{A}^* \mathbf{A}\|^{1/2} = \|\mathbf{A} \mathbf{A}^*\|^{1/2} \quad \text{für alle } \mathbf{A} \in \mathbb{K}^{n \times m}. \quad (3.8d)$$



*Beweis.* Seien  $\mathbf{A} \in \mathbb{K}^{n \times m}$  und  $\mathbf{z} \in \mathbb{K}^m$ . Falls  $\mathbf{z} = \mathbf{0}$  gilt, folgt  $\|\mathbf{Az}\| = 0 = \|\mathbf{A}\|\|\mathbf{z}\|$  unmittelbar.

Ansonsten gilt nach Definition

$$\frac{\|\mathbf{Az}\|}{\|\mathbf{z}\|} \leq \|\mathbf{A}\|,$$

und Multiplikation mit  $\|\mathbf{z}\|$  führt zu (3.8a).

Indem wir (3.8a) zweimal anwenden erhalten wir

$$\|\mathbf{ABz}\| \leq \|\mathbf{A}\| \|\mathbf{Bz}\| \leq \|\mathbf{A}\| \|\mathbf{B}\| \|\mathbf{z}\| \quad \text{für alle } \mathbf{z} \in \mathbb{K}^k,$$

so dass sich aus Definition 3.19 unmittelbar die Ungleichung (3.8b) ergibt.

Sei  $\mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\}$  gegeben. Mit der Cauchy-Schwarz-Ungleichung (3.7) erhalten wir

$$\frac{|\langle \mathbf{y}, \mathbf{Az} \rangle|}{\|\mathbf{y}\| \|\mathbf{z}\|} \leq \frac{\|\mathbf{y}\| \|\mathbf{Az}\|}{\|\mathbf{y}\| \|\mathbf{z}\|} = \frac{\|\mathbf{Az}\|}{\|\mathbf{z}\|} \quad \text{für alle } \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}. \quad (3.9)$$

Falls  $\mathbf{y} := \mathbf{Az} \neq \mathbf{0}$  gilt, haben wir

$$\frac{|\langle \mathbf{y}, \mathbf{Az} \rangle|}{\|\mathbf{y}\| \|\mathbf{z}\|} = \frac{|\langle \mathbf{Az}, \mathbf{Az} \rangle|}{\|\mathbf{Az}\| \|\mathbf{z}\|} = \frac{\|\mathbf{Az}\|^2}{\|\mathbf{Az}\| \|\mathbf{z}\|} = \frac{\|\mathbf{Az}\|}{\|\mathbf{z}\|}$$

und gemeinsam mit (3.9) ergibt sich

$$\frac{\|\mathbf{Az}\|}{\|\mathbf{z}\|} \leq \max \left\{ \frac{|\langle \mathbf{y}, \mathbf{Az} \rangle|}{\|\mathbf{y}\| \|\mathbf{z}\|} : \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\} \right\} \leq \frac{\|\mathbf{Az}\|}{\|\mathbf{z}\|}.$$

Diese Abschätzung bleibt offenbar auch korrekt, falls  $\mathbf{Az} = \mathbf{0}$  gilt. Indem wir zu dem Maximum über alle  $\mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\}$  übergehen folgt (3.8c).

Aus dieser Gleichung folgt mit Lemma 3.17 und (3.6c) direkt

$$\begin{aligned} \|\mathbf{A}\| &= \max \left\{ \frac{|\langle \mathbf{y}, \mathbf{Az} \rangle|}{\|\mathbf{y}\| \|\mathbf{z}\|} : \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}, \mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\} \right\} \\ &= \max \left\{ \frac{|\langle \mathbf{A}^* \mathbf{y}, \mathbf{z} \rangle|}{\|\mathbf{y}\| \|\mathbf{z}\|} : \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}, \mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\} \right\} \\ &= \max \left\{ \frac{|\langle \mathbf{z}, \mathbf{A}^* \mathbf{y} \rangle|}{\|\mathbf{z}\| \|\mathbf{y}\|} : \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}, \mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\} \right\} = \|\mathbf{A}^*\|. \end{aligned}$$

Mit dieser Gleichung in Kombination mit Lemma 3.17, (3.8c) sowie (3.8b) erhalten wir

$$\begin{aligned} \|\mathbf{A}\|^2 &= \max \left\{ \frac{\|\mathbf{Az}\|^2}{\|\mathbf{z}\|^2} : \mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\} \right\} \\ &= \max \left\{ \frac{|\langle \mathbf{Az}, \mathbf{Az} \rangle|}{\|\mathbf{z}\|^2} : \mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\} \right\} \\ &= \max \left\{ \frac{|\langle \mathbf{z}, \mathbf{A}^* \mathbf{Az} \rangle|}{\|\mathbf{z}\|^2} : \mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\} \right\} \end{aligned}$$

### 3 Theoretische Grundlagen

$$\begin{aligned} &\leq \max \left\{ \frac{|\langle \mathbf{y}, \mathbf{A}^* \mathbf{A} \mathbf{z} \rangle|}{\|\mathbf{y}\| \|\mathbf{z}\|} : \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}, \mathbf{z} \in \mathbb{K}^m \setminus \{\mathbf{0}\} \right\} \\ &= \|\mathbf{A}^* \mathbf{A}\| \leq \|\mathbf{A}^*\| \|\mathbf{A}\| = \|\mathbf{A}\|^2. \end{aligned}$$

Indem wir dieses Resultat auf  $\mathbf{A}^*$  anstelle von  $\mathbf{A}$  anwenden folgt die letzte Gleichung. ■

In einigen der folgenden Beweise benötigen wir eine Möglichkeit, rechteckige Matrizen zu „invertieren“. Diese Pseudo-Inverse können wir mit Hilfe der Adjungierten definieren. Ein erster Schritt in dieser Richtung ist das folgende Lemma, mit dem sich die Identität von Vektoren im Bild der Matrix überprüfen lässt.

**Lemma 3.21 (Orthogonalität)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times m}$ . Es gilt

$$\langle \mathbf{x}, \mathbf{y} \rangle = 0 \quad \text{für alle } \mathbf{x} \in \text{Bild}(\mathbf{A}), \mathbf{y} \in \text{Kern}(\mathbf{A}^*).$$

Insbesondere haben wir

$$\mathbf{A} \mathbf{z} = \mathbf{0} \iff \mathbf{A}^* \mathbf{A} \mathbf{z} = \mathbf{0} \quad \text{für alle } \mathbf{z} \in \mathbb{K}^m.$$

*Beweis.* Seien  $\mathbf{x} \in \text{Bild}(\mathbf{A})$  und  $\mathbf{y} \in \text{Kern}(\mathbf{A}^*)$  gegeben. Nach Definition finden wir  $\mathbf{z} \in \mathbb{K}^m$  mit  $\mathbf{x} = \mathbf{A} \mathbf{z}$ . Mit Lemma 3.17 folgt

$$\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{A} \mathbf{z}, \mathbf{y} \rangle = \langle \mathbf{z}, \mathbf{A}^* \mathbf{y} \rangle = \langle \mathbf{z}, \mathbf{0} \rangle = 0.$$

Sei nun  $\mathbf{z} \in \mathbb{K}^m$ . Offenbar folgt aus  $\mathbf{A} \mathbf{z} = \mathbf{0}$  unmittelbar auch  $\mathbf{A}^* \mathbf{A} \mathbf{z} = \mathbf{A}^* \mathbf{0} = \mathbf{0}$ .

Für den Nachweis der Umkehrung setzen wir voraus, dass  $\mathbf{A}^* \mathbf{A} \mathbf{z} = \mathbf{0}$  gilt. Also liegt der Vektor  $\mathbf{x} := \mathbf{A} \mathbf{z}$  sowohl im Bild der Matrix  $\mathbf{A}$  als auch im Kern der Adjungierten  $\mathbf{A}^*$ . Mit dem ersten Teil unserer Aussage folgt

$$\|\mathbf{x}\|^2 = \langle \mathbf{x}, \mathbf{x} \rangle = 0,$$

also  $\mathbf{A} \mathbf{z} = \mathbf{x} = \mathbf{0}$ . ■

Aus dieser Eigenschaft folgt bereits, dass zwei Vektoren aus dem Bild einer Matrix  $\mathbf{A}$  genau dann identisch sind, falls ihre Produkte mit der Adjungierten  $\mathbf{A}^*$  identisch sind. Um auch passende Urbilder rekonstruieren zu können, benötigen wir zusätzliche Eigenschaften der Matrix.

**Definition 3.22 (Positiv definit)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Die Matrix  $\mathbf{A}$  heißt positiv definit, falls

$$\langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle > 0 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\} \quad (3.10)$$

gilt. Sie heißt positiv semidefinit, falls die Ungleichung

$$\langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle \geq 0 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n \quad (3.11)$$

erfüllt ist.

Eine positiv definite Matrix ist immer invertierbar, denn aus (3.10) folgt unmittelbar, dass nur der Nullvektor im Kern der Matrix liegen kann, dass  $\mathbf{A}$  also injektiv ist. Da  $\mathbf{A}$  eine quadratische Matrix ist, muss sie nach dem Dimensionssatz für lineare Abbildungen auch invertierbar sein.

Wenn wir für einen Vektor  $\mathbf{x}$  aus dem Bild einer Matrix  $\mathbf{A}$  ein Urbild rekonstruieren wollen, also einen Vektor  $\mathbf{y}$  mit  $\mathbf{x} = \mathbf{A}\mathbf{y}$  suchen, können wir uns dem Problem nähern, indem wir beide Seiten der Gleichung mit der Adjungierten  $\mathbf{A}^*$  multiplizieren und

$$\mathbf{A}^*\mathbf{x} = \mathbf{A}^*\mathbf{A}\mathbf{y}$$

erhalten. Falls die *Gramsche Matrix*  $\mathbf{A}^*\mathbf{A}$  auf der rechten Seite invertierbar ist, können wir die Gleichung mit deren Inversen multiplizieren und eine explizite Formel erhalten, mit der sich  $\mathbf{y}$  aus  $\mathbf{x}$  berechnen lässt.

**Lemma 3.23 (Gramsche Matrizen)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times m}$ . Dann sind  $\mathbf{A}^*\mathbf{A}$  und  $\mathbf{A}\mathbf{A}^*$  positiv semidefinite Matrizen.

$\mathbf{A}$  ist genau dann injektiv, wenn  $\mathbf{A}^*\mathbf{A}$  positiv definit ist.

$\mathbf{A}$  ist genau dann surjektiv, wenn  $\mathbf{A}\mathbf{A}^*$  positiv definit ist.

*Beweis.* Mit Lemma 3.17 erhalten wir

$$\langle \mathbf{A}^*\mathbf{A}\mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{x} \rangle = \|\mathbf{A}\mathbf{x}\|^2 \geq 0, \quad (3.12a)$$

$$\langle \mathbf{A}\mathbf{A}^*\mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{A}^*\mathbf{x}, \mathbf{A}^*\mathbf{x} \rangle = \|\mathbf{A}^*\mathbf{x}\|^2 \geq 0 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n, \quad (3.12b)$$

also sind  $\mathbf{A}^*\mathbf{A}$  und  $\mathbf{A}\mathbf{A}^*$  positiv semidefinit.

Falls  $\mathbf{A}$  injektiv ist, gilt  $\|\mathbf{A}\mathbf{x}\| > 0$  für alle  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$ , also ist nach (3.12a) die Matrix  $\mathbf{A}^*\mathbf{A}$  positiv definit.

Falls umgekehrt  $\mathbf{A}^*\mathbf{A}$  positiv definit ist, gilt nach (3.12a) die Ungleichung  $\|\mathbf{A}\mathbf{x}\| > 0$  für alle  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$ , also kann der Kern der Matrix  $\mathbf{A}$  nur den Nullvektor enthalten. Damit muss  $\mathbf{A}$  injektiv sein.

Falls  $\mathbf{A}$  surjektiv ist, falls wir also  $\text{Bild}(\mathbf{A}) = \mathbb{K}^n$  haben, folgt aus Lemma 3.21 bereits  $\text{Kern}(\mathbf{A}^*) = \{\mathbf{0}\}$ , also gilt insbesondere  $\|\mathbf{A}^*\mathbf{x}\| = 0$  genau dann, wenn  $\mathbf{x} = \mathbf{0}$  gilt. Also muss nach (3.12b) die Matrix  $\mathbf{A}\mathbf{A}^*$  positiv definit sein.

Falls umgekehrt  $\mathbf{A}\mathbf{A}^*$  positiv definit, also insbesondere invertierbar ist, können wir zu jedem  $\mathbf{x} \in \mathbb{K}^n$  den Vektor  $\mathbf{y} := \mathbf{A}^*(\mathbf{A}\mathbf{A}^*)^{-1}\mathbf{x}$  definieren und erhalten

$$\mathbf{A}\mathbf{y} = \mathbf{A}\mathbf{A}^*(\mathbf{A}\mathbf{A}^*)^{-1}\mathbf{x} = \mathbf{x},$$

also folgt  $\mathbf{x} \in \text{Bild}(\mathbf{A})$  für beliebige  $\mathbf{x} \in \mathbb{K}^n$ , so dass  $\mathbf{A}$  surjektiv sein muss. ■

## 3.4 Invariante Unterräume

Wie bereits gesehen, lässt sich der Raum  $\mathbb{K}^n$  im allgemeinen nicht in eine direkte Summe von Eigenräumen zerlegen. Es ist allerdings möglich, die Eigenräume durch allgemeinere Teilräume von  $\mathbb{K}^n$  zu ersetzen, die dann eine direkte Summe bilden:

### 3 Theoretische Grundlagen

**Definition 3.24 (Invariante Unterräume)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Ein Teilraum  $\mathcal{V} \subseteq \mathbb{K}^n$  heißt bezüglich  $\mathbf{A}$  invarianter Unterraum, falls für alle  $\mathbf{x} \in \mathcal{V}$  die Gleichung

$$\mathbf{A}\mathbf{x} \in \mathcal{V}$$

gilt, falls also  $\mathbf{A}\mathcal{V} \subseteq \mathcal{V}$  erfüllt ist.

**Beispiel 3.25** Sei  $p \in \mathbb{N}$  und  $(\mathbf{x}_i)_{i=1}^p$  eine Familie von Eigenvektoren der Eigenwerte  $(\lambda_i)_{i=1}^p$ . Dann ist

$$\mathcal{V} := \text{span}\{\mathbf{x}_i : i \in [1 : p]\}$$

ein bezüglich  $\mathbf{A}$  invarianter Unterraum.

*Beweis.* Offensichtlich gilt  $\mathbf{A}\mathbf{x}_i = \lambda_i\mathbf{x}_i \in \mathcal{V}$  für alle  $i \in [1 : p]$ , also gilt dasselbe auch für den Aufspann dieser Vektoren. ■

**Beispiel 3.26** Sei  $\mathbf{R} \in \mathbb{K}^{n \times n}$  eine obere Dreiecksmatrix. Dann ist für jedes  $p \in [1 : n]$  der Raum

$$\mathcal{V} := \text{span}\{\delta^{(i)} : i \in [1 : p]\} = \mathbb{K}^p \times \{\mathbf{0}\} \subseteq \mathbb{K}^n$$

invariant bezüglich  $\mathbf{R}$ .

Für Eigenräume gilt die Gleichung  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  mit dem Eigenwert  $\lambda$ . Für invariante Unterräume wird diese Gleichung etwas verallgemeinert, indem wir den Eigenwert  $\lambda$  durch eine kleine quadratische Matrix ersetzen und den Eigenvektor  $\mathbf{x}$  durch eine aus mehreren Vektoren gebildete Matrix.

**Lemma 3.27** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und sei  $\mathcal{V} \subseteq \mathbb{K}^n$  ein bezüglich  $\mathbf{A}$  invarianter Unterraum. Sei  $\mathbf{X} \in \mathbb{K}^{n \times p}$  eine Matrix, deren Spalten den Raum  $\mathcal{V}$  aufspannen, die also die Gleichung

$$\text{Bild}(\mathbf{X}) = \mathcal{V}$$

erfüllt. Dann gibt es eine Matrix  $\mathbf{\Lambda} \in \mathbb{K}^{p \times p}$  mit

$$\mathbf{A}\mathbf{X} = \mathbf{X}\mathbf{\Lambda}. \tag{3.13}$$

Falls  $\mathbf{X}$  injektiv ist, ist  $\mathbf{\Lambda}$  durch die obige Gleichung eindeutig festgelegt.

*Beweis.* Sei  $j \in [1 : p]$ . Offenbar gilt  $\mathbf{X}\delta^{(j)} \in \mathcal{V}$ , und aus der Invarianz folgt  $\mathbf{A}\mathbf{X}\delta^{(j)} \in \mathcal{V} = \text{Bild}(\mathbf{X})$ . Also existiert ein Vektor  $\mathbf{z}^{(j)} \in \mathbb{K}^p$  mit

$$\mathbf{A}\mathbf{X}\delta^{(j)} = \mathbf{X}\mathbf{z}^{(j)}.$$

Wir definieren die Matrix

$$\mathbf{\Lambda} := (\mathbf{z}^{(1)} \quad \dots \quad \mathbf{z}^{(p)}),$$

so dass gerade

$$\mathbf{\Lambda}\delta^{(j)} = \mathbf{z}^{(j)} \quad \text{für alle } j \in [1 : p]$$

### 3.5 Selbstadjungierte und unitäre Matrizen

gilt. Nun folgt

$$\mathbf{A}\mathbf{X}\delta^{(j)} = \mathbf{X}\mathbf{z}^{(j)} = \mathbf{X}\mathbf{\Lambda}\delta^{(j)} \quad \text{für alle } j \in [1 : p],$$

so dass unmittelbar (3.13) folgt.

Sei nun  $\mathbf{X}$  injektiv, und sei  $\mathbf{\Lambda} \in \mathbb{K}^{p \times p}$  eine Matrix, die (3.13) erfüllt. Nach Lemma 3.23 ist  $\mathbf{X}^*\mathbf{X}$  invertierbar, so dass aus

$$\begin{aligned} \mathbf{A}\mathbf{X} &= \mathbf{X}\mathbf{\Lambda}, \\ \mathbf{X}^*\mathbf{A}\mathbf{X} &= \mathbf{X}^*\mathbf{X}\mathbf{\Lambda} \end{aligned}$$

bereits

$$\mathbf{\Lambda} = (\mathbf{X}^*\mathbf{X})^{-1}\mathbf{X}^*\mathbf{A}\mathbf{X}$$

folgt. Durch diese Gleichung ist  $\mathbf{\Lambda}$  eindeutig festgelegt. ■

**Bemerkung 3.28** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Falls Matrizen  $\mathbf{X} \in \mathbb{K}^{n \times p}$  und  $\mathbf{\Lambda} \in \mathbb{K}^{p \times p}$  so gegeben sind, dass (3.13) gilt, muss  $\mathcal{V} := \text{Bild}(\mathbf{X})$  bereits ein invarianter Teilraum sein: Für jedes  $\mathbf{x} \in \text{Bild}(\mathbf{X})$  existiert ein Urbild  $\mathbf{y} \in \mathbb{K}^p$  mit  $\mathbf{x} = \mathbf{X}\mathbf{y}$ , so dass

$$\mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{X}\mathbf{y} = \mathbf{X}\mathbf{\Lambda}\mathbf{y} \in \text{Bild}(\mathbf{X})$$

unmittelbar folgt.

**Bemerkung 3.29** Seien  $\mathbf{A}, \widehat{\mathbf{A}} \in \mathbb{K}^{n \times n}$  ähnliche Matrizen mit  $\mathbf{B}\mathbf{A}\mathbf{B}^{-1} = \widehat{\mathbf{A}}$  für eine reguläre Matrix  $\mathbf{B}$ . Dann ist für jeden bezüglich  $\mathbf{A}$  invarianten Unterraum  $\mathcal{V} \subseteq \mathbb{K}^n$  die Menge

$$\widehat{\mathcal{V}} := \mathbf{B}^{-1}\mathcal{V}$$

ein bezüglich  $\widehat{\mathbf{A}}$  invarianter Unterraum.

**Bemerkung 3.30** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ , und seien Matrizen  $\mathbf{X} \in \mathbb{K}^{n \times p}$  und  $\mathbf{\Lambda} \in \mathbb{K}^{p \times p}$  so gegeben sind, dass (3.13) gilt. Falls  $\mathbf{X}$  injektiv ist, ist jeder Eigenwert der Matrix  $\mathbf{\Lambda}$  auch ein Eigenwert der Matrix  $\mathbf{A}$ : Für jedes  $\lambda \in \sigma(\mathbf{\Lambda})$  existiert ein Eigenvektor  $\mathbf{y} \in \mathbb{K}^p \setminus \{\mathbf{0}\}$ . Da  $\mathbf{X}$  injektiv ist, ist  $\mathbf{x} := \mathbf{X}\mathbf{y}$  nicht der Nullvektor. Es folgt

$$\mathbf{A}\mathbf{x} = \mathbf{A}\mathbf{X}\mathbf{y} = \mathbf{X}\mathbf{\Lambda}\mathbf{y} = \mathbf{X}\lambda\mathbf{y} = \lambda\mathbf{X}\mathbf{y} = \lambda\mathbf{x}.$$

## 3.5 Selbstadjungierte und unitäre Matrizen

Einen invarianten Unterraum können wir durch eine beliebige Basis darstellen. Aus der Perspektive der Numerik ist es ratsam, eine *Orthonormalbasis* zu verwenden, denn derartige Basen zeichnen sich durch eine besondere Unempfindlichkeit gegenüber Rundungsfehlern aus.

### 3 Theoretische Grundlagen

**Definition 3.31 (Isometrisch und unitär)** Sei  $\mathbf{Q} \in \mathbb{K}^{n \times m}$ . Falls die Gleichung

$$\mathbf{Q}^* \mathbf{Q} = \mathbf{I}$$

gilt, nennen wir  $\mathbf{Q}$  isometrisch.

Falls  $\mathbf{Q}$  isometrisch und quadratisch ist, falls also auch noch  $n = m$  gilt, nennen wir  $\mathbf{Q}$  unitär.

Um die Bezeichnung „isometrische Matrix“ zu rechtfertigen müssen wir auf eine auch noch für andere Zwecke nützliche Klasse von Matrizen zurückgreifen, nämlich auf die *selbstadjungierten Matrizen*.

**Definition 3.32 (Selbstadjungierte Matrix)** Eine Matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  heißt selbstadjungiert, falls  $\mathbf{A} = \mathbf{A}^*$  gilt.

Neben vielen anderen vorteilhaften Eigenschaften bieten selbstadjungierte Matrizen den Vorteil, dass sich viele wichtige Eigenschaften bereits an dem Skalarprodukt  $\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle$  ablesen lassen. Für unseren Zweck genügt zunächst die folgende Identitätsaussage:

**Lemma 3.33 (Identität selbstadjungierter Matrizen)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix. Falls

$$\langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle = 0 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n \quad (3.14)$$

gilt, folgt bereits  $\mathbf{A} = \mathbf{0}$ .

*Beweis.* Wir nehmen an, dass (3.14) gilt. Seien  $\mathbf{x}, \mathbf{y} \in \mathbb{K}^n$ . Nach Voraussetzung folgt mit Lemma 3.17 die Gleichung

$$\begin{aligned} 0 &= \langle \mathbf{A}(\mathbf{x} + \mathbf{y}), \mathbf{x} + \mathbf{y} \rangle = \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle + \langle \mathbf{A}\mathbf{y}, \mathbf{x} \rangle + \langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle + \langle \mathbf{A}\mathbf{y}, \mathbf{y} \rangle \\ &= \langle \mathbf{A}\mathbf{y}, \mathbf{x} \rangle + \langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{A}^* \mathbf{x} \rangle + \langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{A}\mathbf{x} \rangle + \langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle. \end{aligned}$$

Nun können wir  $\mathbf{y} = \mathbf{A}\mathbf{x}$  setzen und erhalten  $0 = 2\|\mathbf{A}\mathbf{x}\|^2$ , also  $\mathbf{A}\mathbf{x} = \mathbf{0}$  für jeden beliebigen Vektor  $\mathbf{x} \in \mathbb{K}^n$ , und damit insbesondere  $\mathbf{A} = \mathbf{0}$ . ■

Mit diesem Hilfsmittel können wir nun eine alternative Charakterisierung isometrischer Matrizen angeben:

**Lemma 3.34 (Isometrische Matrix)** Sei  $\mathbf{Q} \in \mathbb{K}^{n \times m}$ .  $\mathbf{Q}$  ist genau dann isometrisch, wenn

$$\|\mathbf{Q}\mathbf{x}\| = \|\mathbf{x}\| \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n \quad (3.15)$$

gilt, wenn die Multiplikation mit  $\mathbf{Q}$  also die Norm unverändert lässt.

*Beweis.* „ $\Rightarrow$ “: Sei zunächst  $\mathbf{Q}$  isometrisch. Nach Lemma 3.17 folgt

$$\|\mathbf{Q}\mathbf{x}\|^2 = \langle \mathbf{Q}\mathbf{x}, \mathbf{Q}\mathbf{x} \rangle = \langle \mathbf{Q}^* \mathbf{Q}\mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{x}, \mathbf{x} \rangle = \|\mathbf{x}\|^2 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n,$$

also gerade (3.15).

„ $\Leftarrow$ “: Gelte nun umgekehrt (3.15). Es folgt

$$\begin{aligned} 0 &= \|\mathbf{Q}\mathbf{x}\|^2 - \|\mathbf{x}\|^2 = \langle \mathbf{Q}\mathbf{x}, \mathbf{Q}\mathbf{x} \rangle - \langle \mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{Q}^* \mathbf{Q}\mathbf{x}, \mathbf{x} \rangle - \langle \mathbf{x}, \mathbf{x} \rangle \\ &= \langle (\mathbf{Q}^* \mathbf{Q} - \mathbf{I})\mathbf{x}, \mathbf{x} \rangle \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n. \end{aligned}$$

Da die Matrix  $\mathbf{Q}^* \mathbf{Q} - \mathbf{I}$  selbstadjungiert ist, erhalten wir mit Lemma 3.33 bereits die Gleichung  $\mathbf{Q}^* \mathbf{Q} - \mathbf{I} = \mathbf{0}$ , also muss  $\mathbf{Q}$  isometrisch sein. ■

In Hinblick auf die für die Behandlung von Eigenwertproblemen sehr wichtigen Ähnlichkeitstransformationen bieten unitäre Matrizen den Vorteil, dass sich ihre Inversen besonders einfach berechnen lassen.

**Lemma 3.35 (Unitäre Matrix)** Sei  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  unitär.

Dann gilt  $\mathbf{Q}\mathbf{Q}^* = \mathbf{I}$ , also  $\mathbf{Q}^* = \mathbf{Q}^{-1}$ , die Adjungierte ist die Inverse der Matrix  $\mathbf{Q}$ ,

*Beweis.* Da die Matrix  $\mathbf{Q}$  nach (3.15) insbesondere injektiv ist, folgt mit der Dimensionsformel, dass sie als quadratische Matrix auch surjektiv sein muss, es gilt also  $\text{Bild}(\mathbf{Q}) = \mathbb{K}^n$ .

Sei  $\mathbf{x} \in \mathbb{K}^n$ . Dank der Surjektivität gilt  $\mathbf{x} \in \text{Bild}(\mathbf{Q})$ , also finden wir ein Urbild  $\mathbf{y} \in \mathbb{K}^n$  mit  $\mathbf{x} = \mathbf{Q}\mathbf{y}$ . Mit Definition 3.31 erhalten wir

$$\mathbf{x} = \mathbf{Q}\mathbf{y} = \mathbf{Q}(\mathbf{Q}^* \mathbf{Q})\mathbf{y} = \mathbf{Q}\mathbf{Q}^* \mathbf{Q}\mathbf{y} = \mathbf{Q}\mathbf{Q}^* \mathbf{x}.$$

Da  $\mathbf{x}$  beliebig gewählt wurde, folgt daraus bereits  $\mathbf{I} = \mathbf{Q}\mathbf{Q}^*$ . Da wegen Definition 3.31 auch  $\mathbf{I} = \mathbf{Q}^* \mathbf{Q}$  gilt, muss  $\mathbf{Q}^*$  die Inverse der Matrix  $\mathbf{Q}$  sein. ■

## 3.6 Schur-Zerlegung

Nun können wir daran gehen, nach einer Möglichkeit zu suchen, um eine Matrix durch unitäre Ähnlichkeitstransformationen auf obere Dreiecksgestalt zu bringen.

Wir werden dabei induktiv vorgehen und zunächst versuchen, die erste Spalte der Matrix auf ein Vielfaches des ersten kanonischen Einheitsvektors abzubilden. Für derartige Aufgaben bietet die numerische lineare Algebra ein nützliches Hilfsmittel: *Householder-Spiegelungen* sind unitäre Abbildungen, die einen beliebigen Vektor in den Aufspann eines beliebigen anderen (von null verschiedenen) Vektors, beispielsweise eines kanonischen Einheitsvektors, überführen.

**Erinnerung 3.36 (Householder-Spiegelung)** Zu jedem Vektor  $\mathbf{v} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  ist die Householder-Spiegelung

$$\mathbf{Q}_v := \mathbf{I} - 2 \frac{\mathbf{v}\mathbf{v}^*}{\mathbf{v}^* \mathbf{v}}$$

### 3 Theoretische Grundlagen

eine unitäre Matrix.

Für beliebige Vektoren  $\mathbf{x} \in \mathbb{K}^n$  und  $\mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  können wir einen Householder-Vektor  $\mathbf{v} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  so finden, dass

$$\mathbf{Q}_v \mathbf{x} = \alpha \mathbf{y}$$

mit einem  $\alpha \in \mathbb{K}$  gilt, dass also  $\mathbf{x}$  in den von  $\mathbf{y}$  aufgespannten Raum abgebildet wird.

**Übungsaufgabe 3.37 (Allgemeine Householder-Spiegelung)** Für  $\mathbf{v} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  und  $\tau \in \mathbb{K}$  können wir die verallgemeinerte Householder-Spiegelung

$$\mathbf{Q}_{v,\tau} := \mathbf{I} - \tau \mathbf{v} \mathbf{v}^*$$

definieren.

(a) Beweisen Sie, dass  $\mathbf{Q}_{v,\tau}$  genau dann unitär ist, wenn  $|\tau|^2 \|\mathbf{v}\|^2 = 2 \operatorname{Re} \tau$  gilt.

(b) Seien  $\mathbf{a} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  und  $\alpha \in \mathbb{K}$  mit  $|\alpha| = \|\mathbf{a}\|$  gegeben.

Wir setzen  $\mathbf{v} := \mathbf{a} - \alpha \delta^{(1)}$  und nehmen an, dass  $\mathbf{v} \neq \mathbf{0}$  gilt.

Beweisen Sie, dass  $\mathbf{Q}_{v,\tau} \mathbf{a} = \alpha \delta^{(1)}$  genau dann gilt, wenn  $\tau \langle \mathbf{v}, \mathbf{a} \rangle = 1$ .

(c) Beweisen Sie unter den Voraussetzungen des Teils (b), dass  $\tau = \frac{1}{\langle \mathbf{v}, \mathbf{a} \rangle}$  die Gleichung  $|\tau|^2 \|\mathbf{v}\|^2 = 2 \operatorname{Re} \tau$  erfüllt, wir also eine unitäre Matrix konstruieren können, die  $\mathbf{a}$  auf  $\alpha \delta^{(1)}$  abbildet.

(d) Sei  $\mathbf{a} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$ . Beweisen Sie, dass die Wahl  $\alpha := -\operatorname{sgn}(\operatorname{Re} a_1) \|\mathbf{a}\|$  die Voraussetzungen der Teile (b) und (c) erfüllt.

Wir finden also verallgemeinerte Householder-Spiegelungen, die einen Vektor auf ein reelles Vielfaches des ersten kanonischen Einheitsvektors abbilden.

Damit können wir beispielsweise dafür sorgen, dass bei komplexen QR-Zerlegungen die Diagonale immer reell ist.

Hinweis: Für  $z \in \mathbb{C} \setminus \{0\}$  gilt  $\operatorname{Re}(1/z) = \operatorname{Re}(z)/|z|^2$ .

**Erinnerung 3.38 (QR-Zerlegung)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times m}$ . Dann existieren eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  und eine obere Dreiecksmatrix  $\mathbf{R} \in \mathbb{K}^{n \times m}$  mit  $\mathbf{A} = \mathbf{QR}$ .

Mit Hilfe der QR-Zerlegung können wir eine isometrische Matrix zu einer unitären Matrix ergänzen.

**Lemma 3.39 (Basisergänzung)** Sei  $\widehat{\mathbf{Q}} \in \mathbb{K}^{n \times p}$  eine isometrische Matrix. Dann existiert eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  mit  $\mathbf{Q}|_{n \times p} = \widehat{\mathbf{Q}}$ .

*Beweis.* Nach Erinnerung 3.38 existiert eine QR-Zerlegung  $\mathbf{Q}_0 \mathbf{R}_0 = \widehat{\mathbf{Q}}$  der Matrix  $\widehat{\mathbf{Q}}$ .

Nach Lemma 3.34 ist  $\widehat{\mathbf{Q}}$  injektiv, also muss  $p \leq n$  gelten. Da  $\mathbf{R}_0$  eine obere Dreiecksmatrix mit mehr Zeilen als Spalten ist, existiert eine Matrix  $\mathbf{R} \in \mathbb{K}^{p \times p}$  mit

$$\mathbf{R}_0 = \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix}.$$



Es gilt

$$\mathbf{R}^* \mathbf{R} = \mathbf{R}_0^* \mathbf{R}_0 = \mathbf{R}_0^* \mathbf{Q}_0^* \mathbf{Q}_0 \mathbf{R}_0 = \widehat{\mathbf{Q}}^* \widehat{\mathbf{Q}} = \mathbf{I},$$

also ist  $\mathbf{R}$  eine unitäre Dreiecksmatrix. Wir definieren

$$\mathbf{Q} := \mathbf{Q}_0 \begin{pmatrix} \mathbf{R} & \\ & \mathbf{I} \end{pmatrix}$$

und stellen fest, dass

$$\mathbf{Q}|_{n \times p} = \mathbf{Q}_0 \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix} = \mathbf{Q}_0 \mathbf{R}_0 = \widehat{\mathbf{Q}}$$

gilt. Also haben wir die gewünschte Fortsetzung gefunden.  $\blacksquare$

Tatsächlich können wir sogar beweisen, dass sich die ersten  $p$  Spalten der Matrix  $\mathbf{Q}_0$  lediglich durch ihre Vorzeichen von den ersten Spalten der Matrix  $\widehat{\mathbf{Q}}$  unterscheiden: Man kann sich überlegen, dass jede unitäre Dreiecksmatrix bereits eine Diagonalmatrix sein muss, und dieses Resultat lässt sich auf die Matrix  $\mathbf{R}$  anwenden.

**Lemma 3.40 (Deflation)** *Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ , sei  $p \in \mathbb{N}$  und sei  $\mathbf{X} \in \mathbb{K}^{n \times p}$  eine isometrische Matrix, die  $\mathbf{A}\mathbf{X} = \mathbf{X}\mathbf{\Lambda}$  (siehe (3.13)) mit einer geeigneten Matrix  $\mathbf{\Lambda} \in \mathbb{K}^{p \times p}$  erfüllt. Dann existieren eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  und weitere Matrizen  $\mathbf{D} \in \mathbb{K}^{(n-p) \times (n-p)}$ ,  $\mathbf{R} \in \mathbb{K}^{p \times (n-p)}$  derart, dass die Gleichung*

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \begin{pmatrix} \mathbf{\Lambda} & \mathbf{R} \\ \mathbf{0} & \mathbf{D} \end{pmatrix} \quad (3.16)$$

erfüllt ist und  $\mathbf{Q}|_{n \times p} = \mathbf{X}$  gilt.

*Beweis.* Gemäß Lemma 3.39 können wir  $\mathbf{X}$  zu einer unitären Matrix  $\mathbf{Q}$  mit  $\mathbf{Q}|_{n \times p} = \mathbf{X}$  ergänzen.

Wir zerlegen  $\mathbf{Q}$  dementsprechend in

$$\mathbf{Q} = \begin{pmatrix} \mathbf{X} & \mathbf{Y} \end{pmatrix}$$

und stellen fest, dass auch  $\mathbf{Y} \in \mathbb{K}^{n \times (n-p)}$  isometrisch ist und aus

$$\begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} = \mathbf{I} = \mathbf{Q}^* \mathbf{Q} = \begin{pmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{pmatrix} \begin{pmatrix} \mathbf{X} & \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^* \mathbf{X} & \mathbf{X}^* \mathbf{Y} \\ \mathbf{Y}^* \mathbf{X} & \mathbf{Y}^* \mathbf{Y} \end{pmatrix}$$

die Identität  $\mathbf{Y}^* \mathbf{X} = \mathbf{0}$  folgt.

Aus der Gleichung (3.13) erhalten wir

$$\begin{aligned} \mathbf{Q}^* \mathbf{A} \mathbf{Q} &= \begin{pmatrix} \mathbf{X}^* \\ \mathbf{Y}^* \end{pmatrix} \mathbf{A} \begin{pmatrix} \mathbf{X} & \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathbf{X}^* \mathbf{A} \mathbf{X} & \mathbf{X}^* \mathbf{A} \mathbf{Y} \\ \mathbf{Y}^* \mathbf{A} \mathbf{X} & \mathbf{Y}^* \mathbf{A} \mathbf{Y} \end{pmatrix} \\ &= \begin{pmatrix} \mathbf{X}^* \mathbf{X} \mathbf{\Lambda} & \mathbf{X}^* \mathbf{A} \mathbf{Y} \\ \mathbf{Y}^* \mathbf{X} \mathbf{\Lambda} & \mathbf{Y}^* \mathbf{A} \mathbf{Y} \end{pmatrix} = \begin{pmatrix} \mathbf{\Lambda} & \mathbf{X}^* \mathbf{A} \mathbf{Y} \\ \mathbf{0} & \mathbf{Y}^* \mathbf{A} \mathbf{Y} \end{pmatrix}, \end{aligned}$$

### 3 Theoretische Grundlagen

also können wir den Beweis abschließen, indem wir

$$\mathbf{R} := \mathbf{X}^* \mathbf{A} \mathbf{Y}, \quad \mathbf{D} := \mathbf{Y}^* \mathbf{A} \mathbf{Y}$$

setzen. ■

Da nach Satz 3.9 eine komplexe Matrix mindestens einen Eigenwert besitzt, ist der zugehörige Eigenraum nach Beispiel 3.25 ein invarianter Unterraum, der mittels Lemma 3.40 Anlass zu einer vereinfachenden Ähnlichkeitstransformation gibt. Bei wiederholter Anwendung erhält man eine sogenannte *Schur-Zerlegung* der Matrix:

**Satz 3.41 (Schur-Zerlegung)** *Sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Dann existieren eine unitäre Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und eine obere Dreiecksmatrix  $\mathbf{R} \in \mathbb{C}^{n \times n}$  so, dass die Gleichung*

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{R} \tag{3.17}$$

*erfüllt ist.*

*Beweis.* Per Induktion über die Dimension  $n$ . Der Induktionsanfang  $n = 1$  ist trivial.

Wir nehmen an, dass  $n > 1$  gilt und die Aussage für  $n - 1$  bewiesen ist. Nach Satz 3.9 besitzt  $\mathbf{A}$  mindestens einen Eigenwert  $\lambda \in \mathbb{C}$ . Sei  $\mathbf{x} \in \mathbb{C}^n \setminus \{\mathbf{0}\}$  ein zugehöriger Eigenvektor, den wir auf  $\|\mathbf{x}\| = 1$  normieren. Indem wir  $\mathbf{x}$  als Matrix  $\mathbf{X} \in \mathbb{C}^{n \times 1}$  und  $\lambda$  als Matrix  $\mathbf{\Lambda} \in \mathbb{C}^{1 \times 1}$  interpretieren, können wir Lemma 3.40 anwenden, um eine unitäre Matrix  $\mathbf{Q}_1 \in \mathbb{C}^{n \times n}$  und Matrizen  $\mathbf{R}_1 \in \mathbb{C}^{1 \times n}$  und  $\mathbf{A}_1 \in \mathbb{C}^{(n-1) \times (n-1)}$  zu erhalten, für die

$$\mathbf{Q}_1^* \mathbf{A} \mathbf{Q}_1 = \begin{pmatrix} \lambda & \mathbf{R}_1 \\ \mathbf{0} & \mathbf{A}_1 \end{pmatrix}$$

gilt. Wir wenden die Induktionsvoraussetzung auf  $\mathbf{A}_1$  an und erhalten eine orthonormale Matrix  $\mathbf{Q}_2 \in \mathbb{C}^{(n-1) \times (n-1)}$  und eine obere Dreiecksmatrix  $\mathbf{R}_2 \in \mathbb{C}^{(n-1) \times (n-1)}$  mit

$$\mathbf{Q}_2^* \mathbf{A}_1 \mathbf{Q}_2 = \mathbf{R}_2.$$

Nun können wir

$$\mathbf{Q} := \mathbf{Q}_1 \begin{pmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2 \end{pmatrix}, \quad \mathbf{R} := \begin{pmatrix} \lambda & \mathbf{R}_1 \mathbf{Q}_2 \\ \mathbf{0} & \mathbf{R}_2 \end{pmatrix}$$

definieren und erhalten

$$\begin{aligned} \mathbf{Q}^* \mathbf{A} \mathbf{Q} &= \begin{pmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2^* \end{pmatrix} \mathbf{Q}_1^* \mathbf{A} \mathbf{Q}_1 \begin{pmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2 \end{pmatrix} = \begin{pmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2^* \end{pmatrix} \begin{pmatrix} \lambda & \mathbf{R}_1 \\ \mathbf{0} & \mathbf{A}_1 \end{pmatrix} \begin{pmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{Q}_2 \end{pmatrix} \\ &= \begin{pmatrix} \lambda & \mathbf{R}_1 \mathbf{Q}_2 \\ \mathbf{0} & \mathbf{Q}_2^* \mathbf{A}_1 \mathbf{Q}_2 \end{pmatrix} = \begin{pmatrix} \lambda & \mathbf{R}_1 \mathbf{Q}_2 \\ \mathbf{0} & \mathbf{R}_2 \end{pmatrix} = \mathbf{R}. \end{aligned}$$

Da  $\mathbf{R}_2$  eine obere Dreiecksmatrix ist, muss auch  $\mathbf{R}$  eine obere Dreiecksmatrix sein, also ist der Beweis vollständig. ■

**Bemerkung 3.42** Im Beweis von Satz 3.41 können wir offenbar die Reihenfolge der Eigenwerte beliebig wählen und so beispielsweise auf der Diagonalen von  $\mathbf{R}$  jede Anordnung erreichen.

Die Gleichung  $\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{R}$  ist äquivalent zu  $\mathbf{A} \mathbf{Q} = \mathbf{Q} \mathbf{R}$ , und indem wir diese Gleichung auf die ersten  $p$  Spalten einschränken und die Dreiecksstruktur von  $\mathbf{R}$  ausnutzen, folgt sofort, dass der Aufspann der ersten  $p$  Spalten von  $\mathbf{Q}$  jeweils ein invarianter Teilraum von  $\mathbf{A}$  sein muss.

**Übungsaufgabe 3.43 (Reelle Schur-Zerlegung)** Sei  $\mathbf{A} \in \mathbb{R}^{n \times n}$  eine reelle Matrix.

- (a) Wenn wir  $\mathbf{A}$  als komplexe Matrix interpretieren, besitzt sie mindestens einen Eigenwert  $\lambda \in \mathbb{C}$  mit einem Eigenvektor  $\mathbf{e} \in \mathbb{C}^n \setminus \{\mathbf{0}\}$ .

Beweisen Sie, dass dann  $\bar{\lambda}$  ein Eigenwert mit dem Eigenvektor  $\bar{\mathbf{e}}$  ist.

- (b) Beweisen Sie, dass im Fall  $\lambda \neq \bar{\lambda}$  eine reelle injektive Matrix  $\mathbf{X} \in \mathbb{R}^{n \times 2}$  und eine reelle Matrix  $\mathbf{\Lambda} \in \mathbb{R}^{2 \times 2}$  mit  $\mathbf{A} \mathbf{X} = \mathbf{X} \mathbf{\Lambda}$  existieren.

- (c) Beweisen Sie, dass eine reelle unitäre Matrix  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  und  $m \in \mathbb{N}$ ,  $n_1, \dots, n_m \in \{1, 2\}$  sowie reelle Matrizen  $\mathbf{R}_{ij} \in \mathbb{R}^{n_i \times n_j}$  mit

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \begin{pmatrix} \mathbf{R}_{11} & \cdots & \mathbf{R}_{1m} \\ & \ddots & \vdots \\ & & \mathbf{R}_{mm} \end{pmatrix}$$

existieren. Diese Zerlegung nennt man reelle Schur-Zerlegung der Matrix  $\mathbf{A}$ .

## 3.7 Diagonalisierbarkeit durch unitäre Transformationen

Die Matrix auf obere Dreiecksgestalt bringen zu können hilft uns zwar bei der Bestimmung der Eigenwerte, allerdings nur sehr bedingt bei der Untersuchung der Eigenvektoren. Insbesondere wird für viele Beweise eine *Orthonormalbasis* von Eigenvektoren benötigt, die wir Satz 3.41 nicht unmittelbar entnehmen können.

Für reelle Eigenwerte lässt sich relativ leicht klären, welche Eigenschaften eine Matrix aufweisen muss, um sich unitär diagonalisieren zu lassen: Wenn eine reelle Diagonalmatrix  $\mathbf{D} \in \mathbb{R}^{n \times n}$  und einer unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  mit

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{D}$$

existieren, folgt mit Lemma 3.35

$$\mathbf{A} = \mathbf{Q} \mathbf{D} \mathbf{Q}^*$$

und damit bereits

$$\mathbf{A}^* = \mathbf{Q}^{**} \mathbf{D}^* \mathbf{Q}^* = \mathbf{Q} \mathbf{D} \mathbf{Q}^* = \mathbf{A}, \quad (3.18)$$

also muss  $\mathbf{A}$  mindestens selbstadjungiert sein.

### 3 Theoretische Grundlagen

Wir können beweisen, dass diese Eigenschaft nicht nur notwendig, sondern auch hinreichend ist. Im Prinzip könnte man dazu auf Satz 3.41 zurückgreifen, allerdings würden wir damit nur die Existenz einer *komplexen* unitären Ähnlichkeitstransformation erhalten. Indem wir den Fundamentalsatz der Algebra durch eine auf selbstadjungierte Matrizen zugeschnittene Aussage ersetzen, können wir beweisen, dass für reelle selbstadjungierte Matrizen auch eine reelle Ähnlichkeitstransformation genügt.

Der Ausgangspunkt unserer Betrachtung ist der *Rayleigh-Quotient*: Falls  $\mathbf{e} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  ein Eigenvektor zu einem Eigenwert  $\lambda \in \mathbb{K}$  einer Matrix  $\mathbf{A}$  ist, können wir den Eigenwert dank (3.6b) aus dem Eigenvektor rekonstruieren:

$$\begin{aligned}\lambda \mathbf{e} &= \mathbf{A} \mathbf{e}, \\ \lambda \langle \mathbf{e}, \mathbf{e} \rangle &= \langle \mathbf{e}, \lambda \mathbf{e} \rangle = \langle \mathbf{e}, \mathbf{A} \mathbf{e} \rangle, \\ \lambda &= \frac{\langle \mathbf{e}, \mathbf{A} \mathbf{e} \rangle}{\langle \mathbf{e}, \mathbf{e} \rangle}.\end{aligned}$$

Den Ausdruck auf der rechten Seite bezeichnet man als einen *Rayleigh-Quotienten*.

**Definition 3.44 (Rayleigh-Quotient)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Die Rayleigh-Quotientenabbildung ist gegeben durch

$$\Lambda_A: \mathbb{K}^n \setminus \{\mathbf{0}\} \rightarrow \mathbb{K}, \quad \mathbf{x} \mapsto \frac{\langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle}.$$

Falls  $\mathbf{A}$  selbstadjungiert ist, gilt wegen Lemma 3.17 und (3.6c) für alle  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  die Gleichung

$$\langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle = \langle \mathbf{A}^* \mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{A} \mathbf{x}, \mathbf{x} \rangle = \overline{\langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle},$$

also folgt  $\langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle \in \mathbb{R}$ . Der Rayleigh-Quotient ist für selbstadjungierte Matrizen also immer reellwertig.

**Satz 3.45 (Courant-Fischer)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  selbstadjungiert. Die Rayleigh-Quotientenabbildung besitzt ein Maximum, dieses Maximum ist gerade der größte Eigenwert der Matrix  $\mathbf{A}$ , und jedes seiner Urbilder ist ein Eigenvektor zu diesem Eigenwert.

*Beweis.* Die Rayleigh-Quotientenabbildung ist invariant unter Skalierung, denn für alle  $\alpha \in \mathbb{K} \setminus \{0\}$  gilt wegen (3.6a) und (3.6b)

$$\Lambda_A(\alpha \mathbf{x}) = \frac{\langle \alpha \mathbf{x}, \alpha \mathbf{A} \mathbf{x} \rangle}{\langle \alpha \mathbf{x}, \alpha \mathbf{x} \rangle} = \frac{|\alpha|^2 \langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle}{|\alpha|^2 \langle \mathbf{x}, \mathbf{x} \rangle} = \Lambda_A(\mathbf{x}), \quad (3.19)$$

also dürfen wir uns auf der Suche nach dem Maximum auf die Einheitskugel

$$\mathcal{S} := \{\mathbf{x} \in \mathbb{K}^n : \|\mathbf{x}\| = 1\}$$

beschränken. Diese Menge ist beschränkt und abgeschlossen, also nach dem Satz von Heine-Borel kompakt. Als stetige Abbildung nimmt  $\Lambda_A$  auf der kompakten Menge  $\mathcal{S}$  ein Maximum an, das wir mit  $\lambda \in \mathbb{R}$  bezeichnen. Wir fixieren  $\mathbf{e} \in \mathcal{S}$  mit  $\lambda = \Lambda_A(\mathbf{e})$ .

### 3.7 Diagonalisierbarkeit durch unitäre Transformationen

Seien nun  $\mathbf{y} \in \mathbb{K}^n$  und  $\alpha \in (0, 1/\|\mathbf{y}\|)$  (mit  $1/0 = \infty$ ) gegeben. Dann gilt

$$\|\mathbf{e} + \alpha\mathbf{y}\| \geq \|\mathbf{e}\| - |\alpha|\|\mathbf{y}\| > 1 - 1 = 0,$$

also  $\mathbf{e} + \alpha\mathbf{y} \neq \mathbf{0}$ . Da  $\lambda$  das Maximum der Rayleigh-Quotientenabbildung ist, gilt

$$\frac{\langle \mathbf{e} + \alpha\mathbf{y}, \mathbf{A}(\mathbf{e} + \alpha\mathbf{y}) \rangle}{\langle \mathbf{e} + \alpha\mathbf{y}, \mathbf{e} + \alpha\mathbf{y} \rangle} = \Lambda_A(\mathbf{e} + \alpha\mathbf{y}) \leq \lambda.$$

Mit (3.6) folgt

$$\begin{aligned} \lambda\|\mathbf{e}\|^2 + 2\lambda\alpha \operatorname{Re}\langle \mathbf{y}, \mathbf{e} \rangle + \lambda\alpha^2\|\mathbf{y}\|^2 &= \lambda\langle \mathbf{e}, \mathbf{e} \rangle + \lambda\alpha\langle \mathbf{y}, \mathbf{e} \rangle + \lambda\alpha\langle \mathbf{e}, \mathbf{y} \rangle + \lambda\alpha^2\langle \mathbf{y}, \mathbf{y} \rangle \\ &= \lambda(\langle \mathbf{e}, \mathbf{e} \rangle + \langle \alpha\mathbf{y}, \mathbf{e} \rangle + \langle \mathbf{e}, \alpha\mathbf{y} \rangle + \langle \alpha\mathbf{y}, \alpha\mathbf{y} \rangle) \\ &= \lambda\langle \mathbf{e} + \alpha\mathbf{y}, \mathbf{e} + \alpha\mathbf{y} \rangle \geq \langle \mathbf{e} + \alpha\mathbf{y}, \mathbf{A}(\mathbf{e} + \alpha\mathbf{y}) \rangle \\ &= \langle \mathbf{e}, \mathbf{A}\mathbf{e} \rangle + \langle \alpha\mathbf{y}, \mathbf{A}\mathbf{e} \rangle + \langle \mathbf{e}, \alpha\mathbf{A}\mathbf{y} \rangle + \langle \alpha\mathbf{y}, \alpha\mathbf{A}\mathbf{y} \rangle \\ &= \langle \mathbf{e}, \mathbf{A}\mathbf{e} \rangle + \alpha\langle \mathbf{y}, \mathbf{A}\mathbf{e} \rangle + \alpha\langle \mathbf{A}\mathbf{e}, \mathbf{y} \rangle + \alpha^2\langle \mathbf{y}, \mathbf{A}\mathbf{y} \rangle \\ &= \langle \mathbf{e}, \mathbf{A}\mathbf{e} \rangle + \alpha\langle \mathbf{y}, \mathbf{A}\mathbf{e} \rangle + \alpha\overline{\langle \mathbf{y}, \mathbf{A}\mathbf{e} \rangle} + \alpha^2\langle \mathbf{y}, \mathbf{A}\mathbf{y} \rangle \\ &= \lambda\|\mathbf{e}\|^2 + 2\alpha \operatorname{Re}\langle \mathbf{y}, \mathbf{A}\mathbf{e} \rangle + \alpha^2\langle \mathbf{y}, \mathbf{A}\mathbf{y} \rangle. \end{aligned}$$

Wir dürfen  $\lambda\|\mathbf{e}\|^2$  auf beiden Seiten streichen und durch  $\alpha$  dividieren, um

$$\begin{aligned} \alpha^2(\langle \mathbf{y}, \lambda\mathbf{y} \rangle - \langle \mathbf{y}, \mathbf{A}\mathbf{y} \rangle) &\geq 2\alpha(\operatorname{Re}\langle \mathbf{y}, \mathbf{A}\mathbf{e} \rangle - \operatorname{Re}\langle \mathbf{y}, \lambda\mathbf{e} \rangle), \\ \alpha\langle \mathbf{y}, \lambda\mathbf{y} - \mathbf{A}\mathbf{y} \rangle &\geq 2\operatorname{Re}\langle \mathbf{y}, \mathbf{A}\mathbf{e} - \lambda\mathbf{e} \rangle \end{aligned}$$

zu erhalten. Wir setzen  $\mathbf{y} := \mathbf{A}\mathbf{e} - \lambda\mathbf{e}$  und erhalten

$$\alpha\langle \mathbf{y}, \lambda\mathbf{y} - \mathbf{A}\mathbf{y} \rangle \geq 2\|\mathbf{A}\mathbf{e} - \lambda\mathbf{e}\|^2.$$

Da wir  $\alpha$  beliebig klein wählen dürfen, folgt  $0 \geq \|\mathbf{A}\mathbf{e} - \lambda\mathbf{e}\|^2$ , also muss  $\mathbf{A}\mathbf{e} = \lambda\mathbf{e}$  gelten. Damit ist  $\mathbf{e}$  ein Eigenvektor zu dem Eigenwert  $\lambda$ .

Sei nun  $\tilde{\lambda} \in \mathbb{K}$  ein beliebiger weiterer Eigenwert der Matrix  $\mathbf{A}$ , und sei  $\tilde{\mathbf{e}} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  ein zugehöriger Eigenvektor, der Einfachheit halber so skaliert, dass  $\|\tilde{\mathbf{e}}\| = 1$  gilt. Dann gilt wegen  $\tilde{\mathbf{e}} \in \mathcal{S}$  und der Wahl des Maximums  $\lambda$  die Ungleichung

$$\lambda \geq \Lambda_A(\tilde{\mathbf{e}}) = \frac{\langle \tilde{\mathbf{e}}, \mathbf{A}\tilde{\mathbf{e}} \rangle}{\langle \tilde{\mathbf{e}}, \tilde{\mathbf{e}} \rangle} = \frac{\langle \tilde{\mathbf{e}}, \tilde{\lambda}\tilde{\mathbf{e}} \rangle}{\langle \tilde{\mathbf{e}}, \tilde{\mathbf{e}} \rangle} = \tilde{\lambda} \frac{\langle \tilde{\mathbf{e}}, \tilde{\mathbf{e}} \rangle}{\langle \tilde{\mathbf{e}}, \tilde{\mathbf{e}} \rangle} = \tilde{\lambda},$$

also ist  $\lambda$  nicht nur das Maximum der Rayleigh-Quotientenabbildung, sondern tatsächlich der größte Eigenwert. ■

**Bemerkung 3.46 (Reelle Eigenvektoren)** *Im Fall  $\mathbb{K} = \mathbb{R}$  besagt der Satz 3.45 von Courant und Fischer, dass selbstadjungierte Matrizen  $\mathbf{A} \in \mathbb{R}^{n \times n}$  immer mindestens einen reellen Eigenwert haben, der dem Maximum des Rayleigh-Quotienten entspricht und zu dem ein reeller Eigenvektor  $\mathbf{e} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$  existiert.*

### 3 Theoretische Grundlagen

Mit Hilfe des Satzes 3.45 von Courant und Fischer können wir den Rest des Beweises des Satzes 3.41 ohne „Umweg über  $\mathbb{C}$ “ nachvollziehen.

**Satz 3.47 (Hauptachsentransformation)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  und eine reelle Diagonalmatrix  $\mathbf{D} \in \mathbb{R}^{n \times n}$  mit

$$\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^* \quad (3.20)$$

existieren genau dann, wenn  $\mathbf{A}$  selbstadjungiert ist.

*Beweis.* „ $\Rightarrow$ “: Es gelte (3.20). Aus (3.18) folgt dann  $\mathbf{A} = \mathbf{A}^*$ .

„ $\Leftarrow$ “: Sei  $\mathbf{A}$  selbstadjungiert. Wie schon bei der Schur-Zerlegung gehen wir induktiv vor. Der Fall  $n = 1$  ist trivial.

Sei  $n \in \mathbb{N}$  so gegeben, dass die Behauptung für Matrizen der Dimension  $n - 1$  gilt. Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix. Nach Satz 3.45 besitzt  $\mathbf{A}$  einen Eigenwert  $\lambda \in \mathbb{R}$  mit einem zugehörigen Eigenvektor  $\mathbf{e} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$ .

Sei  $\mathbf{Q}_1 \in \mathbb{K}^{n \times n}$  eine Householder-Matrix, die den ersten kanonischen Einheitsvektor  $\delta^{(1)}$  auf ein Vielfaches des Eigenvektors  $\mathbf{e}$  abbildet, es gelte also  $\mathbf{Q}_1 \delta^{(1)} = \alpha \mathbf{e}$  mit  $\alpha \in \mathbb{K}$ . Dann folgt

$$\mathbf{Q}_1^* \mathbf{A} \mathbf{Q}_1 \delta^{(1)} = \alpha \mathbf{Q}_1^* \mathbf{A} \mathbf{e} = \alpha \lambda \mathbf{Q}_1^* \mathbf{e} = \lambda \delta^{(1)}.$$

Da  $\mathbf{Q}_1^* \mathbf{A} \mathbf{Q}_1$  wieder eine selbstadjungierte Matrix ist, muss sie die Form

$$\mathbf{Q}_1^* \mathbf{A} \mathbf{Q}_1 = \begin{pmatrix} \lambda & \mathbf{0} \\ \mathbf{0} & \widehat{\mathbf{A}} \end{pmatrix}$$

mit einer selbstadjungierten Matrix  $\widehat{\mathbf{A}} \in \mathbb{R}^{(n-1) \times (n-1)}$  aufweisen.

Also können wir die Induktionsvoraussetzung auf  $\widehat{\mathbf{A}}$  anwenden und wie im Beweis von Satz 3.41 fortfahren. ■

**Übungsaufgabe 3.48 (Ableitung des Rayleigh-Quotienten)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix, und seien  $\mathbf{x}, \mathbf{p} \in \mathbb{K}^n$  Vektoren mit  $\|\mathbf{x}\| = 1$  und  $\|\mathbf{p}\| \leq 1$ .

Da  $\|\mathbf{x} + t\mathbf{p}\| \geq \|\mathbf{x}\| - |t|\|\mathbf{p}\| \geq 1 - |t| > 0$  für alle  $t \in (-1, 1)$  gilt, ist

$$f: (-1, 1) \rightarrow \mathbb{R}, \quad t \mapsto \Lambda_A(\mathbf{x} + t\mathbf{p}) = \frac{\langle \mathbf{x} + t\mathbf{p}, \mathbf{A}(\mathbf{x} + t\mathbf{p}) \rangle}{\langle \mathbf{x} + t\mathbf{p}, \mathbf{x} + t\mathbf{p} \rangle},$$

eine wohldefinierte stetig differenzierbare Abbildung.

(a) Beweisen Sie

$$f'(0) = 2 \operatorname{Re} \frac{\langle \mathbf{p}, \mathbf{A}\mathbf{x} - \Lambda_A(\mathbf{x})\mathbf{x} \rangle}{\|\mathbf{x}\|^2}.$$

(b) Beweisen Sie, dass alle lokalen Extremstellen des Rayleigh-Quotienten Eigenvektoren der Matrix  $\mathbf{A}$  sind.

### 3.7 Diagonalisierbarkeit durch unitäre Transformationen

Wir wissen, dass wir selbstadjungierte Matrizen unitär diagonalisieren können, allerdings benötigen wir dafür das Maximum des Rayleigh-Quotienten, dessen Berechnung in der Regel nur iterativ erfolgen kann.

Wir können allerdings mit Householder-Spiegelungen eine selbstadjungierte Matrix direkt auf Tridiagonalgestalt bringen, und mit den in Übungsaufgabe 3.37 eingeführten verallgemeinerten Householder-Spiegelungen können wir sogar sicherstellen, dass wir eine *reelle* Tridiagonalmatrix erhalten. Mit diesem Vorbereitungsschritt lässt sich die Berechnung der Eigenwerte erheblich vereinfachen.

**Übungsaufgabe 3.49 (Tridiagonalgestalt)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  selbstadjungiert.

- (a) Beweisen Sie, dass im Fall  $n > 1$  eine verallgemeinerte Householder-Spiegelung  $\mathbf{Q}_1 \in \mathbb{K}^{n \times n}$  (siehe Übungsaufgabe 3.37) und reelle Zahlen  $d, \ell \in \mathbb{R}$  sowie eine selbstadjungierte Matrix  $\hat{\mathbf{A}} \in \mathbb{K}^{(n-1) \times (n-1)}$  existieren mit

$$\mathbf{Q}_1^* \mathbf{A} \mathbf{Q}_1 = \begin{pmatrix} d & \ell & 0 & \cdots & 0 \\ \ell & \hat{a}_{11} & \hat{a}_{12} & \cdots & \hat{a}_{1,n-1} \\ 0 & \hat{a}_{21} & \ddots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & \hat{a}_{n-1,1} & \cdots & \cdots & \hat{a}_{n-1,n-1} \end{pmatrix}, \quad \mathbf{Q}_1 \delta^{(1)} = \delta^{(1)}.$$

- (b) Verwenden Sie (a), um zu beweisen, dass es eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  und reelle Zahlen  $d_1, \dots, d_n \in \mathbb{R}$ ,  $\ell_1, \dots, \ell_{n-1} \in \mathbb{R}$  gibt mit

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \begin{pmatrix} d_1 & \ell_1 & & & \\ \ell_1 & \ddots & \ddots & & \\ & \ddots & \ddots & \ell_{n-1} & \\ & & & \ell_{n-1} & d_n \end{pmatrix},$$

dass wir also jede selbstadjungierte Matrix mit einer unitären Ähnlichkeitstransformation in eine reelle Tridiagonalmatrix überführen können.

Nun wollen wir untersuchen, unter welchen Bedingungen sich Matrizen mit unitären Ähnlichkeitstransformationen in eine *komplexe* Diagonalgestalt überführen lassen.

Sei also  $\mathbf{A} \in \mathbb{C}^{n \times n}$  gegeben. Falls eine unitäre Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und eine Diagonalmatrix  $\mathbf{D} \in \mathbb{C}^{n \times n}$  mit

$$\mathbf{A} = \mathbf{Q}^* \mathbf{D} \mathbf{Q}$$

existieren, haben wir mit Lemma 3.35

$$\begin{aligned} \mathbf{A}^* \mathbf{A} &= (\mathbf{Q}^* \mathbf{D} \mathbf{Q})^* \mathbf{Q}^* \mathbf{D} \mathbf{Q} = \mathbf{Q}^* \mathbf{D}^* \mathbf{Q} \mathbf{Q}^* \mathbf{D} \mathbf{Q} = \mathbf{Q}^* \mathbf{D}^* \mathbf{D} \mathbf{Q} \\ &= \mathbf{Q}^* \mathbf{D} \mathbf{D}^* \mathbf{Q} = \mathbf{Q}^* \mathbf{D} \mathbf{Q} \mathbf{Q}^* \mathbf{D}^* \mathbf{Q} = \mathbf{A} \mathbf{A}^*, \end{aligned} \quad (3.21)$$

da für Diagonalmatrizen  $\mathbf{D}^* \mathbf{D} = \mathbf{D} \mathbf{D}^*$  gilt. Diese Gleichung muss also gelten, falls  $\mathbf{A}$  diagonalisierbar ist.

Wir werden nun beweisen, dass diese Eigenschaft nicht nur notwendig, sondern auch schon hinreichend ist.

### 3 Theoretische Grundlagen

**Definition 3.50 (Normale Matrizen)** Eine Matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  heißt normal, falls

$$\mathbf{A}^* \mathbf{A} = \mathbf{A} \mathbf{A}^*.$$

**Beispiel 3.51 (Drehstreckungen)** Seien  $c, s, r \in \mathbb{K}$  mit  $|c|^2 + |s|^2 = 1$  gegeben (beispielsweise  $c = \cos(\varphi)$  und  $s = \sin(\varphi)$  mit  $\varphi \in \mathbb{R}$ ).

Die zugehörige Drehstreckung

$$\mathbf{G} := r \begin{pmatrix} c & s \\ -\bar{s} & \bar{c} \end{pmatrix}$$

erfüllt

$$\mathbf{G}^* \mathbf{G} = |r|^2 \begin{pmatrix} \bar{c}c + \bar{s}s & \bar{c}s - s\bar{c} \\ \bar{s}c - c\bar{s} & \bar{s}s + c\bar{c} \end{pmatrix} = |r|^2 \mathbf{I}, \quad \mathbf{G} \mathbf{G}^* = |r|^2 \begin{pmatrix} c\bar{c} + s\bar{s} & -c\bar{s} + s\bar{c} \\ -sc + cs & s\bar{r} + c\bar{c} \end{pmatrix} = |r|^2 \mathbf{I},$$

ist also eine normale Matrix.

**Lemma 3.52 (Metrische Äquivalenz)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine Matrix.  $\mathbf{A}$  ist genau dann normal, wenn

$$\|\mathbf{A}\mathbf{x}\| = \|\mathbf{A}^*\mathbf{x}\| \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n. \quad (3.22)$$

*Beweis.* Zunächst nutzen wir Lemma 3.17 und erhalten für alle  $\mathbf{x} \in \mathbb{K}^n$  die Gleichungen

$$\|\mathbf{A}\mathbf{x}\|^2 = \langle \mathbf{A}\mathbf{x}, \mathbf{A}\mathbf{x} \rangle = \langle \mathbf{A}^* \mathbf{A}\mathbf{x}, \mathbf{x} \rangle, \quad \|\mathbf{A}^*\mathbf{x}\|^2 = \langle \mathbf{A}^*\mathbf{x}, \mathbf{A}^*\mathbf{x} \rangle = \langle \mathbf{A} \mathbf{A}^* \mathbf{x}, \mathbf{x} \rangle.$$

„ $\Rightarrow$ “: Sei  $\mathbf{A}$  normal. Dann folgt aus diesen Gleichungen bereits

$$\|\mathbf{A}\mathbf{x}\|^2 = \langle \mathbf{A}^* \mathbf{A}\mathbf{x}, \mathbf{x} \rangle = \langle \mathbf{A} \mathbf{A}^* \mathbf{x}, \mathbf{x} \rangle = \|\mathbf{A}^*\mathbf{x}\|^2$$

für alle  $\mathbf{x} \in \mathbb{K}^n$ .

„ $\Leftarrow$ “: Gelte nun die Gleichheit der Normen (3.22). Dann folgt

$$\langle (\mathbf{A}^* \mathbf{A} - \mathbf{A} \mathbf{A}^*) \mathbf{x}, \mathbf{x} \rangle = \|\mathbf{A}\mathbf{x}\|^2 - \|\mathbf{A}^*\mathbf{x}\|^2 = 0$$

für alle  $\mathbf{x} \in \mathbb{K}^n$ , und da  $\mathbf{A}^* \mathbf{A} - \mathbf{A} \mathbf{A}^*$  eine selbstadjungierte Matrix ist, können wir Lemma 3.33 anwenden, um zu folgern, dass sie gleich null sein muss. ■

**Lemma 3.53 (Normale Dreiecksmatrizen)** Sei  $\mathbf{R} \in \mathbb{K}^{n \times n}$  eine obere Dreiecksmatrix. Dann ist  $\mathbf{R}$  genau dann normal, wenn die Matrix diagonal ist.

*Beweis.* Falls  $\mathbf{R}$  diagonal ist, stehen auf der Diagonalen der Produkte  $\mathbf{R}\mathbf{R}^*$  und  $\mathbf{R}^*\mathbf{R}$  jeweils die Einträge  $|r_{ii}|^2$  für  $i \in [1 : n]$ , also ist  $\mathbf{R}$  auch normal.

Die Gegenrichtung beweisen wir per Induktion. Für  $n = 1$  ist die Aussage trivial.

Sei nun  $n \in \mathbb{N}$  so gegeben, dass jede normale obere Dreiecksmatrix  $\mathbf{R} \in \mathbb{K}^{n \times n}$  bereits diagonal ist.



### 3.7 Diagonalisierbarkeit durch unitäre Transformationen

Sei  $\mathbf{R} \in \mathbb{K}^{(n+1) \times (n+1)}$  eine normale obere Dreiecksmatrix. Mit Lemma 3.52 folgt

$$|r_{11}|^2 = \|\mathbf{R}\delta^{(1)}\|^2 = \|\mathbf{R}^*\delta^{(1)}\|^2 = \sum_{j=1}^{n+1} |r_{1j}|^2,$$

also  $|r_{1j}| = 0$  für alle  $j \in [2 : n+1]$ . Demnach ist  $\mathbf{R}$  von der Gestalt

$$\mathbf{R} = \begin{pmatrix} r_{11} & \mathbf{0} \\ & \widehat{\mathbf{R}} \end{pmatrix}, \quad \widehat{\mathbf{R}} \in \mathbb{K}^{n \times n},$$

und wir haben

$$\begin{pmatrix} |r_{11}|^2 & \\ & \widehat{\mathbf{R}}^*\widehat{\mathbf{R}} \end{pmatrix} = \mathbf{R}^*\mathbf{R} = \mathbf{R}\mathbf{R}^* = \begin{pmatrix} |r_{11}|^2 & \\ & \widehat{\mathbf{R}}\widehat{\mathbf{R}}^* \end{pmatrix},$$

also folgt insbesondere  $\widehat{\mathbf{R}}^*\widehat{\mathbf{R}} = \widehat{\mathbf{R}}\widehat{\mathbf{R}}^*$ ,  $\widehat{\mathbf{R}}$  ist normal. Wir können die Induktionsvoraussetzung anwenden und folgern, dass  $\widehat{\mathbf{R}}$  diagonal sein muss, und damit auch  $\mathbf{R}$ . ■

**Folgerung 3.54 (Komplexe Diagonalisierbarkeit)** Sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Eine unitäre Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und eine Diagonalmatrix  $\mathbf{D} \in \mathbb{C}^{n \times n}$  mit

$$\mathbf{Q}^*\mathbf{A}\mathbf{Q} = \mathbf{D} \tag{3.23}$$

existieren genau dann, wenn  $\mathbf{A}$  normal ist.

*Beweis.* „ $\Rightarrow$ “: Setze zunächst voraus, dass  $\mathbf{Q}, \mathbf{D}$  wie in (3.23) existieren. Dann folgt mit (3.21), dass  $\mathbf{A}$  normal sein muss.

„ $\Leftarrow$ “: Sei  $\mathbf{A}$  normal. Wir nutzen Satz 3.41, um eine unitäre Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und eine obere Dreiecksmatrix  $\mathbf{R} \in \mathbb{C}^{n \times n}$  mit  $\mathbf{R} = \mathbf{Q}^*\mathbf{A}\mathbf{Q}$  zu finden.

Da  $\mathbf{A}$  normal ist, gilt mit Lemma 3.18 und Lemma 3.35

$$\begin{aligned} \mathbf{R}^*\mathbf{R} &= \mathbf{Q}^*\mathbf{A}^*\mathbf{Q}\mathbf{Q}^*\mathbf{A}\mathbf{Q} = \mathbf{Q}^*\mathbf{A}^*\mathbf{A}\mathbf{Q} \\ &= \mathbf{Q}^*\mathbf{A}\mathbf{A}^*\mathbf{Q} = \mathbf{Q}^*\mathbf{A}\mathbf{Q}\mathbf{Q}^*\mathbf{A}^*\mathbf{Q} = \mathbf{R}\mathbf{R}^*, \end{aligned}$$

also ist  $\mathbf{R}$  normal. Nach Lemma 3.53 muss  $\mathbf{R}$  damit bereits eine Diagonalmatrix sein. ■

Der Name *Spektralnorm* legt nahe, dass diese Norm in enger Beziehung zu dem Spektrum stehen dürfte. Das folgende Lemma bestätigt diese Vermutung.

**Lemma 3.55 (Spektralradius)** Sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$  eine Matrix. Das Maximum der Beträge ihrer Eigenwerte

$$\varrho(\mathbf{A}) := \max\{|\lambda| : \lambda \in \sigma(\mathbf{A})\}$$

nennen wir ihren Spektralradius. Falls  $\mathbf{A}$  normal ist, gilt  $\varrho(\mathbf{A}) = \|\mathbf{A}\|$ .

### 3 Theoretische Grundlagen

*Beweis.* Sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$  eine normale Matrix.

Nach Folgerung 3.54 finden wir eine unitäre Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und eine Diagonalmatrix  $\mathbf{D} \in \mathbb{C}^{n \times n}$  mit  $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^*$ .

Mit Definition 3.19 und (3.15) sowie der Substitution  $\widehat{\mathbf{z}} = \mathbf{Q}^*\mathbf{z}$ ,  $\mathbf{Q}\widehat{\mathbf{z}} = \mathbf{z}$  erhalten wir

$$\begin{aligned} \|\mathbf{A}\| &= \max \left\{ \frac{\|\mathbf{A}\mathbf{z}\|}{\|\mathbf{z}\|} : \mathbf{z} \in \mathbb{K}^n \setminus \{\mathbf{0}\} \right\} = \max \left\{ \frac{\|\mathbf{Q}\mathbf{D}\mathbf{Q}^*\mathbf{z}\|}{\|\mathbf{z}\|} : \mathbf{z} \in \mathbb{K}^n \setminus \{\mathbf{0}\} \right\} \\ &= \max \left\{ \frac{\|\mathbf{D}\mathbf{Q}^*\mathbf{z}\|}{\|\mathbf{z}\|} : \mathbf{z} \in \mathbb{K}^n \setminus \{\mathbf{0}\} \right\} = \max \left\{ \frac{\|\mathbf{D}\widehat{\mathbf{z}}\|}{\|\mathbf{Q}\widehat{\mathbf{z}}\|} : \widehat{\mathbf{z}} \in \mathbb{K}^n \setminus \{\mathbf{0}\} \right\} \\ &= \max \left\{ \frac{\|\mathbf{D}\widehat{\mathbf{z}}\|}{\|\widehat{\mathbf{z}}\|} : \widehat{\mathbf{z}} \in \mathbb{K}^n \setminus \{\mathbf{0}\} \right\} = \|\mathbf{D}\|. \end{aligned}$$

Da  $\mathbf{A}$  und  $\mathbf{D}$  ähnliche Matrizen sind, muss nach Lemma 3.13 ein  $j \in [1 : n]$  mit  $|d_{jj}| = \varrho(\mathbf{A})$  existieren. Für jeden Vektor  $\widehat{\mathbf{z}} \in \mathbb{K}^n$  gilt

$$\|\mathbf{D}\widehat{\mathbf{z}}\|^2 = \sum_{i=1}^n |d_{ii}|^2 |\widehat{z}_i|^2 \leq \sum_{i=1}^n |d_{jj}|^2 |\widehat{z}_i|^2 = \varrho(\mathbf{A})^2 \sum_{i=1}^n |\widehat{z}_i|^2 = \varrho(\mathbf{A})^2 \|\widehat{\mathbf{z}}\|^2.$$

Andererseits gilt auch

$$\|\mathbf{D}\delta^{(j)}\| = \|d_{jj}\delta^{(j)}\| = |d_{jj}| \|\delta^{(j)}\| = \varrho(\mathbf{A}) \|\delta^{(j)}\|,$$

also folgt  $\varrho(\mathbf{A}) = \|\mathbf{D}\|$  und der Beweis ist abgeschlossen.  $\blacksquare$

Es stellt sich die Frage, ob sich eine Matrix über die Schur-Zerlegung hinaus durch Ähnlichkeitstransformationen weiter vereinfachen lässt. Um die „Einfachheit“ messen zu können, benötigen wir eine geeignete Norm:

**Definition 3.56 (Frobenius-Norm)** *Die Abbildung*

$$\|\cdot\|_F : \mathbb{K}^{n \times m} \rightarrow \mathbb{R}_{\geq 0}, \quad \mathbf{A} \mapsto \left( \sum_{i=1}^n \sum_{j=1}^m |a_{ij}|^2 \right)^{1/2},$$

*ist eine Norm auf  $\mathbb{K}^{n \times m}$  und wird als Frobenius-Norm bezeichnet.*

Eine für unsere Untersuchungen wesentliche Eigenschaft der Frobenius-Norm ist ihre Invarianz unter unitären Transformationen:

**Lemma 3.57 (Invarianz der Frobenius-Norm)** *Sei  $\mathbf{A} \in \mathbb{K}^{n \times m}$ . Seien  $\mathbf{P} \in \mathbb{K}^{k \times n}$  und  $\mathbf{Q} \in \mathbb{K}^{k \times m}$  isometrische Matrizen. Dann gilt*

$$\|\mathbf{P}\mathbf{A}\|_F = \|\mathbf{A}\|_F = \|\mathbf{A}\mathbf{Q}^*\|_F.$$

### 3.7 Diagonalisierbarkeit durch unitäre Transformationen

*Beweis.* Wir bezeichnen die Spalten der Matrix  $\mathbf{A}$  mit

$$\mathbf{a}_j := \mathbf{A}\delta^{(j)} \quad \text{für alle } j \in [1 : m].$$

Nach Definition 3.56 und Lemma 3.34 gilt

$$\|\mathbf{A}\|_F^2 = \sum_{j=1}^m \|\mathbf{a}_j\|_2^2 = \sum_{j=1}^m \|\mathbf{P}\mathbf{a}_j\|_2^2 = \|\mathbf{P}\mathbf{A}\|_F^2.$$

Wegen

$$\|\mathbf{A}\|_F = \|\mathbf{A}^*\|_F$$

folgt aus der Orthogonalität von  $\mathbf{Q}^*$  und Lemma 3.18 auch

$$\|\mathbf{A}\|_F = \|\mathbf{A}^*\|_F = \|\mathbf{Q}\mathbf{A}^*\|_F = \|(\mathbf{Q}\mathbf{A}^*)^*\|_F = \|\mathbf{A}\mathbf{Q}^*\|_F.$$

■

Als Maß für die „Diagonalität“ einer Matrix verwenden wir die Norm ihrer Außerdiagonaleinträge:

**Definition 3.58** Für  $\mathbf{A} \in \mathbb{K}^{n \times n}$  definieren wir

$$\text{off}(\mathbf{A}) := \left( \sum_{i=1}^n \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|^2 \right)^{1/2}.$$

Offenbar gilt  $\text{off}^2(\mathbf{A}) = \|\mathbf{A}\|_F^2 - \sum_{i=1}^n |a_{ii}|^2$ .

Jetzt lässt sich zeigen, dass diese Größe für jede Schur-Zerlegung einer festen Matrix dieselbe ist, sich also eine Matrix mit unitären Transformationen nicht beliebig weit einer Diagonalmatrix annähern lässt:

**Satz 3.59 (Invariante)** Sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Seien  $\mathbf{Q}_1, \mathbf{Q}_2 \in \mathbb{C}^{n \times n}$  unitäre Matrizen und  $\mathbf{R}_1, \mathbf{R}_2 \in \mathbb{C}^{n \times n}$  obere Dreiecksmatrizen mit

$$\mathbf{Q}_1 \mathbf{R}_1 \mathbf{Q}_1^* = \mathbf{A} = \mathbf{Q}_2 \mathbf{R}_2 \mathbf{Q}_2^*.$$

Dann gilt

$$\text{off}(\mathbf{R}_1) = \text{off}(\mathbf{R}_2).$$

*Beweis.* Im Körper  $\mathbb{C}$  zerfällt das charakteristische Polynom  $p_A$  in Linearfaktoren, wir finden also Nullstellen  $\lambda_1, \dots, \lambda_n \in \mathbb{C}$  mit

$$p_A(\lambda) = \prod_{i=1}^n (\lambda - \lambda_i).$$

### 3 Theoretische Grundlagen

Da  $\mathbf{R}_1$  und  $\mathbf{R}_2$  Dreiecksmatrizen sind, liegen ihre Eigenwerte auf ihren Diagonalen. Da beide Matrizen der Matrix  $\mathbf{A}$  ähnlich sind, folgt

$$\sum_{i=1}^n |r_{1,ii}|^2 = \sum_{i=1}^n |\lambda_i|^2 = \sum_{i=1}^n |r_{2,ii}|^2.$$

Es gelten

$$\begin{aligned} \text{off}^2(\mathbf{R}_1) &= \|\mathbf{R}_1\|_F^2 - \sum_{i=1}^n |r_{1,ii}|^2 = \|\mathbf{R}_1\|_F^2 - \sum_{i=1}^n |\lambda_i|^2, \\ \text{off}^2(\mathbf{R}_2) &= \|\mathbf{R}_2\|_F^2 - \sum_{i=1}^n |r_{2,ii}|^2 = \|\mathbf{R}_2\|_F^2 - \sum_{i=1}^n |\lambda_i|^2. \end{aligned}$$

Mit Lemma 3.57 folgen

$$\begin{aligned} \text{off}^2(\mathbf{R}_1) &= \|\mathbf{R}_1\|_F^2 - \sum_{i=1}^n |\lambda_i|^2 = \|\mathbf{Q}_1 \mathbf{R}_1 \mathbf{Q}_1^*\|_F^2 - \sum_{i=1}^n |\lambda_i|^2 = \|\mathbf{A}\|_F^2 - \sum_{i=1}^n |\lambda_i|^2, \\ \text{off}^2(\mathbf{R}_2) &= \|\mathbf{R}_2\|_F^2 - \sum_{i=1}^n |\lambda_i|^2 = \|\mathbf{Q}_2 \mathbf{R}_2 \mathbf{Q}_2^*\|_F^2 - \sum_{i=1}^n |\lambda_i|^2 = \|\mathbf{A}\|_F^2 - \sum_{i=1}^n |\lambda_i|^2, \end{aligned}$$

und daraus ergibt sich unmittelbar die gewünschte Identität. ■

Für alle Schur-Zerlegungen  $\mathbf{A} = \mathbf{Q}\mathbf{R}\mathbf{Q}^*$  ist also der Außerdiagonalteil der Matrix  $\mathbf{R}$  gleich groß. Wir wussten bereits, dass er nur verschwindet, wenn  $\mathbf{A}$  eine normale Matrix ist. Nun wissen wir, dass wir ihn auch durch noch so geschickte unitäre Transformationen nicht reduzieren können.

## 3.8 Nicht-unitäre Transformationen

Lässt man statt der unitären Transformationen allgemeinere Ähnlichkeitstransformationen zu, kann man die in der Schur-Zerlegung auftretende Dreiecksmatrix durch eine Block-Diagonalmatrix ersetzen.

Dazu betrachten wir eine Blockmatrix

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ & \mathbf{A}_{22} \end{pmatrix}$$

mit  $\mathbf{A} \in \mathbb{K}^{n \times n}$ ,  $\mathbf{A}_{11} \in \mathbb{K}^{m \times m}$ ,  $\mathbf{A}_{22} \in \mathbb{K}^{(n-m) \times (n-m)}$  und  $\mathbf{A}_{12} \in \mathbb{K}^{m \times (n-m)}$  und suchen nach einer Ähnlichkeitstransformation der Form

$$\mathbf{B} = \begin{pmatrix} \mathbf{I} & \mathbf{X} \\ & \mathbf{I} \end{pmatrix}, \quad \mathbf{B}^{-1} = \begin{pmatrix} \mathbf{I} & -\mathbf{X} \\ & \mathbf{I} \end{pmatrix}$$

mit  $\mathbf{X} \in \mathbb{K}^{m \times (n-m)}$ , die  $\mathbf{A}$  in Block-Diagonalform überführt. Aufgrund der Gleichung

$$\mathbf{B}^{-1}\mathbf{A}\mathbf{B} = \begin{pmatrix} \mathbf{I} & -\mathbf{X} \\ & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{11}\mathbf{X} + \mathbf{A}_{12} \\ & \mathbf{A}_{22} \end{pmatrix} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{11}\mathbf{X} - \mathbf{X}\mathbf{A}_{22} + \mathbf{A}_{12} \\ & \mathbf{A}_{22} \end{pmatrix} \quad (3.24)$$

ist das äquivalent dazu, ein  $\mathbf{X} \in \mathbb{K}^{m \times (n-m)}$  so zu finden, dass

$$\mathbf{A}_{11}\mathbf{X} - \mathbf{X}\mathbf{A}_{22} = -\mathbf{A}_{12}$$

gilt. Gleichungen dieser Form nennt man *Sylvester-Gleichungen*.

**Satz 3.60 (Sylvester-Gleichung)** *Seien  $\mathbf{A} \in \mathbb{C}^{n \times n}$  und  $\mathbf{B} \in \mathbb{C}^{m \times m}$  gegeben. Die Sylvester-Gleichung*

$$\mathbf{A}\mathbf{X} - \mathbf{X}\mathbf{B} = \mathbf{C} \quad (3.25)$$

*besitzt genau dann für alle  $\mathbf{C} \in \mathbb{C}^{n \times m}$  eine Lösung, wenn  $\sigma(\mathbf{A}) \cap \sigma(\mathbf{B}) = \emptyset$  gilt.*

*Beweis.* Wir werden zeigen, dass die lineare Abbildung

$$\mathcal{L}: \mathbb{C}^{n \times m} \rightarrow \mathbb{C}^{n \times m}, \quad \mathbf{X} \mapsto \mathbf{A}\mathbf{X} - \mathbf{X}\mathbf{B},$$

genau dann injektiv ist, wenn  $\sigma(\mathbf{A}) \cap \sigma(\mathbf{B}) = \emptyset$  gilt. Da  $\mathcal{L}$  ein Automorphismus zwischen endlich-dimensionalen Räumen ist, sind Injektivität und Surjektivität dank des Dimensionssatzes äquivalent.

Gelte zunächst  $\sigma(\mathbf{A}) \cap \sigma(\mathbf{B}) \neq \emptyset$ . Dann existiert ein Eigenwert  $\lambda \in \sigma(\mathbf{A}) \cap \sigma(\mathbf{B})$ . Da  $\lambda\mathbf{I} - \mathbf{B}$  nicht invertierbar ist, kann nach Lemma 3.18 auch die adjungierte Matrix  $\bar{\lambda}\mathbf{I} - \mathbf{B}^*$  nicht invertierbar sein, also ist  $\bar{\lambda}$  ein Eigenwert der Adjungierten  $\mathbf{B}^*$ .

Also finden wir  $\mathbf{x} \in \mathbb{C}^n \setminus \{\mathbf{0}\}$  und  $\mathbf{y} \in \mathbb{C}^m \setminus \{\mathbf{0}\}$  mit  $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$  und  $\mathbf{B}^*\mathbf{y} = \bar{\lambda}\mathbf{y}$ . Die Matrix  $\mathbf{X} := \mathbf{x}\mathbf{y}^*$  ist dann ungleich null und erfüllt

$$\mathcal{L}[\mathbf{X}] = \mathbf{A}\mathbf{X} - \mathbf{X}\mathbf{B} = \mathbf{A}\mathbf{x}\mathbf{y}^* - \mathbf{x}\mathbf{y}^*\mathbf{B} = \mathbf{A}\mathbf{x}\mathbf{y}^* - \mathbf{x}(\mathbf{B}^*\mathbf{y})^* = \lambda\mathbf{x}\mathbf{y}^* - \mathbf{x}(\bar{\lambda}\mathbf{y})^* = \mathbf{0}.$$

Wir haben also ein  $\mathbf{X} \in \mathbb{C}^{n \times m} \setminus \{\mathbf{0}\}$  mit  $\mathcal{L}[\mathbf{X}] = \mathbf{0}$  gefunden, somit ist  $\mathcal{L}$  nicht injektiv.

Sei nun umgekehrt ein  $\mathbf{X} \in \mathbb{C}^{n \times m} \setminus \{\mathbf{0}\}$  mit  $\mathcal{L}[\mathbf{X}] = \mathbf{0}$  gegeben. Dann gilt

$$\mathbf{A}\mathbf{X} = \mathbf{X}\mathbf{B}.$$

Wir nehmen zunächst an, dass  $\mathbf{B}$  eine obere Dreiecksmatrix ist. Sei  $j \in [1 : m]$  der kleinste Index mit  $\mathbf{X}\delta^{(j)} \neq \mathbf{0}$ , also der Index der ersten Spalte der Matrix  $\mathbf{X}$ , die nicht gleich null ist. Da wir  $\mathbf{X} \neq \mathbf{0}$  vorausgesetzt haben, existiert ein derartiger Index. Wir setzen  $\mathbf{e} := \mathbf{X}\delta^{(j)} \neq \mathbf{0}$ . Da  $\mathbf{B}$  eine obere Dreiecksmatrix ist, gilt

$$\mathbf{A}\mathbf{e} = \mathbf{A}\mathbf{X}\delta^{(j)} = \mathbf{X}\mathbf{B}\delta^{(j)} = \mathbf{X} \sum_{i=1}^j b_{ij}\delta^{(i)} = \sum_{i=1}^j b_{ij}\mathbf{X}\delta^{(i)}.$$

Aufgrund unserer Wahl des Index  $j$  gilt  $\mathbf{X}\delta^{(i)} = \mathbf{0}$  für alle  $i \in [1 : j - 1]$ , so dass wir

$$\mathbf{A}\mathbf{e} = b_{jj}\mathbf{X}\delta^{(j)} = b_{jj}\mathbf{e}$$

### 3 Theoretische Grundlagen

erhalten. Also ist  $\mathbf{e}$  ein Eigenvektor der Matrix  $\mathbf{A}$  zu dem Eigenwert  $b_{jj}$ . Da  $\mathbf{B}$  eine obere Dreiecksmatrix ist, ist das Diagonalelement  $b_{jj}$  auch ein Eigenwert der Matrix  $\mathbf{B}$ , es gilt also  $b_{jj} \in \sigma(\mathbf{A}) \cap \sigma(\mathbf{B})$ .

Den allgemeinen Fall können wir mit dem Satz 3.41 über die Schur-Zerlegung auf den bereits behandelten Sonderfall zurückführen: Wir finden eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{m \times m}$  und eine obere Dreiecksmatrix  $\mathbf{R} \in \mathbb{K}^{m \times m}$  mit

$$\mathbf{B} = \mathbf{Q}\mathbf{R}\mathbf{Q}^*,$$

und es gilt mit Lemma 3.35

$$\mathbf{A}\mathbf{X} = \mathbf{X}\mathbf{B} \iff \mathbf{A}\mathbf{X} = \mathbf{X}\mathbf{Q}\mathbf{R}\mathbf{Q}^* \iff \mathbf{A}\mathbf{X}\mathbf{Q} = \mathbf{X}\mathbf{Q}\mathbf{R}.$$

Wir setzen  $\widehat{\mathbf{X}} := \mathbf{X}\mathbf{Q}$  und erhalten

$$\mathbf{A}\widehat{\mathbf{X}} = \widehat{\mathbf{X}}\mathbf{R},$$

und da  $\mathbf{R}$  nun eine obere Dreiecksmatrix ist, können wir mit dem bereits Gezeigten folgern, dass  $\mathbf{A}$  und  $\mathbf{R}$  einen gemeinsamen Eigenwert besitzen müssen. Da  $\mathbf{R}$  und  $\mathbf{B}$  ähnlich sind, folgt auch  $\sigma(\mathbf{A}) \cap \sigma(\mathbf{B}) \neq \emptyset$ . ■

**Satz 3.61 (Block-Diagonalform)** Sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$ . Dann existieren eine reguläre Matrix  $\mathbf{B} \in \mathbb{C}^{n \times n}$  und  $m := |\sigma(\mathbf{A})|$ ,  $n_1, \dots, n_m \in \mathbb{N}$  sowie

$$\mathbf{R}_i \in \mathbb{C}^{n_i \times n_i}, \quad \text{für alle } i \in [1 : m]$$

mit

$$\mathbf{B}^{-1}\mathbf{A}\mathbf{B} = \begin{pmatrix} \mathbf{R}_1 & & \\ & \ddots & \\ & & \mathbf{R}_m \end{pmatrix}.$$

Jede der Matrizen  $\mathbf{R}_i$  ist in oberer Dreiecksform und besitzt genau einen Eigenwert.

*Beweis.* Wir führen den Beweis per Induktion über  $|\sigma(\mathbf{A})|$ .

*Induktionsanfang:* Falls eine Matrix  $\mathbf{A} \in \mathbb{C}^{n \times n}$  nur einen Eigenwert besitzt, ist die Aussage trivial, denn wir können  $\mathbf{A}_1 = \mathbf{A}$  setzen.

*Induktionsvoraussetzung:* Sei  $m \in \mathbb{N}$  so gegeben, dass für alle Matrizen  $\mathbf{A} \in \mathbb{C}^{n \times n}$  mit  $|\sigma(\mathbf{A})| = m$  die Aussage gilt.

*Induktionsschritt:* Sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$  mit  $|\sigma(\mathbf{A})| = m + 1$  gegeben.

Wir wählen einen Eigenwert  $\lambda \in \sigma(\mathbf{A})$  und konstruieren mit Satz 3.41 eine Schur-Zerlegung

$$\mathbf{Q}^*\mathbf{A}\mathbf{Q} = \begin{pmatrix} \mathbf{R}_1 & \mathbf{A}_{12} \\ & \mathbf{A}_{22} \end{pmatrix}.$$

Dabei wählen wir die Reihenfolge der Eigenwerte auf der Diagonalen so, dass  $\mathbf{R}_1 \in \mathbb{C}^{n_1 \times n_1}$  mit  $\sigma(\mathbf{R}_1) = \{\lambda\}$  und  $n_1 = \mu_A^a(\lambda)$  gilt. Es folgt  $\lambda \notin \sigma(\mathbf{A}_{22})$ , also insbesondere

$\sigma(\mathbf{R}_1) \cap \sigma(\mathbf{A}_{22}) = \emptyset$ , so dass wir Satz 3.60 anwenden können, um eine Matrix  $\mathbf{X} \in \mathbb{C}^{n_1 \times (n-n_1)}$  mit

$$\mathbf{R}_1 \mathbf{X} - \mathbf{X} \mathbf{A}_{22} = -\mathbf{A}_{12}$$

zu finden. Gemäß (3.24) gilt

$$\begin{pmatrix} \mathbf{I} & \mathbf{X} \\ & \mathbf{I} \end{pmatrix} \mathbf{Q}^* \mathbf{A} \mathbf{Q} \begin{pmatrix} \mathbf{I} & -\mathbf{X} \\ & \mathbf{I} \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \mathbf{X} \\ & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{R}_1 & \mathbf{A}_{12} \\ & \mathbf{A}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{I} & -\mathbf{X} \\ & \mathbf{I} \end{pmatrix} = \begin{pmatrix} \mathbf{R}_1 & \\ & \mathbf{A}_{22} \end{pmatrix}.$$

Wegen  $\sigma(\mathbf{A}) = \{\lambda\} \cup \sigma(\mathbf{A}_{22})$  können wir die Induktionsvoraussetzung auf  $\mathbf{A}_{22}$  anwenden, um eine invertierbare Matrix  $\widehat{\mathbf{B}} \in \mathbb{C}^{(n-n_1) \times (n-n_1)}$  und obere Dreiecksmatrizen  $\mathbf{R}_2, \dots, \mathbf{R}_{m+1}$  mit

$$\widehat{\mathbf{B}}^{-1} \mathbf{A}_{22} \widehat{\mathbf{B}} = \begin{pmatrix} \mathbf{R}_2 & & \\ & \ddots & \\ & & \mathbf{R}_{m+1} \end{pmatrix}$$

zu finden. Wir setzen

$$\mathbf{B} := \mathbf{Q} \begin{pmatrix} \mathbf{I} & -\mathbf{X} \\ & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{I} & \\ & \widehat{\mathbf{B}} \end{pmatrix}, \quad \mathbf{B}^{-1} = \begin{pmatrix} \mathbf{I} & \\ & \widehat{\mathbf{B}}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{I} & \mathbf{X} \\ & \mathbf{I} \end{pmatrix} \mathbf{Q}^*,$$

und erhalten

$$\mathbf{B}^{-1} \mathbf{A} \mathbf{B} = \begin{pmatrix} \mathbf{I} & \\ & \widehat{\mathbf{B}}^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{R}_1 & \\ & \mathbf{A}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{I} & \\ & \widehat{\mathbf{B}} \end{pmatrix} = \begin{pmatrix} \mathbf{R}_1 & & \\ & \mathbf{R}_2 & \\ & & \ddots \\ & & & \mathbf{R}_{m+1} \end{pmatrix}.$$

■

Wenn wir also darauf bereit sind, allgemeine Ähnlichkeitstransformationen zuzulassen, können wir zumindest im Komplexen jede beliebige Matrix immerhin auf Block-Diagonalform bringen, wobei jeder Block eine obere Dreiecksmatrix mit genau einem Eigenwert ist, also von der Form

$$\mathbf{R}_i = \begin{pmatrix} \lambda_i & r_{12} & \cdots & r_{1,n_i} \\ & \lambda_i & \ddots & \vdots \\ & & \ddots & r_{n_i-1,n_i} \\ & & & \lambda_i \end{pmatrix}.$$

Diese Diagonalblöcke lassen sich durch geeignete Transformationen noch weiter vereinfachen.

**Satz 3.62 (Jordan-Normalform)** Sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$  eine Matrix. Dann existieren eine invertierbare Matrix  $\mathbf{B} \in \mathbb{C}^{n \times n}$ ,  $m \in \mathbb{N}$ ,  $n_1, \dots, n_m \in \mathbb{N}$ ,  $\lambda_1, \dots, \lambda_m \in \mathbb{C}$  derart, dass

$$\mathbf{B}^{-1} \mathbf{A} \mathbf{B} = \begin{pmatrix} \mathbf{J}_1 & & \\ & \ddots & \\ & & \mathbf{J}_m \end{pmatrix},$$

### 3 Theoretische Grundlagen

mit Jordan-Blöcken der Gestalt

$$\mathbf{J}_i = \begin{pmatrix} \lambda_i & 1 & & \\ & \ddots & \ddots & \\ & & \lambda_i & 1 \\ & & & \lambda_i \end{pmatrix} \in \mathbb{C}^{n_i \times n_i} \quad \text{für alle } i \in [1 : m].$$

Es ist allerdings zu beachten, dass nicht-unitäre Ähnlichkeitstransformationen numerisch nicht ungefährlich sind, wie das folgende Beispiel zeigt:

**Beispiel 3.63** Wir betrachten die Matrizen

$$\mathbf{A}_\epsilon = \begin{pmatrix} 1 & 1 \\ 0 & 1 + \epsilon \end{pmatrix}$$

mit  $\epsilon \in \mathbb{R}_{>0}$ . Die Ähnlichkeitstransformation mit

$$\mathbf{B}_\epsilon := \begin{pmatrix} 1 & 1 \\ 0 & \epsilon \end{pmatrix}, \quad \mathbf{B}_\epsilon^{-1} = \begin{pmatrix} 1 & -1/\epsilon \\ 0 & 1/\epsilon \end{pmatrix}$$

führt zu

$$\mathbf{B}_\epsilon^{-1} \mathbf{A}_\epsilon \mathbf{B}_\epsilon = \begin{pmatrix} 1 & 0 \\ 0 & 1 + \epsilon \end{pmatrix}.$$

Also ist  $\mathbf{A}_\epsilon$  für jedes  $\epsilon > 0$  diagonalisierbar.

Vom numerischen Standpunkt aus gesehen ist diese Diagonalisierbarkeit allerdings fragwürdig, da die Kondition beispielsweise in der Zeilensummennorm durch

$$\kappa_\infty(\mathbf{B}_\epsilon) = \|\mathbf{B}_\epsilon\|_\infty \|\mathbf{B}_\epsilon^{-1}\|_\infty = 2 \left(1 + \frac{1}{\epsilon}\right) = 2 + \frac{2}{\epsilon}$$

gegeben ist und deshalb für kleiner werdendes  $\epsilon$  schnell gegen unendlich strebt.

**Übungsaufgabe 3.64 (Singularwertzerlegung)** Seien  $n, m \in \mathbb{N}$  und eine Matrix  $\mathbf{A} \in \mathbb{K}^{m \times n}$  gegeben.

(a) Beweisen Sie, dass es Vektoren  $\mathbf{u} \in \mathbb{K}^m$  und  $\mathbf{v} \in \mathbb{K}^n$  gibt mit  $\|\mathbf{u}\| = \|\mathbf{v}\| = 1$  und

$$|\langle \mathbf{u}, \mathbf{A}\mathbf{v} \rangle| = \max\{|\langle \mathbf{y}, \mathbf{A}\mathbf{x} \rangle| : \mathbf{x} \in \mathbb{K}^n, \mathbf{y} \in \mathbb{K}^m, \|\mathbf{x}\| = \|\mathbf{y}\| = 1\}.$$

(b) Folgern Sie daraus, dass es ein  $\sigma \in \mathbb{K}$  mit  $\mathbf{A}\mathbf{v} = \sigma\mathbf{u}$  und  $\mathbf{A}^*\mathbf{u} = \bar{\sigma}\mathbf{v}$  gibt.

(c) Folgern Sie daraus, dass es unitäre Matrizen  $\mathbf{U}_1 \in \mathbb{K}^{m \times m}$  und  $\mathbf{V}_1 \in \mathbb{K}^{n \times n}$  gibt mit

$$\mathbf{U}_1^* \mathbf{A} \mathbf{V}_1 = \begin{pmatrix} \sigma & \\ & \hat{\mathbf{A}} \end{pmatrix}, \quad \hat{\mathbf{A}} \in \mathbb{K}^{(m-1) \times (n-1)}.$$

(d) Folgern Sie daraus, dass es unitäre Matrizen  $\mathbf{U} \in \mathbb{K}^{m \times m}$  und  $\mathbf{V} \in \mathbb{K}^{n \times n}$  und eine Diagonalmatrix  $\mathbf{D} \in \mathbb{K}^{m \times n}$  gibt mit

$$\mathbf{U}^* \mathbf{A} \mathbf{V} = \mathbf{D}.$$

Hinweis: Der Satz von Heine-Borel und die Cauchy-Schwarz-Ungleichung können helfen.



### 3.9 Eigenwerte nicht-negativer Matrizen\*

Die Matrix (2.10) des „Minipoly“-Beispiels weist die Besonderheit auf, ausschließlich nicht-negative Einträge zu besitzen. Da wir den Eigenvektor zu dem maximalen Eigenwert 1 als invariantes Wahrscheinlichkeitsmaß interpretieren wollen, sind wir daran interessiert, einen Vektor zu erhalten, der keine negativen Einträge enthält.

Derartige Aufgabenstellungen lassen sich mit einer von Oskar Perron [3] und Georg Frobenius [1] entwickelten Methode untersuchen, die naheliegenderweise in der Literatur als *Perron-Frobenius-Theorie* bezeichnet wird.

**Definition 3.65 (Nicht-negative Matrizen und Vektoren)** Seien  $\mathbf{A} \in \mathbb{R}^{n \times n}$  und  $\mathbf{x} \in \mathbb{R}^n$ . Wir definieren

$$\mathbf{A} \geq 0 : \iff \forall i, j \in [1 : n] : a_{ij} \geq 0,$$

$$\mathbf{A} > 0 : \iff \forall i, j \in [1 : n] : a_{ij} > 0,$$

$$\mathbf{x} \geq 0 : \iff \forall i \in [1 : n] : x_i \geq 0,$$

$$\mathbf{x} > 0 : \iff \forall i \in [1 : n] : x_i > 0.$$

Mit  $K := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \geq 0\}$  bezeichnen wir den Kegel der nicht-negativen Vektoren.

Wir gehen im Folgenden davon aus, dass eine Matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  mit  $\mathbf{A} \geq 0$  gegeben ist.

Für unsere Untersuchung verwenden wir einen auf Helmut Wieland [4] zurückgehenden Ansatz, der auf der Abbildung

$$r : K \setminus \{\mathbf{0}\} \rightarrow \mathbb{R}_{\geq 0}, \quad \mathbf{x} \mapsto \min \left\{ \frac{(\mathbf{A}\mathbf{x})_i}{x_i} : i \in [1 : n], x_i \neq 0 \right\} \quad (3.26)$$

beruht. Für einen gegebenen Vektor  $\mathbf{x} \in K$  gilt dann

$$r(\mathbf{x})x_i \leq (\mathbf{A}\mathbf{x})_i \quad \text{für alle } i \in [1 : n],$$

wobei für den Fall  $x_i = 0$  die Nicht-Negativität und für den Fall  $x_i \neq 0$  die Definition der Abbildung  $r$  herangezogen wird. Diese Eigenschaft können wir kompakt als

$$\mathbf{A}\mathbf{x} - r(\mathbf{x})\mathbf{x} \geq 0 \quad (3.27)$$

schreiben. Wenn  $\mathbf{x}$  ein Eigenvektor wäre, würde in dieser Formel Gleichheit gelten, also bietet es sich an, nach Vektoren zu suchen, für die die linke Seite möglichst klein, also  $r(\mathbf{x})$  möglichst groß wird.

**Lemma 3.66 (Extremalvektoren)** Es gibt mindestens einen Vektor  $\mathbf{x}^* \in K \setminus \{\mathbf{0}\}$ , für den

$$r(\mathbf{x}) \leq r(\mathbf{x}^*) \quad \text{für alle } \mathbf{x} \in K \setminus \{\mathbf{0}\} \quad (3.28)$$

gilt, für den also  $r$  sein Maximum annimmt.

### 3 Theoretische Grundlagen

*Beweis.* Sei  $\mathbf{x} \in K \setminus \{\mathbf{0}\}$ , und sei  $\mathbf{e} \in K$  derjenige Vektor, dessen Einträge alle gleich eins sind. Aus (3.27) folgt

$$\langle \mathbf{e}, \mathbf{Ax} \rangle_2 - r(\mathbf{x}) \langle \mathbf{e}, \mathbf{x} \rangle_2 = \sum_{i=1}^n (\mathbf{Ax})_i - r(\mathbf{x})x_i \geq 0,$$

also insbesondere auch

$$r(\mathbf{x}) \leq \frac{\langle \mathbf{e}, \mathbf{Ax} \rangle_2}{\langle \mathbf{e}, \mathbf{x} \rangle_2}.$$

Wenn wir mir  $\alpha$  den größten Eintrag des Vektors  $\mathbf{A}^* \mathbf{e}$  bezeichnen, gilt

$$\langle \mathbf{e}, \mathbf{Ax} \rangle_2 = \langle \mathbf{A}^* \mathbf{e}, \mathbf{x} \rangle_2 = \sum_{i=1}^n (\mathbf{A}^* \mathbf{e})_i x_i \leq \sum_{i=1}^n \alpha x_i = \alpha \langle \mathbf{e}, \mathbf{x} \rangle_2,$$

also erhalten wir

$$r(\mathbf{x}) \leq \frac{\langle \mathbf{e}, \mathbf{Ax} \rangle_2}{\langle \mathbf{e}, \mathbf{x} \rangle_2} \leq \frac{\alpha \langle \mathbf{e}, \mathbf{x} \rangle_2}{\langle \mathbf{e}, \mathbf{x} \rangle_2} = \alpha,$$

so dass der Wertebereich der Abbildung  $r$  in  $[0, \alpha]$  enthalten sein muss. Demzufolge ist das Supremum

$$\varrho := \sup \{r(\mathbf{x}) : \mathbf{x} \in K_{\geq} \setminus \{\mathbf{0}\}\}$$

wohldefiniert und nicht-negativ.

Nach Definition des Supremums finden wir eine Folge von Vektoren  $(\mathbf{x}^{(m)})_{m=1}^{\infty}$  in  $K \setminus \{\mathbf{0}\}$  derart, dass

$$r(\mathbf{x}^{(m)}) \geq \varrho - 1/m \quad \text{für alle } m \in \mathbb{N}$$

gilt. Ein Blick auf die Definition (3.26) zeigt

$$r(\alpha \mathbf{x}) = r(\mathbf{x}) \quad \text{für alle } \alpha \in \mathbb{R}_{>0}, \mathbf{x} \in K \setminus \{\mathbf{0}\},$$

die Skalierung der Vektoren spielt also keine Rolle, so dass wir ohne Beschränkung der Allgemeinheit davon ausgehen können, dass

$$\langle \mathbf{e}, \mathbf{x}^{(m)} \rangle_2 = 1 \quad \text{für alle } m \in \mathbb{N}$$

gilt. Die Menge

$$T := \{\mathbf{x} \in K : \langle \mathbf{e}, \mathbf{x} \rangle_2 = 1\}$$

ist abgeschlossen und beschränkt, also nach dem Satz von Heine-Borel auch kompakt. Die Elemente der Folge  $(\mathbf{x}^{(m)})_{m=1}^{\infty}$  liegen in dieser Menge, also muss die Folge mindestens einen Häufungspunkt  $\mathbf{x}^* \in T$  besitzen.

Wenn  $r$  stetig wäre, könnten wir unmittelbar schließen, dass  $r(\mathbf{x}^*) = \varrho$  gelten muss. Leider ist  $r$  „nicht ganz“ stetig, da sich die Menge, über die in (3.26) das Minimum gebildet wird, abhängig von den Nulleinträgen des Vektors  $\mathbf{x}$  ändert.

### 3.9 Eigenwerte nicht-negativer Matrizen\*

Deshalb müssen wir die gewünschte Gleichung etwas ausführlicher nachprüfen. Dazu wählen wir ein beliebiges  $\epsilon \in \mathbb{R}_{>0}$  und setzen

$$\delta_1 := \min\{x_i^* : x_i^* \neq 0\}.$$

Wegen  $\mathbf{x}^* \in K \setminus \{\mathbf{0}\}$  dürfen wir  $\delta_1 > 0$  festhalten.

Aufgrund der Stetigkeit der Funktionen

$$\mathbf{x} \mapsto \frac{(\mathbf{A}\mathbf{x})_i}{x_i} \quad \text{für alle } i \in [1:n], x_i^* \neq 0$$

in der Nähe des Vektors  $\mathbf{x}^*$  können wir ein  $\delta_2 > 0$  so finden, dass

$$\frac{(\mathbf{A}\mathbf{x})_i}{x_i} \geq \frac{(\mathbf{A}\mathbf{x}^*)_i}{x_i} - \epsilon/2 \quad \text{für alle } \mathbf{x} \in K \text{ mit } \|\mathbf{x} - \mathbf{x}^*\|_2 \leq \delta_2$$

$$\text{und alle } i \in [1:n] \text{ mit } x_i^* \neq 0$$

erfüllt ist. Wir setzen  $\delta := \min\{\delta_1, \delta_2\}$ .

Da  $\mathbf{x}^*$  ein Häufungspunkt ist, finden wir ein  $m \in \mathbb{N}$  derart, dass

$$\|\mathbf{x}^{(m)} - \mathbf{x}^*\|_2 < \delta, \quad r(\mathbf{x}^{(m)}) \geq \varrho - \epsilon/2$$

gelten. Daraus folgt insbesondere

$$x_i^{(m)} = x_i^* + x_i^{(m)} - x_i^* \geq x_i^* - |x_i^{(m)} - x_i^*|$$

$$> x_i^* - \delta \geq 0 \quad \text{für alle } i \in [1:n] \text{ mit } x_i^* \neq 0,$$

so dass wir

$$\{i \in [1:n] : x_i^* \neq 0\} \subseteq \{i \in [1:n] : x_i^{(m)} \neq 0\}$$

erhalten. Daraus folgt insbesondere

$$r(\mathbf{x}^*) = \min \left\{ \frac{(\mathbf{A}\mathbf{x}^*)_i}{x_i^*} : i \in [1:n], x_i^* \neq 0 \right\}$$

$$\geq \min \left\{ \frac{(\mathbf{A}\mathbf{x}^*)_i}{x_i^*} : i \in [1:n], x_i^{(m)} \neq 0 \right\}$$

$$\geq \min \left\{ \frac{(\mathbf{A}\mathbf{x}^{(m)})_i}{x_i^{(m)}} - \epsilon/2 : i \in [1:n], x_i^{(m)} \neq 0 \right\}$$

$$= r(\mathbf{x}^{(m)}) - \epsilon/2 \geq \varrho - \epsilon.$$

Da  $\epsilon$  beliebig gewählt wurde, haben wir  $r(\mathbf{x}^*) = \varrho$  bewiesen. ■

Wir sind daran interessiert, Eigenvektoren zu finden, bei denen alle Komponenten positiv sind. Leider ist das nicht bei allen nicht-negativen Matrizen möglich, beispielsweise verschwindet bei allen Eigenvektoren der Matrix

$$\begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix} \quad (3.29)$$

die zweite Komponente. Um diesen Fall auszuschließen, müssen wir eine zusätzliche Bedingung an die Matrix  $\mathbf{A}$  stellen.

### 3 Theoretische Grundlagen

**Definition 3.67 (Reduzibel und irreduzible Matrizen)** Eine Matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  nennen wir *reduzibel*, falls ein  $k \in [1 : n - 1]$ , eine Permutationsmatrix  $\mathbf{P} \in \mathbb{R}^{n \times n}$  und Matrizen  $\mathbf{A}_{11} \in \mathbb{K}^{k \times k}$ ,  $\mathbf{A}_{12} \in \mathbb{K}^{k \times (n-k)}$  und  $\mathbf{A}_{22} \in \mathbb{K}^{(n-k) \times (n-k)}$  so existieren, dass

$$\mathbf{PAP}^{-1} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ & \mathbf{A}_{22} \end{pmatrix} \quad (3.30)$$

gilt. Falls eine Matrix nicht *reduzibel* ist, nennen wir sie *irreduzibel*.

Irreduzible Matrizen sind also gerade diejenigen, die sich nicht durch das Umsortieren von Indizes auf Block-Dreiecksform bringen lassen.

**Lemma 3.68 (Irreduzible Matrix)** Falls  $\mathbf{A} \in \mathbb{R}^{n \times n}$  irreduzibel ist und  $\mathbf{A} \geq 0$  gilt, existiert für jeden Vektor  $\mathbf{x} \in K \setminus \{\mathbf{0}\}$  ein  $m \in \mathbb{N}_0$  mit

$$(\mathbf{A} + \mathbf{I})^m \mathbf{x} > 0.$$

*Beweis.* Sei  $\mathbf{x} \in K \setminus \{\mathbf{0}\}$  gegeben. Wir definieren

$$\mathbf{x}^{(m)} := (\mathbf{A} + \mathbf{I})^m \mathbf{x} \quad \text{für alle } m \in \mathbb{N}_0.$$

Mit  $\mathbf{A} \geq 0$  folgt daraus unmittelbar

$$\mathbf{x}^{(m+1)} = \mathbf{Ax}^{(m)} + \mathbf{x}^{(m)} \geq \mathbf{x}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0.$$

In dem Vektor  $\mathbf{x}^{(m+1)}$  können also höchstens diejenigen Koeffizienten gleich null sein, die auch in  $\mathbf{x}^{(m)}$  bereits gleich null waren. Da  $\mathbf{x}^{(0)} = \mathbf{x} \neq \mathbf{0}$  gilt, kann insbesondere keiner der Vektoren  $\mathbf{x}^{(m)}$  der Nullvektor sein.

Wären für ein  $m \in \mathbb{N}_0$  dieselben Koeffizienten in  $\mathbf{x}^{(m)}$  und  $\mathbf{x}^{(m+1)}$  gleich null, könnten wir diese Koeffizienten mit einer Permutationsmatrix  $\mathbf{P}$  in die letzten  $n - k$  Komponenten umsordern und so

$$\mathbf{Px}^{(m)} = \begin{pmatrix} \widehat{\mathbf{x}}^{(m)} \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{Px}^{(m+1)} = \begin{pmatrix} \widehat{\mathbf{x}}^{(m+1)} \\ \mathbf{0} \end{pmatrix}$$

mit  $\widehat{\mathbf{x}}^{(m)}, \widehat{\mathbf{x}}^{(m+1)} \in \mathbb{R}^k$ ,  $k \in [1 : n - 1]$  sowie  $\widehat{\mathbf{x}}^{(m)} > 0$  und  $\widehat{\mathbf{x}}^{(m+1)} > 0$  zu erhalten.

Wir zerlegen die permutierte Matrix in der Form

$$\begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix} := \mathbf{PAP}^{-1}, \quad \mathbf{A}_{11} \in \mathbb{R}^{k \times k}, \quad \mathbf{A}_{22} \in \mathbb{R}^{(n-k) \times (n-k)}$$

und erhalten

$$\begin{aligned} \begin{pmatrix} \widehat{\mathbf{x}}^{(m+1)} \\ \mathbf{0} \end{pmatrix} &= \mathbf{Px}^{(m+1)} = \mathbf{P}(\mathbf{x}^{(m)} + \mathbf{Ax}^{(m)}) = \mathbf{Px}^{(m)} + \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix} \mathbf{Px}^{(m)} \\ &= \begin{pmatrix} \widehat{\mathbf{x}}^{(m)} \\ \mathbf{0} \end{pmatrix} + \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{x}}^{(m)} \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \widehat{\mathbf{x}}^{(m)} + \mathbf{A}_{11} \widehat{\mathbf{x}}^{(m)} \\ \mathbf{A}_{21} \widehat{\mathbf{x}}^{(m)} \end{pmatrix} \end{aligned}$$

folgen. Dank  $\widehat{\mathbf{x}}^{(m)} > 0$  und  $\mathbf{A}_{21} \geq 0$  müsste dann bereits  $\mathbf{A}_{21} = \mathbf{0}$  gelten, also wäre (3.30) erfüllt, im Widerspruch zur Voraussetzung.

Damit kann für jedes  $m \in \mathbb{N}$  der Vektor  $\mathbf{x}^{(m+1)}$  nicht dieselben Nulleinträge wie  $\mathbf{x}^{(m)}$  enthalten. Wir haben bereits gesehen, dass er auch nicht weitere Nulleinträge enthalten kann, also muss die Zahl der Nulleinträge sinken. Da  $\mathbf{x}^{(0)} = \mathbf{x} \neq \mathbf{0}$  höchstens  $n - 1$  Nulleinträge aufweisen kann, muss deshalb  $\mathbf{x}^{(m)} > 0$  spätestens für  $m = n - 1$  gelten. ■

**Bemerkung 3.69 (Irreduzible Matrix)** *Ein Blick auf den vorangehenden Beweis zeigt, dass  $(\mathbf{I} + \mathbf{A})^{n-1}\mathbf{x} > 0$  für alle Vektoren  $\mathbf{x} \in K \setminus \{\mathbf{0}\}$  gilt. Indem wir für  $\mathbf{x}$  die kanonischen Einheitsvektoren einsetzen folgt, dass jede Spalte der Matrix  $(\mathbf{I} + \mathbf{A})^{n-1}$  in jeder Komponente echt größer als Null ist, also haben wir sogar  $(\mathbf{I} + \mathbf{A})^{n-1} > 0$  erhalten.*

Falls  $\mathbf{A}$  irreduzibel ist, können wir wie bereits angedeutet folgern, dass ein Vektor  $\mathbf{x}^* \in K \setminus \{\mathbf{0}\}$  mit der in (3.28) gegebenen Extremaleigenschaft ein Eigenvektor sein muss.

**Lemma 3.70 (Eigenvektor)** *Sei  $\mathbf{A} \in \mathbb{R}^{n \times n}$  irreduzibel mit  $\mathbf{A} \geq 0$ . Sei  $\mathbf{x}^* \in K \setminus \{\mathbf{0}\}$  ein Vektor mit der Eigenschaft (3.28) und sei  $\varrho := r(\mathbf{x}^*)$ . Dann gelten  $\mathbf{x}^* > 0$  und*

$$\mathbf{A}\mathbf{x}^* = \varrho\mathbf{x}^*,$$

$\mathbf{x}^*$  ist also ein Eigenvektor zu dem Eigenwert  $\varrho$ .

*Beweis.* Aufgrund der Ungleichung (3.27) gilt

$$\mathbf{y} := \mathbf{A}\mathbf{x}^* - \varrho\mathbf{x}^* \geq 0.$$

Wir müssen nachweisen, dass  $\mathbf{y} = \mathbf{0}$  gilt.

Dazu verwenden wir einen Widerspruchsbeweis: Wir nehmen  $\mathbf{y} \neq \mathbf{0}$  an. Nach Lemma 3.68 existiert dann ein  $m \in \mathbb{N}_0$  so, dass

$$(\mathbf{I} + \mathbf{A})^m \mathbf{y} > 0$$

erfüllt ist. Für den Vektor

$$\mathbf{z} := (\mathbf{I} + \mathbf{A})^m \mathbf{x}^*$$

folgt daraus

$$\begin{aligned} \mathbf{A}\mathbf{z} - \varrho\mathbf{z} &= \mathbf{A}(\mathbf{I} + \mathbf{A})^m \mathbf{x}^* - \varrho(\mathbf{I} + \mathbf{A})^m \mathbf{x}^* \\ &= (\mathbf{I} + \mathbf{A})(\mathbf{I} + \mathbf{A})^{m-1} \mathbf{x}^* - (\varrho + 1)(\mathbf{I} + \mathbf{A})^{m-1} \mathbf{x}^* \\ &= (\mathbf{I} + \mathbf{A})^{m-1} (\mathbf{I} + \mathbf{A}) \mathbf{x}^* - (\varrho + 1)(\mathbf{I} + \mathbf{A})^{m-1} \mathbf{x}^* \\ &= (\mathbf{I} + \mathbf{A})^{m-1} (\mathbf{A}\mathbf{x}^* - \varrho\mathbf{x}^*) = (\mathbf{I} + \mathbf{A})^{m-1} \mathbf{y} > 0. \end{aligned}$$

Da  $\varrho \geq r(\mathbf{z})$  nach (3.28) gilt, erhalten wir insbesondere

$$\mathbf{A}\mathbf{z} - r(\mathbf{z})\mathbf{z} \geq \mathbf{A}\mathbf{z} - \varrho\mathbf{z} > 0.$$

### 3 Theoretische Grundlagen

Nach Definition des Minimums (3.26) muss aber ein  $i \in [1 : n]$  mit  $z_i \neq 0$  und

$$r(\mathbf{z}) = \frac{(\mathbf{Az})_i}{z_i}, \quad r(\mathbf{z})z_i = (\mathbf{Az})_i, \quad (\mathbf{Az} - r(\mathbf{z})\mathbf{z})_i = 0$$

existieren, im Widerspruch zu der obigen Ungleichung.

Also muss  $\mathbf{y} = \mathbf{0}$  gelten, somit ist  $\mathbf{x}^*$  ein Eigenvektor zu dem Eigenwert  $\varrho$ . ■

In Kombination mit Lemma 3.66 folgt also, dass eine nicht-negative irreduzible Matrix mindestens einen Eigenvektor zu dem Eigenwert  $\varrho$  besitzt und dass die Koeffizienten dieses Eigenvektors alle echt positiv sind.

Der Eigenwert  $\varrho$  lässt sich sogar als im Betrag maximaler Eigenwert der gesamten Matrix  $\mathbf{A}$  identifizieren:

**Lemma 3.71 (Maximaler Eigenwert)** Sei  $\mathbf{A} \in \mathbb{R}^{n \times n}$  eine irreduzible Matrix mit  $\mathbf{A} \geq 0$ . Sei  $\varrho$  wie in Lemma 3.70 definiert, und sei  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  ein beliebiger Eigenvektor der Matrix  $\mathbf{A}$  zu dem Eigenwert  $\lambda \in \mathbb{K}$ . Dann gilt  $|\lambda| \leq \varrho$ .

*Beweis.* Wir definieren den Vektor  $\mathbf{y} \in K \setminus \{\mathbf{0}\}$  durch

$$y_i := |x_i| \quad \text{für alle } i \in [1 : n]$$

und folgern mit der Dreiecksungleichung

$$|\lambda|y_i = |\lambda x_i| = |(\mathbf{Ax})_i| = \left| \sum_{j=1}^n a_{ij}x_j \right| \leq \sum_{j=1}^n a_{ij}|x_j| = (\mathbf{Ay})_i \quad \text{für alle } i \in [1 : n].$$

Für jedes  $i \in [1 : n]$  mit  $y_i \neq 0$  folgt daraus

$$|\lambda| \leq \frac{(\mathbf{Ay})_i}{y_i},$$

also nach (3.26) insbesondere

$$|\lambda| \leq r(\mathbf{y}) \leq \varrho,$$

und damit ist die gewünschte Aussage bewiesen. ■

Wir können die Ergebnisse dieses Abschnitts in dem folgenden Satz zusammenfassen:

**Satz 3.72 (Perron-Frobenius)** Sei  $\mathbf{A} \in \mathbb{R}^{n \times n}$  irreduzibel mit  $\mathbf{A} \geq 0$ . Der Spektralradius der Matrix ist durch

$$\varrho(\mathbf{A}) := \max\{|\lambda| : \lambda \in \sigma(\mathbf{A})\}$$

gegeben. Er ist ein Eigenwert der Matrix  $\mathbf{A}$ , zu dem ein Eigenvektor  $\mathbf{x} \in \mathbb{R}^n$  mit  $\mathbf{x} > \mathbf{0}$  existiert.

### 3.9 Eigenwerte nicht-negativer Matrizen\*

*Beweis.* Sei  $\mathbf{x}^* \in K \setminus \{\mathbf{0}\}$  der nach Lemma 3.66 existierende Vektor mit der Eigenschaft (3.28). Nach Lemma 3.70 ist er ein Eigenvektor zu dem Eigenwert  $\varrho = r(\mathbf{x}^*)$  und erfüllt  $\mathbf{x}^* > 0$ . Nach Lemma 3.71 muss auch  $\varrho = \varrho(\mathbf{A})$  gelten. ■

**Bemerkung 3.73 (Einfacher Eigenwert)** *Es lässt sich beweisen, dass  $\varrho(\mathbf{A})$  unter den im vorangehenden Satz gegebenen Bedingungen sogar ein einfacher Eigenwert der Matrix  $\mathbf{A}$  ist. Der betreffende Beweis ist in [4, Beweis von I.f-g] zu finden.*





## 4 Die Jacobi-Iteration

In diesem Kapitel soll eines der einfachsten Iterationsverfahren zur Bestimmung der Schur-Zerlegung einer selbstadjungierte Matrix eingeführt werden: Die Jacobi-Iteration.

### 4.1 Iterierte Ähnlichkeitstransformationen

Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix. Unser Ziel ist die Bestimmung der Schur-Zerlegung

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{D}$$

für unitäre Matrizen  $\mathbf{Q}$  und eine Diagonalmatrix  $\mathbf{D}$ .

Wenn es möglich wäre, eine derartige Zerlegung in endlich vielen Rechenschritten für eine beliebige Dimension  $n$  zu berechnen, könnte man auch sämtliche Nullstellen von Polynomen beliebig hoher Ordnung mit endlich vielen Schritten berechnen. Da bekannt ist, dass diese Aufgabe nicht lösbar ist, dürfen wir auch nicht darauf hoffen, die Schur-Zerlegung mit endlich vielen Rechenoperationen bestimmen zu können.

Stattdessen müssen wir auf Näherungsverfahren zurückgreifen. Unser Ziel ist es, die Matrix  $\mathbf{A}$  mit Hilfe einer unitären Ähnlichkeitstransformation auf Diagonalgestalt zu bringen. Die Diagonalgestalt ist durch  $\text{off}(\mathbf{D}) = 0$  charakterisiert, also ist die Idee nahe liegend, nach Verfahren zu suchen, die diese Größe reduzieren: Wir beginnen mit  $\mathbf{A}^{(0)} := \mathbf{A}$  und bestimmen zu jedem  $m \in \mathbb{N}_0$  ein unitäres  $\mathbf{Q}^{(m)}$ , das  $\text{off}((\mathbf{Q}^{(m)})^* \mathbf{A}^{(m)} \mathbf{Q}^{(m)})$  verkleinert und setzen dann

$$\mathbf{A}^{(m+1)} := (\mathbf{Q}^{(m)})^* \mathbf{A}^{(m)} \mathbf{Q}^{(m)}.$$

Wir brechen ab, sobald  $\text{off}(\mathbf{A}^{(m+1)})$  klein genug ist und verwenden

$$\tilde{\mathbf{Q}} := \mathbf{Q}^{(0)} \mathbf{Q}^{(1)} \dots \mathbf{Q}^{(m)}, \quad \tilde{\mathbf{D}} := \mathbf{A}^{(m+1)} = \tilde{\mathbf{Q}}^* \mathbf{A} \tilde{\mathbf{Q}}$$

als Approximationen von  $\mathbf{Q}$  und  $\mathbf{D}$ .

Da jeder einzelne Schritt des Verfahrens eine unitäre Ähnlichkeitstransformation ist, muss  $\tilde{\mathbf{Q}}$  ebenfalls unitär sein, und an der Größe  $\text{off}(\tilde{\mathbf{D}})$  lässt sich direkt ablesen, wie nahe wir einer durch  $\text{off}(\mathbf{D}) = 0$  charakterisierten exakten Schur-Zerlegung gekommen sind.

Ein Schritt des Verfahrens, also die Berechnung von  $\mathbf{A}^{(m+1)} = (\mathbf{Q}^{(m)})^* \mathbf{A}^{(m)} \mathbf{Q}^{(m)}$ , kann in zwei Teilschritte zerlegt werden:

$$\mathbf{B}^{(m)} := (\mathbf{Q}^{(m)})^* \mathbf{A}^{(m)}, \quad \mathbf{A}^{(m+1)} := \mathbf{B}^{(m)} \mathbf{Q}^{(m)} = ((\mathbf{Q}^{(m)})^* (\mathbf{B}^{(m)})^*)^*.$$

Der erste Schritt entspricht der Anwendung von  $(\mathbf{Q}^{(m)})^*$  auf jede Spalte von  $\mathbf{A}^{(m)}$ , der zweite der Anwendung derselben Matrix auf jede Zeile von  $\mathbf{B}^{(m)}$ . Sofern sich also die Multiplikation mit  $(\mathbf{Q}^{(m)})^*$  effizient gestalten lässt, kann auch die Iteration effizient durchgeführt werden.

## 4.2 Zweidimensionaler Fall

Wie bei vielen anderen numerischen Verfahren empfiehlt es sich, die Strategie zur Lösung eines großen und komplizierten Problems aus Lösungen einer Folge kleinerer und einfacherer Probleme zu konstruieren.

Wir untersuchen zunächst den Fall einer zweidimensionalen selbstadjungierten Matrix

$$\mathbf{A} = \begin{pmatrix} a & b \\ \bar{b} & d \end{pmatrix}$$

mit  $a, d \in \mathbb{R}$  und  $b \in \mathbb{K} \setminus \{0\}$ . Die Eigenwerte dieser Matrix können wir mit Hilfe des charakteristischen Polynoms einfach bestimmen: Für jeden Eigenwert  $\lambda \in \sigma(\mathbf{A})$  muss

$$\begin{aligned} 0 &= \det(\lambda \mathbf{I} - \mathbf{A}) = \det \begin{pmatrix} \lambda - a & -b \\ -\bar{b} & \lambda - d \end{pmatrix} = (\lambda - a)(\lambda - d) - |b|^2 \\ &= \lambda^2 - (d + a)\lambda + ad - |b|^2 = \lambda^2 - (d + a)\lambda + \frac{(d + a)^2}{4} - \frac{(d + a)^2}{4} + \frac{4ad}{4} - |b|^2 \\ &= \left( \lambda - \frac{d + a}{2} \right)^2 - \frac{(d - a)^2 + 4|b|^2}{4} \end{aligned}$$

gelten, also folgt

$$\lambda = \frac{d + a}{2} + \sigma \sqrt{\frac{(d - a)^2}{4} + |b|^2} \quad \text{für ein } \sigma \in \{-1, 1\}. \quad (4.1)$$

Nun kennen wir die Eigenwerte, also müssen wir als nächstes die Eigenvektoren bestimmen. Wir verwenden den Ansatz

$$\mathbf{x} = \begin{pmatrix} c \\ tc \end{pmatrix}$$

für  $t, c \in \mathbb{K}$  mit  $c \neq 0$ . Der Eigenvektor  $\mathbf{x}$  zu dem Eigenwert  $\lambda$  aus (4.1) ergibt sich als Lösung der definierenden Gleichung

$$\mathbf{0} = (\lambda \mathbf{I} - \mathbf{A})\mathbf{x} = \begin{pmatrix} \frac{d-a}{2} + \sigma \sqrt{\frac{(d-a)^2}{4} + |b|^2} & -b \\ -\bar{b} & \frac{a-d}{2} + \sigma \sqrt{\frac{(d-a)^2}{4} + |b|^2} \end{pmatrix} \begin{pmatrix} c \\ tc \end{pmatrix}. \quad (4.2)$$

Aus der ersten Zeile dieser Gleichung folgt

$$btc = \left( \frac{d-a}{2} + \sigma \sqrt{\frac{(d-a)^2}{4} + |b|^2} \right) c, \quad t = \frac{d-a}{2b} + \frac{\sigma}{b} \sqrt{\frac{(d-a)^2}{4} + |b|^2}.$$

Zur Abkürzung verwenden wir

$$\tau := \frac{d-a}{2b}, \quad t = \tau + \sigma \frac{|b|}{b} \sqrt{|\tau|^2 + 1}. \quad (4.3)$$

Nun müssen wir die zweite Zeile der Gleichung (4.2) prüfen. Mit der dritten binomischen Formel erhalten wir

$$\begin{aligned}
\left(\frac{a-d}{2} + \sigma\sqrt{\frac{(d-a)^2}{4} + |b|^2}\right)t &= b\left(-\frac{d-a}{2b} + \sigma\frac{|b|}{b}\sqrt{|\tau|^2+1}\right)t \\
&= b\left(-\tau + \sigma\frac{|b|}{b}\sqrt{|\tau|^2+1}\right)\left(\tau + \sigma\frac{|b|}{b}\sqrt{|\tau|^2+1}\right) \\
&= b\left(\frac{|b|^2}{b^2}(|\tau|^2+1) - \tau^2\right) = \frac{1}{b}(|b|^2|\tau|^2 + |b|^2 - b^2\tau^2) \\
&= \frac{1}{b}\left(\frac{(d-a)^2}{4} + |b|^2 - \frac{(d-a)^2}{4}\right) = \frac{|b|^2}{b} = \bar{b} \quad (4.4)
\end{aligned}$$

und damit

$$-\bar{b}c + \left(\frac{a-d}{2} + \sigma\sqrt{\frac{(d-a)^2}{4} + |b|^2}\right)tc = -\bar{b}c + \bar{b}c = 0,$$

also ist  $\mathbf{x}$  für das in (4.3) definierte  $t$  tatsächlich ein Eigenvektor. Da wir an einem normierten Eigenvektor interessiert sind, also  $\|\mathbf{x}\|_2 = 1$  gelten soll, müssen wir  $c$  so bestimmen, dass

$$1 = |c|^2 + |t|^2|c|^2 = |c|^2(1 + |t|^2)$$

erfüllt ist, also setzen wir

$$c = \frac{1}{\sqrt{1 + |t|^2}}.$$

Zur Abkürzung der Notation verwenden wir

$$s := tc, \quad \mathbf{x} = \begin{pmatrix} c \\ s \end{pmatrix}.$$

Die Bestimmung des zweiten Eigenvektors gestaltet sich wesentlich einfacher: Da  $\mathbf{A}$  selbstadjungiert ist, müssen Eigenvektoren zu verschiedenen Eigenwerten senkrecht aufeinander stehen, und im zweidimensionalen Raum ist es leicht, einen auf  $\mathbf{x}$  senkrecht stehenden Vektor zu berechnen. Eine naheliegende Wahl ist

$$\mathbf{x}^\perp = \begin{pmatrix} -\bar{s} \\ \bar{c} \end{pmatrix},$$

und da dieser Vektor ebenfalls normiert ist, muss die gesuchte unitäre Transformation  $\mathbf{Q}$  durch

$$\mathbf{Q} = \begin{pmatrix} c & -\bar{s} \\ s & \bar{c} \end{pmatrix}$$

gegeben sein. Da  $|c|^2 + |s|^2 = 1$  gilt, können wir diese Matrix im Fall  $\mathbb{K} = \mathbb{R}$  als Rotation um die Null interpretieren, wobei der Rotationswinkel  $\varphi$  durch  $c = \cos(\varphi)$  und  $s = \sin(\varphi)$  gegeben ist. Dann folgt  $t = s/c = \frac{\sin(\varphi)}{\cos(\varphi)} = \tan(\varphi)$ .

#### 4 Die Jacobi-Iteration

**Lemma 4.1 (Zweidimensionale Schur-Zerlegung)** Für eine beliebige selbstadjungierte Matrix

$$\mathbf{A} = \begin{pmatrix} a & b \\ \bar{b} & d \end{pmatrix}$$

mit  $a, d \in \mathbb{R}$  und  $b \in \mathbb{K} \setminus \{0\}$  setzen wir

$$\tau := \frac{d-a}{2b},$$

wählen ein Vorzeichen  $\sigma \in \{-1, 1\}$  und verwenden

$$t := \tau + \sigma \overline{\operatorname{sgn}(b)} \sqrt{|\tau|^2 + 1}, \quad c := \frac{1}{\sqrt{1 + |t|^2}}, \quad s := tc \quad (4.5)$$

zur Definition von

$$\mathbf{Q} := \begin{pmatrix} c & -\bar{s} \\ s & \bar{c} \end{pmatrix}. \quad (4.6)$$

Dann gilt

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \begin{pmatrix} \frac{d+a}{2} + \sigma |b| \sqrt{|\tau|^2 + 1} & 0 \\ 0 & \frac{d+a}{2} - \sigma |b| \sqrt{|\tau|^2 + 1} \end{pmatrix},$$

wir haben also eine explizite Formel für die Schur-Zerlegung gefunden.

*Beweis.* Die Eigenwerte kürzen wir mit

$$\lambda_1 := \frac{a+d}{2} + \sigma |b| \sqrt{|\tau|^2 + 1}, \quad \lambda_2 := \frac{a+d}{2} - \sigma |b| \sqrt{|\tau|^2 + 1} \quad (4.7)$$

ab und berechnen zunächst

$$\mathbf{A} \mathbf{Q} = \begin{pmatrix} a & b \\ \bar{b} & d \end{pmatrix} \begin{pmatrix} c & -\bar{s} \\ s & \bar{c} \end{pmatrix} = \begin{pmatrix} ac + bs & -a\bar{s} + b\bar{c} \\ \bar{b}c + ds & -\bar{b}\bar{s} + d\bar{c} \end{pmatrix} = \begin{pmatrix} (a+bt)c & (b-a\bar{t})\bar{c} \\ (\bar{b}+dt)c & (d-\bar{b}\bar{t})\bar{c} \end{pmatrix}.$$

Um diesen Ausdruck zu vereinfachen verwenden wir (4.4) und die Definitionen, um

$$\begin{aligned} a + bt &= a + \frac{d-a}{2} + \sigma |b| \sqrt{|\tau|^2 + 1} = \lambda_1, \\ \bar{b} + dt &= \left( \frac{a-d}{2} + \sigma |b| \sqrt{|\tau|^2 + 1} \right) t + dt = \lambda_1 t, \\ b - a\bar{t} &= \overline{\left( \frac{a-d}{2} + \sigma |b| \sqrt{|\tau|^2 + 1} \right) \bar{t}} - a\bar{t} = -\lambda_2 \bar{t}, \\ d - \bar{b}\bar{t} &= d - \frac{d-a}{2} - \sigma |b| \sqrt{|\tau|^2 + 1} = \lambda_2 \end{aligned}$$

zu erhalten und so zu

$$\mathbf{A} \mathbf{Q} = \begin{pmatrix} (a+bt)c & (b-a\bar{t})\bar{c} \\ (\bar{b}+dt)c & (d-\bar{b}\bar{t})\bar{c} \end{pmatrix} = \begin{pmatrix} \lambda_1 c & -\lambda_2 \bar{s} \\ \lambda_1 s & \lambda_2 \bar{c} \end{pmatrix}$$

zu gelangen. Multiplikation mit  $\mathbf{Q}^*$  ergibt

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \begin{pmatrix} \bar{c} & \bar{s} \\ -s & c \end{pmatrix} \begin{pmatrix} \lambda_1 c & -\lambda_2 \bar{s} \\ \lambda_1 s & \lambda_2 \bar{c} \end{pmatrix} = \begin{pmatrix} \lambda_1(|c|^2 + |s|^2) & \lambda_2(-\bar{c}\bar{s} + \bar{s}\bar{c}) \\ \lambda_1(-sc + cs) & \lambda_2(|s|^2 + |c|^2) \end{pmatrix} = \begin{pmatrix} \lambda_1 & \\ & \lambda_2 \end{pmatrix},$$

so dass die gewünschte Aussage bewiesen ist.  $\blacksquare$

Falls  $\mathbf{A}$  schon „fast“ diagonal ist, könnte es passieren, dass durch die Transformation  $\mathbf{Q}$  die Diagonalelemente den Platz tauschen. Das würde zu unerwünschten Problemen bei der Untersuchung der Konvergenz des Verfahrens führen und sollte deshalb vermieden werden. Glücklicherweise lässt sich dieser Effekt vermeiden, indem das Vorzeichen  $\sigma$  geschickt gewählt wird.

**Lemma 4.2** *Seien  $\mathbf{A}, \mathbf{Q}, \tau, t, c$  und  $s$  wie in Lemma 4.1 gegeben. Dann gilt*

$$\|\mathbf{A} - \mathbf{Q}^* \mathbf{A} \mathbf{Q}\|_F^2 = \frac{2}{c^2} |b|^2.$$

*Beweis.* Wir kürzen die Eigenwerte wie in (4.7) ab und erhalten

$$\begin{aligned} |a - \lambda_1|^2 &= \left| a - \frac{d+a}{2} - \sigma |b| \sqrt{|\tau|^2 + 1} \right|^2 = \left| \frac{a-d}{2} - \sigma |b| \sqrt{|\tau|^2 + 1} \right|^2 \\ &= \left| \frac{d-a}{2} + \sigma |b| \sqrt{|\tau|^2 + 1} \right|^2 = |b|^2 \left| \frac{d-a}{2b} + \sigma \frac{|b|}{b} \sqrt{|\tau|^2 + 1} \right|^2 \\ &= |b|^2 \left| \tau + \sigma \overline{\operatorname{sgn}(b)} \sqrt{|\tau|^2 + 1} \right|^2 = |b|^2 |t|^2. \end{aligned}$$

Für den zweiten Diagonaleintrag ergibt sich

$$|d - \lambda_2|^2 = \left| d - \frac{a+d}{2} + \sigma |b| \sqrt{|\tau|^2 + 1} \right|^2 = \left| \frac{d-a}{2} + \sigma |b| \sqrt{|\tau|^2 + 1} \right|^2 = |a - \lambda_1|^2,$$

und da die beiden Außerdiagonaleinträge von  $\mathbf{Q}^* \mathbf{A} \mathbf{Q}$  verschwinden, erhalten wir

$$\|\mathbf{A} - \mathbf{Q}^* \mathbf{A} \mathbf{Q}\|_F^2 = 2|b|^2 |t|^2 + 2|b|^2 = 2|b|^2 (|t|^2 + 1) = \frac{2|b|^2}{c^2},$$

also die gewünschte Gleichung.  $\blacksquare$

Wenn wir sicher stellen wollen, dass die Matrix so wenig wie möglich verändert wird, müssen wir also darauf achten, dass  $c$  möglichst groß ist, also  $|t|$  möglichst klein. Falls  $\tau = 0$  gilt, spielt das Vorzeichen von  $t$  keine Rolle. Anderenfalls muss das Vorzeichen von  $\sigma \overline{\operatorname{sgn}(b)}$  gerade dem Vorzeichen von  $\tau$  entgegengesetzt sein, so dass wir

$$\begin{aligned} \sigma \overline{\operatorname{sgn}(b)} &= \frac{\sigma}{\operatorname{sgn}(b)} \stackrel{!}{=} -\operatorname{sgn}(\tau), \\ \sigma &= -\operatorname{sgn}(b) \operatorname{sgn}(\tau) = -\operatorname{sgn}(b\tau) = \operatorname{sgn}\left(\frac{d-a}{2}\right) = -\operatorname{sgn}(d-a) = \operatorname{sgn}(a-d) \end{aligned}$$

#### 4 Die Jacobi-Iteration

erhalten. Dann entsteht allerdings ein Problem: Bei der Berechnung von

$$t = \tau - \operatorname{sgn}(\tau)\sqrt{|\tau|^2 + 1}$$

kann der Algorithmus instabil werden, falls  $|\tau|$  groß ist, denn dann gilt  $\sqrt{|\tau|^2 + 1} \approx |\tau|$ , also

$$\operatorname{sgn}(\tau)\sqrt{|\tau|^2 + 1} \approx \operatorname{sgn}(\tau)|\tau| = \tau,$$

so dass sich die beiden Summanden näherungsweise auslöschen.

Das Problem lässt sich lösen, indem wir  $t$  indirekt berechnen: Wir wählen das entgegengesetzte Vorzeichen  $-\sigma$  und berechnen die zweite Nullstelle

$$\hat{t} := \tau + \operatorname{sgn}(\tau)\sqrt{|\tau|^2 + 1}.$$

Offenbar sind  $t$  und  $\hat{t}$  Nullstellen des Polynoms

$$\begin{aligned} z \mapsto (z - t)(z - \hat{t}) &= z^2 - (t + \hat{t})z + t\hat{t} = z^2 - 2\tau z + \tau^2 - \operatorname{sgn}(\tau)^2(|\tau|^2 + 1) \\ &= z^2 - 2\tau z + \tau^2 - \operatorname{sgn}(\tau)^2|\tau|^2 - \operatorname{sgn}(\tau)^2 \\ &= z^2 - 2\tau z - \operatorname{sgn}(\tau)^2, \end{aligned}$$

und mittels eines Koeffizientenvergleichs folgt

$$t\hat{t} = -\operatorname{sgn}(\tau)^2, \quad t = -\frac{\operatorname{sgn}(\tau)^2}{\hat{t}} = -\frac{\operatorname{sgn}(\tau)}{|\tau| + \sqrt{|\tau|^2 + 1}} = \frac{\operatorname{sgn}(a - d)\overline{\operatorname{sgn}(b)}}{|\tau| + \sqrt{|\tau|^2 + 1}}.$$

Diese Gleichung erlaubt es uns, die benötigte Größe  $t$  auch für große Werte von  $\tau$  noch stabil zu berechnen.

### 4.3 Höherdimensionaler Fall

Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix. Zwei Indizes  $i, j \in [1 : n]$  mit  $i < j$  legen einen Außerdiagonaleintrag  $a_{ij}$  in  $\mathbf{A}$  fest, genauer gesagt einen Eintrag oberhalb der Diagonalen. Gesucht ist eine unitäre Ähnlichkeitstransformation  $\mathbf{Q} \in \mathbb{K}^{n \times n}$ , die  $a_{ij}$  zu 0 macht, die also für  $\mathbf{B} := \mathbf{Q}^* \mathbf{A} \mathbf{Q}$  die Gleichung  $b_{ij} = 0$  erfüllt.

Wenn  $n = 2$  gelten würde, ließe sich  $\mathbf{Q}$  mit Lemma 4.1 bestimmen. Es stellt sich also die Frage, ob sich das  $n$ -dimensionale Problem auf das zweidimensionale reduzieren lässt. Dazu setzen wir

$$\hat{\mathbf{A}} = \begin{pmatrix} a_{ii} & a_{ij} \\ a_{ji} & a_{jj} \end{pmatrix} \tag{4.8}$$

und wählen

$$\hat{\mathbf{Q}} = \begin{pmatrix} c & -\bar{s} \\ s & \bar{c} \end{pmatrix}$$

gemäß Lemma 4.1. Es folgt, dass

$$\hat{\mathbf{B}} := \hat{\mathbf{Q}}^* \hat{\mathbf{A}} \hat{\mathbf{Q}}$$

eine Diagonalmatrix ist.

Wir definieren  $\mathbf{Q} \in \mathbb{K}^{n \times n}$ , indem wir  $\widehat{\mathbf{Q}}$  auf den von dem  $i$ -ten und  $j$ -ten kanonischen Einheitsvektor aufgespannten Unterraum  $\widehat{\mathcal{E}} := \text{span}\{\delta^{(i)}, \delta^{(j)}\}$  anwenden und seinen Senkrechtraum  $\widehat{\mathcal{E}}^\perp := \{\mathbf{x} \in \mathbb{K}^n : x_i = x_j = 0\}$  unverändert lassen.

$\widehat{\mathcal{E}}$  wird von der isometrischen Matrix

$$\widehat{\mathbf{E}} := (\delta^{(i)} \quad \delta^{(j)})$$

aufgespannt und wir haben

$$\widehat{\mathbf{A}} = \widehat{\mathbf{E}}^* \mathbf{A} \widehat{\mathbf{E}}.$$

Wir suchen eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$ , die in  $\widehat{\mathcal{E}}$  dieselbe Wirkung wie  $\widehat{\mathbf{Q}}$  entfaltet, die also

$$\widehat{\mathbf{E}}^* \mathbf{Q} \widehat{\mathbf{E}} = \widehat{\mathbf{Q}}$$

erfüllt. Da  $\widehat{\mathbf{E}}$  isometrisch ist, wäre

$$\mathbf{Q} = \widehat{\mathbf{E}} \widehat{\mathbf{Q}} \widehat{\mathbf{E}}^*$$

ein erster Ansatz, diese Matrix ist allerdings nicht unitär, die Matrix  $\widehat{\mathbf{E}}^*$  besitzt im Fall  $n > 2$  einen nicht-trivialen Kern.

Das Problem lässt sich lösen, indem wir den „Kern wieder hinzuaddieren“: Die Matrix  $\mathbf{I} - \widehat{\mathbf{E}} \widehat{\mathbf{E}}^*$  ist eine orthogonale Projektion auf den Kern.

**Lemma 4.3 (Jacobi-Givens-Rotation)** *Die Matrix*

$$\mathbf{Q} := (\mathbf{I} - \widehat{\mathbf{E}} \widehat{\mathbf{E}}^*) + \widehat{\mathbf{E}} \widehat{\mathbf{Q}} \widehat{\mathbf{E}}^* \tag{4.9}$$

ist unitär und erfüllt

$$\widehat{\mathbf{E}}^* \mathbf{Q} \widehat{\mathbf{E}} = \widehat{\mathbf{Q}}, \quad \widehat{\mathbf{E}}^* \mathbf{Q}^* \mathbf{A} \mathbf{Q} \widehat{\mathbf{E}} = \widehat{\mathbf{Q}}^* \widehat{\mathbf{A}} \widehat{\mathbf{Q}} = \widehat{\mathbf{B}}. \tag{4.10}$$

*Beweis.* Wir halten zunächst

$$\widehat{\mathbf{E}}^* \mathbf{Q} = \widehat{\mathbf{E}}^* (\mathbf{I} - \widehat{\mathbf{E}} \widehat{\mathbf{E}}^* + \widehat{\mathbf{E}} \widehat{\mathbf{Q}} \widehat{\mathbf{E}}^*) = \widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}^* \widehat{\mathbf{E}} \widehat{\mathbf{E}}^* + \widehat{\mathbf{E}}^* \widehat{\mathbf{E}} \widehat{\mathbf{Q}} \widehat{\mathbf{E}}^* = \widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}^* + \widehat{\mathbf{Q}} \widehat{\mathbf{E}}^* = \widehat{\mathbf{Q}} \widehat{\mathbf{E}}^*$$

fest. Analog können dank Lemma 3.18 auch

$$\widehat{\mathbf{E}}^* \mathbf{Q}^* = \widehat{\mathbf{E}}^* (\mathbf{I} - \widehat{\mathbf{E}} \widehat{\mathbf{E}}^* + \widehat{\mathbf{E}} \widehat{\mathbf{Q}} \widehat{\mathbf{E}}^*) = \widehat{\mathbf{Q}}^* \widehat{\mathbf{E}}^*, \quad \mathbf{Q} \widehat{\mathbf{E}} = \widehat{\mathbf{E}} \widehat{\mathbf{Q}}$$

bewiesen werden. Daraus folgen unmittelbar

$$\widehat{\mathbf{E}}^* \mathbf{Q} \widehat{\mathbf{E}} = \widehat{\mathbf{Q}} \widehat{\mathbf{E}}^* \widehat{\mathbf{E}} = \widehat{\mathbf{Q}}, \quad \widehat{\mathbf{E}}^* \mathbf{Q}^* \mathbf{A} \mathbf{Q} \widehat{\mathbf{E}} = \widehat{\mathbf{Q}}^* \widehat{\mathbf{E}}^* \mathbf{A} \widehat{\mathbf{E}} \widehat{\mathbf{Q}} = \widehat{\mathbf{Q}}^* \widehat{\mathbf{A}} \widehat{\mathbf{Q}} = \widehat{\mathbf{B}}.$$

Es bleibt noch zu zeigen, dass  $\mathbf{Q}$  unitär ist. Wir rechnen die Gleichung mit Lemma 3.18 direkt nach:

$$\mathbf{Q}^* \mathbf{Q} = (\mathbf{I} - \widehat{\mathbf{E}} \widehat{\mathbf{E}}^* + \widehat{\mathbf{E}} \widehat{\mathbf{Q}} \widehat{\mathbf{E}}^*)^* (\mathbf{I} - \widehat{\mathbf{E}} \widehat{\mathbf{E}}^* + \widehat{\mathbf{E}} \widehat{\mathbf{Q}} \widehat{\mathbf{E}}^*)$$

#### 4 Die Jacobi-Iteration

$$\begin{aligned}
&= \mathbf{I} - \widehat{\mathbf{E}}\widehat{\mathbf{E}}^* + \widehat{\mathbf{E}}\widehat{\mathbf{Q}}\widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}\widehat{\mathbf{E}}^* + \widehat{\mathbf{E}}\widehat{\mathbf{E}}^*\widehat{\mathbf{E}}\widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}\widehat{\mathbf{E}}^*\widehat{\mathbf{E}}\widehat{\mathbf{Q}}\widehat{\mathbf{E}}^* \\
&\quad + \widehat{\mathbf{E}}\widehat{\mathbf{Q}}^*\widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}\widehat{\mathbf{Q}}^*\widehat{\mathbf{E}}^*\widehat{\mathbf{E}}\widehat{\mathbf{E}}^* + \widehat{\mathbf{E}}\widehat{\mathbf{Q}}^*\widehat{\mathbf{E}}^*\widehat{\mathbf{E}}\widehat{\mathbf{Q}}\widehat{\mathbf{E}}^* \\
&= \mathbf{I} - \widehat{\mathbf{E}}\widehat{\mathbf{E}}^* + \widehat{\mathbf{E}}\widehat{\mathbf{Q}}\widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}\widehat{\mathbf{E}}^* + \widehat{\mathbf{E}}\widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}\widehat{\mathbf{Q}}\widehat{\mathbf{E}}^* \\
&\quad + \widehat{\mathbf{E}}\widehat{\mathbf{Q}}^*\widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}\widehat{\mathbf{Q}}^*\widehat{\mathbf{E}}^* + \widehat{\mathbf{E}}\widehat{\mathbf{Q}}^*\widehat{\mathbf{Q}}\widehat{\mathbf{E}}^* \\
&= \mathbf{I} - \widehat{\mathbf{E}}\widehat{\mathbf{E}}^* + \widehat{\mathbf{E}}\widehat{\mathbf{Q}}\widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}\widehat{\mathbf{E}}^* + \widehat{\mathbf{E}}\widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}\widehat{\mathbf{Q}}\widehat{\mathbf{E}}^* \\
&\quad + \widehat{\mathbf{E}}\widehat{\mathbf{Q}}^*\widehat{\mathbf{E}}^* - \widehat{\mathbf{E}}\widehat{\mathbf{Q}}^*\widehat{\mathbf{E}}^* + \widehat{\mathbf{E}}\widehat{\mathbf{E}}^* = \mathbf{I}.
\end{aligned}$$

■

Die Tatsache, dass  $\mathbf{Q}$  nur auf die  $i$ -te und  $j$ -te Komponente eines Vektors wirkt, wird in der „Punktchennotation“ der Matrix erkennbar:

$$\mathbf{Q} = \begin{pmatrix} 1 & \cdots & 0 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots & & \vdots \\ 0 & \cdots & c & \cdots & -\bar{s} & \cdots & 0 \\ \vdots & & \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & s & \cdots & \bar{c} & \cdots & 0 \\ \vdots & & \vdots & & \vdots & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & 0 & \cdots & 1 \end{pmatrix}. \quad (4.11)$$

Da  $\widehat{\mathbf{Q}}$  eine Rotation ist, ist  $\mathbf{Q}$  eine Rotation in der  $(i, j)$ -Koordinatenebene.

Die für uns wichtige Eigenschaft dieser Rotation ist, dass sie die Einträge  $a_{ij}$  und  $a_{ji}$  eliminiert: Für  $\mathbf{B} := \mathbf{Q}^*\mathbf{A}\mathbf{Q}$  folgt mit (4.10)

$$\begin{pmatrix} b_{ii} & b_{ij} \\ b_{ji} & b_{jj} \end{pmatrix} = \widehat{\mathbf{E}}^*\widehat{\mathbf{B}}\widehat{\mathbf{E}} = \widehat{\mathbf{E}}^*\widehat{\mathbf{Q}}^*\widehat{\mathbf{A}}\widehat{\mathbf{Q}}\widehat{\mathbf{E}} = \widehat{\mathbf{B}},$$

und da  $\widehat{\mathbf{B}}$  eine Diagonalmatrix ist, muss  $b_{ij} = b_{ji} = 0$  gelten.

Es bleibt die Frage, ob die Ähnlichkeitstransformation mit  $\mathbf{Q}$  tatsächlich den Außerdiagonalanteil von  $\mathbf{B}$  gegenüber  $\mathbf{A}$  reduziert hat. Diese Frage beantwortet das folgende Lemma:

**Lemma 4.4** Sei  $\mathbf{Q}$  wie in (4.9) für  $i, j \in [1 : n]$  mit  $i < j$  definiert. Dann gilt für die Matrix  $\mathbf{B} := \mathbf{Q}^*\mathbf{A}\mathbf{Q}$  die Gleichung

$$\text{off}^2(\mathbf{B}) = \text{off}^2(\mathbf{A}) - 2|a_{ij}|^2.$$

*Beweis.* Wir wenden Lemma 3.57 auf die Matrizen  $\widehat{\mathbf{A}}$  und  $\widehat{\mathbf{B}}$  an und erhalten

$$|a_{ii}|^2 + 2|a_{ij}|^2 + |a_{jj}|^2 = \|\widehat{\mathbf{A}}\|_F^2 = \|\widehat{\mathbf{Q}}^*\widehat{\mathbf{A}}\widehat{\mathbf{Q}}\|_F^2 = \|\widehat{\mathbf{B}}\|_F^2 = |b_{ii}|^2 + |b_{jj}|^2.$$



Die einzigen Diagonalelemente, in denen sich  $\mathbf{B}$  von  $\mathbf{A}$  unterscheidet, sind  $b_{ii}$  und  $b_{jj}$ , so dass sich wieder mit Lemma 3.57 für die Norm die Außerdiagonalelemente die Gleichung

$$\begin{aligned} \text{off}^2(\mathbf{B}) &= \|\mathbf{B}\|_F^2 - \sum_{k=1}^n |b_{kk}|^2 = \|\mathbf{Q}^* \mathbf{A} \mathbf{Q}\|_F^2 - \sum_{\substack{k=1 \\ k \neq i, k \neq j}}^n |b_{kk}|^2 - |b_{ii}|^2 - |b_{jj}|^2 \\ &= \|\mathbf{A}\|_F^2 - \sum_{\substack{k=1 \\ k \neq i, k \neq j}}^n |a_{kk}|^2 - |a_{ii}|^2 - |a_{jj}|^2 - 2|a_{ij}|^2 \\ &= \|\mathbf{A}\|_F^2 - \sum_{k=1}^n |a_{kk}|^2 - 2|a_{ij}|^2 = \text{off}^2(\mathbf{A}) - 2|a_{ij}|^2 \end{aligned}$$

ergibt. Also bewirkt die Ähnlichkeitstransformation von  $\mathbf{A}$  mit  $\mathbf{Q}$  dieselbe Reduktion der Außerdiagonalelemente wie die von  $\hat{\mathbf{A}}$  mit  $\hat{\mathbf{Q}}$ . ■

## 4.4 Algorithmus

Seien  $i, j \in [1 : n]$  mit  $i < j$  gegeben. Wir berechnen  $c, s \in \mathbb{K}$  entsprechend Lemma 4.1 und sind daran interessiert, die in (4.9) gegebene Transformation  $\mathbf{Q}$  möglichst effizient anzuwenden. Zunächst betrachten wir dazu die Matrix  $\mathbf{M} = \mathbf{A} \mathbf{Q}$ , deren Einträge durch

$$\mathbf{M} = \begin{pmatrix} a_{11} & \cdots & ca_{1i} + sa_{1j} & \cdots & -\bar{s}a_{1i} + \bar{c}a_{1j} & \cdots & a_{1n} \\ \vdots & & \vdots & & \vdots & & \vdots \\ a_{i1} & \cdots & ca_{ii} + sa_{ij} & \cdots & -\bar{s}a_{ii} + \bar{c}a_{ij} & \cdots & a_{in} \\ \vdots & & \vdots & & \vdots & & \vdots \\ a_{j1} & \cdots & ca_{ji} + sa_{jj} & \cdots & -\bar{s}a_{ji} + \bar{c}a_{jj} & \cdots & a_{jn} \\ \vdots & & \vdots & & \vdots & & \vdots \\ a_{n1} & \cdots & ca_{ni} + sa_{nj} & \cdots & -\bar{s}a_{ni} + \bar{c}a_{nj} & \cdots & a_{nn} \end{pmatrix}$$

gegeben sind: Nur die  $i$ -te und  $j$ -te Spalte wurden verändert, ein Algorithmus könnte diesen Schritt also in  $6n$  Operationen durchführen.

Die Matrix  $\mathbf{B} = \mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{Q}^* \mathbf{M}$  ist entsprechend durch

$$\mathbf{B} = \begin{pmatrix} m_{11} & \cdots & m_{1i} & \cdots & m_{1j} & \cdots & m_{1n} \\ \vdots & & \vdots & & \vdots & & \vdots \\ \bar{c}m_{i1} + \bar{s}m_{j1} & \cdots & \bar{c}m_{ii} + \bar{s}m_{ji} & \cdots & \bar{c}m_{ij} + \bar{s}m_{jj} & \cdots & \bar{c}m_{in} + \bar{s}m_{jn} \\ \vdots & & \vdots & & \vdots & & \vdots \\ -sm_{i1} + cm_{j1} & \cdots & -sm_{ii} + cm_{ji} & \cdots & -sm_{ij} + cm_{jj} & \cdots & -sm_{in} + cm_{jn} \\ \vdots & & \vdots & & \vdots & & \vdots \\ m_{n1} & \cdots & m_{ni} & \cdots & m_{nj} & \cdots & m_{nn} \end{pmatrix}$$

beschrieben: Hier wurden nur die  $i$ -te und die  $j$ -te Zeile verändert, also kann auch dieser Schritt in  $6n$  Operationen durchgeführt werden. Wir erhalten den folgenden Algorithmus:

**Algorithmus 4.5 (Jacobi-Iteration)** Der folgende Algorithmus überschreibt  $\mathbf{A}$  mit einer zu  $\mathbf{A}$  ähnlichen Matrix  $\mathbf{B}$  mit  $\text{off}(\mathbf{B}) \leq \epsilon$  und speichert die korrespondierende Ähnlichkeitstransformation in der Matrix  $\tilde{\mathbf{Q}}$ :

```

 $\tilde{\mathbf{Q}} \leftarrow \mathbf{I};$ 
 $\alpha \leftarrow \text{off}(\mathbf{A})^2;$ 
while  $\alpha > \epsilon^2$  do begin
  Wähle  $i, j \in [1 : n]$  mit  $i < j$ ;
   $\alpha \leftarrow \alpha - 2|a_{ij}|^2;$ 
   $\tau \leftarrow \frac{a_{jj} - a_{ii}}{2a_{ij}}; \quad t \leftarrow \frac{\text{sgn}(a_{ii} - a_{jj})\overline{\text{sgn}(a_{ij})}}{|\tau| + \sqrt{|\tau|^2 + 1}};$ 
   $c \leftarrow 1/\sqrt{1 + |t|^2}; \quad s \leftarrow tc;$ 
  for  $k = 1$  to  $n$  do begin    { Berechne  $\mathbf{A} \leftarrow \mathbf{A}\mathbf{Q}$  und  $\tilde{\mathbf{Q}} \leftarrow \tilde{\mathbf{Q}}\mathbf{Q}$  }
     $h \leftarrow a_{ki}; \quad a_{ki} \leftarrow hc + a_{kj}s; \quad a_{kj} \leftarrow -h\bar{s} + a_{kj}\bar{c};$ 
     $h \leftarrow \tilde{q}_{ki}; \quad q_{ki} \leftarrow hc + \tilde{q}_{kj}s; \quad q_{kj} \leftarrow -h\bar{s} + \tilde{q}_{kj}\bar{c}$ 
  end;
  for  $k = 1$  to  $n$  do begin    { Berechne  $\mathbf{A} \leftarrow \mathbf{Q}^*\mathbf{A}$  }
     $h \leftarrow a_{ik}; \quad a_{ik} \leftarrow \bar{c}h + \bar{s}a_{jk}; \quad a_{jk} \leftarrow -sh + ca_{jk}$ 
  end end

```

Falls wir in diesem Algorithmus das Paar  $1 \leq i < j \leq n$  so wählen, dass  $|a_{ij}|$  den maximalen Wert annimmt, können wir eine einfache Konvergenzaussage herleiten:

**Satz 4.6 (Konvergenz)** Seien  $i, j \in [1 : n]$  mit  $i < j$  so gewählt, dass

$$|a_{k\ell}| \leq |a_{ij}| \quad \text{für alle } k, \ell \in [1 : n] \text{ mit } k < \ell \quad (4.12)$$

gilt. Dann erfüllt die Matrix  $\mathbf{B}$  aus Lemma 4.4 die Abschätzung

$$\text{off}^2(\mathbf{B}) \leq \left(1 - \frac{2}{n(n-1)}\right) \text{off}^2(\mathbf{A}),$$

also konvergiert der Algorithmus mit mindestens linearer Geschwindigkeit.

*Beweis.* Seien  $i, j \in [1 : n]$  mit  $i < j$  so gewählt, dass (4.12) gilt. Dann folgt

$$\text{off}^2(\mathbf{A}) = \sum_{\substack{k, \ell=1 \\ k \neq \ell}}^n |a_{k\ell}|^2 = 2 \sum_{\substack{k, \ell=1 \\ k < \ell}}^n |a_{k\ell}|^2 \leq 2 \sum_{\substack{k, \ell=1 \\ k < \ell}}^n |a_{ij}|^2 = 2 \frac{n(n-1)}{2} |a_{ij}|^2 = n(n-1) |a_{ij}|^2.$$

und damit nach Lemma 4.4

$$\text{off}^2(\mathbf{B}) = \text{off}^2(\mathbf{A}) - 2|a_{ij}|^2 \leq \text{off}^2(\mathbf{A}) - \frac{2}{n(n-1)} \text{off}^2(\mathbf{A}) = \left(1 - \frac{2}{n(n-1)}\right) \text{off}^2(\mathbf{A}).$$

Das ist die gesuchte Abschätzung. ■

Lemma 4.1 legt die Rotation  $\mathbf{Q}$  nur bis auf das Vorzeichen  $\sigma \in \{1, -1\}$  fest. Es stellt sich die Frage, ob eine geschickte Wahl des Vorzeichens Vorteile für das Verfahren bietet. Das folgende Lemma deutet ein mögliches Auswahlkriterium an:

**Lemma 4.7** Sei  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  wie in (4.9) definiert, sei  $\mathbf{B} = \mathbf{Q}^* \mathbf{A} \mathbf{Q}$ . Dann gilt

$$\|\mathbf{A} - \mathbf{B}\|_F^2 = 4(1 - c) \sum_{k \neq i, j} (|a_{ki}|^2 + |a_{kj}|^2) + \frac{2}{c^2} |a_{ij}|^2.$$

*Beweis.* Sei  $\mathbf{M} := \mathbf{A} \mathbf{Q}$ . Sei  $k \in [1 : n]$ . Dann gilt

$$\begin{aligned} |m_{ki} - a_{ki}|^2 &= (m_{ki} - a_{ki})(\bar{m}_{ki} - \bar{a}_{ki}) = |m_{ki}|^2 + |a_{ki}|^2 - m_{ki}\bar{a}_{ki} - a_{ki}\bar{m}_{ki} \\ &= |m_{ki}|^2 + |a_{ki}|^2 - 2 \operatorname{Re}(m_{ki}\bar{a}_{ki}) \\ &= |m_{ki}|^2 + |a_{ki}|^2 - 2 \operatorname{Re}((a_{ki}c + a_{kj}s)\bar{a}_{ki}), \\ |m_{kj} - a_{kj}|^2 &= |m_{kj}|^2 + |a_{kj}|^2 - 2 \operatorname{Re}(a_{kj}\bar{m}_{kj}) \\ &= |m_{kj}|^2 + |a_{kj}|^2 - 2 \operatorname{Re}(a_{kj}(-\bar{a}_{ki}s + \bar{a}_{kj}c)) \end{aligned}$$

und durch Addition beider Gleichungen erhalten wir

$$|m_{ki} - a_{ki}|^2 + |m_{kj} - a_{kj}|^2 = |m_{ki}|^2 + |m_{kj}|^2 + |a_{ki}|^2 + |a_{kj}|^2 - 2(|a_{ki}|^2 + |a_{kj}|^2) \operatorname{Re} c.$$

Da die Transformation mit  $\hat{\mathbf{Q}}$  unitär ist, muss  $|m_{ki}|^2 + |m_{kj}|^2 = |a_{ki}|^2 + |a_{kj}|^2$  gelten und wir erhalten

$$|m_{ki} - a_{ki}|^2 + |m_{kj} - a_{kj}|^2 = 2(1 - \operatorname{Re} c)(|a_{ki}|^2 + |a_{kj}|^2).$$

Indem wir dieselbe Argumentation auf  $\mathbf{B} = \mathbf{Q}^* \mathbf{M}$  anwenden, erhalten wir

$$|b_{ik} - m_{ik}|^2 + |b_{jk} - m_{jk}|^2 = 2(1 - \operatorname{Re} c)(|m_{ik}|^2 + |m_{jk}|^2).$$

Da die Spaltentransformation nur die Spalten  $i, j$  betrifft und die Zeilentransformation nur die Zeilen  $i, j$  verändert, folgt

$$\begin{aligned} |b_{ki} - a_{ki}|^2 + |b_{kj} - a_{kj}|^2 &= 2(1 - \operatorname{Re} c)(|a_{ki}|^2 + |a_{kj}|^2) && \text{für alle } k \notin \{i, j\}, \\ |b_{ik} - a_{ik}|^2 + |b_{jk} - a_{jk}|^2 &= 2(1 - \operatorname{Re} c)(|a_{ki}|^2 + |a_{kj}|^2) && \text{für alle } k \notin \{i, j\}, \end{aligned}$$

wobei im letzten Schritt ausgenutzt wurde, dass  $\mathbf{A}$  selbstadjungiert ist. Da  $c$  in unserem Fall immer reell und positiv ist, erhalten wir den ersten Summanden unserer Gleichung.

Es fehlt nur noch die Betrachtung der Elemente  $a_{ii}$ ,  $a_{ij}$ ,  $a_{ji}$  und  $a_{jj}$ . Dazu greifen wir auf Lemma 4.2 zurück und erhalten

$$|a_{ii} - b_{ii}|^2 + |a_{ij} - b_{ij}|^2 + |a_{ji} - b_{ji}|^2 + |a_{jj} - b_{jj}|^2 = \frac{2}{c^2} |a_{ij}|^2,$$

und der Beweis ist vollständig. ■

Wir können Sprünge in der Folge der Iterierten vermeiden, indem wir dafür sorgen, dass  $c$  möglichst große Werte annimmt. Nach Konstruktion ist das gerade dann der Fall, wenn  $|t|$  möglichst klein ist, wenn also das Vorzeichen von  $a_{ii} - a_{jj}$  dem von  $\sigma |a_{ij}| \sqrt{\tau^2 + 1}$  entspricht. Oder kürzer: Wir müssen  $\sigma$  negativ wählen, falls  $\tau < 0$  gilt, und ansonsten positiv. Beispielsweise wird auf diese Weise verhindert, dass bei bereits kleinen Außerdiagonaleinträgen die Diagonaleinträge die Plätze tauschen und so zwar  $\operatorname{off}(\mathbf{B})^2$  konvergiert, nicht aber die Folge der Matrizen selbst.

**Bemerkung 4.8 (Quadratische Konvergenz)** *Wählt man die Vorzeichen wie oben erläutert, so kann man unter einigen zusätzlichen Voraussetzungen nachweisen, dass das Jacobi-Verfahren lokal quadratisch konvergiert.*

**Bemerkung 4.9 (Parallelisierung)** *Falls zwei Indexpaare  $(i, j)$  und  $(i', j')$  die Bedingung  $\{i, j\} \cap \{i', j'\} = \emptyset$  erfüllen, können die Berechnung der Zeilen zu beiden Rotationen parallel erfolgen: Die Rotation zu  $(i, j)$  betrifft nur die  $i$ -te und die  $j$ -te Zeile, während die Rotation zu  $(i', j')$  nur die  $i'$ -te und die  $j'$ -te Zeile betrifft. Entsprechend können wir auch mit den Spaltenrotationen verfahren.*

## 5 Die Vektoriteration

Die im letzten Kapitel vorgestellte Jacobi-Iteration berechnet im Falle einer symmetrischen Matrix sämtliche Eigenwerte und Eigenvektoren. Es gibt viele Fälle, in denen man lediglich an einen Teil des Spektrums und den zugehörigen Eigenvektoren interessiert ist, beispielsweise bei der Bestimmung von Eigenschwingungen oder eines invarianten Wahrscheinlichkeitsmaßes.

Zur Lösung derartiger Aufgaben werden sehr oft Verfahren auf der Basis der sogenannten *Vektoriteration* (engl. „power iteration“, weil Potenzen der Matrix eine entscheidende Rolle spielen) verwendet. Einigen dieser Methoden ist dieses Kapitel gewidmet.

### 5.1 Grundidee

Wir erinnern uns an das zweite Beispiel aus Kapitel 2, in dem die Aufgabe darin bestand, eine Wahrscheinlichkeitsverteilung  $\mathbf{y} \in \mathbb{R}_{\geq 0}^4 \setminus \{0\}$  zu finden, die stabil bleibt, sich also bei der „Durchführung eines Spielzugs“ nicht ändert. Ein derartiger Vektor ist durch die Gleichung

$$\mathbf{M}\mathbf{y} = \mathbf{y} \tag{5.1}$$

beschrieben. Ein Ansatz zur Bestimmung eines derartigen Vektors besteht darin, zu untersuchen, unter welchen Bedingungen die Folge  $(\mathbf{M}^m \mathbf{z})_{m \in \mathbb{N}}$  für eine Startverteilung  $\mathbf{z} \in \mathbb{R}_{\geq 0}^4$  konvergiert.

Falls wir annehmen, dass die Folge gegen einen Vektor  $\mathbf{z}^* \in \mathbb{R}^4$  konvergiert, finden wir zu jedem  $\epsilon \in \mathbb{R}_{>0}$  ein  $n_0 \in \mathbb{N}$  mit

$$\|\mathbf{M}^m \mathbf{z} - \mathbf{z}^*\| < \epsilon \quad \text{für alle } m \in \mathbb{N}_{\geq n_0},$$

so dass wir

$$\begin{aligned} \|\mathbf{M}\mathbf{z}^* - \mathbf{z}^*\| &= \|\mathbf{M}\mathbf{z}^* - \mathbf{M}^{m+1}\mathbf{z} + \mathbf{M}^{m+1}\mathbf{z} - \mathbf{z}^*\| \leq \|\mathbf{M}(\mathbf{z}^* - \mathbf{M}^m \mathbf{z})\| + \|\mathbf{M}^{m+1}\mathbf{z} - \mathbf{z}^*\| \\ &\leq \|\mathbf{M}\| \|\mathbf{z}^* - \mathbf{M}^m \mathbf{z}\| + \|\mathbf{M}^{m+1}\mathbf{z} - \mathbf{z}^*\| < (\|\mathbf{M}\| + 1)\epsilon \end{aligned}$$

für alle  $m \in \mathbb{N}_{\geq n_0}$  erhalten. Da diese Abschätzung für beliebige  $\epsilon \in \mathbb{R}_{>0}$  gilt, folgt

$$\mathbf{M}\mathbf{z}^* = \mathbf{z}^*,$$

also ist der Grenzwert der Folge, sofern er existiert, auch Lösung der Gleichung (5.1).

Unsere Aufgabe besteht also darin, für eine Matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und einen Startvektor  $\mathbf{z}^{(0)}$  die Folge  $(\mathbf{A}^m \mathbf{z}^{(0)})_{m \in \mathbb{N}_0}$  zu untersuchen.

## 5 Die Vektoriteration

Dazu untersuchen wir zunächst eine Diagonalmatrix

$$\mathbf{D} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix},$$

mit Eigenwerten  $\lambda_1, \dots, \lambda_n \in \mathbb{K}$ . Ohne Beschränkung der Allgemeinheit können wir annehmen, dass die Eigenwerte nach absteigendem Betrag sortiert sind, dass also

$$|\lambda_1| \geq |\lambda_2| \geq |\lambda_3| \geq \dots \geq |\lambda_n| \quad (5.2)$$

gilt. Wenn wir mit einem Vektor  $\hat{\mathbf{x}}^{(0)} \in \mathbb{K}^n$  anfangen, erhalten wir

$$\hat{\mathbf{x}}^{(m)} = \mathbf{D}^m \hat{\mathbf{x}}^{(0)} = \begin{pmatrix} \lambda_1^m \hat{x}_1^{(0)} \\ \vdots \\ \lambda_n^m \hat{x}_n^{(0)} \end{pmatrix} \quad \text{für alle } m \in \mathbb{N}_0.$$

Wir können sehen, dass die erste Komponente dieses Vektors schneller als alle anderen wachsen wird, falls  $|\lambda_1| > |\lambda_2|$  und  $\hat{x}_1^{(0)} \neq 0$  gelten.

Falls allerdings  $\lambda_1 \neq 1$  gilt, dürfen wir nicht erwarten, dass die Folge  $(\hat{\mathbf{x}}^{(m)})_{m=0}^\infty$  im konventionellen Sinn konvergiert, da asymptotisch  $\hat{\mathbf{x}}^{(m+1)} \approx \lambda_1 \hat{\mathbf{x}}^{(m)}$  gelten wird. Es ist deshalb hilfreich, lediglich die von den Vektoren aufgespannten Räume miteinander zu vergleichen statt die Vektoren selbst. Dafür ist der *Winkel* zwischen Vektoren nützlich.

**Definition 5.1 (Winkel)** Seien  $\mathbf{x}, \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$ . Der Winkel zwischen den Vektoren ist definiert durch

$$\angle(\mathbf{x}, \mathbf{y}) = \arccos \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|}{\|\mathbf{x}\| \|\mathbf{y}\|}.$$

Diese Definition hat den Vorteil, dass eine Skalierung der Vektoren  $\mathbf{x}$  und  $\mathbf{y}$  mit beliebigen von null verschiedenen Faktoren den Winkel nicht ändert, so dass wir auf Konvergenz hoffen dürfen.

In der Praxis ist es häufig handlicher, mit von dem Winkel abgeleiteten trigonometrischen Funktionen zu arbeiten, die durch

$$\begin{aligned} \cos \angle(\mathbf{x}, \mathbf{y}) &= \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|}{\|\mathbf{x}\| \|\mathbf{y}\|}, \\ \sin \angle(\mathbf{x}, \mathbf{y}) &:= \sqrt{1 - \cos^2 \angle(\mathbf{x}, \mathbf{y})}, \\ \tan \angle(\mathbf{x}, \mathbf{y}) &:= \frac{\sin \angle(\mathbf{x}, \mathbf{y})}{\cos \angle(\mathbf{x}, \mathbf{y})} \end{aligned} \quad \text{für alle } \mathbf{x}, \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$$

gegeben sind. Ein Blick auf die Taylor-Entwicklung zeigt, dass für kleine Winkel sowohl der Sinus als auch der Tangens ungefähr gleich dem Winkel sind, so dass Konvergenzaussagen über diese Funktionen auch Aussagen über den Winkel zulassen.

**Lemma 5.2 (Konvergenz für Diagonalmatrizen)** Sei  $\hat{\mathbf{x}}^{(0)} \in \mathbb{K}^n$  mit  $\hat{x}_1^{(0)} \neq 0$  gegeben. Dann gilt

$$\tan \angle(\delta^{(1)}, \hat{\mathbf{x}}^{(m)}) \leq \left( \frac{|\lambda_2|}{|\lambda_1|} \right)^m \tan \angle(\delta^{(1)}, \hat{\mathbf{x}}^{(0)}) \quad \text{für alle } m \in \mathbb{N}_0.$$

*Beweis.* Da  $\delta^{(1)} \in \mathbb{K}^n$  der erste kanonische Einheitsvektor ist, gelten

$$\begin{aligned} \cos^2 \angle(\delta^{(1)}, \hat{\mathbf{x}}^{(m)}) &= \frac{|\hat{x}_1^{(m)}|^2}{\|\hat{\mathbf{x}}^{(m)}\|^2} = \frac{|\hat{x}_1^{(m)}|^2}{\sum_{i=1}^n |\hat{x}_i^{(m)}|^2}, \\ \sin^2 \angle(\delta^{(1)}, \hat{\mathbf{x}}^{(m)}) &= 1 - \cos^2 \angle(\delta^{(1)}, \hat{\mathbf{x}}^{(m)}) = \frac{\sum_{i=1}^n |\hat{x}_i^{(m)}|^2 - |\hat{x}_1^{(m)}|^2}{\sum_{i=1}^n |\hat{x}_i^{(m)}|^2} = \frac{\sum_{i=2}^n |\hat{x}_i^{(m)}|^2}{\sum_{i=1}^n |\hat{x}_i^{(m)}|^2}, \\ \tan^2 \angle(\delta^{(1)}, \hat{\mathbf{x}}^{(m)}) &= \frac{\sin^2 \angle(\delta^{(1)}, \hat{\mathbf{x}}^{(m)})}{\cos^2 \angle(\delta^{(1)}, \hat{\mathbf{x}}^{(m)})} = \frac{\sum_{i=2}^n |\hat{x}_i^{(m)}|^2}{|\hat{x}_1^{(m)}|^2} \quad \text{für alle } m \in \mathbb{N}_0. \end{aligned}$$

Sei  $m \in \mathbb{N}_0$ . Mit (5.2) folgt

$$\begin{aligned} \tan^2 \angle(\delta^{(1)}, \hat{\mathbf{x}}^{(m)}) &= \frac{\sum_{i=2}^n |\hat{x}_i^{(m)}|^2}{|\hat{x}_1^{(m)}|^2} = \frac{\sum_{i=2}^n |\lambda_i^{2m}| |\hat{x}_i^{(0)}|^2}{|\lambda_1^{2m}| |\hat{x}_1^{(0)}|^2} \leq \frac{\sum_{i=2}^n |\lambda_2^{2m}| |\hat{x}_i^{(0)}|^2}{|\lambda_1^{2m}| |\hat{x}_1^{(0)}|^2} \\ &= \left( \frac{|\lambda_2|}{|\lambda_1|} \right)^{2m} \frac{\sum_{i=2}^n |\hat{x}_i^{(0)}|^2}{|\hat{x}_1^{(0)}|^2} = \left( \frac{|\lambda_2|}{|\lambda_1|} \right)^{2m} \tan^2 \angle(\delta^{(1)}, \hat{\mathbf{x}}^{(0)}). \end{aligned}$$

Die Behauptung folgt, indem wir auf beiden Seiten die Wurzel ziehen. ■

Auf Konvergenz des Tangens, und damit des Winkels, gegen null dürfen wir demnach hoffen, falls  $|\lambda_1| > |\lambda_2|$  gilt. Gemäß unserer Definition bedeutet das gerade, dass ein Eigenwert im Betrag echt größer ist als alle anderen.

**Definition 5.3 (Dominanter Eigenwert)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Ein Eigenwert  $\lambda_1 \in \sigma(\mathbf{A})$  heißt dominant, falls

$$|\lambda_1| > |\lambda| \quad \text{für alle } \lambda \in \sigma(\mathbf{A}) \setminus \{\lambda_1\}$$

gilt, falls der Betrag von  $\lambda_1$  also echt größer als die Beträge aller anderen Eigenwerte der Matrix  $\mathbf{A}$  ist (Es sei daran erinnert, dass das Spektrum grundsätzlich über dem Körper  $\mathbb{C}$  der komplexen Zahlen definiert wird, so dass  $\sigma(\mathbf{A})$  nicht leer sein kann).

Für die Praxis ist ein Verfahren, das sich nur auf Diagonalmatrizen anwenden lässt, natürlich uninteressant. Als erste Verallgemeinerung untersuchen wir eine normale Matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  (vgl. Definition 3.50). Nach Folgerung 3.54 existieren eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  und eine Diagonalmatrix  $\mathbf{D} \in \mathbb{C}^{n \times n}$  mit

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{D}.$$

## 5 Die Vektoriteration

Nach Bemerkung 3.42 können wir dafür sorgen, dass

$$\mathbf{D} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}, \quad |\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n| \quad (5.3)$$

gelten. Im Folgenden gehen wir davon aus, dass die Matrix  $\mathbf{D}$  diese Eigenschaft besitzt.

Den durch die Gleichung

$$\mathbf{x}^{(m)} = \mathbf{A}^m \mathbf{x}^{(0)} \quad \text{für alle } m \in \mathbb{N}_0$$

definierten Iterationsvektoren können wir nun transformierte Vektoren

$$\widehat{\mathbf{x}}^{(m)} := \mathbf{Q}^* \mathbf{x}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0$$

zuordnen. Dann gilt

$$\begin{aligned} \widehat{\mathbf{x}}^{(m)} &= \mathbf{Q}^* \mathbf{x}^{(m)} = \mathbf{Q}^* \mathbf{A}^m \mathbf{x}^{(0)} = \mathbf{Q}^* \mathbf{A}^m \mathbf{Q} \widehat{\mathbf{x}}^{(0)} \\ &= (\mathbf{Q}^* \mathbf{A} \mathbf{Q})^m \widehat{\mathbf{x}}^{(0)} = \mathbf{D}^m \widehat{\mathbf{x}}^{(0)} \quad \text{für alle } m \in \mathbb{N}_0. \end{aligned}$$

Damit können wir Lemma 5.2 auf die Vektoren  $(\widehat{\mathbf{x}}^{(m)})_{m=0}^{\infty}$  anwenden und müssen lediglich untersuchen, wie sich die Winkel unter unitären Transformationen verändern.

**Satz 5.4 (Konvergenz für normale Matrizen)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine normale Matrix. Sie besitzt eine Schur-Zerlegung  $\mathbf{A} = \mathbf{Q} \mathbf{D} \mathbf{Q}^*$  mit einer unitären Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und einer Diagonalmatrix  $\mathbf{D} \in \mathbb{C}^{n \times n}$  der Form (5.3).

Dann ist  $\mathbf{e} := \mathbf{Q} \delta^{(1)}$  ein Eigenvektor zu dem betragsgrößten Eigenwert  $\lambda_1$ .

Sei  $\mathbf{x}^{(0)} \in \mathbb{K}^n$  mit  $\langle \mathbf{e}, \mathbf{x}^{(0)} \rangle \neq 0$  gegeben. Dann gilt

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(m)}) \leq \left( \frac{|\lambda_2|}{|\lambda_1|} \right)^m \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \quad \text{für alle } m \in \mathbb{N}_0.$$

*Beweis.* Mit Lemma 3.17 gilt

$$\hat{x}_1^{(0)} = \langle \delta^{(1)}, \widehat{\mathbf{x}}^{(0)} \rangle = \langle \mathbf{Q}^* \mathbf{Q} \delta^{(1)}, \widehat{\mathbf{x}}^{(0)} \rangle = \langle \mathbf{Q} \delta^{(1)}, \mathbf{Q} \widehat{\mathbf{x}}^{(0)} \rangle = \langle \mathbf{e}, \mathbf{x}^{(0)} \rangle \neq 0.$$

Mit den Lemmas 3.17 und 3.34 erhalten wir

$$\begin{aligned} \cos \angle(\mathbf{e}, \mathbf{x}^{(m)}) &= \frac{|\langle \mathbf{e}, \mathbf{x}^{(m)} \rangle|}{\|\mathbf{e}\| \|\mathbf{x}^{(m)}\|} = \frac{|\langle \mathbf{Q} \delta^{(1)}, \mathbf{Q} \widehat{\mathbf{x}}^{(m)} \rangle|}{\|\mathbf{Q} \delta^{(1)}\| \|\mathbf{Q} \widehat{\mathbf{x}}^{(m)}\|} = \frac{|\langle \delta^{(1)}, \mathbf{Q}^* \mathbf{Q} \widehat{\mathbf{x}}^{(m)} \rangle|}{\|\delta^{(1)}\| \|\widehat{\mathbf{x}}^{(m)}\|} \\ &= \frac{|\langle \delta^{(1)}, \widehat{\mathbf{x}}^{(m)} \rangle|}{\|\delta^{(1)}\| \|\widehat{\mathbf{x}}^{(m)}\|} = \cos \angle(\delta^{(1)}, \widehat{\mathbf{x}}^{(m)}) \quad \text{für alle } m \in \mathbb{N}_0. \end{aligned}$$

Nach Definition folgt daraus unmittelbar

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(m)}) = \tan \angle(\delta^{(1)}, \widehat{\mathbf{x}}^{(m)}) \quad \text{für alle } m \in \mathbb{N}_0.$$



Nun können wir Lemma 5.2 anwenden und erhalten

$$\begin{aligned}\tan \angle(\mathbf{e}, \mathbf{x}^{(m)}) &= \tan \angle(\delta^{(1)}, \widehat{\mathbf{x}}^{(m)}) \leq \left( \frac{|\lambda_2|}{|\lambda_1|} \right)^m \tan \angle(\delta^{(1)}, \widehat{\mathbf{x}}^{(0)}) \\ &= \left( \frac{|\lambda_2|}{|\lambda_1|} \right)^m \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \quad \text{für alle } m \in \mathbb{N}_0.\end{aligned}$$

■

Falls wir Aussagen für Matrizen erhalten wollen, die nicht normal sind, können wir nicht länger auf unitäre Ähnlichkeitstransformationen zurückgreifen. In dieser Situation kann es nützlich sein, eine alternative Charakterisierung des Winkels zu verwenden.

**Lemma 5.5 (Sinus als Minimum)** *Es gilt*

$$\sin \angle(\mathbf{x}, \mathbf{y}) = \min \left\{ \frac{\|\mathbf{x} - \alpha \mathbf{y}\|}{\|\mathbf{x}\|} : \alpha \in \mathbb{K} \right\} \quad \text{für alle } \mathbf{x}, \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}.$$

*Beweis.* Seien  $\mathbf{x}, \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  gegeben. Wir setzen  $\beta := \langle \mathbf{y}, \mathbf{x} \rangle / \|\mathbf{y}\|^2$ . Sei  $\alpha \in \mathbb{K}$ . Es gilt

$$\begin{aligned}\|\mathbf{x} - \alpha \mathbf{y}\|^2 &= \|\mathbf{x} - \beta \mathbf{y} + (\beta - \alpha) \mathbf{y}\|^2 \\ &= \langle \mathbf{x} - \beta \mathbf{y} + (\beta - \alpha) \mathbf{y}, \mathbf{x} - \beta \mathbf{y} + (\beta - \alpha) \mathbf{y} \rangle \\ &= \|\mathbf{x} - \beta \mathbf{y}\|^2 + \overline{(\beta - \alpha)} \langle \mathbf{y}, \mathbf{x} - \beta \mathbf{y} \rangle + (\beta - \alpha) \langle \mathbf{x} - \beta \mathbf{y}, \mathbf{y} \rangle + |\beta - \alpha|^2 \|\mathbf{y}\|^2.\end{aligned}$$

Wir haben mit (3.6c)

$$\begin{aligned}\langle \mathbf{x} - \beta \mathbf{y}, \mathbf{y} \rangle &= \langle \mathbf{x}, \mathbf{y} \rangle - \bar{\beta} \|\mathbf{y}\|^2 = \langle \mathbf{x}, \mathbf{y} \rangle - \overline{\langle \mathbf{y}, \mathbf{x} \rangle} = 0, \\ \langle \mathbf{y}, \mathbf{x} - \beta \mathbf{y} \rangle &= \overline{\langle \mathbf{x} - \beta \mathbf{y}, \mathbf{y} \rangle} = 0,\end{aligned}$$

so dass wir

$$\|\mathbf{x} - \alpha \mathbf{y}\|^2 = \|\mathbf{x} - \beta \mathbf{y}\|^2 + |\beta - \alpha|^2 \|\mathbf{y}\|^2$$

erhalten. Offenbar nimmt dieser Ausdruck sein Minimum für  $\alpha = \beta$  an, und dieses Minimum ist gerade

$$\begin{aligned}\|\mathbf{x} - \beta \mathbf{y}\|^2 &= \|\mathbf{x}\|^2 - \bar{\beta} \langle \mathbf{y}, \mathbf{x} \rangle - \beta \langle \mathbf{x}, \mathbf{y} \rangle + |\beta|^2 \|\mathbf{y}\|^2 \\ &= \|\mathbf{x}\|^2 - \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|^2}{\|\mathbf{y}\|^2} - \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|^2}{\|\mathbf{y}\|^2} + \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|^2}{\|\mathbf{y}\|^4} \|\mathbf{y}\|^2 = \|\mathbf{x}\|^2 - \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|^2}{\|\mathbf{y}\|^2}.\end{aligned}$$

Es folgt

$$\sin^2 \angle(\mathbf{x}, \mathbf{y}) = 1 - \cos^2 \angle(\mathbf{x}, \mathbf{y}) = \frac{\|\mathbf{x}\|^2 \|\mathbf{y}\|^2 - |\langle \mathbf{x}, \mathbf{y} \rangle|^2}{\|\mathbf{x}\|^2 \|\mathbf{y}\|^2} = \frac{\|\mathbf{x} - \beta \mathbf{y}\|^2}{\|\mathbf{x}\|^2},$$

und damit die Behauptung. ■

Diese alternative Charakterisierung des Winkels ist nützlich, weil sie es uns erlaubt, auch allgemeine Transformationen in Betracht zu ziehen.

## 5 Die Vektoriteration

**Lemma 5.6 (Transformierter Sinus)** Sei  $\mathbf{B} \in \mathbb{K}^{n \times n}$  eine invertierbare Matrix. Es gilt

$$\sin \angle(\mathbf{B}\mathbf{x}, \mathbf{B}\mathbf{y}) \leq \|\mathbf{B}\| \|\mathbf{B}^{-1}\| \sin \angle(\mathbf{x}, \mathbf{y}) \quad \text{für alle } \mathbf{x}, \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}.$$

*Beweis.* Seien  $\mathbf{x}, \mathbf{y} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$ . Mit Lemma 5.5 finden wir ein  $\alpha \in \mathbb{K}$  so, dass

$$\sin \angle(\mathbf{x}, \mathbf{y}) = \frac{\|\mathbf{x} - \alpha\mathbf{y}\|}{\|\mathbf{x}\|}$$

gilt. Mit (3.8a) folgen

$$\begin{aligned} \|\mathbf{B}\mathbf{x} - \alpha\mathbf{B}\mathbf{y}\| &= \|\mathbf{B}(\mathbf{x} - \alpha\mathbf{y})\| \leq \|\mathbf{B}\| \|\mathbf{x} - \alpha\mathbf{y}\|, \\ \|\mathbf{B}^{-1}\| \|\mathbf{B}\mathbf{x}\| &\geq \|\mathbf{B}^{-1}\mathbf{B}\mathbf{x}\| = \|\mathbf{x}\|, \end{aligned}$$

und wir erhalten schließlich

$$\sin \angle(\mathbf{B}\mathbf{x}, \mathbf{B}\mathbf{y}) \leq \frac{\|\mathbf{B}\mathbf{x} - \alpha\mathbf{B}\mathbf{y}\|}{\|\mathbf{B}\mathbf{x}\|} \leq \|\mathbf{B}\| \|\mathbf{B}^{-1}\| \frac{\|\mathbf{x} - \alpha\mathbf{y}\|}{\|\mathbf{x}\|} = \|\mathbf{B}\| \|\mathbf{B}^{-1}\| \sin \angle(\mathbf{x}, \mathbf{y}).$$

■

Wenden wir uns also dem allgemeinen Fall zu: Für einen Startvektor  $\mathbf{x}^{(0)} \in \mathbb{K}^n$  und die allgemeine Matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  definieren wir die Folge der Iterierten durch

$$\mathbf{x}^{(m)} := \mathbf{A}^m \mathbf{x}^{(0)} = \mathbf{A} \mathbf{x}^{(m-1)} \quad \text{für alle } m \in \mathbb{N}. \quad (5.4)$$

Wir setzen voraus, dass  $\mathbf{A}$  diagonalisierbar ist und die Eigenwerte nach ihrem Betrag absteigend sortiert sind, dass also eine reguläre Matrix  $\mathbf{T} \in \mathbb{K}^{n \times n}$  mit

$$\mathbf{T}^{-1} \mathbf{A} \mathbf{T} = \mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$$

und der bereits aus (5.2) bekannten Anordnung

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

existiert. Mit Hilfe des Lemmas 5.6 können wir den Satz 5.4 wie folgt verallgemeinern:

**Folgerung 5.7 (Konvergenz)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  diagonalisierbar mit  $\mathbf{A} = \mathbf{T}\mathbf{D}\mathbf{T}^{-1}$  und  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  sowie (5.2). Gelte  $\gamma := \langle \delta^{(1)}, \mathbf{T}^{-1} \mathbf{x}^{(0)} \rangle \neq 0$ .

Dann ist  $\mathbf{e} := \mathbf{T}\delta^{(1)}$  ein Eigenvektor zu dem betragsgrößten Eigenwert  $\lambda_1$  und es gilt

$$\sin \angle(\mathbf{e}, \mathbf{x}^{(m)}) \leq \left( \frac{|\lambda_2|}{|\lambda_1|} \right)^m \|\mathbf{T}\| \|\mathbf{T}^{-1}\| \tan \angle(\delta^{(1)}, \mathbf{T}^{-1} \mathbf{x}^{(0)}) \quad \text{für alle } m \in \mathbb{N}_0.$$

*Beweis.* Wir definieren

$$\widehat{\mathbf{x}}^{(m)} := \mathbf{T}^{-1} \mathbf{x}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0$$

und stellen fest, dass

$$\widehat{\mathbf{x}}^{(m)} = \mathbf{T}^{-1} \mathbf{x}^{(m)} = \mathbf{T}^{-1} \mathbf{A}^m \mathbf{x}^{(0)} = \mathbf{T}^{-1} \mathbf{A}^m \mathbf{T} \widehat{\mathbf{x}}^{(0)} = (\mathbf{T}^{-1} \mathbf{A} \mathbf{T})^m \widehat{\mathbf{x}}^{(0)} = \mathbf{D}^m \widehat{\mathbf{x}}^{(0)}$$

für alle  $m \in \mathbb{N}_0$  gilt. Mit Lemma 5.6 erhalten wir

$$\sin \angle(\mathbf{e}, \mathbf{x}^{(m)}) = \sin \angle(\mathbf{T} \delta^{(1)}, \mathbf{T} \widehat{\mathbf{x}}^{(m)}) \leq \|\mathbf{T}\| \|\mathbf{T}^{-1}\| \sin \angle(\delta^{(1)}, \widehat{\mathbf{x}}^{(m)}).$$

Wegen  $\hat{x}_1^{(0)} = \langle \delta^{(1)}, \widehat{\mathbf{x}}^{(0)} \rangle = \langle \delta^{(1)}, \mathbf{T}^{-1} \mathbf{x}^{(0)} \rangle = \gamma \neq 0$  können wir wie bisher mit Lemma 5.2 abschätzen, um den Sinus abzuschätzen und die gewünschte Aussage zu erhalten. ■

In der Praxis kann die Bestimmung der Iterierten  $\mathbf{x}^{(m)}$  zu Schwierigkeiten führen, weil die Komponenten der Vektoren durch Maschinenzahlen dargestellt werden, die nicht beliebig groß oder klein werden können: Der Vektor wird asymptotisch in jedem Schritt ungefähr mit dem Faktor  $\lambda_1$  multipliziert werden. Falls  $|\lambda_1| > 1$  gilt, wird er also exponentiell wachsen, bis die Menge der Maschinenzahlen ausgeschöpft ist. Im IEEE-754-Standard würden dann die „übergelaufenen“ Koeffizienten gleich unendlich gesetzt werden und wären für unsere Zwecke unbrauchbar. Falls  $|\lambda_1| < 1$  gilt, werden die Vektoren exponentiell schrumpfen, bis sie so nahe an der Null sind, dass sie zu null abgerundet werden und damit ebenfalls unbrauchbar werden.

Um zu verhindern, dass die Koeffizienten der Iterationsvektoren die Menge der Maschinenzahlen verlassen, empfiehlt es sich, eine Normierung einzuführen, also die Vektoren im Zuge des Verfahrens so zu skalieren, dass Über- und Unterläufe ausgeschlossen werden.

Da gemäß Lemma 5.5 auch

$$\begin{aligned} \sin \angle(\beta \mathbf{x}, \mathbf{y}) &= \min \left\{ \frac{\|\beta \mathbf{x} - \gamma \mathbf{y}\|}{\|\beta \mathbf{x}\|} : \gamma \in \mathbb{K} \right\} = \min \left\{ \frac{\|\beta \mathbf{x} - \beta \gamma' \mathbf{y}\|}{\|\beta \mathbf{x}\|} : \gamma' \in \mathbb{K} \right\} \\ &= \min \left\{ \frac{\|\mathbf{x} - \gamma' \mathbf{y}\|}{\|\mathbf{x}\|} : \gamma' \in \mathbb{K} \right\} = \sin \angle(\mathbf{x}, \mathbf{y}) \end{aligned}$$

für alle  $\beta \in \mathbb{K} \setminus \{0\}$  gilt, beeinflusst eine beliebige Skalierung die Konvergenz der Vektoren nicht im Geringsten. Häufig wählt man die Skalierung so, dass die Iterierten Einheitsvektoren bezüglich einer geeigneten Norm sind. Der korrespondierende Algorithmus nimmt dann die folgende Form an:

**Algorithmus 5.8 (Vektoriteration)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und  $\mathbf{x}^{(0)} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  gegeben. Der folgende Algorithmus führt die Vektoriteration aus.

```

m ← 0
x(m) ← x(m) / ||x(m)||
while „Fehler zu groß“ do begin
  w(m+1) ← Ax(m)
  x(m+1) ← w(m+1) / ||w(m+1)||
  m ← m + 1
end

```

## 5 Die Vektoriteration

Selbstverständlich können wir den Algorithmus nicht unendlich lange arbeiten lassen, sondern wir müssen den Fehler der Approximation des Eigenvektors einschätzen und die Iteration abbrechen, sobald der Fehler klein genug geworden ist.

Eine einfache Strategie zur Schätzung des Fehlers der Vektoriteration besteht darin, in jedem Schritt des Verfahrens das *Residuum* der exakten Eigenwertgleichung  $\mathbf{Ax} = \lambda\mathbf{x}$  zu bestimmen, also ein Kriterium der Gestalt

$$\|\mathbf{Ax}^{(m)} - \lambda\mathbf{x}^{(m)}\| \leq \epsilon \|\mathbf{x}^{(m)}\|$$

zu verwenden, wobei  $\|\mathbf{x}^{(m)}\|$  auf der rechten Seite auftritt, um sicher zu stellen, dass das Kriterium invariant unter Skalierungen des Vektors  $\mathbf{x}^{(m)}$  ist.

Falls wir den Eigenwert  $\lambda$  kennen, ist dieses Kriterium sicherlich sinnvoll. Falls wir den Eigenwert hingegen nicht kennen, können wir ihn durch eine Näherung ersetzen: Falls  $\mathbf{x}$  ein Eigenvektor zu einem Eigenwert  $\lambda$  ist, bietet der in Definition 3.44 eingeführte *Rayleigh-Quotient*

$$\Lambda_A(\mathbf{x}) = \frac{\langle \mathbf{x}, \mathbf{Ax} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\},$$

die Möglichkeit, den Eigenwert  $\lambda$  aus dem Eigenvektor  $\mathbf{x}$  zu berechnen. Falls  $\mathbf{x}$  lediglich eine *Approximation* eines Eigenvektors ist, lässt sich beweisen, dass der Rayleigh-Quotient eine Approximation des Eigenwerts bietet.

**Lemma 5.9 (Eigenwert-Approximation)** *Sei  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  eine Näherung eines Eigenvektors  $\mathbf{e} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  der Matrix  $\mathbf{A}$  zu dem Eigenwert  $\lambda$ . Dann gilt*

$$|\Lambda_A(\mathbf{x}) - \lambda| \leq \|\mathbf{A} - \lambda\mathbf{I}\| \sin \angle(\mathbf{e}, \mathbf{x}) \leq \|\mathbf{A} - \lambda\mathbf{I}\| \frac{\|\mathbf{x} - \alpha\mathbf{e}\|}{\|\mathbf{x}\|} \quad \text{für alle } \alpha \in \mathbb{K}.$$

Falls  $\mathbf{A}$  eine normale Matrix ist, gilt sogar

$$|\Lambda_A(\mathbf{x}) - \lambda| \leq \|\mathbf{A} - \lambda\mathbf{I}\| \sin^2 \angle(\mathbf{e}, \mathbf{x}) \leq \|\mathbf{A} - \lambda\mathbf{I}\| \left( \frac{\|\mathbf{x} - \alpha\mathbf{e}\|}{\|\mathbf{x}\|} \right)^2 \quad \text{für alle } \alpha \in \mathbb{K}.$$

*Beweis.* Sei  $\alpha \in \mathbb{K}$ . Nach Definition gilt  $(\mathbf{A} - \lambda\mathbf{I})\mathbf{e} = \mathbf{0}$ , also folgt die Abschätzung

$$\begin{aligned} |\Lambda_A(\mathbf{x}) - \lambda| &= \left| \frac{\langle \mathbf{x}, \mathbf{Ax} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} - \frac{\langle \mathbf{x}, \lambda\mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \right| = \left| \frac{\langle \mathbf{x}, (\mathbf{A} - \lambda\mathbf{I})\mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \right| = \left| \frac{\langle \mathbf{x}, (\mathbf{A} - \lambda\mathbf{I})(\mathbf{x} - \alpha\mathbf{e}) \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \right| \\ &\leq \frac{\|\mathbf{x}\| \|(\mathbf{A} - \lambda\mathbf{I})(\mathbf{x} - \alpha\mathbf{e})\|}{\|\mathbf{x}\|^2} \leq \frac{\|\mathbf{A} - \lambda\mathbf{I}\| \|\mathbf{x} - \alpha\mathbf{e}\|}{\|\mathbf{x}\|} \end{aligned}$$

Da wir diese Ungleichung für jedes  $\alpha \in \mathbb{K}$  bewiesen haben, folgt aus Lemma 5.5 die erste Aussage. Sei nun  $\mathbf{A}$  eine normale Matrix. Nach Lemma 3.52 gilt

$$\|(\mathbf{A}^* - \bar{\lambda}\mathbf{I})\mathbf{e}\| = \|(\mathbf{A} - \lambda\mathbf{I})^*\mathbf{e}\| = \|(\mathbf{A} - \lambda\mathbf{I})\mathbf{e}\| = 0,$$

also insbesondere auch  $(\mathbf{A}^* - \bar{\lambda}\mathbf{I})\mathbf{e} = \mathbf{0}$ ,  $\mathbf{e}$  ist also ein Eigenvektor von  $\mathbf{A}^*$  zu dem Eigenwert  $\bar{\lambda}$ . Mit Hilfe dieser Gleichung können wir Lemma 3.17 verwenden, um unsere Abschätzung zu verbessern:

$$\begin{aligned} |\Lambda_A(\mathbf{x}) - \lambda| &= \left| \frac{\langle \mathbf{x}, (\mathbf{A} - \lambda\mathbf{I})(\mathbf{x} - \alpha\mathbf{e}) \rangle}{\|\mathbf{x}\|^2} \right| = \left| \frac{\langle (\mathbf{A} - \lambda\mathbf{I})^*\mathbf{x}, \mathbf{x} - \alpha\mathbf{e} \rangle}{\|\mathbf{x}\|^2} \right| \\ &= \left| \frac{\langle (\mathbf{A} - \lambda\mathbf{I})^*(\mathbf{x} - \alpha\mathbf{e}), \mathbf{x} - \alpha\mathbf{e} \rangle}{\|\mathbf{x}\|^2} \right| = \left| \frac{\langle \mathbf{x} - \alpha\mathbf{e}, (\mathbf{A} - \lambda\mathbf{I})(\mathbf{x} - \alpha\mathbf{e}) \rangle}{\|\mathbf{x}\|^2} \right| \\ &\leq \frac{\|\mathbf{x} - \alpha\mathbf{e}\| \|\mathbf{A} - \lambda\mathbf{I}\| \|\mathbf{x} - \alpha\mathbf{e}\|}{\|\mathbf{x}\|^2} = \|\mathbf{A} - \lambda\mathbf{I}\| \left( \frac{\|\mathbf{x} - \alpha\mathbf{e}\|}{\|\mathbf{x}\|} \right)^2, \end{aligned}$$

und da auch diese Abschätzung für beliebiges  $\alpha \in \mathbb{K}$  gezeigt wurde, folgt die zweite Fehlerabschätzung. ■

Es ist bemerkenswert, dass bei normalen Matrizen der mit Hilfe des Rayleigh-Quotienten berechnete Eigenwert wesentlich schneller als der Eigenvektor konvergieren kann. Unter gewissen Voraussetzungen lässt sich diese Eigenschaft verallgemeinern.

**Übungsaufgabe 5.10 (Verallgemeinerter Rayleigh-Quotient)** Sei  $n \in \mathbb{N}$ , und sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine Matrix. Wir definieren die Menge

$$\mathcal{W} := \{(\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{K}^n \times \mathbb{K}^n : \langle \mathbf{x}_1, \mathbf{x}_2 \rangle \neq 0\}$$

und den verallgemeinerten Rayleigh-Quotienten

$$\Lambda_A: \mathcal{W} \rightarrow \mathbb{K}, \quad (\mathbf{x}_1, \mathbf{x}_2) \mapsto \frac{\langle \mathbf{x}_1, \mathbf{A}\mathbf{x}_2 \rangle}{\langle \mathbf{x}_1, \mathbf{x}_2 \rangle}.$$

Sei  $\lambda \in \mathbb{K}$  ein Eigenwert der Matrix  $\mathbf{A}$ . Dann ist  $\bar{\lambda}$  ein Eigenwert der Matrix  $\mathbf{A}^*$ .

- (a) Sei  $\mathbf{e}_2 \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  ein Eigenvektor der Matrix  $\mathbf{A}$  zu dem Eigenwert  $\lambda$ . Beweisen Sie

$$\Lambda_A(\mathbf{x}_1, \mathbf{e}_2) = \lambda \quad \text{für alle } \mathbf{x}_1 \in \mathbb{K}^n \text{ mit } \langle \mathbf{x}_1, \mathbf{e}_2 \rangle \neq 0.$$

- (b) Sei  $\mathbf{e}_1 \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  ein Eigenvektor der Matrix  $\mathbf{A}^*$  zu dem Eigenwert  $\bar{\lambda}$ . Beweisen Sie

$$\Lambda_A(\mathbf{e}_1, \mathbf{x}_2) = \bar{\lambda} \quad \text{für alle } \mathbf{x}_2 \in \mathbb{K}^n \text{ mit } \langle \mathbf{e}_1, \mathbf{x}_2 \rangle \neq 0.$$

- (c) Seien  $\mathbf{e}_1, \mathbf{e}_2 \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  Eigenvektoren der Matrizen  $\mathbf{A}^*$  und  $\mathbf{A}$  zu den Eigenwerten  $\bar{\lambda}$  und  $\lambda$ . Beweisen Sie

$$|\lambda - \Lambda_A(\mathbf{x}_1, \mathbf{x}_2)| \leq \frac{\|\lambda\mathbf{I} - \mathbf{A}\|}{\cos \angle(\mathbf{x}_1, \mathbf{x}_2)} \sin \angle(\mathbf{x}_1, \mathbf{e}_1) \sin \angle(\mathbf{x}_2, \mathbf{e}_2) \quad \text{für alle } (\mathbf{x}_1, \mathbf{x}_2) \in \mathcal{W}.$$

- (d) Beweisen oder widerlegen Sie, dass die in Teil (c) eingeführten Eigenvektoren immer  $\cos \angle(\mathbf{e}_1, \mathbf{e}_2) > 0$  erfüllen.

## 5 Die Vektoriteration

Mit der durch den Rayleigh-Quotienten zur Verfügung gestellten Näherung

$$\lambda^{(m)} := \Lambda_A(\mathbf{x}^{(m)}) \quad \text{für alle } m \in \mathbb{N}_0$$

können wir

$$\epsilon^{(m)} := \|\mathbf{Ax}^{(m)} - \lambda^{(m)}\mathbf{x}^{(m)}\| \quad \text{für alle } m \in \mathbb{N}_0 \quad (5.5)$$

als Maß des Fehlers verwenden.

In Algorithmus 5.8 steht uns der Vektor  $\mathbf{w}^{(m+1)} = \mathbf{Ax}^{(m)}$  zur Verfügung und  $\mathbf{x}^{(m)}$  ist ein Einheitsvektor, so dass der Rayleigh-Quotient sich in der Form

$$\lambda^{(m)} = \Lambda_A(\mathbf{x}^{(m)}) = \frac{\langle \mathbf{x}^{(m)}, \mathbf{Ax}^{(m)} \rangle}{\langle \mathbf{x}^{(m)}, \mathbf{x}^{(m)} \rangle} = \langle \mathbf{x}^{(m)}, \mathbf{w}^{(m+1)} \rangle \quad \text{für alle } m \in \mathbb{N}_0$$

darstellen lässt, die ohne weitere Matrix-Vektor-Multiplikationen ausgewertet werden kann.

Der Vektor  $\mathbf{w}^{(m+1)}$  lässt sich auch benutzen, um die Berechnung der Größe  $\epsilon^{(m)}$  zu beschleunigen, denn es gilt

$$\epsilon^{(m)} := \|\mathbf{Ax}^{(m)} - \lambda^{(m)}\mathbf{x}^{(m)}\| = \|\mathbf{w}^{(m+1)} - \lambda^{(m)}\mathbf{x}^{(m)}\| \quad \text{für alle } m \in \mathbb{N}_0,$$

auch hier ist also keine weitere Matrix-Vektor-Multiplikation erforderlich.

Da die im Computer verwendeten Maschinenzahlen eine gewisse *relative* Genauigkeit garantieren, bietet es sich an, Algorithmen auch so zu konstruieren, dass sie diese Eigenschaft erhalten. In unserem Fall bedeutet das, dass das Abbruchkriterium auch von der Skalierung der Matrix unabhängig sein sollte. Da eine Skalierung der Matrix eine Skalierung der Eigenwerte zur Folge hat, leistet die Bedingung  $\epsilon^{(m)} \leq \epsilon|\lambda^{(m)}| \|\mathbf{x}^{(m)}\|$  das Geforderte.

**Algorithmus 5.11 (Vektoriteration mit Abbruchkriterium)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und  $\mathbf{x}^{(0)} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  gegeben. Der folgende Algorithmus führt die Vektoriteration aus und prüft in jedem Schritt, ob das Abbruchkriterium erfüllt ist.

```

m ← 0
x(m) ← x(m) / ||x(m)||
w(m+1) ← Ax(m)
λ(m) ← ⟨x(m), w(m+1)⟩
while ||w(m+1) - λ(m)x(m)|| > ε|λ(m)| do begin
    x(m+1) ← w(m+1) / ||w(m+1)||
    m ← m + 1
    w(m+1) ← Ax(m)
    λ(m) ← ⟨x(m), w(m+1)⟩
end

```

Nach dem Ende der Iteration erfüllen die berechneten Näherungen des Eigenvektors  $\mathbf{x}^{(m)}$  und des Eigenwerts  $\lambda^{(m)}$  die Abschätzung  $\|\mathbf{Ax}^{(m)} - \lambda^{(m)}\mathbf{x}^{(m)}\| \leq \epsilon|\lambda| \|\mathbf{x}^{(m)}\|$ .

**Bemerkung 5.12 (Implementierung)** Für den Algorithmus 5.8 ist es lediglich erforderlich, dass sich die Matrix  $\mathbf{A}$  mit einem Vektor multiplizieren lässt. Diese Eigenschaft ist sehr wichtig, da es sehr viele Fälle gibt, in denen das Auswerten einer Matrix relativ kostengünstig durchzuführen ist, etwa bei Matrizen mit vielen Nulleinträgen (z.B. bei Bandmatrizen wie der Tridiagonalmatrix aus Abschnitt 2.1) oder bei Matrizen, die implizit als Lösungsoperator eines linearen Gleichungssystems gegeben sind.

**Bemerkung 5.13 (Iteration im Teilraum)** Falls  $\mathbf{A}$  eine normale Matrix ist, stehen ihre Eigenvektoren senkrecht aufeinander. Nehmen wir an, dass ein Eigenvektor  $\mathbf{e}^{(1)}$  zu dem Eigenwert  $\lambda_1$  berechnet wurde. Dann können wir einen Anfangsvektor  $\mathbf{x}^{(0)}$  wählen, der senkrecht auf dem bereits berechneten Eigenvektor steht. Dann läuft die gesamte Vektoriteration in dem orthogonalen Komplement dieses Eigenvektors ab.

Wenn wir  $\mathbf{A}$  als Abbildung dieses invarianten Teilraums in sich interpretieren, ist ihr betragsgrößter Eigenwert nun  $\lambda_2$ , und falls  $|\lambda_2| > |\lambda_3|$  gilt, wird die Vektoriteration gegen einen Eigenvektor zu dem Eigenwert  $\lambda_2$  konvergieren.

Falls  $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$  gilt, können wir in dieser Weise nach und nach Eigenvektoren zu allen Eigenwerten berechnen und die Matrix diagonalisieren.

Wir haben bereits in Übungsaufgabe 3.49 gesehen, dass wir beliebige selbstadjungierte Matrizen mit Hilfe einiger Householder-Spiegelungen auf Tridiagonalgestalt bringen können. Deshalb lohnt es sich, diese Klasse von Matrizen etwas genauer zu untersuchen.

Das in Abschnitt 2.1 betrachtete Modellproblem führt beispielsweise direkt zu selbstadjungierten Tridiagonalmatrizen mit konstanten Diagonal- und Nebendiagonaleinträgen, für die sich die Eigenwerte und Eigenvektoren explizit angeben lassen.

**Übungsaufgabe 5.14 (1D-Modellproblem)** Seien  $n \in \mathbb{N}$ ,  $d \in \mathbb{R}$  und  $\ell \in \mathbb{K}$  gegeben. Wir untersuchen die Matrix

$$\mathbf{A} = \begin{pmatrix} d & \bar{\ell} & & \\ \ell & \ddots & \ddots & \\ & \ddots & \ddots & \bar{\ell} \\ & & \ell & d \end{pmatrix}.$$

Wir setzen  $z := -\operatorname{sgn}(\ell)$  und  $h = \frac{1}{n+1}$ .

(a) Sei  $\nu \in [1 : n]$ . Beweisen Sie, dass der durch

$$e_i := z^i \sin(\pi \nu i h) \quad \text{für alle } i \in [1 : n]$$

definierte Vektor  $\mathbf{e} \in \mathbb{K}^n$  ein Eigenvektor der Matrix  $\mathbf{A}$  ist mit dem Eigenwert

$$\lambda = d - 2|\ell| \cos(\pi \nu h).$$

(b) Sei  $\nu \in [1 : n]$ ,  $d = 2h^{-2}$  und  $\ell = -h^{-2}$ . Zeigen Sie, dass der in Teil (a) konstruierte Eigenwert die Gleichung

$$\lambda = 4h^{-2} \sin^2(\pi \nu h/2)$$

erfüllt.

## 5 Die Vektoriteration

Hinweis: Das Additionstheorem  $\sin(\alpha + \beta) = \sin(\alpha)\cos(\beta) + \cos(\alpha)\sin(\beta)$  und die Gleichung  $\cos(2\alpha) = 1 - 2\sin^2(\alpha)$  könnten hilfreich sein.

Tridiagonalmatrizen mit konstanten Einträgen auf der Diagonalen und den Nebendiagonalen lassen sich mit einfachen diagonalen Ähnlichkeitstransformationen in selbstadjungierte Tridiagonalmatrizen überführen. Damit eröffnet sich die Möglichkeit, die für derartige Matrizen entwickelten Algorithmen auch in allgemeineren Fällen einzusetzen.

**Übungsaufgabe 5.15 (Tridiagonale Toeplitz-Matrix)** Sei  $n \in \mathbb{N}$ , seien  $\ell, d, r \in \mathbb{C}$  mit  $\ell r \in \mathbb{R}_{>0}$  gegeben. Wir betrachten die Tridiagonalmatrix

$$\mathbf{A} = \begin{pmatrix} d & r & & \\ \ell & \ddots & \ddots & \\ & \ddots & \ddots & r \\ & & \ell & d \end{pmatrix} \in \mathbb{C}^{n \times n}.$$

Beweisen Sie, dass es eine invertierbare Diagonalmatrix  $\mathbf{D} \in \mathbb{C}^{n \times n}$  und ein  $b \in \mathbb{C}$  gibt mit

$$\mathbf{D}^{-1}\mathbf{A}\mathbf{D} = \begin{pmatrix} d & \bar{b} & & \\ b & d & \ddots & \\ & \ddots & \ddots & \bar{b} \\ & & b & d \end{pmatrix}.$$

Unter Zuhilfenahme der in Abschnitt 3.8 entwickelten Theorie können unsere Konvergenzsätze für die Vektoriteration auf den Fall nicht-diagonalisierbarer Matrizen ausgedehnt werden, solange der dominante Eigenwert einfach ist.

**Übungsaufgabe 5.16 (Allgemeinerer Konvergenzsatz)** Sei  $n \in \mathbb{N}$ ,  $n > 1$ , und sei  $\mathbf{A} \in \mathbb{C}^{n \times n}$  eine Matrix mit einem dominanten Eigenwert  $\lambda_1 \in \mathbb{C}$  der algebraischen Vielfachheit  $n_1 = 1$ .

Mit Satz 3.61 finden wir eine reguläre Matrix  $\mathbf{B} \in \mathbb{C}^{n \times n}$ , Zahlen  $n_2, \dots, n_s \in \mathbb{N}$  für  $s = |\sigma(\mathbf{A})| > 1$ , und obere Dreiecksmatrizen mit

$$\mathbf{R}_i \in \mathbb{C}^{n_i \times n_i}, \quad \sigma(\mathbf{R}_i) = \{\lambda_i\} \quad \text{für alle } i \in [1 : s]$$

und

$$\mathbf{B}^{-1}\mathbf{A}\mathbf{B} = \begin{pmatrix} \mathbf{R}_1 & & \\ & \ddots & \\ & & \mathbf{R}_s \end{pmatrix}.$$

Ohne Beschränkung der Allgemeinheit gelte  $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_s|$ .

Die (theoretische) Vektoriteration zu einem Startvektor  $\mathbf{x}^{(0)} \in \mathbb{K}^n$  ist wie zuvor durch  $\mathbf{x}^{(m)} = \mathbf{A}^m \mathbf{x}^{(0)}$  für alle  $m \in \mathbb{N}$  definiert.



(a) Wir definieren die Abbildung

$$\text{nil}: \mathbb{K}^{n \times n} \rightarrow \mathbb{Z},$$

$$\mathbf{A} \mapsto \max\{\alpha \in \mathbb{Z} : a_{ij} = 0 \text{ für alle } i, j \in [1 : n] \text{ mit } i > j - \alpha\},$$

die beschreibt, wie viele obere Nebendiagonalen einer Matrix gleich null sind.

Seien  $\mathbf{A}, \mathbf{B} \in \mathbb{K}^{n \times n}$  mit  $\text{nil}(\mathbf{A}), \text{nil}(\mathbf{B}) \geq 0$  gegeben.

Beweisen Sie  $\text{nil}(\mathbf{AB}) \geq \text{nil}(\mathbf{A}) + \text{nil}(\mathbf{B})$ .

(b) Folgern Sie mit Teil (a), dass für eine Matrix  $\mathbf{N} \in \mathbb{K}^{n \times n}$  mit  $\text{nil}(\mathbf{N}) > 0$  auch  $\mathbf{N}^n = \mathbf{0}$  gilt, sie also nilpotent ist.

(c) Sei  $i \in [1 : s]$ . Beweisen Sie, dass  $\mathbf{R}_i = \lambda_i \mathbf{I} + \mathbf{N}_i$  mit einer Matrix  $\mathbf{N}_i \in \mathbb{C}^{n_i \times n_i}$  mit  $\text{nil}(\mathbf{N}_i) > 0$  gilt.

Folgern Sie daraus mit Teil (b), dass ein Polynom  $p_i \in \Pi_{n_i}$  existiert mit

$$\|\mathbf{R}_i^m\| \leq |\lambda_i|^{m-n_i} p_i(m) \quad \text{für alle } m \in \mathbb{N}_{\geq n_i}.$$

(d) Gelte  $(\mathbf{B}^{-1} \mathbf{x}^{(0)})_1 \neq 0$ .  $\mathbf{e} := \mathbf{B} \delta^{(1)}$  ist ein Eigenvektor zu dem Eigenwert  $\lambda_1$ .

Sei  $k := \max\{n_2, \dots, n_s\}$ . Zeigen Sie, dass ein Polynom  $p \in \Pi_k$  existiert mit

$$\sin \angle(\mathbf{x}^{(m)}, \mathbf{e}) \leq \left( \frac{|\lambda_2|}{|\lambda_1|} \right)^{m-k} p(m) \quad \text{für alle } m \in \mathbb{N}_{\geq k}.$$

Hinweis: Für Teil (c) kann der binomische Satz hilfreich sein. Für Teil (d) kann mit Hilfe von Teil (c) ähnlich wie bei Lemma 5.2 verfahren werden.

## 5.2 Fehleranalyse

Der Algorithmus 5.11 endet, sobald die Norm des *Residuums*

$$\mathbf{r}_A(\mathbf{x}) := \Lambda_A(\mathbf{x})\mathbf{x} - \mathbf{A}\mathbf{x}$$

den Wert  $\epsilon|\lambda|$  unterschreitet. Es stellt sich die Frage, ob diese Eigenschaft bereits bedeutet, dass wir gute Näherungen des Eigenwerts und des Eigenvektors berechnet haben.

Wir beschränken uns bei der Untersuchung auf den Fall einer selbstadjungierten Matrix  $\mathbf{A} = \mathbf{A}^* \in \mathbb{K}^{n \times n}$ .

Bei der Analyse des Fehlers verwenden wir ein allgemeines Prinzip, das auch in anderen Gebieten der numerischen Mathematik von Bedeutung ist: Die Idee der *Rückwärtsanalyse* beruht darauf, zu einer Näherungslösung eines Problems ein zweites Problem zu konstruieren, das durch die Näherungslösung exakt gelöst wird. Falls Aussagen darüber zur Verfügung stehen, wie sich die Lösungen eines Problems unter Störungen der Problemstellung verändern, kann man dann Rückschlüsse auf die Genauigkeit der Näherungslösung ziehen<sup>1</sup>.

<sup>1</sup>Für die Anregung zu diesem Beweis bedanke ich mich bei Prof. Dr. Volker Mehrmann.

## 5 Die Vektoriteration

In unserem Fall suchen wir einen Eigenvektor  $\mathbf{e}$ , also eine Lösung der Gleichung

$$\mathbf{A}\mathbf{e} = \lambda\mathbf{e}.$$

Uns steht ein genäherter Eigenvektor  $\mathbf{x}$  zur Verfügung, zu dem  $\tilde{\lambda} = \Lambda_A(\mathbf{x})$  eine Näherung des Eigenwerts darstellt. Wir suchen eine Matrix  $\tilde{\mathbf{A}}$  derart, dass

$$\tilde{\mathbf{A}}\mathbf{x} = \tilde{\lambda}\mathbf{x}$$

gilt. Wenn wir das Residuum mit

$$\mathbf{r} := \Lambda_A(\mathbf{x})\mathbf{x} - \mathbf{A}\mathbf{x} = \tilde{\lambda}\mathbf{x} - \mathbf{A}\mathbf{x}$$

bezeichnen, erhalten wir wegen

$$\langle \mathbf{x}, \mathbf{r} \rangle = \Lambda_A(\mathbf{x})\langle \mathbf{x}, \mathbf{x} \rangle - \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle = \frac{\langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} \langle \mathbf{x}, \mathbf{x} \rangle - \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle = 0 \quad (5.6)$$

für die Matrix

$$\tilde{\mathbf{A}} := \mathbf{A} + \frac{\mathbf{x}\mathbf{r}^*}{\|\mathbf{x}\|^2} + \frac{\mathbf{r}\mathbf{x}^*}{\|\mathbf{x}\|^2}$$

die Gleichung

$$\tilde{\mathbf{A}}\mathbf{x} = \mathbf{A}\mathbf{x} + \mathbf{x} \frac{\langle \mathbf{r}, \mathbf{x} \rangle}{\|\mathbf{x}\|^2} + \mathbf{r} \frac{\langle \mathbf{x}, \mathbf{x} \rangle}{\|\mathbf{x}\|^2} = \mathbf{A}\mathbf{x} + \mathbf{r} = \mathbf{A}\mathbf{x} + \tilde{\lambda}\mathbf{x} - \mathbf{A}\mathbf{x} = \tilde{\lambda}\mathbf{x},$$

also erfüllt  $\tilde{\mathbf{A}}$  unsere Anforderungen.

Die Spektralnorm der Störung lässt sich besonders elegant darstellen:

**Lemma 5.17 (Rang-2-Matrix)** Seien  $\mathbf{a}, \mathbf{b} \in \mathbb{K}^n$  mit  $\langle \mathbf{a}, \mathbf{b} \rangle = 0$  gegeben. Sei  $\mathbf{E} := \mathbf{a}\mathbf{b}^* + \mathbf{b}\mathbf{a}^*$ . Dann gilt  $\|\mathbf{E}\| = \|\mathbf{a}\| \|\mathbf{b}\|$ .

*Beweis.* Wir untersuchen zunächst den Sonderfall  $\|\mathbf{a}\| = 1 = \|\mathbf{b}\|$ .

Sei  $\mathbf{x} \in \mathbb{K}^n$ . Wir zerlegen den Vektor in Anteile aus dem Aufspann der Vektoren  $\mathbf{a}$  und  $\mathbf{b}$  sowie einen Rest, der senkrecht auf beiden steht:

$$\alpha := \langle \mathbf{a}, \mathbf{x} \rangle, \quad \beta := \langle \mathbf{b}, \mathbf{x} \rangle, \quad \mathbf{x}_0 := \mathbf{x} - \alpha\mathbf{a} - \beta\mathbf{b}.$$

Dann erhalten wir

$$\begin{aligned} \langle \mathbf{a}, \mathbf{x}_0 \rangle &= \langle \mathbf{a}, \mathbf{x} \rangle - \alpha \langle \mathbf{a}, \mathbf{a} \rangle - \beta \langle \mathbf{a}, \mathbf{b} \rangle = \alpha - \alpha = 0, \\ \langle \mathbf{b}, \mathbf{x}_0 \rangle &= \langle \mathbf{b}, \mathbf{x} \rangle - \alpha \langle \mathbf{b}, \mathbf{a} \rangle - \beta \langle \mathbf{b}, \mathbf{b} \rangle = \beta - \beta = 0, \\ \mathbf{E}\mathbf{x}_0 &= \mathbf{a}\langle \mathbf{b}, \mathbf{x}_0 \rangle + \mathbf{b}\langle \mathbf{a}, \mathbf{x}_0 \rangle = 0, \end{aligned}$$

also liegt  $\mathbf{x}_0$  im Kern der Matrix  $\mathbf{E}$ . Es folgt

$$\mathbf{E}\mathbf{x} = \mathbf{E}(\mathbf{x}_0 + \alpha\mathbf{a} + \beta\mathbf{b}) = \alpha\mathbf{a}\langle \mathbf{b}, \mathbf{a} \rangle + \alpha\mathbf{b}\langle \mathbf{a}, \mathbf{a} \rangle + \beta\mathbf{a}\langle \mathbf{b}, \mathbf{b} \rangle + \beta\mathbf{b}\langle \mathbf{a}, \mathbf{b} \rangle = \alpha\mathbf{b} + \beta\mathbf{a},$$

und infolge der Orthogonalität der Vektoren  $\mathbf{x}_0$ ,  $\mathbf{a}$  und  $\mathbf{b}$  ergibt sich

$$\begin{aligned}\|\mathbf{E}\mathbf{x}\|^2 &= \|\alpha\mathbf{b} + \beta\mathbf{a}\|^2 = |\alpha|^2\|\mathbf{b}\|^2 + |\beta|^2\|\mathbf{a}\|^2 = |\alpha|^2 + |\beta|^2 \leq \|\mathbf{x}_0\|^2 + |\alpha|^2 + |\beta|^2 \\ &= \|\mathbf{x}_0\|^2 + \|\alpha\mathbf{a}\|^2 + \|\beta\mathbf{b}\|^2 = \|\mathbf{x}_0 + \alpha\mathbf{a} + \beta\mathbf{b}\|^2 = \|\mathbf{x}\|^2,\end{aligned}$$

also  $\|\mathbf{E}\| \leq 1 = \|\mathbf{a}\| \|\mathbf{b}\|$  nach Definition der Spektralnrm. Wegen

$$\|\mathbf{E}\mathbf{a}\| = \|\mathbf{b}\langle\mathbf{a}, \mathbf{a}\rangle + \mathbf{a}\langle\mathbf{b}, \mathbf{a}\rangle\| = \|\mathbf{b}\| = 1$$

muss auch  $\|\mathbf{E}\| \geq 1$  gelten, so dass wir  $\|\mathbf{E}\| = 1$  bewiesen haben.

Widmen wir uns nun dem allgemeinen Fall. Sollte  $\mathbf{a} = \mathbf{0}$  oder  $\mathbf{b} = \mathbf{0}$  gelten, so folgt  $\mathbf{E} = \mathbf{0}$  und damit die Behauptung.

Anderenfalls setzen wir  $\hat{\mathbf{a}} := \mathbf{a}/\|\mathbf{a}\|$  und  $\hat{\mathbf{b}} := \mathbf{b}/\|\mathbf{b}\|$  und wenden den bereits bewiesenen Sonderfall auf  $\hat{\mathbf{E}} := \hat{\mathbf{a}}\hat{\mathbf{b}}^* + \hat{\mathbf{b}}\hat{\mathbf{a}}^*$  an, um  $\|\hat{\mathbf{E}}\| = 1$  zu erhalten. Mit  $\|\mathbf{E}\| = \|\mathbf{a}\| \|\mathbf{b}\| \|\hat{\mathbf{E}}\|$  folgt die Behauptung. ■

Die Anwendung dieses Lemmas auf unseren Fall führt zu

$$\|\mathbf{A} - \tilde{\mathbf{A}}\| = \frac{\|\mathbf{x}\| \|\mathbf{r}\|}{\|\mathbf{x}\|^2} = \frac{\|\mathbf{r}\|}{\|\mathbf{x}\|}, \quad (5.7)$$

die Norm der Störung der Matrix ist also gleich der relativen Norm des Residuums, die wir in unserem Algorithmus explizit berechnen können.

Unser Ziel ist es, aus dieser Abschätzung Rückschlüsse auf die Genauigkeit der Approximation des Eigenwerts und des Eigenvektors zu gewinnen. Als Hilfsmittel verwenden wir eine Variante des Satzes 3.45 von Courant und Fischer:

**Folgerung 5.18 (Rayleigh-Minimum)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix. Der Rayleigh-Quotient  $\Lambda_A$  besitzt in  $\mathbb{K}^n \setminus \{\mathbf{0}\}$  ein Minimum, und dieses Minimum ist der kleinste Eigenwert der Matrix  $\mathbf{A}$ .

*Beweis.* Es gilt

$$\Lambda_{-\mathbf{A}}(\mathbf{x}) = \frac{\langle \mathbf{x}, -\mathbf{A}\mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} = -\frac{\langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} = -\Lambda_A(\mathbf{x}) \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}.$$

Da die Matrix  $-\mathbf{A}$  selbstadjungiert ist, garantiert der Satz 3.45 von Courant und Fischer, dass der zugehörige Rayleigh-Quotient  $\Lambda_{-\mathbf{A}}$  ein Maximum annimmt, das dem größten Eigenwert der Matrix entspricht.

Also nimmt  $\Lambda_A = -\Lambda_{-\mathbf{A}}$  an derselben Stelle ein Minimum an, das dem kleinsten Eigenwert der Matrix  $\mathbf{A}$  entspricht. ■

Neben seiner Bedeutung für den folgenden Störungssatz kann das Courant-Fischer-Minimierungsprinzip auch als Ausgangspunkt bei der Konstruktion numerischer Näherungsverfahren dienen: Statt unmittelbar nach einem Eigenvektor zu suchen, können wir uns auch um ein Minimum des Rayleigh-Quotienten bemühen.

Zunächst formulieren wir allerdings den *Bauer-Fike-Störungssatz* für selbstadjungierte Matrizen, der Aussagen über die Konvergenz der Eigenwerte ermöglicht.

## 5 Die Vektoriteration

**Satz 5.19 (Bauer-Fike)** Seien  $\mathbf{A}, \tilde{\mathbf{A}} \in \mathbb{K}^{n \times n}$  selbstadjungierte Matrizen und sei  $\tilde{\lambda} \in \sigma(\tilde{\mathbf{A}})$ . Dann existiert ein Eigenwert  $\lambda \in \sigma(\mathbf{A})$  mit

$$|\lambda - \tilde{\lambda}| \leq \|\mathbf{A} - \tilde{\mathbf{A}}\|.$$

*Beweis.* Wir wählen ein  $\lambda \in \sigma(\mathbf{A})$  mit minimalem Abstand zu  $\tilde{\lambda}$ , es soll also

$$|\lambda - \tilde{\lambda}| \leq |\mu - \tilde{\lambda}| \quad \text{für alle } \mu \in \sigma(\mathbf{A})$$

gelten. Mit Satz 3.47 finden wir eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  und eine reelle Diagonalmatrix  $\mathbf{D} \in \mathbb{R}^{n \times n}$  mit  $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^*$ . Es folgt

$$(\mathbf{A} - \tilde{\lambda}\mathbf{I})^2 = (\mathbf{Q}\mathbf{D}\mathbf{Q}^* - \tilde{\lambda}\mathbf{Q}\mathbf{Q}^*)^2 = \mathbf{Q}(\mathbf{D} - \tilde{\lambda}\mathbf{I})^2\mathbf{Q}^*,$$

und da die Eigenwerte der Matrix  $\mathbf{A}$  die Diagonalelemente der Matrix  $\mathbf{D}$  sind, muss  $(\lambda - \tilde{\lambda})^2$  der kleinste Eigenwert der Matrix  $(\mathbf{A} - \tilde{\lambda}\mathbf{I})^2$  sein.

Sei  $\tilde{\mathbf{e}} \in \mathbb{K}^n$  ein Eigenvektor der Matrix  $\tilde{\mathbf{A}}$  zu dem Eigenwert  $\tilde{\lambda}$  mit  $\|\tilde{\mathbf{e}}\| = 1$ . Dann gilt nach Satz 5.18 und (3.8a) die Abschätzung

$$\begin{aligned} (\lambda - \tilde{\lambda})^2 &= \min\{\langle (\mathbf{A} - \tilde{\lambda}\mathbf{I})^2 \mathbf{z}, \mathbf{z} \rangle : \mathbf{z} \in \mathbb{K}^n, \|\mathbf{z}\| = 1\} \\ &\leq \langle (\mathbf{A} - \tilde{\lambda}\mathbf{I})^2 \tilde{\mathbf{x}}, \tilde{\mathbf{x}} \rangle = \langle (\mathbf{A} - \tilde{\lambda}\mathbf{I})\tilde{\mathbf{e}}, (\mathbf{A} - \tilde{\lambda}\mathbf{I})\tilde{\mathbf{e}} \rangle \\ &= \langle (\mathbf{A} - \tilde{\mathbf{A}})\tilde{\mathbf{e}}, (\mathbf{A} - \tilde{\mathbf{A}})\tilde{\mathbf{e}} \rangle = \|(\mathbf{A} - \tilde{\mathbf{A}})\tilde{\mathbf{e}}\|^2 \leq \|\mathbf{A} - \tilde{\mathbf{A}}\|^2 \|\tilde{\mathbf{e}}\|^2 = \|\mathbf{A} - \tilde{\mathbf{A}}\|^2, \end{aligned}$$

und wir müssen nur noch die Wurzel ziehen, um den Beweis abzuschließen.  $\blacksquare$

Indem wir diese Abschätzung mit (5.7) kombinieren, erhalten wir für ein  $\lambda \in \sigma(\mathbf{A})$  schon die Ungleichung

$$|\lambda - \tilde{\lambda}| \leq \|\mathbf{A} - \tilde{\mathbf{A}}\| = \frac{\|\mathbf{r}\|}{\|\mathbf{x}\|},$$

können also die Genauigkeit des genäherten Eigenwerts durch die praktisch berechenbaren Größen  $\|\mathbf{r}\|$  und  $\|\mathbf{x}\|$  beschreiben.

Natürlich sind wir auch daran interessiert, Aussagen über die Qualität der Approximation des Eigenvektors zu gewinnen. Diese Aufgabe erweist sich als schwieriger:

**Beispiel 5.20 (Keine Konvergenz)** Wir untersuchen die Matrizen

$$\mathbf{A} := \begin{pmatrix} 1 + 2\epsilon & \\ & 1 \end{pmatrix}, \quad \tilde{\mathbf{A}} := \begin{pmatrix} 1 & \epsilon \\ \epsilon & 1 \end{pmatrix} = \frac{1}{2} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} 1 + \epsilon & \\ & 1 - \epsilon \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}.$$

Für  $\epsilon \rightarrow 0$  konvergieren beide gegen die Einheitsmatrix, also insbesondere auch gegen einander, wir können die Differenz mit Lemma 3.55 sogar explizit berechnen, indem wir die Nullstellen des charakteristischen Polynoms bestimmen:

$$\begin{aligned} \|\mathbf{A} - \tilde{\mathbf{A}}\| &= \left\| \begin{pmatrix} 2\epsilon & -\epsilon \\ -\epsilon & 0 \end{pmatrix} \right\| = \max\{|\lambda| : (\lambda - 2\epsilon)\lambda - \epsilon^2 = 0\} \\ &= \max\{(\sqrt{2} - 1)\epsilon, (\sqrt{2} + 1)\epsilon\} = (\sqrt{2} + 1)\epsilon. \end{aligned}$$

Der Winkel zwischen den kanonischen Einheitsvektoren, die die Eigenvektoren der Matrix  $\mathbf{A}$  sind, und den Eigenvektoren der Matrix  $\tilde{\mathbf{A}}$  dagegen ist durch

$$\cos \angle(\mathbf{x}, \tilde{\mathbf{x}}) = \frac{\left| \left\langle \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\rangle \right|}{\sqrt{2}} = 1/\sqrt{2}$$

gegeben, beträgt also unabhängig von  $\epsilon$  immer  $\pi/4$ . Die Eigenvektoren der Matrizen konvergieren demzufolge nicht gegeneinander.

Für  $\epsilon = 0$  sind in unserem Beispiel *alle* von null verschiedenen Vektoren Eigenvektoren, so dass die Eigenschaft, Eigenvektor zu sein, keine Aussagen über einen Vektor mehr zulässt. Diesen Sonderfall können wir ausschließen, indem wir messen, wie nahe die Eigenwerte beieinander liegen, um so eine „Durchmischung der Eigenräume“ zu vermeiden.

**Definition 5.21 (Spektrallücke)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine normale Matrix. Für alle  $\lambda \in \mathbb{K}$  definieren wir die Spektrallücke zu  $\lambda$  in Bezug auf  $\mathbf{A}$  durch

$$\gamma_A(\lambda) := \inf\{|\mu - \lambda| : \mu \in \sigma(\mathbf{A}) \setminus \{\lambda\}\},$$

also gerade als den Abstand von  $\lambda$  zu dem nächstgelegenen Eigenwert. Wie üblich setzen wir  $\inf \emptyset = \infty$  und verwenden in den folgenden Argumenten die Konvention  $1/\infty = 0$ .

Ausgehend von einer Abschätzung für  $\|\mathbf{A} - \tilde{\mathbf{A}}\|$  ist es relativ leicht, Aussagen im Bildbereich der Matrizen zu formulieren. Da wir daran interessiert sind, eine Aussage im Definitionsbereich, nämlich über die Störung der Eigenvektoren, zu erhalten, brauchen wir eine Möglichkeit, aus ersteren letztere zu gewinnen. Falls  $\lambda$  ein Eigenwert der Matrix  $\mathbf{A}$  ist, kann  $\lambda\mathbf{I} - \mathbf{A}$  nicht injektiv sein, also müssen wir den Kern der Matrix ausblenden.

**Lemma 5.22 (Urbild-Abschätzung)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix, sei  $\lambda \in \mathbb{K}$ . Dann existiert eine orthogonale Projektion  $\mathbf{\Pi} \in \mathbb{K}^{n \times n}$  derart, dass

$$\|\mathbf{z} - \mathbf{\Pi}\mathbf{z}\| \leq \frac{1}{\gamma_A(\lambda)} \|(\lambda\mathbf{I} - \mathbf{A})\mathbf{z}\| \quad \text{für alle } \mathbf{z} \in \mathbb{K}^n \quad (5.8)$$

und  $(\lambda\mathbf{I} - \mathbf{A})\mathbf{\Pi} = \mathbf{0}$  gelten. Letztere Gleichung bedeutet, dass das Bild der Matrix  $\mathbf{\Pi}$  im Kern der Matrix  $\lambda\mathbf{I} - \mathbf{A}$  enthalten ist.

*Beweis.* Da  $\mathbf{A}$  selbstadjungiert ist, finden wir nach Satz 3.47 eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  und eine Diagonalmatrix  $\mathbf{D} \in \mathbb{R}^{n \times n}$  mit

$$\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^*.$$

Falls  $\gamma_A(\lambda) = \infty$  gilt, also  $\sigma(\mathbf{A}) = \{\lambda\}$ , folgt  $\gamma_D(\lambda) = \infty$ , also  $\mathbf{D} = \lambda\mathbf{I}$  und damit  $\lambda\mathbf{I} - \mathbf{A} = \mathbf{0}$ . In diesem Fall setzen wir  $\mathbf{\Pi} = \mathbf{I}$  und sind fertig.

## 5 Die Vektoriteration

Anderenfalls definieren wir eine Diagonalmatrix  $\widehat{\mathbf{\Pi}} \in \mathbb{R}^{n \times n}$  durch

$$\widehat{\pi}_{ij} := \begin{cases} 1 & \text{falls } i = j \text{ und } d_{ii} = \lambda, \\ 0 & \text{ansonsten} \end{cases} \quad \text{für alle } i, j \in [1 : n]$$

und setzen  $\mathbf{\Pi} := \mathbf{Q}\widehat{\mathbf{\Pi}}\mathbf{Q}^*$ . Offenbar gelten  $\mathbf{\Pi}^* = \mathbf{Q}^*\widehat{\mathbf{\Pi}}^*\mathbf{Q} = \mathbf{\Pi}$  und  $\mathbf{\Pi}^2 = \mathbf{Q}\widehat{\mathbf{\Pi}}^2\mathbf{Q}^* = \mathbf{\Pi}$ .

Aus der Definition folgt  $(\lambda\mathbf{I} - \mathbf{D})\widehat{\mathbf{\Pi}} = \mathbf{0}$ , und wir erhalten

$$(\lambda\mathbf{I} - \mathbf{A})\mathbf{\Pi} = \mathbf{Q}(\lambda\mathbf{I} - \mathbf{D})\mathbf{Q}^*\mathbf{Q}\widehat{\mathbf{\Pi}}\mathbf{Q}^* = \mathbf{Q}(\lambda\mathbf{I} - \mathbf{D})\widehat{\mathbf{\Pi}}\mathbf{Q}^* = \mathbf{Q}\mathbf{0}\mathbf{Q}^* = \mathbf{0}.$$

Sei  $\mathbf{z} \in \mathbb{K}^n$ . Wir definieren  $\widehat{\mathbf{z}} := \mathbf{Q}^*\mathbf{z}$  und halten fest, dass nach Lemma 3.34 die Gleichung

$$\|(\lambda\mathbf{I} - \mathbf{A})\mathbf{z}\| = \|\mathbf{Q}(\lambda\mathbf{I} - \mathbf{D})\mathbf{Q}^*\mathbf{z}\| = \|(\lambda\mathbf{I} - \mathbf{D})\widehat{\mathbf{z}}\|$$

gilt. Die Norm auf der rechten Seite können wir einfach berechnen und erhalten

$$\begin{aligned} \|(\lambda\mathbf{I} - \mathbf{D})\widehat{\mathbf{z}}\|^2 &= \sum_{i=1}^n |\lambda - d_{ii}|^2 |\widehat{z}_i|^2 = \sum_{\substack{i=1 \\ d_{ii} \neq \lambda}}^n |\lambda - d_{ii}|^2 |\widehat{z}_i|^2 \geq \sum_{\substack{i=1 \\ d_{ii} \neq \lambda}}^n \gamma_A(\lambda)^2 |\widehat{z}_i|^2 \\ &= \gamma_A(\lambda)^2 \sum_{i=1}^n (1 - \widehat{\pi}_{ii})^2 |\widehat{z}_i|^2 = \gamma_A(\lambda)^2 \|(\mathbf{I} - \widehat{\mathbf{\Pi}})\widehat{\mathbf{z}}\|^2. \end{aligned}$$

Mit Lemma 3.34 ergibt sich

$$\begin{aligned} \gamma_A(\lambda) \|\mathbf{z} - \mathbf{\Pi}\mathbf{z}\| &= \gamma_A(\lambda) \|\mathbf{Q}^*(\mathbf{z} - \mathbf{\Pi}\mathbf{z})\| = \gamma_A(\lambda) \|\widehat{\mathbf{z}} - \widehat{\mathbf{\Pi}}\widehat{\mathbf{z}}\| \\ &\leq \|(\lambda\mathbf{I} - \mathbf{D})\widehat{\mathbf{z}}\| = \|(\lambda\mathbf{I} - \mathbf{A})\mathbf{z}\|, \end{aligned}$$

und die Division durch  $\gamma_A(\lambda) > 0$  führt zu der gewünschten Abschätzung.  $\blacksquare$

**Bemerkung 5.23 (Projektion)** Die in Lemma 5.22 definierte Matrix  $\mathbf{\Pi}$  ist eine orthogonale Projektion, erfüllt also  $\mathbf{\Pi}^2 = \mathbf{\Pi} = \mathbf{\Pi}^*$ .

Aus (5.8) folgt, dass für jedes  $\mathbf{z} \in \text{Kern}(\lambda\mathbf{I} - \mathbf{A})$  die rechte Seite gleich null ist, also auch die linke, so dass  $\mathbf{\Pi}\mathbf{z} = \mathbf{z}$  folgt. Damit ist das Bild der Matrix  $\mathbf{\Pi}$  nicht nur im Kern der Matrix enthalten, sondern es ist bereits der gesamte Kern.

Damit ist  $\mathbf{\Pi}$  eine orthogonale Projektion auf den Kern der Matrix  $\lambda\mathbf{I} - \mathbf{A}$ . Für jedes  $\lambda \in \mathbb{K}$  gibt es genau eine solche Projektion, und sie ordnet jedem Vektor  $\mathbf{z}$  denjenigen Vektor  $\mathbf{\Pi}\mathbf{z}$  aus dem Kern zu, der ihm in der euklidischen Norm am nächsten kommt.

In unserem Fall wird  $\lambda$  ein Eigenwert der Matrix  $\mathbf{A}$  sein, also ist  $\mathbf{\Pi}$  die Projektion auf den zugehörigen Eigenraum  $\mathcal{E}_A(\lambda)$ .

**Satz 5.24 (Gestörtes Eigenwertproblem)** Seien  $\mathbf{A}, \widetilde{\mathbf{A}} \in \mathbb{K}^{n \times n}$  selbstadjungierte Matrizen, sei  $\tilde{\lambda} \in \sigma(\widetilde{\mathbf{A}})$  ein Eigenwert der Matrix  $\widetilde{\mathbf{A}}$  und  $\tilde{\mathbf{e}} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  ein zugehöriger Eigenvektor.

Dann existieren ein Eigenwert  $\lambda \in \sigma(\mathbf{A})$  und ein Vektor  $\mathbf{e} \in \mathcal{E}_A(\lambda)$  aus dem zugehörigen Eigenraum der Matrix  $\mathbf{A}$  mit

$$|\lambda - \tilde{\lambda}| \leq \|\mathbf{A} - \widetilde{\mathbf{A}}\|, \quad \frac{\|\mathbf{e} - \tilde{\mathbf{e}}\|}{\|\tilde{\mathbf{e}}\|} \leq \frac{2}{\gamma_A(\lambda)} \|\mathbf{A} - \widetilde{\mathbf{A}}\|.$$

*Beweis.* Mit dem Satz 5.19 finden wir einen Eigenwert  $\lambda \in \sigma(\mathbf{A})$  mit

$$|\lambda - \tilde{\lambda}| \leq \|\mathbf{A} - \tilde{\mathbf{A}}\|.$$

Mit der Dreiecksungleichung,  $\tilde{\lambda}\tilde{\mathbf{e}} = \tilde{\mathbf{A}}\tilde{\mathbf{e}}$  und (3.8a) erhalten wir

$$\begin{aligned} \|(\lambda\mathbf{I} - \mathbf{A})\tilde{\mathbf{e}}\| &= \|(\tilde{\lambda}\mathbf{I} - \tilde{\mathbf{A}})\tilde{\mathbf{e}} - (\mathbf{A} - \tilde{\mathbf{A}})\tilde{\mathbf{e}} + (\lambda - \tilde{\lambda})\tilde{\mathbf{e}}\| \\ &\leq \|(\tilde{\lambda}\mathbf{I} - \tilde{\mathbf{A}})\tilde{\mathbf{e}}\| + \|(\mathbf{A} - \tilde{\mathbf{A}})\tilde{\mathbf{e}}\| + |\lambda - \tilde{\lambda}| \|\tilde{\mathbf{e}}\| \\ &\leq \|\mathbf{A} - \tilde{\mathbf{A}}\| \|\tilde{\mathbf{e}}\| + \|\mathbf{A} - \tilde{\mathbf{A}}\| \|\tilde{\mathbf{e}}\| = 2\|\mathbf{A} - \tilde{\mathbf{A}}\| \|\tilde{\mathbf{e}}\|. \end{aligned}$$

Wir wenden Lemma 5.22 an und erhalten eine orthogonale Projektion  $\mathbf{\Pi}$  mit  $(\lambda\mathbf{I} - \mathbf{A})\mathbf{\Pi} = \mathbf{0}$  und

$$\|\tilde{\mathbf{e}} - \mathbf{\Pi}\tilde{\mathbf{e}}\| \leq \frac{1}{\gamma_A(\lambda)} \|(\lambda\mathbf{I} - \mathbf{A})\tilde{\mathbf{e}}\| \leq \frac{2}{\gamma_A(\lambda)} \|\mathbf{A} - \tilde{\mathbf{A}}\| \|\tilde{\mathbf{e}}\|.$$

Wir setzen  $\mathbf{e} := \mathbf{\Pi}\tilde{\mathbf{e}}$ . Aus

$$(\lambda\mathbf{I} - \mathbf{A})\mathbf{e} = (\lambda\mathbf{I} - \mathbf{A})\mathbf{\Pi}\tilde{\mathbf{e}} = \mathbf{0}\tilde{\mathbf{e}} = \mathbf{0}$$

folgt  $\mathbf{e} \in \mathcal{E}_A(\lambda)$ . ■

Indem wir diesen Satz auf die im Rahmen der Rückwärtsanalyse konstruierte Matrix  $\tilde{\mathbf{A}}$  anwenden und das Abbruchkriterium der in Algorithmus 5.11 gegebenen Vektoriteration einsetzen, erhalten wir die folgende Aussage über das von ihr berechnete Ergebnis:

**Folgerung 5.25 (Ergebnis der Vektoriteration)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  selbstadjungiert.

Sei  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  ein Vektor, der die Abbruchbedingung

$$\|\mathbf{A}\mathbf{x} - \tilde{\lambda}\mathbf{x}\| \leq \epsilon|\tilde{\lambda}| \|\mathbf{x}\|$$

mit  $\tilde{\lambda} := \Lambda_A(\mathbf{x})$  erfüllt.

Dann existieren ein Eigenwert  $\lambda \in \sigma(\mathbf{A})$  der Matrix  $\mathbf{A}$  und ein Vektor  $\mathbf{e} \in \mathcal{E}_A(\lambda)$  aus dem zugehörigen Eigenraum mit

$$|\lambda - \tilde{\lambda}| \leq |\tilde{\lambda}|\epsilon, \quad |\lambda - \tilde{\lambda}| \leq \frac{4\|\lambda\mathbf{I} - \mathbf{A}\| |\tilde{\lambda}|^2}{\gamma_A(\lambda)^2} \epsilon^2, \quad \frac{\|\mathbf{x} - \mathbf{e}\|}{\|\mathbf{x}\|} \leq \frac{2|\tilde{\lambda}|}{\gamma_A(\lambda)} \epsilon.$$

*Beweis.* Wir bezeichnen das Residuum wieder mit

$$\mathbf{r} := \tilde{\lambda}\mathbf{x} - \mathbf{A}\mathbf{x}$$

und setzen wie zuvor

$$\tilde{\mathbf{A}} := \mathbf{A} + \frac{\mathbf{x}\mathbf{r}^*}{\|\mathbf{x}\|^2} + \frac{\mathbf{r}\mathbf{x}^*}{\|\mathbf{x}\|^2}.$$

Nach Lemma 5.17 und der Voraussetzung gilt

$$\|\mathbf{A} - \tilde{\mathbf{A}}\| = \left\| \frac{\mathbf{x}\mathbf{r}^*}{\|\mathbf{x}\|^2} + \frac{\mathbf{r}\mathbf{x}^*}{\|\mathbf{x}\|^2} \right\| = \frac{\|\mathbf{x}\| \|\mathbf{r}\|}{\|\mathbf{x}\|^2} = \frac{\|\mathbf{r}\|}{\|\mathbf{x}\|} \leq \epsilon|\tilde{\lambda}|.$$

## 5 Die Vektoriteration

Mit Satz 5.24 finden wir  $\lambda \in \sigma(\mathbf{A})$  und  $\mathbf{e} \in \mathcal{E}_A(\lambda)$  mit

$$|\lambda - \tilde{\lambda}| \leq \epsilon |\tilde{\lambda}|, \quad \frac{\|\mathbf{e} - \mathbf{x}\|}{\|\mathbf{x}\|} \leq \frac{2}{\gamma_A(\lambda)} \epsilon |\tilde{\lambda}|.$$

Indem wir die Abschätzung für die Eigenvektoren in Lemma 5.9 einsetzen, folgt

$$|\lambda - \tilde{\lambda}| = |\Lambda_A(\mathbf{x}) - \lambda| \leq \|\lambda \mathbf{I} - \mathbf{A}\| \frac{\|\mathbf{e} - \mathbf{x}\|^2}{\|\mathbf{x}\|^2} \leq \frac{4\|\lambda \mathbf{I} - \mathbf{A}\| |\tilde{\lambda}|^2}{\gamma_A(\lambda)^2} \epsilon^2,$$

also die noch fehlende Abschätzung. ■

**Bemerkung 5.26 (Relativer Fehler)** *Wir können aus den Ergebnissen der vorangehenden Folgerung auch Abschätzungen des relativen Fehlers gewinnen, falls wir  $\lambda \neq 0$  und  $\epsilon < 1$  voraussetzen: Für den Eigenwert gilt*

$$\frac{|\lambda - \tilde{\lambda}|}{|\lambda|} = \frac{|\lambda - \tilde{\lambda}|}{|\tilde{\lambda} + \lambda - \tilde{\lambda}|} \leq \frac{|\lambda - \tilde{\lambda}|}{|\tilde{\lambda}| - |\lambda - \tilde{\lambda}|} = \frac{|\lambda - \tilde{\lambda}|/|\tilde{\lambda}|}{1 - |\lambda - \tilde{\lambda}|/|\tilde{\lambda}|} \leq \frac{\epsilon}{1 - \epsilon}.$$

Man beachte, dass die rechte Seite nun unabhängig von  $\tilde{\lambda}$  und  $\mathbf{x}$  ist.

Für den Eigenvektor erhalten wir

$$\begin{aligned} \frac{\|\mathbf{x} - \mathbf{e}\|}{\|\mathbf{x}\|} &\leq \frac{2|\tilde{\lambda}|}{\gamma_A(\lambda)} \epsilon \leq \frac{2(|\lambda| + |\tilde{\lambda} - \lambda|)}{\gamma_A(\lambda)} \epsilon = \frac{2|\lambda|}{\gamma_A(\lambda)} \left(1 + \frac{|\tilde{\lambda} - \lambda|}{|\lambda|}\right) \epsilon \\ &\leq \frac{2|\lambda|}{\gamma_A(\lambda)} \left(1 + \frac{\epsilon}{1 - \epsilon}\right) \epsilon = \frac{2|\lambda|}{\gamma_A(\lambda)} \frac{1}{1 - \epsilon} \epsilon = \frac{2|\lambda|}{\gamma_A(\lambda)} \frac{\epsilon}{1 - \epsilon}. \end{aligned}$$

Bei dieser Abschätzung ist von besonderem Interesse, dass die relative Spektrallücke  $\gamma_A(\lambda)/|\lambda|$  ausreicht.

Mit Lemma 5.9 folgt daraus

$$|\tilde{\lambda} - \lambda| \leq 4\|\mathbf{A} - \lambda \mathbf{I}\| \frac{|\lambda|^2}{\gamma_A(\lambda)^2} \frac{\epsilon^2}{(1 - \epsilon)^2}.$$

Falls der Fehler des Eigenvektors klein genug ist, falls nämlich

$$\frac{\|\mathbf{x} - \mathbf{e}\|}{\|\mathbf{e}\|} \leq \frac{2|\lambda|}{\gamma_A(\lambda)} \frac{\epsilon}{1 - \epsilon} < 1$$

gilt, erhalten wir mit

$$\begin{aligned} \frac{\|\mathbf{x} - \mathbf{e}\|}{\|\mathbf{e}\|} &= \frac{\|\mathbf{x} - \mathbf{e}\|}{\|\mathbf{x} - \mathbf{x} + \mathbf{e}\|} \leq \frac{\|\mathbf{x} - \mathbf{e}\|}{\|\mathbf{x}\| - \|\mathbf{x} - \mathbf{e}\|} \leq \frac{\|\mathbf{x} - \mathbf{e}\|}{\|\mathbf{x}\| - \frac{2|\lambda|}{\gamma_A(\lambda)} \frac{\epsilon}{1 - \epsilon} \|\mathbf{x}\|} \\ &= \frac{1}{1 - \frac{2|\lambda|}{\gamma_A(\lambda)} \frac{\epsilon}{1 - \epsilon}} \frac{\|\mathbf{x} - \mathbf{e}\|}{\|\mathbf{x}\|} \leq \frac{1}{1 - \frac{2|\lambda|}{\gamma_A(\lambda)} \frac{\epsilon}{1 - \epsilon}} \frac{2|\lambda|}{\gamma_A(\lambda)} \frac{\epsilon}{1 - \epsilon} \end{aligned}$$

eine Abschätzung des relativen Fehlers des Eigenvektors, die lediglich von der Genauigkeit  $\epsilon$  und der relativen Spektrallücke  $\gamma_A(\lambda)/|\lambda|$  abhängt.



### 5.3 Inverse Iteration mit und ohne Shift

Im Beispiel 2.1 ist die Anwendung der Vektoriteration in der Form des Algorithmus 5.8 nicht sinnvoll, da die Eigenvektoren zu den größten Eigenwerten gerade diejenigen sind, die wegen des Diskretisierungsfehlers besonders wenig mit den Eigenvektoren des kontinuierlichen Problems zu tun haben.

Interessanter sind in diesem Fall die Eigenvektoren zu den *kleinsten* Eigenwerten, da diese Vektoren in der Regel relativ gut approximiert werden und auch für ingenieurtechnische Anwendungen von weitaus größerem Interesse sind, schließlich will man häufig die *niedrigste* Frequenz kennen, bei der es zu Oszillationen kommen kann.

Falls die Matrix  $\mathbf{A}$  regulär und  $\lambda$  einer ihrer Eigenwerte mit einem zugehörigen Eigenvektor  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  ist, gilt

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{x}, \quad \mathbf{x} = \lambda\mathbf{A}^{-1}\mathbf{x}, \quad \frac{1}{\lambda}\mathbf{x} = \mathbf{A}^{-1}\mathbf{x},$$

also ist der Kehrwert jedes Eigenwerts von  $\mathbf{A}$  auch ein Eigenwert von  $\mathbf{A}^{-1}$ . Da die Inverse von  $\mathbf{A}^{-1}$  wieder  $\mathbf{A}$  ist, folgt, dass das Spektrum von  $\mathbf{A}^{-1}$  nur aus den Kehrwerten der Eigenwerte von  $\mathbf{A}$  besteht.

Insbesondere ist der Kehrwert des betragskleinsten Eigenwerts von  $\mathbf{A}$  gerade der betragsgrößte Eigenwert von  $\mathbf{A}^{-1}$ . Falls wir an dem betragskleinsten Eigenwert interessiert sind, liegt es also nahe, den Algorithmus 5.8 auf die inverse Matrix  $\mathbf{A}^{-1}$  anzuwenden.

Wir definieren also die  $m$ -te Iterierte unseres neuen Verfahrens durch

$$\mathbf{x}^{(m)} := \mathbf{A}^{-m}\mathbf{x}^{(0)} = \mathbf{A}^{-1}\mathbf{x}^{(m-1)} \quad \text{für alle } m \in \mathbb{N}. \quad (5.9)$$

Da es aus der Anwendung der Vektoriteration auf die Inverse entsteht, trägt dieses Iterationsverfahren den Namen *inverse Iteration*.

Für die Konvergenzuntersuchung müssen wir die Voraussetzungen so wählen, dass sich Satz 5.4 oder Folgerung 5.7 auf  $\mathbf{A}^{-1}$  statt  $\mathbf{A}$  anwenden lassen.

Wir beschränken uns auf den einfacheren der beiden Fälle: Sei  $\mathbf{A}$  im Folgenden eine normale Matrix. Dann existieren nach Folgerung 3.54 eine unitäre Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und eine Diagonalmatrix  $\mathbf{D} \in \mathbb{C}^n$  mit

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{D}.$$

Nach Bemerkung 3.42 können wir die Reihenfolge der Eigenwerte frei wählen, so dass wir

$$\mathbf{D} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}, \quad |\lambda_1| \leq |\lambda_2| \leq \dots \leq |\lambda_n| \quad (5.10)$$

erhalten. Daraus folgt

$$\left| \frac{1}{\lambda_1} \right| \geq \left| \frac{1}{\lambda_2} \right| \geq \dots \geq \left| \frac{1}{\lambda_n} \right|.$$

Mit  $\mathbf{e} := \mathbf{Q}\delta^{(1)}$  bezeichnen wir wieder einen Eigenvektor, der zu dem Eigenwert  $\lambda_1$  der Matrix  $\mathbf{A}$  gehört.

## 5 Die Vektoriteration

**Satz 5.27 (Konvergenz)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine invertierbare normale Matrix. Sie besitzt eine Schur-Zerlegung  $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^*$  mit einer unitären Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und einer Diagonalmatrix  $\mathbf{D} \in \mathbb{C}^{n \times n}$  der Form (5.10).

Dann ist  $\mathbf{e} := \mathbf{Q}\delta^{(1)}$  ein Eigenvektor zu dem betragskleinsten Eigenwert  $\lambda_1$ .

Sei  $\mathbf{x}^{(0)} \in \mathbb{K}^n$  mit  $\langle \mathbf{e}, \mathbf{x}^{(0)} \rangle \neq 0$  gegeben. Dann gilt

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(m)}) \leq \left( \frac{|\lambda_1|}{|\lambda_2|} \right)^m \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \quad \text{für alle } m \in \mathbb{N}_0.$$

*Beweis.* Wir setzen  $\widehat{\mathbf{A}} := \mathbf{A}^{-1}$ . Es gilt mit Lemma 3.35

$$\mathbf{Q}^* \widehat{\mathbf{A}} \mathbf{Q} = \mathbf{Q}^* \mathbf{A}^{-1} \mathbf{Q} = \mathbf{Q}^* (\mathbf{Q}\mathbf{D}\mathbf{Q}^*)^{-1} \mathbf{Q} = \mathbf{Q}^* \mathbf{Q} \mathbf{D}^{-1} \mathbf{Q}^* \mathbf{Q} = \mathbf{D}^{-1},$$

also ist insbesondere  $\widehat{\mathbf{A}}$  diagonalisierbar mit derselben Transformation  $\mathbf{Q}$ . Die Eigenwerte von  $\widehat{\mathbf{A}}$  sind offenbar gerade die Diagonalelemente von

$$\mathbf{D}^{-1} = \begin{pmatrix} 1/\lambda_1 & & \\ & \ddots & \\ & & 1/\lambda_n \end{pmatrix},$$

und sie sind nach Voraussetzung dem Betrag nach absteigend sortiert. Also dürfen wir Satz 5.4 anwenden und erhalten

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(m)}) \leq \left( \frac{1/|\lambda_2|}{1/|\lambda_1|} \right)^m \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) = \left( \frac{|\lambda_1|}{|\lambda_2|} \right)^m \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \quad \text{für alle } m \in \mathbb{N}_0. \quad \blacksquare$$

**Bemerkung 5.28** Mit der inversen Iteration lässt sich der Eigenvektor zu dem kleinsten Eigenwert der Matrix aus Beispiel 2.1 berechnen. Da der kontinuierliche Operator die Eigenwerte

$$\lambda_k = c \frac{\pi^2}{\ell^2} k^2$$

besitzt und die Eigenwerte der diskreten Matrix gegen diese Werte konvergieren, wird die Konvergenzgeschwindigkeit der inversen Iteration gegen

$$\frac{|\lambda_1|}{|\lambda_2|} = 1/4$$

streben, also unabhängig von der Auflösung der Diskretisierung sein.

Während bei der Jacobi-Iteration die Konvergenzgeschwindigkeit potentiell von der Dimension der Matrix abhängt, ist sie also bei der inversen Iteration davon unabhängig. Das ist insbesondere bei hohen Auflösungen ein sehr großer Vorteil.

An diesem Beispiel zeigt sich auch die Wichtigkeit einer dem Problem angemessenen Implementierung des Verfahrens: Falls man bei der Lösung des tridiagonalen Gleichungssystems die Bandstruktur (etwa mit Hilfe einer LR-Zerlegung) ausnutzt, benötigt

der gesamte Schleifenrumpf lediglich  $\mathcal{O}(n)$  Operationen, so dass sich mit fast linearem Aufwand der gesuchte Eigenvektor bestimmen lässt.

Würde man stattdessen  $\mathbf{A}^{-1}$  explizit berechnen, wäre dafür im Allgemeinen ein Aufwand von  $\mathcal{O}(n^3)$  Operationen erforderlich, während die Multiplikation mit  $\mathbf{A}^{-1}$  in jedem Schritt der inversen Iteration einen Aufwand von  $\mathcal{O}(n^2)$  nach sich ziehen würde.

Wie wir gesehen haben steht uns  $\mathbf{A}^{-1}$  in der Praxis oft nicht zur Verfügung, oder der Aufwand für die Berechnung der Inversen ist inakzeptabel, deshalb empfiehlt es sich, den Schritt

$$\mathbf{w}^{(m+1)} := \mathbf{A}^{-1}\mathbf{x}^{(m)}$$

des Originalverfahrens durch das Lösen des Gleichungssystems

$$\mathbf{A}\mathbf{w}^{(m+1)} = \mathbf{x}^{(m)}$$

zu ersetzen, das sich in vielen praktischen Anwendungen mit Hilfe einer Faktorisierung oder eines iterativen Lösungsverfahrens effizient durchführen lässt.

**Algorithmus 5.29 (Inverse Iteration)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und  $\mathbf{x}^{(0)} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  gegeben. Der folgende Algorithmus führt die inverse Iteration aus.

```

m ← 0
x(m) ← x(m) / ||x(m)||
while „Fehler zu groß“ do begin
  Löse A w(m+1) = x(m)
  x(m+1) ← w(m+1) / ||w(m+1)||
  m ← m + 1
end

```

Auch dieser Algorithmus lässt sich wieder mit dem Rayleigh-Quotienten kombinieren, um eine Approximation des betragskleinsten Eigenwerts und damit auch eine Schätzung des Iterationsfehlers zu gewinnen. Wir können den im Zuge der Iteration ohnehin berechneten Vektor  $\mathbf{w}^{(m+1)}$  verwenden, um  $\lambda^{(m)} := \Lambda_{\mathbf{A}^{-1}}(\mathbf{x}^{(m)})$  effizient zu berechnen und

$$\epsilon^{(m)} := \|\mathbf{A}^{-1}\mathbf{x}^{(m)} - \lambda^{(m)}\mathbf{x}^{(m)}\| = \|\mathbf{w}^{(m+1)} - \langle \mathbf{x}^{(m)}, \mathbf{w}^{(m+1)} \rangle \mathbf{x}^{(m)}\|$$

als Maß des Fehlers benutzen. Dabei sollte man natürlich nicht vergessen, dass  $\lambda^{(m)}$  nun gegen den Kehrwert des kleinsten Eigenwerts konvergieren wird, nicht mehr gegen den Eigenwert selbst. Entsprechend kann auch Folgerung 5.25 nur in modifizierter Form angewendet werden, nämlich mit  $\mathbf{A}^{-1}$  statt  $\mathbf{A}$ .

**Bemerkung 5.30 (Fehlerschätzung)** Falls die Multiplikation der Matrix  $\mathbf{A}$  mit einem Vektor effizient durchgeführt werden kann, können wir natürlich auch den Hilfsvektor  $\mathbf{a}^{(m)} := \mathbf{A}\mathbf{x}^{(m)}$  berechnen, damit den Rayleigh-Quotienten  $\lambda^{(m)} = \Lambda_{\mathbf{A}}(\mathbf{x}^{(m)}) = \langle \mathbf{x}^{(m)}, \mathbf{a}^{(m)} \rangle$  bestimmen, und damit das von der Vektoriteration bekannte Residuum  $\mathbf{r}^{(m)} = \lambda^{(m)}\mathbf{x}^{(m)} - \mathbf{a}^{(m)}$ , so dass sich die bereits entwickelte Fehlertheorie auch für die inverse Iteration verwenden lässt.

## 5 Die Vektoriteration

Bei der Betrachtung der inversen Iteration stellen sich zwei Fragen:

1. Die Forderung nach der Regularität von  $\mathbf{A}$  ist im Kontext von Eigenwertverfahren unerwartet, schließlich kam sie bei keiner der bisherigen theoretischen Aussagen vor und schränkt die Anwendbarkeit der inversen Iteration ein. Lässt sie sich vermeiden?
2. Mit der Vektoriteration und der inversen Iteration stehen Verfahren zur Bestimmung der größten und kleinsten Eigenwerte zur Verfügung. Lassen sich auch „mittlere“ Eigenwerte berechnen?

Beide Fragen lassen sich positiv beantworten, wenn man von der Inversen  $\mathbf{A}^{-1}$  von  $\mathbf{A}$  zu der einer um einen gewissen Betrag  $\mu \in \mathbb{K}$  „verschobenen“ Matrix übergeht, also  $\mathbf{A}^{-1}$  durch  $(\mathbf{A} - \mu\mathbf{I})^{-1}$  ersetzt.

Falls  $\mu \notin \sigma(\mathbf{A})$  gilt, folgt für  $\lambda \in \mathbb{K}$ ,  $\mathbf{x} \in \mathbb{K}^n \setminus \{0\}$  die Äquivalenz

$$\mathbf{Ax} = \lambda\mathbf{x} \iff (\mathbf{A} - \mu\mathbf{I})\mathbf{x} = (\lambda - \mu)\mathbf{x} \iff (\mathbf{A} - \mu\mathbf{I})^{-1}\mathbf{x} = \frac{1}{\lambda - \mu}\mathbf{x},$$

so dass alle Eigenvektoren von  $\mathbf{A}$  auch Eigenvektoren von  $(\mathbf{A} - \mu\mathbf{I})^{-1}$  sind. Man kann einfach nachweisen, dass jeder Eigenwert der letzteren Matrix sich als  $1/(\lambda - \mu)$  mit einem  $\lambda \in \sigma(\mathbf{A})$  darstellen lässt, wir gewinnen oder verlieren also keine Eigenwerte, ähnlich wie bei der inversen Iteration.

Die Konvergenz der inversen Iteration hängt von dem Verhältnis des vom Betrag her kleinsten zum vom Betrag her zweitkleinsten Eigenwerts ab. Verwendet man die Matrix  $(\mathbf{A} - \mu\mathbf{I})^{-1}$  anstelle der Matrix  $\mathbf{A}^{-1}$ , so ist der vom Betrag her größte Eigenwert derjenige, der dem Wert  $\mu$  am nächsten liegt. Durch die Wahl des sogenannten *Shift-Parameters*  $\mu$  lässt sich also festlegen, welche Eigenwerte und Eigenvektoren berechnet werden sollen.

Für die *inverse Iteration mit Shift*  $\mu$  definieren wir die Folge  $(\mathbf{x}^{(m)})_{m \in \mathbb{N}_0}$  der Iterierten durch

$$\mathbf{x}^{(m)} := (\mathbf{A} - \mu\mathbf{I})^{-m}\mathbf{x}^{(0)} = (\mathbf{A} - \mu\mathbf{I})^{-1}\mathbf{x}^{(m-1)} \quad \text{für alle } m \in \mathbb{N}. \quad (5.11)$$

Auch in diesem Fall lässt sich Satz 5.4 anwenden, um eine Konvergenzaussage zu erhalten. Sei  $\mu \in \mathbb{K} \setminus \sigma(\mathbf{A})$ . Wie gehabt soll  $\mathbf{A}$  normal sein, so dass wir mit Folgerung 3.54 eine unitäre Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und eine Diagonalmatrix  $\mathbf{D} \in \mathbb{C}^{n \times n}$  mit

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{D}$$

finden. Dank Bemerkung 3.42 können wir diesmal die Eigenwerte so anordnen, dass

$$\mathbf{D} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}, \quad |\lambda_1 - \mu| \leq |\lambda_2 - \mu| \leq \dots \leq |\lambda_n - \mu| \quad (5.12)$$

gilt. Daraus folgt nun

$$\left| \frac{1}{\lambda_1 - \mu} \right| \geq \left| \frac{1}{\lambda_2 - \mu} \right| \geq \dots \geq \left| \frac{1}{\lambda_n - \mu} \right|.$$

Mit  $\mathbf{e} := \mathbf{Q}\delta^{(1)}$  bezeichnen wir wieder einen Eigenvektor, der zu dem Eigenwert  $\lambda_1$  der Matrix  $\mathbf{A}$  gehört. Dann erhalten wir die folgende Konvergenzaussage:

**Satz 5.31 (Konvergenz)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine normale Matrix, sei  $\mu \in \mathbb{K} \setminus \sigma(\mathbf{A})$ . Die Matrix  $\mathbf{A}$  besitzt eine Schurzerlegung  $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^*$  mit einer unitären Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und einer Diagonalmatrix  $\mathbf{D} \in \mathbb{C}^{n \times n}$  der Form (5.12).

Dann ist  $\mathbf{e} := \mathbf{Q}\delta^{(1)}$  ein Eigenvektor zu dem Eigenwert  $\lambda_1$ , der  $\mu$  am nächsten liegt. Sei  $\mathbf{x}^{(0)} \in \mathbb{K}^n$  mit  $\langle \mathbf{e}, \mathbf{x}^{(0)} \rangle \neq 0$  gegeben. Dann gilt

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(m)}) \leq \left( \frac{|\lambda_1 - \mu|}{|\lambda_2 - \mu|} \right)^m \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \quad \text{für alle } m \in \mathbb{N}_0.$$

*Beweis.* Wir setzen  $\widehat{\mathbf{A}} := (\mathbf{A} - \mu\mathbf{I})^{-1}$ . Es gilt

$$\mathbf{Q}^* \widehat{\mathbf{A}} \mathbf{Q} = \mathbf{Q}^* (\mathbf{A} - \mu\mathbf{I})^{-1} \mathbf{Q} = (\mathbf{D} - \mu\mathbf{I})^{-1} = \begin{pmatrix} 1/(\lambda_1 - \mu) & & \\ & \ddots & \\ & & 1/(\lambda_n - \mu) \end{pmatrix},$$

also ist  $\widehat{\mathbf{A}}$  diagonalisierbar und die Eigenwerte von  $(\mathbf{D} - \mu\mathbf{I})^{-1}$  sind dem Betrag nach absteigend sortiert. Also können wir den Satz 5.4 anwenden, um die gesuchte Abschätzung zu erhalten. ■

Die Folge der normierten Iterierten lässt sich einfach mit Hilfe des folgenden Algorithmus berechnen:

**Algorithmus 5.32 (Inverse Iteration mit Shift)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und  $\mathbf{x}^{(0)} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  gegeben. Der folgende Algorithmus führt die inverse Iteration mit dem Shift-Wert  $\mu \in \mathbb{K} \setminus \sigma(\mathbf{A})$  aus.

```

m ← 0
x(0) ← x(0) / ||x(0)||
while „Fehler zu groß“ do begin
  Löse (A - μI)w(m+1) = x(m)
  x(m+1) ← w(m+1) / ||w(m+1)||
  m ← m + 1
end

```

Der Shift-Parameter ermöglicht es uns nicht nur, beliebige Eigenwerte zu approximieren, er kann bei geschickter Wahl auch zu einer erheblichen Beschleunigung der Konvergenz führen, wie die folgenden Beispiele zeigen:

**Beispiel 5.33 (Trennung von Eigenwerten)** Wir untersuchen die Matrix

$$\mathbf{A}_\epsilon := \begin{pmatrix} 1 & 0 \\ 0 & 1 + \epsilon \end{pmatrix}$$

## 5 Die Vektoriteration

für  $\epsilon \in \mathbb{R}_{>0}$ . Offenbar gilt  $\sigma(\mathbf{A}_\epsilon) = \{1, 1 + \epsilon\}$ . Wendet man die inverse Iteration auf diese Matrix und einen Startvektor aus  $(\mathbb{R} \setminus \{0\}) \times (\mathbb{R} \setminus \{0\})$  an, so konvergiert sie gemäß Satz 5.27 mit einer Geschwindigkeit von  $1/(1 + \epsilon)$ . Falls  $\epsilon$  klein ist, erhalten wir sehr langsame Konvergenz.

Wendet man die inverse Iteration mit einem Shift von  $\mu = 1 + \epsilon/3$  auf die Matrix und den Startvektor an, so konvergiert sie mit einer Geschwindigkeit von

$$\left| \frac{1 - \mu}{1 + \epsilon - \mu} \right| = \left| \frac{-\epsilon/3}{2\epsilon/3} \right| = 1/2.$$

Durch geschickte Wahl des Shift-Parameters lässt sich also auch bei nahe beieinander liegenden (auch als „schlecht separiert“ bezeichneten) Eigenwerten eine gute Konvergenzrate erzielen.

Betrachten wir als nächstes die Matrix

$$\mathbf{B} := \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}.$$

Ihr charakteristisches Polynom ist  $p_B(\lambda) = \lambda^2 + 1$ , also folgt  $\sigma(\mathbf{B}) = \{i, -i\}$ . Wegen  $|i| = |-i| = 1$  ist nicht zu erwarten, dass die Vektoriteration oder die inverse Iteration bei dieser Matrix konvergieren. Selbst in diesem Fall lässt sich durch Verwendung eines Shift-Wertes von  $\mu = i/2$  noch die gute Konvergenzrate von  $1/3$  erzielen:

$$\left| \frac{i - \mu}{-i - \mu} \right| = \left| \frac{i/2}{-3i/2} \right| = 1/3.$$

Zum Abschluß dieses Abschnittes sollen noch einmal die Eigenschaften der einzelnen vorgestellten Verfahren zusammengefasst werden:

- Vektoriteration:  
 Berechne  $\mathbf{x}^{(m+1)} = \mathbf{A}\mathbf{x}^{(m)}$ ,  
 konvergent, falls  $|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_n|$ ,  
 Konvergenzrate  $|\lambda_2|/|\lambda_1|$ .
- Inverse Iteration:  
 Löse  $\mathbf{A}\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)}$ ,  
 konvergent, falls  $|\lambda_1| < |\lambda_2| \leq \dots \leq |\lambda_n|$ ,  
 Konvergenzrate  $|\lambda_1|/|\lambda_2|$ .
- Inverse Iteration mit Shift:  
 Löse  $(\mathbf{A} - \mu\mathbf{I})\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)}$ ,  
 konvergent, falls  $|\lambda_1 - \mu| < |\lambda_2 - \mu| \leq \dots \leq |\lambda_n - \mu|$ ,  
 Konvergenzrate  $|\lambda_1 - \mu|/|\lambda_2 - \mu|$ .

## 5.4 Inverse Iteration mit Rayleigh-Shift

Die geschickte Wahl des Shift-Parameters  $\mu$  ist offensichtlich von großer Bedeutung für die Geschwindigkeit des Verfahrens. Da die Konvergenz desto besser wird, je näher  $\mu$  an dem gesuchten Eigenwert liegt (beziehungsweise desto weiter es von allen anderen entfernt ist), sind wir daran interessiert,  $\mu$  als Approximation des uns interessierenden Eigenwerts zu wählen. Lemma 5.9 legt die Idee nahe, für die Berechnung von  $\mu$  den Rayleigh-Quotienten zu verwenden.

Wenn wir davon ausgehen, dass  $\mathbf{x}^{(0)}$  bereits eine relativ gute Approximation eines Eigenvektors  $\mathbf{e}$  zu dem Eigenwert  $\lambda_1$  ist, wird nach Lemma 5.9 die Zahl  $\lambda^{(0)} := \Lambda_A(\mathbf{x}^{(0)})$  eine gute Approximation von  $\lambda_1$  sein. Wenn wir nun  $\lambda^{(0)}$  als Shift verwenden und

$$\mathbf{x}^{(1)} := (\mathbf{A} - \lambda^{(0)}\mathbf{I})^{-1}\mathbf{x}^{(0)}$$

berechnen, erhalten wir nach Satz 5.31 eine Abschätzung der Form

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) \leq \frac{|\lambda_1 - \lambda^{(0)}|}{|\lambda_2 - \lambda^{(0)}|} \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}),$$

wobei die Eigenwerte wieder in der Form

$$|\lambda_1 - \lambda^{(0)}| \leq |\lambda_2 - \lambda^{(0)}| \leq \dots \leq |\lambda_n - \lambda^{(0)}|$$

angeordnet sind und

$$\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \quad (5.13)$$

mit einer geeigneten unitären Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  gilt.

Aus Lemma 5.9 erhalten wir eine Abschätzung für  $|\lambda_1 - \lambda^{(0)}|$ , so dass wir

$$\begin{aligned} \tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) &\leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{|\lambda_2 - \lambda^{(0)}|} \sin \angle(\mathbf{e}, \mathbf{x}^{(0)}) \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \\ &\leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{|\lambda_2 - \lambda^{(0)}|} \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \end{aligned} \quad (5.14)$$

bewiesen haben. Diese Formel suggeriert, dass der Fehler der neuen Iterierten  $\mathbf{x}^{(1)}$  sich wie das Quadrat des Fehlers der alten Iterierten  $\mathbf{x}^{(0)}$  verhält, dass wir also auf *quadratische Konvergenz* hoffen dürfen. Das würde bedeuten, dass das Verfahren sehr schnell konvergiert, sobald  $\mathbf{x}^{(0)}$  dem gesuchten Eigenraum hinreichend nahe ist.

Ausgehend von  $\mathbf{x}^{(1)}$  können wir dann einen neuen Shift-Parameter  $\lambda^{(1)} := \Lambda_A(\mathbf{x}^{(1)})$  bestimmen und den Vorgang wiederholen, um eine weitere Näherung  $\mathbf{x}^{(2)}$  zu gewinnen. Da die Berechnung des Rayleigh-Quotienten nichtlinear ist, wird die so definierte *inverse Iteration mit Rayleigh-Shift*, kurz *Rayleigh-Iteration*, anders als die Vektoriteration oder die konventionelle inverse Iteration, ein nichtlineares Verfahren für die Approximation des Eigenvektors sein.

**Algorithmus 5.34 (Rayleigh-Iteration)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und  $\mathbf{x}^{(0)} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  gegeben. Der folgende Algorithmus führt die Rayleigh-Iteration aus.

```

 $m \leftarrow 0$ 
 $\mathbf{x}^{(m)} \leftarrow \mathbf{x}^{(0)} / \|\mathbf{x}^{(0)}\|$ 
while „Fehler zu groß“ do begin
   $\lambda^{(m)} \leftarrow \langle \mathbf{x}^{(m)}, \mathbf{A}\mathbf{x}^{(m)} \rangle$ 
  Löse  $(\mathbf{A} - \lambda^{(m)}\mathbf{I})\mathbf{w}^{(m+1)} = \mathbf{x}^{(m)}$ 
   $\mathbf{x}^{(m+1)} \leftarrow \mathbf{w}^{(m+1)} / \|\mathbf{w}^{(m+1)}\|$ 
   $m \leftarrow m + 1$ 
end

```

Auf den ersten Blick ist die inverse Iteration mit Rayleigh-Shift nicht viel aufwendiger als die inverse Iteration mit konstantem Shift. In praktischen Implementierungen ist das leider häufig nicht der Fall:

**Bemerkung 5.35 (Rechenaufwand)** In der Praxis werden die Gleichungssysteme

$$(\mathbf{A} - \mu\mathbf{I})\mathbf{w}^{(m+1)} = \mathbf{x}^{(m)}, \quad (\mathbf{A} - \lambda^{(m)}\mathbf{I})\mathbf{w}^{(m+1)} = \mathbf{x}^{(m)},$$

die bei der inversen Iteration mit konstantem bzw. mit Rayleigh-Shift auftreten, häufig mit Hilfe einer Faktorisierung gelöst, etwa mit einer LR-Faktorisierung mit Pivotsuche.

Im Falle eines konstanten Shift-Parameters kann die Berechnung der Faktorisierung vor dem Eintritt in die zentrale Schleife der inversen Iteration stattfinden, so dass ein Iterationsschritt lediglich das Vorwärts- und Rückwärtseinsetzen in die bereits berechneten Faktoren erfordert.

Für den Rayleigh-Shift dagegen ändert sich die Matrix potentiell in jedem Schritt, so dass für jede Iterierte die Faktorisierung erneut berechnet werden muss. Deshalb kann die inverse Iteration mit Rayleigh-Shift unter Umständen wesentlich aufwendiger als die Variante mit konstantem Shift werden.

Ein einfacher Ausweg besteht darin, den Shift-Parameter nicht in jedem Schritt zu aktualisieren, sondern in größeren Abständen, um die Anzahl der Faktorisierungen zu reduzieren. Dann verliert man zwar potentiell die quadratische Konvergenz, gewinnt aber ein schnelleres Verfahren.

Alternativ kann wie in Übungsaufgabe 3.49 die Matrix zunächst (mit potentiell kubischem Aufwand) in Tridiagonalgestalt gebracht werden. Für Tridiagonalmatrizen lassen sich Faktorisierungen mit linearem Aufwand berechnen, so dass die Durchführung der Rayleigh-Iteration sehr effizient wird.

Um aus der Ungleichung (5.14) eine quadratische Konvergenzaussage zu gewinnen, müssen wir den Nenner  $|\lambda_2 - \lambda^{(m)}|$  durch eine von  $m$  unabhängige Konstante abschätzen. Das gelingt uns, falls wir voraussetzen, dass  $\lambda^{(m)}$  hinreichend nahe an  $\lambda_1$  liegt, und das folgt, falls der Winkel zwischen  $\mathbf{x}^{(m)}$  und dem gesuchten Eigenvektor  $\mathbf{e}$  hinreichend klein ist. Falls uns also gute Startvektoren zur Verfügung stehen, dürfen wir auf schnelle Konvergenz hoffen.



**Satz 5.36 (Quadratische Konvergenz)** Sei  $\delta \in [0, 1)$ . Sei  $\mathbf{x}^{(0)} \in \mathbb{K}^n$  mit

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq \delta \frac{|\lambda_2 - \lambda_1|}{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}$$

gegeben, und sei  $\lambda^{(0)} := \Lambda_A(\mathbf{x}^{(0)})$ . Für den Vektor  $\mathbf{x}^{(1)} := (\mathbf{A} - \lambda^{(0)} \mathbf{I})^{-1} \mathbf{x}^{(0)}$  gilt dann

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) \leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{(1 - \delta)|\lambda_2 - \lambda_1|} \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}).$$

Falls  $\delta \leq 1/2$  gilt, erhalten wir insbesondere

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) \leq \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}),$$

so dass auch  $\mathbf{x}^{(1)}$  die Voraussetzungen des Satzes erfüllt und die Rayleigh-Iteration fortgeführt werden kann.

*Beweis.* Den größten Teil der Arbeit haben wir bereits in (5.14) geleistet, wir müssen lediglich noch den Nenner abschätzen. Mit der Dreiecksungleichung erhalten wir

$$|\lambda_2 - \lambda^{(0)}| \geq |\lambda_2 - \lambda_1| - |\lambda_1 - \lambda^{(0)}|,$$

und mit Lemma 5.9 und der Voraussetzung folgt

$$|\lambda_1 - \lambda^{(0)}| \leq \|\mathbf{A} - \lambda_1 \mathbf{I}\| \sin \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq \|\mathbf{A} - \lambda_1 \mathbf{I}\| \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq \delta |\lambda_2 - \lambda_1|,$$

so dass wir insgesamt zu

$$|\lambda_2 - \lambda^{(0)}| \geq |\lambda_2 - \lambda_1| - \delta |\lambda_2 - \lambda_1| = (1 - \delta) |\lambda_2 - \lambda_1|$$

gelangen. Aus (5.14) folgt

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) \leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{|\lambda_2 - \lambda^{(0)}|} \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{(1 - \delta) |\lambda_2 - \lambda_1|} \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}).$$

Gelte nun  $\delta \leq 1/2$ . Dann erhalten wir

$$\begin{aligned} \tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) &\leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{(1 - \delta) |\lambda_2 - \lambda_1|} \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \\ &\leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{(1 - \delta) |\lambda_2 - \lambda_1|} \delta \frac{|\lambda_2 - \lambda_1|}{\|\mathbf{A} - \lambda_1 \mathbf{I}\|} \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \\ &= \frac{\delta}{1 - \delta} \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}). \end{aligned}$$

■

Diese Aussage gilt auch noch für allgemeine diagonalisierbare Matrizen, wenn man wie in Folgerung 5.7 die Konditionszahl der diagonalisierenden Ähnlichkeitstransformation an geeigneter Stelle berücksichtigt.

Für normale Matrizen lässt sich sogar *kubische* Konvergenz beweisen:

## 5 Die Vektoriteration

**Satz 5.37 (Kubische Konvergenz)** Sei  $\mathbf{A}$  normal. Sei  $\delta \in [0, 1)$ . Sei  $\mathbf{x}^{(0)} \in \mathbb{K}^n$  mit

$$\tan^2 \angle(\mathbf{e}, \mathbf{x}^{(1)}) \leq \delta \frac{|\lambda_2 - \lambda_1|}{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}$$

gegeben, und sei  $\lambda^{(0)} := \Lambda_A(\mathbf{x}^{(0)})$ . Für den Vektor  $\mathbf{x}^{(1)} := (\mathbf{A} - \lambda^{(0)} \mathbf{I})^{-1} \mathbf{x}^{(0)}$  gilt dann

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) \leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{(1 - \delta)|\lambda_2 - \lambda_1|} \tan^3 \angle(\mathbf{e}, \mathbf{x}^{(0)}).$$

Falls  $\delta \leq 1/2$  gilt, erhalten wir insbesondere

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) \leq \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}),$$

so dass auch  $\mathbf{x}^{(1)}$  die Voraussetzungen des Satzes erfüllt und die Rayleigh-Iteration fortgeführt werden kann.

*Beweis.* Nach Lemma 5.9 gilt

$$|\lambda_1 - \lambda^{(0)}| \leq \|\mathbf{A} - \lambda_1 \mathbf{I}\| \sin^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq \|\mathbf{A} - \lambda_1 \mathbf{I}\| \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}),$$

und Einsetzen in die Abschätzung des Satzes 5.31 ergibt

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) \leq \frac{|\lambda_1 - \lambda^{(0)}|}{|\lambda_2 - \lambda^{(0)}|} \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{|\lambda_2 - \lambda^{(0)}|} \tan^3 \angle(\mathbf{e}, \mathbf{x}^{(0)}).$$

Mit der umgekehrten Dreiecksungleichung folgt

$$\begin{aligned} |\lambda_2 - \lambda^{(0)}| &\geq |\lambda_2 - \lambda_1| - |\lambda_1 - \lambda^{(0)}| \geq |\lambda_2 - \lambda_1| - \|\mathbf{A} - \lambda_1 \mathbf{I}\| \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \\ &\geq |\lambda_2 - \lambda_1| - \delta |\lambda_2 - \lambda_1| = (1 - \delta) |\lambda_2 - \lambda_1|. \end{aligned}$$

Das gewünschte Resultat folgt durch Einsetzen in den Nenner der vorangehenden Ungleichung.

Gelte nun  $\delta \leq 1/2$ . Dann erhalten wir

$$\begin{aligned} \tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) &\leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{(1 - \delta) |\lambda_2 - \lambda_1|} \tan^3 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \\ &\leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{(1 - \delta) |\lambda_2 - \lambda_1|} \delta \frac{|\lambda_2 - \lambda_1|}{\|\mathbf{A} - \lambda_1 \mathbf{I}\|} \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \\ &= \frac{\delta}{1 - \delta} \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq \tan \angle(\mathbf{e}, \mathbf{x}^{(0)}). \end{aligned}$$

■

Als Beispiel setzen wir  $\delta = 1/2$  und definieren die Konstante

$$C := \frac{2\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{|\lambda_2 - \lambda_1|}.$$

Die Voraussetzung des Satzes 5.36 nimmt dann die Form

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq 1/C \quad (5.15)$$

an, und seine Aussage kann als

$$\begin{aligned} \tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) &\leq \frac{\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{(1 - \delta)|\lambda_2 - \lambda_1|} \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \\ &= \frac{2\|\mathbf{A} - \lambda_1 \mathbf{I}\|}{|\lambda_2 - \lambda_1|} \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}) = C \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \end{aligned}$$

formuliert werden. Wir nehmen an, dass unser Startvektor etwas besser als unbedingt nötig ist, dass nämlich

$$\tan \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq \frac{1}{2C}$$

gilt. Dann erhalten wir

$$\begin{aligned} \tan \angle(\mathbf{e}, \mathbf{x}^{(1)}) &\leq C \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq C \frac{1}{4C^2} = \frac{1}{4C} = \frac{1}{2^{2^1} C}, \\ \tan \angle(\mathbf{e}, \mathbf{x}^{(2)}) &\leq C \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(1)}) \leq C \frac{1}{16C^2} = \frac{1}{16C} = \frac{1}{2^{2^2} C}, \\ \tan \angle(\mathbf{e}, \mathbf{x}^{(3)}) &\leq C \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(2)}) \leq C \frac{1}{256C^2} = \frac{1}{256C} = \frac{1}{2^{2^3} C}, \\ \tan \angle(\mathbf{e}, \mathbf{x}^{(m)}) &\leq \frac{1}{2^{2^m} C} \quad \text{für alle } m \in \mathbb{N}_0. \end{aligned}$$

Für den Fall einer normalen Matrix können wir entsprechend mit Satz 5.37 verfahren: Wir setzen

$$\tan^2 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq \frac{1}{2C}$$

voraus und erhalten

$$\begin{aligned} \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(1)}) &\leq C^2 \tan^6 \angle(\mathbf{e}, \mathbf{x}^{(0)}) \leq C^2 \frac{1}{2^3 C^3} = \frac{1}{2^{3^1} C}, \\ \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(2)}) &\leq C^2 \tan^6 \angle(\mathbf{e}, \mathbf{x}^{(1)}) \leq C^2 \frac{1}{2^{3^2} C^3} = \frac{1}{2^{3^2} C}, \\ \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(3)}) &\leq C^2 \tan^6 \angle(\mathbf{e}, \mathbf{x}^{(2)}) \leq C^2 \frac{1}{2^{3^3} C^3} = \frac{1}{2^{3^3} C}, \\ \tan^2 \angle(\mathbf{e}, \mathbf{x}^{(m)}) &\leq \frac{1}{2^{3^m} C}, \\ \tan \angle(\mathbf{e}, \mathbf{x}^{(m)}) &\leq \frac{1}{2^{3^m/2} \sqrt{C}} \quad \text{für alle } m \in \mathbb{N}_0. \end{aligned}$$

Da die Rayleigh-Iteration höhere Ansprüche an den Startvektor stellt als die inverse Iteration, kann es in der Praxis durchaus sinnvoll sein, zunächst einige Schritte der inversen Iteration mit einem festen Shift-Parameter durchzuführen, bevor man zu Rayleigh-Shifts wechselt.

## 5.5 Orthogonale Iteration

Die Vektoriteration und die von ihr abgeleiteten Verfahren dienen der Berechnung eines Eigenvektors zu einem bestimmten Eigenwert, der in geeigneter Weise dominant sein muss: Bei der einfachen Vektoriteration mussten wir voraussetzen, dass  $|\lambda_1|$  echt größer als die Beträge aller anderen Eigenwerte ist. Insbesondere mussten wir dabei mehrfache Eigenwerte ausschließen, denn im Falle  $|\lambda_1| = |\lambda_2|$  würde unser Konvergenzsatz 5.4 keine brauchbare Fehlerabschätzung mehr zur Verfügung stellen. Falls immerhin noch  $\lambda_1 \neq \lambda_2$  gilt, können wir dieses Problem mit einem geeigneten Shift-Parameter beheben, im Falle  $\lambda_1 = \lambda_2$  dagegen, also bei einem doppelten Eigenwert, erhalten wir mit der bisherigen Theorie keine Konvergenzaussage.

Wenn wir eine Diagonalmatrix  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  mit

$$|\lambda_1| = |\lambda_2| > |\lambda_3| \geq \dots \geq |\lambda_n|$$

näher untersuchen, stellen wir fest, dass in den Vektoren

$$\mathbf{x}^{(m)} := \frac{1}{|\lambda_1|^m} \mathbf{D}^m \mathbf{x}^{(0)} \quad \text{für } m \in \mathbb{N}_0$$

immer noch alle Komponenten außer den *ersten beiden* gegen null konvergieren. Die Vektoren werden also zwar nicht gegen einen Vektor des Eigenwerts  $\lambda_1$  konvergieren, aber immer noch gegen einen Vektor aus dem von den Eigenvektoren zu  $\lambda_1$  und  $\lambda_2$  aufgespannten *invarianten Teilraum*.

Konvergenzaussagen wie Satz 5.4 basieren darauf, dass die Iterierte  $\mathbf{x}^{(m)}$  mit einem geeigneten Vielfachen des ersten Eigenvektors  $\mathbf{e}$  verglichen wird. Wenn wir die Konvergenz gegen einen Teilraum untersuchen wollen, liegt es also nahe  $\mathbf{x}^{(m)}$  mit einem Vektor aus dem Teilraum zu vergleichen, der ihm möglichst nahe liegt.

Ein einfacher Zugang besteht darin, die Approximation einer Iterierten durch eine Matrix  $\mathbf{P}$  zu beschreiben, deren Bild der gewünschte Teilraum ist. Dann vergleichen wir die Iterierte  $\mathbf{x}^{(m)}$  mit dem Element  $\mathbf{P}\mathbf{x}^{(m)}$  des Teilraums, und falls die Differenz klein ist, liegt  $\mathbf{x}^{(m)}$  „fast“ im Teilraum. Besonders günstig ist es, wenn  $\mathbf{P}$  eine Projektion auf den Teilraum ist, also  $\mathbf{P}^2 = \mathbf{P}$  gilt. Für unsere Zwecke ideal unter derartigen Matrizen sind die *orthogonalen Projektionen*.

**Definition 5.38 (Orthogonale Projektion)** Sei  $\mathbf{P} \in \mathbb{K}^{n \times n}$ . Falls  $\mathbf{P}^2 = \mathbf{P} = \mathbf{P}^*$  gilt, nennen wir  $\mathbf{P}$  eine orthogonale Projektion.

**Lemma 5.39 (Orthogonale Projektion)** Sei  $\mathbf{P} \in \mathbb{K}^{n \times n}$  eine orthogonale Projektion. Dann gilt

$$\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 + \|\mathbf{P}\mathbf{x} - \mathbf{y}\|^2 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n, \mathbf{y} \in \text{Bild}(\mathbf{P}). \quad (5.16a)$$

Daraus folgen

$$\|\mathbf{x} - \mathbf{P}\mathbf{x}\| \leq \|\mathbf{x} - \mathbf{y}\| \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n, \mathbf{y} \in \text{Bild}(\mathbf{P}), \quad (5.16b)$$

$$\|\mathbf{P}\mathbf{x}\|^2 = \|\mathbf{x}\|^2 - \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 \leq \|\mathbf{x}\|^2 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n, \quad (5.16c)$$

die Projektion bildet also jeden Vektor auf seine beste Approximation in ihrem Bild ab und vergrößert dabei nicht seine Norm.

*Beweis.* Seien  $\mathbf{x} \in \mathbb{K}^n$  und  $\mathbf{y} \in \text{Bild}(\mathbf{P})$  gegeben. Dann existiert ein  $\mathbf{z} \in \mathbb{K}^n$  mit  $\mathbf{y} = \mathbf{P}\mathbf{z}$ . Wir haben

$$\begin{aligned} \|\mathbf{x} - \mathbf{y}\|^2 &= \langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle \\ &= \langle (\mathbf{x} - \mathbf{P}\mathbf{x}) + (\mathbf{P}\mathbf{x} - \mathbf{y}), (\mathbf{x} - \mathbf{P}\mathbf{x}) + (\mathbf{P}\mathbf{x} - \mathbf{y}) \rangle \\ &= \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 + \langle \mathbf{x} - \mathbf{P}\mathbf{x}, \mathbf{P}\mathbf{x} - \mathbf{y} \rangle + \langle \mathbf{P}\mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{P}\mathbf{x} \rangle + \|\mathbf{P}\mathbf{x} - \mathbf{y}\|^2 \\ &= \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 + \langle \mathbf{x} - \mathbf{P}\mathbf{x}, \mathbf{P}(\mathbf{x} - \mathbf{z}) \rangle + \langle \mathbf{P}(\mathbf{x} - \mathbf{z}), \mathbf{x} - \mathbf{P}\mathbf{x} \rangle + \|\mathbf{P}\mathbf{x} - \mathbf{y}\|^2 \\ &= \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 + \langle \mathbf{P}^*(\mathbf{x} - \mathbf{P}\mathbf{x}), \mathbf{x} - \mathbf{z} \rangle + \langle \mathbf{x} - \mathbf{z}, \mathbf{P}^*(\mathbf{x} - \mathbf{P}\mathbf{x}) \rangle + \|\mathbf{P}\mathbf{x} - \mathbf{y}\|^2 \\ &= \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 + \langle \mathbf{P}(\mathbf{x} - \mathbf{P}\mathbf{x}), \mathbf{x} - \mathbf{z} \rangle + \langle \mathbf{x} - \mathbf{z}, \mathbf{P}(\mathbf{x} - \mathbf{P}\mathbf{x}) \rangle + \|\mathbf{P}\mathbf{x} - \mathbf{y}\|^2 \\ &= \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 + \langle \mathbf{P}\mathbf{x} - \mathbf{P}^2\mathbf{x}, \mathbf{x} - \mathbf{z} \rangle + \langle \mathbf{x} - \mathbf{z}, \mathbf{P}\mathbf{x} - \mathbf{P}^2\mathbf{x} \rangle + \|\mathbf{P}\mathbf{x} - \mathbf{y}\|^2 \\ &= \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 + \langle \mathbf{0}, \mathbf{x} - \mathbf{z} \rangle + \langle \mathbf{x} - \mathbf{z}, \mathbf{0} \rangle + \|\mathbf{P}\mathbf{x} - \mathbf{y}\|^2 \\ &= \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 + \|\mathbf{P}\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

Daraus folgt unmittelbar (5.16b). Durch Einsetzen von  $\mathbf{y} = \mathbf{0}$  folgt (5.16c).  $\blacksquare$

**Lemma 5.40 (Existenz und Eindeutigkeit)** Sei  $\mathcal{V} \subseteq \mathbb{K}^n$  ein Teilraum. Dann existiert genau eine orthogonale Projektion  $\mathbf{P} \in \mathbb{K}^{n \times n}$  mit  $\text{Bild}(\mathbf{P}) = \mathcal{V}$ . Wir nennen sie die orthogonale Projektion auf  $\mathcal{V}$ .

*Beweis.* Zunächst zeigen wir, dass eine orthogonale Projektion existiert. Falls  $\mathcal{V} = \{\mathbf{0}\}$  gilt, setzen wir  $\mathbf{P} = \mathbf{0}$  und sind fertig.

Anderenfalls sei  $k \in [1 : n]$  die Dimension des Raums  $\mathcal{V}$ . Indem wir eine Basis  $(\mathbf{v}_i)_{i=1}^k$  des Raums wählen und ihre Elemente als Spalten einer Matrix  $\mathbf{V} \in \mathbb{K}^{n \times k}$  verwenden, erhalten wir eine injektive Matrix mit  $\text{Bild}(\mathbf{V}) = \mathcal{V}$ .

Nach Lemma 3.23 ist  $\mathbf{V}^*\mathbf{V}$  positiv definit, also insbesondere invertierbar. Wir definieren nun

$$\mathbf{P} := \mathbf{V}(\mathbf{V}^*\mathbf{V})^{-1}\mathbf{V}^*$$

und stellen mit Lemma 3.18 fest, dass  $\mathbf{P} = \mathbf{P}^*$  gilt. Es gilt auch

$$\mathbf{P}^2 = \mathbf{V} \underbrace{(\mathbf{V}^*\mathbf{V})^{-1}\mathbf{V}^*\mathbf{V}}_{=\mathbf{I}} (\mathbf{V}^*\mathbf{V})^{-1}\mathbf{V}^* = \mathbf{V}(\mathbf{V}^*\mathbf{V})^{-1}\mathbf{V}^* = \mathbf{P},$$

also haben wir eine orthogonale Projektion gefunden.

Sei nun  $\widehat{\mathbf{P}} \in \mathbb{K}^{n \times n}$  eine weitere orthogonale Projektion auf  $\mathcal{V}$ . Sei  $\mathbf{x} \in \mathbb{K}^n$ . Wir setzen

$$\mathbf{y} := \mathbf{P}\mathbf{x}, \quad \widehat{\mathbf{y}} := \widehat{\mathbf{P}}\mathbf{x}$$

## 5 Die Vektoriteration

und erhalten mit Lemma 5.39

$$\begin{aligned}\|\mathbf{x} - \widehat{\mathbf{y}}\|^2 &= \|\mathbf{x} - \mathbf{y}\|^2 + \|\mathbf{y} - \widehat{\mathbf{y}}\|^2, \\ \|\mathbf{x} - \mathbf{y}\|^2 &= \|\mathbf{x} - \widehat{\mathbf{y}}\|^2 + \|\widehat{\mathbf{y}} - \mathbf{y}\|^2.\end{aligned}$$

Indem wir die erste Gleichung in die zweite einsetzen, folgt

$$\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x} - \widehat{\mathbf{y}}\|^2 + 2\|\mathbf{y} - \widehat{\mathbf{y}}\|^2,$$

also muss  $\mathbf{P}\mathbf{x} = \mathbf{y} = \widehat{\mathbf{y}} = \widehat{\mathbf{P}}\mathbf{x}$  gelten. Da wir diese Gleichung für beliebige  $\mathbf{x} \in \mathbb{K}^n$  bewiesen haben, folgt  $\mathbf{P} = \widehat{\mathbf{P}}$ . ■

**Lemma 5.41 (Konvergenz für diagonale Matrizen)** Sei  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  gegeben und gelte

$$|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_k| > |\lambda_{k+1}| \geq \dots \geq |\lambda_n| \quad (5.17)$$

für ein  $k \in [1 : n - 1]$ . Sei  $\mathbf{P} \in \mathbb{K}^{n \times n}$  die orthogonale Projektion auf den Teilraum  $\mathbb{K}^k \times \{\mathbf{0}\}$ , der von den ersten  $k$  Eigenvektoren der Matrix  $\mathbf{D}$  aufgespannt wird.

Sei  $\mathbf{x}^{(0)} \in \mathbb{K}^n$  mit  $\mathbf{P}\mathbf{x} \neq \mathbf{0}$  gegeben. Dann gilt für die Iterierten

$$\mathbf{x}^{(m)} := \mathbf{D}^m \mathbf{x}^{(0)} \quad \text{für alle } m \in \mathbb{N}.$$

die Abschätzung

$$\frac{\|\mathbf{x}^{(m)} - \mathbf{P}\mathbf{x}^{(m)}\|}{\|\mathbf{P}\mathbf{x}^{(m)}\|} \leq \left( \frac{|\lambda_{k+1}|}{|\lambda_k|} \right)^m \frac{\|\mathbf{x}^{(0)} - \mathbf{P}\mathbf{x}^{(0)}\|}{\|\mathbf{P}\mathbf{x}^{(0)}\|} \quad \text{für alle } m \in \mathbb{N}_0.$$

*Beweis.* Sei  $\mathbf{I}_k \in \mathbb{K}^{k \times k}$  die Identität auf  $\mathbb{K}^k$ , und sei

$$\mathbf{P} = \begin{pmatrix} \mathbf{I}_k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \in \mathbb{K}^{n \times n}$$

ihre Fortsetzung durch null. Dann gilt

$$\mathbf{P}\mathbf{x}^{(m)} = \begin{pmatrix} \lambda_1^m x_1^{(0)} \\ \dots \\ \lambda_k^m x_k^{(0)} \\ 0 \\ \dots \\ 0 \end{pmatrix}, \quad \mathbf{x}^{(m)} - \mathbf{P}\mathbf{x}^{(m)} = \begin{pmatrix} 0 \\ \dots \\ 0 \\ \lambda_{k+1}^m x_{k+1}^{(0)} \\ \dots \\ \lambda_n^m x_n^{(0)} \end{pmatrix} \quad \text{für alle } m \in \mathbb{N}_0.$$

Für die Normen erhalten wir

$$\|\mathbf{x}^{(m)} - \mathbf{P}\mathbf{x}^{(m)}\|^2 = \sum_{i=k+1}^n |\lambda_i^m x_i^{(0)}|^2 \leq |\lambda_{k+1}|^{2m} \sum_{i=k+1}^n |x_i^{(0)}|^2 \leq |\lambda_{k+1}|^{2m} \|\mathbf{x}^{(0)} - \mathbf{P}\mathbf{x}^{(0)}\|^2,$$

$$\|\mathbf{P}\mathbf{x}^{(m)}\|^2 = \sum_{i=1}^k |\lambda_i^m x_i^{(0)}|^2 \geq |\lambda_k|^{2m} \sum_{i=1}^k |x_i^{(0)}|^2 = |\lambda_k|^{2m} \|\mathbf{P}\mathbf{x}^{(0)}\|^2,$$

also folgt

$$\frac{\|\mathbf{x}^{(m)} - \mathbf{P}\mathbf{x}^{(m)}\|}{\|\mathbf{P}\mathbf{x}^{(m)}\|} \leq \left( \frac{|\lambda_{k+1}|}{|\lambda_k|} \right)^m \frac{\|\mathbf{x}^{(0)} - \mathbf{P}\mathbf{x}^{(0)}\|}{\|\mathbf{P}\mathbf{x}^{(0)}\|} \quad \text{für alle } m \in \mathbb{N}_0,$$

und das ist die zu zeigende Abschätzung.  $\blacksquare$

Auch dieses Ergebnis können wir durch einen Winkelbegriff ausdrücken: Analog zu Lemma 5.5 definieren wir für einen Vektor  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  und einen Teilraum  $\mathcal{Y} \subseteq \mathbb{K}^n$  den Winkel durch

$$\sin \angle(\mathbf{x}, \mathcal{Y}) := \min \left\{ \frac{\|\mathbf{x} - \mathbf{y}\|}{\|\mathbf{x}\|} : \mathbf{y} \in \mathcal{Y} \right\}.$$

Falls  $\mathbf{P} \in \mathbb{K}^{n \times n}$  die orthogonale Projektion auf  $\mathcal{Y}$  ist, können wir mit Lemma 5.39 das Minimum explizit darstellen und erhalten

$$\sin \angle(\mathbf{x}, \mathcal{Y}) = \frac{\|\mathbf{x} - \mathbf{P}\mathbf{x}\|}{\|\mathbf{x}\|}.$$

Indem wir (5.16a) auf  $\mathbf{y} = \mathbf{0}$  anwenden, erhalten wir  $\|\mathbf{x}\|^2 = \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2 + \|\mathbf{P}\mathbf{x}\|^2$  und können den Cosinus des Winkels durch

$$\cos^2 \angle(\mathbf{x}, \mathcal{Y}) = 1 - \sin^2 \angle(\mathbf{x}, \mathcal{Y}) = \frac{\|\mathbf{x}\|^2 - \|\mathbf{x} - \mathbf{P}\mathbf{x}\|^2}{\|\mathbf{x}\|^2} = \frac{\|\mathbf{P}\mathbf{x}\|^2}{\|\mathbf{x}\|^2}$$

einführen. Insgesamt haben wir

$$\sin \angle(\mathbf{x}, \mathcal{Y}) = \frac{\|\mathbf{x} - \mathbf{P}\mathbf{x}\|}{\|\mathbf{x}\|}, \quad \cos \angle(\mathbf{x}, \mathcal{Y}) = \frac{\|\mathbf{P}\mathbf{x}\|}{\|\mathbf{x}\|}, \quad \tan \angle(\mathbf{x}, \mathcal{Y}) = \frac{\|\mathbf{x} - \mathbf{P}\mathbf{x}\|}{\|\mathbf{P}\mathbf{x}\|}. \quad (5.18)$$

Damit können wir die Aussage des Lemmas 5.41 in die uns vertraute Form bringen:

**Folgerung 5.42 (Konvergenz für diagonale Matrizen)** Sei  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  gegeben mit (5.17) für ein  $k \in [1 : n - 1]$ . Sei  $\mathcal{E} := \mathbb{K}^k \times \{\mathbf{0}\}$  der von den ersten  $k$  Eigenvektoren der Matrix  $\mathbf{D}$  aufgespannte invariante Teilraum.

Sei  $\mathbf{x}^{(0)} \in \mathbb{K}^n$  mit  $\cos \angle(\mathbf{x}^{(0)}, \mathcal{E}) > 0$  gegeben. Dann gilt für die Iterierten  $(\mathbf{x}^{(m)})_{m=0}^\infty$  der Vektoriteration die Abschätzung

$$\tan \angle(\mathbf{x}^{(m)}, \mathcal{E}) \leq \left( \frac{|\lambda_{k+1}|}{|\lambda_k|} \right)^m \tan \angle(\mathbf{x}^{(0)}, \mathcal{E}) \quad \text{für alle } m \in \mathbb{N}_0.$$

*Beweis.* Wir kombinieren Lemma 5.41 mit (5.18).  $\blacksquare$

Falls nicht zufällig  $\lambda_1 = \dots = \lambda_k$  gilt, erhalten wir in diesem Fall *keine* Konvergenz gegen einen Eigenraum, sondern nur gegen den invarianten Teilraum  $\mathcal{E}$ .

## 5 Die Vektoriteration

Unser Ziel ist es nun, wenigstens diesen Teilraum vollständig zu beschreiben, indem wir eine Basis konstruieren. Die Idee ist einfach: Wenn die Vektoriteration zu *einem* Startvektor gegen *einen* Vektor aus dem Unterraum konvergiert, dann wird die Vektoriteration zu *k* Startvektoren gegen *k* Vektoren aus dem Unterraum konvergieren, und falls wir sicherstellen können, dass diese Vektoren linear unabhängig sind, erhalten wir eine Basis.

Zur Abkürzung der Notation fassen wir die Vektoren in einer Matrix zusammen: Die *k* Spalten einer Matrix  $\mathbf{X}^{(m)} \in \mathbb{K}^{n \times k}$  interpretieren wir als die Iterierten von *k* simultan ausgeführten Vektoriterationen mit den *k* Spalten der Matrix  $\mathbf{X}^{(0)} \in \mathbb{K}^{n \times k}$  als Startvektoren. Die Iterierten sind dann durch

$$\mathbf{X}^{(m)} := \mathbf{D}^m \mathbf{X}^{(0)} \quad \text{für alle } m \in \mathbb{N} \quad (5.19)$$

definiert. Es stellt sich die Frage, unter welchen Bedingungen wir darauf hoffen dürfen, dass die Spalten von  $\mathbf{X}^{(m)}$  linear unabhängig sind. Offenbar müssen dafür zumindest die Spalten von  $\mathbf{X}^{(0)}$  linear unabhängig sein, die Matrix muss also vollen Rang besitzen. Das reicht allerdings noch nicht: Falls eine Linearkombination der Spalten in den Kern der Matrix  $\mathbf{D}^m$  fallen sollte, würde  $\mathbf{X}^{(m)}$  trotzdem keinen vollen Rang besitzen.

Dieser Fall lässt sich einfach ausschließen, indem wir die in Lemma 5.41 eingeführte Projektion  $\mathbf{P}$  auf den Unterraum verwenden: Statt zu fordern, dass  $\mathbf{X}^{(0)}$  vollen Rang hat, fordern wir diese Eigenschaft von  $\mathbf{P}\mathbf{X}^{(0)}$ . Diese Voraussetzung ist ausreichend:

**Lemma 5.43 (Basis)** *Seien  $\mathbf{D}, \mathbf{P}, \lambda_1, \dots, \lambda_n$  wie in Lemma 5.41 gegeben. Sei  $\mathbf{X}^{(0)} \in \mathbb{K}^{n \times k}$  so gegeben, dass  $\mathbf{P}\mathbf{X}^{(0)}$  vollen Rang hat.*

*Dann haben für jedes  $m \in \mathbb{N}_0$  auch die Matrizen  $\mathbf{P}\mathbf{X}^{(m)}$  und  $\mathbf{X}^{(m)}$  vollen Rang.*

*Beweis.* Sei  $m \in \mathbb{N}_0$ . Dazu führen wir

$$\widehat{\mathbf{D}} := \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_k \end{pmatrix}, \quad \widehat{\mathbf{X}}^{(m)} := \begin{pmatrix} x_{11}^{(m)} & \dots & x_{1k}^{(m)} \\ \vdots & \ddots & \vdots \\ x_{k1}^{(m)} & \dots & x_{kk}^{(m)} \end{pmatrix}$$

ein und können  $\mathbf{P}^2 = \mathbf{P}$  sowie  $\mathbf{P}\mathbf{D} = \mathbf{D}\mathbf{P}$  ausnutzen, um

$$\mathbf{P}\mathbf{X}^{(m)} = \mathbf{P}\mathbf{D}^m \mathbf{X}^{(0)} = \mathbf{P}^2 \mathbf{D}^m \mathbf{X}^{(0)} = \mathbf{P}\mathbf{D}^m \mathbf{P}\mathbf{X}^{(0)} = \begin{pmatrix} \widehat{\mathbf{D}}^m & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{X}}^{(0)} \\ \mathbf{0} \end{pmatrix} = \begin{pmatrix} \widehat{\mathbf{D}}^m \widehat{\mathbf{X}}^{(0)} \\ \mathbf{0} \end{pmatrix}$$

zu erhalten. Nach Voraussetzung hat

$$\mathbf{P}\mathbf{X}^{(0)} = \begin{pmatrix} \widehat{\mathbf{X}}^{(0)} \\ \mathbf{0} \end{pmatrix}$$

vollen Rang, also muss  $\widehat{\mathbf{X}}^{(0)}$  regulär sein. Aus (5.17) folgt, dass  $\widehat{\mathbf{D}}$  regulär ist, also muss auch  $\widehat{\mathbf{D}}^m$  regulär sein, und somit auch  $\widehat{\mathbf{D}}^m \widehat{\mathbf{X}}^{(0)}$ . Damit ist bewiesen, dass  $\mathbf{P}\mathbf{X}^{(m)}$  vollen Rang hat.



Die Dimension des Bildraums dieses Produkts ist also  $k$ , und damit müssen auch die Dimensionen der Bildräume der Faktoren mindestens  $k$  sein. Also hat insbesondere  $\mathbf{X}^{(m)}$  vollen Rang. ■

Indem wir Lemma 5.41 und Lemma 5.43 kombinieren, können wir folgern, dass die Matrizen  $\mathbf{X}^{(m)}$  gegen Basen des für uns interessanten invarianten Unterraums konvergieren werden.

Dieses Konvergenzverhalten können wir auch quantifizieren: Nach Lemma 3.27 bilden die Spalten einer Matrix  $\mathbf{X} \in \mathbb{K}^{n \times k}$  genau dann eine Basis eines invarianten Unterraums, wenn es eine Matrix  $\mathbf{\Lambda} \in \mathbb{K}^{k \times k}$  gibt, mit der  $\mathbf{DX} = \mathbf{X}\mathbf{\Lambda}$  gilt. Wenn  $\mathbf{X}^{(m)}$  „fast“ eine solche Basis ist, sollte eine Eigenschaft der Form

$$\|\mathbf{DX}^{(m)} - \mathbf{X}^{(m)}\mathbf{\Lambda}\| \leq \epsilon$$

für ein geeignetes  $\epsilon \in \mathbb{R}_{>0}$  gelten. Wenn wir die Skalierung der Matrix  $\mathbf{X}^{(m)}$  geeignet berücksichtigen, erhalten wir das folgende Resultat:

**Lemma 5.44 (Konvergenz gegen Teilraum)** *Seien  $\mathbf{D}, \mathbf{P}, \lambda_1, \dots, \lambda_n$  wie in Lemma 5.41 gegeben. Sei  $\mathbf{X}^{(0)} \in \mathbb{K}^{n \times k}$  so gegeben, dass  $\mathbf{PX}^{(0)}$  vollen Rang hat.*

*Dann gibt es eine Matrix  $\mathbf{\Lambda} \in \mathbb{K}^{k \times k}$  so, dass für alle  $\mathbf{y} \in \mathbb{K}^k \setminus \{\mathbf{0}\}$  und alle  $m \in \mathbb{N}$  die Abschätzungen*

$$\begin{aligned} \frac{\|(\mathbf{DX}^{(m)} - \mathbf{X}^{(m)}\mathbf{\Lambda})\mathbf{y}\|}{\|\mathbf{PX}^{(m)}\mathbf{y}\|} &\leq \left(\frac{|\lambda_{k+1}|}{|\lambda_k|}\right)^m \frac{\|(\mathbf{DX}^{(0)} - \mathbf{X}^{(0)}\mathbf{\Lambda})\mathbf{y}\|}{\|\mathbf{PX}^{(0)}\mathbf{y}\|}, \\ \frac{\|(\mathbf{X}^{(m)} - \mathbf{PX}^{(m)})\mathbf{y}\|}{\|\mathbf{PX}^{(m)}\mathbf{y}\|} &\leq \left(\frac{|\lambda_{k+1}|}{|\lambda_k|}\right)^m \frac{\|(\mathbf{X}^{(0)} - \mathbf{PX}^{(0)})\mathbf{y}\|}{\|\mathbf{PX}^{(0)}\mathbf{y}\|} \end{aligned}$$

erfüllt sind.

*Beweis.* Wir definieren Hilfsmatrizen

$$\begin{aligned} \widehat{\mathbf{X}} &:= \begin{pmatrix} x_{11}^{(0)} & \cdots & x_{1k}^{(0)} \\ \vdots & \ddots & \vdots \\ x_{k1}^{(0)} & \cdots & x_{kk}^{(0)} \end{pmatrix}, & \mathbf{X}_\perp &:= \begin{pmatrix} x_{k+1,1}^{(0)} & \cdots & x_{k+1,k}^{(0)} \\ \vdots & \ddots & \vdots \\ x_{n,1}^{(0)} & \cdots & x_{n,k}^{(0)} \end{pmatrix}, \\ \widehat{\mathbf{D}} &:= \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_k \end{pmatrix}, & \mathbf{D}_\perp &:= \begin{pmatrix} \lambda_{k+1} & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \end{aligned}$$

und erhalten

$$\mathbf{X}^{(0)} = \begin{pmatrix} \widehat{\mathbf{X}} \\ \mathbf{X}_\perp \end{pmatrix}, \quad \mathbf{PX}^{(0)} = \begin{pmatrix} \widehat{\mathbf{X}} \\ \mathbf{0} \end{pmatrix}, \quad \mathbf{D} = \begin{pmatrix} \widehat{\mathbf{D}} & \\ & \mathbf{D}_\perp \end{pmatrix}.$$

Nach Voraussetzung hat  $\mathbf{PX}^{(0)}$  vollen Rang, also muss  $\widehat{\mathbf{X}} \in \mathbb{K}^{k \times k}$  regulär sein. Wir definieren

$$\mathbf{\Lambda} := \widehat{\mathbf{X}}^{-1}\widehat{\mathbf{D}}\widehat{\mathbf{X}}$$

## 5 Die Vektoriteration

und erhalten

$$\widehat{\mathbf{X}}\mathbf{\Lambda} = \widehat{\mathbf{X}}\widehat{\mathbf{X}}^{-1}\widehat{\mathbf{D}}\widehat{\mathbf{X}} = \widehat{\mathbf{D}}\widehat{\mathbf{X}}.$$

Nun können wir die Terme unserer Abschätzung untersuchen. Sei  $\mathbf{y} \in \mathbb{K}^k \setminus \{\mathbf{0}\}$  und  $m \in \mathbb{N}$ . Es gilt

$$\begin{aligned} \|(\mathbf{D}\mathbf{X}^{(m)} - \mathbf{X}^{(m)}\mathbf{\Lambda})\mathbf{y}\| &= \|\mathbf{D}^{m+1}\mathbf{X}^{(0)}\mathbf{y} - \mathbf{D}^m\mathbf{X}^{(0)}\mathbf{\Lambda}\mathbf{y}\| \\ &= \left\| \begin{pmatrix} (\widehat{\mathbf{D}}^{m+1}\widehat{\mathbf{X}} - \widehat{\mathbf{D}}^m\widehat{\mathbf{X}}\mathbf{\Lambda})\mathbf{y} \\ (\mathbf{D}_\perp^{m+1}\mathbf{X}_\perp - \mathbf{D}_\perp^m\mathbf{X}_\perp\mathbf{\Lambda})\mathbf{y} \end{pmatrix} \right\| = \left\| \begin{pmatrix} (\widehat{\mathbf{D}}^{m+1}\widehat{\mathbf{X}} - \widehat{\mathbf{D}}^{m+1}\widehat{\mathbf{X}})\mathbf{y} \\ (\mathbf{D}_\perp^{m+1}\mathbf{X}_\perp - \mathbf{D}_\perp^m\mathbf{X}_\perp\mathbf{\Lambda})\mathbf{y} \end{pmatrix} \right\| \\ &= \left\| \begin{pmatrix} \mathbf{0} \\ \mathbf{D}_\perp^m(\mathbf{D}_\perp\mathbf{X}_\perp - \mathbf{X}_\perp\mathbf{\Lambda})\mathbf{y} \end{pmatrix} \right\| \leq |\lambda_{k+1}|^m \|(\mathbf{D}_\perp\mathbf{X}_\perp - \mathbf{X}_\perp\mathbf{\Lambda})\mathbf{y}\| \\ &= |\lambda_{k+1}|^m \left\| \begin{pmatrix} (\widehat{\mathbf{D}}\widehat{\mathbf{X}} - \widehat{\mathbf{X}}\mathbf{\Lambda})\mathbf{y} \\ (\mathbf{D}_\perp\mathbf{X}_\perp - \mathbf{X}_\perp\mathbf{\Lambda})\mathbf{y} \end{pmatrix} \right\| = |\lambda_{k+1}|^m \|(\mathbf{D}\mathbf{X}^{(0)} - \mathbf{X}^{(0)}\mathbf{\Lambda})\mathbf{y}\|, \end{aligned}$$

und wir erhalten außerdem

$$\|\mathbf{P}\mathbf{X}^{(m)}\mathbf{y}\| = \left\| \begin{pmatrix} \widehat{\mathbf{D}}^m\widehat{\mathbf{X}}\mathbf{y} \\ \mathbf{0} \end{pmatrix} \right\| \geq |\lambda_k|^m \left\| \begin{pmatrix} \widehat{\mathbf{X}}\mathbf{y} \\ \mathbf{0} \end{pmatrix} \right\| = |\lambda_k|^m \|\mathbf{P}\mathbf{X}^{(0)}\mathbf{y}\|. \quad (5.20)$$

Einsetzen dieser Ungleichungen in den zu untersuchenden Bruch führt zu der ersten Abschätzung.

Für die zweite Abschätzung verwenden wir einfach

$$\begin{aligned} \|(\mathbf{X}^{(m)} - \mathbf{P}\mathbf{X}^{(m)})\mathbf{y}\| &= \left\| \begin{pmatrix} \widehat{\mathbf{D}}^m\widehat{\mathbf{X}} - \widehat{\mathbf{D}}^m\widehat{\mathbf{X}} \\ \mathbf{D}_\perp^m\mathbf{X}_\perp \end{pmatrix} \mathbf{y} \right\| = \|\mathbf{D}_\perp^m\mathbf{X}_\perp\mathbf{y}\| \leq |\lambda_{k+1}|^m \|\mathbf{X}_\perp\mathbf{y}\| \\ &= |\lambda_{k+1}|^m \left\| \begin{pmatrix} \widehat{\mathbf{X}} - \widehat{\mathbf{X}} \\ \mathbf{X}_\perp \end{pmatrix} \mathbf{y} \right\| = |\lambda_{k+1}|^m \|(\mathbf{X}^{(0)} - \mathbf{P}\mathbf{X}^{(0)})\mathbf{y}\| \end{aligned}$$

in Kombination mit (5.20). ■

Wie üblich können wir diese Konvergenzaussage auf den Fall normaler Matrizen  $\mathbf{A} \in \mathbb{K}^{n \times n}$  übertragen, indem wir sie mit Hilfe der Folgerung 3.54 diagonalisieren.

**Satz 5.45 (Konvergenz)** *Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine normale Matrix. Sie besitzt eine Schurzerlegung  $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^*$  mit einer unitären Matrix  $\mathbf{Q} \in \mathbb{C}^{n \times n}$  und einer Diagonalmatrix*

$$\mathbf{D} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \quad |\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_n|.$$

*Sei  $k \in [1 : n - 1]$  so gegeben, dass  $|\lambda_k| > |\lambda_{k+1}|$  gilt, sei  $\mathcal{E} := \mathbf{Q}(\mathbb{K}^k \times \{\mathbf{0}\})$  der von den ersten  $k$  Eigenvektoren aufgespannte invariante Teilraum von  $\mathbb{K}^n$ , und sei  $\mathbf{P} \in \mathbb{K}^{n \times n}$  die Projektion auf diesen Teilraum.*

*Sei  $\mathbf{X}^{(0)} \in \mathbb{K}^{n \times k}$  so gegeben, dass die Matrix  $\mathbf{P}\mathbf{X}^{(0)}$  vollen Rang hat.*

Dann existiert eine Matrix  $\mathbf{\Lambda} \in \mathbb{K}^{k \times k}$  so, dass

$$\frac{\|(\mathbf{A}\mathbf{X}^{(m)} - \mathbf{X}^{(m)}\mathbf{\Lambda})\mathbf{y}\|}{\|\mathbf{P}\mathbf{X}^{(m)}\mathbf{y}\|} \leq \left(\frac{|\lambda_{k+1}|}{|\lambda_k|}\right)^m \frac{\|(\mathbf{A}\mathbf{X}^{(0)} - \mathbf{X}^{(0)}\mathbf{\Lambda})\mathbf{y}\|}{\|\mathbf{P}\mathbf{X}^{(0)}\mathbf{y}\|}, \quad (5.21a)$$

$$\frac{\|(\mathbf{X}^{(m)} - \mathbf{P}\mathbf{X}^{(m)})\mathbf{y}\|}{\|\mathbf{P}\mathbf{X}^{(m)}\mathbf{y}\|} \leq \left(\frac{|\lambda_{k+1}|}{|\lambda_k|}\right)^m \frac{\|(\mathbf{X}^{(0)} - \mathbf{P}\mathbf{X}^{(0)})\mathbf{y}\|}{\|\mathbf{P}\mathbf{X}^{(0)}\mathbf{y}\|} \quad (5.21b)$$

für alle  $m \in \mathbb{N}_0$  und  $\mathbf{y} \in \mathbb{K}^k \setminus \{\mathbf{0}\}$  gelten. Insbesondere haben die Matrizen  $\mathbf{P}\mathbf{X}^{(m)}$  und  $\mathbf{X}^{(m)}$  vollen Rang.

*Beweis.* Sei  $\widehat{\mathbf{P}} \in \mathbb{K}^{n \times n}$  die Projektion aus Lemma 5.41. Durch Rücktransformation erhalten wir die orthogonale Projektion

$$\mathbf{P} := \mathbf{Q}\widehat{\mathbf{P}}\mathbf{Q}^*$$

auf  $\mathcal{E}$ . Wie schon im Beweis von Satz 5.4 definieren wir durch

$$\widehat{\mathbf{X}}^{(m)} := \mathbf{Q}^*\mathbf{X}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0$$

die transformierte Folge der Iterierten, für die wegen

$$\widehat{\mathbf{P}}\widehat{\mathbf{X}}^{(0)} = \mathbf{Q}^*\mathbf{P}\mathbf{Q}\widehat{\mathbf{X}}^{(0)} = \mathbf{Q}^*\mathbf{P}\mathbf{X}^{(0)}$$

und der Gleichung

$$\widehat{\mathbf{X}}^{(m)} = \mathbf{Q}^*\mathbf{X}^{(m)} = \mathbf{Q}^*\mathbf{A}^m\mathbf{X}^{(0)} = \mathbf{Q}^*\mathbf{A}^m\mathbf{Q}\widehat{\mathbf{X}}^{(0)} = \mathbf{D}^m\widehat{\mathbf{X}}^{(0)} \quad \text{für alle } m \in \mathbb{N}_0$$

die Voraussetzungen des Lemmas 5.44 erfüllt sind. Also existiert ein  $\mathbf{\Lambda} \in \mathbb{K}^{k \times k}$  derart, dass die Abschätzungen

$$\begin{aligned} \frac{\|(\mathbf{D}\widehat{\mathbf{X}}^{(m)} - \widehat{\mathbf{X}}^{(m)}\mathbf{\Lambda})\mathbf{y}\|}{\|\widehat{\mathbf{P}}\widehat{\mathbf{X}}^{(m)}\mathbf{y}\|} &\leq \left(\frac{|\lambda_{k+1}|}{|\lambda_k|}\right)^m \frac{\|(\mathbf{D}\widehat{\mathbf{X}}^{(0)} - \widehat{\mathbf{X}}^{(0)}\mathbf{\Lambda})\mathbf{y}\|}{\|\widehat{\mathbf{P}}\widehat{\mathbf{X}}^{(0)}\mathbf{y}\|}, \\ \frac{\|(\widehat{\mathbf{X}}^{(m)} - \widehat{\mathbf{P}}\widehat{\mathbf{X}}^{(m)})\mathbf{y}\|}{\|\widehat{\mathbf{P}}\widehat{\mathbf{X}}^{(m)}\mathbf{y}\|} &\leq \left(\frac{|\lambda_{k+1}|}{|\lambda_k|}\right)^m \frac{\|(\widehat{\mathbf{X}}^{(0)} - \widehat{\mathbf{P}}\widehat{\mathbf{X}}^{(0)})\mathbf{y}\|}{\|\widehat{\mathbf{P}}\widehat{\mathbf{X}}^{(0)}\mathbf{y}\|} \end{aligned}$$

für alle  $\mathbf{y} \in \mathbb{K}^k \setminus \{\mathbf{0}\}$  und alle  $m \in \mathbb{N}$  gelten. Mit (3.15) erhalten wir

$$\begin{aligned} \|(\mathbf{D}\widehat{\mathbf{X}}^{(m)} - \widehat{\mathbf{X}}^{(m)}\mathbf{\Lambda})\mathbf{y}\| &= \|(\mathbf{Q}\mathbf{D}\mathbf{Q}^*\mathbf{X}^{(m)} - \mathbf{Q}\widehat{\mathbf{X}}^{(m)}\mathbf{\Lambda})\mathbf{y}\| = \|(\mathbf{A}\mathbf{X}^{(m)} - \mathbf{X}^{(m)}\mathbf{\Lambda})\mathbf{y}\|, \\ \|(\widehat{\mathbf{X}}^{(m)} - \widehat{\mathbf{P}}\widehat{\mathbf{X}}^{(m)})\mathbf{y}\| &= \|(\mathbf{Q}\widehat{\mathbf{X}}^{(m)} - \mathbf{Q}\widehat{\mathbf{P}}\mathbf{Q}^*\mathbf{X}^{(m)})\mathbf{y}\| = \|(\mathbf{X}^{(m)} - \mathbf{P}\mathbf{X}^{(m)})\mathbf{y}\|, \\ \|\widehat{\mathbf{P}}\widehat{\mathbf{X}}^{(m)}\mathbf{y}\| &= \|\mathbf{Q}\widehat{\mathbf{P}}\mathbf{Q}^*\mathbf{X}^{(m)}\mathbf{y}\| = \|\mathbf{P}\mathbf{X}^{(m)}\mathbf{y}\| \end{aligned}$$

für alle  $m \in \mathbb{N}_0$  und alle  $\mathbf{y} \in \mathbb{K}^k$ , und damit folgen die gewünschten Aussagen.  $\blacksquare$

Wie schon im Falle der konventionellen Vektoriteration kann auch bei dieser Iteration das Problem auftreten, dass es durch die wiederholte Multiplikation mit  $\mathbf{A}$  zu Über- oder

## 5 Die Vektoriteration

Unterläufen kommt. Sehr viel schlimmer ist allerdings, dass in der Regel alle Spalten der Matrizen  $\mathbf{X}^{(m)}$  gegen Vielfache des Eigenvektors eines dominanten Eigenwerts konvergieren werden, so dass sie zwar theoretisch linear unabhängig bleiben, in der numerischen Praxis aber fast linear abhängig werden können.

Beide Probleme lassen sich wieder durch eine geschickte Normierung beheben: Da wir nur an dem invarianten Unterraum, also dem Bild der Matrix  $\mathbf{X}^{(m)}$ , interessiert sind, können wir die Matrix fast beliebig skalieren und Linearkombinationen ihrer Spalten bilden, ohne diesen Unterraum zu verändern. Insbesondere können wir, beispielsweise mit Hilfe der Gram-Schmidt-Orthonormalisierung, dafür sorgen, dass die von den Spalten beschriebene Basis orthonormal ist, und damit insbesondere aus linear unabhängigen Einheitsvektoren besteht.

Da die Gram-Schmidt-Orthonormalisierung numerisch instabil sein kann, verwenden wir stattdessen allerdings die sicherere QR-Zerlegung: Zu jeder Matrix  $\mathbf{X}^{(m)} \in \mathbb{K}^{n \times k}$  existieren eine unitäre Matrix  $\widehat{\mathbf{Q}}^{(m)} \in \mathbb{K}^{n \times n}$  und eine obere Dreiecksmatrix  $\widehat{\mathbf{R}}^{(m)} \in \mathbb{K}^{n \times k}$  so, dass

$$\mathbf{X}^{(m)} = \widehat{\mathbf{Q}}^{(m)} \widehat{\mathbf{R}}^{(m)}$$

gilt. Wenn wir die ersten  $k$  Zeilen von  $\widehat{\mathbf{R}}^{(m)}$  mit  $\mathbf{R}^{(m)}$  sowie die ersten  $k$  Spalten von  $\widehat{\mathbf{Q}}^{(m)}$  mit  $\mathbf{Q}^{(m)}$  sowie die restlichen Spalten mit  $\widetilde{\mathbf{Q}}^{(m)}$  bezeichnen, erhalten wir

$$\mathbf{X}^{(m)} = \widehat{\mathbf{Q}}^{(m)} \widehat{\mathbf{R}}^{(m)} = \begin{pmatrix} \mathbf{Q}^{(m)} & \widetilde{\mathbf{Q}}^{(m)} \end{pmatrix} \begin{pmatrix} \mathbf{R}^{(m)} \\ \mathbf{0} \end{pmatrix} = \mathbf{Q}^{(m)} \mathbf{R}^{(m)},$$

da  $\widehat{\mathbf{R}}^{(m)}$  eine obere Dreiecksmatrix ist. Wir können also zu jedem  $\mathbf{X}^{(m)} \in \mathbb{K}^{n \times k}$  eine isometrische Matrix  $\mathbf{Q}^{(m)} \in \mathbb{K}^{n \times k}$  und eine obere Dreiecksmatrix  $\mathbf{R}^{(m)} \in \mathbb{K}^{k \times k}$  mit

$$\mathbf{X}^{(m)} = \mathbf{Q}^{(m)} \mathbf{R}^{(m)} \quad (5.22)$$

konstruieren, eine sogenannte *dünne QR-Zerlegung*. Die Spalten der Matrix  $\mathbf{Q}^{(m)}$  beschreiben dann die gesuchte Orthonormalbasis.

Für derartige Matrizen vereinfacht sich die Aussage von Satz 5.45 wesentlich, wenn wir berücksichtigen, dass durch die Orthogonalisierung in jedem Schritt die Basis verändert wird.

**Folgerung 5.46 (Orthogonale Iteration)** *Unter den Voraussetzungen von Satz 5.45 und mit (5.22) finden wir eine Familie  $(\widetilde{\boldsymbol{\Lambda}}^{(m)})_{m=0}^{\infty}$  in  $\mathbb{K}^{k \times k}$  mit der folgenden Eigenschaft: Für jedes  $m \in \mathbb{N}_0$  und jeden Vektor  $\mathbf{y} \in \mathbb{K}^k \setminus \{\mathbf{0}\}$  existiert ein Vektor  $\mathbf{z} \in \mathbb{K}^k \setminus \{\mathbf{0}\}$  mit*

$$\frac{\|(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\widetilde{\boldsymbol{\Lambda}}^{(m)})\mathbf{y}\|}{\|\mathbf{P}\mathbf{Q}^{(m)}\mathbf{y}\|} \leq \left( \frac{|\lambda_{k+1}|}{|\lambda_k|} \right)^m \frac{\|(\mathbf{A}\mathbf{Q}^{(0)} - \mathbf{Q}^{(0)}\widetilde{\boldsymbol{\Lambda}}^{(0)})\mathbf{z}\|}{\|\mathbf{P}\mathbf{Q}^{(0)}\mathbf{z}\|}, \quad (5.23a)$$

$$\frac{\|(\mathbf{Q}^{(m)} - \mathbf{P}\mathbf{Q}^{(m)})\mathbf{y}\|}{\|\mathbf{P}\mathbf{Q}^{(m)}\mathbf{y}\|} \leq \left( \frac{|\lambda_{k+1}|}{|\lambda_k|} \right)^m \frac{\|(\mathbf{Q}^{(0)} - \mathbf{P}\mathbf{Q}^{(0)})\mathbf{z}\|}{\|\mathbf{P}\mathbf{Q}^{(0)}\mathbf{z}\|}. \quad (5.23b)$$

*Beweis.* Sei  $\mathbf{X}^{(0)} \in \mathbb{K}^{n \times k}$  mit  $\text{rank}(\mathbf{P}\mathbf{X}^{(0)}) = k$  gegeben. Da  $\mathbf{X}^{(m)}$  nach Satz 5.45 für alle  $m \in \mathbb{N}_0$  vollen Rang hat, müssen nach (5.22) auch die Matrizen  $\mathbf{R}^{(m)}$  vollen Rang haben, also regulär sein. Damit ist

$$\widetilde{\boldsymbol{\Lambda}}^{(m)} := \mathbf{R}^{(m)} \boldsymbol{\Lambda} (\mathbf{R}^{(m)})^{-1} \quad \text{für alle } m \in \mathbb{N}_0$$

mit der Matrix  $\mathbf{\Lambda}$  aus demselben Satz wohldefiniert.

Seien nun  $m \in \mathbb{N}$  und ein Vektor  $\mathbf{y} \in \mathbb{K}^k \setminus \{\mathbf{0}\}$  fixiert. Wir definieren

$$\tilde{\mathbf{y}} := (\mathbf{R}^{(m)})^{-1}\mathbf{y}, \quad \mathbf{z} := \mathbf{R}^{(0)}\tilde{\mathbf{y}},$$

und erhalten mit (5.22) die Gleichungen

$$\begin{aligned} \|(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\tilde{\mathbf{\Lambda}}^{(m)})\mathbf{y}\| &= \|(\mathbf{A}\mathbf{Q}^{(m)}\mathbf{R}^{(m)}(\mathbf{R}^{(m)})^{-1} - \mathbf{Q}^{(m)}\mathbf{R}^{(m)}\mathbf{\Lambda}(\mathbf{R}^{(m)})^{-1})\mathbf{y}\| \\ &= \|(\mathbf{A}\mathbf{X}^{(m)} - \mathbf{X}^{(m)}\mathbf{\Lambda})(\mathbf{R}^{(m)})^{-1}\mathbf{y}\| \\ &= \|(\mathbf{A}\mathbf{X}^{(m)} - \mathbf{X}^{(m)}\mathbf{\Lambda})\tilde{\mathbf{y}}\|, \\ \|(\mathbf{A}\mathbf{Q}^{(0)} - \mathbf{Q}^{(0)}\tilde{\mathbf{\Lambda}}^{(0)})\mathbf{z}\| &= \|(\mathbf{A}\mathbf{X}^{(0)} - \mathbf{X}^{(0)}\mathbf{\Lambda})(\mathbf{R}^{(0)})^{-1}\mathbf{z}\| \\ &= \|(\mathbf{A}\mathbf{X}^{(0)} - \mathbf{X}^{(0)}\mathbf{\Lambda})\tilde{\mathbf{y}}\|, \\ \|(\mathbf{Q}^{(m)} - \mathbf{P}\mathbf{Q}^{(m)})\mathbf{y}\| &= \|(\mathbf{Q}^{(m)}\mathbf{R}^{(m)} - \mathbf{P}\mathbf{Q}^{(m)}\mathbf{R}^{(m)})(\mathbf{R}^{(m)})^{-1}\mathbf{y}\| \\ &= \|(\mathbf{X}^{(m)} - \mathbf{P}\mathbf{X}^{(m)})\tilde{\mathbf{y}}\|, \\ \|(\mathbf{Q}^{(0)} - \mathbf{P}\mathbf{Q}^{(0)})\mathbf{z}\| &= \|(\mathbf{X}^{(0)} - \mathbf{P}\mathbf{X}^{(0)})\tilde{\mathbf{y}}\|, \\ \|\mathbf{P}\mathbf{Q}^{(m)}\mathbf{y}\| &= \|\mathbf{P}\mathbf{Q}^{(m)}\mathbf{R}^{(m)}\tilde{\mathbf{y}}\| = \|\mathbf{P}\mathbf{X}^{(m)}\tilde{\mathbf{y}}\|, \\ \|\mathbf{P}\mathbf{Q}^{(0)}\mathbf{z}\| &= \|\mathbf{P}\mathbf{Q}^{(0)}\mathbf{R}^{(0)}\tilde{\mathbf{y}}\| = \|\mathbf{P}\mathbf{X}^{(0)}\tilde{\mathbf{y}}\|, \end{aligned}$$

mit denen die Aussagen des Satzes 5.45 die gewünschte Form annehmen.  $\blacksquare$

Auch die Konvergenz von Teilräumen können wir durch Winkel ausdrücken: Für zwei nicht leere Teilräume  $\mathcal{X}, \mathcal{Y} \subseteq \mathbb{K}^n$  definieren wir

$$\begin{aligned} \sin \angle(\mathcal{X}, \mathcal{Y}) &:= \max\{\sin \angle(\mathbf{x}, \mathcal{Y}) : \mathbf{x} \in \mathcal{X} \setminus \{\mathbf{0}\}\}, \\ \cos \angle(\mathcal{X}, \mathcal{Y}) &:= \min\{\cos \angle(\mathbf{x}, \mathcal{Y}) : \mathbf{x} \in \mathcal{X} \setminus \{\mathbf{0}\}\}, \\ \tan \angle(\mathcal{X}, \mathcal{Y}) &:= \max\{\tan \angle(\mathbf{x}, \mathcal{Y}) : \mathbf{x} \in \mathcal{X} \setminus \{\mathbf{0}\}\}. \end{aligned}$$

Dann erhalten wir mit (5.18) insbesondere für den Tangens

$$\begin{aligned} \tan \angle(\text{Bild } \mathbf{X}^{(m)}, \mathcal{E}) &= \max \left\{ \frac{\|\mathbf{x} - \mathbf{P}\mathbf{x}\|}{\|\mathbf{P}\mathbf{x}\|} : \mathbf{x} \in \text{Bild } \mathbf{X}^{(m)} \setminus \{\mathbf{0}\} \right\} \\ &= \max \left\{ \frac{\|(\mathbf{X}^{(m)} - \mathbf{P}\mathbf{X}^{(m)})\mathbf{y}\|}{\|\mathbf{P}\mathbf{X}^{(m)}\mathbf{y}\|} : \mathbf{y} \in \mathbb{K}^k \setminus \{\mathbf{0}\} \right\}. \end{aligned}$$

Indem wir auf beiden Seiten der Abschätzungen (5.21b) und (5.23b) zu dem Maximum übergehen, erhalten wir

$$\begin{aligned} \tan \angle(\text{Bild } \mathbf{X}^{(m)}, \mathcal{E}) &\leq \left( \frac{|\lambda_{k+1}|}{|\lambda_k|} \right)^m \tan \angle(\text{Bild } \mathbf{X}^{(0)}, \mathcal{E}), \\ \tan \angle(\text{Bild } \mathbf{Q}^{(m)}, \mathcal{E}) &\leq \left( \frac{|\lambda_{k+1}|}{|\lambda_k|} \right)^m \tan \angle(\text{Bild } \mathbf{Q}^{(0)}, \mathcal{E}) \quad \text{für alle } m \in \mathbb{N}_0. \end{aligned}$$

**Übungsaufgabe 5.47 (Winkel zwischen Teilräumen)** Seien  $n \in \mathbb{N}$  und  $k \in [1 : n]$  gegeben, und seien  $\mathcal{X}, \mathcal{Y} \subseteq \mathbb{K}^n$   $k$ -dimensionale Teilräume. Der Winkel zwischen den Teilräumen ist gegeben durch

$$\cos \angle(\mathcal{X}, \mathcal{Y}) = \min \left\{ \max \left\{ \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|}{\|\mathbf{x}\| \|\mathbf{y}\|} : \mathbf{y} \in \mathcal{Y} \setminus \{\mathbf{0}\} \right\} : \mathbf{x} \in \mathcal{X} \setminus \{\mathbf{0}\} \right\}.$$

(a) Sei eine isometrische Matrix  $\mathbf{Y} \in \mathbb{K}^{n \times k}$  mit  $\text{Bild } \mathbf{Y} = \mathcal{Y}$  gegeben. Beweisen Sie

$$\max \left\{ \frac{|\langle \mathbf{x}, \mathbf{y} \rangle|}{\|\mathbf{y}\|} : \mathbf{y} \in \mathcal{Y} \setminus \{\mathbf{0}\} \right\} = \|\mathbf{Y}^* \mathbf{x}\| \quad \text{für alle } \mathbf{x} \in \mathcal{X}.$$

(b) Seien isometrische Matrizen  $\mathbf{X}, \mathbf{Y} \in \mathbb{K}^{n \times k}$  mit  $\text{Bild } \mathbf{X} = \mathcal{X}$  und  $\text{Bild } \mathbf{Y} = \mathcal{Y}$  gegeben. Beweisen Sie, dass  $\cos \angle(\mathcal{X}, \mathcal{Y}) > 0$  genau dann gilt, wenn  $\mathbf{Y}^* \mathbf{X}$  invertierbar ist, und dass in diesem Fall die folgende Gleichung gilt:

$$\cos \angle(\mathcal{X}, \mathcal{Y}) = \frac{1}{\|(\mathbf{Y}^* \mathbf{X})^{-1}\|}.$$

(c) Beweisen Sie  $\cos \angle(\mathcal{X}, \mathcal{Y}) = \cos \angle(\mathcal{Y}, \mathcal{X})$ .

Hinweis: Die Cauchy-Schwarz-Ungleichung sowie die Lemmas 3.18 und 3.20 können hilfreich sein.

Die Fehlerabschätzungen (5.21a) und (5.23a) sind noch etwas unhandlich, da sie auf den Matrizen  $\tilde{\mathbf{\Lambda}}^{(m)}$  beruhen, die sich eventuell in der Praxis nur schwierig berechnen lassen. Wir können allerdings diese Matrizen durch eine geeignete Verallgemeinerung des Rayleigh-Quotienten ersetzen: Wir definieren

$$\mathbf{\Lambda}^{(m)} := (\mathbf{Q}^{(m)})^* \mathbf{A} \mathbf{Q}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0. \quad (5.24)$$

Wenn wir die Einheitsvektoren  $\mathbf{x}^{(m)} \in \mathbb{K}^n$  der ursprünglichen Vektoriteration als isometrische Matrizen  $\mathbf{Q}^{(m)} \in \mathbb{K}^{n \times 1}$  interpretieren, ist sofort ersichtlich, dass  $\lambda^{(m)} = \Lambda_A(\mathbf{x}^{(m)}) = \mathbf{\Lambda}^{(m)}$  gilt, dass also beide Definitionen zusammenfallen. Die Orthogonalitätsbeziehung (5.6) lässt sich wie folgt verallgemeinern:

**Lemma 5.48** Sei  $\mathbf{\Lambda} \in \mathbb{K}^{k \times k}$  eine beliebige Matrix. Es gilt

$$\|(\mathbf{A} \mathbf{Q}^{(m)} - \mathbf{Q}^{(m)} \mathbf{\Lambda}) \mathbf{y}\|^2 = \|(\mathbf{A} \mathbf{Q}^{(m)} - \mathbf{Q}^{(m)} \mathbf{\Lambda}^{(m)}) \mathbf{y}\|^2 + \|(\mathbf{\Lambda}^{(m)} - \mathbf{\Lambda}) \mathbf{y}\|^2$$

für alle  $m \in \mathbb{N}$  und  $\mathbf{y} \in \mathbb{K}^k$ .

Insbesondere minimiert die Wahl  $\mathbf{\Lambda} = \mathbf{\Lambda}^{(m)}$  die rechte Seite der Gleichung, und damit natürlich auch linke.

*Beweis.* Seien  $m \in \mathbb{N}$ ,  $\mathbf{y} \in \mathbb{K}^k$  und  $\mathbf{\Lambda} \in \mathbb{K}^{k \times k}$  gegeben. Dann folgt mit der binomischen Gleichung und den Lemmas 3.17 und 3.34

$$\|(\mathbf{A} \mathbf{Q}^{(m)} - \mathbf{Q}^{(m)} \mathbf{\Lambda}) \mathbf{y}\|^2 = \|(\mathbf{A} \mathbf{Q}^{(m)} - \mathbf{Q}^{(m)} \mathbf{\Lambda}^{(m)}) \mathbf{y} + \mathbf{Q}^{(m)} (\mathbf{\Lambda}^{(m)} - \mathbf{\Lambda}) \mathbf{y}\|^2$$

$$\begin{aligned}
 &= \|(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\boldsymbol{\Lambda}^{(m)})\mathbf{y}\|^2 + \|\mathbf{Q}^{(m)}(\boldsymbol{\Lambda}^{(m)} - \boldsymbol{\Lambda})\mathbf{y}\|^2 \\
 &\quad + 2\operatorname{Re}\langle(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\boldsymbol{\Lambda}^{(m)})\mathbf{y}, \mathbf{Q}^{(m)}(\boldsymbol{\Lambda}^{(m)} - \boldsymbol{\Lambda})\mathbf{y}\rangle_2^2 \\
 &= \|(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\boldsymbol{\Lambda}^{(m)})\mathbf{y}\|^2 + \|(\boldsymbol{\Lambda}^{(m)} - \boldsymbol{\Lambda})\mathbf{y}\|^2 \\
 &\quad + 2\operatorname{Re}\langle(\mathbf{Q}^{(m)})^*(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\boldsymbol{\Lambda}^{(m)})\mathbf{y}, (\boldsymbol{\Lambda}^{(m)} - \boldsymbol{\Lambda})\mathbf{y}\rangle_2^2 \\
 &= \|(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\boldsymbol{\Lambda}^{(m)})\mathbf{y}\|^2 + \|(\boldsymbol{\Lambda}^{(m)} - \boldsymbol{\Lambda})\mathbf{y}\|^2 \\
 &\quad + 2\operatorname{Re}\langle(\boldsymbol{\Lambda}^{(m)} - \boldsymbol{\Lambda}^{(m)})\mathbf{y}, (\boldsymbol{\Lambda}^{(m)} - \boldsymbol{\Lambda})\mathbf{y}\rangle_2^2 \\
 &= \|(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\boldsymbol{\Lambda}^{(m)})\mathbf{y}\|^2 + \|(\boldsymbol{\Lambda}^{(m)} - \boldsymbol{\Lambda})\mathbf{y}\|^2,
 \end{aligned}$$

und das ist die gewünschte Darstellung des Fehlers.  $\blacksquare$

Wir können also auf der linken Seite der Folgerung 5.46 einfach die unhandliche Matrix  $\tilde{\boldsymbol{\Lambda}}^{(m)}$  durch den verallgemeinerten Rayleigh-Quotienten  $\boldsymbol{\Lambda}^{(m)}$  ersetzen, den wir praktisch berechnen können.

Als Vorbereitung auf das folgende Kapitel halten wir fest, dass die Folgerung 5.46 auch eine Aussage über die Konvergenz der Spektralnorm enthält.

**Folgerung 5.49 (Spektralnorm)** *Unter den Voraussetzungen von Satz 5.45 und mit (5.22) gilt mit der dort eingeführten Matrix  $\tilde{\boldsymbol{\Lambda}}^{(0)}$  die Abschätzung*

$$\|\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\boldsymbol{\Lambda}^{(m)}\| \leq \left(\frac{|\lambda_{k+1}|}{|\lambda_k|}\right)^m \max \left\{ \frac{\|(\mathbf{A}\mathbf{Q}^{(0)} - \mathbf{Q}^{(0)}\tilde{\boldsymbol{\Lambda}}^{(0)})\mathbf{z}\|}{\|\mathbf{P}\mathbf{Q}^{(0)}\mathbf{z}\|} : \mathbf{z} \in \mathbb{K}^k \setminus \{\mathbf{0}\} \right\}$$

für alle  $m \in \mathbb{N}_0$ .

*Beweis.* Da  $\mathbf{P}$  eine orthogonale Projektion ist, gilt mit (5.16c) und Lemma 3.34 die Abschätzung

$$\|\mathbf{P}\mathbf{Q}^{(m)}\mathbf{y}\| \leq \|\mathbf{Q}^{(m)}\mathbf{y}\| = \|\mathbf{y}\| \quad \text{für alle } \mathbf{y} \in \mathbb{K}^k.$$

Aus Lemma 5.48 erhalten wir

$$\|(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\boldsymbol{\Lambda}^{(m)})\mathbf{y}\| \leq \|(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\tilde{\boldsymbol{\Lambda}}^{(m)})\mathbf{y}\| \quad \text{für alle } \mathbf{y} \in \mathbb{K}^k.$$

Die Kombination beider Abschätzungen mit der Folgerung 5.46 führt zu

$$\begin{aligned}
 \frac{\|(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\boldsymbol{\Lambda}^{(m)})\mathbf{y}\|}{\|\mathbf{y}\|} &\leq \frac{\|(\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\tilde{\boldsymbol{\Lambda}}^{(m)})\mathbf{y}\|}{\|\mathbf{P}\mathbf{Q}^{(m)}\mathbf{y}\|} \\
 &\leq \left(\frac{|\lambda_{k+1}|}{|\lambda_k|}\right)^m \frac{\|(\mathbf{A}\mathbf{Q}^{(0)} - \mathbf{Q}^{(0)}\tilde{\boldsymbol{\Lambda}}^{(0)})\mathbf{z}\|}{\|\mathbf{P}\mathbf{Q}^{(0)}\mathbf{z}\|},
 \end{aligned}$$

mit einem Vektor  $\mathbf{z} \in \mathbb{K}^k$ . Indem wir auf beiden Seiten dieser Abschätzung zu dem Maximum übergehen, folgt das gewünschte Ergebnis.  $\blacksquare$

**Übungsaufgabe 5.50 (Konvergenz des Rayleigh-Quotienten)** Seien  $n \in \mathbb{N}$ ,  $k \in [1 : n]$  gegeben. Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und  $\tilde{\mathbf{\Lambda}} \in \mathbb{K}^{k \times k}$  Matrizen. Sei  $\mathbf{Q} \in \mathbb{K}^{n \times k}$  isometrisch, und sei  $\mathbf{\Lambda} := \mathbf{Q}^* \mathbf{A} \mathbf{Q}$ .

(a) Beweisen Sie

$$\|(\mathbf{\Lambda} - \tilde{\mathbf{\Lambda}})\mathbf{y}\| = \|\mathbf{Q}\mathbf{Q}^*(\mathbf{A}\mathbf{Q} - \mathbf{Q}\tilde{\mathbf{\Lambda}})\mathbf{y}\| \quad \text{für alle } \mathbf{y} \in \mathbb{K}^k.$$

(b) Folgern Sie daraus

$$\|(\mathbf{A}\mathbf{Q} - \mathbf{Q}\mathbf{\Lambda})\mathbf{y}\| = \|(\mathbf{I} - \mathbf{Q}\mathbf{Q}^*)(\mathbf{A}\mathbf{Q} - \mathbf{Q}\tilde{\mathbf{\Lambda}})\mathbf{y}\| \quad \text{für alle } \mathbf{y} \in \mathbb{K}^k.$$

Für die Praxis ist es nicht sinnvoll, die Matrizen  $\mathbf{X}^{(m)}$  zu berechnen, da sie aus den bereits erwähnten Gründen zunehmend schlecht konditioniert sein werden. Stattdessen wollen wir möglichst mit den isometrischen Matrizen  $\mathbf{Q}^{(m)}$  arbeiten, bei denen diese Gefahr nicht besteht.

Wir verwenden (5.22), um

$$\mathbf{X}^{(m+1)} = \mathbf{A}\mathbf{X}^{(m)} = \mathbf{A}\mathbf{Q}^{(m)}\mathbf{R}^{(m)} \quad (5.25)$$

zu erhalten. Nun berechnen wir die, da  $\mathbf{Q}^{(m)}$  isometrisch ist, hoffentlich gut konditionierte Matrix

$$\mathbf{W}^{(m+1)} := \mathbf{A}\mathbf{Q}^{(m)} \in \mathbb{K}^{n \times k}$$

und bestimmen ihre (dünne) QR-Zerlegung

$$\mathbf{W}^{(m+1)} = \mathbf{Q}^{(m+1)}\tilde{\mathbf{R}}^{(m+1)}$$

mit einer isometrischen Matrix  $\mathbf{Q}^{(m+1)} \in \mathbb{K}^{n \times k}$  sowie einer oberen Dreiecksmatrix  $\tilde{\mathbf{R}}^{(m+1)} \in \mathbb{K}^{k \times k}$ . Wir setzen diese Zerlegung in die Gleichung (5.25) ein, um

$$\mathbf{X}^{(m+1)} = \mathbf{W}^{(m+1)}\mathbf{R}^{(m)} = \mathbf{Q}^{(m+1)}\tilde{\mathbf{R}}^{(m+1)}\mathbf{R}^{(m)}$$

zu erhalten, und da  $\tilde{\mathbf{R}}^{(m+1)}$  und  $\mathbf{R}^{(m)}$  obere Dreiecksmatrizen sind, muss auch

$$\mathbf{R}^{(m+1)} := \tilde{\mathbf{R}}^{(m+1)}\mathbf{R}^{(m)}$$

eine obere Dreiecksmatrix sein. Wir können also eine Zerlegung der gewünschten Form (5.22) berechnen, indem wir die Matrix  $\mathbf{A}$  mit der *isometrischen* — und deshalb gut konditionierten — Matrix  $\mathbf{Q}^{(m)}$  multiplizieren und eine QR-Zerlegung des Ergebnisses berechnen.

Die resultierende Verallgemeinerung der Vektoriteration bezeichnet man als *orthogonale Iteration*:

**Algorithmus 5.51 (Orthogonale Iteration)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und eine isometrische Matrix  $\mathbf{Q}^{(0)} \in \mathbb{K}^{n \times k}$  gegeben. Der folgende Algorithmus führt die orthogonale Iteration aus und berechnet die verallgemeinerten Rayleigh-Quotienten  $(\mathbf{\Lambda}^{(m)})_{m=0}^{\infty}$ .



```

m ← 0
W(m+1) ← AQ(m)
Λ(m) ← (Q(m))*W(m+1)
while „Fehler zu groß“ do begin
  Berechne eine dünne QR-Zerlegung Q(m+1)R(m+1) = W(m+1)
  m ← m + 1
  W(m+1) ← AQ(m)
  Λ(m) ← (Q(m))*W(m+1)
end

```

Selbstverständlich müssen wir uns auch bei der orthogonalen Iteration Gedanken über ein geeignetes Abbruchkriterium machen. In Anlehnung an das Vorgehen bei der Vektoriteration bietet es sich an, dass verallgemeinerte Residuum

$$\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\mathbf{\Lambda}^{(m)}$$

zu verwenden. Beispielsweise wissen wir dank Folgerung 5.49, dass die Spektralnorm dieser Matrix gegen null konvergieren dürfte, so dass die Bedingung

$$\|\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\mathbf{\Lambda}^{(m)}\| \leq \epsilon$$

geeignet erscheint. Um dafür zu sorgen, dass das Kriterium von der Skalierung der Matrix  $\mathbf{A}$  unabhängig ist, könnte man auch

$$\|\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\mathbf{\Lambda}^{(m)}\| \leq \epsilon \|\mathbf{\Lambda}^{(m)}\|$$

in Erwägung ziehen.

Wir können auf die in Abschnitt 5.2 entwickelten Fehlerabschätzungen zurückgreifen, indem wir uns an die Bemerkung 3.30 erinnern: Wenn  $\tilde{\lambda}$  ein Eigenwert der Matrix  $\mathbf{\Lambda}^{(m)}$  mit einem Eigenvektor  $\hat{\mathbf{x}} \in \mathbb{K}^k \setminus \{\mathbf{0}\}$  ist, erfüllt  $\mathbf{x} := \mathbf{Q}^{(m)}\hat{\mathbf{x}}$  wegen Lemma 3.34 sowohl  $\|\mathbf{x}\| = \|\hat{\mathbf{x}}\| > 0$  als auch

$$\begin{aligned} \|\mathbf{A}\mathbf{x} - \tilde{\lambda}\mathbf{x}\| &= \|\mathbf{A}\mathbf{Q}^{(m)}\hat{\mathbf{x}} - \tilde{\lambda}\mathbf{Q}^{(m)}\hat{\mathbf{x}}\| = \|\mathbf{A}\mathbf{Q}^{(m)}\hat{\mathbf{x}} - \mathbf{Q}^{(m)}\mathbf{\Lambda}^{(m)}\hat{\mathbf{x}}\| \\ &\leq \|\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\mathbf{\Lambda}^{(m)}\| \|\hat{\mathbf{x}}\| = \|\mathbf{A}\mathbf{Q}^{(m)} - \mathbf{Q}^{(m)}\mathbf{\Lambda}^{(m)}\| \|\mathbf{x}\|. \end{aligned}$$

Falls also die Spektralnorm des verallgemeinerten Residuums klein ist, finden wir zu jedem Eigenwert  $\tilde{\lambda}$  der Matrix  $\mathbf{\Lambda}^{(m)}$  einen Vektor  $\mathbf{x}$ , der einen Eigenvektor der Matrix  $\mathbf{A}$  zu demselben Eigenwert approximiert. Da wegen Lemma 3.17 auch

$$\begin{aligned} \tilde{\lambda} = \Lambda_{\mathbf{\Lambda}^{(m)}}(\hat{\mathbf{x}}) &= \frac{\langle \hat{\mathbf{x}}, \mathbf{\Lambda}^{(m)}\hat{\mathbf{x}} \rangle}{\langle \hat{\mathbf{x}}, \hat{\mathbf{x}} \rangle} = \frac{\langle \hat{\mathbf{x}}, (\mathbf{Q}^{(m)})^* \mathbf{A}\mathbf{Q}^{(m)}\hat{\mathbf{x}} \rangle}{\langle \hat{\mathbf{x}}, (\mathbf{Q}^{(m)})^* \mathbf{Q}^{(m)}\hat{\mathbf{x}} \rangle} \\ &= \frac{\langle \mathbf{Q}^{(m)}\hat{\mathbf{x}}, \mathbf{A}\mathbf{Q}^{(m)}\hat{\mathbf{x}} \rangle}{\langle \mathbf{Q}^{(m)}\hat{\mathbf{x}}, \mathbf{Q}^{(m)}\hat{\mathbf{x}} \rangle} = \frac{\langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{x} \rangle} = \Lambda_A(\mathbf{x}) \end{aligned}$$

gilt, können wir die Theorie des Abschnitts 5.2 auf die Eigenwertapproximationen  $\tilde{\lambda}$  und die zugehörigen Eigenvektorapproximationen anwenden.

**Bemerkung 5.52 (Orthogonale inverse Iteration)** *Selbstverständlich können wir auch die inverse Iteration mit und ohne Shift in ähnlicher Weise modifizieren, um Näherungen invarianter Unterräume zu berechnen.*

*Dann müssen für die Berechnung der Matrizen  $\mathbf{W}^{(m+1)}$  in Algorithmus 5.51 lineare Gleichungssysteme mit mehreren rechten Seiten gelöst werden.*

**Bemerkung 5.53 (Orthogonale Rayleigh-Iteration)** *Auf der Diagonalen der Matrix  $\mathbf{\Lambda}^{(m)}$  finden sich die Rayleigh-Quotienten zu den Spaltenvektoren der Matrix  $\mathbf{Q}^{(m)}$ .*

*Wir können eine Variante der Rayleigh-Iteration konstruieren, indem wir die Spalten der Matrix  $\mathbf{W}^{(m+1)}$  berechnen, indem wir geschiftete lineare Gleichungssysteme lösen: Für die  $\ell$ -te Spalte nehmen wir den  $\ell$ -ten Diagonaleintrag der Matrix  $\mathbf{\Lambda}^{(m)}$  als Shift-Parameter und verwenden die  $\ell$ -te Spalte der Matrix  $\mathbf{Q}^{(m)}$  als rechte Seite.*

*Unter geeigneten Voraussetzungen können wir dann die schnelle (quadratische) Konvergenz der Rayleigh-Iteration auch für die Berechnung mehrerer Eigenvektoren erhalten.*

**Übungsaufgabe 5.54 (Dünne QR-Zerlegung)** *Seien  $n \in \mathbb{N}$  und  $k \in [1 : n]$  gegeben, sei  $\mathbf{A} \in \mathbb{K}^{n \times k}$  eine Matrix mit vollem Rang.*

*Seien nun dünne QR-Zerlegungen*

$$\mathbf{A} = \mathbf{Q}_1 \mathbf{R}_1, \quad \mathbf{A} = \mathbf{Q}_2 \mathbf{R}_2$$

*mit isometrischen Matrizen  $\mathbf{Q}_1, \mathbf{Q}_2 \in \mathbb{K}^{n \times k}$  und rechten oberen Dreiecksmatrizen  $\mathbf{R}_1, \mathbf{R}_2 \in \mathbb{K}^{k \times k}$  gegeben.*

*Beweisen Sie, dass eine unitäre Diagonalmatrix  $\mathbf{D} \in \mathbb{K}^{k \times k}$  existiert mit*

$$\mathbf{Q}_2 = \mathbf{Q}_1 \mathbf{D}, \quad \mathbf{R}_2 = \mathbf{D}^* \mathbf{R}_1.$$

*Hinweise: Ohne Beweis kann verwendet werden, dass die Inverse einer Dreiecksmatrix und auch das Produkt zweier Dreiecksmatrizen jeweils wieder Dreiecksmatrizen sind.*

*Ein Blick auf die Matrix  $\mathbf{Q}_1 \mathbf{Q}_1^* \mathbf{Q}_2$  könnte sich lohnen.*

*Lemma 3.53 kann in einem Beweisschritt hilfreich sein, wenn man bedenkt, dass unitäre Matrizen auch normal sind.*

**Übungsaufgabe 5.55 (Norm des Rayleigh-Quotienten)** *Seien  $n \in \mathbb{N}$  und  $\mathbf{A} \in \mathbb{K}^{n \times n}$  gegeben.*

- (a) *Seien  $k \in [1 : n]$ , eine Matrix  $\mathbf{\Lambda} \in \mathbb{K}^{k \times k}$  und eine injektive Matrix  $\mathbf{X} \in \mathbb{K}^{n \times k}$  gegeben mit*

$$\mathbf{A} \mathbf{X} = \mathbf{X} \mathbf{\Lambda}.$$

*Beweisen oder widerlegen Sie, dass dann  $\|\mathbf{\Lambda}\| \leq \|\mathbf{A}\|$  gelten muss.*

- (b) *Seien  $k \in [1 : n]$ , eine Matrix  $\mathbf{\Lambda} \in \mathbb{K}^{k \times k}$  und eine isometrische Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times k}$  gegeben mit*

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

*Beweisen Sie, dass dann  $\|\mathbf{\Lambda}\| \leq \|\mathbf{A}\|$  gilt.*

Hinweis: Die Lemmas 3.20 und 3.55 können bei der Berechnung der Spektralnorm nützlich sein.

**Übungsaufgabe 5.56 (Konvergenz des Residuums)** Seien  $n \in \mathbb{N}$ ,  $k \in [1 : n]$  gegeben, sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine Matrix und  $\mathcal{E} \subseteq \mathbb{K}^n$  ein  $k$ -dimensionaler invarianter Teilraum.

Sei  $\mathbf{E} \in \mathbb{K}^{n \times k}$  eine isometrische Matrix mit  $\text{Bild } \mathbf{E} = \mathcal{E}$ .

Sei  $\mathbf{Q} \in \mathbb{K}^{n \times k}$  eine beliebige isometrische Matrix.

(a) Beweisen Sie

$$\|\mathbf{Q} - \mathbf{E}\mathbf{E}^*\mathbf{Q}\| = \sin \angle(\text{Bild } \mathbf{Q}, \mathcal{E}), \quad \|\mathbf{E} - \mathbf{Q}\mathbf{Q}^*\mathbf{E}\| = \sin \angle(\text{Bild } \mathbf{Q}, \mathcal{E}).$$

(b) Folgern Sie daraus, dass eine Matrix  $\mathbf{\Lambda} \in \mathbb{K}^{k \times k}$  existiert mit

$$\|\mathbf{A}\mathbf{Q} - \mathbf{Q}\mathbf{\Lambda}\| \leq 2\|\mathbf{A}\| \sin \angle(\text{Bild } \mathbf{Q}, \mathcal{E}).$$

Hinweise: Übungsaufgabe 5.47 kann verwendet werden. Die Lemmas 3.20, 3.27 und 5.39 sind einen Blick wert.



## 6 Die QR-Iteration

Wir wenden uns zunächst wieder der Frage zu, wie sich *alle* Eigenvektoren einer gegebenen Matrix bestimmen lassen. Für selbstadjungierte Matrizen haben wir mit der Jacobi-Iteration bereits ein erstes Verfahren zur Lösung dieser Aufgabe kennengelernt, allerdings zeigt sich, dass diese Iteration den Nachteil hat, dass bei bereits „fast“ diagonalen Matrizen, etwa Tridiagonalmatrizen, diese Struktur wieder zunichte gemacht wird. In diesem Kapitel entwickeln wir einen Algorithmus, der diesen Nachteil nicht aufweist.

Darüber hinaus können wir bei diesem Ansatz Shift-Strategien besonders einfach umsetzen, die Konvergenz überwachen, und die Dimension reduzieren, falls einzelne Eigenräume bereits konvergiert sind.

### 6.1 Grundidee

Am Beispiel der orthogonalen Iteration haben wir gesehen, dass sich mehrere Eigenvektoren simultan berechnen lassen. Also liegt es nahe, die Iteration so zu erweitern, dass *alle* Eigenvektoren berechnet werden, indem man sie mit einer vollständigen Basis startet, also beispielsweise mit  $\mathbf{Q}^{(0)} = \mathbf{I}$ .

Wir erhalten damit das folgende Verfahren:

```

 $\mathbf{Q}^{(0)} \leftarrow \mathbf{I}$ 
 $m \leftarrow 0$ 
while „Fehler zu groß“ do begin
   $\mathbf{W}^{(m+1)} \leftarrow \mathbf{A}\mathbf{Q}^{(m)}$ 
  Berechne die QR-Zerlegung  $\mathbf{Q}^{(m+1)}\mathbf{R}^{(m+1)} = \mathbf{W}^{(m+1)}$ 
   $m \leftarrow m + 1$ 
end

```

Bisher haben wir bei der Untersuchung der orthogonalen Iteration nicht ausgenutzt, dass die Matrizen  $\mathbf{R}^{(m)}$  obere Dreiecksmatrizen sind. Diese Eigenschaft hat nützliche Konsequenzen, denen wir uns jetzt widmen werden.

Wir wählen ein  $k \in [1 : n - 1]$  und zerlegen die im Rahmen der orthogonalen Iteration auftretenden Matrizen in der Form

$$\mathbf{Q}^{(m)} = \begin{pmatrix} \mathbf{Q}_k^{(m)} & \mathbf{Q}_*^{(m)} \end{pmatrix}, \quad \mathbf{Q}_k^{(m)} \in \mathbb{K}^{n \times k},$$

$$\mathbf{R}^{(m)} = \begin{pmatrix} \mathbf{R}_{kk}^{(m)} & \mathbf{R}_{k*}^{(m)} \\ & \mathbf{R}_{**}^{(m)} \end{pmatrix}, \quad \mathbf{R}_{kk}^{(m)} \in \mathbb{K}^{k \times k} \quad \text{für alle } m \in \mathbb{N}_0.$$

## 6 Die QR-Iteration

Aus den definierenden Gleichungen der Iteration folgt dann

$$\begin{aligned}
 & \left( \mathbf{Q}_k^{(m+1)} \mathbf{R}_{kk}^{(m+1)} \quad \mathbf{Q}_k^{(m+1)} \mathbf{R}_{k*}^{(m+1)} + \mathbf{Q}_*^{(m+1)} \mathbf{R}_{**}^{(m+1)} \right) \\
 &= \left( \mathbf{Q}_k^{(m+1)} \quad \mathbf{Q}_*^{(m+1)} \right) \begin{pmatrix} \mathbf{R}_{kk}^{(m+1)} & \mathbf{R}_{k*}^{(m+1)} \\ & \mathbf{R}_{**}^{(m+1)} \end{pmatrix} = \mathbf{Q}^{(m+1)} \mathbf{R}^{(m+1)} \\
 &= \mathbf{W}^{(m+1)} = \mathbf{A} \mathbf{Q}^{(m)} = \mathbf{A} \begin{pmatrix} \mathbf{Q}_k^{(m)} & \mathbf{Q}_*^{(m)} \end{pmatrix} \\
 &= \begin{pmatrix} \mathbf{A} \mathbf{Q}_k^{(m)} & \mathbf{A} \mathbf{Q}_*^{(m)} \end{pmatrix} \quad \text{für alle } m \in \mathbb{N}_0,
 \end{aligned}$$

also insbesondere

$$\mathbf{Q}_k^{(m+1)} \mathbf{R}_{kk}^{(m+1)} = \mathbf{A} \mathbf{Q}_k^{(m)} \quad \text{für alle } m \in \mathbb{N}_0.$$

Die ersten  $k$  Spalten  $\mathbf{Q}_k^{(m)}$  der Matrizen  $\mathbf{Q}^{(m)}$  können demnach auch als Ergebnis einer orthogonalen Iteration mit der  $k$ -spaltigen Anfangsmatrix  $\mathbf{Q}_k^{(0)}$  interpretiert werden. Also lassen sich unsere Konvergenzresultate auch auf diese Teilmatrizen anwenden.

Falls insbesondere

$$|\lambda_1| \geq |\lambda_2| \geq \dots \geq |\lambda_k| > |\lambda_{k+1}| \geq \dots \geq |\lambda_n|$$

gilt, können wir Lemma 5.49 verwenden, um eine Abschätzung der Form

$$\|\mathbf{A} \mathbf{Q}_k^{(m)} - \mathbf{Q}_k^{(m)} \mathbf{\Lambda}_{kk}^{(m)}\|_2 \leq C \left( \frac{|\lambda_{k+1}|}{|\lambda_k|} \right)^m \quad \text{für alle } m \in \mathbb{N} \quad (6.1)$$

zu erhalten, indem wir entsprechend Lemma 5.48 die Matrix

$$\mathbf{\Lambda}_{kk}^{(m)} := (\mathbf{Q}_k^{(m)})^* \mathbf{A} \mathbf{Q}_k^{(m)}$$

einsetzen. Um diese Abschätzung auszunutzen, untersuchen wir, wie sich die Matrix  $\mathbf{A}$  in der durch  $\mathbf{Q}^{(m)}$  gegebenen Basis darstellt, wir sind also an den Matrizen

$$\mathbf{A}^{(m)} := (\mathbf{Q}^{(m)})^* \mathbf{A} \mathbf{Q}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0 \quad (6.2)$$

interessiert. Aus der Blockdarstellung der Matrizen  $\mathbf{Q}^{(m)}$  folgt

$$\begin{aligned}
 \mathbf{A}^{(m)} &= (\mathbf{Q}^{(m)})^* \mathbf{A} \mathbf{Q}^{(m)} = \begin{pmatrix} (\mathbf{Q}_k^{(m)})^* \\ (\mathbf{Q}_*^{(m)})^* \end{pmatrix} \mathbf{A} \begin{pmatrix} \mathbf{Q}_k^{(m)} & \mathbf{Q}_*^{(m)} \end{pmatrix} \\
 &= \begin{pmatrix} (\mathbf{Q}_k^{(m)})^* \mathbf{A} \mathbf{Q}_k^{(m)} & (\mathbf{Q}_k^{(m)})^* \mathbf{A} \mathbf{Q}_*^{(m)} \\ (\mathbf{Q}_*^{(m)})^* \mathbf{A} \mathbf{Q}_k^{(m)} & (\mathbf{Q}_*^{(m)})^* \mathbf{A} \mathbf{Q}_*^{(m)} \end{pmatrix}.
 \end{aligned}$$

Wir interessieren uns für den linken unteren Block dieser Matrix und möchten beweisen, dass er gegen null konvergiert. Aus (6.1) folgt  $\mathbf{A} \mathbf{Q}_k^{(m)} \approx \mathbf{Q}_k^{(m)} \mathbf{\Lambda}_{kk}^{(m)}$ , also

$$(\mathbf{Q}_*^{(m)})^* \mathbf{A} \mathbf{Q}_k^{(m)} \approx (\mathbf{Q}_*^{(m)})^* \mathbf{Q}_k^{(m)} \mathbf{\Lambda}_{kk}^{(m)}.$$

Da die Spalten der Matrix  $\mathbf{Q}^{(m)}$  eine Orthonormalbasis sind, müssen sie senkrecht aufeinander stehen. Konkret können wir nachrechnen, dass

$$\begin{aligned} \begin{pmatrix} (\mathbf{Q}_k^{(m)})^* \mathbf{Q}_k^{(m)} & (\mathbf{Q}_k^{(m)})^* \mathbf{Q}_*^{(m)} \\ (\mathbf{Q}_*^{(m)})^* \mathbf{Q}_k^{(m)} & (\mathbf{Q}_*^{(m)})^* \mathbf{Q}_*^{(m)} \end{pmatrix} &= \begin{pmatrix} (\mathbf{Q}_k^{(m)})^* \\ (\mathbf{Q}_*^{(m)})^* \end{pmatrix} \begin{pmatrix} \mathbf{Q}_k^{(m)} & \mathbf{Q}_*^{(m)} \end{pmatrix} \\ &= (\mathbf{Q}^{(m)})^* \mathbf{Q}^{(m)} = \mathbf{I} = \begin{pmatrix} \mathbf{I} & \\ & \mathbf{I} \end{pmatrix} \end{aligned}$$

gilt, also insbesondere  $(\mathbf{Q}_*^{(m)})^* \mathbf{Q}_k^{(m)} = \mathbf{0}$ . Wegen  $(\mathbf{Q}_k^{(m)})^* \mathbf{Q}_k^{(m)} = \mathbf{I}$  und  $(\mathbf{Q}_*^{(m)})^* \mathbf{Q}_*^{(m)} = \mathbf{I}$  sind die Matrizen  $\mathbf{Q}_k^{(m)}$  und  $\mathbf{Q}_*^{(m)}$  außerdem isometrisch. Für den rechten unteren Block der Matrix  $\mathbf{A}^{(m)}$  erhalten wir mit (6.1) und Lemma 3.20 die Abschätzung

$$\begin{aligned} \|(\mathbf{Q}_*^{(m)})^* \mathbf{A} \mathbf{Q}_k^{(m)}\| &= \|(\mathbf{Q}_*^{(m)})^* \mathbf{Q}_k^{(m)} \Lambda_{kk}^{(m)} + (\mathbf{Q}_*^{(m)})^* (\mathbf{A} \mathbf{Q}_k^{(m)} - \mathbf{Q}_k^{(m)} \Lambda_{kk}^{(m)})\| \\ &= \|(\mathbf{Q}_*^{(m)})^* (\mathbf{A} \mathbf{Q}_k^{(m)} - \mathbf{Q}_k^{(m)} \Lambda_{kk}^{(m)})\| \\ &\leq \|(\mathbf{Q}_*^{(m)})^*\| \|\mathbf{A} \mathbf{Q}_k^{(m)} - \mathbf{Q}_k^{(m)} \Lambda_{kk}^{(m)}\| \\ &= \|\mathbf{Q}_*^{(m)}\| \|\mathbf{A} \mathbf{Q}_k^{(m)} - \mathbf{Q}_k^{(m)} \Lambda_{kk}^{(m)}\| \\ &\leq C \left( \frac{|\lambda_{k+1}|}{|\lambda_k|} \right)^m \quad \text{für alle } m \in \mathbb{N}_0. \end{aligned}$$

Die Matrizen  $\mathbf{A}^{(m)}$  werden sich also einer oberen Block-Dreiecksform annähern.

Falls wir dieses Argument auf alle  $k \in [1 : n - 1]$  anwenden können, falls also

$$|\lambda_1| > |\lambda_2| > \dots > |\lambda_{n-1}| > |\lambda_n|$$

gilt, erhalten wir sogar, dass die Matrix gegen obere Dreiecksgestalt konvergiert, dass wir uns also iterativ der Schur-Zerlegung (vgl. Satz 3.41) nähern.

In diesem Fall bietet uns die orthogonale Iteration, angewendet auf eine vollständige Basis, eine Möglichkeit, die Schur-Zerlegung zu approximieren. Da  $\mathbf{Q}^{(m)}$  in diesem Fall immer eine vollständige Basis ist, ist  $\mathbf{A}^{(m)}$  das Ergebnis einer unitären Ähnlichkeitstransformation der Ausgangsmatrix  $\mathbf{A}$ , wir berechnen also eine Folge unitär ähnlicher Matrizen, die gegen eine obere Dreiecksmatrix konvergieren.

Wir würden diese Matrizen gerne berechnen, ohne die vollen orthogonalen Basen  $\mathbf{Q}^{(m)}$  mitführen zu müssen, denn es gibt Anwendungen, bei denen wir nur an den Eigenwerten, aber nicht an den Eigenvektoren interessiert sind. Nach Konstruktion haben wir

$$\begin{aligned} \mathbf{A} \mathbf{Q}^{(m)} &= \mathbf{W}^{(m+1)} = \mathbf{Q}^{(m+1)} \mathbf{R}^{(m+1)}, \\ \mathbf{A}^{(m)} &= (\mathbf{Q}^{(m)})^* \mathbf{A} \mathbf{Q}^{(m)} = (\mathbf{Q}^{(m)})^* \mathbf{Q}^{(m+1)} \mathbf{R}^{(m+1)} = \widehat{\mathbf{Q}}^{(m+1)} \mathbf{R}^{(m+1)}, \end{aligned} \quad (6.3)$$

mit der Matrix

$$\widehat{\mathbf{Q}}^{(m+1)} := (\mathbf{Q}^{(m)})^* \mathbf{Q}^{(m+1)},$$

die nach Lemma 3.35 äquivalent mit

$$\mathbf{Q}^{(m)} \widehat{\mathbf{Q}}^{(m+1)} = \mathbf{Q}^{(m+1)} \quad (6.4)$$

## 6 Die QR-Iteration

ist, also den Wechsel von  $\mathbf{Q}^{(m)}$  zu  $\mathbf{Q}^{(m+1)}$  beschreibt.

Aus den Gleichungen (6.4) und (6.3) folgt

$$\begin{aligned}\mathbf{A}^{(m+1)} &= (\mathbf{Q}^{(m+1)})^* \mathbf{A} \mathbf{Q}^{(m+1)} = (\widehat{\mathbf{Q}}^{(m+1)})^* (\mathbf{Q}^{(m)})^* \mathbf{A} \mathbf{Q}^{(m)} \widehat{\mathbf{Q}}^{(m+1)} \\ &= (\widehat{\mathbf{Q}}^{(m+1)})^* \mathbf{A}^{(m)} \widehat{\mathbf{Q}}^{(m+1)} = (\widehat{\mathbf{Q}}^{(m+1)})^* \widehat{\mathbf{Q}}^{(m+1)} \mathbf{R}^{(m+1)} \widehat{\mathbf{Q}}^{(m+1)} \\ &= \mathbf{R}^{(m+1)} \widehat{\mathbf{Q}}^{(m+1)},\end{aligned}$$

so dass wir  $\mathbf{A}^{(m+1)}$  direkt bestimmen können, ohne auf  $\mathbf{A}$  oder  $\mathbf{Q}^{(m+1)}$  zurückgreifen zu müssen. Zusammen mit (6.3) erhalten wir die Gleichungen

$$\mathbf{A}^{(m)} = \widehat{\mathbf{Q}}^{(m+1)} \mathbf{R}^{(m+1)}, \quad \mathbf{A}^{(m+1)} = \mathbf{R}^{(m+1)} \widehat{\mathbf{Q}}^{(m+1)}, \quad (6.5)$$

mit denen die vollständige Iteration bereits beschrieben ist. Wir können also die orthogonale Iteration durchführen, ohne die Matrizen  $\mathbf{Q}^{(m)}$  explizit aufzustellen und ohne die ursprüngliche Matrix  $\mathbf{A}$  zu verwenden.

Dieser Zugang ähnelt wegen

$$\mathbf{A}^{(m+1)} = \mathbf{R}^{(m+1)} \widehat{\mathbf{Q}}^{(m+1)} = (\widehat{\mathbf{Q}}^{(m+1)})^* \mathbf{A}^{(m)} \widehat{\mathbf{Q}}^{(m+1)}$$

dem bereits bei der Jacobi-Iteration verwendeten Ansatz: Alle Iterierten sind unitär ähnlich, und wir versuchen, die unitären Transformationen so zu wählen, dass sie gegen die obere Dreiecksgestalt konvergieren.

**Algorithmus 6.1 (QR-Iteration)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Der folgende Algorithmus berechnet die Folge  $(\mathbf{A}^{(m)})_{m \in \mathbb{N}}$  aus (6.2).

```

 $\mathbf{A}^{(0)} \leftarrow \mathbf{A}$ 
 $m \leftarrow 0$ 
while „Fehler zu groß“ do begin
  Bestimme die QR-Zerlegung  $\widehat{\mathbf{Q}}^{(m+1)} \mathbf{R}^{(m+1)} = \mathbf{A}^{(m)}$ 
   $\mathbf{A}^{(m+1)} \leftarrow \mathbf{R}^{(m+1)} \widehat{\mathbf{Q}}^{(m+1)}$ 
   $m \leftarrow m + 1$ 
end

```

**Bemerkung 6.2** Nach der durch Gleichung (6.2) gegebenen Definition sind alle Matrizen  $\mathbf{A}^{(m)}$  unitär ähnlich zu  $\mathbf{A}$ , besitzen also dieselben Eigenwerte.

Wir können die Gleichung (6.4) verwenden, um während der Durchführung der QR-Iteration ausgehend von  $\mathbf{Q}^{(0)} = \mathbf{I}$  auch die Matrizen  $\mathbf{Q}^{(m)}$  zu berechnen.

Falls  $\mathbf{A}$  eine normale Matrix ist und die Matrizen  $\mathbf{A}^{(m)}$  gegen Diagonalform konvergieren, bilden die Spalten der Matrizen  $\mathbf{Q}^{(m)}$  eine Orthonormalbasis aus genäherten Eigenvektoren.

**Bemerkung 6.3** Die Wahl  $\mathbf{A}^{(0)} = \mathbf{A}$ ,  $\mathbf{Q}^{(0)} = \mathbf{I}$ , in Algorithmus 6.1 liegt zwar nahe, ist aber häufig nicht die beste. In Abschnitt 6.3 werden wir sehen, dass sich durch eine geschickte Vorgehensweise die Effizienz des Algorithmus erheblich verbessern lässt.



## 6.2 Shift-Strategien und Deflation

Zur Verbesserung der Konvergenzgeschwindigkeit der Vektoriteration haben wir in Kapitel 5 die inverse Iteration mit Shift verwendet: Statt mit der Matrix  $\mathbf{A}$  haben wir mit der Matrix  $(\mathbf{A} - \mu\mathbf{I})^{-1}$  gearbeitet, die dieselben Eigenvektoren wie  $\mathbf{A}$  besitzt, deren Eigenwerte aber durch

$$\lambda \in \sigma(\mathbf{A}) \quad \Longleftrightarrow \quad \frac{1}{\lambda - \mu} \in \sigma((\mathbf{A} - \mu\mathbf{I})^{-1})$$

gegeben sind, so dass wir durch Wahl des Shift-Parameters  $\mu$  in der Nähe eines Eigenwerts für schnelle Konvergenz sorgen können. Dieser Ansatz lässt sich verfeinern, indem  $\mu$  mit Hilfe des Rayleigh-Quotienten automatisch gewählt wird, und in diesem Fall konvergiert die entsprechende Iteration quadratisch, für normale Matrizen sogar kubisch.

Der Einsatz eines Shift-Parameters ist also von großem Vorteil zur Beschleunigung des Verfahrens. Leider macht er es bei der inversen Iteration auch erforderlich, mit der Inversen von  $\mathbf{A} - \mu\mathbf{I}$  zu arbeiten, die uns bei der QR-Iteration nicht ohne weiteres zur Verfügung steht.

Ein genauerer Blick auf die die Konvergenz des Verfahrens beschreibende Abschätzung (6.1) legt nahe, dass wir auch eine andere Shift-Strategie verwenden können: Wenn  $\mu$  hinreichend nahe an einem  $\alpha$ -fachen Eigenwert von  $\mathbf{A}$  liegt, können wir die Eigenwerte der Matrix  $\mathbf{A} - \mu\mathbf{I}$  in die Reihenfolge

$$|\lambda_1 - \mu| \geq \dots \geq |\lambda_{n-\alpha} - \mu| > |\lambda_{n-\alpha+1} - \mu| = \dots = |\lambda_n - \mu|$$

mit  $\lambda_{n-\alpha+1} = \dots = \lambda_n$  bringen. Angewendet auf diese Matrix sollte also der linke untere Matrixblock mit  $\alpha$  Zeilen und  $n - \alpha$  Spalten gegen null konvergieren, und der Iterationsfehler sollte sich durch

$$\left( \frac{|\lambda_n - \mu|}{|\lambda_{n-\alpha} - \mu|} \right)^m$$

abschätzen lassen. Insbesondere sollte die Konvergenz um so schneller werden, je geringer der Abstand zwischen  $\mu$  und dem Eigenwert  $\lambda_{n-\alpha+1} = \lambda_n$  ist.

Eine Verschiebung des Spektrums von  $\mathbf{A}$  kann also auch dann von Vorteil sein, wenn wir nicht die inverse Iteration verwenden. Da bei diesem Ansatz keine Berechnung der Inversen erforderlich ist, können wir den Shift-Parameter praktisch ohne zusätzlichen Rechenaufwand in jedem Schritt anpassen und so beispielsweise die Rayleigh-Quotienten-Strategie verwenden.

Für unsere Variante der QR-Iteration wollen wir weiterhin möglichst mit den Matrizen  $\mathbf{A}^{(m)}$  statt  $\mathbf{A}^{(m)} - \mu\mathbf{I}$  arbeiten. Dieses Ziel können wir erreichen, indem wir (6.5) durch

$$\mathbf{A}^{(m)} - \mu\mathbf{I} = \widehat{\mathbf{Q}}^{(m+1)}\mathbf{R}^{(m+1)}, \quad \mathbf{A}^{(m+1)} = \mathbf{R}^{(m+1)}\widehat{\mathbf{Q}}^{(m+1)} + \mu\mathbf{I} \quad (6.6)$$

ersetzen, denn dann gilt weiterhin

$$\mathbf{A}^{(m+1)} = \mathbf{R}^{(m+1)}\widehat{\mathbf{Q}}^{(m+1)} + \mu\mathbf{I}$$

## 6 Die QR-Iteration

$$\begin{aligned}
&= (\widehat{\mathbf{Q}}^{(m+1)})^* (\mathbf{A}^{(m)} - \mu \mathbf{I}) \widehat{\mathbf{Q}}^{(m+1)} + \mu (\widehat{\mathbf{Q}}^{(m+1)})^* \widehat{\mathbf{Q}}^{(m+1)} \\
&= (\widehat{\mathbf{Q}}^{(m+1)})^* (\mathbf{A} - \mu \mathbf{I} + \mu \mathbf{I}) \widehat{\mathbf{Q}}^{(m+1)} = (\widehat{\mathbf{Q}}^{(m+1)})^* \mathbf{A}^{(m)} \widehat{\mathbf{Q}}^{(m+1)},
\end{aligned}$$

wir berechnen also wie gehabt eine Ähnlichkeitstransformation von  $\mathbf{A}^{(m)}$ , geändert hat sich nur die Wahl der unitären Transformation  $\widehat{\mathbf{Q}}^{(m+1)}$ .

Es stellt sich die Frage nach der geeigneten Wahl des Shift-Parameters  $\mu$ . Da wir hoffen, dass alle Außerdiagonalelemente der  $n$ -ten Zeile der Matrix  $\mathbf{A}^{(m)}$  gegen null konvergieren, können wir annehmen, dass das letzte Diagonalelement  $a_{nn}^{(m)}$  gegen einen Eigenwert konvergieren wird. Unter diesen Annahmen ist also  $\mu = a_{nn}^{(m)}$  eine gute Wahl für den Shift-Parameter.

Diese Wahl lässt sich auch begründen, indem wir auf den Rayleigh-Quotienten zurückgreifen, der sich bereits bei der inversen Iteration als nützlich erwiesen hat: Wir bezeichnen den  $n$ -ten kanonischen Einheitsvektor mit  $\delta^{(n)} \in \mathbb{K}^n$  und die letzte Spalte von  $\mathbf{Q}^{(m)}$  mit  $\mathbf{q}^{(m)} = \mathbf{Q}^{(m)} \delta^{(n)} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  und stellen fest, dass

$$\begin{aligned}
\mu = a_{nn}^{(m)} &= \frac{\langle \delta^{(n)}, \mathbf{A}^{(m)} \delta^{(n)} \rangle_2}{\langle \delta^{(n)}, \delta^{(n)} \rangle_2} = \frac{\langle \delta^{(n)}, (\mathbf{Q}^{(m)})^* \mathbf{A} \mathbf{Q}^{(m)} \delta^{(n)} \rangle_2}{\langle \delta^{(n)}, (\mathbf{Q}^{(m)})^* \mathbf{Q}^{(m)} \delta^{(n)} \rangle_2} \\
&= \frac{\langle \mathbf{Q}^{(m)} \delta^{(n)}, \mathbf{A} \mathbf{Q}^{(m)} \delta^{(n)} \rangle_2}{\langle \mathbf{Q}^{(m)} \delta^{(n)}, \mathbf{Q}^{(m)} \delta^{(n)} \rangle_2} = \frac{\langle \mathbf{q}^{(m)}, \mathbf{A} \mathbf{q}^{(m)} \rangle_2}{\langle \mathbf{q}^{(m)}, \mathbf{q}^{(m)} \rangle_2} = \Lambda_{\mathbf{A}}(\mathbf{q}^{(m)})
\end{aligned}$$

gilt. Damit entspricht unsere Wahl des Shift-Parameters gerade der Verwendung des Rayleigh-Quotienten, und analog zu Satz 5.36 dürfen wir auf Konvergenz gegen einen Eigenwert hoffen, falls  $\mathbf{q}^{(m)}$  gegen einen Eigenvektor konvergiert.

Falls  $\mathbf{A}^{(m)}$  näherungsweise von oberer Dreiecksgestalt ist, ist die adjungierte Matrix  $(\mathbf{A}^{(m)})^*$  näherungsweise von unterer Dreiecksgestalt, und es gilt

$$(\mathbf{A}^{(m)})^* \delta^{(n)} \approx \bar{a}_{nn}^{(m)} \delta^{(n)},$$

und damit wegen der Lemmas 3.18 und 3.35 auch

$$\begin{aligned}
\mathbf{A}^* \mathbf{q}^{(m)} &= \mathbf{Q}^{(m)} (\mathbf{Q}^{(m)})^* \mathbf{A}^* \mathbf{Q}^{(m)} \delta^{(n)} = \mathbf{Q}^{(m)} ((\mathbf{Q}^{(m)})^* \mathbf{A} \mathbf{Q}^{(m)})^* \delta^{(n)} \\
&= \mathbf{Q}^{(m)} (\mathbf{A}^{(m)})^* \delta^{(n)} \approx \bar{a}_{nn}^{(m)} \mathbf{Q}^{(m)} \delta^{(n)} = \bar{a}_{nn}^{(m)} \mathbf{q}^{(m)},
\end{aligned}$$

also approximiert  $\mathbf{q}^{(m)}$  einen Eigenvektor der adjungierten Matrix  $\mathbf{A}^*$ . Indem wir Lemma 5.9 auf  $\mathbf{A}^*$  statt  $\mathbf{A}$  anwenden, folgt die gewünschte Konvergenzaussage, und wir dürfen auf quadratische oder, falls  $\mathbf{A}$  normal ist, sogar auf kubische Konvergenz hoffen.

Mit Hilfe dieser Modifikation können wir also darauf hoffen, dass die Matrizen  $\mathbf{A}^{(m)}$  schnell gegen die Form

$$\begin{pmatrix} \tilde{\mathbf{A}} & \mathbf{C} \\ & \lambda \end{pmatrix}$$

konvergieren werden, dass also

$$\mathbf{A} \approx \mathbf{Q}^{(m)} \begin{pmatrix} \tilde{\mathbf{A}} & \mathbf{C} \\ & \lambda \end{pmatrix} (\mathbf{Q}^{(m)})^*$$

für ein relativ kleines  $m$  gelten wird. Sobald die Einträge im linken unteren Block klein genug sind, können wir uns darauf konzentrieren, nur noch den linken oberen Block  $\tilde{\mathbf{A}}$  auf obere Dreiecksgestalt zu bringen. Falls wir nämlich eine näherungsweise Schur-Zerlegung

$$\tilde{\mathbf{A}} \approx \tilde{\mathbf{Q}}\tilde{\mathbf{R}}\tilde{\mathbf{Q}}^*$$

gefunden haben, ist durch

$$\mathbf{Q}^{(m)} \begin{pmatrix} \tilde{\mathbf{Q}} & \\ & \mathbf{1} \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{R}} & \tilde{\mathbf{Q}}^* \mathbf{C} \\ & \lambda \end{pmatrix} \begin{pmatrix} \tilde{\mathbf{Q}}^* & \\ & \mathbf{1} \end{pmatrix} (\mathbf{Q}^{(m)})^* \approx \mathbf{Q}^{(m)} \begin{pmatrix} \tilde{\mathbf{A}} & \mathbf{C} \\ & \lambda \end{pmatrix} (\mathbf{Q}^{(m)})^* \approx \mathbf{A}$$

eine näherungsweise Schur-Zerlegung der Gesamtmatrix  $\mathbf{A}$  gegeben. Indem wir das Konvergenzverhalten der Blöcke unterhalb der Diagonalen kontrollieren, können wir also bereits konvergierte Teile der Matrix abspalten und so die Problemdimension nach und nach reduzieren.

Dieser Ansatz, also das Entfernen bereits konvergierter Diagonalblöcke aus dem Verfahren, trägt den Namen *Deflation* und sorgt einerseits für eine Reduktion des Rechenaufwands, während andererseits auch dafür gesorgt wird, dass bereits konvergierte Eigenwerte nicht mehr weiter als Shift-Parameter verwendet werden.

In der Praxis zeigt sich, dass mit Rayleigh-Shift in der Regel bereits nach sehr wenigen Schritten ein Eigenwert hinreichend gut approximiert ist, um die Deflation anwenden zu können. Insgesamt werden bei diesem Ansatz dann nur  $\sim n$  Iterationen benötigt, um die Matrix auf obere Dreiecksgestalt zu bringen.

**Übungsaufgabe 6.4 (Deflation)** Seien  $n \in \mathbb{N}$  und  $m \in [1 : n - 1]$  gegeben, und sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine Matrix der Gestalt

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{0} & \mathbf{A}_{22} \end{pmatrix}, \quad \mathbf{A}_{11} \in \mathbb{K}^{m \times m}, \quad \mathbf{A}_{22} \in \mathbb{K}^{(n-m) \times (n-m)}.$$

- (a) Sei  $\lambda \in \mathbb{K}$  ein Eigenwert der Matrix  $\mathbf{A}_{11}$ , und sei  $\hat{\mathbf{e}} \in \mathbb{K}^m \setminus \{\mathbf{0}\}$  ein zugehöriger Eigenvektor. Zeigen Sie, dass  $\lambda$  auch ein Eigenwert der Matrix  $\mathbf{A}$  ist, und geben Sie einen zugehörigen Eigenvektor an.
- (b) Sei  $\lambda \in \mathbb{K}$  ein Eigenwert der Matrix  $\mathbf{A}_{22}$ , und sei  $\hat{\mathbf{e}} \in \mathbb{K}^{n-m} \setminus \{\mathbf{0}\}$  ein zugehöriger Eigenvektor. Zeigen Sie, dass  $\lambda$  auch ein Eigenwert der Matrix  $\mathbf{A}$  ist, und geben Sie einen zugehörigen Eigenvektor an.

Hinweis: Teil (b) ist etwas schwieriger als Teil (a). Um einen Eigenvektor anzugeben, dürfen Sie auf die Eigenvektoren der Teilmatrizen zurückgreifen und lineare Gleichungssysteme lösen.

## 6.3 Hessenberg-Form

Bisher haben wir uns auf die Konvergenzrate der QR-Iteration konzentriert und den Aufwand für die Durchführung eines QR-Schrittes vernachlässigt. Dieser Aufwand ist sehr

## 6 Die QR-Iteration

hoch: Die Berechnung einer QR-Zerlegung erfordert in der Regel  $\sim n^3$  Operationen, und die Berechnung des Produkks  $\mathbf{R}^{(m+1)}\widehat{\mathbf{Q}}^{(m+1)}$  hat einen ähnlich hohen Rechenaufwand. Wenn wir davon ausgehen, dass wir  $\sim n$  Iterationen benötigen, um die Matrix auf obere Dreiecksgestalt zu bringen, würde unser Algorithmus  $\sim n^4$  Operationen benötigen und wäre damit sehr aufwendig.

**Definition 6.5 (Hessenberg-Form)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Falls

$$a_{ij} = 0 \quad \text{für alle } i, j \in [1 : n] \text{ mit } i > j + 1$$

gilt, bezeichnet man  $\mathbf{A}$  als Matrix in (oberer) Hessenberg-Form oder als (obere) Hessenberg-Matrix.

Jede obere Dreiecksmatrix ist auch eine Hessenberg-Matrix, aber bei einer Hessenberg-Matrix sind auch noch Einträge in der unteren Nebendiagonalen der Matrix erlaubt: Hessenberg-Matrizen haben die Gestalt

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ & a_{32} & a_{33} & \dots & a_{3n} \\ & & \ddots & \ddots & \vdots \\ & & & a_{n,n-1} & a_{nn} \end{pmatrix}. \quad (6.7)$$

Diese Struktur bietet den Vorteil, dass sich die QR-Zerlegung einer Hessenberg-Matrix besonders einfach berechnen lässt: Wir wählen eine Givens-Rotation  $\mathbf{G}_1 \in \mathbb{K}^{n \times n}$ , die den Eintrag  $a_{21}$  eliminiert, indem die erste und zweite Zeile miteinander kombiniert werden, beispielsweise als

$$\mathbf{G}_1 = \begin{pmatrix} \bar{c} & \bar{s} & & & \\ -s & c & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{pmatrix}, \quad c = \frac{a_{11}}{r}, \quad s = \frac{a_{21}}{r}, \quad r = \sqrt{|a_{11}|^2 + |a_{21}|^2}.$$

Entsprechend wählen wir weitere Rotationen  $\mathbf{G}_2, \dots, \mathbf{G}_{n-1} \in \mathbb{K}^{n \times n}$ , die der Reihe nach die Einträge  $a_{32}, \dots, a_{n,n-1}$  eliminieren. Zu Illustration stellen wir potentiell von null verschiedene Einträge in der folgenden Skizze abkürzend mit  $\times$  dar und als  $\boxtimes$ , falls sie im aktuellen Schritt verändert wurden:

$$\mathbf{A} \rightsquigarrow \mathbf{G}_1 \mathbf{A} = \begin{pmatrix} \boxtimes & \boxtimes & \boxtimes & \dots \\ \mathbf{0} & \boxtimes & \boxtimes & \dots \\ & \times & \times & \dots \\ & & \times & \ddots \\ & & & \ddots \end{pmatrix} \rightsquigarrow \mathbf{G}_2 \mathbf{G}_1 \mathbf{A} = \begin{pmatrix} \times & \times & \times & \dots \\ 0 & \boxtimes & \boxtimes & \dots \\ & \mathbf{0} & \boxtimes & \dots \\ & & \times & \ddots \\ & & & \ddots \end{pmatrix} \rightsquigarrow \dots \rightsquigarrow \mathbf{R}.$$

Damit überführen wir die Matrix in eine obere Dreiecksmatrix  $\mathbf{R} \in \mathbb{K}^{n \times n}$ , und es gelten

$$\mathbf{G}_{n-1} \dots \mathbf{G}_1 \mathbf{A} = \mathbf{R}, \quad \mathbf{A} = \mathbf{G}_1^* \dots \mathbf{G}_{n-1}^* \mathbf{R},$$

also haben wir die QR-Zerlegung in  $\sim n^2$  Operationen berechnet, und der Faktor  $\mathbf{Q}$  ist durch

$$\mathbf{Q} = \mathbf{G}_1^* \dots \mathbf{G}_{n-1}^* \quad (6.8)$$

gegeben. Für die zweite Hälfte des QR-Schritts müssen wir nun

$$\mathbf{RQ} = \mathbf{RG}_1^* \mathbf{G}_2^* \dots \mathbf{G}_{n-1}^*$$

berechnen. Das entspricht gerade der Anwendung der Givens-Rotationen auf die Spalten der Matrix  $\mathbf{R}$ , die Reihenfolge der Rotationen ist dieselbe wie bei der Berechnung der Zerlegung. Deshalb benötigt auch diese Berechnung  $\sim n^2$  Operationen, so dass sich ein vollständiger Schritt der QR-Iteration für eine Hessenberg-Matrix in  $\sim n^2$  Operationen durchführen lässt.

Besonders wichtig ist, dass  $\mathbf{RQ}$  wieder eine Matrix in Hessenberg-Form ist: Da  $\mathbf{R}$  eine obere Dreiecksmatrix ist, wird in  $\mathbf{RG}_1^*$  lediglich in der ersten Spalte ein Eintrag unterhalb der Diagonalen entstehen, in  $\mathbf{RG}_1^* \mathbf{G}_2^*$  entsteht einer in der zweiten Spalte, und so weiter:

$$\begin{aligned} \mathbf{R} = \begin{pmatrix} \times & \times & \times & \times & \dots \\ 0 & \times & \times & \times & \dots \\ & 0 & \times & \times & \dots \\ & & 0 & \times & \dots \\ & & & 0 & \ddots \end{pmatrix} \rightsquigarrow \mathbf{RG}_1^* = \begin{pmatrix} \boxtimes & \boxtimes & \times & \times & \dots \\ \boxtimes & \boxtimes & \times & \times & \dots \\ & 0 & \times & \times & \dots \\ & & 0 & \times & \dots \\ & & & 0 & \ddots \end{pmatrix} \\ \rightsquigarrow \mathbf{RG}_1^* \mathbf{G}_2^* = \begin{pmatrix} \times & \boxtimes & \boxtimes & \times & \dots \\ \times & \boxtimes & \boxtimes & \times & \dots \\ & \boxtimes & \boxtimes & \times & \dots \\ & & 0 & \times & \dots \\ & & & 0 & \ddots \end{pmatrix} \rightsquigarrow \dots \rightsquigarrow \mathbf{RG}_1^* \dots \mathbf{G}_{n-1}^*. \end{aligned}$$

Diese Eigenschaft ist von entscheidender Bedeutung: Sofern die Ausgangsmatrix  $\mathbf{A}^{(0)}$  in Hessenberg-Form vorliegt, werden alle gemäß der oben beschriebenen Vorgehensweise berechneten Matrizen  $\mathbf{A}^{(m)}$  der QR-Iteration ebenfalls in Hessenberg-Form sein, so dass *jeder* Iterationsschritt nur  $\sim n^2$  Operationen erfordert. Das ist wesentlich effizienter als die  $\sim n^3$  Operationen, die ohne Ausnutzung der Hessenberg-Form nötig wären.

Unser Ziel sollte es nun also sein, die erste Iterierte  $\mathbf{A}^{(0)}$  möglichst in Hessenberg-Form zu überführen, indem wir die Basis  $\mathbf{Q}^{(0)}$  geschickter als bisher wählen. Auch dieses Ziel lässt sich mit Hilfe geeigneter Givens-Rotationen erreichen: Um eine beliebige Matrix

$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ a_{31} & a_{32} & \dots & a_{3n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{pmatrix}$$

## 6 Die QR-Iteration

in die Hessenberg-Form zu überführen, müssen wir eine unitäre Ähnlichkeitstransformation finden, die die Elemente  $a_{31}, \dots, a_{n1}, a_{42}, \dots, a_{n2}, \dots$  eliminiert. Falls wir, wie zuvor beschrieben, das Element  $a_{31}$  durch Kombination der ersten und dritten Zeile eliminieren würden, müssten wir, da wir jetzt eine Ähnlichkeitstransformation benötigen, auch die erste und dritte Spalte kombinieren, und dadurch könnte der Eintrag  $a_{31}$  wieder einen anderen Wert als null erhalten.

Also kombinieren wir die *zweite* mit der dritten Zeile, um  $a_{31}$  zu eliminieren. Die korrespondierende Transformation der Spalten beeinflusst dann nur die zweite und dritte Spalte, so dass die in  $a_{31}$  eingeführte Null erhalten bleibt.

Im nächsten Schritt verwenden wir die zweite und die vierte Zeile, um  $a_{41}$  zu eliminieren, dann die zweite und die fünfte, um  $a_{51}$  zu null werden zu lassen. Da keine der Spaltentransformationen die erste Spalte betrifft, bleiben die von den Zeilentransformationen eingeführten Nullen dort erhalten:

$$\begin{aligned}
 \mathbf{A} \rightsquigarrow \mathbf{G}_{31}\mathbf{A} &= \begin{pmatrix} \times & \times & \times & \times & \times & \dots \\ \boxtimes & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ \mathbf{0} & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ \times & \times & \times & \times & \times & \dots \\ \times & \times & \times & \times & \times & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix} \rightsquigarrow \mathbf{G}_{31}\mathbf{A}\mathbf{G}_{31}^* = \begin{pmatrix} \times & \boxtimes & \boxtimes & \times & \times & \dots \\ \times & \boxtimes & \boxtimes & \times & \times & \dots \\ 0 & \boxtimes & \boxtimes & \times & \times & \dots \\ \times & \boxtimes & \boxtimes & \times & \times & \dots \\ \times & \boxtimes & \boxtimes & \times & \times & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix} \\
 \rightsquigarrow \mathbf{G}_{41}\mathbf{G}_{31}\mathbf{A}\mathbf{G}_{31}^* &= \begin{pmatrix} \times & \times & \times & \times & \times & \dots \\ \boxtimes & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ 0 & \times & \times & \times & \times & \dots \\ \mathbf{0} & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ \times & \times & \times & \times & \times & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix} \\
 \rightsquigarrow \mathbf{G}_{41}\mathbf{G}_{31}\mathbf{A}\mathbf{G}_{31}^*\mathbf{G}_{41}^* &= \begin{pmatrix} \times & \boxtimes & \times & \boxtimes & \times & \dots \\ \times & \boxtimes & \times & \boxtimes & \times & \dots \\ 0 & \boxtimes & \times & \boxtimes & \times & \dots \\ 0 & \boxtimes & \times & \boxtimes & \times & \dots \\ \times & \boxtimes & \times & \boxtimes & \times & \dots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \ddots \end{pmatrix} \\
 \rightsquigarrow \mathbf{G}_{n,n-2} \dots \mathbf{G}_{31}\mathbf{A}\mathbf{G}_{31}^* \dots \mathbf{G}_{n,n-2}^* &= (\mathbf{Q}^{(0)})^* \mathbf{A} \mathbf{Q}^{(0)}.
 \end{aligned}$$

Ein gelegentlich sehr nützlicher Nebeneffekt der hier verwendeten Givens-Rotationen besteht darin, dass die entstehenden Einträge unterhalb der Diagonalen reell und nicht-negativ sind, denn Operationen mit reellen Zahlen benötigen erheblich weniger Zeit als solche mit komplexen.

Falls wir auch an den Eigenvektoren interessiert sind, können wir die Matrix  $\mathbf{Q}^{(0)}$  konstruieren, indem wir ausgehend von der Einheitsmatrix die Adjungierte jeder verwendeten Givens-Rotation auf die Spalten der Matrix anwenden, um

$$\mathbf{I} \rightsquigarrow \mathbf{I}\mathbf{G}_{31}^* \rightsquigarrow \mathbf{I}\mathbf{G}_{31}^*\mathbf{G}_{41}^* \rightsquigarrow \dots \rightsquigarrow \mathbf{Q}^{(0)} = \mathbf{G}_{31}^*\mathbf{G}_{41}^* \dots \mathbf{G}_{n,n-2}^*$$

zu erhalten.

**Algorithmus 6.6** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Der folgende Algorithmus überschreibt  $\mathbf{A}$  mit einer Matrix  $\mathbf{A}^{(0)}$  in Hessenberg-Form, die die Gleichung  $\mathbf{A}^{(0)} = (\mathbf{Q}^{(0)})^* \mathbf{A} \mathbf{Q}^{(0)}$  erfüllt.

```

 $\mathbf{Q}^{(0)} \leftarrow \mathbf{I}$ 
for  $j = 1$  to  $n - 2$  do
  for  $i = j + 2$  to  $n$  do begin
     $r \leftarrow \sqrt{|a_{j+1,j}|^2 + |a_{ij}|^2}$ ;  $c \leftarrow a_{j+1,j}/r$ ;  $s \leftarrow a_{ij}/r$ 
     $a_{j+1,j} \leftarrow r$ ;  $a_{ij} \leftarrow 0$ 
    for  $k \in [j + 1 : n]$  do begin
       $h \leftarrow a_{j+1,k}$ ;  $a_{j+1,k} \leftarrow \bar{c}h + \bar{s}a_{ik}$ ;  $a_{ik} \leftarrow -sh + ca_{ik}$ 
    end
    for  $k \in [1 : n]$  do begin
       $h \leftarrow a_{k,j+1}$ ;  $a_{k,j+1} \leftarrow ch + sa_{ki}$ ;  $a_{ki} \leftarrow -\bar{s}h + \bar{c}a_{ki}$ 
       $h \leftarrow q_{k,j+1}^{(0)}$ ;  $q_{k,j+1}^{(0)} \leftarrow ch + sq_{ki}^{(0)}$ ;  $q_{ki}^{(0)} \leftarrow -\bar{s}h + \bar{c}q_{ki}^{(0)}$ 
    end
  end
end

```

Mit diesem Algorithmus können wir eine beliebige Matrix in  $\sim n^3$  Operationen in Hessenberg-Form überführen, so dass sich dann die weiteren Schritte der QR-Iteration effizient durchführen lassen.

**Bemerkung 6.7 (Householder-Spiegelungen)** Selbstverständlich können wir statt Givens-Rotationen auch Householder-Spiegelungen verwenden, um eine beliebige Matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  in Hessenberg-Gestalt zu transformieren.

Im ersten Schritt zerlegen wir die Matrix  $\mathbf{A}$  in der Form

$$\mathbf{A} = \begin{pmatrix} a_{11} & \mathbf{A}_{1*} \\ \mathbf{A}_{*1} & \mathbf{A}_{**} \end{pmatrix}, \quad \mathbf{A}_{**} \in \mathbb{K}^{(n-1) \times (n-1)},$$

und wählen eine Householder-Spiegelung  $\mathbf{H}_1 \in \mathbb{K}^{(n-1) \times (n-1)}$ , die  $\mathbf{A}_{*1}$  auf ein Vielfaches des ersten kanonischen Einheitsvektors abbildet.

Dann erhalten wir

$$\begin{pmatrix} 1 & \\ & \mathbf{H}_1 \end{pmatrix} \begin{pmatrix} a_{11} & \mathbf{A}_{1*} \\ \mathbf{A}_{*1} & \mathbf{A}_{**} \end{pmatrix} \begin{pmatrix} 1 & \\ & \mathbf{H}_1^* \end{pmatrix} = \begin{pmatrix} a_{11} & \mathbf{A}_{1*} \mathbf{H}_1^* \\ \alpha \delta^{(1)} & \mathbf{H}_1 \mathbf{A}_{**} \mathbf{H}_1^* \end{pmatrix},$$

die erste Spalte weist also die gewünschte Form auf. Wir können induktiv fortfahren, um eine unitäre Matrix  $\widehat{\mathbf{H}} \in \mathbb{K}^{(n-2) \times (n-2)}$  so zu finden, dass die unitäre Matrix

$$\begin{pmatrix} 1 & \\ & \widehat{\mathbf{H}} \end{pmatrix} \in \mathbb{K}^{(n-1) \times (n-1)}$$

die Matrix  $\widehat{\mathbf{A}} := \mathbf{H}_1 \mathbf{A}_{**} \mathbf{H}_1^*$  in Hessenberg-Gestalt bringt. Die Kombination beider Transformationen leistet dann das Gewünschte.

In praktischen Implementierungen sind Householder-Spiegelungen in der Regel erheblich effizienter als Givens-Rotationen. In diesem Skript wurde den Givens-Rotationen nur der Vorzug gegeben, um den Schreibaufwand zu reduzieren.

## 6 Die QR-Iteration

Falls zusätzliche Informationen über die Struktur von  $\mathbf{A}$  vorliegen, können wir sie ausnutzen, um den Algorithmus effizienter zu gestalten: Bei Bandmatrizen der Bandbreite  $k$  etwa genügen  $\sim n^2k$  Operationen für die Transformation. Besonders günstig ist die Transformation auf Hessenberg-Form, falls  $\mathbf{A}$  selbstadjungiert ist:

**Bemerkung 6.8** Falls  $\mathbf{A}$  selbstadjungiert ist, muss auch  $\mathbf{A}^{(0)}$  selbstadjungiert sein. Damit folgt aus  $a_{ij}^{(0)} = 0$  auch  $a_{ji}^{(0)} = 0$  für alle  $1 < j + 1 < i < n$ , die Matrix  $\mathbf{A}^{(0)}$  ist also tridiagonal. Man kann sich einfach überlegen, dass der beschriebene QR-Schritt für eine tridiagonale Matrix sogar nur  $\sim n$  Operationen erfordert, die Iteration wird also im Fall selbstadjungierter Matrizen noch wesentlich effizienter als im allgemeinen Fall sein.

Bei komplexwertigen selbstadjungierten Matrizen kommt hinzu, dass die von uns gewählten Givens-Rotationen zu reellen Nebendiagonaleinträgen führen. Da bei einer selbstadjungierten Matrix die Diagonaleinträge ohnehin reell sind, entsteht eine vollständig reelle Tridiagonalmatrix, so dass wir bei der Durchführung der QR-Iteration viel Rechenzeit sparen können.

Da eine Hessenberg-Matrix schon „fast“ von oberer Dreiecksgestalt ist, können wir besonders einfach feststellen, wann eine Teilmatrix unterhalb der Diagonalen konvergiert ist und wir die Dimension reduzieren können:

**Bemerkung 6.9** Sobald  $|a_{i+1,i}|$  für ein  $i \in [1 : n - 1]$  hinreichend klein geworden ist, können wir per Deflation zu einer kleineren Teilmatrix übergehen. In der Praxis verwendet man skalierungsinvariante Kriterien der Form

$$|a_{i+1,i}| \leq \epsilon(|a_{ii}| + |a_{i+1,i+1}|)$$

mit einer Fehlerschranke  $\epsilon \in \mathbb{R}_{>0}$ , um zu erkennen, wann eine Teilmatrix unterhalb der Diagonalen konvergiert ist.

Die Hessenberg-Form bietet uns auch die Möglichkeit, den Shift-Parameter geschickter als bisher zu wählen: Die Idee des *Wilkinson-Shifts* besteht darin, nicht nur den rechten unteren Diagonaleintrag, also das Gegenstück des Rayleigh-Quotienten, zu verwenden, sondern stattdessen die Eigenwerte des rechten unteren  $2 \times 2$ -Diagonalblocks zu benutzen.

Im selbstadjungierten Fall besitzt  $\mathbf{A}^{(m)}$  Tridiagonalgestalt, und der Wilkinson-Shift-Parameter ergibt sich aus der Untersuchung der Eigenwerte der  $2 \times 2$ -Matrix

$$\mathbf{S} := \begin{pmatrix} a_{n-1,n-1}^{(m)} & a_{n-1,n}^{(m)} \\ a_{n-1,n}^{(m)} & a_{nn}^{(m)} \end{pmatrix}.$$

Wir können analog zu Abschnitt 4.2 vorgehen und erhalten die Eigenwerte

$$\lambda_1 = m + \sqrt{d^2 + |a_{n-1,n}^{(m)}|^2}, \quad \lambda_2 = m - \sqrt{d^2 + |a_{n-1,n}^{(m)}|^2},$$

wobei wir die Hilfsgrößen

$$m := \frac{a_{n-1,n-1}^{(m)} + a_{nn}^{(m)}}{2}, \quad d := \frac{a_{n-1,n-1}^{(m)} - a_{nn}^{(m)}}{2}$$



für den Mittelwert und die halbe Differenz der Diagonalelemente verwenden.

Da wir darauf hoffen, dass das letzte Diagonalelement  $a_{nn}^{(m)}$  gegen einen Eigenwert konvergiert, bietet es sich an, denjenigen Eigenwert als Shift-Parameter  $\mu$  zu wählen, der dem derzeitigen Wert von  $a_{nn}^{(m)}$  am nächsten liegt. Wir haben

$$\begin{aligned}\lambda_1 - a_{nn}^{(m)} &= m - a_{nn}^{(m)} + \sqrt{d^2 + |a_{n-1,n}^{(m)}|^2} \\ &= \frac{a_{n-1,n-1}^{(m)} + a_{nn}^{(m)}}{2} - a_{nn}^{(m)} + \sqrt{d^2 + |a_{n-1,n}^{(m)}|^2} \\ &= d + \sqrt{d^2 + |a_{n-1,n}^{(m)}|^2}, \\ \lambda_2 - a_{nn}^{(m)} &= d - \sqrt{d^2 + |a_{n-1,n}^{(m)}|^2},\end{aligned}$$

also müssen wir  $\lambda_1$  wählen, falls  $d$  negativ ist, und ansonsten  $\lambda_2$ :

$$\mu = \begin{cases} m - \sqrt{d^2 + |a_{n-1,n}^{(m)}|^2} & \text{falls } d \geq 0, \\ m + \sqrt{d^2 + |a_{n-1,n}^{(m)}|^2} & \text{ansonsten.} \end{cases} \quad (6.9)$$

Indem wir dieses  $\mu$  als Shift-Parameter in der QR-Iteration verwenden, ergibt sich die *QR-Iteration mit Wilkinson-Shift*, bei der wir auf bessere Konvergenz als bei der Verwendung des einfachen Rayleigh-Quotienten  $\mu = a_{nn}^{(m)}$  hoffen dürfen.

**Bemerkung 6.10 (Stabile Berechnung)** Falls  $|a_{n-1,n}^{(m)}|$  klein ist, wird die Wurzel ungefähr gleich  $|d|$  sein, und falls auch  $|a_{nn}^{(m)}|$  klein ist, könnte es bei der Berechnung des Shift-Werts zu Auslöschungseffekten kommen. Da  $\mu$  „nur“ als Shift-Wert in der QR-Iteration auftritt, ist das weniger kritisch als bei der Jacobi-Iteration.

Bei Bedarf können wir die Gleichung

$$\lambda_1 \lambda_2 = m^2 - (d^2 + |a_{n-1,n}^{(m)}|^2) = a_{n-1,n-1}^{(m)} a_{nn}^{(m)} - |a_{n-1,n}^{(m)}|^2$$

ausnutzen, um aus dem stabil berechenbaren betragsgrößeren der beiden Eigenwerte den gewünschten betragskleineren zu bestimmen.

Auch im Fall nicht selbstadjungierter Matrizen lassen sich verfeinerte Shift-Strategien entwickeln, allerdings muss hier berücksichtigt werden, dass komplexe Shift-Parameter auftreten können. Falls wir ohnehin mit einer komplexen Matrix arbeiten, ist das nicht weiter problematisch.

Falls wir allerdings mit einer reellen Matrix rechnen, wäre ein „Umweg“ über die komplexen Zahlen unattraktiv. Leider lässt er sich aber oft nicht vermeiden, denn reelle Matrizen können durchaus „echt komplexe“ Eigenwerte besitzen, wie wir im Beispiel 3.10 gesehen haben. Bei derartigen Matrizen können wir mit reellen Shift-Parametern die komplexen Eigenwerte nicht beliebig gut approximieren, müssten also potentiell erhebliche Einschränkungen der Konvergenzgeschwindigkeit in Kauf nehmen.

## 6 Die QR-Iteration

Dieses Problem lässt sich lösen, indem man zwei QR-Schritte kombiniert: Der erste verwendet einen komplexen Shift-Parameter  $\mu \in \mathbb{C}$ , der zweite verwendet den konjugierten Parameter  $\bar{\mu}$ . Man kann beweisen, dass nach dem zweiten QR-Schritt die Matrix wieder reell ist.

**Übungsaufgabe 6.11 (Doppelshift)** Sei  $n \in \mathbb{N}$ .

(a) Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ , und seien  $\mu_1, \mu_2 \in \mathbb{K}$ . Eine konventionelle QR-Iteration berechnet

$$\begin{aligned} \mathbf{Q}_1 \mathbf{R}_1 &= \mathbf{A} - \mu_1 \mathbf{I}, & \mathbf{A}_1 &:= \mathbf{R}_1 \mathbf{Q}_1 + \mu_1 \mathbf{I}, \\ \mathbf{Q}_2 \mathbf{R}_2 &= \mathbf{A}_1 - \mu_2 \mathbf{I}, & \mathbf{A}_2 &:= \mathbf{R}_2 \mathbf{Q}_2 + \mu_2 \mathbf{I} \end{aligned}$$

mit QR-Zerlegungen  $(\mathbf{Q}_1, \mathbf{R}_1)$  und  $(\mathbf{Q}_2, \mathbf{R}_2)$ .

Wir können alternativ einen Doppelschritt

$$\mathbf{Q}_d \mathbf{R}_d = (\mathbf{A} - \mu_2 \mathbf{I})(\mathbf{A} - \mu_1 \mathbf{I}), \quad \mathbf{A}_d := \mathbf{Q}_d^* \mathbf{A} \mathbf{Q}_d$$

mit einer QR-Zerlegung  $(\mathbf{Q}_d, \mathbf{R}_d)$  durchführen.

Unter diesen Voraussetzungen und der Annahme, dass  $\mu_1$  und  $\mu_2$  keine Eigenwerte der Matrix  $\mathbf{A}$  sind, beweisen Sie, dass eine unitäre Diagonalmatrix  $\mathbf{D} \in \mathbb{K}^{n \times n}$  mit  $\mathbf{A}_d = \mathbf{D}^* \mathbf{A}_2 \mathbf{D}$  existiert.

(b) Sei  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , sei  $\mu \in \mathbb{C}$  kein Eigenwert der Matrix  $\mathbf{A}$ .

Beweisen Sie, dass die QR-Zerlegung bei einem Doppelschritt mit  $\mu_1 = \mu$  und  $\mu_2 = \bar{\mu}$  so gewählt werden kann, dass  $\mathbf{A}_d$  wieder reell ist.

(c) Sei  $\mathbf{A}$  in Hessenberg-Form. Beweisen Sie für das Produkt

$$\mathbf{B} := (\mathbf{A} - \mu_2 \mathbf{I})(\mathbf{A} - \mu_1 \mathbf{I})$$

die Eigenschaft

$$b_{ij} = 0 \quad \text{für alle } i, j \in [1 : n] \text{ mit } i > j + 2.$$

Hinweis: Sie dürfen voraussetzen, dass zwei QR-Zerlegungen  $(\mathbf{Q}, \mathbf{R})$  und  $(\tilde{\mathbf{Q}}, \tilde{\mathbf{R}})$  derselben invertierbaren Matrix  $\mathbf{X} \in \mathbb{K}^{n \times n}$  sich nur durch eine unitäre Diagonalmatrix  $\mathbf{D} \in \mathbb{K}^{n \times n}$  unterscheiden, die  $\mathbf{Q} = \tilde{\mathbf{Q}} \mathbf{D}$  und  $\mathbf{R} = \mathbf{D}^* \tilde{\mathbf{R}}$  erfüllt.

**Bemerkung 6.12 (Eigenvektoren)** Wir können den Rechenaufwand der QR-Iteration deutlich reduzieren, indem wir auf die Akkumulation der Transformationsmatrizen  $\mathbf{Q}^{(m)}$  verzichten: Für Hessenberg-Matrizen reduziert sich so der Aufwand eines QR-Schritts von  $\mathcal{O}(n^3)$  auf  $\mathcal{O}(n^2)$ , für tridiagonale Matrizen sogar auf  $\mathcal{O}(n)$ .

Um die Eigenvektoren zu rekonstruieren, bietet sich die Verwendung einer inversen Iteration an: Da alle Eigenwerte bekannt sind, können wir sehr gute Shift-Parameter wählen und also auf sehr schnelle Konvergenz hoffen. Wenn wir die inverse Iteration auf die Hessenberg-Matrix  $\mathbf{A}^{(0)}$  statt direkt auf  $\mathbf{A}$  anwenden, können wir die einzelnen Schritte der inversen Iteration effizient durchführen.

**Übungsaufgabe 6.13 (Matrixfreie Hessenberg-Transformation)** Sei  $n \in \mathbb{N}$ , und sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine beliebige Matrix.

- (a) Sei  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  eine unitäre Matrix derart, dass  $\mathbf{H} := \mathbf{Q}^* \mathbf{A} \mathbf{Q}$  in Hessenberg-Form ist. Wir bezeichnen mit  $\mathbf{q}^{(i)} := \mathbf{Q} \delta^{(i)}$  für alle  $i \in [1 : n]$  die  $i$ -te Spalte der Matrix  $\mathbf{Q}$ . Beweisen Sie

$$h_{ij} = \langle \mathbf{q}^{(i)}, \mathbf{A} \mathbf{q}^{(j)} \rangle \quad \text{für alle } i, j \in [1 : n],$$

$$h_{j+1,j} \mathbf{q}^{(j+1)} = \mathbf{A} \mathbf{q}^{(j)} - \sum_{i=1}^j h_{ij} \mathbf{q}^{(i)} \quad \text{für alle } j \in [1 : n-1].$$

- (b) Gegeben seien Einheitsvektoren  $\mathbf{q}^{(1)}, \dots, \mathbf{q}^{(n)} \in \mathbb{K}^n$  mit  $\|\mathbf{q}^{(1)}\| = 1$ , die

$$\langle \mathbf{q}^{(j+1)}, \mathbf{A} \mathbf{q}^{(j)} \rangle \neq 0,$$

$$\langle \mathbf{q}^{(j+1)}, \mathbf{A} \mathbf{q}^{(j)} \rangle \mathbf{q}^{(j+1)} = \mathbf{A} \mathbf{q}^{(j)} - \sum_{i=1}^j \langle \mathbf{q}^{(i)}, \mathbf{A} \mathbf{q}^{(j)} \rangle \mathbf{q}^{(i)} \quad \text{für alle } j \in [1 : n-1]$$

erfüllen. Sei  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  die Matrix, deren  $i$ -te Spalte für alle  $i \in [1 : n]$  gerade  $\mathbf{q}^{(i)}$  ist, und sei  $\mathbf{H} := \mathbf{Q}^* \mathbf{A} \mathbf{Q}$ .

Beweisen Sie, dass  $\mathbf{Q}$  unitär und  $\mathbf{H}$  in Hessenberg-Form ist.

**Übungsaufgabe 6.14 (Shifts aus Teilmatrizen)** Sei  $n \in \mathbb{N}_{>1}$ , und sei  $\mathbf{F} \in \mathbb{K}^{n \times n}$  definiert durch

$$f_{ij} = \begin{cases} 1 & \text{falls } j - i \in \{1, 1 - n\}, \\ 0 & \text{ansonsten} \end{cases} \quad \text{für alle } i, j \in [1 : n],$$

die Matrix hat also die Gestalt

$$\mathbf{F} = \begin{pmatrix} 0 & 1 & & & \\ & \ddots & \ddots & & \\ & & 0 & 1 & \\ 1 & & & & 0 \end{pmatrix}.$$

- (a) Sei  $\omega := \exp(2\pi i/n)$  die  $n$ -te komplexe Einheitswurzel ( $i$  ist hier die komplexe Einheit mit  $i^2 = -1$ ). Beweisen Sie  $\sigma(\mathbf{F}) = \{\omega^k : k \in [1 : n]\}$ .
- (b) Bei Shift-Strategien wie dem Rayleigh- oder dem Wilkinson-Shift werden Eigenwerte von Teilmatrizen verwendet. Für ein  $m \in [1 : n-1]$  betrachten wir die Teilmatrix  $\hat{\mathbf{F}} \in \mathbb{K}^{m \times m}$ , die sich aus den letzten  $m$  Zeilen und Spalten der Matrix  $\mathbf{F}$  zusammensetzt, also

$$\hat{f}_{ij} = f_{i+(n-m), j+(n-m)} \quad \text{für alle } i, j \in [1 : m]$$

erfüllt. Beweisen Sie  $\sigma(\hat{\mathbf{F}}) = \{0\}$ .

- (c) Sind Eigenwerte der Matrix  $\hat{\mathbf{F}}$  als Shift-Parameter für die Matrix  $\mathbf{F}$  geeignet?

## 6.4 Implizite Verfahren

Bisher haben wir einen QR-Schritt explizit durchgeführt: Erst wird die QR-Zerlegung  $\widehat{\mathbf{Q}}^{(m+1)}\mathbf{R}^{(m+1)} = \mathbf{A}^{(m)}$  berechnet, dann wird daraus  $\mathbf{A}^{(m+1)} = \mathbf{R}^{(m+1)}\widehat{\mathbf{Q}}^{(m+1)}$  konstruiert. Bei dieser Vorgehensweise müssen wir sämtliche Givens-Rotationen während der Zerlegung speichern, um sie anschließend bei der Multiplikation verwenden zu können.

Wesentlich eleganter wäre es, wenn wir es vermeiden könnten, die Givens-Rotationen explizit zu speichern, wenn wir also den Schritt von  $\mathbf{A}^{(m+1)}$  zu  $\mathbf{A}^{(m)}$  mit einer *impliziten* QR-Zerlegung durchführen könnten.

Diese Möglichkeit besteht, wenn wir mit Hessenberg-Matrizen arbeiten: Im QR-Schritt verwenden wir unitäre Ähnlichkeitstransformationen, die Hessenberg-Matrizen wieder in Hessenberg-Matrizen überführen. Derartige Transformationen unterliegen Gesetzmäßigkeiten, die sich praktisch ausnutzen lassen.

**Bemerkung 6.15** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  in Hessenberg-Form. Die Matrix ist genau dann irreduzibel (siehe Definition 3.67), falls

$$a_{i+1,i} \neq 0 \quad \text{für alle } i \in [1 : n - 1] \text{ gilt.}$$

**Satz 6.16 (Eindeutigkeit)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine Matrix. Seien  $\mathbf{Q}, \widehat{\mathbf{Q}} \in \mathbb{K}^{n \times n}$  unitäre Matrizen derart, dass  $\mathbf{H} := \mathbf{Q}^* \mathbf{A} \mathbf{Q}$  und  $\widehat{\mathbf{H}} := \widehat{\mathbf{Q}}^* \mathbf{A} \widehat{\mathbf{Q}}$  in Hessenberg-Form sind und  $\mathbf{H}$  irreduzibel ist.

Falls die ersten Spalten der Matrizen  $\mathbf{Q}$  und  $\widehat{\mathbf{Q}}$  übereinstimmen, existiert eine unitäre Diagonalmatrix  $\mathbf{D} \in \mathbb{K}^{n \times n}$  mit  $\mathbf{Q} = \widehat{\mathbf{Q}} \mathbf{D}$ , die beiden Matrizen können sich also allenfalls in den Vorzeichen ihrer Spalten unterscheiden.

*Beweis.* Wir bezeichnen mit

$$\mathbf{q}^{(i)} := \mathbf{Q} \delta^{(i)}, \quad \widehat{\mathbf{q}}^{(i)} := \widehat{\mathbf{Q}} \delta^{(i)} \quad \text{für alle } i \in [1 : n]$$

die Spalten der Matrizen  $\mathbf{Q}$  und  $\widehat{\mathbf{Q}}$ . Wir zeigen per abschnittsweiser Induktion, dass Konstanten  $\alpha_1, \dots, \alpha_n \in \mathbb{K}$  existieren mit

$$|\alpha_i| = 1, \quad \mathbf{q}^{(i)} = \alpha_i \widehat{\mathbf{q}}^{(i)} \quad \text{für alle } i \in [1 : n].$$

Daraus folgt die Behauptung mit

$$\mathbf{D} := \begin{pmatrix} \alpha_1 & & \\ & \ddots & \\ & & \alpha_n \end{pmatrix}.$$

*Induktionsanfang:* Nach Voraussetzung gilt  $\mathbf{q}^{(1)} = \widehat{\mathbf{q}}^{(1)}$ , also setzen wir  $\alpha_1 = 1$ .

*Induktionsvoraussetzung:* Sei  $m \in [1 : n - 1]$  so gegeben, dass  $\alpha_1, \dots, \alpha_m \in \mathbb{K}$  existieren mit

$$|\alpha_i| = 1, \quad \mathbf{q}^{(i)} = \alpha_i \widehat{\mathbf{q}}^{(i)} \quad \text{für alle } i \in [1 : m].$$

*Induktionsschritt:* Mit Lemma 3.35 und der Hessenberg-Form der Matrizen  $\mathbf{H}$  und  $\widehat{\mathbf{H}}$  erhalten wir

$$\mathbf{A}\mathbf{q}^{(m)} = \mathbf{Q}\mathbf{Q}^*\mathbf{A}\mathbf{Q}\delta^{(m)} = \mathbf{Q}\mathbf{H}\delta^{(m)} = \mathbf{Q}\sum_{j=1}^{m+1} h_{jm}\delta^{(j)} = \sum_{j=1}^{m+1} h_{jm}\mathbf{q}^{(j)}, \quad (6.10a)$$

$$\mathbf{A}\widehat{\mathbf{q}}^{(m)} = \widehat{\mathbf{Q}}\widehat{\mathbf{Q}}^*\mathbf{A}\widehat{\mathbf{Q}}\delta^{(m)} = \widehat{\mathbf{Q}}\widehat{\mathbf{H}}\delta^{(m)} = \widehat{\mathbf{Q}}\sum_{j=1}^{m+1} \hat{h}_{jm}\delta^{(j)} = \sum_{j=1}^{m+1} \hat{h}_{jm}\widehat{\mathbf{q}}^{(j)}. \quad (6.10b)$$

Mit der Induktionsvoraussetzung folgt wegen Lemma 3.17 einerseits

$$\begin{aligned} h_{jm} &= \langle \delta^{(j)}, \mathbf{H}\delta^{(m)} \rangle = \langle \delta^{(j)}, \mathbf{Q}^*\mathbf{A}\mathbf{Q}\delta^{(m)} \rangle = \langle \mathbf{Q}\delta^{(j)}, \mathbf{A}\mathbf{Q}\delta^{(m)} \rangle = \langle \mathbf{q}^{(j)}, \mathbf{A}\mathbf{q}^{(m)} \rangle \\ &= \langle \alpha_j \widehat{\mathbf{q}}^{(j)}, \mathbf{A}\alpha_m \widehat{\mathbf{q}}^{(m)} \rangle = \bar{\alpha}_j \alpha_m \langle \widehat{\mathbf{q}}^{(j)}, \mathbf{A}\widehat{\mathbf{q}}^{(m)} \rangle = \bar{\alpha}_j \alpha_m \hat{h}_{jm} \quad \text{für alle } j \in [1 : m], \end{aligned}$$

andererseits damit aber auch

$$\begin{aligned} h_{m+1,m}\mathbf{q}^{(m+1)} &= \mathbf{A}\mathbf{q}^{(m)} - \sum_{j=1}^m h_{jm}\mathbf{q}^{(j)} = \alpha_m \mathbf{A}\widehat{\mathbf{q}}^{(m)} - \sum_{j=1}^m h_{jm}\alpha_j \widehat{\mathbf{q}}^{(j)} \\ &= \alpha_m \sum_{j=1}^{m+1} \hat{h}_{jm}\widehat{\mathbf{q}}^{(j)} - \sum_{j=1}^m \bar{\alpha}_j \alpha_m \hat{h}_{jm}\alpha_j \widehat{\mathbf{q}}^{(j)} \\ &= \alpha_m \sum_{j=1}^{m+1} \hat{h}_{jm}\widehat{\mathbf{q}}^{(j)} - \alpha_m \sum_{j=1}^m |\alpha_j|^2 \hat{h}_{jm}\widehat{\mathbf{q}}^{(j)} = \alpha_m \hat{h}_{m+1,m}\widehat{\mathbf{q}}^{(m+1)}. \end{aligned}$$

Da  $\mathbf{H}$  irreduzibel ist, gilt insbesondere  $h_{m+1,m} \neq 0$ , so dass wir

$$\mathbf{q}^{(m+1)} = \alpha_m \frac{\hat{h}_{m+1,m}}{h_{m+1,m}} \widehat{\mathbf{q}}^{(m+1)}$$

bewiesen haben. Wir setzen also

$$\alpha_{m+1} := \alpha_m \frac{\hat{h}_{m+1,m}}{h_{m+1,m}}$$

und halten fest, dass auch

$$1 = \|\mathbf{q}^{(m+1)}\| = \|\alpha_{m+1}\widehat{\mathbf{q}}^{(m+1)}\| = |\alpha_{m+1}| \|\widehat{\mathbf{q}}^{(m+1)}\| = |\alpha_{m+1}|$$

gelten muss, da  $\mathbf{q}^{(m+1)}$  und  $\widehat{\mathbf{q}}^{(m+1)}$  Einheitsvektoren sind.  $\blacksquare$

Tatsächlich lässt sich der Beweis auch auf den nicht vollständig irreduziblen Fall anwenden, um ein für unsere Zwecke völlig ausreichendes Teilresultat zu erzielen:

**Bemerkung 6.17 (Nicht-irreduzibler Fall)** Falls  $\mathbf{A}$  in Satz 6.16 nicht irreduzibel ist, sondern lediglich

$$h_{i+1,i} \neq 0 \quad \text{für alle } i \in [1 : k-1]$$

## 6 Die QR-Iteration

und  $h_{k+1,k} = 0$  mit einem  $k < n$  erfüllt, können wir immerhin  $k - 1$  Schritte der Induktion durchführen, um zu zeigen, dass die ersten  $k$  Spalten der Matrizen  $\mathbf{Q}$  und  $\widehat{\mathbf{Q}}$  sich nur in ihren Vorzeichen unterscheiden können.

Aus  $h_{k+1,k} = 0$  folgt mit (6.10a) die Gleichung

$$\mathbf{A}\widehat{\mathbf{q}}^{(k)} = \bar{\alpha}_k \mathbf{A}\mathbf{q}^{(k)} = \bar{\alpha}_k \sum_{j=1}^k h_{jk} \mathbf{q}^{(j)} = \bar{\alpha}_k \sum_{j=1}^k h_{jk} \alpha_j \widehat{\mathbf{q}}^{(j)} = \sum_{j=1}^k \hat{h}_{jk} \widehat{\mathbf{q}}^{(j)},$$

und daraus mit (6.10b) bereits  $\hat{h}_{k+1,k} = 0$ . Falls also in  $\mathbf{H}$  die Irreduzibilität verletzt ist, ist sie es auch in  $\widehat{\mathbf{H}}$ , und zwar in demselben Nebendiagonalelement.

Bisher haben wir Schritte der QR-Iteration durchgeführt, indem wir die QR-Zerlegung  $\widehat{\mathbf{Q}}^{(m+1)} \mathbf{R}^{(m+1)} = \mathbf{A}^{(m)}$  berechnet haben, um dann

$$\mathbf{A}^{(m+1)} = \mathbf{R}^{(m+1)} \widehat{\mathbf{Q}}^{(m+1)} = (\widehat{\mathbf{Q}}^{(m+1)})^* \mathbf{A}^{(m)} \widehat{\mathbf{Q}}^{(m+1)}$$

zu konstruieren. Praktisch haben wir  $\widehat{\mathbf{Q}}^{(m+1)}$  als Folge von  $n - 1$  Givens-Rotationen dargestellt (siehe (6.8)), also als

$$\widehat{\mathbf{Q}}^{(m+1)} = \mathbf{G}_1^* \mathbf{G}_2^* \dots \mathbf{G}_{n-1}^*,$$

bei denen  $\mathbf{G}_i$  beziehungsweise  $\mathbf{G}_i^*$  jeweils nur auf die  $i$ -te und die  $(i + 1)$ -te Komponente eines Vektors wirkt. Daraus folgt, dass die erste Spalte der Matrix  $\widehat{\mathbf{Q}}^{(m+1)}$  ausschließlich von  $\mathbf{G}_1$  abhängt, aber von keiner der anderen Rotationen.

Falls wir also eine zweite unitäre Transformation

$$\widetilde{\mathbf{Q}}^{(m+1)} = \mathbf{G}_1^* \widetilde{\mathbf{G}}_2^* \dots \widetilde{\mathbf{G}}_{n-1}^*$$

aus Givens-Rotationen  $\widetilde{\mathbf{G}}_i^*$  konstruieren, die auch jeweils nur die  $i$ -te und die  $(i + 1)$ -te Spalte beeinflussen, müssen beide dieselbe erste Spalte besitzen.

Falls  $\mathbf{A}^{(m+1)}$  irreduzibel ist und falls wir  $\widetilde{\mathbf{G}}_2^*, \dots, \widetilde{\mathbf{G}}_{n-1}^*$  so gewählt haben, dass

$$\widetilde{\mathbf{A}}^{(m+1)} := (\widetilde{\mathbf{Q}}^{(m+1)})^* \mathbf{A}^{(m)} \widetilde{\mathbf{Q}}^{(m+1)}$$

wieder eine Hessenberg-Matrix ist, erhalten wir mit Satz 6.16, dass sich beide Matrizen nur durch eine unitäre Diagonalskalierung unterscheiden.

Derartige Skalierungen sind für die Konvergenz der QR-Iteration ohne Belang, wir können also  $\widetilde{\mathbf{A}}^{(m+1)}$  anstelle von  $\mathbf{A}^{(m+1)}$  verwenden. Falls  $\mathbf{A}^{(m+1)}$  nicht irreduzibel sein sollte, folgt aus Bemerkung 6.17, dass wir auf einen Teil der Matrix eine nicht näher festgelegte unitäre Transformation angewendet haben. Da in diesem Fall aber ohnehin als nächstes eine Deflation durchgeführt und diese Teilmatrix separat weiter behandelt werden würde, schadet diese überflüssige Transformation nicht.

Der Vorteil dieses Zugangs besteht darin, dass wir  $\widetilde{\mathbf{A}}^{(m+1)}$  direkt aus  $\mathbf{A}^{(m)}$  konstruieren können, ohne eine Folge von Givens-Rotationen zu speichern: Zunächst wählen wir die

Givens-Rotation  $\mathbf{G}_1$  so, dass das Nebendiagonalelement  $a_{21}^{(m)}$  der geshifteten Matrix  $\mathbf{A}^{(m)} - \mu\mathbf{I}$  eliminiert wird:

$$\mathbf{G}_1 = \begin{pmatrix} \bar{c} & \bar{s} & & & & \\ -s & c & & & & \\ & & 1 & & & \\ & & & \ddots & & \\ & & & & \ddots & \\ & & & & & 1 \end{pmatrix}, \quad c = \frac{a_{11}^{(m)} - \mu}{r}, \quad s = \frac{a_{21}^{(m)}}{r}, \quad r = \sqrt{|a_{11}^{(m)} - \mu|^2 + |a_{21}^{(m)}|^2}.$$

Die durch  $\mathbf{G}_1$  definierte unitäre Ähnlichkeitstransformation wenden wir auf die erste und zweite Zeile und Spalte an:

$$\begin{aligned} \mathbf{A}^{(m)} \rightsquigarrow \mathbf{G}_1 \mathbf{A}^{(m)} &= \begin{pmatrix} \boxtimes & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ \boxtimes & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ & \times & \times & \times & \times & \dots \\ & & \times & \times & \times & \dots \\ & & & \times & \times & \dots \\ & & & & \ddots & \ddots \end{pmatrix} \\ \rightsquigarrow \mathbf{G}_1 \mathbf{A}^{(m)} \mathbf{G}_1^* &= \begin{pmatrix} \boxtimes & \boxtimes & \times & \times & \times & \dots \\ \boxtimes & \boxtimes & \times & \times & \times & \dots \\ \boxtimes & \boxtimes & \times & \times & \times & \dots \\ & & \times & \times & \times & \dots \\ & & & \times & \times & \dots \\ & & & & \ddots & \ddots \end{pmatrix} \end{aligned}$$

Da bei der Spaltentransformation (also der Multiplikation mit  $\mathbf{G}_1^*$  von rechts) die erste und zweite Spalte kombiniert werden, wird im Allgemeinen in der dritten Zeile der ersten Spalte der Matrix  $\mathbf{B} := \mathbf{G}_1 \mathbf{A} \mathbf{G}_1^*$  keine Null mehr stehen, die Matrix wird also nicht mehr in Hessenberg-Gestalt sein. Um diese Gestalt wieder herzustellen, eliminieren wir diesen Eintrag mit einer Givens-Rotation, die wir auf die zweite und dritte Zeile anwenden:

$$\tilde{\mathbf{G}}_2 = \begin{pmatrix} 1 & & & & & \\ & \bar{c} & \bar{s} & & & \\ & -s & c & & & \\ & & & 1 & & \\ & & & & \ddots & \\ & & & & & 1 \end{pmatrix}, \quad c = \frac{b_{21}}{r}, \quad s = \frac{b_{31}}{r}, \quad r = \sqrt{|b_{21}|^2 + |b_{31}|^2}.$$

Um eine Ähnlichkeitstransformation zu erhalten, müssen wir mit  $\tilde{\mathbf{G}}_2$  von links und mit  $\tilde{\mathbf{G}}_2^*$  von rechts multiplizieren, also die Rotation auf die Zeilen und ihre Adjungierte auf

## 6 Die QR-Iteration

die Spalten anwenden:

$$\mathbf{B} \rightsquigarrow \tilde{\mathbf{G}}_2 \mathbf{B} = \begin{pmatrix} \times & \times & \times & \times & \times & \dots \\ \boxtimes & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ \mathbf{0} & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ & & \times & \times & \times & \dots \\ & & & \times & \times & \dots \\ & & & & \ddots & \ddots \end{pmatrix} \rightsquigarrow \tilde{\mathbf{G}}_2 \mathbf{B} \tilde{\mathbf{G}}_2^* = \begin{pmatrix} \times & \boxtimes & \boxtimes & \times & \times & \dots \\ \times & \boxtimes & \boxtimes & \times & \times & \dots \\ & \boxtimes & \boxtimes & \times & \times & \dots \\ & & \boxtimes & \boxtimes & \times & \dots \\ & & & \times & \times & \dots \\ & & & & \ddots & \ddots \end{pmatrix}$$

Die Ähnlichkeitstransformation mit  $\tilde{\mathbf{G}}_2$  führt zwar dazu, dass der Eintrag  $b_{31}$  eliminiert wird, allerdings sorgt auch hier die Transformation der Spalten dafür, dass in dem Ergebnis  $\mathbf{C} := \tilde{\mathbf{G}}_2 \mathbf{B} \tilde{\mathbf{G}}_2^*$  ein von Null verschiedener Eintrag außerhalb der Hessenberg-Gestalt auftreten kann, diesmal in der vierten Zeile der zweiten Spalte.

Diesen Eintrag können wir mit einer weiteren Givens-Rotation  $\tilde{\mathbf{G}}_3$  eliminieren, die auf die dritte und vierte Zeile beziehungsweise Spalte wirkt:

$$\mathbf{C} \rightsquigarrow \tilde{\mathbf{G}}_3 \mathbf{C} = \begin{pmatrix} \times & \times & \times & \times & \times & \dots \\ \times & \times & \times & \times & \times & \dots \\ & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ \mathbf{0} & \boxtimes & \boxtimes & \boxtimes & \dots \\ & & \times & \times & \dots \\ & & & \ddots & \ddots \end{pmatrix} \rightsquigarrow \tilde{\mathbf{G}}_3 \mathbf{C} \tilde{\mathbf{G}}_3^* = \begin{pmatrix} \times & \times & \boxtimes & \boxtimes & \times & \dots \\ \times & \times & \boxtimes & \boxtimes & \times & \dots \\ & \times & \boxtimes & \boxtimes & \times & \dots \\ & & \boxtimes & \boxtimes & \times & \dots \\ & & & \boxtimes & \boxtimes & \times & \dots \\ & & & & \ddots & \ddots \end{pmatrix}$$

Nach diesem Schritt ist potentiell in der fünften Zeile der dritten Spalte ein von Null verschiedener Eintrag entstanden. Mit jeder Givens-Rotation verschiebt sich also das störende Element weiter nach rechts unten, bis es mit einer letzten Rotation  $\tilde{\mathbf{G}}_{n-1}$  endgültig eliminiert werden kann und die Matrix

$$\tilde{\mathbf{A}}^{(m+1)} = \tilde{\mathbf{G}}_{n-1} \dots \tilde{\mathbf{G}}_2 \mathbf{G}_1 \mathbf{A}^{(m)} \mathbf{G}_1^* \tilde{\mathbf{G}}_2^* \dots \tilde{\mathbf{G}}_{n-1}^*$$

wieder Hessenberg-Form hat. Gemäß Satz 6.16 unterscheidet sich diese Matrix von  $\mathbf{A}^{(m+1)}$  nur durch eine unitäre Diagonalskalierung, also können wir fortfahren, als wäre sie die neue Iterierte.

Bei der praktischen Implementierung bietet es sich an, den „überschüssigen“ Eintrag, der durch die erste Givens-Rotation entsteht, separat zu behandeln, schließlich ist denkbar, dass eine Darstellung der Matrix  $\mathbf{A}$  verwendet wird, bei der keine Speicherplatz für über die Hessenberg-Form hinausgehende Einträge vorgesehen ist. In dem folgenden Beispiel wird der kritische Eintrag in einer separaten Variablen  $\gamma$  aufbewahrt, die sich aus der jeweils vorangehenden Givens-Rotation ergibt und in die Berechnung der nächsten Givens-Rotation eingeht.

**Algorithmus 6.18** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine Hessenberg-Matrix. Der folgende Algorithmus überschreibt  $\mathbf{A}$  mit einer Matrix, die im Wesentlichen (siehe Bemerkung 6.17) mit der nächsten Iterierten der QR-Iteration mit Shift  $\mu$  übereinstimmt.



```

 $r \leftarrow \sqrt{|a_{11} - \mu|^2 + |a_{21}|^2}; \quad c \leftarrow (a_{11} - \mu)/r; \quad s \leftarrow a_{21}/r$ 
for  $j \in [1 : n]$  do begin
   $h \leftarrow a_{1j}; \quad a_{1j} \leftarrow \bar{c}h + \bar{s}a_{2j}; \quad a_{2j} \leftarrow -sh + ca_{2j}$ 
end
 $h \leftarrow a_{11}; \quad a_{11} \leftarrow ch + sa_{12}; \quad a_{12} \leftarrow -\bar{s}h + \bar{c}a_{12}$ 
 $h \leftarrow a_{21}; \quad a_{21} \leftarrow ch + sa_{22}; \quad a_{22} \leftarrow -\bar{s}h + \bar{c}a_{22}$ 
for  $i = 2$  to  $n - 1$  do begin
   $\gamma \leftarrow sa_{i+1,i}; \quad a_{i+1,i} \leftarrow \bar{c}a_{i+1,i}$ 
   $r \leftarrow \sqrt{|a_{i,i-1}|^2 + |\gamma|^2}; \quad c \leftarrow a_{i,i-1}/r; \quad s \leftarrow \gamma/r$ 
   $a_{i,i-1} \leftarrow r$ 
  for  $j \in [i : n]$  do begin
     $h \leftarrow a_{ij}; \quad a_{ij} \leftarrow \bar{c}h + \bar{s}a_{i+1,j}; \quad a_{i+1,j} \leftarrow -sh + ca_{i+1,j}$ 
  end
  for  $k \in [1 : i + 1]$  do begin
     $h \leftarrow a_{ki}; \quad a_{ki} \leftarrow ch + sa_{k,i+1}; \quad a_{k,i+1} \leftarrow -\bar{s}h + \bar{c}a_{k,i+1}$ 
  end
end
end

```

Neben dem geringeren Speicherbedarf besteht ein weiterer Vorteil des impliziten Verfahrens darin, dass es es erlaubt, *Multi-Shift-Strategien* elegant zu realisieren. Angenommen, wir führen einen Schritt der QR-Iteration mit einem Shift-Wert  $\mu_1$  durch und erhalten

$$\mathbf{A}^{(m+1)} = \mathbf{G}_{n-1} \dots \mathbf{G}_1 \mathbf{A}^{(m)} \mathbf{G}_1^* \dots \mathbf{G}_{n-1}^*$$

mit Givens-Rotationen  $\mathbf{G}_i$ , die jeweils auf die  $i$ -te und  $(i + 1)$ -Zeile wirken.

Nun führen wir einen weiteren Schritt mit einem zweiten Shift-Wert  $\mu_2$  durch und gelangen zu

$$\mathbf{A}^{(m+2)} = \widehat{\mathbf{G}}_{n-1} \dots \widehat{\mathbf{G}}_1 \mathbf{A}^{(m+1)} \widehat{\mathbf{G}}_1^* \dots \widehat{\mathbf{G}}_{n-1}^*$$

mit neuen Givens-Rotationen  $\widehat{\mathbf{G}}_i$ , die ebenfalls nur auf die  $i$ -te und  $(i + 1)$ -te Zeile wirken.

Wenn wir direkt von  $\mathbf{A}^{(m)}$  zu  $\mathbf{A}^{(m+2)}$  gelangen wollen, erhalten wir zunächst

$$\begin{aligned} \mathbf{A}^{(m+2)} &= \widehat{\mathbf{G}}_{n-1} \dots \widehat{\mathbf{G}}_1 \mathbf{A}^{(m+1)} \widehat{\mathbf{G}}_1^* \dots \widehat{\mathbf{G}}_{n-1}^* \\ &= \widehat{\mathbf{G}}_{n-1} \dots \widehat{\mathbf{G}}_1 \mathbf{G}_{n-1} \dots \mathbf{G}_1 \mathbf{A}^{(m)} \mathbf{G}_1^* \dots \mathbf{G}_{n-1}^* \widehat{\mathbf{G}}_1^* \dots \widehat{\mathbf{G}}_{n-1}^*. \end{aligned}$$

Da  $\widehat{\mathbf{G}}_1$  nur auf die erste und zweite Zeile wirkt, während die Rotationen  $\mathbf{G}_{n-1}, \dots, \mathbf{G}_3$  diese Zeilen nicht verändern, kommutieren die Matrizen, so dass wir

$$\mathbf{A}^{(m+2)} = \underbrace{\widehat{\mathbf{G}}_{n-1} \dots \widehat{\mathbf{G}}_2 \mathbf{G}_{n-1} \dots \mathbf{G}_3 \widehat{\mathbf{G}}_1 \mathbf{G}_2 \mathbf{G}_1}_{=: \mathbf{Q}_{12}^*} \mathbf{A}^{(m)} \underbrace{\mathbf{G}_1^* \mathbf{G}_2^* \widehat{\mathbf{G}}_1^* \mathbf{G}_3^* \dots \mathbf{G}_{n-1}^* \widehat{\mathbf{G}}_2^* \dots \widehat{\mathbf{G}}_{n-1}^*}_{=: \mathbf{Q}_{12}}$$

erhalten. Da nur  $\widehat{\mathbf{G}}_1 \mathbf{G}_2 \mathbf{G}_1$  die erste Spalte der Matrix  $\mathbf{Q}_{12}$  beeinflussen und sowohl  $\mathbf{A}^{(m)}$  als auch  $\mathbf{A}^{(m+2)}$  Hessenberg-Matrizen sind, können wir wieder Satz 6.16 einsetzen, um

## 6 Die QR-Iteration

den Schritt von  $\mathbf{A}^{(m)}$  zu  $\mathbf{A}^{(m+2)}$  implizit durchzuführen: Zunächst berechnen wir

$$\begin{aligned} \mathbf{A}^{(m)} \rightsquigarrow \mathbf{Q}_{12}^* \mathbf{A}^{(m)} &= \begin{pmatrix} \boxtimes & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ \boxtimes & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ \boxtimes & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ & & \times & \times & \times & \dots \\ & & & \times & \times & \dots \\ & & & & \ddots & \ddots \end{pmatrix} \\ \rightsquigarrow \mathbf{B} := \mathbf{Q}_{12}^* \mathbf{A}^{(m)} \mathbf{Q}_{12} &= \begin{pmatrix} \boxtimes & \boxtimes & \boxtimes & \times & \times & \dots \\ \boxtimes & \boxtimes & \boxtimes & \times & \times & \dots \\ \boxtimes & \boxtimes & \boxtimes & \times & \times & \dots \\ \boxtimes & \boxtimes & \boxtimes & \times & \times & \dots \\ & & & \times & \times & \dots \\ & & & & \ddots & \ddots \end{pmatrix}. \end{aligned}$$

Nun müssen wir nur noch unitäre Transformationen finden, die die Hessenberg-Form wiederherstellen, also die störenden Einträge in der ersten und zweiten Spalte eliminieren. Wir berechnen eine QR-Zerlegung  $(\widehat{\mathbf{Q}}, \widehat{\mathbf{R}})$  der Teilmatrix

$$\widehat{\mathbf{B}} := \begin{pmatrix} b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \\ b_{41} & b_{42} & b_{43} \end{pmatrix},$$

halten  $\widehat{\mathbf{Q}}^* \widehat{\mathbf{B}} = \widehat{\mathbf{R}}$  fest, und setzen  $\widehat{\mathbf{Q}}$  zu der unitären Matrix

$$\mathbf{Q}_2 := \begin{pmatrix} 1 & & \\ & \widehat{\mathbf{Q}}^* & \\ & & \mathbf{I} \end{pmatrix} \in \mathbb{K}^{n \times n}$$

fort, die die Anwendung von  $\widehat{\mathbf{Q}}$  auf die Zeilen zwei bis vier eines Vektors darstellt. Es folgt

$$\mathbf{B} \rightsquigarrow \mathbf{Q}_2 \mathbf{B} = \begin{pmatrix} \times & \times & \times & \times & \times & \dots \\ \boxtimes & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ \mathbf{0} & \boxtimes & \boxtimes & \boxtimes & \boxtimes & \dots \\ \mathbf{0} & \mathbf{0} & \boxtimes & \boxtimes & \boxtimes & \dots \\ & & & \times & \times & \dots \\ & & & & \ddots & \ddots \end{pmatrix} \rightsquigarrow \mathbf{Q}_2 \mathbf{B} \mathbf{Q}_2^* = \begin{pmatrix} \times & \boxtimes & \boxtimes & \boxtimes & \times & \dots \\ \times & \boxtimes & \boxtimes & \boxtimes & \times & \dots \\ & \boxtimes & \boxtimes & \boxtimes & \times & \dots \\ & \boxtimes & \boxtimes & \boxtimes & \times & \dots \\ & \boxtimes & \boxtimes & \boxtimes & \times & \dots \\ & & & & \ddots & \ddots \end{pmatrix},$$

nun treten also störende Einträge in den Spalten zwei und drei auf. Wie zuvor können wir sie mit einer weiteren QR-Zerlegung in die Spalten drei und vier verschieben und diesen Prozess wiederholen, bis sie aus der Matrix verschwunden sind und die Hessenberg-Form vorliegt. Mit Satz 6.16 folgt dann wieder, dass sich die resultierende Matrix nur in ihren Vorzeichen von  $\mathbf{A}^{(m+2)}$  unterscheiden kann. Wir können also mit *einem* Durchlauf durch die Matrix das Äquivalent von *zwei* QR-Schritten mit bei Bedarf sogar unterschiedlichen Shifts reproduzieren. Dieser Ansatz lässt sich auch auf  $k$  Schritte mit  $k$  Shifts erweitern.

## 6.5 Singulärwertzerlegung\*

Offensichtlich können nur quadratische Matrizen über Eigenwerte verfügen. Im Fall rechteckiger Matrizen kann gelegentlich die *Singulärwertzerlegung* an die Stelle der Schur-Zerlegung treten.

**Satz 6.19 (Singulärwertzerlegung)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times m}$ , sei  $k := \min\{n, m\}$ .

Dann existieren zwei isometrische Matrizen  $\mathbf{U} \in \mathbb{K}^{n \times k}$  und  $\mathbf{V} \in \mathbb{K}^{m \times k}$  und reelle Zahlen  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_k \geq 0$  mit

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*, \quad \mathbf{\Sigma} = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_k \end{pmatrix}. \quad (6.11)$$

Eine solche Faktorisierung nennt man eine Singulärwertzerlegung der Matrix  $\mathbf{A}$ . Die Zahlen  $\sigma_1, \dots, \sigma_k$  werden die Singulärwerte der Matrix  $\mathbf{A}$  genannt, die Spalten der Matrizen  $\mathbf{U}$  und  $\mathbf{V}$  linke und rechte Singulärvektoren.

*Beweis.* Wir führen den Beweis per Induktion über  $k \in \mathbb{N}$ .

*Induktionsanfang:* Sei  $\mathbf{A} \in \mathbb{K}^{n \times m}$  mit  $k = \min\{n, m\} = 1$  gegeben.

Falls  $\mathbf{A} = \mathbf{0}$  gilt, können wir  $\sigma_1 = 0$  verwenden und beliebige isometrische Matrizen  $\mathbf{U}$  und  $\mathbf{V}$  wählen.

Anderenfalls setzen wir  $\sigma_1 = \|\mathbf{A}\| > 0$ .

Falls  $m = 1$  gilt, setzen wir  $v_{11} = 1$  und  $\mathbf{U} = \mathbf{A}/\sigma_1$ . Dann ist  $\mathbf{U}$  ein Einheitsvektor, also isometrisch.

Anderenfalls gilt  $n = 1$  und wir setzen  $u_{11} = 1$  und  $\mathbf{V} = \mathbf{A}^*/\sigma_1$ . Laut Lemma 3.20 gilt  $\|\mathbf{A}\| = \|\mathbf{A}^*\|$ , also ist  $\mathbf{V}$  ein Einheitsvektor, und damit isometrisch.

*Induktionsvoraussetzung:* Sei  $k \in \mathbb{N}$  so gegeben, dass jede Matrix  $\mathbf{A} \in \mathbb{K}^{n \times m}$  mit  $k = \min\{n, m\}$  eine Singulärwertzerlegung der Form (6.11) besitzt.

*Induktionsschritt:* Sei  $\mathbf{A} \in \mathbb{K}^{n \times m}$  mit  $\min\{n, m\} = k + 1$  gegeben. Wir bezeichnen mit

$$\mathcal{S}_n := \{\mathbf{x} \in \mathbb{K}^n : \|\mathbf{x}\| = 1\}, \quad \mathcal{S}_m := \{\mathbf{y} \in \mathbb{K}^m : \|\mathbf{y}\| = 1\}$$

die  $n$ - und die  $m$ -dimensionale Einheitssphäre und halten fest, dass  $\mathcal{S}_n$  und  $\mathcal{S}_m$  nach dem Satz von Heine-Borel kompakte Mengen sind.

Die stetige Funktion

$$f: \mathcal{S}_n \times \mathcal{S}_m \rightarrow \mathbb{R}_{\geq 0}, \quad (\mathbf{x}, \mathbf{y}) \mapsto |\langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle|,$$

muss deshalb ein Maximum besitzen, das wir mit  $\sigma_1$  bezeichnen. Wir fixieren  $\mathbf{u} \in \mathcal{S}_n$  und  $\mathbf{v} \in \mathcal{S}_m$  mit  $\sigma_1 = f(\mathbf{u}, \mathbf{v}) = |\langle \mathbf{u}, \mathbf{A}\mathbf{v} \rangle|$ .

Aus der Cauchy-Schwarz-Ungleichung (3.7) folgt

$$|\langle \mathbf{u}, \mathbf{A}\mathbf{v} \rangle| \leq \|\mathbf{u}\| \|\mathbf{A}\mathbf{v}\|, \quad (6.12)$$

und Gleichheit gilt genau dann, wenn  $\mathbf{u}$  und  $\mathbf{A}\mathbf{v}$  linear abhängig sind. Da

$$\sigma_1 = \max\{\max\{|\langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle| : \mathbf{y} \in \mathcal{S}_m\} : \mathbf{x} \in \mathcal{S}_n\}$$

## 6 Die QR-Iteration

gilt, muss in (6.12) Gleichheit gelten, also müssen  $\mathbf{u}$  und  $\mathbf{A}\mathbf{v}$  linear abhängig sein. Dann finden wir ein  $\alpha \in \mathbb{K}$  mit

$$\begin{aligned}\mathbf{A}\mathbf{v} &= \alpha\mathbf{u}, \\ |\alpha| &= |\alpha| |\langle \mathbf{u}, \mathbf{u} \rangle| = |\langle \mathbf{u}, \alpha\mathbf{u} \rangle| = |\langle \mathbf{u}, \mathbf{A}\mathbf{v} \rangle| = \sigma_1.\end{aligned}$$

Indem wir das Vorzeichen von  $\mathbf{u}$  anpassen können wir

$$\mathbf{A}\mathbf{v} = \sigma_1\mathbf{u}$$

sicher stellen. Analog folgt aus Lemma 3.17 und der Cauchy-Schwarz-Ungleichung auch

$$|\langle \mathbf{u}, \mathbf{A}\mathbf{v} \rangle| = |\langle \mathbf{A}^*\mathbf{u}, \mathbf{v} \rangle| \leq \|\mathbf{A}^*\mathbf{u}\| \|\mathbf{v}\|,$$

und wir können wie zuvor folgern, dass  $\mathbf{A}^*\mathbf{u}$  und  $\mathbf{v}$  linear abhängig sind, dass also ein  $\beta \in \mathbb{K}$  mit

$$\mathbf{A}^*\mathbf{u} = \beta\mathbf{v}$$

existiert. Mit unserer modifizierten Wahl des Vektors  $\mathbf{u}$  folgt

$$\beta = \beta \langle \mathbf{v}, \mathbf{v} \rangle = \langle \mathbf{v}, \beta\mathbf{v} \rangle = \langle \mathbf{v}, \mathbf{A}^*\mathbf{u} \rangle = \langle \mathbf{A}\mathbf{v}, \mathbf{u} \rangle = \langle \sigma_1\mathbf{u}, \mathbf{u} \rangle = \bar{\sigma}_1 = \sigma_1,$$

also haben wir insgesamt

$$\mathbf{A}\mathbf{v} = \sigma_1\mathbf{u}, \quad \mathbf{A}^*\mathbf{u} = \sigma_1\mathbf{v}$$

bewiesen. Wir wählen Householder-Spiegelungen  $\widehat{\mathbf{U}}_1 \in \mathbb{K}^{n \times n}$  und  $\widehat{\mathbf{V}}_1 \in \mathbb{K}^{m \times m}$  derart, dass

$$\nu\mathbf{u} = \widehat{\mathbf{U}}_1\delta^{(1)}, \quad \mu\mathbf{v} = \widehat{\mathbf{V}}_1\delta^{(1)}$$

mit geeigneten  $\nu, \mu \in \mathbb{K}$  gelten. Wegen  $|\nu| = \|\nu\mathbf{u}\| = \|\delta^{(1)}\| = 1$  und  $|\mu| = \|\mu\mathbf{v}\| = \|\delta^{(1)}\| = 1$  sind auch  $\mathbf{U}_1 := \bar{\nu}\widehat{\mathbf{U}}_1$  und  $\mathbf{V}_1 := \bar{\mu}\widehat{\mathbf{V}}_1$  unitäre Matrizen mit

$$\mathbf{u} = \mathbf{U}_1\delta^{(1)}, \quad \mathbf{v} = \mathbf{V}_1\delta^{(1)}$$

Es folgen

$$\begin{aligned}\mathbf{U}_1^*\mathbf{A}\mathbf{V}_1\delta^{(1)} &= \mathbf{U}_1^*\mathbf{A}\mathbf{v} = \sigma_1\mathbf{U}_1^*\mathbf{u} = \sigma_1\delta^{(1)}, \\ (\mathbf{U}_1^*\mathbf{A}\mathbf{V}_1)^*\delta^{(1)} &= \mathbf{V}_1^*\mathbf{A}^*\mathbf{U}_1\delta^{(1)} = \mathbf{V}_1^*\mathbf{A}^*\mathbf{u} = \sigma_1\mathbf{V}_1^*\mathbf{v} = \sigma_1\delta^{(1)},\end{aligned}$$

so dass die transformierte Matrix die Gestalt

$$\mathbf{U}_1^*\mathbf{A}\mathbf{V}_1 = \begin{pmatrix} \sigma_1 & \\ & \widehat{\mathbf{A}} \end{pmatrix}, \quad \widehat{\mathbf{A}} \in \mathbb{K}^{(n-1) \times (m-1)}$$

aufweist. Nach der Induktionsvoraussetzung finden wir eine Singulärwertzerlegung der Matrix  $\hat{\mathbf{A}}$ , wir finden also isometrische Matrizen  $\hat{\mathbf{U}} \in \mathbb{K}^{(n-1) \times k}$  und  $\hat{\mathbf{V}} \in \mathbb{K}^{(m-1) \times k}$  sowie  $\sigma_2 \geq \sigma_3 \geq \dots \geq \sigma_{k+1}$  mit

$$\hat{\mathbf{A}} = \hat{\mathbf{U}} \hat{\mathbf{\Sigma}} \hat{\mathbf{V}}^*, \quad \hat{\mathbf{\Sigma}} = \begin{pmatrix} \sigma_2 & & \\ & \ddots & \\ & & \sigma_{k+1} \end{pmatrix}.$$

Zusammengesetzt erhalten wir

$$\mathbf{A} = \mathbf{U}_1 \begin{pmatrix} \sigma_1 & \\ & \hat{\mathbf{A}} \end{pmatrix} \mathbf{V}_1^* = \mathbf{U}_1 \begin{pmatrix} \sigma_1 & \\ & \hat{\mathbf{U}} \hat{\mathbf{\Sigma}} \hat{\mathbf{V}}^* \end{pmatrix} \mathbf{V}_1^* = \mathbf{U}_1 \begin{pmatrix} 1 & \\ & \hat{\mathbf{U}} \end{pmatrix} \begin{pmatrix} \sigma_1 & \\ & \hat{\mathbf{\Sigma}} \end{pmatrix} \begin{pmatrix} 1 & \\ & \hat{\mathbf{V}} \end{pmatrix}^* \mathbf{V}_1^*,$$

so dass wir mit

$$\mathbf{U} := \mathbf{U}_1 \begin{pmatrix} 1 & \\ & \hat{\mathbf{U}} \end{pmatrix}, \quad \mathbf{V} := \mathbf{V}_1 \begin{pmatrix} 1 & \\ & \hat{\mathbf{V}} \end{pmatrix}, \quad \mathbf{\Sigma} := \begin{pmatrix} \sigma_1 & \\ & \hat{\mathbf{\Sigma}} \end{pmatrix}$$

die gewünschte Zerlegung gefunden haben. Mit  $\mathbf{x} := \mathbf{U} \delta^{(2)}$  und  $\mathbf{y} := \mathbf{V} \delta^{(2)}$  erhalten wir

$$\sigma_2 = |\sigma_2| = |\langle \delta^{(2)}, \mathbf{\Sigma} \delta^{(2)} \rangle| = |\langle \mathbf{U}^* \mathbf{x}, \mathbf{\Sigma} \mathbf{V}^* \mathbf{y} \rangle| = |\langle \mathbf{x}, \mathbf{U} \mathbf{\Sigma} \mathbf{V}^* \mathbf{y} \rangle| = |\langle \mathbf{x}, \mathbf{A} \mathbf{y} \rangle| \leq \sigma_1,$$

so dass die Singulärwerte auch bereits die gewünschte Reihenfolge aufweisen.  $\blacksquare$

Eine solche Singulärwertzerlegung ist beispielsweise sehr nützlich, um lineare Ausgleichsprobleme zu lösen (vgl. Übungsaufgabe 6.20) oder approximative Faktorisierungen der Matrix  $\mathbf{A}$  (vgl. Übungsaufgabe 6.21) zu konstruieren.

Die Singulärwertzerlegung ist eng verwandt mit der Schur-Zerlegung: Falls

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*$$

mit isometrischen Matrizen  $\mathbf{U}$  und  $\mathbf{V}$  gilt, haben wir

$$\mathbf{A} \mathbf{A}^* = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^* \mathbf{V} \mathbf{\Sigma} \mathbf{U}^* = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^*, \quad \mathbf{A}^* \mathbf{A} = \mathbf{V} \mathbf{\Sigma} \mathbf{U}^* \mathbf{U} \mathbf{\Sigma} \mathbf{V}^* = \mathbf{V} \mathbf{\Sigma}^2 \mathbf{V}^*,$$

die linken und rechten Singulärvektoren sind also Eigenvektoren der positiv semidefiniten und selbstadjungierten Matrizen  $\mathbf{A} \mathbf{A}^*$  und  $\mathbf{A}^* \mathbf{A}$ .

Aus dieser Beobachtung lässt sich ein Algorithmus für die iterative Berechnung einer Singulärwertzerlegung gewinnen: Wir führen die QR-Iteration für die *Gramsche Matrix*  $\mathbf{G} := \mathbf{A} \mathbf{A}^*$  durch und stellen fest, dass jedes dabei auftretende Zwischenergebnis

$$\mathbf{G}^{(m)} = (\mathbf{Q}^{(m)})^* \mathbf{G} \mathbf{Q}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0$$

wieder in faktorisierter Form vorliegt: Es gilt

$$\mathbf{G}^{(m)} = (\mathbf{Q}^{(m)})^* \mathbf{G} \mathbf{Q}^{(m)} = (\mathbf{Q}^{(m)})^* \mathbf{A} ((\mathbf{Q}^{(m)})^* \mathbf{A})^* = \mathbf{A}^{(m)} (\mathbf{A}^{(m)})^* \quad \text{für alle } m \in \mathbb{N}_0$$

## 6 Die QR-Iteration

mit  $\mathbf{A}^{(m)} = (\mathbf{Q}^{(m)})^* \mathbf{A}$ . Um einen Schritt der QR-Iteration für die Matrix  $\mathbf{G}$  auszuführen, müssen wir also lediglich die korrespondierenden Transformationen auf die Zeilen der Matrix  $\mathbf{A}$  anwenden.

Wenn die QR-Iteration konvergiert ist, wenn also (bis auf eine gegebene Genauigkeit) eine Diagonalmatrix  $\mathbf{G}^{(m)}$  erreicht wurde, können wir eine Singulärwertzerlegung rekonstruieren: Da  $\mathbf{G}^{(m)} = \mathbf{A}^{(m)}(\mathbf{A}^{(m)})$  diagonal ist, stehen die Zeilen der Matrix  $\mathbf{A}^{(m)}$  senkrecht aufeinander. Wir berechnen eine LQ-Zerlegung

$$\mathbf{A}^{(m)} = \mathbf{L}^{(m)} \mathbf{P}^{(m)}$$

mit einer unteren Dreiecksmatrix  $\mathbf{L}^{(m)} \in \mathbb{K}^{n \times m}$  und einer unitären Matrix  $\mathbf{P}^{(m)} \in \mathbb{K}^{m \times m}$  und stellen fest, dass wegen Lemma 3.35 auch

$$\mathbf{G}^{(m)} = \mathbf{A}^{(m)}(\mathbf{A}^{(m)})^* = \mathbf{L}^{(m)} \mathbf{P}^{(m)} (\mathbf{P}^{(m)})^* (\mathbf{L}^{(m)})^* = \mathbf{L}^{(m)} (\mathbf{L}^{(m)})^*$$

gilt, auch die Zeilen der Matrix  $\mathbf{L}^{(m)}$  stehen also senkrecht aufeinander. Bei einer Dreiecksmatrix bedeutet das aber, dass die Matrix diagonal ist, wir haben also unitäre Transformationen gefunden, die  $\mathbf{A}^{(m)}$  auf Diagonalgestalt bringen. Indem wir die Vorzeichen anpassen und die Diagonalelemente sortieren, erhalten wir die gewünschte Singulärwertzerlegung.

In dieser allgemeinen Form wäre die Berechnung allerdings zu aufwendig. Analog zu der für die Effizienz der QR-Iteration entscheidenden Hessenberg-Transformation empfiehlt es sich deshalb, in einem ersten Schritt die Matrix  $\mathbf{A}$  zu *bidiagonalisieren*, sie also mit unitären Transformationen  $\mathbf{U}^{(0)} \in \mathbb{K}^{n \times k}$  und  $\mathbf{V}^{(0)} \in \mathbb{K}^{m \times k}$  in die Form

$$(\mathbf{U}^{(0)})^* \mathbf{A} \mathbf{V}^{(0)} = \begin{pmatrix} \alpha_1 & & & & \\ \beta_1 & \alpha_2 & & & \\ & \ddots & \ddots & & \\ & & & \beta_{k-1} & \alpha_k \end{pmatrix}$$

zu bringen. Das gelingt dem *Golub-Kahan-Bidiagonalisierungsalgorithmus*<sup>1</sup>, mit Hilfe von Householder-Spiegelungen: Zunächst wenden wir eine Spiegelung auf die Spalten der Matrix an, um die erste Zeile in die gewünschte Form zu bringen:

$$\mathbf{A} \mathbf{H}_1 = \begin{pmatrix} \alpha_1 & \mathbf{0} & \mathbf{0} & \dots & \mathbf{0} \\ \boxtimes & \boxtimes & \boxtimes & \dots & \boxtimes \\ \boxtimes & \boxtimes & \boxtimes & \dots & \boxtimes \\ \vdots & \vdots & \vdots & \ddots & \times \\ \boxtimes & \boxtimes & \boxtimes & \dots & \boxtimes \end{pmatrix}.$$

In einem zweiten Schritt wenden wir eine Spiegelung auf die Zeilen der Matrix an, die

<sup>1</sup>G. Golub und W. Kahan, *Calculating the singular values and pseudo-inverse of a matrix*, J. SIAM Num. Anal. B 2(2), 205–224 (1965)

die erste Zeile unverändert lässt und Nullen im Rest der ersten Spalte einführt:

$$\mathbf{H}_2^* \mathbf{A} \mathbf{H}_1 = \begin{pmatrix} \alpha_1 & 0 & 0 & \dots & 0 \\ \beta_1 & \boxtimes & \boxtimes & \dots & \boxtimes \\ \mathbf{0} & \boxtimes & \boxtimes & \dots & \boxtimes \\ \mathbf{0} & \boxtimes & \boxtimes & \dots & \boxtimes \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \boxtimes & \boxtimes & \dots & \boxtimes \end{pmatrix}$$

Im nächsten Schritt wird wieder eine Spiegelung auf die Spalten angewendet, die die erste Spalte unverändert lässt und die zweite Zeile in die gewünschte Form bringt:

$$\mathbf{H}_2^* \mathbf{A} \mathbf{H}_1 \mathbf{H}_3 = \begin{pmatrix} \alpha_1 & 0 & 0 & \dots & 0 \\ \beta_1 & \alpha_2 & \mathbf{0} & \dots & \mathbf{0} \\ 0 & \boxtimes & \boxtimes & \dots & \boxtimes \\ 0 & \boxtimes & \boxtimes & \dots & \boxtimes \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \boxtimes & \boxtimes & \dots & \boxtimes \end{pmatrix}.$$

Nun ist es wieder Zeit für eine Zeilentransformation, bei der die ersten *beiden* Zeilen nicht angefasst und Nulleinträge in der zweite Spalte hergestellt werden.

$$\mathbf{H}_4^* \mathbf{H}_2^* \mathbf{A} \mathbf{H}_1 \mathbf{H}_3 = \begin{pmatrix} \alpha_1 & 0 & 0 & \dots & 0 \\ \beta_1 & \alpha_2 & \mathbf{0} & \dots & \mathbf{0} \\ 0 & \beta_2 & \boxtimes & \dots & \boxtimes \\ 0 & \mathbf{0} & \boxtimes & \dots & \boxtimes \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & \mathbf{0} & \boxtimes & \dots & \boxtimes \end{pmatrix}.$$

In dieser Weise können wir fortfahren, bis die gewünschte Bidiagonalform erreicht ist.

Wenn die Matrix  $\mathbf{A}^{(0)} := (\mathbf{U}^{(0)})^* \mathbf{A} \mathbf{V}^{(0)}$  in Bidiagonalgestalt gegeben ist, ist das korrespondierende Produkt  $\mathbf{G}^{(0)} = \mathbf{A}^{(0)} (\mathbf{A}^{(0)})^*$  eine Tridiagonalmatrix der Form

$$\mathbf{G}^{(0)} = \begin{pmatrix} |\alpha_1|^2 & \alpha_1 \bar{\beta}_1 & & & & \\ \bar{\alpha}_1 \beta_1 & |\alpha_2|^2 + |\beta_1|^2 & \alpha_2 \bar{\beta}_2 & & & \\ & \ddots & \ddots & \ddots & & \\ & & \bar{\alpha}_{k-2} \beta_{k-2} & |\alpha_{k-1}|^2 + |\beta_{k-2}|^2 & \alpha_{k-1} \bar{\beta}_{k-1} & \\ & & & \bar{\alpha}_{k-1} \beta_{k-1} & |\alpha_k|^2 + |\beta_{k-1}|^2 & \end{pmatrix}, \quad (6.13)$$

so dass sich ein Schritt der QR-Iteration besonders effizient durchführen lässt. Sehr elegant wird der Algorithmus, wenn wir eine *implizite* QR-Iteration verwenden: Wir bestimmen einen geeigneten Shift-Wert und wenden die erste Givens-Rotation auf die

## 6 Die QR-Iteration

ersten beiden Zeilen der Matrix  $\mathbf{A}^{(m)}$  an:

$$\mathbf{A}^{(m)} = \begin{pmatrix} \times & & & & \\ \times & \times & & & \\ & \times & \times & & \\ & & \times & \times & \\ & & & \ddots & \ddots \end{pmatrix} \rightsquigarrow \mathbf{G}_1 \mathbf{A}^{(m)} = \begin{pmatrix} \boxtimes & \boxtimes & & & \\ \boxtimes & \boxtimes & & & \\ & \times & \times & & \\ & & \times & \times & \\ & & & \ddots & \ddots \end{pmatrix}$$

Dadurch wird die Bidiagonalstruktur verletzt, in der ersten Zeile entsteht ein Eintrag oberhalb der Diagonalen. Diesen Eintrag können wir mit einer auf die erste und zweite Spalte angewandten Givens-Rotation eliminieren, erhalten dabei aber einen unerwünschten Eintrag in der dritten Zeile:

$$\mathbf{G}_1 \mathbf{A}^{(m)} = \begin{pmatrix} \times & \times & & & \\ \times & \times & & & \\ & \times & \times & & \\ & & \times & \times & \\ & & & \ddots & \ddots \end{pmatrix} \rightsquigarrow \mathbf{G}_1 \mathbf{A}^{(m)} \mathbf{G}_2^* = \begin{pmatrix} \boxtimes & \mathbf{0} & & & \\ \boxtimes & \boxtimes & & & \\ \boxtimes & \boxtimes & \times & & \\ & & \times & \times & \\ & & & \ddots & \ddots \end{pmatrix}$$

Eine Givens-Rotation der zweiten und dritten Zeile beseitigt ihn, erzeugt aber einen Eintrag oberhalb der Diagonalen in der dritten Spalte:

$$\mathbf{G}_1 \mathbf{A}^{(m)} \mathbf{G}_2^* = \begin{pmatrix} \times & & & & \\ \times & \times & & & \\ \times & \times & \times & & \\ & & \times & \times & \\ & & & \ddots & \ddots \end{pmatrix} \rightsquigarrow \mathbf{G}_3 \mathbf{G}_1 \mathbf{A}^{(m)} \mathbf{G}_2^* = \begin{pmatrix} \times & & & & \\ \boxtimes & \boxtimes & \boxtimes & & \\ \mathbf{0} & \boxtimes & \boxtimes & & \\ & & \times & \times & \\ & & & \ddots & \ddots \end{pmatrix}$$

Diesen Eintrag können wir mit einer auf die zweite und dritte Spalte angewandten Givens-Rotation eliminieren, die aber wieder zu einem von null verschiedenen Eintrag in der vierten Zeile führt:

$$\mathbf{G}_3 \mathbf{G}_1 \mathbf{A}^{(m)} \mathbf{G}_2^* = \begin{pmatrix} \times & & & & \\ \times & \times & \times & & \\ & \times & \times & & \\ & & \times & \times & \\ & & & \ddots & \ddots \end{pmatrix} \rightsquigarrow \mathbf{G}_3 \mathbf{G}_1 \mathbf{A}^{(m)} \mathbf{G}_2^* \mathbf{G}_4^* = \begin{pmatrix} \times & & & & \\ \times & \boxtimes & \mathbf{0} & & \\ & \boxtimes & \boxtimes & & \\ & \boxtimes & \boxtimes & \times & \\ & & & \ddots & \ddots \end{pmatrix}$$

In dieser Weise können wir die problematischen von null verschiedenen Einträge „aus der Matrix heraus schieben“, wie wir es schon im Fall des impliziten QR-Verfahrens getan haben.

Damit liegen zwei unitäre Transformationen der Matrix  $\mathbf{G}^{(m)}$  auf Hessenberg-Gestalt vor, so dass wir Satz 6.16 anwenden könnten, um zu beweisen, dass sie bis auf Vorzeichen identisch sind. Dafür wäre es aber erforderlich, zu zeigen, dass die ersten Spalten der Transformationen identisch sind. Glücklicherweise brauchen wir dabei nur





## 6 Die QR-Iteration

Die Matrix  $\mathbf{A}^+ := \mathbf{V}\boldsymbol{\Sigma}^+\mathbf{U}^*$  nennen wir die Pseudoinverse der Matrix  $\mathbf{A}$ .

- (a) Beweisen Sie  $\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$ ,  $\mathbf{A}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+$  und dass  $\mathbf{A}^+\mathbf{A}$  und  $\mathbf{A}\mathbf{A}^+$  selbstadjungiert sind.
- (b) Sei  $\mathbf{b} \in \mathbb{K}^n$ . Mit

$$\mathcal{S} := \{\mathbf{x} \in \mathbb{K}^m : \|\mathbf{A}\mathbf{x} - \mathbf{b}\| \text{ ist minimal}\}$$

bezeichnen wir die Menge der Lösungen des zu  $\mathbf{A}$  und  $\mathbf{b}$  gehörenden linearen Ausgleichsproblems.

Beweisen Sie, dass  $\mathbf{A}^+\mathbf{b} \in \mathcal{S}$  und

$$\|\mathbf{A}^+\mathbf{b}\| = \min\{\|\mathbf{y}\| : \mathbf{y} \in \mathcal{S}\}$$

gelten. Aufgrund dieser Eigenschaft nennt man  $\mathbf{A}^+\mathbf{b}$  die Minimumnormlösung des linearen Ausgleichsproblems.

**Übungsaufgabe 6.21 (Niedrigrangapproximation)** Seien  $n, m \in \mathbb{N}$  gegeben, sei  $k := \min\{n, m\}$ . Sei  $\mathbf{A} \in \mathbb{K}^{n \times m}$ . Sei  $\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^* = \mathbf{A}$  eine Singulärwertzerlegung der Matrix mit den Singulärwerten  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_k \geq 0$ .

Für alle  $\ell \in [0 : k - 1]$  definieren wir

$$\boldsymbol{\Sigma}_\ell := \begin{pmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_\ell & \\ & & & \mathbf{0} \end{pmatrix} \in \mathbb{K}^{k \times k}$$

und  $\mathbf{A}_\ell := \mathbf{U}\boldsymbol{\Sigma}_\ell\mathbf{V}^*$ .

- (a) Beweisen Sie  $\|\mathbf{A} - \mathbf{A}_\ell\| = \sigma_{\ell+1}$  für alle  $\ell \in [0 : k - 1]$ .
- (b) Sei  $\ell \in [0 : k - 1]$ , und sei  $\mathbf{R} \in \mathbb{K}^{n \times m}$  eine beliebige Matrix mit  $\text{Rang } \ell$ , also mit  $\dim(\text{Bild } \mathbf{R}) = \ell$ .

Beweisen Sie, dass ein Vektor  $\mathbf{z} \in \mathbb{K}^m$  existiert mit

$$\mathbf{R}\mathbf{z} = \mathbf{0}, \quad \|\mathbf{A}\mathbf{z}\| \geq \sigma_{\ell+1}, \quad \|\mathbf{z}\| = 1.$$

- (c) Folgern Sie aus Teil (b), dass für alle Matrizen mit  $\text{Rang } \ell \in [0 : k - 1]$  die Ungleichung  $\|\mathbf{A} - \mathbf{R}\| \geq \sigma_{\ell+1}$  gilt.

Unter allen Approximationen der Matrix  $\mathbf{A}$  mit Rang höchstens  $\ell$  erreicht also  $\mathbf{A}_\ell$  den minimalen Fehler in der Spektralnrm.

## 7 Verfahren für Tridiagonalmatrizen

Im vorigen Kapitel haben wir gesehen, dass sich eine beliebige selbstadjungierte Matrix mit Hilfe von Householder-Spiegelungen in eine Tridiagonalmatrix überführen lässt. Für Tridiagonalmatrizen lassen sich nicht nur die verschiedenen Varianten der QR-Iteration effizient durchführen, es existieren auch eine Reihe weiterer Verfahren, die die Tridiagonalgestalt ausnutzen können.

Ein Beispiel sind Bisektionsverfahren, mit deren Hilfe sich nach Nullstellen des charakteristischen Polynoms  $p_A$  suchen lässt. Nach Lemma 3.5 sind diese Nullstellen gerade die Eigenwerte der Matrix  $\mathbf{A}$ . Ein besonderer Vorteil dieser Verfahren besteht darin, dass wir mit ihrer Hilfe nicht nur nach dem kleinsten oder größten Eigenwert suchen können, sondern auch nach dem  $k$ -ten.

### 7.1 Auswertung des charakteristischen Polynoms

Wir fixieren eine selbstadjungierte Tridiagonalmatrix

$$\mathbf{A} = \begin{pmatrix} \alpha_1 & \bar{\beta}_1 & & & \\ \beta_1 & \alpha_2 & \bar{\beta}_2 & & \\ & \ddots & \ddots & \ddots & \\ & & \beta_{n-2} & \alpha_{n-1} & \bar{\beta}_{n-1} \\ & & & \beta_{n-1} & \alpha_n \end{pmatrix}. \quad (7.1)$$

Für ein  $i \in [1 : n]$  und ein  $x \in \mathbb{R}$  bezeichnen wir die Determinante der  $i$ -ten Haupttermatrix von  $x\mathbf{I} - \mathbf{A}$  mit

$$p_i(x) := \det \begin{pmatrix} x - \alpha_1 & -\bar{\beta}_1 & & & \\ -\beta_1 & x - \alpha_2 & -\bar{\beta}_2 & & \\ & \ddots & \ddots & \ddots & \\ & & -\beta_{i-2} & x - \alpha_{i-1} & -\bar{\beta}_{i-1} \\ & & & -\beta_{i-1} & x - \alpha_i \end{pmatrix}.$$

Offenbar ist  $p_n$  gerade das charakteristische Polynom  $p_A$ , während  $p_1(x) = x - \alpha_1$  sehr einfach zu berechnen ist. Unser Ziel ist es, mit möglichst geringem Rechenaufwand von  $p_1$  zu  $p_n$  zu gelangen. Dazu verwenden wir den Laplaceschen Entwicklungssatz für Determinanten: Für  $i > 2$  entwickeln wir die Determinante erst nach der letzten Spalte und

## 7 Verfahren für Tridiagonalmatrizen

anschließend nach der letzten Zeile, um den Zusammenhang

$$\begin{aligned}
 p_i(x) &:= \det \left( \begin{array}{cccc|c}
 x - \alpha_1 & -\bar{\beta}_1 & & & \\
 -\beta_1 & x - \alpha_2 & -\bar{\beta}_2 & & \\
 & \ddots & \ddots & \ddots & \\
 & & -\beta_{i-3} & x - \alpha_{i-2} & -\bar{\beta}_{i-2} \\
 & & & -\beta_{i-2} & x - \alpha_{i-1} & -\bar{\beta}_{i-1} \\
 & & & & -\beta_{i-1} & x - \alpha_i
 \end{array} \right) \\
 &= (x - \alpha_i) \det \left( \begin{array}{cccc}
 x - \alpha_1 & -\bar{\beta}_1 & & \\
 -\bar{\beta}_1 & x - \alpha_2 & -\bar{\beta}_2 & \\
 & \ddots & \ddots & \ddots \\
 & & -\beta_{i-3} & x - \alpha_{i-2} & -\bar{\beta}_{i-2} \\
 & & & -\beta_{i-2} & x - \alpha_{i-1}
 \end{array} \right) \\
 &\quad + \bar{\beta}_{i-1} \det \left( \begin{array}{cccc}
 x - \alpha_1 & -\bar{\beta}_1 & & \\
 -\beta_1 & x - \alpha_2 & -\bar{\beta}_2 & \\
 & \ddots & \ddots & \ddots \\
 & & -\beta_{i-3} & x - \alpha_{i-2} & -\bar{\beta}_{i-2} \\
 & & & & -\beta_{i-1}
 \end{array} \right) \\
 &= (x - \alpha_i) \det \left( \begin{array}{cccc}
 x - \alpha_1 & -\bar{\beta}_1 & & \\
 -\beta_1 & x - \alpha_2 & -\bar{\beta}_2 & \\
 & \ddots & \ddots & \ddots \\
 & & -\beta_{i-2} & x - \alpha_{i-2} & -\bar{\beta}_{i-2} \\
 & & & -\beta_{i-2} & x - \alpha_{i-1}
 \end{array} \right) \\
 &\quad - |\beta_{i-1}|^2 \det \left( \begin{array}{cccc}
 x - \alpha_1 & -\bar{\beta}_1 & & \\
 -\beta_1 & x - \alpha_2 & -\bar{\beta}_2 & \\
 & \ddots & \ddots & \ddots \\
 & & -\beta_{i-3} & x - \alpha_{i-3} & -\bar{\beta}_{i-3} \\
 & & & -\beta_{i-3} & x - \alpha_{i-2}
 \end{array} \right) \\
 &= (x - \alpha_i)p_{i-1}(x) - |\beta_{i-1}|^2 p_{i-2}(x)
 \end{aligned}$$

zu erhalten. Wenn wir zur Vereinfachung  $p_0 = 1$  einführen, ergibt sich wegen

$$p_2(x) = \det \begin{pmatrix} x - \alpha_1 & -\bar{\beta}_1 \\ -\beta_1 & x - \alpha_2 \end{pmatrix} = (x - \alpha_1)(x - \alpha_2) - |\beta_1|^2 \quad \text{für alle } x \in \mathbb{R}$$

die Rekursionsformel

$$p_i(x) = \begin{cases} 1 & \text{falls } i = 0, \\ x - \alpha_1 & \text{falls } i = 1, \\ (x - \alpha_i)p_{i-1}(x) - |\beta_{i-1}|^2 p_{i-2}(x) & \text{ansonsten} \end{cases} \quad \text{für alle } i \in [0 : n], x \in \mathbb{R}, \quad (7.2)$$

an der wir unmittelbar den folgenden Algorithmus für die Berechnung des Tupels  $(p_n(x), \dots, p_0(x))$ , also insbesondere auch die Auswertung des charakteristischen Polynoms, ablesen können.

**Algorithmus 7.1 (Auswertung von  $p_i$ )** *Der folgende Algorithmus berechnet das Tupel  $\mathbf{p} = (p_0(x), \dots, p_n(x))$  für ein beliebiges  $x \in \mathbb{R}$ :*

```

 $p_0 \leftarrow 1; \quad p_1 \leftarrow x - \alpha_1$ 
for  $i = 2$  to  $n$  do
   $p_i \leftarrow (x - \alpha_i)p_{i-1} - |\beta_{i-1}|^2 p_{i-2}$ 

```

Insbesondere lässt sich mit Hilfe des Algorithmus 7.1 das charakteristische Polynom  $p_A = p_n$  in  $\mathcal{O}(n)$  Operationen auswerten, und nebenbei werden die charakteristischen Polynome aller Hauptuntermatrizen berechnet, die sich für bestimmte Algorithmen als sehr nützlich erweisen können.

Da Eigenwerte von  $\mathbf{A}$  Nullstellen des charakteristischen Polynoms  $p_A$  sind, das wir mit Algorithmus 7.1 elegant und effizient auswerten können, bietet es sich an, Standardverfahren zur Nullstellenberechnung auf  $p_A$  anzuwenden.

Eines der einfachsten und trotzdem zuverlässigsten Verfahren ist die Bisektion, bei der eine Folge von Intervallen berechnet wird, die jeweils mindestens eine Nullstelle enthalten: Wir beginnen mit einem Intervall  $[a, b]$  und fordern, dass  $p_A(a)$  und  $p_A(b)$  unterschiedliche Vorzeichen besitzen. Dann muss nach Mittelwertsatz in dem Intervall  $[a, b]$  eine Nullstelle von  $p_A$  enthalten sein. Wir unterteilen  $[a, b]$  in zwei Teilintervall  $[a, c]$  und  $[c, b]$  mit  $c = (b + a)/2$  und fahren mit demjenigen Intervall fort, für das die Vorzeichenbedingung immer noch erfüllt ist.

**Algorithmus 7.2 (Bisektion)** *Seien  $a, b \in \mathbb{R}$  mit  $a < b$  und  $p_A(a)p_A(b) < 0$  gegeben. Dann approximiert der folgende Algorithmus einen Eigenwert von  $\mathbf{A}$  im Intervall  $[a, b]$ :*

```

 $p_a \leftarrow p_A(a); \quad p_b \leftarrow p_A(b)$ 
while  $|b - a| > \epsilon$  do begin
   $c \leftarrow (b + a)/2$ 
   $p_c \leftarrow p_A(c)$ 
  if  $p_a p_c < 0$  then begin
     $b \leftarrow c; \quad p_b \leftarrow p_c$ 
  end else begin
     $a \leftarrow c; \quad p_a \leftarrow p_c$ 
  end
end

```

Dieser Algorithmus benötigt pro Iterationsschritt nur eine Auswertung des Polynoms  $p_A$ , die sich, wie wir bereits gesehen haben, in  $\mathcal{O}(n)$  Operationen durchführen lässt. Da sich mit jedem Schritt das Intervall halbiert, können wir in  $\mathcal{O}(\log_2(1/\epsilon))$  Iterationen eine Genauigkeit von  $\epsilon \in \mathbb{R}_{>0}$  erreichen.

Ein Vorteil dieses Algorithmus besteht darin, dass wir das zu untersuchende Intervall explizit vorgeben können und dass er sehr stabil arbeitet, falls die Auswertung von  $p_A$

## 7 Verfahren für Tridiagonalmatrizen

stabil erfolgt. Ein Nachteil besteht darin, dass nicht klar ist, wie man ein Startintervall  $[a, b]$  finden kann, das die benötigte Vorzeichenbedingung  $p_A(a)p_A(b) < 0$  erfüllt.

Durch Differenzieren der Rekursionsformel für  $p_A$  können wir die Rekursionsformel

$$p'_i(x) = \begin{cases} 0 & \text{falls } i = 0, \\ 1 & \text{falls } i = 1, \\ p_{i-1}(x) + (x - \alpha_i)p'_{i-1}(x) & \text{ansonsten} \\ - |\beta_{i-1}|^2 p'_{i-2}(x) & \end{cases} \quad \text{für alle } i \in \mathbb{N}_0, x \in \mathbb{R} \quad (7.3)$$

gewinnen, mit der sich die erste Ableitung von  $p_A$  ebenfalls effizient berechnen lässt, so dass wir statt der Bisektion auch das Newton-Verfahren verwenden können. Zwar konvergiert das Newton-Verfahren unter Umständen wesentlich schneller als das Bisektionsverfahren, aber das passiert in der Regel nur, wenn ein guter Startwert vorliegt. Auch bei diesem Zugang stellt sich also die Frage nach geeigneten Startwerten.

## 7.2 Sturmsche Ketten

Wir sind daran interessiert, das einfache Bisektionsverfahren so zu modifizieren, dass wir nicht mehr auf eine Vorzeichenbedingung angewiesen sind, sondern eine Voraussetzung finden, die einfacher zu erfüllen ist.

Dazu gehen wir zunächst davon aus, dass alle Nullstellen von  $p_A$  einfach sind und wir sie deshalb in die Reihenfolge

$$\lambda_1 < \lambda_2 < \dots < \lambda_n$$

bringen können. Nach dem Satz von Rolle liegt zwischen zwei dieser Nullstellen jeweils mindestens eine Nullstelle der Ableitung  $p'_A$  des charakteristischen Polynoms, wir können also für jedes  $i \in [1 : n - 1]$  ein  $\lambda_i^{(1)} \in (\lambda_i, \lambda_{i+1})$  mit  $p'_A(\lambda_i^{(1)}) = 0$  finden. Wir erhalten

$$\lambda_1 < \lambda_1^{(1)} < \lambda_2 < \lambda_2^{(1)} < \lambda_3 < \dots < \lambda_{n-1} < \lambda_{n-1}^{(1)} < \lambda_n.$$

Da  $p'_A$  nur noch ein Polynom der Ordnung  $n - 1$  ist, kann es keine weiteren Nullstellen außer  $\lambda_1^{(1)}, \dots, \lambda_{n-1}^{(1)}$  besitzen.

In ähnlicher Weise können wir beweisen, dass die Nullstellen der  $(i + 1)$ -ten Ableitung von  $p_A$  gerade die der  $i$ -ten Ableitung trennen, solange  $i < n$  gilt.

Falls  $\xi$  eine einfache Nullstelle von  $p_A$  ist, können wir ein  $\epsilon > 0$  so finden, dass  $p'_A$  in  $[\xi - \epsilon, \xi + \epsilon]$  keine Nullstelle aufweist. Mit dem Mittelwertsatz der Differentialrechnung finden wir  $\eta_+ \in [\xi, \xi + \epsilon]$  und  $\eta_- \in [\xi - \epsilon, \xi]$  so, dass

$$\begin{aligned} p_A(\xi + \epsilon) &= p_A(\xi + \epsilon) - p_A(\xi) = \epsilon p'_A(\eta_+), \\ p_A(\xi - \epsilon) &= p_A(\xi - \epsilon) - p_A(\xi) = -\epsilon p'_A(\eta_-) \end{aligned}$$

gilt, und es folgt

$$p_A(\xi + \epsilon)p_A(\xi - \epsilon) = -\epsilon^2 p'_A(\eta_+)p'_A(\eta_-) < 0,$$

weil  $\xi$  eine Nullstelle des Polynoms  $p_A$  ist und  $p'_A(\eta_+)$  und  $p'_A(\eta_-)$  dasselbe Vorzeichen besitzen. An jeder einfachen Nullstelle wechselt  $p_A$  also das Vorzeichen. Wenn wir Nullstellen zählen wollen, bietet es sich demnach an, nach Vorzeichenwechseln zu suchen.

Es lässt sich leicht nachprüfen, dass der führende Koeffizient des charakteristischen Polynoms gerade 1 ist, so dass

$$\lim_{x \rightarrow \infty} p_A^{(m)}(x) = \infty \quad \text{für alle } m \in [0 : n - 1]$$

gilt. Wir haben bereits gesehen, dass die Nullstellen der Ableitung  $p_A^{(m+1)}$  zwischen denen der Funktion  $p_A^{(m)}$  liegen, also muss insbesondere

$$p_A^{(m)}(x) > 0 \quad \text{für alle } x > \lambda_n, \quad m \in [0 : n]$$

gelten. Für hinreichend großes  $x$  sind also die Vorzeichen aller Ableitungen des Polynoms  $p_A$  identisch.

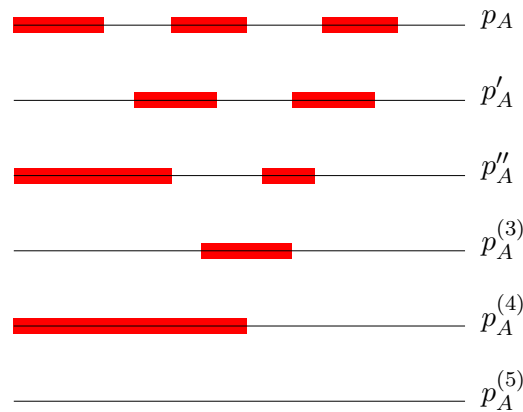


Abbildung 7.1: Vorzeichen eines charakteristischen Polynoms und seiner Ableitungen im Fall  $n = 5$ . Rot markiert Bereiche, in denen das Vorzeichen negativ ist.

In Abbildung 7.1 sind die Vorzeichen eines charakteristischen Polynoms und seiner Ableitungen für den Fall  $n = 5$  dargestellt. Die Bereiche, in denen die einzelnen Funktionen negativ sind, sind rot markiert. Für unsere Zwecke von Bedeutung ist die Beobachtung, dass in jedem Punkt  $x$  die Anzahl der Vorzeichenwechsel *zwischen den Ableitungen*  $p_A(x), p'_A(x), p''_A(x), \dots, p_A^{(n)}(x)$  gerade die Anzahl der Nullstellen größer als  $x$  angibt.

Diese Eigenschaft verdanken wir der Tatsache, dass zwischen zwei einfachen Nullstellen unserer Polynome jeweils genau eine einfache Nullstelle seiner Ableitung liegt, so dass „rechts“ von einer Nullstelle von  $p_A^{(m)}$  immer die Vorzeichen von  $p_A^{(m)}$  und  $p_A^{(m+1)}$  übereinstimmen. Falls  $m > 0$  gilt, sind die Vorzeichen von  $p_A^{(m)}$  und  $p_A^{(m-1)}$  ungleich. In diesem Fall wird also ein Vorzeichenwechsel zwischen  $p_A^{(m)}$  und  $p_A^{(m-1)}$  durch einen

## 7 Verfahren für Tridiagonalmatrizen

zwischen  $p_A^{(m)}$  und  $p_A^{(m+1)}$  ersetzt, lediglich für  $m = 0$  reduziert sich die Gesamtzahl der Vorzeichenwechsel. Falls eine exakte Null auftreten sollte, können wir festlegen, ob sie als positiv oder negativ gelten soll.

Die Ableitungen des charakteristischen Polynoms können wir zwar im Prinzip berechnen, sehr viel eleganter ist es allerdings, festzustellen, dass die bei seiner Auswertung berechneten Hilfspolynome  $p_m$  die Rolle der Ableitungen  $p_A^{(n-m)}$  übernehmen können. Dadurch können wir die Vorzeichenwechsel mit sehr geringem zusätzlichem Aufwand zählen.

Die für uns interessanten Eigenschaften eines Tupels von Funktionen fasst die folgende Definition zusammen:

**Definition 7.3 (Sturmsche Kette)** Ein Tupel  $(p_0, p_1, \dots, p_n)$  von reellen Polynomen heißt Sturmsche Kette, wenn

1.  $\overline{p_{n-1}(\xi)p_n'(\xi)} > 0$  für alle Nullstellen  $\xi \in \mathbb{K}$  von  $p_n$  erfüllt ist,
2.  $\overline{p_{i+1}(\xi)p_{i-1}(\xi)} < 0$  für alle  $i \in [1 : n - 1]$  und Nullstellen  $\xi \in \mathbb{K}$  von  $p_i$  gilt, sowie
3.  $p_0$  keine Nullstelle besitzt.

Bedingung 1 stellt sicher, dass in den Nullstellen des Polynom  $p_n$  die Vorzeichen der Ableitung  $p_n'$  und des Polynoms  $p_{n-1}$  übereinstimmen. Insbesondere müssen diese Nullstellen einfach sein, da  $p_n'$  nicht gleich null sein kann.

Bedingung 2 sorgt dafür, dass bei einer Nullstelle eines Polynoms  $p_i$  die „benachbarten“ Polynome  $p_{i-1}$  und  $p_{i+1}$  entgegengesetzte Vorzeichen aufweisen, so dass die Gesamtzahl der Vorzeichenwechsel unverändert bleibt.

Bedingung 3 schließlich benötigen wir, um ein „Referenzvorzeichen“ festzulegen, an dem wir die anderen Vorzeichen messen.

Nun stehen wir vor der Aufgabe, nachzuweisen, dass die durch die Rekursionsformel (7.2) definierten Polynome eine Sturmsche Kette bilden.

**Lemma 7.4 (Sturmsche Kette)** Falls  $\mathbf{A}$  irreduzibel ist, falls also  $\beta_i \neq 0$  für alle  $i \in [1 : n - 1]$  gilt, sind die durch (7.2) definierten Polynome  $p_0, \dots, p_n$  eine Sturmsche Kette.

*Beweis.* Der Beweis beruht darauf, Eigenvektoren zu den Eigenwerten der Matrix  $\mathbf{A}$  zu konstruieren. Für ein  $x \in \mathbb{K}$  suchen wir einen Vektor  $\mathbf{e}(x) \in \mathbb{K}^n$ , der „möglichst nahe“ an dem Kern von  $x\mathbf{I} - \mathbf{A}$  liegt, für den nämlich die ersten  $n - 1$  Komponenten des Vektors  $(x\mathbf{I} - \mathbf{A})\mathbf{e}(x)$  verschwinden.

Damit  $\mathbf{e}(x)$  nicht der Nullvektor wird, setzen wir  $e_1(x) = 1$ . Durch Einsetzen in die erste Zeile der Matrix erhalten wir nun

$$(x - \alpha_1)e_1(x) - \bar{\beta}_1 e_2(x) = 0,$$

$$e_2(x) = \frac{x - \alpha_1}{\bar{\beta}_1} e_1(x) = \frac{p_1(x)}{\bar{\beta}_1}$$



während Einsetzen in die  $i$ -te Zeile für  $i \in [2 : n - 1]$  die Gleichung

$$\begin{aligned} (-\beta_{i-1}e_{i-1}(x) + (x - \alpha_i)e_i(x) - \bar{\beta}_i e_{i+1}(x)) &= 0, \\ e_{i+1}(x) &= \frac{(x - \alpha_i)e_i(x) - \beta_{i-1}e_{i-1}(x)}{\bar{\beta}_i} \end{aligned}$$

für  $i \in [1 : n - 1]$  ergibt. Im Zähler dieses Terms erkennen wir Teile der Rekursionsformel (7.2) wieder und wählen den Ansatz

$$e_i(x) = \frac{p_{i-1}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{i-1}} \quad \text{für } i \in [2 : n].$$

Mit ihm nimmt unsere Gleichung die Form

$$\begin{aligned} e_{i+1}(x) &= \frac{\frac{(x - \alpha_i)p_{i-1}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{i-1}} - \frac{\beta_{i-1}p_{i-2}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{i-2}}}{\bar{\beta}_i} \\ &= \frac{(x - \alpha_i)\frac{p_{i-1}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{i-1}} - |\beta_{i-1}|^2 \frac{p_{i-2}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{i-1}}}{\bar{\beta}_i} \\ &= \frac{(x - \alpha_i)p_{i-1}(x) - |\beta_{i-1}|^2 p_{i-2}(x)}{\bar{\beta}_1 \dots \bar{\beta}_i} = \frac{p_i(x)}{\bar{\beta}_1 \dots \bar{\beta}_i} \end{aligned}$$

an. Mit einer einfachen Induktion folgt, dass der durch

$$e_i(x) = \begin{cases} 1 & \text{falls } i = 1, \\ \frac{p_{i-1}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{i-1}} & \text{ansonsten} \end{cases} \quad \text{für alle } i \in [1 : n]$$

definierte Vektor die gewünschten Gleichungen erfüllt.

Für die  $n$ -te Zeile schließlich erhalten wir

$$\begin{aligned} -\beta_{n-1}e_{n-1}(x) + (x - \alpha_n)e_n(x) &= -\beta_{n-1}\frac{p_{n-2}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{n-2}} + (x - \alpha_n)\frac{p_{n-1}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{n-1}} \\ &= -|\beta_{n-1}|^2 \frac{p_{n-2}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{n-1}} + (x - \alpha_n)\frac{p_{n-1}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{n-1}} \\ &= \frac{(x - \alpha_n)p_{n-1}(x) - |\beta_{n-1}|^2 p_{n-2}(x)}{\bar{\beta}_1 \dots \bar{\beta}_{n-1}} \\ &= \frac{p_n(x)}{\bar{\beta}_1 \dots \bar{\beta}_{n-1}} =: \gamma(x). \end{aligned}$$

Falls  $p_n(x) = 0$  gilt, haben wir auch  $\gamma(x) = 0$ , also folgt die Gleichung  $(x\mathbf{I} - \mathbf{A})\mathbf{e}(x) = \mathbf{0}$  und  $\mathbf{e}(x)$  ist ein Eigenvektor der Matrix  $\mathbf{A}$  zu dem Eigenwert  $x$ . Im allgemeinen Fall haben wir

$$(x\mathbf{I} - \mathbf{A})\mathbf{e}(x) = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \gamma(x) \end{pmatrix} \quad \text{für alle } x \in \mathbb{K} \quad (7.4)$$

## 7 Verfahren für Tridiagonalmatrizen

bewiesen. Indem wir diese Funktion nach  $x$  differenzieren, erhalten wir

$$\mathbf{e}(x) + (x\mathbf{I} - \mathbf{A})\mathbf{e}'(x) = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ \gamma'(x) \end{pmatrix} \quad \text{für alle } x \in \mathbb{K}. \quad (7.5)$$

Sei nun  $\xi \in \mathbb{K}$  eine Nullstelle des Polynoms  $p_n$ , also auch des Polynoms  $\gamma$ . Indem wir (7.5) mit  $\mathbf{e}(\xi)$  multiplizieren erhalten wir wegen  $(x\mathbf{I} - \mathbf{A})\mathbf{e}(x) = \mathbf{0}$  und Lemma 3.17 die Gleichung

$$\begin{aligned} 0 < \|\mathbf{e}(\xi)\|^2 &= \|\mathbf{e}(\xi)\|^2 + \langle (\xi\mathbf{I} - \mathbf{A})\mathbf{e}(\xi), \mathbf{e}'(\xi) \rangle = \|\mathbf{e}(\xi)\|^2 + \langle \mathbf{e}(\xi), (\xi\mathbf{I} - \mathbf{A})\mathbf{e}'(\xi) \rangle \\ &= \langle \mathbf{e}(\xi), \mathbf{e}(\xi) + (\xi\mathbf{I} - \mathbf{A})\mathbf{e}'(\xi) \rangle = \overline{e_n(\xi)}\gamma'(\xi) = \frac{\overline{p_{n-1}(\xi)}}{\beta_1 \dots \beta_{n-1}} \frac{p'_n(\xi)}{\bar{\beta}_1 \dots \bar{\beta}_{n-1}} \\ &= \frac{\overline{p_{n-1}(\xi)}p'_n(\xi)}{|\beta_1|^2 \dots |\beta_{n-1}|^2}, \end{aligned}$$

also insbesondere Bedingung 1 der Definition 7.3.

Für den Nachweis der Bedingung 2 dieser Definition setzen wir die Rekursionsformel (7.2) ein. Sei  $i \in [1 : n - 1]$ , und sei  $\xi \in \mathbb{K}$  eine Nullstelle von  $p_i$ . Dann folgt aus (7.2) die Gleichung

$$p_{i+1}(\xi) = -|\beta_i|^2 p_{i-1}(\xi),$$

also

$$\overline{p_{i+1}(\xi)}p_{i-1}(\xi) = -|\beta_i|^2 \overline{p_{i-1}(\xi)}p_{i-1}(\xi) = -|\beta_i|^2 |p_{i-1}(\xi)|^2 \leq 0.$$

Falls nun  $p_{i+1}(\xi) = 0$  gelten würde, hätten wir wegen  $\beta_i \neq 0$  auch  $p_{i-1}(\xi) = 0$  und könnten mit der Rekurrenzform (7.2) induktiv fortfahren, um zu dem Widerspruch  $0 = p_{i+1}(\xi) = \overline{p_i(\xi)} = \dots = p_0(\xi) = 1$  zu gelangen. Also müssen  $p_{i+1}(\xi), p_{i-1}(\xi) \neq 0$  gelten, und damit  $\overline{p_{i+1}(\xi)}p_{i-1}(\xi) < 0$ . ■

Aus diesem Lemma lassen sich bereits erste nützliche Aussagen über selbstadjungierte Tridiagonalmatrizen gewinnen.

**Folgerung 7.5 (Einfache Eigenwerte)** *Eine irreduzible selbstadjungierte Tridiagonalmatrix besitzt nur einfache Eigenwerte.*

*Beweis.* Mit Lemma 7.4 folgt die Aussage unmittelbar aus der ersten Bedingung in Definition 7.3. ■

Nach Satz 3.47 wissen wir bereits, dass selbstadjungierte Matrizen reell diagonalisierbar sind, dass also charakteristischen Polynome in reelle Linearfaktoren zerfallen. Bei einer irreduziblen selbstadjungierten Tridiagonalmatrix sind diese Linearfaktoren alle unterschiedlich, so dass wir reelle Eigenwerte  $\lambda_1 < \lambda_2 < \dots < \lambda_n$  finden.

**Folgerung 7.6 (Trennungseigenschaft)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine irreduzible selbstadjungierte Tridiagonalmatrix. Seien  $\lambda_1 < \lambda_2 < \dots < \lambda_n$  ihre Eigenwerte und  $\lambda'_1 < \lambda'_2 < \dots < \lambda'_{n-1}$  die Eigenwerte ihrer  $(n-1)$ -ten Hauptuntermatrix  $\mathbf{A}_{n-1}$ . Dann gilt

$$\lambda_1 < \lambda'_1 < \lambda_2 < \dots < \lambda_{n-1} < \lambda'_{n-1} < \lambda_n.$$

*Beweis.* Sei  $i \in [1 : n-1]$ . Nach der Bedingung 1 der Definition 7.3 können  $p'_A(\lambda_i)$  und  $p'_A(\lambda_{i+1})$  nicht gleich null sein,  $\lambda_i$  und  $\lambda_{i+1}$  sind also einfache Nullstellen des Polynoms  $p_A$ . Nach unserer Vorbetrachtung besitzt  $p'_A$  dann genau eine einfache Nullstelle in  $(\lambda_i, \lambda_{i+1})$ , also müssen sich die Vorzeichen von  $p'_A(\lambda_i)$  und  $p'_A(\lambda_{i+1})$  unterscheiden.

Nach der Bedingung 1 der Definition 7.3 weist  $p_{n-1}$  in  $\lambda_i$  und  $\lambda_{i+1}$  dieselben Vorzeichen wie  $p'_A = p'_n$  auf, also muss auch  $p_{n-1}$  mindestens eine Nullstelle  $\lambda'_i$  in  $(\lambda_i, \lambda_{i+1})$  besitzen.

Damit haben wir  $n-1$  Nullstellen von  $p_{n-1}$  gefunden, und da  $p_{n-1}$  höchstens den Grad  $n-1$  aufweist und nicht das Nullpolynom ist, können nach dem Identitätssatz für Polynome keine weiteren Nullstellen existieren. ■

Die für uns entscheidende Eigenschaft der Sturmschen Kette besteht darin, dass ihre Vorzeichenwechsel eine Beziehung zu der Anzahl der Nullstellen besitzen, so dass wir Aussagen darüber treffen können, wieviele Nullstellen in einem Intervall liegen, ohne sie explizit berechnen zu müssen.

**Satz 7.7 (Nullstellenzähler)** Sei  $(p_0, p_1, \dots, p_n)$  eine Sturmsche Kette. Wir definieren für alle  $x \in \mathbb{R}$

$$W_x := \{i \in [0 : n-1] : p_i(x)p_{i+1}(x) < 0 \text{ oder } p_i(x) = 0\}$$

und die Funktion

$$w : \mathbb{R} \rightarrow \mathbb{N}_0, \quad x \mapsto |W_x|,$$

die  $x$  die Mächtigkeit der Menge  $W_x$  zuordnet. Seien  $a, b \in \mathbb{R}$  mit  $a < b$  gegeben. Dann besitzt  $p_n$  genau  $w(a) - w(b)$  Nullstellen im Intervall  $(a, b]$ .

*Beweis.* Offensichtlich ändert sich  $w$  nur, wenn eines der Polynome  $(p_i)_{i=0}^n$  sein Vorzeichen ändert, also eine Nullstelle passiert. Wir definieren für alle  $x \in \mathbb{R}$  die Menge

$$N_x := \{i \in [0 : n] : p_i(x) = 0\}.$$

Infolge der Bedingung 3 in Definition 7.3 gilt  $0 \notin N_x$  für alle  $x \in \mathbb{R}$ . Wir werden nun nachweisen, dass die Mächtigkeit von  $W_x$  genau dann um eins sinkt, wenn  $x$  „von links nach rechts“ eine Nullstelle von  $p_n$  passiert.

Sei  $\xi \in \mathbb{R}$  mit  $N_\xi \neq \emptyset$ . Da nicht-konstante Polynome nur endlich viele Nullstellen besitzen können, können sich die Nullstellen nicht häufen, also muss ein  $\epsilon \in \mathbb{R}_{>0}$  so existieren, dass  $N_x = \emptyset$  für alle  $x \in [\xi - \epsilon, \xi + \epsilon] \setminus \{\xi\}$  gilt.

Sei  $i \in N_\xi$ . Falls  $i < n$  gilt, folgt aus Bedingung 2 in Definition 7.3 die Ungleichung  $p_{i-1}(\xi)p_{i+1}(\xi) < 0$ , also gelten  $i-1 \notin N_\xi$  und  $i+1 \notin N_\xi$ . Nach Wahl von  $\epsilon$  besitzen dann  $p_{i-1}$  und  $p_{i+1}$  keine Nullstellen in  $[\xi - \epsilon, \xi + \epsilon]$ , es folgt

$$p_{i-1}(x)p_{i+1}(x) < 0 \quad \text{für alle } x \in [\xi - \epsilon, \xi + \epsilon].$$

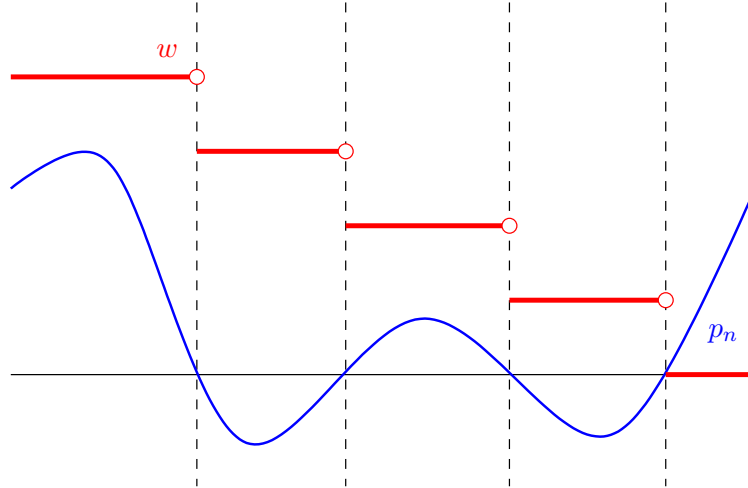


Abbildung 7.2: Die Funktion  $p_n$  (blau) und ihr „Nullstellenzähler“  $w$  (rot).

Sei  $x \in [\xi - \epsilon, \xi + \epsilon]$ . Falls  $p_{i-1}(x)p_i(x) < 0$  gilt, folgt  $p_{i+1}(x)p_i(x) > 0$ , also  $|W_x \cap \{i-1, i\}| = |\{i-1\}| = 1$ . Falls  $p_{i-1}(x)p_i(x) > 0$  gilt, folgt  $p_{i+1}(x)p_i(x) < 0$ , also  $|W_x \cap \{i-1, i\}| = |\{i\}| = 1$ . Falls schließlich  $p_{i-1}(x)p_i(x) = 0$  gilt, folgt  $p_i(x) = 0$ , also  $|W_x \cap \{i-1, i\}| = |\{i\}| = 1$ . Also ändern Nullstellen von  $p_i$  für  $i > 0$  nichts an der Mächtigkeit der Menge  $W_x$ .

Falls  $i = n$  gilt, folgt aus der Bedingung 1 in Definition 7.3 die Ungleichung  $p_{n-1}(\xi)p'_n(\xi) > 0$ , also insbesondere  $p_{n-1}(\xi) \neq 0$  und  $p'_n(\xi) \neq 0$ . Wir wählen  $\delta \in (0, \epsilon]$  so, dass  $p'_n$  keine Nullstellen in  $[\xi - \delta, \xi + \delta]$  besitzt. Mit dem Mittelwertsatz der Differentialrechnung finden wir  $\eta_+ \in (\xi, \xi + \delta]$  und  $\eta_- \in [\xi - \delta, \xi)$  mit

$$\begin{aligned} p_n(\xi + \delta) - p_n(\xi) &= p'_n(\eta_+) \delta, \\ p_n(\xi) - p_n(\xi - \delta) &= -p'_n(\eta_-) \delta. \end{aligned}$$

Aus  $p_{n-1}(\xi)p'_n(\xi) > 0$  und der Tatsache, dass beide Polynome aufgrund ihrer Stetigkeit ihre Vorzeichen in  $[\xi - \delta, \xi + \delta]$  nicht ändern können, folgen

$$\begin{aligned} p_{n-1}(\xi)p_n(\xi + \delta) &= p_{n-1}(\xi)p'_n(\eta_+) \delta > 0, \\ p_{n-1}(\xi)p_n(\xi - \delta) &= -p_{n-1}(\xi)p'_n(\eta_-) \delta < 0. \end{aligned}$$

Auf dem Intervall  $[\xi - \epsilon, \xi + \epsilon]$  kann  $p_{n-1}$  sein Vorzeichen nicht ändern, und  $p_n$  kann es nur an der Nullstelle  $\xi$ , damit erhalten wir

$$\begin{aligned} p_{n-1}(x)p_n(x) &< 0, \text{ also } W_x \cap \{n-1, n\} = \{n-1\} && \text{für alle } x \in [\xi - \delta, \xi), \\ p_{n-1}(x)p_n(x) &> 0, \text{ also } W_x \cap \{n-1, n\} = \emptyset && \text{für alle } x \in (\xi, \xi + \delta], \\ p_{n-1}(\xi)p_n(\xi) &= 0, \text{ also } W_x \cap \{n-1, n\} = \emptyset. \end{aligned}$$

Also reduzieren Nullstellen des Polynoms  $p_n$  die Mächtigkeit der Menge  $W_x$  um eins. ■

Der Satz 7.7 enthält lediglich eine Aussage über die Differenzen  $w(b) - w(a)$  der Funktion  $w$ . Das folgende Lemma bestimmt  $w$  näher:

**Lemma 7.8** *Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine irreduzible selbstadjungierte Tridiagonalmatrix, sei  $\lambda_1$  ihr kleinster Eigenwert und  $\lambda_n$  ihr größter. Sei  $p_0, \dots, p_n$  die gemäß (7.2) definierte Sturmsche Kette, und sei  $w: \mathbb{R} \rightarrow \mathbb{N}_0$  wie in Satz 7.7 definiert.*

*Für alle  $x \in \mathbb{R}_{<\lambda_1}$  gilt  $w(x) = n$ , und für alle  $x \in \mathbb{R}_{\geq\lambda_n}$  gilt  $w(x) = 0$ .*

*Insbesondere ist  $w(x)$  für jedes  $x \in \mathbb{R}$  gerade die Anzahl der Eigenwerte, die echt größer als  $x$  sind.*

*Beweis.* Nach Konstruktion des Polynoms  $p_i$  hat sein führender Koeffizient ein positives Vorzeichen. Für  $x \rightarrow \infty$  muss also  $p_i(x) \rightarrow \infty$  gelten, somit existiert ein  $x_0 \in \mathbb{R}$  mit

$$p_i(x) > 0 \quad \text{für alle } i \in [0 : n], \quad x \in \mathbb{R}_{\geq x_0}.$$

Es folgt  $w(x_0) = 0$ . Da  $w$  seinen Wert nur bei Nullstellen von  $p_n$  ändert, folgt  $w(x) = 0$  für alle  $x \in \mathbb{R}_{\geq\lambda_n}$ . Da  $p_n$  genau  $n$  einfache Nullstellen besitzt, folgt  $w(x) = n$  für alle  $x \in \mathbb{R}_{<\lambda_1}$  direkt aus Satz 7.7, indem wir  $b = x_0$  einsetzen. ■

Auf dieser Grundlage können wir einen Algorithmus konstruieren, der einen beliebigen Eigenwert bestimmt:

Sei  $k \in [1 : n]$ , und seien  $a, b \in \mathbb{R}$  mit  $a < b$  gegeben.

Falls  $w(b) \leq n - k$  gilt, enthält das unendliche Intervall  $(-\infty, b]$  gerade  $n - w(b) \geq k$  Eigenwerte, also insbesondere auch den  $k$ -ten Eigenwert, den wir suchen.

Falls  $w(a) > n - k$  gilt, enthält das unendliche Intervall  $(-\infty, a]$  gerade  $n - w(a) < k$  Eigenwerte, also gerade *nicht* den  $k$ -ten Eigenwert, den wir suchen.

Also muss der  $k$ -te Eigenwert in dem Intervall  $(a, b] = (-\infty, b] \setminus (-\infty, a]$  liegen.

**Algorithmus 7.9** *Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine irreduzible selbstadjungierte Tridiagonalmatrix, seien  $\lambda_1 < \dots < \lambda_n$  ihre Eigenwerte. Sei  $k \in [1, n]$ , und seien Intervallgrenzen  $a, b \in \mathbb{R}$  mit  $a < b$ ,  $w(b) \leq n - k < w(a)$  gegeben. Dann bestimmt der folgende Algorithmus eine Folge von  $(a, b]$ , die  $\lambda_k$  enthalten:*

```

while |b - a| > ε do begin
  c ← (b + a)/2
  if w(c) ≤ n - k then
    b ← c
  else
    a ← c
end

```

Auch dieser Bisektionsalgorithmus halbiert den Fehler in jedem Schritt und benötigt dank Algorithmus 7.1 nur  $\mathcal{O}(n)$  Operationen dafür. Allerdings ist auch er auf geeignete Startwerte angewiesen und nur auf selbstadjungierte irreduzible Tridiagonalmatrizen anwendbar.

## 7 Verfahren für Tridiagonalmatrizen

Wir haben bereits gesehen, dass wir jede beliebige selbstadjungierte Matrix mit unitären Ähnlichkeitstransformationen auf Tridiagonalgestalt bringen können. Falls die Tridiagonalmatrix nicht irreduzibel ist, können wir sie in eine Blockdiagonalmatrix mit irreduziblen Diagonalblöcken zerlegen.

Damit bleibt nur noch zu klären, wie wir ein geeignetes Anfangsintervall finden. Einen einfachen Zugang bieten *Gerschgorin-Kreise*, die uns eine Möglichkeit bieten, das Spektrum einer beliebigen Matrix einzugrenzen.

**Satz 7.10 (Gerschgorin-Kreise)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$ . Die abgeschlossenen Kreisscheiben

$$\mathcal{D}_i := \{z \in \mathbb{C} : |z - a_{ii}| \leq r_i\}, \quad r_i := \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|,$$

bezeichnen wir als die Gerschgorin-Kreise zu der Matrix  $\mathbf{A}$ . Es gilt

$$\sigma(\mathbf{A}) \subseteq \bigcup_{i=1}^n \mathcal{D}_i,$$

jeder Eigenwert ist also in mindestens einer der Kreisscheiben enthalten.

*Beweis.* Sei  $\lambda \in \sigma(\mathbf{A})$  ein Eigenwert der Matrix  $\mathbf{A}$ . Sei  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  ein korrespondierender Eigenvektor. Wir bezeichnen den Index des betragsgrößten Koeffizienten mit  $i \in [1 : n]$ , es gilt also

$$|x_j| \leq |x_i| \quad \text{für alle } j \in [1 : n].$$

Da  $\mathbf{x}$  ein Eigenvektor ist, folgt mit der Dreiecksungleichung

$$\begin{aligned} \lambda x_i &= (\lambda \mathbf{x})_i = (\mathbf{A} \mathbf{x})_i = \sum_{j=1}^n a_{ij} x_j, \\ (\lambda - a_{ii}) x_i &= \sum_{\substack{j=1 \\ j \neq i}}^n a_{ij} x_j, \\ |\lambda - a_{ii}| |x_i| &\leq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| |x_j| \leq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| |x_i| = r_i |x_i|, \\ |\lambda - a_{ii}| &\leq r_i, \end{aligned}$$

also bereits  $\lambda \in \mathcal{D}_i$ . ■

Im von uns untersuchten Fall selbstadjungierter Tridiagonalmatrizen  $\mathbf{A}$  sind die *Gerschgorin-Kreise* besonders einfach zu bestimmen: Wenn wir  $\beta_0 = \beta_n = 0$  setzen, erhalten wir

$$\mathcal{D}_i = \{z \in \mathbb{C} : |z - \alpha_i| \leq |\beta_{i-1}| + |\beta_i|\}.$$

Da  $\mathbf{A}$  selbstadjungiert ist, ist das Spektrum reell, es gilt also

$$\sigma(\mathbf{A}) \subseteq \bigcup_{i=1}^n [\alpha_i - (|\beta_{i-1}| + |\beta_i|), \alpha_i + (|\beta_{i-1}| + |\beta_i|)],$$

so dass wir folgern können, dass das Spektrum in dem Intervall  $[a, b]$  mit

$$\begin{aligned} a &:= \min\{\alpha_i - |\beta_{i-1}| - |\beta_i| : i \in [1 : n]\}, \\ b &:= \max\{\alpha_i + |\beta_{i-1}| + |\beta_i| : i \in [1 : n]\} \end{aligned}$$

enthalten ist. Dieses Intervall erfüllt also die Voraussetzungen von Algorithmus 7.9 für jedes  $k \in [1 : n]$ , so dass wir gezielt jeden beliebigen Eigenwert berechnen können.

### 7.3 Trägheitssatz und Dreieckszerlegungen

Der Einsatz der Sturmschen Ketten kann zu Schwierigkeiten führen, falls Rundungsfehler die Vorzeichen der einzelnen Polynome verfälschen. Es gibt allerdings alternative Möglichkeiten, um festzustellen, wie viele Eigenwerte kleiner oder größer als eine gegebene Zahl sind: Wir untersuchen, wieviele negative und positive Eigenwerte die „spektral verschobene“ Matrix  $\mathbf{A} - \mu\mathbf{I}$  besitzt. Die Anzahl der Eigenwerte, die echt kleiner als  $\mu$  sind, ist gerade die Anzahl der negativen Eigenwerte der Matrix  $\mathbf{A} - \mu\mathbf{I}$ .

Die Anzahl der negativen Eigenwerte lässt sich bestimmen, ohne sie explizit zu berechnen, indem wir geeignete Transformationen auf die Matrix anwenden, die wesentlich einfacher handzuhaben sind als die bisher verwendeten Ähnlichkeitstransformationen.

**Definition 7.11 (Kongruenztransformation)** Sei  $\mathbf{T} \in \mathbb{K}^{n \times n}$  eine reguläre Matrix. Die Abbildung

$$\mathbb{K}^{n \times n} \rightarrow \mathbb{K}^{n \times n}, \quad \mathbf{A} \mapsto \mathbf{T}^* \mathbf{A} \mathbf{T},$$

nennen wir die zu  $\mathbf{T}$  gehörende Kongruenztransformation.

Für unitäre Matrizen  $\mathbf{T}$  ist die Kongruenztransformation auch eine Ähnlichkeitstransformation, im allgemeinen Fall gilt das allerdings nicht. Trotzdem können Kongruenztransformationen sehr nützlich sein: Der *Trägheitssatz von Sylvester* besagt, dass die Anzahl der positiven und negativen Eigenwerte unter Kongruenztransformationen unverändert bleibt.

**Satz 7.12 (Trägheitssatz)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix. Sei  $\mathbf{T} \in \mathbb{K}^{n \times n}$  eine reguläre Matrix und

$$\widehat{\mathbf{A}} := \mathbf{T}^* \mathbf{A} \mathbf{T}.$$

Dann besitzen  $\mathbf{A}$  und  $\widehat{\mathbf{A}}$  jeweils gleich viele echt positive und echt negative Eigenwerte.

## 7 Verfahren für Tridiagonalmatrizen

*Beweis.* Nach Satz 3.47 existieren unitäre Matrizen  $\mathbf{Q}, \widehat{\mathbf{Q}} \in \mathbb{K}^{n \times n}$  und reelle Diagonalmatrizen  $\mathbf{D}, \widehat{\mathbf{D}} \in \mathbb{R}^{n \times n}$  mit

$$\begin{aligned} \mathbf{A} &= \mathbf{Q}\mathbf{D}\mathbf{Q}^*, & \mathbf{D} &= \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}, & \lambda_1 &\geq \lambda_2 \geq \dots \geq \lambda_n, \\ \widehat{\mathbf{A}} &= \widehat{\mathbf{Q}}\widehat{\mathbf{D}}\widehat{\mathbf{Q}}^*, & \widehat{\mathbf{D}} &= \begin{pmatrix} \hat{\lambda}_1 & & \\ & \ddots & \\ & & \hat{\lambda}_n \end{pmatrix}, & \hat{\lambda}_1 &\geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_n. \end{aligned}$$

Wir bezeichnen mit  $p, \hat{p} \in [0 : n]$  die Anzahl der echt positiven Eigenwerte der Matrizen  $\mathbf{A}$  und  $\widehat{\mathbf{A}}$ , es sollen also

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p > 0 \geq \lambda_{p+1}, \quad \hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_{\hat{p}} > 0 \geq \hat{\lambda}_{\hat{p}+1}$$

gelten. Unsere Aufgabe ist es, nachzuweisen, dass  $p = \hat{p}$  gilt.

Dazu definieren wir die Spaltenvektoren

$$\mathbf{q}^{(j)} := \mathbf{Q}\delta^{(j)}, \quad \widehat{\mathbf{q}}^{(j)} := \widehat{\mathbf{Q}}\delta^{(j)} \quad \text{für alle } j \in [1 : n].$$

Die von den Eigenvektoren zu echt positiven Eigenwerten aufgespannten Vektorräume bezeichnen wir mit

$$\mathcal{P} := \text{span}\{\mathbf{q}^{(1)}, \dots, \mathbf{q}^{(p)}\}, \quad \widehat{\mathcal{P}} := \text{span}\{\widehat{\mathbf{q}}^{(1)}, \dots, \widehat{\mathbf{q}}^{(\hat{p})}\}.$$

Wenn wir zeigen können, dass beide Räume dieselbe Dimension besitzen, sind wir fertig.

Als Hilfsmittel führen wir die Räume

$$\mathcal{N} := \text{span}\{\mathbf{q}^{(p+1)}, \dots, \mathbf{q}^{(n)}\}, \quad \widehat{\mathcal{N}} := \text{span}\{\widehat{\mathbf{q}}^{(\hat{p}+1)}, \dots, \widehat{\mathbf{q}}^{(n)}\}$$

zu den nicht-positiven Eigenwerten ein und untersuchen die Abbildungen

$$\begin{aligned} f: \mathbb{K}^n &\rightarrow \mathbb{R}, & \mathbf{x} &\mapsto \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle, \\ \hat{f}: \mathbb{K}^n &\rightarrow \mathbb{R}, & \widehat{\mathbf{x}} &\mapsto \langle \widehat{\mathbf{x}}, \widehat{\mathbf{A}}\widehat{\mathbf{x}} \rangle. \end{aligned}$$

Für einen Vektor  $\mathbf{x} \in \mathcal{P} \setminus \{\mathbf{0}\}$  finden wir Koeffizienten  $\alpha_1, \dots, \alpha_j \in \mathbb{K}$  mit

$$\mathbf{x} = \sum_{j=1}^p \alpha_j \mathbf{q}^{(j)},$$

die wegen  $\mathbf{x} \neq \mathbf{0}$  nicht alle gleich null sein können, und erhalten

$$f(\mathbf{x}) = \sum_{i=1}^p \sum_{j=1}^p \bar{\alpha}_i \alpha_j \langle \mathbf{q}^{(i)}, \mathbf{A}\mathbf{q}^{(j)} \rangle = \sum_{i=1}^p |\alpha_i|^2 \langle \mathbf{q}^{(i)}, \lambda_i \mathbf{q}^{(i)} \rangle = \sum_{i=1}^p |\alpha_i|^2 \lambda_i > 0.$$



Für  $\hat{f}$  und  $\hat{\mathcal{N}}$  erhalten wir ein entsprechendes Resultat und dürfen

$$f(\mathbf{x}) > 0, \quad \hat{f}(\hat{\mathbf{x}}) \leq 0 \quad \text{für alle } \mathbf{x} \in \mathcal{P} \setminus \{\mathbf{0}\}, \hat{\mathbf{x}} \in \hat{\mathcal{N}}$$

festhalten. Sei nun  $\hat{\mathbf{x}} \in \hat{\mathcal{N}}$  gegeben. Wir stellen fest, dass

$$0 \geq \hat{f}(\hat{\mathbf{x}}) = \langle \hat{\mathbf{x}}, \hat{\mathbf{A}}\hat{\mathbf{x}} \rangle = \langle \hat{\mathbf{x}}, \mathbf{T}^* \mathbf{A} \mathbf{T} \hat{\mathbf{x}} \rangle = \langle \mathbf{T}\hat{\mathbf{x}}, \mathbf{A} \mathbf{T}\hat{\mathbf{x}} \rangle = f(\mathbf{T}\hat{\mathbf{x}})$$

gilt. Also folgt

$$f(\mathbf{x}) \leq 0 \quad \text{für alle } \mathbf{x} \in \mathbf{T}\hat{\mathcal{N}},$$

also kann der Schnitt der Teilräume  $\mathcal{P}$  und  $\mathbf{T}\hat{\mathcal{N}}$  nur den Nullvektor enthalten.

Da  $\mathbf{T}$  invertierbar ist, folgt daraus

$$n \geq \dim(\mathcal{P}) + \dim(\mathbf{T}\hat{\mathcal{N}}) = \dim(\mathcal{P}) + \dim(\hat{\mathcal{N}}) = p + n - \hat{p},$$

also  $0 \geq p - \hat{p}$  und damit  $\hat{p} \geq p$ .

Indem wir entsprechend mit den Räumen  $\hat{\mathcal{P}}$  und  $\mathcal{N}$  verfahren, folgt auch  $p \geq \hat{p}$ , so dass  $p = \hat{p}$  bewiesen ist.

Wir können dieselbe Argumentation auf die Matrizen  $-\mathbf{A}$  und  $-\hat{\mathbf{A}}$  anwenden, um zu zeigen, dass auch die Anzahl der echt negativen Eigenwerte identisch ist. ■

Unser Ziel ist es, die Anzahl der negativen Eigenwerte der Matrix  $\mathbf{A} - \mu\mathbf{I}$  zu berechnen. Dank des Trägheitssatzes 7.12 ändert sich diese Anzahl nicht, wenn wir Kongruenztransformationen auf die Matrix anwenden, also bietet es sich an, nach Kongruenztransformationen zu suchen, die die Matrix in eine Form bringen, an der wir die Anzahl der negativen Eigenwerte unmittelbar ablesen können. Ideal geeignet wäre eine Diagonalmatrix, wir suchen also eine reguläre Matrix  $\mathbf{T} \in \mathbb{K}^{n \times n}$  derart, dass  $\mathbf{T}^*(\mathbf{A} - \mu\mathbf{I})\mathbf{T}$  eine Diagonalmatrix ist, denn dann stehen die Eigenwerte auf der Diagonalen.

Diese Aufgabe lässt sich lösen, indem wir  $\mathbf{T}$  als Dreiecksmatrix wählen und eine verallgemeinerte Form der *Cholesky-Zerlegung* berechnen. Zur Abkürzung setzen wir  $\mathbf{B} := \mathbf{A} - \mu\mathbf{I}$  und suchen nach einer normierten unteren Dreiecksmatrix  $\mathbf{L} \in \mathbb{K}^{n \times n}$  und einer reellen Diagonalmatrix  $\mathbf{D} \in \mathbb{R}^{n \times n}$  mit

$$\mathbf{B} = \mathbf{L}\mathbf{D}\mathbf{L}^* \quad \iff \quad \mathbf{L}^{-1}(\mathbf{A} - \mu\mathbf{I})\mathbf{L}^{-*} = \mathbf{D}.$$

Wir zerlegen die auftretenden Matrizen in Teilmatrizen:

$$\begin{aligned} \mathbf{B} &= \begin{pmatrix} b_{11} & \mathbf{B}_{1*} \\ \mathbf{B}_{*1} & \mathbf{B}_{**} \end{pmatrix} & \text{mit } \mathbf{B}_{1*} \in \mathbb{K}^{1 \times (n-1)}, \mathbf{B}_{**} \in \mathbb{K}^{(n-1) \times (n-1)}, \\ \mathbf{L} &= \begin{pmatrix} 1 & \\ \mathbf{L}_{*1} & \mathbf{L}_{**} \end{pmatrix} & \text{mit } \mathbf{L}_{*1} \in \mathbb{K}^{1 \times (n-1)}, \mathbf{L}_{**} \in \mathbb{K}^{(n-1) \times (n-1)}, \\ \mathbf{D} &= \begin{pmatrix} d_1 & \\ & \mathbf{D}_{**} \end{pmatrix} & \text{mit } \mathbf{D}_{**} \in \mathbb{R}^{(n-1) \times (n-1)}. \end{aligned}$$

## 7 Verfahren für Tridiagonalmatrizen

Durch Einsetzen in die definierende Gleichung erhalten wir

$$\begin{aligned} \begin{pmatrix} b_{11} & \mathbf{B}_{1*} \\ \mathbf{B}_{*1} & \mathbf{B}_{**} \end{pmatrix} = \mathbf{B} = \mathbf{L}\mathbf{D}\mathbf{L}^* &= \begin{pmatrix} 1 & \\ \mathbf{L}_{*1} & \mathbf{L}_{**} \end{pmatrix} \begin{pmatrix} d_1 & \\ & \mathbf{D}_{**} \end{pmatrix} \begin{pmatrix} 1 & \mathbf{L}_{*1}^* \\ & \mathbf{L}_{**}^* \end{pmatrix} \\ &= \begin{pmatrix} d_1 & \\ \mathbf{L}_{*1}d_1 & \mathbf{L}_{**}\mathbf{D}_{**} \end{pmatrix} \begin{pmatrix} 1 & \mathbf{L}_{*1}^* \\ & \mathbf{L}_{**}^* \end{pmatrix} = \begin{pmatrix} d_1 & d_1\mathbf{L}_{*1}^* \\ \mathbf{L}_{*1}d_1 & \mathbf{L}_{*1}d_1\mathbf{L}_{*1}^* + \mathbf{L}_{**}\mathbf{D}_{**}\mathbf{L}_{**}^* \end{pmatrix}. \end{aligned}$$

Aus  $\mathbf{B}^* = \mathbf{B}$  folgt  $\mathbf{B}_{1*} = \mathbf{B}_{*1}^*$ , so dass lediglich die folgenden drei Gleichungen zu erfüllen sind:

$$\begin{aligned} b_{11} &= d_1, \\ \mathbf{B}_{*1} &= \mathbf{L}_{*1}d_1, \\ \mathbf{B}_{**} &= \mathbf{L}_{*1}d_1\mathbf{L}_{*1}^* + \mathbf{L}_{**}\mathbf{D}_{**}\mathbf{L}_{**}^*. \end{aligned}$$

Falls  $b_{11} \neq 0$  gilt, können wir sie umformen und finden

$$\begin{aligned} d_1 &= b_{11}, \\ \mathbf{L}_{*1} &= \mathbf{B}_{*1}/d_1, \\ \mathbf{L}_{**}\mathbf{D}_{**}\mathbf{L}_{**}^* &= \widehat{\mathbf{B}} := \mathbf{B}_{**} - \mathbf{L}_{*1}d_1\mathbf{L}_{*1}^*. \end{aligned}$$

Die letzte Gleichung ist von der Form des ursprünglichen Problems, so dass wir sie per Induktion behandeln können.

In unserem Fall ist  $\mathbf{B}$  eine Tridiagonalmatrix der Form

$$\mathbf{B} = \begin{pmatrix} \alpha_1 - \mu & \bar{\beta}_1 & & \\ \beta_1 & \alpha_2 - \mu & \ddots & \\ & \ddots & \ddots & \bar{\beta}_{n-1} \\ & & \beta_{n-1} & \alpha_n - \mu \end{pmatrix},$$

so dass die Gleichungen

$$\mathbf{L}_{*1} = \begin{pmatrix} \beta_1/(\alpha_1 - \mu) \\ 0 \\ \vdots \\ 0 \end{pmatrix}, \quad \widehat{\mathbf{B}} = \begin{pmatrix} \alpha_2 - \mu - \frac{|\beta_1|^2}{(\alpha_1 - \mu)^2} & \bar{\beta}_2 & & \\ \beta_2 & \alpha_3 - \mu & \ddots & \\ & \ddots & \ddots & \bar{\beta}_{n-1} \\ & & \beta_{n-1} & \alpha_n - \mu \end{pmatrix}$$

gelten und sich demnach der Induktionsschritt mit wenigen Rechenoperationen vollziehen lässt.

Ein weiteres Verfahren für Tridiagonalmatrizen soll hier nur kurz in Form einer Übungsaufgabe angerissen werden: Wir können eine Tridiagonalmatrix  $\mathbf{A}$  in zwei ungefähr gleich große unabhängige Diagonalblöcke zerlegen, indem wir  $k = \lfloor n/2 \rfloor$  wählen und

$$\mathbf{R} := (|\beta_k|\delta^{(k)} + \beta_k\delta^{(k+1)})(\delta^{(k)} + \frac{\beta_k}{|\beta_k|}\delta^{(k+1)})^*$$

$$= |\beta_k| \delta^{(k)} (\delta^{(k)})^* + \beta_k \delta^{(k+1)} (\delta^{(k)})^* + \bar{\beta}_k \delta^{(k)} (\delta^{(k+1)})^* + |\beta_k| \delta^{(k+1)} (\delta^{(k+1)})^*$$

subtrahieren, um den  $k$ -ten Nebendiagonaleintrag zu eliminieren.

Die Diagonalblöcke diagonalisieren wir rekursiv. Um die Eigenwerte der Originalmatrix zu erhalten, müssen wir nun  $\mathbf{R}$  wieder hinzuaddieren. Dabei können wir ausnutzen, dass es sich um eine Matrix mit Rang eins handelt.

**Übungsaufgabe 7.13 (Niedrigrangstörung)** Seien  $n, m \in \mathbb{N}$ .

- (a) Seien  $\mathbf{A} \in \mathbb{K}^{n \times m}$  und  $\mathbf{B} \in \mathbb{K}^{m \times n}$ . Beweisen Sie  $\sigma(\mathbf{AB}) \cup \{0\} = \sigma(\mathbf{BA}) \cup \{0\}$ .
- (b) Seien  $\mathbf{A} \in \mathbb{K}^{n \times m}$  und  $\mathbf{B} \in \mathbb{K}^{m \times n}$ . Beweisen Sie  $\det(\mathbf{I} + \mathbf{AB}) = \det(\mathbf{I} + \mathbf{BA})$ .
- (c) Sei  $\mathbf{D} = \text{diag}(\lambda_1, \dots, \lambda_n)$  eine reelle Diagonalmatrix mit  $\lambda_1 < \lambda_2 < \dots < \lambda_n$ , seien  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ , und sei  $\tilde{\mathbf{D}} := \mathbf{D} + \mathbf{ab}^*$ . Beweisen Sie

$$p_{\tilde{\mathbf{D}}}(x) = p_{\mathbf{D}}(x) \left( 1 - \sum_{i=1}^n \frac{a_i b_i}{x - \lambda_i} \right) \quad \text{für alle } x \in \mathbb{K} \setminus \sigma(\mathbf{D}).$$

- (d) Beweisen Sie unter den Voraussetzungen des Aufgabenteils (c): Falls  $a_i b_i > 0$ ,  $a_{i+1} b_{i+1} > 0$  für ein  $i \in [1 : n - 1]$  gelten, besitzt  $\tilde{\mathbf{D}}$  einen Eigenwert im Intervall  $(\lambda_i, \lambda_{i+1})$ .

Hinweise: Sind die Matrizen  $\begin{pmatrix} \alpha \mathbf{I} + \mathbf{AB} & \mathbf{A} \\ \mathbf{0} & \alpha \mathbf{I} \end{pmatrix}$  und  $\begin{pmatrix} \alpha \mathbf{I} & \mathbf{A} \\ \mathbf{0} & \alpha \mathbf{I} + \mathbf{BA} \end{pmatrix}$  ähnlich? Der Determinanten-Multiplikationssatz ist hilfreich. Es darf verwendet werden, dass die Determinante einer oberen Block-Dreiecksmatrix das Produkt der Determinanten der Diagonalblöcke ist.

Neben Tridiagonalmatrizen gibt es noch weitere spezielle Matrixtypen, für die die Behandlung des Eigenwertproblems besonders einfach ist.

**Übungsaufgabe 7.14 (Zirkulante Matrizen)** Sei  $n \in \mathbb{N}$ . Wir zerlegen  $\mathbb{Z}$  in die Äquivalenzklassen

$$[k]_n := \{k + nz : z \in \mathbb{Z}\} \quad \text{für alle } k \in \mathbb{Z}.$$

Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine zirkulante Matrix, es gelte also

$$[j - i]_n = [k - \ell]_n \Rightarrow a_{ij} = a_{k\ell} \quad \text{für alle } i, j, k, \ell \in [1 : n].$$

Wir bezeichnen mit  $\omega_n := \exp(2\pi i/n)$  die  $n$ -te Einheitswurzel. Hier ist  $i$  die imaginäre Einheit mit  $i^2 = -1$ , und es gilt  $\omega_n^n = 1$ . Die diskrete Fourier-Transformation  $\mathbf{F} \in \mathbb{C}^{n \times n}$  ist definiert durch

$$f_{ij} = \frac{1}{\sqrt{n}} \omega_n^{ij} \quad \text{für alle } i, j \in [1 : n].$$

## 7 Verfahren für Tridiagonalmatrizen

- (a) Beweisen Sie  $\mathbf{F}^*\mathbf{F} = \mathbf{I}$ .
- (b) Beweisen Sie, dass  $\alpha_0, \dots, \alpha_{n-1} \in \mathbb{K}$  existieren mit

$$\mathbf{A} = \begin{pmatrix} \alpha_0 & \alpha_1 & \dots & \alpha_{n-1} \\ \alpha_{n-1} & \alpha_0 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \alpha_1 \\ \alpha_1 & \dots & \alpha_{n-1} & \alpha_0 \end{pmatrix}.$$

- (c) Beweisen Sie, dass  $\mathbf{F}^*\mathbf{A}\mathbf{F}$  eine Diagonalmatrix ist.

Hinweise: Bei Teil (a) mag die geometrische Summenformel nützlich sein, bei Teil (c) könnte man  $\mathbf{A}\mathbf{F} = \mathbf{F}\mathbf{D}$  mit einer geeigneten Diagonalmatrix  $\mathbf{D}$  zeigen und Teil (a) verwenden.

**Übungsaufgabe 7.15 (Tensorprodukte)** Seien  $n, m \in \mathbb{N}$ , und seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und  $\mathbf{B} \in \mathbb{K}^{m \times m}$ . Das Tensorprodukt der Matrizen ist definiert durch

$$\mathbf{A} \otimes \mathbf{B} := \begin{pmatrix} a_{11}\mathbf{B} & \dots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{n1}\mathbf{B} & \dots & a_{nn}\mathbf{B} \end{pmatrix} \in \mathbb{K}^{(nm) \times (nm)}.$$

- (a) Seien  $\lambda, \mu \in \mathbb{K}$  Eigenwerte der Matrizen  $\mathbf{A}$  und  $\mathbf{B}$  mit Eigenvektoren  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  und  $\mathbf{y} \in \mathbb{K}^m \setminus \{\mathbf{0}\}$ . Beweisen Sie, dass  $\lambda\mu$  ein Eigenwert der Matrix  $\mathbf{A} \otimes \mathbf{B}$  ist und geben Sie einen Eigenvektor an.
- (b) Seien  $\lambda, \mu \in \mathbb{K}$  Eigenwerte der Matrizen  $\mathbf{A}$  und  $\mathbf{B}$  mit Eigenvektoren  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  und  $\mathbf{y} \in \mathbb{K}^m \setminus \{\mathbf{0}\}$ . Beweisen Sie, dass  $\lambda + \mu$  ein Eigenwert der Matrix  $\mathbf{A} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{B}$  ist und geben Sie einen Eigenvektor an.
- (c) Seien  $\mathbf{C} \in \mathbb{K}^{n \times n}$  und  $\mathbf{D} \in \mathbb{K}^{m \times m}$ . Beweisen Sie  $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{A}\mathbf{C}) \otimes (\mathbf{B}\mathbf{D})$ .
- (d) Seien  $\mathbf{A}$  und  $\mathbf{B}$  diagonalisierbar. Beweisen Sie, dass  $\mathbf{A} \otimes \mathbf{B}$  diagonalisierbar ist.

# 8 Lanczos-Verfahren für schwachbesetzte Matrizen

In der Praxis treten sehr häufig Matrizen auf, die besondere Eigenschaften aufweisen, die zur Beschleunigung bestimmter Operationen ausgenutzt werden können: Bei einer Tridiagonalmatrix lassen sich in  $\mathcal{O}(n)$  Operationen die Matrix-Vektor-Multiplikation ausführen, eine QR-Zerlegung konstruieren oder das charakteristische Polynom auswerten. Bei einer Hessenberg-Matrix kann eine QR-Zerlegung immerhin in  $\mathcal{O}(n^2)$  Operationen berechnet werden, während bei einer allgemeinen Matrix ein kubischer Aufwand erforderlich ist.

Wir werden uns in diesem Kapitel auf Matrizen konzentrieren, bei denen die Matrix-Vektor-Multiplikation  $\mathbf{x} \mapsto \mathbf{Ax}$  effizient ausgeführt werden kann. Ein prominentes Beispiel sind *schwachbesetzte Matrizen*, die sich dadurch auszeichnen, dass die meisten ihrer Koeffizienten gleich null sind und deshalb bei der Multiplikation keine Rolle spielen.

## 8.1 Zweidimensionales Modellproblem

Wir untersuchen als Modellproblem die numerische Approximation des zweidimensionalen *Poisson-Eigenwertproblems*

$$-\Delta u(x, y) = \lambda u(x, y)$$

mit  $(x, y) \in ]0, 1[^2$ , der Randbedingung

$$u(0, \cdot) = u(1, \cdot) = u(\cdot, 0) = u(\cdot, 1) = 0$$

und dem Laplace-Operator

$$\Delta u(x, y) := \frac{\partial^2 u}{\partial x^2}(x, y) + \frac{\partial^2 u}{\partial y^2}(x, y). \tag{8.1}$$

Zur Diskretisierung wählen wir ein  $N \in \mathbb{N}$  und setzen

$$h := \frac{1}{N+1}, \quad \mathcal{I} := [1 : N]^2, \quad \Omega_h := \{(hi, hj) : (i, j) \in \mathcal{I}\}.$$

Analog zu dem Beispiel in Abschnitt 2.1 approximieren wir die Differentialquotienten in (8.1) durch Differenzenquotienten:

$$-\Delta u(x, y) \approx \frac{4u(x, y) - u(x-h, y) - u(x+h, y) - u(x, y-h) - u(x, y+h)}{h^2}.$$

Wie im eindimensionalen Fall schränken wir den Definitionsbereich auf  $\Omega_h$  ein und erhalten ein lineares Gleichungssystem im Raum  $\mathbb{R}^{\mathcal{I}}$ : Wir ersetzen den Differentialoperator  $-\Delta$  durch eine Matrix  $\mathbf{A} \in \mathbb{R}^{\mathcal{I} \times \mathcal{I}}$  und die Funktion  $u$  durch einen Vektor  $\mathbf{x} \in \mathbb{R}^{\mathcal{I}}$ , die durch

$$a_{ij} := \begin{cases} 4h^{-2} & \text{falls } i_x = j_x \text{ und } i_y = j_y, \\ -h^{-2} & \text{falls } |i_x - j_x| = 1, \\ -h^{-2} & \text{falls } |i_y - j_y| = 1, \\ 0 & \text{ansonsten} \end{cases},$$

$$x_j := u(hj_x, hj_y) \quad \text{für alle } i = (i_x, i_y) \in \mathcal{I}, j = (j_x, j_y) \in \mathcal{I}$$

definiert sind, und erhalten das Gleichungssystem

$$\mathbf{A}\mathbf{x} \approx \lambda\mathbf{x}.$$

Wie im eindimensionalen Fall ist die Matrix  $\mathbf{A}$  symmetrisch und positiv definit, und auch in diesem Fall sind wir an ihren kleinsten Eigenwerten interessiert. Unsere Aufgabe besteht also darin, ein  $n$ -dimensionales Eigenwertproblem zu lösen, wobei  $n = N^2 \approx h^{-2}$  gilt und sich beweisen lässt, dass die Eigenwerte und Eigenvektoren des diskreten Systems mit einem Fehler proportional zu  $h^2$  gegen die des ursprünglichen kontinuierlichen Problems konvergieren. Um eine brauchbare Genauigkeit zu erreichen, müssen wir also darauf vorbereitet sein, sehr große Eigenwertprobleme zu behandeln.

Während allerdings im eindimensionalen Fall eine Tridiagonalmatrix entstand, lässt sich im zweidimensionalen Fall keine Anordnung der Indexmenge  $\mathcal{I}$  finden, die aus  $\mathbf{A}$  eine Bandmatrix mit von  $N$  unabhängiger Bandbreite macht: Man kann beweisen, dass die Bandbreite immer mindestens  $N/2$  betragen muss. In der Praxis verwendet man häufig die *lexikographische* Anordnung der Indexmenge  $\mathcal{I}$ , bei der die Matrix  $\mathbf{A}$  die Block-Tridiagonaldarstellung

$$\mathbf{A} = h^{-2} \begin{pmatrix} \mathbf{T} & -\mathbf{I} & & & \\ -\mathbf{I} & \ddots & \ddots & & \\ & \ddots & \ddots & -\mathbf{I} & \\ & & & -\mathbf{I} & \mathbf{T} \end{pmatrix} \in \mathbb{R}^{n \times n}, \quad \mathbf{T} = \begin{pmatrix} 4 & -1 & & & \\ -1 & \ddots & \ddots & & \\ & \ddots & \ddots & -1 & \\ & & & -1 & 4 \end{pmatrix} \in \mathbb{R}^{N \times N},$$

besitzt, die eine Bandbreite von  $N$  aufweist.

Eine Besonderheit von Tridiagonalmatrizen weist allerdings auch die Matrix  $\mathbf{A}$  unabhängig von der Anordnung der Indizes auf: Da pro Zeile dieser Matrizen lediglich höchstens fünf von null verschiedene Einträge auftreten, lässt sich das Matrix-Vektor-Produkt  $\mathbf{y} := \mathbf{A}\mathbf{x}$  für  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^{\mathcal{I}}$  in  $\mathcal{O}(n)$  Operationen auswerten. Matrizen mit dieser Eigenschaft werden als *schwachbesetzt* bezeichnet.

Theoretisch können wir alle bereits diskutierten Verfahren auf die Matrix  $\mathbf{A}$  anwenden, praktisch stoßen wir dabei allerdings auf Schwierigkeiten: Sobald wir Linearkombinationen von Zeilen oder Spalten berechnen, beispielsweise bei der Anwendung einer Givens-Rotation, entstehen dadurch in der Regel neue von null verschiedene Einträge in

der Matrix. Das bedeutet, dass der Speicherbedarf und der Rechenaufwand sehr stark zunehmen. Für hohe Problemdimensionen ist dieser Ansatz deshalb nicht mehr sinnvoll durchführbar.

Günstig dagegen wäre die Vektoriteration, da sie lediglich Matrix-Vektor-Multiplikationen erfordert, die sich, wie gesagt, sehr effizient durchführen lassen. Leider eignet sie sich nur für die Berechnung der größten Eigenwerte, während wir bei dieser Anwendung an den kleinsten interessiert sind.

**Bemerkung 8.1** *Mit der in Übungsaufgabe 7.15 eingeführten Notation erhalten wir*

$$\mathbf{A} = \mathbf{L} \otimes \mathbf{I} + \mathbf{I} \otimes \mathbf{L},$$

wobei

$$\mathbf{L} := h^{-2} \begin{pmatrix} 2 & -1 & & & \\ -1 & \ddots & \ddots & & \\ & \ddots & \ddots & -1 & \\ & & & -1 & 2 \end{pmatrix}$$

die aus Abschnitt 2.1 bekannte Matrix des eindimensionalen Modellproblems ist.

Aus Übungsaufgabe 5.14 sind uns die Eigenwerte der Matrix  $\mathbf{L}$  bekannt, und mit Übungsaufgabe 7.15 können wir damit auch die Eigenwerte der Matrix  $\mathbf{A}$  exakt berechnen.

## 8.2 Krylow-Räume

Ein sinnvoller erster Schritt bei der Berechnung der Eigenwerte einer selbstadjungierten Matrix wird in der Regel darin bestehen, die Matrix auf Hessenberg-, also Tridiagonalgestalt zu bringen, da sich beispielsweise die QR-Iteration für solche Matrizen sehr effizient durchführen lässt.

Bisher haben wir für diesen Zweck Householder-Spiegelungen oder Givens-Rotationen eingesetzt, die bei schwachbesetzten Matrizen die vorhandene Struktur zunichte machen würden und deshalb unattraktiv sind. In Übungsaufgabe 6.13 haben wir bereits gesehen, dass es unter gewissen Umständen möglich ist, eine Matrix nur mit Matrix-Vektor-Multiplikationen auf Hessenberg-Gestalt zu bringen. Dieser Ansatz soll nun ausgearbeitet und genauer untersucht werden.

Insbesondere interessieren wir uns für die Frage, ob wirklich eine vollständige Hessenberg-Transformation nötig ist, oder ob wir die Konstruktion eventuell früher abbrechen können und trotzdem gute Näherungen der gesuchten Eigenwerte erhalten.

Sei im folgenden  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix mit den Eigenwerten  $\lambda_1 \leq \dots \leq \lambda_n$ . Dank des Satzes 3.45 von Courant und Fischer wissen wir, dass der Rayleigh-Quotient

$$\Lambda_A: \mathbb{K}^n \setminus \{\mathbf{0}\} \rightarrow \mathbb{R}, \quad \mathbf{x} \mapsto \frac{\langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle_2}{\langle \mathbf{x}, \mathbf{x} \rangle_2},$$

die Gleichungen

$$\begin{aligned}\lambda_1 &= \min\{\Lambda_A(\mathbf{x}) : \mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}\}, \\ \lambda_n &= \max\{\Lambda_A(\mathbf{x}) : \mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}\}\end{aligned}$$

erfüllt. In ähnlicher Weise lassen sich auch andere Eigenwerte mit Hilfe geeigneter lokaler Minima und Maxima des Rayleigh-Quotienten charakterisieren.

Eine gute Strategie zur Approximation der Eigenwerte könnte also darin bestehen, den Rayleigh-Quotienten zu minimieren oder zu maximieren. Für derartige Aufgaben ist das *Gradientenverfahren* geeignet: Der *Gradient* einer Funktion  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  in einem Punkt  $\mathbf{x} \in \mathbb{R}^n$  ist der Vektor  $\nabla f(\mathbf{x}) \in \mathbb{R}^n$ , für den

$$Df(\mathbf{x}) \cdot \mathbf{y} = \langle \nabla f(\mathbf{x}), \mathbf{y} \rangle_2 \quad \text{für alle } \mathbf{y} \in \mathbb{R}^n$$

gilt, mit dem sich also alle Richtungsableitungen von  $f$  durch ein Skalarprodukt beschreiben lassen. Es lässt sich nachweisen, dass die Funktion  $f$  in Richtung des Gradienten am schnellsten ansteigt und in Richtung des negativen Gradienten am schnellsten abnimmt, also sind diese Richtungen für die Suche nach Minima und Maxima besonders interessant.

Da wir nach Minima und Maxima des Rayleigh-Quotienten suchen, bietet es sich deshalb an, seinen Gradienten zu bestimmen.

**Lemma 8.2 (Gradient)** *Es gilt*

$$\nabla \Lambda_A(\mathbf{x}) = \frac{2}{\langle \mathbf{x}, \mathbf{x} \rangle_2} (\mathbf{A}\mathbf{x} - \Lambda_A(\mathbf{x})\mathbf{x}) \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}.$$

*Insbesondere ist der Gradient  $\nabla \Lambda_A(\mathbf{x})$  genau dann gleich null, wenn  $\mathbf{x}$  ein Eigenvektor von  $\mathbf{A}$  ist.*

*Beweis.* Wir führen die Hilfsfunktionen

$$f(\mathbf{x}) = \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle_2, \quad g(\mathbf{x}) = \langle \mathbf{x}, \mathbf{x} \rangle_2 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n$$

ein und erhalten mit Hilfe der Produktregel die Gleichungen

$$\begin{aligned}Df(\mathbf{x}) \cdot \mathbf{y} &= \langle \mathbf{y}, \mathbf{A}\mathbf{x} \rangle_2 + \langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle_2 = \langle \mathbf{y}, \mathbf{A}\mathbf{x} \rangle_2 + \langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle_2 = 2\langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle_2, \\ Dg(\mathbf{x}) \cdot \mathbf{y} &= 2\langle \mathbf{x}, \mathbf{y} \rangle_2 \quad \text{für alle } \mathbf{x}, \mathbf{y} \in \mathbb{K}^n.\end{aligned}$$

Aus  $\Lambda_A(\mathbf{x}) = f(\mathbf{x})/g(\mathbf{x})$  und der Quotientenregel folgt

$$\begin{aligned}D\Lambda_A(\mathbf{x}) \cdot \mathbf{y} &= \frac{g(\mathbf{x})Df(\mathbf{x}) \cdot \mathbf{y} - f(\mathbf{x})Dg(\mathbf{x}) \cdot \mathbf{y}}{g^2(\mathbf{x})} = \frac{2\langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle_2}{\langle \mathbf{x}, \mathbf{x} \rangle_2} - \frac{f(\mathbf{x})}{g(\mathbf{x})} \frac{2\langle \mathbf{x}, \mathbf{y} \rangle_2}{\langle \mathbf{x}, \mathbf{x} \rangle_2} \\ &= \frac{2}{\langle \mathbf{x}, \mathbf{x} \rangle_2} \langle \mathbf{A}\mathbf{x} - \Lambda_A(\mathbf{x})\mathbf{x}, \mathbf{y} \rangle_2 \quad \text{für alle } \mathbf{y} \in \mathbb{K}^n, \mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\},\end{aligned}$$

und damit die gesuchte Gleichung. ■



Wir beginnen mit der Suche nach einem Minimum mit einem Startvektor  $\mathbf{x}^{(0)} \in \mathbb{K}^n$ . Um die nächste Iterierte zu berechnen, bestimmen wir den Gradienten des Rayleigh-Quotienten  $\Lambda_A$  in  $\mathbf{x}^{(m)}$  und setzen

$$\begin{aligned} \mathbf{x}^{(m+1)} &:= \mathbf{x}^{(m)} - \alpha_m \nabla \Lambda_A(\mathbf{x}^{(m)}) = \mathbf{x}^{(m)} + \frac{2\alpha_m \Lambda_A(\mathbf{x}^{(m)})}{\langle \mathbf{x}^{(m)}, \mathbf{x}^{(m)} \rangle_2} \mathbf{x}^{(m)} - \frac{2\alpha_m}{\langle \mathbf{x}^{(m)}, \mathbf{x}^{(m)} \rangle_2} \mathbf{A} \mathbf{x}^{(m)} \\ &\in \text{span}\{\mathbf{x}^{(m)}, \mathbf{A} \mathbf{x}^{(m)}\} \quad \text{für alle } m \in \mathbb{N}_0 \end{aligned}$$

mit einem geeigneten Skalierungsfaktor  $\alpha_m$ . Die nächste Iterierte liegt also im Aufspann von  $\mathbf{x}^{(m)}$  und  $\mathbf{A} \mathbf{x}^{(m)}$ . Mit einer einfachen Induktion folgt

$$\mathbf{x}^{(m)} \in \text{span}\{\mathbf{x}^{(0)}, \mathbf{A} \mathbf{x}^{(0)}, \dots, \mathbf{A}^m \mathbf{x}^{(0)}\} \quad \text{für alle } m \in \mathbb{N}_0.$$

Die derart durch das wiederholte Multiplizieren eines Startvektors mit der Matrix  $\mathbf{A}$  definierten Räume sind für uns von besonderem Interesse.

**Definition 8.3 (Krylow-Raum)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$ ,  $\mathbf{q}^{(0)} \in \mathbb{K}^n$  und  $m \in \mathbb{N}_0$ . Der Raum

$$\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) := \text{span}\{\mathbf{A}^i \mathbf{q}^{(0)} : i \in [0 : m]\}$$

wird als  $m$ -ter Krylow-Raum zu der Matrix  $\mathbf{A}$  und dem Startvektor  $\mathbf{q}^{(0)}$  bezeichnet.

**Lemma 8.4 (Krylow-Raum)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$ ,  $\mathbf{q}^{(0)} \in \mathbb{K}^n$  und  $m \in \mathbb{N}_0$ .

Es gelten  $\dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) \leq m + 1$  und

$$\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) = \{p(\mathbf{A})\mathbf{q}^{(0)} : p \in \Pi_m\}.$$

Für jeden Vektor  $\mathbf{x} \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m)$  gilt  $\mathbf{A} \mathbf{x} \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m + 1)$ .

*Beweis.* Die Ungleichung  $\dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) \leq m + 1$  folgt unmittelbar aus der Definition des Krylow-Raums, denn er wird von  $m + 1$  Vektoren aufgespannt.

Sei  $\mathbf{x} \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m)$ . Dann existieren  $\alpha_0, \dots, \alpha_m \in \mathbb{K}$  mit

$$\mathbf{x} = \sum_{i=0}^m \alpha_i \mathbf{A}^i \mathbf{q}^{(0)}.$$

Einerseits folgt unmittelbar

$$\mathbf{A} \mathbf{x} = \mathbf{A} \sum_{i=0}^m \alpha_i \mathbf{A}^i \mathbf{q}^{(0)} = \sum_{i=0}^m \alpha_i \mathbf{A}^{i+1} \mathbf{q}^{(0)} \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m + 1),$$

andererseits erhalten wir mit dem durch

$$p(t) := \sum_{i=0}^m \alpha_i t^i \quad \text{für alle } t \in \mathbb{K}$$

definierten Polynom  $p \in \Pi_m$  die Gleichung

$$p(\mathbf{A})\mathbf{q}^{(0)} = \left( \sum_{i=0}^m \alpha_i \mathbf{A}^i \right) \mathbf{q}^{(0)} = \sum_{i=0}^m \alpha_i \mathbf{A}^i \mathbf{q}^{(0)} = \mathbf{x},$$

also die letzte noch fehlende Gleichung. ■

**Übungsaufgabe 8.5 (Invariante Teilräume)** Sei  $n \in \mathbb{N}$ . Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und  $\mathbf{q}^{(0)} \in \mathbb{K}^n$  gegeben.

- (a) Sei  $\mathcal{V} \subseteq \mathbb{K}^n$  ein bezüglich  $\mathbf{A}$  invarianter Unterraum der Dimension  $k \in [0 : n]$ , und gelte  $\mathbf{q}^{(0)} \in \mathcal{V}$ . Beweisen Sie

$$\dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) \leq k \quad \text{für alle } m \in \mathbb{N}_0.$$

- (b) Seien  $k \in [0 : n]$  und  $\mathbf{q}^{(0)} \in \mathbb{K}^n$  gegeben mit

$$\dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) \leq k \quad \text{für alle } m \in \mathbb{N}_0.$$

Beweisen Sie, dass ein bezüglich  $\mathbf{A}$  invarianter Unterraum  $\mathcal{V}$  mit  $\mathbf{q}^{(0)} \in \mathcal{V}$  und  $\dim \mathcal{V} \leq k$  existiert.

**Übungsaufgabe 8.6 (Approximation)** Sei  $n \in \mathbb{N}_{>1}$ . Sei  $\mathbf{F} \in \mathbb{K}^{n \times n}$  wie in Übungsaufgabe 6.14 gegeben durch

$$f_{ij} = \begin{cases} 1 & \text{falls } i - j \in \{1, 1 - n\}, \\ 0 & \text{ansonsten} \end{cases} \quad \text{für alle } i, j \in [1 : n].$$

Seien  $\alpha_0 \geq \alpha_1 \geq \dots \geq \alpha_{n-1} = 0$  gegeben. Zeigen Sie, dass ein  $\mathbf{b} \in \mathbb{K}^n$  existiert mit

$$\min\{\|\mathbf{b} - \mathbf{y}\| : \mathbf{y} \in \mathcal{K}(\mathbf{F}, \delta^{(1)}, m)\} = \alpha_m \quad \text{für alle } m \in [0 : n - 1].$$

Hinweis: Welche Komponenten eines Vektors in  $\mathcal{K}(\mathbf{F}, \delta^{(1)}, m)$  können ungleich null sein?

### 8.3 Arnoldi-Basis

Der Vorteil der Krylow-Räume besteht darin, dass sie sich relativ einfach mit Hilfe von Matrix-Vektor-Multiplikationen konstruieren lassen. Außerdem besteht, wie bereits gesehen, die Hoffnung, dass sie geeignet sind, um den Rayleigh-Quotienten zu minimieren beziehungsweise zu maximieren und so Approximationen des kleinsten beziehungsweise größten Eigenwerts zu erhalten.

Allerdings ist dafür die kanonische Basis  $\mathbf{q}^{(0)}, \mathbf{A}\mathbf{q}^{(0)}, \dots, \mathbf{A}^m\mathbf{q}^{(0)}$  in der Regel nicht gut geeignet: Für große Werte von  $m$  ist zu befürchten, dass die Basisvektoren, wie schon bei der Vektoriteration gesehen, gegen einen Eigenraum konvergieren und damit „numerisch linear abhängig“ werden. Die Lösung besteht, wie schon bei der orthogonalen Iteration, darin, zu einer orthonormalen Basis zu wechseln. Dafür sollten wir zunächst die Dimension des Raums bestimmen.

**Lemma 8.7 (Maximale Dimension)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und  $\mathbf{q}^{(0)} \in \mathbb{K}^n$ . Wir bezeichnen mit

$$m_0 := \max\{\dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) : m \in \mathbb{N}_0\}$$

die maximale Dimension, die ein Krylow-Raum mit dem Anfangsvektor  $\mathbf{q}^{(0)}$  erreichen kann. Dann gilt

$$m_0 = \min\{m \in \mathbb{N}_0 : \dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) \leq m\}.$$

Falls  $\mathbf{q}^{(0)} \neq \mathbf{0}$  gilt, folgt daraus insbesondere  $\dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m_0 - 1) = m_0$ .

*Beweis.* Aus  $\mathbf{q}^{(0)} = \mathbf{0}$  folgt  $\dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) = 0$  für alle  $m \in \mathbb{N}_0$ , so dass die Aussage trivial ist.

Gelte also nun  $\mathbf{q}^{(0)} \neq \mathbf{0}$ . Wir setzen

$$m_1 := \min\{m \in \mathbb{N}_0 : \dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) \leq m\}$$

und beweisen  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) = \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m_1 - 1)$  für alle  $m \in \mathbb{N}_0$  mit  $m \geq m_1$  per Induktion. Da  $m_1$  minimal gewählt ist, gilt  $\dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m_1 - 1) = m_1$ , und damit folgt die Behauptung.

*Induktionsanfang:* Sei  $m = m_1$ . Da  $m_1$  minimal gewählt wurde, muss

$$\dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m - 1) = m \geq \dim \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m)$$

gelten, und aus  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m - 1) \subseteq \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m)$  folgt  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m - 1) = \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m)$ .

*Induktionsvoraussetzung:* Sei  $m \in \mathbb{N}$  mit  $m \geq m_1$  so gegeben, dass  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) = \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m_1 - 1)$  gilt.

*Induktionsschritt:* Nach Induktionsvoraussetzung gilt  $\mathbf{A}^m \mathbf{q}^{(0)} \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) = \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m_1 - 1)$ . Mit Lemma 8.4 und dem Induktionsanfang folgt

$$\mathbf{A}^{m+1} \mathbf{q}^{(0)} = \mathbf{A}(\mathbf{A}^m \mathbf{q}^{(0)}) \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m_1) = \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m_1 - 1).$$

Da  $\mathbf{q}^{(0)}, \mathbf{A} \mathbf{q}^{(0)}, \dots, \mathbf{A}^m \mathbf{q}^{(0)} \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m_1 - 1)$  bereits nach Induktionsvoraussetzung gilt, folgt  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m + 1) \subseteq \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m_1 - 1)$ , und aus  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m_1 - 1) \subseteq \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m + 1)$  auch die Gleichheit der Räume. ■

Da sich die kanonische Basis des Krylow-Raums  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m + 1)$  von der des Raums  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m)$  nur durch den Vektor  $\mathbf{A}^{m+1} \mathbf{q}^{(0)}$  unterscheidet, bietet es sich an, auch orthonormale Basen zu verwenden, die sich lediglich durch einen Vektor unterscheiden. Wir suchen also orthonormale Vektoren  $\mathbf{q}^{(0)}, \mathbf{q}^{(1)}, \dots, \mathbf{q}^{(m_0-1)} \in \mathbb{K}^n$  derart, dass

$$\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) = \text{span}\{\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m)}\} \quad \text{für alle } m \in [0 : m_0 - 1] \quad (8.2)$$

gilt. Aus dieser Gleichung können wir direkt eine induktive Konstruktion für die Vektoren  $\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m_0-1)}$  gewinnen: Wir gehen davon aus, dass ein  $m \in [0 : m_0 - 1]$  so gegeben ist, dass  $\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m)}$  eine Basis des Krylow-Raums  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m)$  ist. Dann spannen die Vektoren

$$\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m)}, \mathbf{A}^{m+1} \mathbf{q}^{(0)}$$

den Raum  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m + 1)$  auf. Da wir damit rechnen müssen, dass die Vektoren  $\mathbf{A}^{m+1} \mathbf{q}^{(0)}$  gegen einen Eigenraum konvergieren und damit die Basis „numerisch linear abhängig“ wird, soll nun  $\mathbf{A}^{m+1} \mathbf{q}^{(0)}$  durch einen „numerisch stabileren“ Vektor ersetzt werden: Wegen

$$\mathbf{A}^m \mathbf{q}^{(0)} \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) = \text{span}\{\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m)}\}$$

existieren  $\alpha_0, \dots, \alpha_m \in \mathbb{K}$  mit

$$\mathbf{A}^m \mathbf{q}^{(0)} = \alpha_0 \mathbf{q}^{(0)} + \dots + \alpha_m \mathbf{q}^{(m)},$$

also folgt

$$\mathbf{A}^{m+1}\mathbf{q}^{(0)} = \alpha_0\mathbf{A}\mathbf{q}^{(0)} + \dots + \alpha_m\mathbf{A}\mathbf{q}^{(m)}.$$

Nach Lemma 8.4 haben wir

$$\alpha_0\mathbf{A}\mathbf{q}^{(0)} + \dots + \alpha_{m-1}\mathbf{A}\mathbf{q}^{(m-1)} \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m),$$

also folgt

$$\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m+1) = \text{span}\{\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m)}, \mathbf{A}\mathbf{q}^{(m)}\}$$

mit dem „numerisch stabileren“ Vektor  $\mathbf{A}\mathbf{q}^{(m)}$  anstelle von  $\mathbf{A}^{m+1}\mathbf{q}^{(0)}$ .

Da nach Lemma 8.7 für alle  $m \in [0 : m_0 - 2]$  der Krylow-Raum  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m+1)$  die Dimension  $m+2$  hat, müssen die Vektoren  $\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m)}, \mathbf{A}\mathbf{q}^{(m)}$  linear unabhängig sein. Insbesondere kann  $\mathbf{A}\mathbf{q}^{(m)}$  nicht der Nullvektor sein.

Beispielsweise per Gram-Schmidt-Orthonormalisierung können wir aus ihm deshalb den gewünschten Vektor  $\mathbf{q}^{(m+1)}$  konstruieren.

**Definition 8.8 (Arnoldi-Basis)** Sei  $\mathbf{q}^{(0)} \in \mathbb{K}^n$  mit  $\|\mathbf{q}^{(0)}\| = 1$  gegeben. Die Arnoldi-Basis  $\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m_0-1)}$  ist durch

$$\begin{aligned} \mathbf{p}^{(m+1)} &:= \mathbf{A}\mathbf{q}^{(m)} - \sum_{i=0}^m \langle \mathbf{q}^{(i)}, \mathbf{A}\mathbf{q}^{(m)} \rangle \mathbf{q}^{(i)} && \text{für alle } m \in [0 : m_0 - 1], \\ \mathbf{q}^{(m+1)} &:= \frac{\mathbf{p}^{(m+1)}}{\|\mathbf{p}^{(m+1)}\|} && \text{für alle } m \in [0 : m_0 - 2] \end{aligned}$$

definiert. Nach Definition von  $m_0$  gelten  $\mathbf{p}^{(m)} \neq \mathbf{0}$  für alle  $m \in [1 : m_0 - 1]$  sowie  $\mathbf{p}^{(m_0)} = \mathbf{0}$ , also ist die Arnoldi-Basis wohldefiniert.

Nach Konstruktion ist die Basis auch orthonormal und besitzt die Eigenschaft (8.2).

Wir sind daran interessiert, die kleinsten und größten Eigenwerte zu approximieren. Falls  $\mathbf{A}$  selbstadjungiert ist, ist kleinste Eigenwert  $\lambda_1$  das Minimum des Rayleigh-Quotienten  $\Lambda_A$  auf dem Raum  $\mathbb{K}^n \setminus \{\mathbf{0}\}$ , und indem wir diesen Raum durch die Teilräume

$$\mathcal{Q}^{(m)} := \text{span}\{\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m)}\} \quad \text{für alle } m \in [0 : m_0 - 1]$$

ersetzen, können wir Approximationen

$$\lambda_1^{(m)} := \min\{\Lambda_A(\mathbf{x}) : \mathbf{x} \in \mathcal{Q}^{(m)} \setminus \{\mathbf{0}\}\} \quad \text{für alle } m \in [0 : m_0 - 1]$$

konstruieren. Die Suche nach dem Minimum in einem Teilraum ist natürlich etwas unhandlich, deshalb führen wir die Berechnung auf Größen zurück, mit denen wir besser arbeiten können.

Wir fassen die Basisvektoren zu isometrischen Matrizen

$$\mathbf{Q}^{(m)} := (\mathbf{q}^{(0)} \quad \dots \quad \mathbf{q}^{(m)}) \in \mathbb{K}^{n \times [0:m]} \quad \text{für alle } m \in [0 : m_0 - 1] \quad (8.3)$$

zusammen, wobei die Notation  $\mathbb{K}^{n \times [0:m]}$  betonen soll, dass die Spalten der Matrix von 0 bis  $m$  statt von 1 bis  $m+1$  nummeriert sind, und stellen fest, dass für  $m \in [0 : m_0 - 1]$  die Gleichung

$$\begin{aligned} \lambda_1^{(m)} &= \min\{\Lambda_A(\mathbf{x}) : \mathbf{x} \in \mathcal{Q}^{(m)} \setminus \{\mathbf{0}\}\} \\ &= \min\{\Lambda_A(\mathbf{Q}^{(m)}\widehat{\mathbf{x}}) : \widehat{\mathbf{x}} \in \mathbb{K}^{[0:m]} \setminus \{\mathbf{0}\}\} \\ &= \min\left\{\frac{\langle \mathbf{Q}^{(m)}\widehat{\mathbf{x}}, \mathbf{A}\mathbf{Q}^{(m)}\widehat{\mathbf{x}} \rangle_2}{\langle \mathbf{Q}^{(m)}\widehat{\mathbf{x}}, \mathbf{Q}^{(m)}\widehat{\mathbf{x}} \rangle_2} : \widehat{\mathbf{x}} \in \mathbb{K}^{[0:m]} \setminus \{\mathbf{0}\}\right\} \\ &= \min\left\{\frac{\langle \widehat{\mathbf{x}}, (\mathbf{Q}^{(m)})^* \mathbf{A}\mathbf{Q}^{(m)}\widehat{\mathbf{x}} \rangle_2}{\langle \widehat{\mathbf{x}}, (\mathbf{Q}^{(m)})^* \mathbf{Q}^{(m)}\widehat{\mathbf{x}} \rangle_2} : \widehat{\mathbf{x}} \in \mathbb{K}^{[0:m]} \setminus \{\mathbf{0}\}\right\} \\ &= \min\{\Lambda_{(\mathbf{Q}^{(m)})^* \mathbf{A}\mathbf{Q}^{(m)}}(\widehat{\mathbf{x}}) : \widehat{\mathbf{x}} \in \mathbb{K}^{[0:m]} \setminus \{\mathbf{0}\}\} \end{aligned}$$

gilt. Nach dem Satz 3.45 von Courant und Fischer sind also die Werte  $\lambda_1^{(m)}$  gerade die kleinsten Eigenwerte der Matrizen

$$\widehat{\mathbf{A}}^{(m)} := (\mathbf{Q}^{(m)})^* \mathbf{A}\mathbf{Q}^{(m)} \in \mathbb{K}^{[0:m] \times [0:m]} \quad \text{für alle } m \in [0 : m_0 - 1], \quad (8.4)$$

können also durch Lösen eines  $(m+1)$ -dimensionalen Eigenwertproblems berechnet werden. Für  $m \ll n$  lässt sich diese Berechnung wesentlich effizienter durchführen als für die ursprüngliche Matrix  $\mathbf{A}$ . Durch die folgende Beobachtung wird die Berechnung sogar noch weiter vereinfacht:

**Lemma 8.9 (Hessenberg-Form)** *Für jedes  $m \in [0 : m_0 - 1]$  besitzt die Matrix  $\widehat{\mathbf{A}}^{(m)}$  Hessenberg-Gestalt, es gilt*

$$\hat{a}_{ij} = \begin{cases} \langle \mathbf{q}^{(i)}, \mathbf{A}\mathbf{q}^{(j)} \rangle & \text{falls } i \leq j + 1, \\ 0 & \text{ansonsten} \end{cases} \quad \text{für alle } i, j \in [0 : m].$$

*Beweis.* Seien  $m \in [0 : m_0 - 1]$  und  $i, j \in [0 : m]$ .

Aus der Gleichung (8.4) folgt mit der Definition der Matrix  $\mathbf{Q}^{(m)}$  direkt die Gleichung  $\hat{a}_{ij} = \langle \mathbf{q}^{(i)}, \mathbf{A}\mathbf{q}^{(j)} \rangle$ .

Mit Lemma 8.4 folgt aus  $\mathbf{q}^{(j)} \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, j)$  unmittelbar

$$\mathbf{A}\mathbf{q}^{(j)} \in \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, j+1) = \text{span}\{\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(j+1)}\}.$$

Falls  $i > j + 1$  gilt, steht  $\mathbf{q}^{(i)}$  senkrecht auf den Vektoren  $\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(j+1)}$ , also auch auf  $\mathbf{A}\mathbf{q}^{(j)}$ , und es folgt  $\hat{a}_{ij} = \langle \mathbf{q}^{(i)}, \mathbf{A}\mathbf{q}^{(j)} \rangle = 0$ . ■

Die Berechnung der Einträge der Matrix  $\widehat{\mathbf{A}}^{(m)}$  kann besonders elegant gestaltet werden: Die Einträge  $\hat{a}_{ij} = \langle \mathbf{q}^{(i)}, \mathbf{A}\mathbf{q}^{(j)} \rangle$  für  $i \in [0 : j]$  werden bei der Gram-Schmidt-Orthonormalisierung ohnehin berechnet, wir brauchen sie also lediglich für die spätere Verwendung abzuspeichern.

Der untere Nebendiagonaleintrag  $\hat{a}_{i+1,i}$  erfüllt nach Definition, und weil  $\mathbf{q}^{(i+1)}$  senkrecht auf  $\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(i)}$  steht, die Gleichung

$$\begin{aligned}\hat{a}_{i+1,i} &= \langle \mathbf{q}^{(i+1)}, \mathbf{A}\mathbf{q}^{(i)} \rangle = \langle \mathbf{q}^{(i+1)}, \mathbf{p}^{(i+1)} + \sum_{\ell=0}^i \hat{a}_{\ell i} \mathbf{q}^{(\ell)} \rangle \\ &= \langle \mathbf{q}^{(i+1)}, \mathbf{p}^{(i+1)} \rangle = \|\mathbf{p}^{(i+1)}\| \quad \text{für alle } i \in [0 : m-1], m \in [0 : m_0-1],\end{aligned}$$

und diese Norm wird bei der Orthonormalisierung ebenfalls bereits berechnet, so dass wir sie nur aufzubewahren brauchen. Insbesondere sind diese Nebendiagonalelemente alle ungleich null, so dass die entstehende Hessenberg-Matrix irreduzibel ist.

Für eine praktische Konstruktion der Arnoldi-Basis wäre es von Vorteil, wenn wir auch den Parameter  $m_0$  berechnen könnten. Nach Definition 8.8 gilt  $\mathbf{p}^{(m_0)} = \mathbf{0}$ , also können wir diesen Vektor verwenden, um  $m_0$  zu ermitteln. Wegen Rundungsfehlern dürfen wir in der Praxis nicht auf einen exakten Nullvektor hoffen. Stattdessen verwenden wir ein Kriterium der Form

$$\|\mathbf{p}^{(m+1)}\|_2 \leq \epsilon_{\text{ir}} \|\mathbf{A}\mathbf{q}^{(m)}\|_2$$

mit einem  $\epsilon_{\text{ir}} \in \mathbb{R}_{>0}$ : Mit der Abkürzung  $\mathbf{a}^{(m)} := \mathbf{A}\mathbf{q}^{(m)}$  gilt die Gleichung

$$\mathbf{p}^{(m+1)} = \mathbf{A}\mathbf{q}^{(m)} - \sum_{i=0}^m \langle \mathbf{q}^{(i)}, \mathbf{A}\mathbf{q}^{(m)} \rangle \mathbf{q}^{(i)} = \mathbf{a}^{(m)} - \mathbf{Q}^{(m)}(\mathbf{Q}^{(m)})^* \mathbf{a}^{(m)},$$

und da  $\mathbf{Q}^{(m)}(\mathbf{Q}^{(m)})^*$  die orthogonale Projektion auf  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m)$  ist, folgt

$$\sin \angle(\mathbf{a}^{(m)}, \mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m)) = \frac{\|\mathbf{a}^{(m)} - \mathbf{Q}^{(m)}(\mathbf{Q}^{(m)})^* \mathbf{a}^{(m)}\|}{\|\mathbf{a}^{(m)}\|} = \frac{\|\mathbf{p}^{(m)}\|}{\|\mathbf{A}\mathbf{q}^{(m)}\|},$$

unser Kriterium begrenzt also gerade den Winkel zwischen dem neu hinzugekommenen Vektor und dem vorigen Krylow-Raum.

**Algorithmus 8.10 (Arnoldi-Basis)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  und  $\mathbf{q}^{(0)} \in \mathbb{K}^n$  mit  $\|\mathbf{q}^{(0)}\|_2 = 1$  gegeben. Der folgende Algorithmus berechnet die Arnoldi-Basis  $\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m_0-1)}$  und die gemäß (8.4) definierten Matrizen  $\hat{\mathbf{A}}^{(m)}$ . Der Algorithmus endet mit  $m = m_0$ .

```

 $\mathbf{p} \leftarrow \mathbf{A}\mathbf{q}^{(0)}; \quad \gamma \leftarrow \|\mathbf{p}\|$ 
 $\hat{a}_{00} \leftarrow \langle \mathbf{q}^{(0)}, \mathbf{p} \rangle; \quad \mathbf{p} \leftarrow \mathbf{p} - \hat{a}_{00}\mathbf{q}^{(0)}$ 
 $\hat{a}_{10} \leftarrow \|\mathbf{p}\|$ 
 $m \leftarrow 1$ 
while  $\hat{a}_{m,m-1} > \epsilon_{\text{ir}}\gamma$  do begin
   $\mathbf{q}^{(m)} \leftarrow \mathbf{p}/\hat{a}_{m,m-1}$ 
   $\mathbf{p} \leftarrow \mathbf{A}\mathbf{q}^{(m)}; \quad \gamma \leftarrow \|\mathbf{p}\|$ 
  for  $i \in [0 : m]$  do begin
     $\hat{a}_{im} \leftarrow \langle \mathbf{q}^{(i)}, \mathbf{p} \rangle; \quad \mathbf{p} \leftarrow \mathbf{p} - \hat{a}_{im}\mathbf{q}^{(i)}$ 
  end
   $\hat{a}_{m+1,m} \leftarrow \|\mathbf{p}\|$ 
   $m \leftarrow m + 1$ 
end

```

Falls  $\mathbf{A}$  selbstadjungiert ist, gilt dasselbe nach Konstruktion auch für  $\widehat{\mathbf{A}}^{(m)}$ , und da diese Matrix auch eine Hessenberg-Matrix ist, muss sie dann sogar tridiagonal sein, so dass sich der Rechenaufwand für die Orthonormalisierung und die Bestimmung von  $\lambda_1^{(m)}$  noch weiter reduzieren lässt: Es gilt

$$\hat{a}_{im} = \overline{\hat{a}_{mi}} = 0 \quad \text{für alle } i \in [0 : m - 2],$$

so dass wir die meisten der für die Orthogonalisierung erforderlichen Skalarprodukte einsparen können. Wenn wir die Tridiagonalmatrix  $\widehat{\mathbf{A}}^{(m)}$  in der bereits aus Kapitel 7 bekannten Form

$$\widehat{\mathbf{A}}^{(m)} = \begin{pmatrix} \alpha_1 & \bar{\beta}_1 & & & \\ \beta_1 & \alpha_2 & \ddots & & \\ & \ddots & \ddots & \ddots & \\ & & \ddots & \ddots & \bar{\beta}_m \\ & & & \beta_m & \alpha_{m+1} \end{pmatrix}$$

darstellen, ergibt sich

$$\hat{a}_{m-1,m} = \overline{\hat{a}_{m,m-1}} = \bar{\beta}_{m-1},$$

so dass ein weiteres Skalarprodukt entfällt. Der resultierende sehr effiziente *Lanczos-Algorithmus* nimmt damit die folgende Form an:

**Algorithmus 8.11 (Lanczos-Algorithmus)** Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  selbstadjungiert und  $\mathbf{q}^{(0)} \in \mathbb{K}^n$  mit  $\|\mathbf{q}^{(0)}\| = 1$  gegeben. Der folgende Algorithmus berechnet die Arnoldi-Basis  $\mathbf{q}^{(0)}, \dots, \mathbf{q}^{(m_0-1)}$  und die Matrizen  $\widehat{\mathbf{A}}^{(m)}$ . Der Algorithmus endet mit  $m = m_0$ .

```

 $\mathbf{p} \leftarrow \mathbf{A}\mathbf{q}^{(0)}; \quad \gamma \leftarrow \|\mathbf{p}\|$ 
 $\alpha_1 \leftarrow \langle \mathbf{q}^{(0)}, \mathbf{p} \rangle; \quad \mathbf{p} \leftarrow \mathbf{p} - \alpha_1 \mathbf{q}^{(0)}$ 
 $\beta_1 \leftarrow \|\mathbf{p}\|$ 
 $m \leftarrow 1$ 
while  $\beta_m > \epsilon_{ir} \gamma$  do begin
   $\mathbf{q}^{(m)} \leftarrow \mathbf{p} / \beta_m$ 
   $\mathbf{p} \leftarrow \mathbf{A}\mathbf{q}^{(m)}; \quad \gamma \leftarrow \|\mathbf{p}\|$ 
   $\alpha_{m+1} \leftarrow \langle \mathbf{q}^{(m)}, \mathbf{p} \rangle$ 
   $\mathbf{p} \leftarrow \mathbf{p} - \bar{\beta}_m \mathbf{q}^{(m-1)} - \alpha_{m+1} \mathbf{q}^{(m)}$ 
   $\beta_{m+1} \leftarrow \|\mathbf{p}\|$ 
   $m \leftarrow m + 1$ 
end

```

Da der Algorithmus abbricht, sobald  $\beta_m = 0$  gilt, berechnet er nicht nur eine selbstadjungierte Tridiagonalmatrix, sondern sogar eine *irreduzible* selbstadjungierte Tridiagonalmatrix. Deshalb lässt sich die Theorie aus Kapitel 7 anwenden, um beispielsweise die Eigenwerte mit Hilfe Sturmscher Ketten zu berechnen.

Es muss leider darauf hingewiesen werden, dass in der Praxis Rundungsfehler häufig zu einem Verlust der Orthogonalität der mit dem Lanczos-Verfahren berechneten Basis führen. Dadurch kann es beispielsweise passieren, dass  $\widehat{\mathbf{A}}^{(m)}$  bestimmte Eigenwerte bis auf Rundungsfehler mehrfach aufweist, die in  $\mathbf{A}$  nur einfach auftreten.

## 8.4 Konvergenz

Bei der folgenden Untersuchung beschränken wir uns auf den Fall selbstadjungierter Matrizen, setzen also  $\mathbf{A}^* = \mathbf{A}$  voraus. Mit Satz 3.47 finden wir eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  und eine reelle Diagonalmatrix

$$\mathbf{D} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}, \quad \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

derart, dass

$$\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^*$$

gilt. Nach Lemma 8.4 gilt

$$\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m) = \{p(\mathbf{A})\mathbf{q}^{(0)} : p \in \Pi_m\},$$

wobei  $\Pi_m$  den Raum der Polynome mit Grad  $\leq m$  bezeichnet.

Da die Matrizen  $\mathbf{Q}$  unitär sind, gilt  $\mathbf{Q}^*\mathbf{Q} = \mathbf{I}$ , und wir erhalten

$$p(\mathbf{A}) = p(\mathbf{Q}\mathbf{D}\mathbf{Q}^*) = \mathbf{Q}p(\mathbf{D})\mathbf{Q}^* \quad \text{für alle } p \in \Pi_m, m \in \mathbb{N}_0. \quad (8.5)$$

Diese Gleichung ist für unsere Konvergenzanalyse von zentraler Bedeutung, denn sie ermöglicht es uns, durch Wahl geeigneter Polynome  $p \in \Pi_m$  bestimmte Eigenwerte in  $p(\mathbf{D})$  zu vergrößern oder zu verkleinern und so Approximationen zu konstruieren.

**Lemma 8.12 (Eigenvektor-Approximation)** Sei  $\mathbf{q}^{(0)} \in \mathbb{K}^n$  mit  $\|\mathbf{q}^{(0)}\| = 1$  gegeben und sei  $\mathbf{e} := \mathbf{Q}\delta^{(1)}$ . Sei  $p \in \Pi_m$  ein Polynom mit

$$p(\lambda_1) = 1, \quad |p(x)| \leq \epsilon \quad \text{für alle } x \in \{\lambda_2, \dots, \lambda_n\}.$$

Dann gilt

$$\tan \angle(p(\mathbf{A})\mathbf{q}^{(0)}, \mathbf{e}) \leq \epsilon \tan \angle(\mathbf{q}^{(0)}, \mathbf{e}).$$

*Beweis.* Mit (8.5) und  $\hat{\mathbf{q}}^{(0)} := \mathbf{Q}^*\mathbf{q}^{(0)}$  folgt

$$\begin{aligned} \tan^2 \angle(p(\mathbf{A})\mathbf{q}^{(0)}, \mathbf{e}) &= \tan^2 \angle(\mathbf{Q}p(\mathbf{D})\mathbf{Q}^*\mathbf{q}^{(0)}, \mathbf{Q}\delta^{(1)}) \\ &= \tan^2 \angle(p(\mathbf{D})\hat{\mathbf{q}}^{(0)}, \delta^{(1)}) = \frac{\sum_{i=2}^n |p(\lambda_i)|^2 |\hat{q}_i^{(0)}|^2}{|p(\lambda_1)|^2 |\hat{q}_1^{(0)}|^2} \\ &\leq \frac{\epsilon^2 \sum_{i=2}^n |\hat{q}_i^{(0)}|^2}{|\hat{q}_1^{(0)}|^2} = \epsilon^2 \tan^2 \angle(\hat{\mathbf{q}}^{(0)}, \delta^{(1)}) \\ &= \epsilon^2 \tan^2 \angle(\mathbf{Q}\hat{\mathbf{q}}^{(0)}, \mathbf{Q}\delta^{(1)}) = \epsilon^2 \tan^2 \angle(\mathbf{q}^{(0)}, \mathbf{e}). \end{aligned}$$

Das ist die gewünschte Abschätzung. ■



**Bemerkung 8.13 (Weitere Eigenwerte)** Um einen Eigenvektor zu dem zweiten Eigenwert zu finden, können wir ein  $q \in \Pi_{m-1}$  mit  $q(\lambda_2) = 1$  suchen, das auf  $\{\lambda_3, \dots, \lambda_n\}$  im Betrag kleiner als ein  $\epsilon \in \mathbb{R}_{>0}$  ist.

Dann setzen wir  $p(x) := \frac{x-\lambda_1}{\lambda_2-\lambda_1}q(x)$  und stellen fest, dass  $p \in \Pi_m$  und

$$p(\lambda_1) = 0, \quad p(\lambda_2) = 1, \quad |p(x)| \leq \epsilon \frac{\lambda_n - \lambda_1}{\lambda_2 - \lambda_1} \quad \text{für alle } x \in \{\lambda_3, \dots, \lambda_n\}$$

gelten. Nun können wir wie in dem Beweis des Lemmas 8.12 vorgehen. Für weitere Eigenwerte können wir entsprechend argumentieren.

Da der Grad des Polynoms  $p$  unmittelbar mit der Anzahl der Schritte des Lanczos-Verfahrens zusammenhängt, sind wir daran interessiert, Polynome zu finden, die auf  $\{\lambda_2, \dots, \lambda_n\}$  möglichst geringe Werte annehmen, aber in  $\lambda_1$  gleich eins sind. Geeignet transformierte Tschebyscheff-Polynome erfüllen diese Anforderungen.

**Definition 8.14 (Tschebyscheff-Polynome)** Wir definieren durch

$$\begin{aligned} c_0 &: \mathbb{K} \rightarrow \mathbb{K}, & x &\mapsto 1, \\ c_1 &: \mathbb{K} \rightarrow \mathbb{K}, & x &\mapsto x, \\ c_{m+2} &: \mathbb{K} \rightarrow \mathbb{K}, & x &\mapsto 2x c_{m+1}(x) - c_m(x) \end{aligned} \quad \text{für alle } m \in \mathbb{N}_0,$$

die Folge der Tschebyscheff-Polynome.

Aus dieser Darstellung lässt sich direkt ableiten, dass  $c_m \in \Pi_m$  gilt, und die Rekurrenzrelation kann sich bei der Implementierung als nützlich erweisen. Für die theoretische Untersuchung sind alternative Darstellungen der Tschebyscheff-Polynome hilfreicher:

**Lemma 8.15 (Alternative Darstellungen)** Für alle  $x \in [-1, 1]$  und alle  $m \in \mathbb{N}_0$  gilt

$$c_m(x) = \cos(m \arccos x). \quad (8.6a)$$

Für alle  $x \in \mathbb{R} \setminus [-1, 1]$  und alle  $m \in \mathbb{N}_0$  gilt

$$c_m(x) = \frac{1}{2} \left( \left( x + \sqrt{x^2 - 1} \right)^m + \left( x + \sqrt{x^2 - 1} \right)^{-m} \right). \quad (8.6b)$$

*Beweis.* Wir verwenden die Kutta-Schukowski-Transformation

$$s: \mathbb{C} \setminus \{0\} \rightarrow \mathbb{C}, \quad z \mapsto \frac{z + 1/z}{2}.$$

Sie ist surjektiv, denn für jedes  $x \in \mathbb{C}$  gilt

$$s(z) = x \iff \frac{z + 1/z}{2} = x \iff z^2 - 2xz + 1 = 0,$$

und diese quadratische Gleichung ist in dem Körper  $\mathbb{C}$  der komplexen Zahlen immer lösbar. Falls  $z$  eine Lösung ist, gilt dasselbe aus Symmetriegründen auch für  $1/z$ .

## 8 Lanczos-Verfahren für schwachbesetzte Matrizen

Wir definieren für alle  $m \in \mathbb{N}_0$  die Funktionen

$$f_m: \mathbb{C} \setminus \{0\} \rightarrow \mathbb{C}, \quad z \mapsto \frac{z^m + z^{-m}}{2}$$

und werden die Gleichungen

$$c_m(x) = f_m(z) \quad \text{für alle } m \in \mathbb{N}_0, z \in \mathbb{C} \setminus \{0\}, x := s(z) \quad (8.7)$$

per abschnittsweiser Induktion beweisen.

*Induktionsanfang:* Für  $m = 0$  gilt  $f_m(z) = 1 = c_m(x)$  für alle  $z \in \mathbb{C} \setminus \{0\}$  und alle  $x \in \mathbb{C}$ . Für  $m = 1$  haben wir  $f_m = s$ , also folgt aus  $x = s(z)$  bereits  $c_1(x) = x = s(z) = f_m(z)$  für alle  $z \in \mathbb{C} \setminus \{0\}$ .

*Induktionsvoraussetzung:* Sei  $m \in \mathbb{N}_{\geq 1}$  so gegeben, dass die Gleichung  $c_n(x) = f_n(z)$  für alle  $n \in [0 : m]$ ,  $z \in \mathbb{C} \setminus \{0\}$ ,  $x = s(z)$  gilt.

*Induktionsschritt:* Sei  $z \in \mathbb{C} \setminus \{0\}$ , sei  $x := s(z)$ . Es gilt

$$\begin{aligned} f_{m+1}(z) &= \frac{z^{m+1} + \frac{1}{z^{m+1}}}{2} = \frac{(z + \frac{1}{z})z^m - z^{m-1} + (z + \frac{1}{z})\frac{1}{z^m} - \frac{1}{z^{m-1}}}{2} \\ &= \left(z + \frac{1}{z}\right) \frac{z^m + \frac{1}{z^m}}{2} - \frac{z^{m-1} + \frac{1}{z^{m-1}}}{2} = 2s(z) f_m(z) - f_{m-1}(z) \\ &= 2x c_m(x) - c_{m-1}(x) = c_{m+1}(x). \end{aligned}$$

Um (8.6a) zu beweisen, wählen wir  $x \in [-1, 1]$  und setzen  $\xi := \arccos(x) \in [0, \pi]$  sowie  $z := e^{i\xi}$ . Dann gilt mit der Eulerschen Formel

$$s(z) = \frac{z + 1/z}{2} = \frac{e^{i\xi} + e^{-i\xi}}{2} = \frac{e^{i\xi} + \overline{e^{i\xi}}}{2} = \operatorname{Re}(e^{i\xi}) = \cos(\xi) = x.$$

Mit (8.7) und der Eulerschen Formel folgt

$$\begin{aligned} c_m(x) = f_m(z) &= \frac{z^m + z^{-m}}{2} = \frac{e^{im\xi} + e^{-im\xi}}{2} \\ &= \frac{e^{im\xi} + \overline{e^{im\xi}}}{2} = \operatorname{Re}(e^{im\xi}) = \cos(m\xi) = \cos(m \arccos(x)). \end{aligned}$$

Um (8.6b) zu zeigen, wählen wir  $x \in \mathbb{R} \setminus [-1, 1]$  und setzen  $z := x + \sqrt{x^2 - 1}$ . Dann gilt

$$\frac{1}{z} = \frac{1}{x + \sqrt{x^2 - 1}} = \frac{x - \sqrt{x^2 - 1}}{(x + \sqrt{x^2 - 1})(x - \sqrt{x^2 - 1})} = \frac{x - \sqrt{x^2 - 1}}{x^2 - (x^2 - 1)} = x - \sqrt{x^2 - 1},$$

und es folgt

$$s(z) = \frac{z + 1/z}{2} = \frac{x + \sqrt{x^2 - 1} + x - \sqrt{x^2 - 1}}{2} = \frac{2x}{2} = x.$$

Indem wir in (8.7) einsetzen, folgt direkt (8.6b). ■

**Lemma 8.16 (Optimalität)** Sei  $x_0 \in \mathbb{R} \setminus [-1, 1]$ . Unter allen Polynomen  $p \in \Pi_m$  mit  $p(x_0) = c_m(x_0)$  besitzt  $c_m$  die kleinste Maximumnorm auf dem Intervall  $[-1, 1]$ .

*Beweis.* Die Maximumnorm von  $c_m$  auf dem Intervall  $[-1, 1]$  beträgt wegen  $c_m(x) = \cos(m \arccos x)$  gerade eins. Sei  $p \in \Pi_m$  ein Polynom mit  $\|p\|_{\infty, [-1, 1]} < 1$ .

Wir betrachten die Differenz  $r := c_m - p \in \Pi_m$ . Es gilt

$$c_m(\xi_\nu) = \cos(\pi\nu) = (-1)^\nu \quad \text{mit } \xi_\nu := \cos(\pi\nu/m) \quad \text{für alle } \nu \in [0 : m].$$

Aus  $\|p\|_{\infty, [-1, 1]} < 1$  folgt

$$\begin{aligned} r(\xi_\nu) &= (-1)^\nu - p(\xi_\nu) \geq 1 - |p(\nu)| > 1 - 1 = 0 && \text{für alle geraden } \nu \in [0 : m], \\ r(\xi_\nu) &= (-1)^\nu - p(\xi_\nu) \leq -1 + |p(\nu)| < -1 + 1 = 0 && \text{für alle ungeraden } \nu \in [0 : m]. \end{aligned}$$

Da das Polynom  $r$  in  $\xi_\nu$  alternierende Vorzeichen annimmt, muss es nach dem Zwischenwertsatz in jedem Intervall  $(\xi_{\nu+1}, \xi_\nu)$  eine Nullstelle besitzen. Da der Cosinus auf  $[0, \pi]$  streng monoton fällt, besitzt  $r$  also mindestens  $m$  Nullstellen in  $[-1, 1] = [\cos(\pi), \cos(0)]$ .

Da  $x_0 \notin [-1, 1]$  gilt, kann  $x_0$  keine dieser  $m$  Nullstellen sein. Würde also  $p(x_0) = c_m(x_0)$  gelten, besäße  $r \in \Pi_m$  mindestens  $m + 1$  Nullstellen, müsste also nach dem Eindeutigkeitsatz für Polynome gleich null sein, wir hätten  $p = c_m$ . Da  $\|p\|_{\infty, [-1, 1]} < 1 = \|c_m\|_{\infty, [-1, 1]}$  gilt, ist das unmöglich, also kann kein Polynom  $p \in \Pi_m$  mit  $p(x_0) = c_m(x_0)$  auf  $[-1, 1]$  eine Maximumnorm echt kleiner als  $1 = \|c_m\|_{\infty, [-1, 1]}$  aufweisen. ■

Um eine Approximation des kleinsten Eigenwerts  $\lambda_1$  zu erhalten, benötigen wir ein Polynom niedriger Ordnung, das auf dem Intervall  $[\lambda_2, \lambda_n]$  möglichst kleine Werte annimmt. Die Tschebyscheff-Polynome sind auf  $[-1, 1]$  in der oben präzisierten Weise minimal, also bietet es sich an, sie so zu transformieren, dass diese Minimalität sich auf  $[\lambda_2, \lambda_n]$  überträgt.

Die Abbildung

$$\Phi : \mathbb{R} \rightarrow \mathbb{R}, \quad x \mapsto \frac{\lambda_n - x}{\lambda_n - \lambda_2} + \frac{x - \lambda_2}{\lambda_n - \lambda_2}(-1) = \frac{\lambda_n + \lambda_2 - 2x}{\lambda_n - \lambda_2}$$

erfüllt  $\Phi(\lambda_2) = 1$  und  $\Phi(\lambda_n) = -1$  sowie

$$x_0 := \Phi(\lambda_1) = \frac{\lambda_n + \lambda_2 - 2\lambda_1}{\lambda_n - \lambda_2} = \frac{\lambda_n - \lambda_2 + 2(\lambda_2 - \lambda_1)}{\lambda_n - \lambda_2} = 1 + 2 \frac{\lambda_2 - \lambda_1}{\lambda_n - \lambda_2} > 1. \quad (8.8)$$

Mit Hilfe dieser Transformation können wir nun das für die Berechnung der Eigenvektoren benötigte  $p_\epsilon$  optimal wählen: Zu einer gegebenen Dimension  $m \in \mathbb{N}$  setzen wir

$$p_m(x) := \frac{1}{c_m(x_0)} c_m(\Phi(x)) \quad (8.9)$$

und erhalten

$$p_m(\lambda_1) = \frac{1}{c_m(x_0)} c_m(\Phi(\lambda_1)) = 1, \quad |p_m(x)| \leq \frac{1}{c_m(x_0)} \quad \text{für alle } x \in [\lambda_2, \lambda_n].$$

Um konkret angeben zu können, wie gut die durch  $p_m$  induzierte Approximation des Eigenvektors ist, benötigen wir eine Abschätzung für  $1/c_m(x_0)$ :

**Lemma 8.17 (Konvergenzrate)** *Wir definieren*

$$\kappa := \frac{\lambda_n - \lambda_1}{\lambda_2 - \lambda_1} > 1, \quad \varrho := \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \quad (8.10)$$

und erhalten

$$\frac{1}{|c_m(x_0)|} = \frac{2\varrho^m}{1 + \varrho^{2m}} \leq 2\varrho^m \quad \text{für alle } m \in \mathbb{N}_0.$$

*Beweis.* Da wir die Tschebyscheff-Polynome in der zweiten Darstellung aus Lemma 8.15 verwenden wollen, müssen wir  $x_0 + \sqrt{x_0^2 - 1}$  geeignet darstellen. Wir beginnen mit

$$\begin{aligned} x_0^2 - 1 &= \left(1 + 2\frac{\lambda_2 - \lambda_1}{\lambda_n - \lambda_2}\right)^2 - 1 = 1 + 4\frac{\lambda_2 - \lambda_1}{\lambda_n - \lambda_2} + 4\frac{(\lambda_2 - \lambda_1)^2}{(\lambda_n - \lambda_2)^2} - 1 \\ &= 4\frac{(\lambda_2 - \lambda_1)(\lambda_n - \lambda_2) + (\lambda_2 - \lambda_1)^2}{(\lambda_n - \lambda_2)^2} = 4\frac{(\lambda_2 - \lambda_1)(\lambda_n - \lambda_1)}{(\lambda_n - \lambda_2)^2}. \end{aligned}$$

Aus dieser Gleichung und (8.8) ergibt sich

$$\begin{aligned} x_0 + \sqrt{x_0^2 - 1} &= \frac{(\lambda_n - \lambda_2) + 2(\lambda_2 - \lambda_1) + 2\sqrt{\lambda_2 - \lambda_1}\sqrt{\lambda_n - \lambda_1}}{\lambda_n - \lambda_2} \\ &= \frac{(\lambda_2 - \lambda_1) + (\lambda_n - \lambda_1) + 2\sqrt{\lambda_2 - \lambda_1}\sqrt{\lambda_n - \lambda_1}}{\lambda_n - \lambda_2} \\ &= \frac{(\sqrt{\lambda_2 - \lambda_1} + \sqrt{\lambda_n - \lambda_1})^2}{\lambda_n - \lambda_2} = \frac{\lambda_2 - \lambda_1}{\lambda_n - \lambda_2} \left(1 + \sqrt{\frac{\lambda_n - \lambda_1}{\lambda_2 - \lambda_1}}\right)^2 \\ &= \frac{\lambda_2 - \lambda_1}{(\lambda_n - \lambda_1) - (\lambda_2 - \lambda_1)} (1 + \sqrt{\kappa})^2 = \left(\frac{\lambda_n - \lambda_1}{\lambda_2 - \lambda_1} - 1\right)^{-1} (1 + \sqrt{\kappa})^2 \\ &= \frac{1}{\kappa - 1} (1 + \sqrt{\kappa})^2 = \frac{(\sqrt{\kappa} + 1)^2}{(\sqrt{\kappa} + 1)(\sqrt{\kappa} - 1)} = \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} = 1/\varrho. \end{aligned}$$

Einsetzen in die zweite Darstellung in Lemma 8.15 ergibt

$$c_m(x_0) = \frac{1}{2}((1/\varrho)^m + (1/\varrho)^{-m}) = \frac{1}{2}(\varrho^{-m} + \varrho^m) = \frac{1 + \varrho^{2m}}{2\varrho^m}$$

und damit

$$\frac{1}{c_m(x_0)} = \frac{2\varrho^m}{\varrho^{2m} + 1}.$$

■

Die für uns wichtigsten Eigenschaften des durch (8.9) gegebenen transformierten Tschebyscheff-Polynoms fassen wir in folgendem Lemma zusammen:

**Lemma 8.18 (Transformiertes Polynom)** Sei  $m \in \mathbb{N}_0$ , und seien  $\lambda_1 < \lambda_2 \leq \lambda_n$  gegeben. Dann existiert ein Polynom  $p_m \in \Pi_m$  mit

$$p_m(\lambda_1) = 1, \quad \max\{|p_m(x)| : x \in [\lambda_2, \lambda_n]\} \leq \frac{2\varrho^m}{1 + \varrho^{2m}}$$

mit  $\varrho$  aus (8.10).

*Beweis.* Wir definieren  $p_m$  durch (8.9) und wenden Lemma 8.17 an. Da  $c_m$  auf  $[-1, 1]$  nur Werte zwischen  $-1$  und  $1$  annehmen kann, kann  $p_m$  auf  $[\lambda_2, \lambda_n]$  nur Werte zwischen  $-1/c_m(x_0)$  und  $1/c_m(x_0)$  annehmen. ■

Aus diesem Lemma folgt unmittelbar die folgende Konvergenzaussage für den Eigenvektor  $\mathbf{e}$ :

**Folgerung 8.19 (Eigenvektor-Approximation)** Seien  $m \in \mathbb{N}_0$  und  $\mathbf{e} := \mathbf{Q}\delta^{(1)}$ . Sei  $\mathbf{q}^{(0)} \in \mathbb{K}^n$  mit  $\|\mathbf{q}^{(0)}\| = 1$  gegeben, und sei  $\varrho$  wie in Lemma 8.17 gewählt. Dann existiert ein Polynom  $p_m \in \Pi_m$  mit

$$\tan \angle(p_m(\mathbf{A})\mathbf{q}^{(0)}, \mathbf{e}) \leq \frac{2\varrho^m}{1 + \varrho^{2m}} \tan \angle(\mathbf{q}^{(0)}, \mathbf{e}) \quad \text{für alle } m \in \mathbb{N}_0.$$

*Beweis.* Wir kombinieren Lemma 8.12 mit Lemma 8.18. ■

Die Größe  $\kappa$  in (8.10) beschreibt, wie groß der Abstand zwischen  $\lambda_1$  und  $\lambda_2$  im Vergleich zu dem „Durchmesser“  $\lambda_n - \lambda_1$  des gesamten Spektrums ist. Je kleiner dieser relative Abstand wird, desto größer wird  $\kappa$  und desto näher rückt die „Konvergenzrate“  $\varrho$  an eins heran, desto langsamer konvergieren also die Eigenvektoren und Eigenwerte. Wie schon bei der Vektoriteration ist es also auch hier von Vorteil, wenn die Eigenwerte, die wir berechnen wollen, einen möglichst großen Abstand zu dem Rest des Spektrums aufweisen.

Man beachte, dass Folgerung 8.19 nur eine Existenzaussage bietet: Wir wissen, dass mit  $p_m(\mathbf{A})\mathbf{q}^{(0)}$  ein Vektor in dem Krylow-Raum  $\mathcal{K}(\mathbf{A}, \mathbf{q}^{(0)}, m)$  existiert, der  $\mathbf{e}$  approximiert, aber wir können ihn mit diesem Zugang nicht berechnen, weil wir die Eigenwerte  $\lambda_1, \lambda_2$  und  $\lambda_n$  nicht kennen, die für die Konstruktion von  $p_m$  erforderlich wären.

Bei Eigenwerten dagegen können wir einen direkten Bezug zwischen den berechenbaren Näherungswerten  $\lambda_1^{(m)}$  und dem Eigenwert  $\lambda_1$  herstellen:

**Satz 8.20** Sei  $\mathbf{A} = \mathbf{A}^* \in \mathbb{K}^{n \times n}$ . Seien  $\lambda_1 < \lambda_2 \leq \dots \leq \lambda_n$  die Eigenwerte von  $\mathbf{A}$ . Dann ist  $\mathbf{e} = \mathbf{Q}\delta^{(1)}$  ein Eigenvektor zu dem Eigenwert  $\lambda_1$ . Sei  $m \in [0 : m_0 - 1]$  und sei  $\lambda_1^{(m)}$  der kleinste Eigenwert von  $\widehat{\mathbf{A}}^{(m)}$ . Dann gilt

$$\lambda_1 \leq \lambda_1^{(m)} \leq \lambda_1 + (\lambda_n - \lambda_1) \left( \frac{2\varrho^m}{1 + \varrho^{2m}} \right)^2 \tan^2 \angle(\mathbf{q}^{(0)}, \mathbf{e}).$$

für die Konstante  $\varrho$  aus (8.10).

*Beweis.* Die linke Abschätzung folgt direkt aus

$$\begin{aligned}\lambda_1^{(m)} &= \min\{\Lambda_{\widehat{A}^{(m)}}(\widehat{\mathbf{x}}) : \widehat{\mathbf{x}} \in \mathbb{R}^k \setminus \{\mathbf{0}\}\} = \min\{\Lambda_A(\mathbf{x}) : \mathbf{x} \in \text{Bild } \mathbf{Q}^{(m)} \setminus \{\mathbf{0}\}\} \\ &\geq \min\{\Lambda_A(\mathbf{x}) : \mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}\} = \lambda_1.\end{aligned}$$

Für die rechte Abschätzung kombinieren wir Folgerung 8.19 mit Lemma 5.9. Wir setzen  $\mathbf{y} := p_m(\mathbf{A})\mathbf{q}^{(0)}$  und erhalten

$$\begin{aligned}\lambda_1^{(m)} - \lambda_1 &\leq \Lambda_A(\mathbf{y}) - \lambda_1 = |\Lambda_A(\mathbf{y}) - \lambda_1| \leq \|\mathbf{A} - \lambda_1\mathbf{I}\| \sin^2(\mathbf{y}, \mathbf{e}) \\ &\leq \|\mathbf{A} - \lambda_1\mathbf{I}\| \tan^2(p_m(\mathbf{A})\mathbf{q}^{(0)}, \mathbf{e}) \\ &\leq \|\mathbf{A} - \lambda_1\mathbf{I}\| \left(\frac{2\varrho^m}{1 + \varrho^{2m}}\right)^2 \tan^2 \angle(\mathbf{q}^{(0)}, \mathbf{e}).\end{aligned}$$

Es bleibt nur noch, die Norm der Matrix  $\mathbf{A} - \lambda_1\mathbf{I}$  abzuschätzen:

$$\|\mathbf{A} - \lambda_1\mathbf{I}\| = \|\mathbf{Q}\mathbf{D}\mathbf{Q}^* - \lambda_1\mathbf{Q}\mathbf{Q}^*\| = \|\mathbf{Q}(\mathbf{D} - \lambda_1\mathbf{I})\mathbf{Q}^*\| = \|\mathbf{D} - \lambda_1\mathbf{I}\|.$$

Mit der Abschätzung

$$\|(\mathbf{D} - \lambda_1\mathbf{I})\mathbf{z}\|^2 = \sum_{i=1}^n (\lambda_i - \lambda_1)^2 |z_i|^2 \leq \sum_{i=1}^n (\lambda_n - \lambda_1)^2 |z_i|^2 = (\lambda_n - \lambda_1)^2 \|\mathbf{z}\|^2$$

folgt  $\|\mathbf{A} - \lambda_1\mathbf{I}\| = \|\mathbf{D} - \lambda_1\mathbf{I}\| \leq \lambda_n - \lambda_1$ , also die Behauptung. ■

**Bemerkung 8.21 (Mehrfache Eigenwerte)** Die Voraussetzung  $\lambda_1 < \lambda_2$  ist in diesem Fall nicht kritisch, sie dient lediglich der Vereinfachung des Beweises. Falls  $\lambda_1 = \dots = \lambda_k < \lambda_{k+1}$  gilt, können wir den Beweis im Wesentlichen wie bisher durchführen und müssen lediglich das Polynom  $p_m$  so wählen, dass es auf  $[\lambda_{k+1}, \lambda_n]$  kleine Werte annimmt. Außerdem müssen wir den Tangens des Winkels zwischen  $\mathbf{q}^{(1)}$  und  $\mathbf{e}$  durch den Tangens des Winkels zwischen  $\mathbf{q}^{(1)}$  und dem von  $\mathbf{e}^{(1)}, \dots, \mathbf{e}^{(k)}$  aufgespannten Eigenraum zu dem Eigenwert  $\lambda_1$  ersetzen.

Die Fehlerabschätzung gilt dann mit  $\lambda_{k+1}$  anstelle von  $\lambda_2$ , entscheidend ist also der Abstand zwischen  $\lambda_1 = \dots = \lambda_k$  und dem Rest des Spektrums.

**Bemerkung 8.22 (Modellproblem)** Im Fall des zweidimensionalen Modellproblems sind wir vor allem daran interessiert, mit kleinen Gitterschrittweiten  $h \in \mathbb{R}_{>0}$  zu arbeiten. Für kleines  $h$  stellt man fest, dass  $\lambda_n \approx ch^{-2}$  für eine Konstante  $c \in \mathbb{R}_{>0}$  gilt, der größte Eigenwert wird also sehr groß werden. Damit folgt auch  $\kappa \approx ch^{-2}$ , und wir erhalten

$$\varrho = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} = 1 - \frac{2}{\sqrt{\kappa} + 1} \approx 1 - \frac{2}{\sqrt{ch^{-2}} + 1} \approx 1 - \frac{2}{\sqrt{c}}h,$$

die Konvergenzgeschwindigkeit unseres Verfahrens wird also für zunehmend feine Gitter zunehmend langsamer werden.

**Übungsaufgabe 8.23 (Ritz-Vektoren)** Seien  $n \in \mathbb{N}$  und eine selbstadjungierte Matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  gegeben. Seien eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  und

$$\mathbf{D} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix} \quad \text{mit} \quad \lambda_1 < \lambda_2 \leq \dots \leq \lambda_n$$

so gegeben, dass  $\mathbf{A} = \mathbf{Q}\mathbf{D}\mathbf{Q}^*$  gilt. Dann ist  $\mathbf{e} := \mathbf{Q}\delta(1)$  ein Eigenvektor der Matrix  $\mathbf{A}$  zu dem Eigenwert  $\lambda_1$ .

(a) Beweisen Sie

$$\frac{\|\widehat{\mathbf{x}} - \hat{x}_1 \delta^{(1)}\|^2}{\|\widehat{\mathbf{x}}\|^2} \leq \frac{\Lambda_D(\widehat{\mathbf{x}}) - \lambda_1}{\lambda_2 - \lambda_1} \quad \text{für alle } \widehat{\mathbf{x}} \in \mathbb{K}^n \setminus \{\mathbf{0}\}.$$

(b) Wir bezeichnen den Eigenraum zu dem kleinsten Eigenwert  $\lambda_1$  mit  $\mathcal{E} := \mathcal{E}_A(\lambda_1) = \text{span}\{\mathbf{e}\}$ . Beweisen Sie

$$\sin^2 \angle(\mathbf{x}, \mathcal{E}) \leq \frac{\Lambda_A(\mathbf{x}) - \lambda_1}{\lambda_2 - \lambda_1} \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}.$$

(c) Sei  $m \in [1 : n]$ , und sei  $\widehat{\mathbf{Q}} \in \mathbb{K}^{n \times m}$  eine isometrische Matrix. Sei  $\widehat{\mathbf{A}} := \widehat{\mathbf{Q}}^* \mathbf{A} \widehat{\mathbf{Q}}$ , und seien  $\hat{\lambda}_1 \in \mathbb{R}$  der kleinste Eigenwert dieser Matrix und  $\widehat{\mathbf{e}} \in \mathbb{K}^m \setminus \{\mathbf{0}\}$  ein zugehöriger Eigenvektor.

Beweisen Sie, dass der Ritz-Vektor  $\tilde{\mathbf{e}} := \widehat{\mathbf{Q}}\widehat{\mathbf{e}}$  die folgende Abschätzung erfüllt:

$$\sin^2 \angle(\tilde{\mathbf{e}}, \mathcal{E}) \leq \frac{\hat{\lambda}_1 - \lambda_1}{\lambda_2 - \lambda_1}.$$





## 9 Eigenwertverfahren für sehr große Matrizen

Die Konvergenz des Lanczos-Verfahrens hängt von dem Verhältnis zwischen dem größten und dem kleinsten Eigenwert der Matrix  $\mathbf{A}$  ab, also von ihrer Konditionszahl. Gerade bei Eigenwertproblemen, die im Kontext partieller Differentialgleichungen auftreten, wird diese Konditionszahl oft sehr groß sein, so dass wir mit einer relativ langsamen Konvergenz rechnen müssen. Deshalb interessiert man sich für Verfahren, die weniger empfindlich auf die Konditionszahl reagieren. Ein Kandidat wäre die inverse Iteration, bei der bekanntlich nur die Eigenwerte in der Nähe des gesuchten einen Einfluss auf die Konvergenzgeschwindigkeit haben, aber die Berechnung der Inversen oder das (exakte) Lösen eines linearen Gleichungssystems ist bei sehr großen Matrizen in der Regel schwierig. Deshalb sucht man nach Verfahren, die sich im Prinzip wie die inverse Iteration verhalten (oder sogar wie die inverse Iteration mit Shift), aber ohne das *exakte* Lösen großer Gleichungssysteme auskommen.

### 9.1 Richardson-Iteration

Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  wieder eine selbstadjungierte Matrix, es gelte also  $\mathbf{A} = \mathbf{A}^*$ . Nach Satz 3.47 finden wir eine unitäre Matrix  $\mathbf{Q} \in \mathbb{K}^{n \times n}$  und eine reelle Diagonalmatrix

$$\mathbf{D} = \begin{pmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{pmatrix}, \quad \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$$

mit  $\mathbf{Q}^* \mathbf{A} \mathbf{Q} = \mathbf{D}$ . Wir interessieren uns für den kleinsten Eigenwert  $\lambda_1$  der Matrix.

Eines der einfachsten iterativen Lösungsverfahren für lineare Gleichungssysteme der Form

$$\mathbf{A} \mathbf{x} = \mathbf{b} \tag{9.1}$$

mit vorgegebener rechter Seite  $\mathbf{b}$  ist die *Richardson-Iteration*

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} - \theta(\mathbf{A} \mathbf{x}^{(m)} - \mathbf{b}) \quad \text{für alle } m \in \mathbb{N}_0,$$

die ausgehend von einem Startvektor  $\mathbf{x}^{(0)}$  eine Folge von Näherungslösungen berechnet. Dabei ist die korrekte Wahl des *Dämpfungsparameters*  $\theta \in \mathbb{R}_{>0}$  entscheidend für das Konvergenzverhalten.

## 9 Eigenwertverfahren für sehr große Matrizen

Wenn wir den kleinsten Eigenwert und einen passenden Eigenvektor zu berechnen, können wir das Verfahren anpassen: Zu diesem Zweck ersetzen wir (9.1) durch die Eigenwertgleichung

$$\mathbf{A}\mathbf{x} = \lambda_1\mathbf{x},$$

die sich in die Form eines linearen Gleichungssystems

$$(\mathbf{A} - \lambda_1\mathbf{I})\mathbf{x} = \mathbf{0}$$

mit der rechten Seite  $\mathbf{b} = \mathbf{0}$  bringen lässt. Da wir den kleinsten Eigenwert  $\lambda_1$  nicht kennen, ersetzen wir ihn durch eine Näherung  $\mu \in \mathbb{R}$  und erhalten die Iterationsvorschrift

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} - \theta(\mathbf{A}\mathbf{x}^{(m)} - \mu\mathbf{x}^{(m)}) \quad \text{für alle } m \in \mathbb{N}_0. \quad (9.2)$$

Für die Analyse führen wir, wie schon häufiger, die transformierten Vektoren

$$\widehat{\mathbf{x}}^{(m)} := \mathbf{Q}^*\mathbf{x}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0$$

ein und erhalten mit  $\mathbf{x}^{(m)} = \mathbf{Q}\widehat{\mathbf{x}}^{(m)}$  aus (9.2) die Gleichung

$$\begin{aligned} \widehat{\mathbf{x}}^{(m+1)} &= \mathbf{Q}^*\mathbf{x}^{(m+1)} = \mathbf{Q}^*\mathbf{x}^{(m)} - \theta(\mathbf{Q}^*\mathbf{A}\mathbf{x}^{(m)} - \mu\mathbf{Q}^*\mathbf{x}^{(m)}) \\ &= \widehat{\mathbf{x}}^{(m)} - \theta(\mathbf{Q}^*\mathbf{A}\mathbf{Q}\widehat{\mathbf{x}}^{(m)} - \mu\widehat{\mathbf{x}}^{(m)}) \\ &= \widehat{\mathbf{x}}^{(m)} - \theta(\mathbf{D}\widehat{\mathbf{x}}^{(m)} - \mu\widehat{\mathbf{x}}^{(m)}) \\ &= (\mathbf{I} - \theta(\mathbf{D} - \mu\mathbf{I}))\widehat{\mathbf{x}}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0. \end{aligned}$$

Für die einzelnen Komponenten des Vektors  $\widehat{\mathbf{x}}^{(m)}$  folgt daraus

$$\widehat{x}_i^{(m)} = (1 - \theta(\lambda_i - \mu))^m \widehat{x}_i^{(0)} \quad \text{für alle } m \in \mathbb{N}_0, i \in [1 : n]. \quad (9.3)$$

Wir möchten erreichen, dass die erste Komponente der Iterationsvektoren gegenüber den anderen dominiert.

Falls wir  $\mu$  mit dem Rayleigh-Quotienten ermitteln, wird nach dem Satz 3.45 von Courant und Fischer immer  $\lambda_1 \leq \mu$  gelten, so dass  $\lambda_1 - \mu$  negativ ist. Wenn wir  $\theta > 0$  sicher stellen, folgt

$$1 - \theta(\lambda_i - \mu) \geq 1 \quad \text{für alle } i \in [1 : n] \text{ mit } \lambda_i \leq \mu.$$

Die Komponenten des Vektors, die zu Eigenwerten größer als  $\mu$  gehören, sollten möglichst reduziert werden, wir brauchen also

$$\begin{aligned} |1 - \theta(\lambda_i - \mu)| &< 1, \\ -1 &< 1 - \theta(\lambda_i - \mu) < 1 \quad \text{für alle } i \in [1 : n] \text{ mit } \lambda_i > \mu. \end{aligned}$$

Da  $\lambda_n$  der größte Eigenwert ist, genügt es,

$$-1 < 1 - \theta(\lambda_n - \mu) \iff 2 > \theta(\lambda_n - \mu) \iff \theta < \frac{2}{\lambda_n - \mu} \leq \frac{2}{\lambda_n - \lambda_1}$$

sicher zu stellen. Der Dämpfungsparemeter  $\theta$  muss also positiv, aber klein genug sein.

**Satz 9.1 (Konvergenz)** Seien  $k \in [1 : n - 1]$  und  $\mu \in \mathbb{R}_{>0}$  so gewählt, dass

$$\lambda_1 \leq \dots \leq \lambda_k \leq \mu < \lambda_{k+1} \leq \dots \leq \lambda_n$$

gilt. Sei  $\mathcal{E}$  der von den Eigenvektoren zu den ersten  $k$  Eigenwerten aufgespannte invariante Teilraum.

Für alle Dämpfungsparameter  $\theta \in \mathbb{R}_{>0}$  erfüllen die Vektoren der Richardson-Iteration die Abschätzung

$$\tan \angle(\mathbf{x}^{(m)}, \mathcal{E}) \leq \varrho^m \tan \angle(\mathbf{x}^{(0)}, \mathcal{E}) \quad \text{für alle } m \in \mathbb{N}_0,$$

wobei die Konvergenzrate durch

$$\varrho := \max\{1 - \theta(\lambda_{k+1} - \mu), \theta(\lambda_n - \mu) - 1\}$$

gegeben ist. Sie nimmt für den optimalen Dämpfungsparameter

$$\theta_{opt} = \frac{2}{(\lambda_n - \mu) + (\lambda_{k+1} - \mu)}$$

ihr Minimum

$$\varrho_{opt} = \frac{(\lambda_n - \mu) - (\lambda_{k+1} - \mu)}{(\lambda_n - \mu) + (\lambda_{k+1} - \mu)} < 1$$

an, ist aber für jedes  $\theta \in (0, \theta_{opt}]$  echt kleiner als eins.

*Beweis.* Wir definieren den Teilraum  $\widehat{\mathcal{E}} := \mathbb{K}^k \times \{\mathbf{0}\} \subseteq \mathbb{K}^n$  und halten  $\mathcal{E} = \mathbf{Q}\widehat{\mathcal{E}}$  fest. Da  $\mathbf{Q}$  unitär ist, gilt

$$\tan \angle(\mathbf{x}^{(m)}, \mathcal{E}) = \tan \angle(\widehat{\mathbf{x}}^{(m)}, \widehat{\mathcal{E}}) \quad \text{für alle } m \in \mathbb{N}_0.$$

Aus  $\theta > 0$  folgt

$$1 - \theta(\lambda_i - \mu) > 1 \quad \text{für alle } i \in [1 : k].$$

Da die Eigenwerte aufsteigend sortiert sind, gilt

$$1 - \theta(\lambda_n - \mu) \leq 1 - \theta(\lambda_i - \mu) \leq 1 - \theta(\lambda_{k+1} - \mu) \quad \text{für alle } i \in [k + 1 : n].$$

Falls  $\theta(\lambda_n - \mu) \leq 1$  gilt, haben wir

$$0 \leq 1 - \theta(\lambda_n - \mu) \leq 1 - \theta(\lambda_{k+1} - \mu),$$

also

$$|1 - \theta(\lambda_i - \mu)| \leq 1 - \theta(\lambda_{k+1} - \mu) \leq \varrho \quad \text{für alle } i \in [k + 1 : n].$$

Anderenfalls gilt

$$|1 - \theta(\lambda_n - \mu)| = \theta(\lambda_n - \mu) - 1$$

## 9 Eigenwertverfahren für sehr große Matrizen

und wir erhalten

$$|1 - \theta(\lambda_i - \mu)| \leq \max\{1 - \theta(\lambda_{k+1} - \mu), \theta(\lambda_n - \mu) - 1\} = \varrho \quad \text{für alle } i \in [k+1 : n].$$

Wie im Beweis des Lemmas 5.41 folgt mit (9.3) die Abschätzung

$$\begin{aligned} \tan^2 \angle(\widehat{\mathbf{x}}^{(m)}, \widehat{\mathcal{E}}) &= \frac{\sum_{i=k+1}^n |\hat{x}_i^{(m)}|^2}{\sum_{i=1}^k |\hat{x}_i^{(m)}|^2} \\ &= \frac{\sum_{i=k+1}^n |1 - \theta(\lambda_i - \mu)|^{2m} |\hat{x}_i^{(0)}|^2}{\sum_{i=1}^k |1 - \theta(\lambda_i - \mu)|^{2m} |\hat{x}_i^{(0)}|^2} \\ &\leq \frac{\sum_{i=k+1}^n \varrho^{2m} |\hat{x}_i^{(0)}|^2}{\sum_{i=1}^k |\hat{x}_i^{(0)}|^2} = \varrho^{2m} \tan^2 \angle(\widehat{\mathbf{x}}^{(0)}, \widehat{\mathcal{E}}) \quad \text{für alle } m \in \mathbb{N}_0. \end{aligned}$$

Für die Bestimmung des optimalen Dämpfungsparameters  $\theta$  gehen wir davon aus, dass  $\varrho$  als Maximum einer monoton fallenden und einer monoton wachsenden Funktion nur dann minimal sein kann, wenn

$$1 - \theta(\lambda_{k+1} - \mu) = \theta(\lambda_n - \mu) - 1$$

gilt, und diese Gleichung ist äquivalent mit

$$2 = \theta((\lambda_{k+1} - \mu) + (\lambda_n - \mu)).$$

Daraus ergibt sich die Formel für  $\theta_{\text{opt}}$ . Durch Einsetzen in die Definition der Konvergenzrate  $\varrho$  erhalten wir

$$\begin{aligned} \varrho_{\text{opt}} &= 1 - \theta_{\text{opt}}(\lambda_{k+1} - \mu) = 1 - \frac{2(\lambda_{k+1} - \mu)}{(\lambda_n - \mu) + (\lambda_{k+1} - \mu)} \\ &= \frac{(\lambda_n - \mu) + (\lambda_{k+1} - \mu) - 2(\lambda_{k+1} - \mu)}{(\lambda_n - \mu) + (\lambda_{k+1} - \mu)} = \frac{(\lambda_n - \mu) - (\lambda_{k+1} - \mu)}{(\lambda_n - \mu) + (\lambda_{k+1} - \mu)}, \end{aligned}$$

und wegen  $\mu < \lambda_{k+1}$  folgt  $\varrho_{\text{opt}} < 1$ . Für ein beliebiges  $\theta \in (0, \theta_{\text{opt}}]$  erhalten wir

$$1 > 1 - \theta(\lambda_{k+1} - \mu) \geq 1 - \theta_{\text{opt}}(\lambda_{k+1} - \mu) = \theta_{\text{opt}}(\lambda_n - \mu) - 1 \geq \theta(\lambda_n - \mu) - 1,$$

also insbesondere  $\varrho < 1$ . ■

Wir sehen, dass auch bei der Richardson-Iteration der größte Eigenwert eine Rolle für die Konvergenzrate spielt: Selbst wenn wir den Dämpfungsparameter  $\theta$  optimal wählen, haben wir

$$\varrho_{\text{opt}} = \frac{(\lambda_n - \mu) - (\lambda_{k+1} - \mu)}{(\lambda_n - \mu) + (\lambda_{k+1} - \mu)} = 1 - \frac{2\lambda_{k+1} - 2\mu}{\lambda_n + \lambda_{k+1} - 2\mu},$$

und dieser Ausdruck wird für große  $\lambda_n$  sehr nahe an eins liegen, so dass eine sehr langsame Konvergenz zu erwarten ist.

## 9.2 Optimale Dämpfung

Es wäre von Vorteil, wenn wir den Dämpfungsparameter  $\theta$  der Richardson-Iteration in optimaler Weise durch einen geeigneten Algorithmus wählen lassen könnten.

Gemäß (9.2) ist die neue Iterierte  $\mathbf{x}^{(m+1)}$  eine Linearkombination zwischen der alten Iterierten  $\mathbf{x}^{(m)}$  und dem *Residuum*

$$\mathbf{r}^{(m)} := \mu \mathbf{x}^{(m)} - \mathbf{A} \mathbf{x}^{(m)},$$

also stellt sich die Frage, wie wir diese Linearkombination wählen sollen, um dem gesuchten Eigenvektor möglichst nahe zu kommen.

Bei der Beantwortung dieser Frage können wir uns wieder von dem Satz 3.45 von Courant und Fischer leiten lassen: Der gewünschte Eigenwert  $\lambda_1$  ist das Minimum des Rayleigh-Quotienten, also bietet es sich an,

$$\mathbf{x}^{(m+1)} \in \mathcal{V}^{(m)} := \text{span}\{\mathbf{x}^{(m)}, \mathbf{r}^{(m)}\}$$

so zu wählen, dass

$$\Lambda_A(\mathbf{x}^{(m+1)}) \leq \Lambda_A(\mathbf{z}) \quad \text{für alle } \mathbf{z} \in \mathcal{V}^{(m)} \quad (9.4)$$

gilt. Dazu stellen wir den zweidimensionalen Teilraum  $\mathcal{V}^{(m)}$  als Bild der Matrix

$$\mathbf{V}^{(m)} := \begin{pmatrix} \mathbf{x}^{(m)} & \mathbf{r}^{(m)} \end{pmatrix}$$

dar und konstruieren mit der dünnen QR-Zerlegung

$$\mathbf{Q}^{(m)} \mathbf{R}^{(m)} = \mathbf{V}^{(m)}$$

eine isometrische Matrix  $\mathbf{Q}^{(m)} \in \mathbb{K}^{n \times 2}$ , deren Bild den Teilraum  $\mathcal{V}^{(m)}$  zumindest enthält.

Wenn wir die nächste Iterierte als

$$\mathbf{x}^{(m+1)} = \mathbf{Q}^{(m)} \widehat{\mathbf{x}}^{(m+1)}$$

setzen, erhalten wir mit Lemma 3.17 die Gleichung

$$\begin{aligned} \Lambda_A(\mathbf{x}^{(m+1)}) &= \frac{\langle \mathbf{x}^{(m+1)}, \mathbf{A} \mathbf{x}^{(m+1)} \rangle}{\langle \mathbf{x}^{(m+1)}, \mathbf{x}^{(m+1)} \rangle} = \frac{\langle \mathbf{Q}^{(m)} \widehat{\mathbf{x}}^{(m+1)}, \mathbf{A} \mathbf{Q}^{(m)} \widehat{\mathbf{x}}^{(m+1)} \rangle}{\langle \mathbf{Q}^{(m)} \widehat{\mathbf{x}}^{(m+1)}, \mathbf{Q}^{(m)} \widehat{\mathbf{x}}^{(m+1)} \rangle} \\ &= \frac{\langle \widehat{\mathbf{x}}^{(m+1)}, (\mathbf{Q}^{(m)})^* \mathbf{A} \mathbf{Q}^{(m)} \widehat{\mathbf{x}}^{(m+1)} \rangle}{\langle \widehat{\mathbf{x}}^{(m+1)}, (\mathbf{Q}^{(m)})^* \mathbf{Q}^{(m)} \widehat{\mathbf{x}}^{(m+1)} \rangle} = \frac{\langle \widehat{\mathbf{x}}^{(m+1)}, (\mathbf{Q}^{(m)})^* \mathbf{A} \mathbf{Q}^{(m)} \widehat{\mathbf{x}}^{(m+1)} \rangle}{\langle \widehat{\mathbf{x}}^{(m+1)}, \widehat{\mathbf{x}}^{(m+1)} \rangle}, \end{aligned}$$

wir können also

$$\widehat{\mathbf{A}}^{(m)} := (\mathbf{Q}^{(m)})^* \mathbf{A} \mathbf{Q}^{(m)} \in \mathbb{K}^{2 \times 2}$$

definieren, um zu der Gleichung

$$\Lambda_A(\mathbf{x}^{(m+1)}) = \Lambda_{\widehat{\mathbf{A}}^{(m)}}(\widehat{\mathbf{x}}^{(m+1)}) \quad (9.5)$$

## 9 Eigenwertverfahren für sehr große Matrizen

zu gelangen. Nach dem Satz 3.45 von Courant und Fischer nimmt die rechte Seite ihr Minimum an, wenn  $\hat{\mathbf{x}}^{(m+1)}$  der Eigenvektor zu dem kleinsten Eigenwert der zweidimensionalen Matrix  $\hat{\mathbf{A}}^{(m)}$  ist.

Also können wir die bestmögliche Linearkombination der Vektoren  $\mathbf{x}^{(m)}$  und  $\mathbf{r}^{(m)}$  bestimmen, indem wir ein zweidimensionales Eigenwertproblem lösen. Wie man das bewerkstelligt, wissen wir bereits seit Kapitel 4.

Der Eigenwert entspricht dank (9.5) dem Rayleigh-Quotienten der neuen Iterierten  $\mathbf{x}^{(m+1)}$ , so dass wir ihn für das aus Kapitel 5 bekannte Abbruchkriterium verwenden können, ohne erneut den Rayleigh-Quotienten explizit auszuwerten.

**Algorithmus 9.2 (Gradientenverfahren)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  selbstadjungiert. Der folgende Algorithmus berechnet eine Folge  $(\mathbf{x}^{(m)})_{m \in \mathbb{N}_0}$ , die unter geeigneten Bedingungen gegen einen Eigenvektor zu dem kleinsten Eigenwert von  $\mathbf{A}$  konvergiert.

```

 $\mathbf{x} \leftarrow \mathbf{x} / \|\mathbf{x}\|;$ 
 $\lambda \leftarrow \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle; \quad \mathbf{r} \leftarrow \mathbf{A}\mathbf{x} - \lambda\mathbf{x};$ 
while  $\|\mathbf{r}\| > \epsilon|\lambda|$  do begin
   $\mathbf{V} \leftarrow (\mathbf{x} \ \mathbf{r}) \in \mathbb{K}^{n \times 2};$ 
  Berechne eine dünne QR-Zerlegung  $\mathbf{QR} = \mathbf{V}$  mit  $\mathbf{Q} \in \mathbb{K}^{n \times 2};$ 
   $\mathbf{B} \leftarrow \mathbf{A}\mathbf{Q};$ 
   $\hat{\mathbf{A}} \leftarrow \mathbf{Q}^* \mathbf{B}; \quad \{ = \mathbf{Q}^* \mathbf{A} \mathbf{Q} \}$ 
  Finde einen normierten Eigenvektor  $\mathbf{y}$  für
    den minimalen Eigenwert  $\lambda$  der Matrix  $\hat{\mathbf{A}}$ ;
   $\mathbf{x} \leftarrow \mathbf{Q}\mathbf{y};$ 
   $\mathbf{r} \leftarrow \mathbf{B}\mathbf{y} - \lambda\mathbf{x} \quad \{ = \mathbf{A}\mathbf{x} - \lambda\mathbf{x} \}$ 
end

```

Die Bezeichnung „Gradientenverfahren“ ist dadurch motiviert, dass wegen Lemma 8.2

$$\text{span}\{\mathbf{x}^{(m)}, \mathbf{r}^{(m)}\} = \text{span}\{\mathbf{x}^{(m)}, \nabla \Lambda_{\mathbf{A}}(\mathbf{x}^{(m)})\}$$

gilt, wir verbessern unsere Lösung also gerade in Richtung des Gradienten des Rayleigh-Quotienten. Da unser Algorithmus dabei die bestmögliche Schrittweite verwendet, entspricht er dem konventionellen Gradientenverfahren für die Minimierung des Rayleigh-Quotienten, allerdings mit der Besonderheit, dass die Iterationsvektoren in jedem Schritt normiert werden. Da der Rayleigh-Quotient invariant unter Skalierung des Arguments ist, ändert sich dadurch das Konvergenzverhalten nicht.

Bei der Analyse des Gradientenverfahrens können wir auf die in Kapitel 8 geleistete Vorarbeit zurückgreifen: In jedem Schritt des Verfahrens wird die bestmögliche Linearkombination zwischen  $\mathbf{x}^{(m)}$  und  $\mathbf{r}^{(m)}$  gewählt, und wegen

$$\text{span}\{\mathbf{x}^{(m)}, \mathbf{r}^{(m)}\} = \text{span}\{\mathbf{x}^{(m)}, \mathbf{A}\mathbf{x}^{(m)}\}$$

können wir einen solchen Schritt wie ein Lanczos-Verfahren behandeln, das nach einem einzigen Schritt abbricht.

**Satz 9.3 (Konvergenz)** *Es gelte  $\lambda_1 < \lambda_2$ . Mit*

$$\varrho := \frac{(\lambda_n - \lambda_1) - (\lambda_2 - \lambda_1)}{(\lambda_n - \lambda_1) + (\lambda_2 - \lambda_1)} = 1 - \frac{2(\lambda_2 - \lambda_1)}{(\lambda_n - \lambda_1) + (\lambda_2 - \lambda_1)} < 1$$

*erhalten wir die Abschätzung*

$$\Lambda_A(\mathbf{x}^{(m+1)}) - \lambda_1 \leq \min \left\{ (\lambda_n - \lambda_1) \varrho^2 \tan^2 \angle(\mathbf{x}^{(m)}, \mathbf{e}), \right. \\ \left. \Lambda_A(\mathbf{x}^{(m)}) - \lambda_1 \right\} \quad \text{für alle } m \in \mathbb{N}_0,$$

*wobei  $\mathbf{e} := \mathbf{Q}\delta^{(1)}$  ein Eigenvektor zu dem Eigenwert  $\lambda_1$  ist.*

*Beweis.* Da das Gradientenverfahren nach Konstruktion den Rayleigh-Quotienten in dem Krylow-Raum  $\mathcal{K}(\mathbf{A}, \mathbf{x}^{(m)}, 1) = \text{span}\{\mathbf{x}^{(m)}, \mathbf{A}\mathbf{x}^{(m)}\}$  minimiert, können wir seine Analyse auf die eines Lanczos-Verfahrens mit einem Schritt zurückführen.

Mit Satz 8.20 erhalten wir die Abschätzung

$$\Lambda_A(\mathbf{x}^{(m+1)}) - \lambda_1 \leq (\lambda_n - \lambda_1) \left( \frac{2\hat{\varrho}}{1 + \hat{\varrho}^2} \right)^2 \tan^2 \angle(\mathbf{x}^{(m)}, \mathbf{e}) \quad \text{für alle } m \in \mathbb{N}_0$$

mit den Konstanten

$$\hat{\varrho} = \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1}, \quad \kappa = \frac{\lambda_n - \lambda_1}{\lambda_2 - \lambda_1}.$$

Dank der Gleichung

$$\begin{aligned} \varrho &= 2 \left( \frac{1 + \hat{\varrho}^2}{\hat{\varrho}} \right)^{-1} = 2 \left( \frac{1}{\hat{\varrho}} + \hat{\varrho} \right)^{-1} = 2 \left( \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} + \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right)^{-1} \\ &= 2 \left( \frac{(\sqrt{\kappa} + 1)^2 + (\sqrt{\kappa} - 1)^2}{(\sqrt{\kappa} + 1)(\sqrt{\kappa} - 1)} \right)^{-1} = 2 \left( \frac{\kappa + 2\sqrt{\kappa} + 1 + \kappa - 2\sqrt{\kappa} + 1}{\kappa - 1} \right)^{-1} \\ &= 2 \left( \frac{2\kappa + 2}{\kappa - 1} \right)^{-1} = \frac{\kappa - 1}{\kappa + 1} = \frac{\frac{\lambda_n - \lambda_1}{\lambda_2 - \lambda_1} - 1}{\frac{\lambda_n - \lambda_1}{\lambda_2 - \lambda_1} + 1} = \frac{(\lambda_n - \lambda_1) - (\lambda_2 - \lambda_1)}{(\lambda_n - \lambda_1) + (\lambda_2 - \lambda_1)} \end{aligned}$$

folgt die erste Abschätzung. Die zweite folgt unmittelbar, da wir den Rayleigh-Quotienten auf einem Teilraum minimieren, der  $\mathbf{x}^{(m)}$  enthält. ■

**Bemerkung 9.4 (Mehrfacher Eigenwert)** *Falls  $\lambda_1 = \lambda_2 = \dots = \lambda_k < \lambda_{k+1}$  gilt, erhalten wir ein leicht modifiziertes Konvergenzresultat: Wir ersetzen  $\mathbf{e}$  durch den Aufspann  $\mathcal{E} = \text{span}\{\mathbf{e}^{(1)}, \dots, \mathbf{e}^{(k)}\}$  der ersten  $k$  Eigenvektoren  $\mathbf{e}^{(i)} := \mathbf{Q}\delta^{(i)}$ ,  $i \in [1 : k]$ , und erhalten*

$$\Lambda_A(\mathbf{x}^{(m+1)}) - \lambda_1 \leq (\lambda_n - \lambda_1) \varrho^2 \tan^2 \angle(\mathbf{x}^{(m)}, \mathcal{E}) \quad \text{für alle } m \in \mathbb{N}_0$$

*mit der Konvergenzrate*

$$\varrho := \frac{(\lambda_n - \lambda_1) - (\lambda_{k+1} - \lambda_1)}{(\lambda_n - \lambda_1) + (\lambda_{k+1} - \lambda_1)} < 1.$$

## 9 Eigenwertverfahren für sehr große Matrizen

In vielen Anwendungen ist  $\lambda_n$  sehr viel größer als  $\lambda_1$ , so dass der Faktor  $\varrho$  relativ nahe bei eins liegen wird und deshalb mit relativ langsamer Konvergenz rechnen müssen.

**Übungsaufgabe 9.5 (Lokales Minimum)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  selbstadjungiert, sei  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$ . Beweisen Sie, dass aus

$$\Lambda_A(\mathbf{x}) \leq \Lambda_A(\mathbf{y}) \quad \text{für alle } \mathbf{y} \in \text{span}\{\mathbf{x}, \mathbf{Ax}\}$$

folgt, dass  $\mathbf{x}$  ein Eigenvektor ist.

Hinweis: Man könnte eine Arnoldi-Basis konstruieren und Folgerung 7.6 auf die Matrizen  $\widehat{\mathbf{A}}^{(0)}$  und  $\widehat{\mathbf{A}}^{(1)}$  anwenden.

### 9.3 Vorkonditionierer

Unser Ziel besteht nun darin, die Geschwindigkeit des Gradientenverfahrens zu verbessern. Insbesondere stellt die Abhängigkeit der Konvergenz von dem größten Eigenwert  $\lambda_n$  ein Problem dar, da bei vielen Aufgabenstellungen, beispielsweise auch bei unserem Modellproblem, dieser Eigenwert sehr groß werden kann.

Die inverse Iteration kennt dieses Problem nicht, für sie konnten wir eine Konvergenzrate von  $|\lambda_2|/|\lambda_1|$  nachweisen. Leider ist es unrealistisch, bei sehr großen Matrizen  $\mathbf{A}$  die Inverse exakt zu berechnen, und die Approximation der Eigenwerte einer genäherten Inversen ist etwas unbefriedigend.

Wir können allerdings die inverse Iteration so umschreiben, dass wir mit Näherungen arbeiten können, aber weiterhin die „richtigen“ Eigenwerte bestimmen: Da wir lediglich an der Richtung der Vektoren interessiert sind, können wir die Iterierten der inversen Iteration beliebig skalieren. Wir wählen die Skalierung mit dem Rayleigh-Quotienten und erhalten

$$\begin{aligned} \mathbf{x}^{(m+1)} &= \Lambda_A(\mathbf{x}^{(m)})\mathbf{A}^{-1}\mathbf{x}^{(m)} = \mathbf{x}^{(m)} - \mathbf{A}^{-1}\mathbf{Ax}^{(m)} + \Lambda_A(\mathbf{x}^{(m)})\mathbf{A}^{-1}\mathbf{x}^{(m)} \\ &= \mathbf{x}^{(m)} - \mathbf{A}^{-1}(\mathbf{Ax}^{(m)} - \Lambda_A(\mathbf{x}^{(m)})\mathbf{x}^{(m)}). \end{aligned}$$

Falls  $\mathbf{x}^{(m)}$  ein Eigenvektor der Matrix  $\mathbf{A}$  ist, gilt  $\mathbf{Ax}^{(m)} = \Lambda_A(\mathbf{x}^{(m)})\mathbf{x}^{(m)}$ , also  $\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)}$ , exakte Eigenvektoren sind also Fixpunkte dieser Iteration.

Diese Eigenschaft bleibt erhalten, wenn wir die exakte Inverse durch eine Näherung  $\mathbf{N}$  ersetzen, um zu der *vorkonditionierten Richardson-Iteration* (auch bekannt als PINVIT, *preconditioned inverse iteration*).

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} - \mathbf{N}(\mathbf{Ax}^{(m)} - \Lambda_A(\mathbf{x}^{(m)})\mathbf{x}^{(m)}) \quad \text{für alle } m \in \mathbb{N}_0$$

zu gelangen. Für jede invertierbare *Vorkonditionierungsmatrix*  $\mathbf{N}$  sind die Fixpunkte dieser Iteration Eigenvektoren der Matrix  $\mathbf{A}$ . Da die Matrix  $\mathbf{A}$  selbstadjungiert ist und  $\mathbf{N}$  ihre Inverse approximieren soll, wird in der Regel vorausgesetzt, dass auch  $\mathbf{N}$  selbstadjungiert ist.



**Bemerkung 9.6 (Vorkonditionierung)** Der Begriff Vorkonditionierung stammt dabei aus der Welt der linearen Gleichungssysteme: Falls wir

$$\mathbf{Ax} = \mathbf{b}$$

lösen wollen, aber die Konditionszahl der Matrix  $\mathbf{A}$  groß ist, werden viele iterative Lösungsverfahren nur sehr langsam konvergieren. Deshalb ersetzt man das System durch das äquivalente System

$$\mathbf{NAx} = \mathbf{Nb},$$

in dem durch eine geeignete Vorkonditionierungsmatrix  $\mathbf{N}$  die Konditionszahl reduziert und so das Konvergenzverhalten verbessert werden kann.

Bei der vorkonditionierten Richardson-Iteration können wir den Dämpfungparameter in der Matrix  $\mathbf{N}$  verschwinden lassen, da er keinen Einfluss auf die Konditionszahl hat, so dass sie die einfache Form

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} - \mathbf{N}(\mathbf{Ax}^{(m)} - \mathbf{b}) \quad \text{für alle } m \in \mathbb{N}_0$$

annimmt. Die Fehler bezüglich der exakten Lösung  $\mathbf{x}$  des linearen Gleichungssystems entwickeln sich gemäß

$$\begin{aligned} \mathbf{x}^{(m+1)} - \mathbf{x} &= \mathbf{x}^{(m)} - \mathbf{x} - \mathbf{N}(\mathbf{Ax}^{(m)} - \mathbf{b}) \\ &= \mathbf{x}^{(m)} - \mathbf{x} - \mathbf{NA}(\mathbf{x}^{(m)} - \mathbf{x}) \\ &= (\mathbf{I} - \mathbf{NA})(\mathbf{x}^{(m)} - \mathbf{x}) \end{aligned} \quad \text{für alle } m \in \mathbb{N}_0.$$

Für die Konvergenz ist demnach die Iterationsmatrix  $\mathbf{I} - \mathbf{NA}$  ausschlaggebend. Gilt beispielsweise  $\|\mathbf{I} - \mathbf{NA}\| < 1$ , so konvergieren die Vektoren  $\mathbf{x}^{(m)}$  gegen die Lösung  $\mathbf{x}$ .

Die vorkonditionierte Richardson-Iteration für Eigenwertprobleme entsteht aus der Fassung für lineare Gleichungssysteme, indem in jedem Schritt die konstante rechte Seite  $\mathbf{b}$  durch den von  $\mathbf{x}^{(m)}$  abhängenden Vektor  $\Lambda_A(\mathbf{x}^{(m)})$  ersetzt wird. Damit wird die Iteration nichtlinear, so dass sich ihr Konvergenzverhalten nicht mehr einfach durch die Iterationsmatrix charakterisieren lässt.

Die vorkonditionierte Richardson-Iteration ähnelt sehr der uns bereits bekannten nicht vorkonditionierten Version, wir haben lediglich das Residuum

$$\mathbf{r}^{(m)} = \Lambda_A(\mathbf{x}^{(m)})\mathbf{x}^{(m)} - \mathbf{Ax}^{(m)}$$

durch ein vorkonditioniertes Residuum

$$\mathbf{p}^{(m)} = \mathbf{Nr}^{(m)} = \mathbf{N}(\Lambda_A(\mathbf{x}^{(m)})\mathbf{x}^{(m)} - \mathbf{Ax}^{(m)})$$

ersetzt. Indem wir wieder orthogonalisieren und ein zweidimensionales Eigenwertproblem lösen, um den optimalen Dämpfungparameter zu wählen, also die bestmögliche Linearkombination des alten Iterationsvektors und des vorkonditionierten Residuums, gelangen wir zu dem vorkonditionierten Gradientenverfahren.

**Algorithmus 9.7 (Vorkonditioniertes Gradientenverfahren)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix. Sei  $\mathbf{N} \in \mathbb{K}^{n \times n}$  ein für  $\mathbf{A}$  geeigneter Vorkonditionierer. Der folgende Algorithmus berechnet eine Folge  $(\mathbf{x}^{(m)})_{m \in \mathbb{N}_0}$ , die unter geeigneten Bedingungen gegen einen Eigenvektor konvergiert.

```

 $\mathbf{x} \leftarrow \mathbf{x} / \|\mathbf{x}\|;$ 
 $\lambda \leftarrow \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle; \quad \mathbf{r} \leftarrow \mathbf{A}\mathbf{x} - \lambda\mathbf{x};$ 
 $\mathbf{p} \leftarrow \mathbf{N}\mathbf{r};$ 
while  $\|\mathbf{r}\| > \epsilon|\lambda|$  do begin
   $\mathbf{V} \leftarrow (\mathbf{x} \ \mathbf{p}) \in \mathbb{K}^{n \times 2};$ 
  Berechne eine dünne QR-Zerlegung  $\mathbf{QR} = \mathbf{V}$  mit  $\mathbf{Q} \in \mathbb{K}^{n \times 2};$ 
   $\mathbf{B} \leftarrow \mathbf{A}\mathbf{Q};$ 
   $\widehat{\mathbf{A}} \leftarrow \mathbf{Q}^*\mathbf{B}; \quad \{ = \mathbf{Q}^*\mathbf{A}\mathbf{Q} \}$ 
  Finde einen normierten Eigenvektor  $\mathbf{y}$  für
    den minimalen Eigenwert  $\lambda$  der Matrix  $\widehat{\mathbf{A}};$ 
   $\mathbf{x} \leftarrow \mathbf{Q}\mathbf{y};$ 
   $\mathbf{r} \leftarrow \mathbf{B}\mathbf{y} - \lambda\mathbf{x}; \quad \{ = \mathbf{A}\mathbf{x} - \lambda\mathbf{x} \}$ 
   $\mathbf{p} \leftarrow \mathbf{N}\mathbf{r}$ 
end

```

Auch dieses Verfahren hat den großen Vorteil, dass in jedem Schritt nur zwei Multiplikationen mit der Matrix  $\mathbf{A}$  und sogar nur eine mit der Matrix  $\mathbf{N}$  benötigt werden, so dass man es auch in Situationen anwenden kann, in denen diese Matrizen nicht explizit im Speicher vorliegen. Beispielsweise wird bei manchen Diskretisierungen die Matrix  $\mathbf{A}$  nicht gespeichert, weil sich Speicherplatz sparen lässt, indem man ihre Koeffizienten nach Bedarf aus der Beschreibung der zugrundeliegenden Geometrie rekonstruiert.

Bei der Behandlung linearer Gleichungssysteme lässt sich das Gradientenverfahren wesentlich verbessern, indem man zu dem *Verfahren der konjugierten Gradienten* (engl. *conjugate gradients method*, deshalb auch im Deutschen als *cg-Verfahren* bekannt) übergeht, das häufig erheblich schneller konvergiert. Also wäre es interessant, auch für Eigenwertprobleme eine entsprechende Variante des Gradientenverfahrens zu entwickeln.

Der Satz 3.45 besagt, dass der kleinste Eigenwert einer selbstadjungierten Matrix  $\mathbf{A}$  das Minimum des Rayleigh-Quotienten ist, das Eigenwertproblem ist also äquivalent mit einem Minimierungsproblem.

Dementsprechend wollen wir nun auch die Suche nach der Lösung eines linearen Gleichungssystems

$$\mathbf{A}\mathbf{x} = \mathbf{b} \tag{9.6}$$

als Minimierungsproblem formulieren. Wir nehmen an, dass die Matrix  $\mathbf{A}$  selbstadjungiert und positiv definit (vgl. Definition 3.22) ist, und führen die Funktion

$$f: \mathbb{K}^n \rightarrow \mathbb{R}, \quad \mathbf{x} \mapsto \frac{1}{2} \langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle - \operatorname{Re} \langle \mathbf{b}, \mathbf{x} \rangle$$

ein. Da  $\mathbf{A}$  positiv definit ist, gilt  $\langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle \in \mathbb{R}_{\geq 0}$ , so dass die Funktion  $f$  tatsächlich reellwertig ist.

Man kann nachrechnen, dass

$$f(\mathbf{x}) \leq f(\mathbf{x} + t\mathbf{p}) \quad \text{für alle } t \in \mathbb{K} \quad (9.7a)$$

und

$$\langle \mathbf{p}, \mathbf{b} - \mathbf{Ax} \rangle = 0. \quad (9.7b)$$

äquivalent sind. Falls  $\mathbf{x}$  Lösung des linearen Gleichungssystems (9.6) ist, gilt offenbar (9.7b) für alle  $\mathbf{p}$ , also ist  $f(\mathbf{x})$  minimal.

Falls umgekehrt  $f(\mathbf{x})$  minimal ist, können wir in (9.7b) auch  $\mathbf{p} := \mathbf{b} - \mathbf{Ax}$  einsetzen und erhalten  $\|\mathbf{b} - \mathbf{Ax}\|^2 = 0$ , also folgt

$$f(\mathbf{x}) \text{ minimal} \quad \iff \quad \mathbf{Ax} = \mathbf{b}.$$

Das Verfahren der konjugierten Gradienten beginnt mit einem Startvektor  $\mathbf{x}^{(0)}$ , berechnet sein Residuum

$$\mathbf{r}^{(0)} = \mathbf{b} - \mathbf{Ax}^{(0)},$$

setzt  $\mathbf{p}^{(0)} := \mathbf{r}^{(0)}$  und sucht dann einen Skalierungsfaktor  $\lambda_0$  so, dass die verbesserte Näherungslösung

$$\mathbf{x}^{(1)} := \mathbf{x}^{(0)} + \lambda_0 \mathbf{p}^{(0)}$$

die Optimalitätsbedingung

$$f(\mathbf{x}^{(1)}) \leq f(\mathbf{x}^{(0)} + \mathbf{y}) \quad \text{für alle } \mathbf{y} \in \text{span}\{\mathbf{p}^{(0)}\}$$

erfüllt. Damit ist der erste Schritt vollzogen.

Wenn nun für  $m \in \mathbb{N}$  eine Lösung  $\mathbf{x}^{(m)}$  mit

$$f(\mathbf{x}^{(m)}) \leq f(\mathbf{x}^{(0)} + \mathbf{y}) \quad \text{für alle } \mathbf{y} \in \text{span}\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(m-1)}\} \quad (9.8)$$

berechnet worden ist, soll sie in einer neuen Richtung  $\mathbf{p}^{(m)}$  verbessert werden, die neue Näherung

$$\mathbf{x}^{(m+1)} := \mathbf{x}^{(m)} + \lambda_m \mathbf{p}^{(m)}$$

soll also

$$f(\mathbf{x}^{(m+1)}) \leq f(\mathbf{x}^{(0)} + \mathbf{y}) \quad \text{für alle } \mathbf{y} \in \text{span}\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(m)}\}$$

erfüllen. Das ist äquivalent mit

$$\langle \mathbf{p}, \mathbf{b} - \mathbf{Ax}^{(m+1)} \rangle = 0 \quad \text{für alle } \mathbf{p} \in \text{span}\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(m)}\}.$$

Für  $\mathbf{p} = \mathbf{p}^{(m)}$  können wir diese Gleichung durch die Wahl des Skalierungsparameters  $\lambda_m$  sicher stellen. Für  $\ell \in [0 : m - 1]$  folgt mit (9.7b) die Gleichung

$$0 = \langle \mathbf{p}^{(\ell)}, \mathbf{b} - \mathbf{Ax}^{(m+1)} \rangle = \langle \mathbf{p}^{(\ell)}, \mathbf{b} - \mathbf{Ax}^{(m)} \rangle - \lambda_m \langle \mathbf{p}^{(\ell)}, \mathbf{Ap}^{(m)} \rangle = -\lambda_m \langle \mathbf{p}^{(\ell)}, \mathbf{Ap}^{(m)} \rangle,$$

## 9 Eigenwertverfahren für sehr große Matrizen

und da wir nicht  $\lambda_m = 0$  wählen wollen, muss

$$\langle \mathbf{p}^{(\ell)}, \mathbf{A}\mathbf{p}^{(m)} \rangle = 0 \quad \text{für alle } \ell \in [0 : m - 1]$$

gelten. Also ist es sinnvoll, dafür zu sorgen, dass alle Richtungen die Eigenschaft

$$\langle \mathbf{p}^{(\ell)}, \mathbf{A}\mathbf{p}^{(k)} \rangle = 0 \quad \text{für alle } k, \ell \in [0 : m] \text{ mit } k \neq \ell$$

besitzen. Für  $\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(m-1)}$  dürfen wir das im Rahmen der Induktion voraussetzen, die Richtung  $\mathbf{p}^{(m)}$  müssen wir entsprechend konstruieren. Dazu gehen wir von dem Residuum

$$\mathbf{r}^{(m)} := \mathbf{b} - \mathbf{A}\mathbf{x}^{(m)}$$

aus und nutzen den Gram-Schmidt-Orthogonalisierungsalgorithmus

$$\mathbf{p}^{(m)} = \mathbf{r}^{(m)} - \sum_{k=0}^{m-1} \frac{\langle \mathbf{p}^{(k)}, \mathbf{A}\mathbf{r}^{(m)} \rangle}{\langle \mathbf{p}^{(k)}, \mathbf{A}\mathbf{p}^{(k)} \rangle} \mathbf{p}^{(k)}, \quad (9.9)$$

der für alle  $\ell \in [0 : m]$  die gewünschte Gleichung

$$\begin{aligned} \langle \mathbf{p}^{(\ell)}, \mathbf{A}\mathbf{p}^{(m)} \rangle &= \langle \mathbf{p}^{(\ell)}, \mathbf{A}\mathbf{r}^{(m)} \rangle - \sum_{k=0}^m \frac{\langle \mathbf{p}^{(k)}, \mathbf{A}\mathbf{r}^{(m)} \rangle}{\langle \mathbf{p}^{(k)}, \mathbf{A}\mathbf{p}^{(k)} \rangle} \langle \mathbf{p}^{(\ell)}, \mathbf{A}\mathbf{p}^{(k)} \rangle \\ &= \langle \mathbf{p}^{(\ell)}, \mathbf{A}\mathbf{r}^{(m)} \rangle - \frac{\langle \mathbf{p}^{(\ell)}, \mathbf{A}\mathbf{r}^{(m)} \rangle}{\langle \mathbf{p}^{(\ell)}, \mathbf{A}\mathbf{p}^{(\ell)} \rangle} \langle \mathbf{p}^{(\ell)}, \mathbf{A}\mathbf{p}^{(\ell)} \rangle = 0 \end{aligned}$$

garantiert. Man kann beweisen (vgl. Übungsaufgabe 9.8), dass

$$\mathbf{A}\mathbf{p}^{(k)} \in \text{span}\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(k+1)}\} \quad \text{für alle } k \in [0 : m - 2]$$

gilt, so dass aus der Optimalitätsbedingung (9.8) und  $\mathbf{A}^* = \mathbf{A}$  mit Lemma 3.17 die Gleichung

$$\langle \mathbf{p}^{(k)}, \mathbf{A}\mathbf{r}^{(m)} \rangle = \langle \mathbf{A}\mathbf{p}^{(k)}, \mathbf{r}^{(m)} \rangle = 0 \quad \text{für alle } k \in [0 : m - 2]$$

folgt und sich die Orthogonalisierung (9.9) vereinfacht zu

$$\mathbf{p}^{(m)} = \mathbf{r}^{(m)} - \frac{\langle \mathbf{p}^{(m-1)}, \mathbf{A}\mathbf{r}^{(m)} \rangle}{\langle \mathbf{p}^{(m-1)}, \mathbf{A}\mathbf{p}^{(m-1)} \rangle} \mathbf{p}^{(m-1)}. \quad (9.10)$$

Insgesamt erhalten wir

$$\begin{aligned} \mathbf{x}^{(m+1)} &= \mathbf{x}^{(m)} + \lambda_m \mathbf{p}^{(m)} \\ &= \mathbf{x}^{(m)} + \lambda_m \mathbf{r}^{(m)} - \lambda_m \frac{\langle \mathbf{p}^{(m-1)}, \mathbf{A}\mathbf{r}^{(m)} \rangle}{\langle \mathbf{p}^{(m-1)}, \mathbf{A}\mathbf{p}^{(m-1)} \rangle} \mathbf{p}^{(m-1)} \in \text{span}\{\mathbf{x}^{(m)}, \mathbf{p}^{(m-1)}, \mathbf{r}^{(m)}\}. \end{aligned}$$

Wir dürfen davon ausgehen, dass  $\lambda_{m-1} \neq 0$  gilt, so dass aus

$$\mathbf{x}^{(m)} = \mathbf{x}^{(m-1)} + \lambda_{m-1} \mathbf{p}^{(m-1)} \iff \lambda_{m-1} \mathbf{p}^{(m-1)} = \mathbf{x}^{(m)} - \mathbf{x}^{(m-1)} \in \text{span}\{\mathbf{x}^{(m-1)}, \mathbf{x}^{(m)}\}$$

auch  $\mathbf{p}^{(m-1)} \in \text{span}\{\mathbf{x}^{(m-1)}, \mathbf{x}^{(m)}\}$  folgt und wir insgesamt

$$\mathbf{x}^{(m+1)} \in \text{span}\{\mathbf{x}^{(m-1)}, \mathbf{x}^{(m)}, \mathbf{r}^{(m)}\}$$

bewiesen haben. Das Verfahren der konjugierten Gradienten bestimmt also seinen nächsten Schritt wie folgt:

Sei  $\mathbf{r}^{(m)} := \mathbf{b} - \mathbf{A}\mathbf{x}^{(m)}$ . Finde  $\mathbf{x}^{(m+1)} \in \text{span}\{\mathbf{x}^{(m-1)}, \mathbf{x}^{(m)}, \mathbf{r}^{(m)}\}$  mit

$$f(\mathbf{x}^{(m+1)}) \leq f(\mathbf{y}) \quad \text{für alle } \mathbf{y} \in \text{span}\{\mathbf{x}^{(m-1)}, \mathbf{x}^{(m)}, \mathbf{r}^{(m)}\}.$$

Für den Schritt von (9.9) zu (9.10) ist das folgende Lemma nützlich:

**Übungsaufgabe 9.8 (Basen der Krylow-Räume)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  selbstadjungiert und positiv definit. Seien  $\mathbf{x}^{(0)}, \mathbf{b} \in \mathbb{K}^n$  gegeben. Wir konstruieren

$$\begin{aligned} \mathbf{r}^{(m)} &:= \mathbf{b} - \mathbf{A}\mathbf{x}^{(m)} && \text{für alle } m \in [0 : m_0], \\ \mathbf{p}^{(m)} &:= \begin{cases} \mathbf{r}^{(0)} & \text{falls } m = 0, \\ \mathbf{r}^{(m)} - \frac{\langle \mathbf{p}^{(m-1)}, \mathbf{A}\mathbf{r}^{(m)} \rangle}{\langle \mathbf{p}^{(m-1)}, \mathbf{A}\mathbf{p}^{(m-1)} \rangle} \mathbf{p}^{(m-1)} & \text{ansonsten} \end{cases} && \text{für alle } m \in [0 : m_0], \\ \mathbf{x}^{(m)} &:= \mathbf{x}^{(m-1)} + \frac{\langle \mathbf{p}^{(m-1)}, \mathbf{r}^{(m-1)} \rangle}{\langle \mathbf{p}^{(m-1)}, \mathbf{A}\mathbf{p}^{(m-1)} \rangle} \mathbf{p}^{(m-1)} && \text{für alle } m \in [1 : m_0], \end{aligned}$$

wobei  $m_0 \in \mathbb{N}_0$  die kleinste Zahl mit  $\mathbf{r}^{(m_0)} = \mathbf{0}$  ist.

(a) Beweisen Sie  $\langle \mathbf{p}^{(m-1)}, \mathbf{r}^{(m)} \rangle = 0$  für alle  $m \in [1 : m_0]$ .

(b) Beweisen Sie, dass die Vektoren  $\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(m_0-1)}$  linear unabhängig sind.

(c) Beweisen Sie für alle  $m \in [0 : m_0 - 1]$

$$\text{span}\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(m)}\} \subseteq \text{span}\{\mathbf{r}^{(0)}, \dots, \mathbf{r}^{(m)}\} \subseteq \mathcal{K}(\mathbf{A}, \mathbf{r}^{(0)}, m)$$

(d) Beweisen Sie  $\mathcal{K}(\mathbf{A}, \mathbf{r}^{(0)}, m) \subseteq \text{span}\{\mathbf{p}^{(0)}, \dots, \mathbf{p}^{(m)}\}$ , dass also alle in (c) genannten Teilräume identisch sind.

Hinweis: Für den Beweis des Teils (d) sind die Teile (b) und (c) nützlich.

Dieser Ansatz lässt sich offensichtlich auf unsere Aufgabenstellung, nämlich die Minimierung des Rayleigh-Quotienten, übertragen: Das Residuum berechnen wir wie in Kapitel 5 unter Verwendung des Rayleigh-Quotienten als  $\mathbf{r}^{(m)} := \Lambda_A(\mathbf{x}^{(m)})\mathbf{x}^{(m)} - \mathbf{A}\mathbf{x}^{(m)}$  und bestimmen dann die nächste Approximation des Eigenvektors als Lösung des folgenden Minimierungsproblems:

Sei  $\mathbf{r}^{(m)} := \Lambda_A(\mathbf{x}^{(m)})\mathbf{x}^{(m)} - \mathbf{A}\mathbf{x}^{(m)}$ .

Finde  $\mathbf{x}^{(m+1)} \in \text{span}\{\mathbf{x}^{(m-1)}, \mathbf{x}^{(m)}, \mathbf{r}^{(m)}\}$  mit

$$\Lambda_A(\mathbf{x}^{(m+1)}) \leq \Lambda_A(\mathbf{y}) \quad \text{für alle } \mathbf{y} \in \text{span}\{\mathbf{x}^{(m-1)}, \mathbf{x}^{(m)}, \mathbf{r}^{(m)}\}.$$

## 9 Eigenwertverfahren für sehr große Matrizen

Natürlich wollen wir uns die Möglichkeit bewahren, einen Vorkonditionierer einzusetzen, deshalb ersetzen wir wieder das Residuum  $\mathbf{r}^{(m)}$  durch das vorkonditionierte Residuum  $\mathbf{q}^{(m)} := \mathbf{N}\mathbf{r}^{(m)}$  und erhalten das folgende Verfahren:

Sei  $\mathbf{r}^{(m)} := \Lambda_A(\mathbf{x}^{(m)})\mathbf{x}^{(m)} - \mathbf{A}\mathbf{x}^{(m)}$ , sei  $\mathbf{q}^{(m)} := \mathbf{N}\mathbf{r}^{(m)}$ .  
 Finde  $\mathbf{x}^{(m+1)} \in \text{span}\{\mathbf{x}^{(m-1)}, \mathbf{x}^{(m)}, \mathbf{q}^{(m)}\}$  mit

$$\Lambda_A(\mathbf{x}^{(m+1)}) \leq \Lambda_A(\mathbf{y}) \quad \text{für alle } \mathbf{y} \in \text{span}\{\mathbf{x}^{(m-1)}, \mathbf{x}^{(m)}, \mathbf{q}^{(m)}\}.$$

Da das Residuum bei Eigenwertproblemen erheblich weniger handlich als bei linearen Gleichungssystemen ist, vererbt sich die Optimalität einer Näherung bezüglich eines Teilraums leider nicht auf im folgenden Schritt berechnete Näherungen. Trotzdem wird die Näherung im Vergleich zum Gradientenverfahren besser sein, da das Minimum in einem größeren Raum gesucht wird. Deshalb spricht man nur von einem *lokal* optimalen cg-Verfahren (LOPCG, *locally optimal preconditioned conjugate gradients*).

**Algorithmus 9.9 (Vorkonditioniertes cg-Verfahren)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix. Sei  $\mathbf{N} \in \mathbb{K}^{n \times n}$  ein für  $\mathbf{A}$  geeigneter Vorkonditionierer. Der folgende Algorithmus berechnet eine Folge  $(\mathbf{x}^{(m)})_{m \in \mathbb{N}_0}$ , die unter geeigneten Bedingungen gegen einen Eigenvektor konvergiert.

```

 $\mathbf{x} \leftarrow \mathbf{x} / \|\mathbf{x}\|;$ 
 $\lambda \leftarrow \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle; \quad \mathbf{r} \leftarrow \mathbf{A}\mathbf{x} - \lambda\mathbf{x};$ 
 $\mathbf{q} \leftarrow \mathbf{N}\mathbf{r};$ 
 $\mathbf{V} \leftarrow (\mathbf{x} \quad \mathbf{q}) \in \mathbb{K}^{n \times 2};$ 
while  $\|\mathbf{r}\| > \epsilon|\lambda|$  do begin
    Berechne eine dünne QR-Zerlegung  $\mathbf{QR} = \mathbf{V};$ 
     $\mathbf{B} \leftarrow \mathbf{AQ};$ 
     $\hat{\mathbf{A}} \leftarrow \mathbf{Q}^*\mathbf{B}; \quad \{ = \mathbf{Q}^*\mathbf{A}\mathbf{Q} \}$ 
    Finde einen normierten Eigenvektor  $\mathbf{y}$  für
    den minimalen Eigenwert  $\lambda$  der Matrix  $\hat{\mathbf{A}};$ 
     $\mathbf{z} \leftarrow \mathbf{x};$ 
     $\mathbf{x} \leftarrow \mathbf{Q}\mathbf{y};$ 
     $\mathbf{r} \leftarrow \mathbf{B}\mathbf{y} - \lambda\mathbf{x} \quad \{ = \mathbf{A}\mathbf{x} - \lambda\mathbf{x} \}$ 
     $\mathbf{q} \leftarrow \mathbf{N}\mathbf{r};$ 
     $\mathbf{V} \leftarrow (\mathbf{z} \quad \mathbf{x} \quad \mathbf{q}) \in \mathbb{K}^{n \times 3}$ 
end

```

Ein *global* optimales Verfahren haben wir bereits kennen gelernt: Der Lanczos-Algorithmus 8.11 berechnet das Minimum jeweils im Aufspann *aller* bisher aufgetretenen Residuen, benötigt allerdings auch sehr viel mehr Speicher als das lokal optimale Verfahren, da für die Konstruktion der Approximation des Eigenvektors alle bis zu dem aktuellen Iterationsschritt konstruierten Vektoren der Arnoldi-Basis aufbewahrt werden müssen. Falls wir nur an dem Eigenwert interessiert sind, lässt sich dieser Speicherplatz natürlich einsparen.

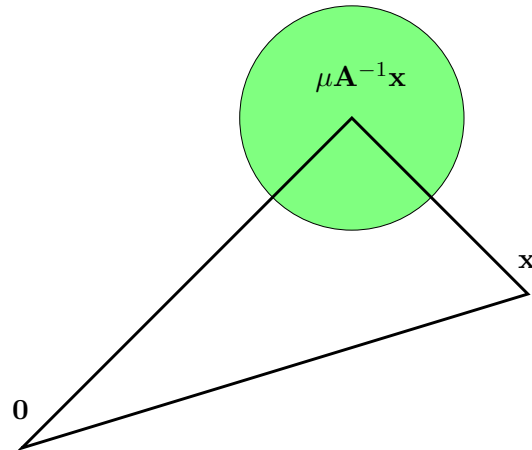


Abbildung 9.1: Geometrie der vorkonditionierten Iteration: Alle möglichen Iterationsvektoren  $\mathbf{y}$  liegen in einem Kreis mit Radius  $\gamma\|\mathbf{x} - \mu\mathbf{A}^{-1}\mathbf{x}\|_A$  um den Iterationsvektor  $\mu\mathbf{A}^{-1}\mathbf{x}$  der inversen Iteration.

**Übungsaufgabe 9.10 (Vorkonditionierte Richardson-Iteration)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte und positiv definite Matrix. Wir bezeichnen mit  $\mathbf{A}^{1/2}$  ihre Quadratwurzel, also diejenige selbstadjungierte und positiv definite Matrix, die  $(\mathbf{A}^{1/2})^2 = \mathbf{A}$  erfüllt. Deren Inverse bezeichnen wir mit  $\mathbf{A}^{-1/2}$ .

Die Energienorm eines Vektors  $\mathbf{x} \in \mathbb{K}^n$  ist durch  $\|\mathbf{x}\|_A := \|\mathbf{A}^{1/2}\mathbf{x}\|$  definiert.

Sei  $\mathbf{N} \in \mathbb{K}^{n \times n}$  selbstadjungiert, sei  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$ . Die erste Iterierte des vorkonditionierten Richardson-Verfahrens bezeichnen wir mit

$$\mathbf{y} := \mathbf{x} - \mathbf{N}(\mathbf{A}\mathbf{x} - \mu\mathbf{x}), \quad \mu := \Lambda_A(\mathbf{x}).$$

(a) Beweisen Sie die Gleichung

$$\|\mathbf{x}\|_A^2 = \|\mu\mathbf{A}^{-1}\mathbf{x}\|_A^2 + \|\mathbf{x} - \mu\mathbf{A}^{-1}\mathbf{x}\|_A^2.$$

(b) Beweisen Sie

$$\|\mathbf{y} - \mu\mathbf{A}^{-1}\mathbf{x}\|_A \leq \|\mathbf{I} - \mathbf{A}^{1/2}\mathbf{N}\mathbf{A}^{1/2}\| \|\mathbf{x} - \mu\mathbf{A}^{-1}\mathbf{x}\|_A.$$

(c) Sei  $\gamma \in [0, 1]$  gegeben. Beweisen Sie, dass aus

$$(1 - \gamma)\langle \mathbf{z}, \mathbf{A}^{-1}\mathbf{z} \rangle \leq \langle \mathbf{z}, \mathbf{N}\mathbf{z} \rangle \leq (1 + \gamma)\langle \mathbf{z}, \mathbf{A}^{-1}\mathbf{z} \rangle \quad \text{für alle } \mathbf{z} \in \mathbb{K}^n$$

die Ungleichung  $\|\mathbf{I} - \mathbf{A}^{1/2}\mathbf{N}\mathbf{A}^{1/2}\| \leq \gamma$  folgt.

Hinweis: Bei Teil (c) können Satz 3.45 und Lemma 3.55 helfen.

## 9.4 Block-Verfahren

Bisher haben wir uns lediglich mit der Frage nach der Berechnung eines Eigenvektors für den kleinsten Eigenwert beschäftigt. Ähnlich wie im Fall der orthogonalen Iteration (vgl. Abschnitt 5.5) können wir den einzelnen Iterationsvektor durch eine Basis solcher Vektoren ersetzen und in dieser Weise versuchen, eine Anzahl der kleinsten Eigenwerte zu ermitteln. Wir führen wieder mehrere Iterationen simultan durch und fassen die dabei entstehenden Vektoren zu isometrischen Matrizen  $\mathbf{X}^{(m)} \in \mathbb{K}^{n \times k}$  zusammen.

Wir hoffen, dass die  $\ell$ -ten Spalten der Matrizen  $\mathbf{X}^{(m)}$  gegen den Eigenraum des  $\ell$ -ten Eigenwerts konvergieren werden. Falls die Eigenwerte unterschiedlich sind, sollten wir bei der Berechnung des Residuums also unterschiedliche Eigenwerte einsetzen. Da wir in diesem Abschnitt häufig mit einzelnen Spalten bestimmter Matrizen arbeiten müssen, führen wir die Notation ein, dass  $\mathbf{M}_\ell = \mathbf{M}\delta^{(\ell)}$  gerade die  $\ell$ -te Spalte einer Matrix  $\mathbf{M}$  bezeichnet.

Das Residuum ist dann die durch

$$\mathbf{R}_\ell^{(m)} := \mathbf{A}\mathbf{X}_\ell^{(m)} - \lambda_\ell^{(m)}\mathbf{X}_\ell^{(m)}, \quad \lambda_\ell^{(m)} := \langle \mathbf{X}_\ell^{(m)}, \mathbf{A}\mathbf{X}_\ell^{(m)} \rangle \quad \text{für alle } \ell \in [1 : k]$$

definierte Matrix  $\mathbf{R}^{(m)} \in \mathbb{K}^{n \times k}$ . Der Rest des Verfahrens lässt sich einfach anpassen, indem alle Vektoren durch Matrizen ersetzt werden:

**Algorithmus 9.11 (Block-Gradientenverfahren)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix. Sei  $\mathbf{N} \in \mathbb{K}^{n \times n}$  ein für  $\mathbf{A}$  geeigneter Vorkonditionierer. Der folgende Algorithmus berechnet eine Folge von Matrizen  $(\mathbf{X}^{(m)})_{m \in \mathbb{N}_0}$ , deren Spalten unter geeigneten Bedingungen gegen Eigenvektoren zu den kleinsten Eigenwerten konvergieren.

```

for  $\ell \in [1 : k]$  do begin
     $\lambda_\ell \leftarrow \langle \mathbf{X}_\ell, \mathbf{A}\mathbf{X}_\ell \rangle$ ;    $\mathbf{R}_\ell \leftarrow \mathbf{A}\mathbf{X}_\ell - \lambda_\ell\mathbf{X}_\ell$ 
end;
 $\mathbf{P} \leftarrow \mathbf{N}\mathbf{R}$ ;
while  $\|\mathbf{R}\|$  zu groß do begin
     $\mathbf{V} \leftarrow (\mathbf{X} \ \mathbf{P}) \in \mathbb{K}^{n \times (2k)}$ ;
    Berechne eine dünne QR-Zerlegung  $\mathbf{Q}\mathbf{R} = \mathbf{V}$  mit  $\mathbf{Q} \in \mathbb{K}^{n \times (2k)}$ ;
     $\mathbf{B} \leftarrow \mathbf{A}\mathbf{Q}$ ;
     $\hat{\mathbf{A}} \leftarrow \mathbf{Q}^*\mathbf{B}$ ;    $\{ = \mathbf{Q}^*\mathbf{A}\mathbf{Q} \}$ 
    Finde eine isometrische Matrix  $\mathbf{Y} \in \mathbb{K}^{(2k) \times k}$ , deren Spalten Eigenvektoren
        zu den  $k$  kleinsten Eigenwerten  $\lambda_1 \leq \dots \leq \lambda_k$  der Matrix  $\hat{\mathbf{A}}$  enthält;
     $\mathbf{X} \leftarrow \mathbf{Q}\mathbf{Y}$ ;
    for  $\ell \in [1 : k]$  do
         $\mathbf{R}_\ell \leftarrow \mathbf{A}\mathbf{X}_\ell - \lambda_\ell\mathbf{X}_\ell$ ;
     $\mathbf{P} \leftarrow \mathbf{N}\mathbf{R}$ 
end

```

Das Block-Gradientenverfahren bietet dieselben Vorteile wie die anderen bisher besprochenen Blockverfahren: Es lassen sich mehrere Eigenvektoren simultan berechnen,



so dass sich beispielsweise zu mehrfachen Eigenwerten die vollständigen Eigenräume konstruieren lassen.

Selbstverständlich können wir wie im vorigen Abschnitt auch eine Blockvariante des lokal optimalen cg-Verfahrens konstruieren: Wir nehmen  $\mathbf{X}^{(m-1)}$  hinzu und erhalten so das *lokal optimale vorkonditionierte Block-cg-Verfahren* (LOBPCG, *locally optimal block preconditioned conjugate gradients*):

**Algorithmus 9.12 (Block-cg-Verfahren)** Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  eine selbstadjungierte Matrix. Sei  $\mathbf{N} \in \mathbb{K}^{n \times n}$  ein für  $\mathbf{A}$  geeigneter Vorkonditionierer. Der folgende Algorithmus berechnet eine Folge von Matrizen  $(\mathbf{X}^{(m)})_{m \in \mathbb{N}_0}$ , deren Spalten unter geeigneten Bedingungen gegen Eigenvektoren zu den kleinsten Eigenwerten konvergieren.

```

for  $\ell \in [1 : k]$  do begin
   $\lambda_\ell \leftarrow \langle \mathbf{X}_\ell, \mathbf{A}\mathbf{X}_\ell \rangle$ ;    $\mathbf{R}_\ell \leftarrow \mathbf{A}\mathbf{X}_\ell - \lambda_\ell \mathbf{X}_\ell$ 
end;
 $\mathbf{P} \leftarrow \mathbf{N}\mathbf{R}$ ;
 $\mathbf{V} \leftarrow (\mathbf{X} \ \mathbf{P}) \in \mathbb{K}^{n \times (2k)}$ ;
while  $\|\mathbf{R}\|$  zu groß do begin
  Berechne eine dünne QR-Zerlegung  $\mathbf{QR} = \mathbf{V}$ ;
   $\mathbf{B} \leftarrow \mathbf{A}\mathbf{Q}$ ;
   $\hat{\mathbf{A}} \leftarrow \mathbf{Q}^*\mathbf{B}$ ;   {  $= \mathbf{Q}^*\mathbf{A}\mathbf{Q}$  }
  Finde eine isometrische Matrix  $\mathbf{Y}$ , deren Spalten Eigenvektoren
    zu den  $k$  kleinsten Eigenwerten  $\lambda_1 \leq \dots \leq \lambda_k$  der Matrix  $\hat{\mathbf{A}}$  enthält;
   $\mathbf{Z} \leftarrow \mathbf{X}$ ;
   $\mathbf{X} \leftarrow \mathbf{Q}\mathbf{Y}$ ;
  for  $\ell \in [1 : k]$  do
     $\mathbf{R}_\ell \leftarrow \mathbf{A}\mathbf{X}_\ell - \lambda_\ell \mathbf{X}_\ell$ ;
   $\mathbf{P} \leftarrow \mathbf{N}\mathbf{R}$ ;
   $\mathbf{V} \leftarrow (\mathbf{Z} \ \mathbf{X} \ \mathbf{P}) \in \mathbb{K}^{n \times (3k)}$ 
end

```

**Bemerkung 9.13 (Residuum)** Falls wir lediglich auf der Suche nach einem invarianten Teilraum sind, können wir wie in Kapitel 5 vorgehen und den verallgemeinerten Rayleigh-Quotienten

$$\mathbf{\Lambda}^{(m)} := (\mathbf{X}^{(m)})^* \mathbf{A} \mathbf{X}^{(m)}$$

verwenden, der auch die Interaktion zwischen verschiedenen Spalten der Matrix  $\mathbf{X}^{(m)}$  erfasst. Wir verwenden dann das modifizierte Residuum

$$\mathbf{R}^{(m)} = \mathbf{A}\mathbf{X}^{(m)} - \mathbf{X}^{(m)}\mathbf{\Lambda}^{(m)}.$$

Die Matrix  $\mathbf{\Lambda}^{(m)}$  wird gegen Diagonalform konvergieren, allerdings möglicherweise langsam, falls die ersten  $k$  Eigenwerte nahe beieinander liegen.

## 9.5 Eigenwert-Mehrgitterverfahren

Bei sehr großen Matrizen führt kein Weg daran vorbei, deren strukturelle Eigenschaften so gut wie möglich auszunutzen. Im Fall partieller Differentialgleichungen besteht eine solche Eigenschaft darin, dass wir sie in unterschiedlichen Auflösungen diskretisieren können, um unterschiedlich große Matrizen zu erhalten. Da diese Matrizen Näherungen desselben kontinuierlichen Problems darstellen, hoffen wir, dass sie zueinander in Beziehung stehen und dass sich diese Beziehung für unsere Zwecke nutzen lässt.

Bei partiellen Differentialgleichungen kommt häufig eine *Finite-Elemente-Diskretisierung* zum Einsatz, bei der das kontinuierliche Eigenwertproblem in die Form eines *verallgemeinerten Eigenwertproblems* übergeht: Neben die Matrix  $\mathbf{A} \in \mathbb{K}^{n \times n}$  tritt eine zweite  $\mathbf{M} \in \mathbb{K}^{n \times n}$ , und  $\lambda \in \mathbb{K}$  ist ein Eigenwert, falls ein Eigenvektor  $\mathbf{e} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  existiert mit

$$\mathbf{A}\mathbf{e} = \lambda\mathbf{M}\mathbf{e}. \quad (9.11)$$

Sofern  $\mathbf{M}$  invertierbar ist, sind solche verallgemeinerten Eigenwertprobleme nicht schwieriger als die bisher betrachteten, weil (9.11) dann äquivalent ist mit

$$\mathbf{M}^{-1}\mathbf{A}\mathbf{e} = \lambda\mathbf{e}.$$

Häufig ist  $\mathbf{M}$  sogar selbstadjungiert und positiv definit, dann können wir die Cholesky-Zerlegung  $\mathbf{M} = \mathbf{L}\mathbf{L}^*$  benutzen, um

$$\mathbf{L}^{-1}\mathbf{A}\mathbf{L}^{-*}\hat{\mathbf{e}} = \lambda\hat{\mathbf{e}}, \quad \hat{\mathbf{e}} = \mathbf{L}^*\mathbf{e}$$

zu erhalten, so dass die verbliebene Matrix sogar wieder selbstadjungiert ist und wir die bisher entwickelte Theorie in vollem Umfang anwenden können.

Insbesondere finden wir dann mit Satz 3.47 eine Orthonormalbasis  $(\hat{\mathbf{e}}^{(i)})_{i=1}^n$  aus Eigenvektoren. Es folgt

$$\begin{aligned} \langle \mathbf{e}^{(i)}, \mathbf{M}\mathbf{e}^{(j)} \rangle &= \langle \mathbf{L}^{-*}\hat{\mathbf{e}}^{(i)}, \mathbf{L}\mathbf{M}\mathbf{L}^{-*}\hat{\mathbf{e}}^{(j)} \rangle = \langle \hat{\mathbf{e}}^{(i)}, \mathbf{L}^{-1}\mathbf{M}\mathbf{L}^{-*}\hat{\mathbf{e}}^{(j)} \rangle \\ &= \langle \hat{\mathbf{e}}^{(i)}, \hat{\mathbf{e}}^{(j)} \rangle = \begin{cases} 1 & \text{falls } i = j, \\ 0 & \text{ansonsten} \end{cases} \quad \text{für alle } i, j \in [1 : n], \end{aligned}$$

die verallgemeinerten Eigenvektoren bilden also bezüglich des *Massen-Skalarprodukts*

$$\langle \mathbf{x}, \mathbf{y} \rangle_M := \langle \mathbf{x}, \mathbf{M}\mathbf{y} \rangle \quad \text{für alle } \mathbf{x}, \mathbf{y} \in \mathbb{K}^n$$

wieder eine Orthonormalbasis. Die von diesem Skalarprodukt induzierte Norm

$$\|\mathbf{x}\|_M := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_M} = \sqrt{\langle \mathbf{x}, \mathbf{M}\mathbf{x} \rangle} \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n$$

nennen wir die *Massen-Norm*.

Im Folgenden setzen wir voraus, dass  $\mathbf{M}$  positiv definit und selbstadjungiert ist. Dann können wir für das verallgemeinerte Eigenwertproblem auch einen Rayleigh-Quotienten

$$\Lambda_{A,M} : \mathbb{K}^n \setminus \{\mathbf{0}\} \rightarrow \mathbb{K}, \quad \mathbf{x} \mapsto \frac{\langle \mathbf{x}, \mathbf{A}\mathbf{x} \rangle}{\langle \mathbf{x}, \mathbf{M}\mathbf{x} \rangle}$$

definieren, der ähnliche Eigenschaften wie der bisher betrachtete aufweist, insbesondere gilt eine Version des Satzes 3.45 von Courant und Fischer.

Wie bereits erwähnt beruhen Mehrgitterverfahren auf der Idee, die Beziehungen zwischen unterschiedlich feinen Diskretisierungen derselben Differentialgleichung auszunutzen. Wir gehen deshalb davon aus, dass es Familien selbstadjungierter Matrizen  $(\mathbf{A}_\ell)_{\ell=0}^L$  und selbstadjungierter positiv definiter Matrizen  $(\mathbf{M}_\ell)_{\ell=0}^L$  mit

$$\mathbf{A}_\ell, \mathbf{M}_\ell \in \mathbb{K}^{n_\ell \times n_\ell} \quad n_\ell \in \mathbb{N}, \quad \text{für alle } \ell \in [0 : L]$$

gibt. Dabei soll jeweils  $\mathbf{A}_\ell$  zu einer feineren Diskretisierung als  $\mathbf{A}_{\ell-1}$  gehören.  $\mathbf{A}_0$  gehört zu der gröbsten Diskretisierung und  $\mathbf{A}_L$  gehört zu der feinsten.

Unsere Aufgabe besteht darin, den kleinsten Eigenwert und seinen Eigenraum für das Eigenwertproblem auf der feinsten Stufe zu berechnen. Um die Darstellung des Algorithmus nicht unnötig kompliziert zu gestalten, gehen wir davon aus, dass der kleinste Eigenwert einfach ist, dass also sein Eigenraum von einem einzigen Eigenvektor aufgespannt wird.

Wir werden gleich sehen, dass es für die Lösung dieser Aufgabe sehr hilfreich ist, auch die kleinsten Eigenwerte und ihre Eigenvektoren auf allen anderen Stufen zu berechnen, also die Eigenwertprobleme

$$\mathbf{A}_\ell \mathbf{e}_\ell = \lambda_\ell \mathbf{M}_\ell \mathbf{e}_\ell, \quad \mathbf{e}_\ell \neq \mathbf{0}, \quad \text{für alle } \ell \in [0 : L] \quad (9.12)$$

zu lösen, wobei  $\lambda_\ell$  jeweils den kleinsten Eigenwert der Matrix  $\mathbf{A}_\ell$  bezeichnet.

Wenn unser Algorithmus ausnutzen soll, dass die Matrizen zu unterschiedlichen Diskretisierungen desselben Differentialoperators gehören, müssen wir eine Möglichkeit haben, eine Beziehung zwischen diesen Diskretisierungen herzustellen. Dazu verwenden wir *Prolongationsmatrizen*

$$\mathbf{p}_\ell \in \mathbb{K}^{n_\ell \times n_{\ell-1}} \quad \text{für alle } \ell \in [1 : L],$$

die die Einbettung eines auf der Stufe  $\ell-1$  gegebenen Vektors in den Raum auf der Stufe  $\ell$  beschreiben. Die adjungierte Matrix  $\mathbf{p}_\ell^*$  wird häufig als *Restriktionsmatrix* bezeichnet.

Bei Finite-Elemente-Verfahren besitzen diese Matrizen in der Regel die *Galerkin-Eigenschaft*

$$\mathbf{A}_{\ell-1} = \mathbf{p}_\ell^* \mathbf{A}_\ell \mathbf{p}_\ell, \quad \mathbf{M}_{\ell-1} = \mathbf{p}_\ell^* \mathbf{M}_\ell \mathbf{p}_\ell \quad \text{für alle } \ell \in [1 : L], \quad (9.13)$$

wir können also die Matrizen der gröberen Diskretisierungen exakt durch ihre Gegenstücke der feineren Diskretisierungen darstellen.

Die grundlegende Idee der Mehrgitterverfahren besteht darin, ein *Glättungsverfahren* einzusetzen, das den Fehler zwischen einer Näherungslösung und dem gesuchten Eigenvektor glättet, so dass er mit einer gröberen Diskretisierung approximiert werden kann.

Ein einfaches Glättungsverfahren ist die uns bereits bekannte Richardson-Iteration, die im Fall eines verallgemeinerten Eigenwertproblems etwas angepasst werden muss.

## 9 Eigenwertverfahren für sehr große Matrizen

Sei  $\ell \in [1 : L]$ , und seien  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  die verallgemeinerten Eigenwerte zu den Matrizen  $\mathbf{A}_\ell$  und  $\mathbf{M}_\ell$ . Wir verwenden den Dämpfungsparameter

$$\theta_\ell := \frac{1}{\lambda_n}$$

und die angepasste Iterationsvorschrift

$$\mathbf{x}_\ell^{(m+1)} = \mathbf{x}_\ell^{(m)} - \theta_\ell \mathbf{M}_\ell^{-1} \left( \mathbf{A}_\ell \mathbf{x}_\ell^{(m)} - \mu \mathbf{M}_\ell \mathbf{x}_\ell^{(m)} \right).$$

Für einen Eigenvektor  $\mathbf{e}^{(i)}$  zu einem großen Eigenwert  $\lambda_i$ ,  $i \in [1 : n]$ , der also

$$\lambda_i - \mu \geq \frac{\lambda_n}{2}$$

erfüllt, erhalten wir

$$\begin{aligned} \mathbf{e}^{(i)} - \theta_\ell \mathbf{M}_\ell^{-1} \left( \mathbf{A}_\ell \mathbf{e}^{(i)} - \mu \mathbf{M}_\ell \mathbf{e}^{(i)} \right) &= \mathbf{e}^{(i)} - \theta_\ell \left( \lambda_i \mathbf{e}^{(i)} - \mu \mathbf{e}^{(i)} \right) = (1 - \theta_\ell (\lambda_i - \mu)) \mathbf{e}^{(i)} \\ &= \left( 1 - \frac{\lambda_i - \mu}{\lambda_n} \right) \mathbf{e}^{(i)} \leq \left( 1 - \frac{\lambda_n}{2\lambda_n} \right) \mathbf{e}^{(i)} = \frac{1}{2} \mathbf{e}^{(i)}, \end{aligned}$$

der Eigenvektor wird also mindestens um einen Faktor  $\frac{1}{2}$  verkürzt. Die Richardson-Iteration reduziert demnach diejenigen Anteile eines Vektors, die zu Eigenvektoren mit großen Eigenwerten gehören. Da wir ohnehin an kleinen Eigenwerten interessiert sind, ist uns diese Wirkung durchaus willkommen.

Bei partiellen Differentialgleichungen gehören große Eigenwerte zu Eigenvektoren, die stark oszillieren, so dass die Richardson-Iteration dazu führt, dass Vektoren in einem geeigneten Sinn „glatter“ werden. Das ist die Motivation der Bezeichnung „Glättungsverfahren“.

Um eine Näherung eines Eigenvektors zu finden, müssen wir allerdings auch glatte Anteile des Fehlers reduzieren. Wir nehmen an, dass wir mit dem Glättungsverfahren eine Näherung  $\mathbf{x}_\ell$  des gesuchten Eigenvektors berechnet haben. Da sie „glatt“ ist, dürfen wir davon ausgehen, dass der verbliebene Fehler sich auch in der größeren Diskretisierung auf der Stufe  $\ell - 1$  noch hinreichend gut approximieren lässt.

Für die Herleitung der Gleichung nehmen wir zunächst an, dass  $\mu = \lambda_\ell$  gilt, dass wir also als Shift den exakten Eigenwert verwenden. Dann erfüllt der gesuchte Eigenvektor  $\mathbf{e}_\ell$  die Gleichung

$$\mathbf{A}_\ell \mathbf{e}_\ell = \lambda_\ell \mathbf{M}_\ell \mathbf{e}_\ell \quad \iff \quad (\mathbf{A}_\ell - \mu \mathbf{M}_\ell) \mathbf{e}_\ell = \mathbf{0},$$

der verbliebene Fehler  $\mathbf{f}_\ell$  ist durch

$$(\mathbf{A}_\ell - \mu \mathbf{M}_\ell) \mathbf{f}_\ell = (\mathbf{A}_\ell - \mu \mathbf{M}_\ell) \mathbf{x}_\ell \tag{9.14}$$

bis auf Anteile aus dem Eigenraum beschrieben. Wie gesagt nehmen wir an, dass  $\mathbf{f}_\ell$  in einem geeigneten Sinn glatt ist, dass also ein Vektor  $\mathbf{x}_{\ell-1}$  mit

$$\mathbf{f}_\ell \approx \mathbf{p}_\ell \mathbf{x}_{\ell-1}$$

existiert. Unser Ziel ist es, diesen Vektor  $\mathbf{x}_{\ell-1}$  effizient zu berechnen, um dann mit seiner Hilfe die Näherung  $\mathbf{x}_\ell$  zu verbessern.

Aus (9.14) erhalten wir

$$(\mathbf{A}_\ell - \mu \mathbf{M}_\ell) \mathbf{p}_\ell \mathbf{x}_{\ell-1} \approx (\mathbf{A}_\ell - \mu \mathbf{M}_\ell) \mathbf{f}_\ell = (\mathbf{A}_\ell - \mu \mathbf{M}_\ell) \mathbf{x}_\ell.$$

Um  $\mathbf{x}_{\ell-1}$  als Lösung eines quadratischen Gleichungssystems berechnen zu können, multiplizieren wir diese Gleichung mit der *Restriktion*  $\mathbf{p}_\ell^*$  und erhalten mit der Galerkin-Eigenschaft (9.13) die Gleichung

$$(\mathbf{A}_{\ell-1} - \mu \mathbf{M}_{\ell-1}) \mathbf{x}_{\ell-1} = \mathbf{p}_\ell^* (\mathbf{A}_\ell - \mu \mathbf{M}_\ell) \mathbf{p}_\ell \mathbf{x}_{\ell-1} \approx \mathbf{p}_\ell^* (\mathbf{A}_\ell - \mu \mathbf{M}_\ell) \mathbf{x}_\ell.$$

Wenn wir also

$$\mathbf{b}_{\ell-1} := \mathbf{p}_\ell^* (\mathbf{A}_\ell - \mu \mathbf{M}_\ell) \mathbf{x}_\ell$$

setzen, müssen wir das lineare Gleichungssystem

$$(\mathbf{A}_{\ell-1} - \mu \mathbf{M}_{\ell-1}) \mathbf{x}_{\ell-1} = \mathbf{b}_{\ell-1} \quad (9.15)$$

lösen. Dabei ergibt sich eine kleine Schwierigkeit: Die meisten Lösungsverfahren für lineare Gleichungssysteme reagieren empfindlich, wenn die Matrix fast singulär ist, und in unserem Fall müssen wir annehmen, dass die kleinen Eigenwerte auf der Stufe  $\ell - 1$  nahe an den kleinen Eigenwerten auf der Stufe  $\ell$  liegen, dass also  $\mathbf{A}_{\ell-1} - \mu \mathbf{M}_{\ell-1}$  in der Tat fast singulär ist.

Je nach Wahl des Lösungsverfahrens kann das dazu führen, dass die Lösung durch Anteile aus dem Eigenraum des betragskleinsten Eigenwerts „verunreinigt“ wird, das wären in diesem Fall Vielfache des Eigenvektors  $\mathbf{e}_{\ell-1}$ .

Wir lösen das Problem, indem wir es zunächst verschlimmern: Auf Stufe  $\ell - 1$  ersetzen wir  $\mu$  durch den exakten Eigenwert  $\lambda_{\ell-1}$ , so dass die Matrix  $\mathbf{A}_{\ell-1} - \lambda_{\ell-1} \mathbf{M}_{\ell-1}$  nun singulär ist. Allerdings kennen wir ihren Kern explizit, es ist gerade der Aufspann des Eigenvektors  $\mathbf{e}_{\ell-1}$ , so dass wir trotzdem ein lösbares lineares Gleichungssystem erreichen können, indem wir die rechte Seite  $\mathbf{b}_{\ell-1}$  in das Bild der Matrix projizieren.

**Lemma 9.14 (Projektionen)** *Wir definieren die Matrix*

$$\mathbf{N}_\ell := \mathbf{I} - \frac{\mathbf{e}_\ell \mathbf{e}_\ell^* \mathbf{M}_\ell}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle}, \quad \mathbf{N}_\ell^* = \mathbf{I} - \frac{\mathbf{M}_\ell \mathbf{e}_\ell \mathbf{e}_\ell^*}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle}.$$

*Sie erfüllt die Gleichungen*

$$\langle \mathbf{e}_\ell, \mathbf{N}_\ell \mathbf{x}_\ell \rangle_{M_\ell} = 0 \quad \text{für alle } \mathbf{x}_\ell \in \mathbb{K}^{n_\ell}, \quad (9.16a)$$

$$\text{Bild}(\mathbf{N}_\ell^*) = \text{Bild}(\mathbf{A}_\ell - \lambda_\ell \mathbf{M}_\ell). \quad (9.16b)$$

$$\mathbf{N}_\ell^2 = \mathbf{N}_\ell, \quad (9.16c)$$

$$\langle \mathbf{x}_\ell, \mathbf{N}_\ell \mathbf{y}_\ell \rangle_{M_\ell} = \langle \mathbf{N}_\ell \mathbf{x}_\ell, \mathbf{y}_\ell \rangle_{M_\ell} \quad \text{für alle } \mathbf{x}_\ell, \mathbf{y}_\ell \in \mathbb{K}^{n_\ell}. \quad (9.16d)$$

## 9 Eigenwertverfahren für sehr große Matrizen

*Beweis.* Sei  $\mathbf{x}_\ell \in \mathbb{K}^{n_\ell}$ . Dann gilt

$$\begin{aligned} \langle \mathbf{e}_\ell, \mathbf{N}_\ell \mathbf{y}_\ell \rangle_{M_\ell} &= \langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{N}_\ell \mathbf{y}_\ell \rangle = \langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{y}_\ell \rangle - \langle \mathbf{e}_\ell, \mathbf{M}_\ell \frac{\mathbf{e}_\ell \mathbf{e}_\ell^* \mathbf{M}_\ell}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} \mathbf{y}_\ell \rangle \\ &= \langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{y}_\ell \rangle - \langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle \frac{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{y}_\ell \rangle}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} = \langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{y}_\ell \rangle - \langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{y}_\ell \rangle = 0, \end{aligned}$$

also (9.16a). Mit Lemma 3.17 und  $\mathbf{M}_\ell^* = \mathbf{M}_\ell$  erhalten wir die Gleichung

$$\begin{aligned} \langle \mathbf{N}_\ell^* \mathbf{x}_\ell, \mathbf{e}_\ell \rangle &= \langle \mathbf{x}_\ell, \mathbf{e}_\ell \rangle - \langle \frac{\mathbf{M}_\ell \mathbf{e}_\ell \mathbf{e}_\ell^*}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} \mathbf{x}_\ell, \mathbf{e}_\ell \rangle \\ &= \langle \mathbf{x}_\ell, \mathbf{e}_\ell \rangle - \frac{\langle \mathbf{x}_\ell, \mathbf{e}_\ell \mathbf{e}_\ell^* \mathbf{M}_\ell \mathbf{e}_\ell \rangle}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} = \langle \mathbf{x}_\ell, \mathbf{e}_\ell \rangle - \langle \mathbf{x}_\ell, \mathbf{e}_\ell \rangle \frac{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} = 0, \end{aligned}$$

also steht das Bild senkrecht auf dem Eigenvektor  $\mathbf{e}_\ell$ . Für jeden Vektor  $\mathbf{x}_\ell \in \mathbb{K}^{n_\ell}$ , der diese Eigenschaft besitzt, also  $\langle \mathbf{x}_\ell, \mathbf{e}_\ell \rangle = 0$  erfüllt, gilt auch

$$\mathbf{N}_\ell^* \mathbf{x}_\ell = \mathbf{x}_\ell - \mathbf{M}_\ell \mathbf{e}_\ell \frac{\mathbf{e}_\ell^* \mathbf{x}_\ell}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} = \mathbf{x}_\ell - \mathbf{M}_\ell \mathbf{e}_\ell \frac{\langle \mathbf{e}_\ell, \mathbf{x}_\ell \rangle}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} = \mathbf{x}_\ell,$$

also besteht das Bild der Matrix  $\mathbf{N}_\ell^*$  aus allen Vektoren, die senkrecht auf  $\mathbf{e}_\ell$  stehen.

Sei nun  $\mathbf{x}_\ell \in \text{Bild}(\mathbf{A}_\ell - \lambda_\ell \mathbf{M}_\ell)$  gegeben. Dann existiert ein  $\mathbf{y}_\ell \in \mathbb{K}^{n_\ell}$  mit  $\mathbf{x}_\ell = (\mathbf{A}_\ell - \lambda_\ell \mathbf{M}_\ell) \mathbf{y}_\ell$ . Da  $\mathbf{A}_\ell$  und  $\mathbf{M}_\ell$  selbstadjungiert sind, folgt mit Lemma 3.17 und  $\lambda_\ell \in \mathbb{R}$

$$\langle \mathbf{x}_\ell, \mathbf{e}_\ell \rangle = \langle (\mathbf{A}_\ell - \lambda_\ell \mathbf{M}_\ell) \mathbf{y}_\ell, \mathbf{e}_\ell \rangle = \langle \mathbf{y}_\ell, (\mathbf{A}_\ell - \lambda_\ell \mathbf{M}_\ell) \mathbf{e}_\ell \rangle = \langle \mathbf{y}_\ell, \mathbf{0} \rangle = 0,$$

also  $\text{Bild}(\mathbf{A}_\ell - \lambda_\ell \mathbf{M}_\ell) \subseteq \text{Bild}(\mathbf{N}_\ell^*)$ . Nach dem Dimensionssatz gilt

$$\dim \text{Bild}(\mathbf{A}_\ell - \lambda_\ell \mathbf{M}_\ell) = n_\ell - \dim \text{Kern}(\mathbf{A}_\ell - \lambda_\ell \mathbf{M}_\ell) = n_\ell - 1 = \dim \text{Bild}(\mathbf{N}_\ell^*),$$

sind  $\text{Bild}(\mathbf{A}_\ell - \lambda_\ell \mathbf{M}_\ell)$  und  $\text{Bild}(\mathbf{N}_\ell^*)$  identisch und wir erhalten (9.16b). Aus

$$\begin{aligned} \mathbf{N}_\ell^2 &= \mathbf{I} - 2 \frac{\mathbf{e}_\ell \mathbf{e}_\ell^* \mathbf{M}_\ell}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} + \frac{\mathbf{e}_\ell \mathbf{e}_\ell^* \mathbf{M}_\ell \mathbf{e}_\ell \mathbf{e}_\ell^* \mathbf{M}_\ell}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle^2} \\ &= \mathbf{I} - 2 \frac{\mathbf{e}_\ell \mathbf{e}_\ell^* \mathbf{M}_\ell}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} + \frac{\mathbf{e}_\ell \langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle \mathbf{e}_\ell^* \mathbf{M}_\ell}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle^2} = \mathbf{I} - 2 \frac{\mathbf{e}_\ell \mathbf{e}_\ell^* \mathbf{M}_\ell}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} + \frac{\mathbf{e}_\ell \mathbf{e}_\ell^* \mathbf{M}_\ell}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} = \mathbf{N}_\ell \end{aligned}$$

folgt direkt (9.16c). Für  $\mathbf{x}_\ell, \mathbf{y}_\ell \in \mathbb{K}^{n_\ell}$  haben wir mit Lemma 3.17 die Gleichung

$$\begin{aligned} \langle \mathbf{x}_\ell, \mathbf{N}_\ell \mathbf{y}_\ell \rangle_{M_\ell} &= \langle \mathbf{x}_\ell, \mathbf{M}_\ell \mathbf{y}_\ell \rangle - \langle \mathbf{x}_\ell, \mathbf{M}_\ell \frac{\mathbf{e}_\ell \mathbf{e}_\ell^* \mathbf{M}_\ell}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} \mathbf{y}_\ell \rangle \\ &= \langle \mathbf{x}_\ell, \mathbf{M}_\ell \mathbf{y}_\ell \rangle - \langle \frac{\mathbf{e}_\ell \mathbf{e}_\ell^*}{\langle \mathbf{e}_\ell, \mathbf{M}_\ell \mathbf{e}_\ell \rangle} \mathbf{M}_\ell \mathbf{x}_\ell, \mathbf{M}_\ell \mathbf{y}_\ell \rangle \\ &= \langle \mathbf{N}_\ell \mathbf{x}_\ell, \mathbf{M}_\ell \mathbf{y}_\ell \rangle = \langle \mathbf{N}_\ell \mathbf{x}_\ell, \mathbf{y}_\ell \rangle_{M_\ell}, \end{aligned}$$

also ist auch (9.16d) bewiesen. ■

Die Gleichung (9.16d) besagt, dass  $\mathbf{N}_\ell$  bezüglich des Massen-Skalarprodukts selbstadjungiert ist. Wegen (9.16c) handelt es sich also bezüglich dieses Skalarprodukts um eine orthogonale Projektion. Das ist sinnvoll, denn die Eigenvektoren unseres verallgemeinerten Eigenwertproblems bilden bezüglich des Massen-Skalarprodukts eine Orthonormalbasis, nicht aber bezüglich des euklidischen.

Mit Hilfe dieser Projektionen können wir das Gleichungssystem (9.15) so modifizieren, dass es lösbar wird: Wir ersetzen  $\mathbf{b}_{\ell-1}$  durch die Projektion  $\mathbf{N}_{\ell-1}^* \mathbf{b}_{\ell-1}$  und erhalten

$$(\mathbf{A}_{\ell-1} - \lambda_{\ell-1} \mathbf{M}_{\ell-1}) \mathbf{x}_{\ell-1} = \mathbf{N}_{\ell-1}^* \mathbf{b}_{\ell-1}.$$

Damit haben wir zwar mit (9.16b) die Existenz einer Lösung  $\mathbf{x}_{\ell-1}$  sicher gestellt, nicht aber deren Eindeutigkeit. Da sich unterschiedliche Lösungen des Gleichungssystems nur durch einen Vektor aus dem Kern der Matrix  $\mathbf{A}_{\ell-1} - \lambda_{\ell-1} \mathbf{M}_{\ell-1}$  unterscheiden können, also ein Vielfaches des Eigenvektors  $\mathbf{e}_{\ell-1}$ , bietet es sich an, mit der Matrix  $\mathbf{N}_{\ell-1}$  dafür zu sorgen, dass die berechnete Lösung im Sinn der Gleichung (9.16a) senkrecht auf dem Kern steht. Da wir davon ausgehen dürfen, dass der Eigenvektor  $\mathbf{e}_{\ell-1}$  auf der Stufe  $\ell-1$  eine Näherung des gesuchten Eigenvektors  $\mathbf{e}_\ell$  auf der Stufe  $\ell$  ist, stehen die Chancen gut, dass die Korrektur

$$\mathbf{x}_\ell \leftarrow \mathbf{x}_\ell - \mathbf{p}_\ell \mathbf{x}_{\ell-1}$$

Anteile des Eigenvektors  $\mathbf{e}_\ell$  in  $\mathbf{x}_\ell$  nicht wesentlich verändert.

Falls wir davon ausgehen, dass  $\mathbf{x}_\ell$  bereits eine passable Näherung des gesuchten Eigenvektors ist, können wir den Rayleigh-Quotienten verwenden, um den Shift-Parameter  $\mu$  zu bestimmen. Da bei dessen Berechnung ohnehin auch das Massen-Skalarprodukt berechnet wird, können wir auch gleich dafür sorgen, dass wir mit bezüglich dieses Skalarprodukts normierten Vektoren arbeiten.

**Algorithmus 9.15 (Eigenwert-Zweigitterverfahren)** *Unter der Annahme, dass  $\mathbf{x}_\ell$  eine Näherung des gesuchten Eigenvektors auf Stufe  $\ell$  ist und auf Stufe  $\ell-1$  der exakte Eigenwert  $\lambda_{\ell-1}$  und ein Eigenvektor  $\mathbf{e}_{\ell-1}$  mit  $\|\mathbf{e}_{\ell-1}\|_{M_{\ell-1}} = 1$  bekannt sind, führt der folgende Algorithmus einen Schritt des Eigenwert-Zweigitterverfahrens aus.*

```

 $\alpha \leftarrow \langle \mathbf{x}_\ell, \mathbf{M}_\ell \mathbf{x}_\ell \rangle;$ 
 $\mathbf{x}_\ell \leftarrow \mathbf{x}_\ell / \alpha;$ 
 $\lambda_\ell \leftarrow \langle \mathbf{x}_\ell, \mathbf{A}_\ell \mathbf{x}_\ell \rangle;$ 
for  $i = 1$  to  $\nu$  do
     $\mathbf{x}_\ell \leftarrow \mathbf{x}_\ell - \theta_\ell \mathbf{M}_\ell^{-1} (\mathbf{A}_\ell \mathbf{x}_\ell - \lambda_\ell \mathbf{M}_\ell \mathbf{x}_\ell);$ 
 $\mathbf{d}_\ell \leftarrow \mathbf{A}_\ell \mathbf{x}_\ell - \lambda_\ell \mathbf{M}_\ell \mathbf{x}_\ell;$ 
 $\mathbf{b}_{\ell-1} \leftarrow \mathbf{p}_\ell^* \mathbf{d}_\ell;$ 
 $\mathbf{b}_{\ell-1} \leftarrow \mathbf{b}_{\ell-1} - \langle \mathbf{e}_{\ell-1}, \mathbf{b}_{\ell-1} \rangle \mathbf{M}_{\ell-1} \mathbf{e}_{\ell-1};$     { =  $\mathbf{N}_{\ell-1}^* \mathbf{b}_{\ell-1}$  }
    Löse  $(\mathbf{A}_{\ell-1} \mathbf{x}_{\ell-1} - \lambda_{\ell-1} \mathbf{M}_{\ell-1} \mathbf{x}_{\ell-1}) \mathbf{x}_{\ell-1} = \mathbf{b}_{\ell-1};$ 
 $\mathbf{x}_{\ell-1} \leftarrow \mathbf{x}_{\ell-1} - \mathbf{e}_{\ell-1} \langle \mathbf{e}_{\ell-1}, \mathbf{M}_{\ell-1} \mathbf{x}_{\ell-1} \rangle;$     { =  $\mathbf{N}_{\ell-1} \mathbf{x}_{\ell-1}$  }
 $\mathbf{x}_\ell \leftarrow \mathbf{x}_\ell - \mathbf{p}_\ell \mathbf{x}_{\ell-1}$ 
    
```

## 9 Eigenwertverfahren für sehr große Matrizen

Das Lösen eines Gleichungssystems auf der Stufe  $\ell - 1$  wird in der Regel immer noch sehr aufwendig sein, deshalb verwenden wir eine Variante unserer bisherigen Vorgehensweise: Wir glätten den Fehler mit einer Richardson-Iteration

$$\mathbf{x}_{\ell-1}^{(m+1)} = \mathbf{x}_{\ell-1}^{(m)} - \theta_{\ell-1} \mathbf{M}_{\ell-1}^{-1} \left( \mathbf{A}_{\ell-1} \mathbf{x}_{\ell-1}^{(m)} - \lambda_{\ell-1} \mathbf{M}_{\ell-1} \mathbf{x}_{\ell-1}^{(m)} - \mathbf{b}_{\ell-1} \right)$$

und approximieren den Rest dann auf der nächstgrößeren Stufe  $\ell - 2$ . In dieser Weise können wir rekursiv fortfahren, bis wir die größte Stufe  $\ell = 0$  erreichen, auf der wir es uns leisten können, das hoffentlich kleine verbliebene Gleichungssystem direkt zu lösen.

**Algorithmus 9.16 (Singuläres Mehrgitterverfahren)** *Unter der Annahme, dass auf allen größeren Stufen die exakten Eigenwerte und bezüglich der Massen-Norm normierte Eigenvektoren vorliegen, führt der folgende Algorithmus einen Schritt des singulären Mehrgitterverfahrens aus. Dabei gibt  $\gamma \in \mathbb{N}$  die Anzahl der rekursiven Aufrufe an, die an die Stelle des exakten Lösens der Grobgittergleichung treten.*

```

procedure SMG( $\ell$ );
  if  $\ell = 0$  then
    Löse  $(\mathbf{A}_0 - \lambda_0 \mathbf{M}_0) \mathbf{x}_0 = \mathbf{b}_0$ 
  else begin
    for  $i = 1$  to  $\nu$  do
       $\mathbf{x}_\ell \leftarrow \mathbf{x}_\ell - \theta_\ell \mathbf{M}_\ell^{-1} (\mathbf{A}_\ell \mathbf{x}_\ell - \lambda_\ell \mathbf{M}_\ell \mathbf{x}_\ell - \mathbf{b}_\ell)$ ;
       $\mathbf{d}_\ell \leftarrow \mathbf{A}_\ell \mathbf{x}_\ell - \lambda_\ell \mathbf{M}_\ell \mathbf{x}_\ell - \mathbf{b}_\ell$ ;
       $\mathbf{b}_{\ell-1} \leftarrow \mathbf{p}_\ell^* \mathbf{d}_\ell$ ;
       $\mathbf{b}_{\ell-1} \leftarrow \mathbf{b}_{\ell-1} - \langle \mathbf{e}_{\ell-1}, \mathbf{b}_{\ell-1} \rangle \mathbf{M}_{\ell-1} \mathbf{e}_{\ell-1}$ ;      { =  $\mathbf{N}_{\ell-1}^* \mathbf{b}_{\ell-1}$  }
       $\mathbf{x}_{\ell-1} \leftarrow \mathbf{0}$ ;
      for  $i = 1$  to  $\gamma$  do
        SMG( $\ell - 1$ );
       $\mathbf{x}_{\ell-1} \leftarrow \mathbf{x}_{\ell-1} - \mathbf{e}_{\ell-1} \langle \mathbf{e}_{\ell-1}, \mathbf{M}_{\ell-1} \mathbf{x}_{\ell-1} \rangle$ ;      { =  $\mathbf{N}_{\ell-1} \mathbf{x}_{\ell-1}$  }
       $\mathbf{x}_\ell \leftarrow \mathbf{x}_\ell - \mathbf{p}_\ell \mathbf{x}_{\ell-1}$ 
  end

```

Beide Algorithmen gehen davon aus, dass auf allen größeren Stufen exakte Eigenwerte und Eigenvektoren bekannt sind. Dieses Ziel lässt sich näherungsweise mit Hilfe einer *geschachtelten Iteration* erreichen: Wir berechnen zunächst den Eigenwert und Eigenvektor exakt auf der größten Stufe  $\ell = 0$ .

Dann können wir die Eigenwert-Iteration einsetzen, um Näherungen auf der Stufe  $\ell = 1$  zu ermitteln. Sobald sie gut genug sind, können wir zu der Stufe  $\ell = 2$  wechseln. Wenn wir dabei den jeweiligen Startvektor per Prolongation aus der Näherung der vorigen Stufe berechnen, stehen die Chancen gut, dass das Verfahren schnell konvergiert und deshalb auf jeder Stufe nur wenige Schritte erforderlich sind.



# 10 Verwandte Fragestellungen

Bisher haben wir uns ausschließlich mit der Frage nach der Berechnung einzelner, mehrerer oder aller Eigenwerte und eventuell der zugehörigen Eigenvektoren befasst. In diesem Kapitel beschäftigen wir uns mit Problemen, die eng mit Eigenwertproblemen verwandt sind, aber modifizierte Lösungsverfahren erfordern.

## 10.1 Verallgemeinerte Eigenwertprobleme

Für zwei Matrizen  $\mathbf{A}, \mathbf{B} \in \mathbb{K}^{n \times n}$  können wir uns die Frage stellen, ob Zahlen  $\lambda \in \mathbb{K}$  und Vektoren  $\mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}$  existieren, die die Gleichung

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{B}\mathbf{x} \quad (10.1)$$

erfüllen. Diese Aufgabe bezeichnet man als *verallgemeinertes Eigenwertproblem*. Falls  $\mathbf{B}$  invertierbar ist, können wir durch Multiplikation mit  $\mathbf{B}^{-1}$  zu dem gewöhnlichen Eigenwertproblem

$$\mathbf{B}^{-1}\mathbf{A}\mathbf{x} = \lambda\mathbf{x} \quad (10.2)$$

übergehen und die bereits diskutierten Verfahren einsetzen.

Von Interesse sind deshalb zwei Typen verallgemeinerter Eigenwertprobleme: Einerseits diejenigen, bei denen  $\mathbf{B}$  nicht invertierbar ist, und andererseits diejenigen, bei denen  $\mathbf{A}$  selbstadjungiert ist und diese nützliche Eigenschaft erhalten bleiben sollte. Diesem zweiten Fall widmen wir uns in einem separaten Abschnitt.

Wir untersuchen zunächst den Fall, dass  $\mathbf{B}$  nicht invertierbar oder sehr schlecht konditioniert ist. In diesem Fall ist man daran interessiert, eine *verallgemeinerte Schur-Zerlegung* zu berechnen, nämlich unitäre Matrizen  $\mathbf{Q}, \mathbf{P} \in \mathbb{K}^{n \times n}$ , für die die Matrizen

$$\hat{\mathbf{A}} := \mathbf{Q}\mathbf{A}\mathbf{P}^*, \quad \hat{\mathbf{B}} := \mathbf{Q}\mathbf{B}\mathbf{P}^*$$

obere Dreiecksmatrizen sind. An den Diagonalelementen der beiden transformierten Matrizen lassen sich dann die verallgemeinerten Eigenwerte direkt ablesen, durch Rückwärtseinsetzen in  $\hat{\mathbf{A}} - \lambda\hat{\mathbf{B}}$  können wir auch Eigenvektoren bestimmen.

Bei der Berechnung der verallgemeinerten Schur-Zerlegung können wir uns an der Vorgehensweise für die gewöhnliche Schur-Zerlegung orientieren: Mit Hilfe geeigneter unitärer Matrizen  $\mathbf{Q}$  und  $\mathbf{P}$  können wir das Eigenwertproblem auf die Form

$$\hat{\mathbf{A}}\mathbf{x} = \lambda\hat{\mathbf{B}}\mathbf{x}, \quad \hat{\mathbf{A}} = \mathbf{Q}\mathbf{A}\mathbf{P}^*, \quad \hat{\mathbf{B}} = \mathbf{Q}\mathbf{B}\mathbf{P}^*$$

bringen, bei der  $\hat{\mathbf{B}}$  bereits eine rechte obere Dreiecksmatrix ist,  $\hat{\mathbf{A}}$  allerdings nur eine Hessenberg-Matrix.

### Implizite Iteration für reguläres $\mathbf{B}$

Die Dreiecksstruktur der Matrix  $\widehat{\mathbf{B}}$  lässt sich ausnutzen, um ähnlich wie im Fall der Berechnung der Singulärwertzerlegung vorzugehen: Falls  $\widehat{\mathbf{B}}$  regulär ist, könnte man *implizit* das Gegenstück der Formulierung (10.2) behandeln, also die Schur-Zerlegung der Matrix  $\mathbf{H}^{(0)} = \widehat{\mathbf{B}}^{-1}\widehat{\mathbf{A}}$  berechnen, indem man die in der QR-Iteration auftretenden Transformationen  $\mathbf{Q}^{(0)}, \mathbf{Q}^{(1)}, \dots$  anwendet:

$$\mathbf{H}^{(0)} := \widehat{\mathbf{B}}^{-1}\widehat{\mathbf{A}}, \quad \mathbf{H}^{(m+1)} := (\mathbf{Q}^{(m)})^* \mathbf{H}^{(m)} \mathbf{Q}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0.$$

Da  $\widehat{\mathbf{B}}$  schlecht konditioniert sein kann, wollen wir die Multiplikation mit  $\widehat{\mathbf{B}}^{-1}$  vermeiden. Wie bei der Singulärwertzerlegung definieren wir

$$\begin{aligned} \widehat{\mathbf{A}}^{(0)} &:= \widehat{\mathbf{A}}, & \widehat{\mathbf{A}}^{(m+1)} &:= (\mathbf{P}^{(m)})^* \widehat{\mathbf{A}}^{(m)} \mathbf{Q}^{(m)}, \\ \widehat{\mathbf{B}}^{(0)} &:= \widehat{\mathbf{B}}, & \widehat{\mathbf{B}}^{(m+1)} &:= (\mathbf{P}^{(m)})^* \widehat{\mathbf{B}}^{(m)} \mathbf{Q}^{(m)} \end{aligned} \quad \text{für alle } m \in \mathbb{N}_0$$

mit später noch genauer zu spezifizierenden unitären Matrizen  $\mathbf{P}^{(m)}$  und erhalten wegen

$$\mathbf{H}^{(0)} = (\widehat{\mathbf{B}}^{(0)})^{-1} \widehat{\mathbf{A}}^{(0)}$$

und

$$\begin{aligned} \mathbf{H}^{(m+1)} &= (\mathbf{Q}^{(m)})^* \mathbf{H}^{(m)} \mathbf{Q}^{(m)} = (\mathbf{Q}^{(m)})^* (\widehat{\mathbf{B}}^{(m)})^{-1} \widehat{\mathbf{A}}^{(m)} \mathbf{Q}^{(m)} \\ &= (\mathbf{Q}^{(m)})^* (\widehat{\mathbf{B}}^{(m)})^{-1} \mathbf{P}^{(m)} (\mathbf{P}^{(m)})^* \widehat{\mathbf{A}}^{(m)} \mathbf{Q}^{(m)} \\ &= ((\mathbf{P}^{(m)})^* \widehat{\mathbf{B}} \mathbf{Q}^{(m)})^{-1} ((\mathbf{P}^{(m)})^* \widehat{\mathbf{A}} \mathbf{Q}^{(m)}) = (\widehat{\mathbf{B}}^{(m+1)})^{-1} \widehat{\mathbf{A}}^{(m+1)} \end{aligned} \quad \text{für alle } m \in \mathbb{N}_0$$

mit einer einfachen Induktion die faktorisierte Darstellung

$$\mathbf{H}^{(m)} = (\widehat{\mathbf{B}}^{(m)})^{-1} \widehat{\mathbf{A}}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0. \quad (10.3)$$

Wir können also die Matrizen  $\mathbf{H}^{(m)}$  während der gesamten Iteration in Produktform darstellen und müssen insbesondere niemals die vollständige Inverse der Matrizen  $\widehat{\mathbf{B}}^{(m)}$  berechnen. Allerdings müssen wir darauf achten, dass während der Iteration weder die Dreiecksform der Matrizen  $\widehat{\mathbf{B}}^{(m)}$  noch die Hessenberg-Form der Matrizen  $\widehat{\mathbf{A}}^{(m)}$  verloren geht. Dazu greifen wir wieder auf Satz 6.16 zurück: Wir führen die erste Givens-Rotation der QR-Zerlegung der Matrix  $\mathbf{H}^{(m)} - \mu \mathbf{I}$  mit einem geeigneten Shift-Parameter  $\mu$  durch und sorgen anschließend mit weiteren Givens-Rotationen dafür, dass die Matrizen wieder die vorgesehene Form haben. Mit Satz 6.16 folgt dann, dass sich das Resultat allenfalls im Vorzeichen von dem der ursprünglichen QR-Iteration unterscheidet. Wir stellen die ursprünglichen Matrizen in der Form

$$\widehat{\mathbf{A}}^{(m)} = \begin{pmatrix} a_{11} & a_{12} & \dots & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ & a_{32} & a_{33} & \dots & a_{3n} \\ & & \ddots & \ddots & \vdots \\ & & & a_{n,n-1} & a_{nn} \end{pmatrix}, \quad \widehat{\mathbf{B}}^{(m)} = \begin{pmatrix} b_{11} & b_{12} & \dots & \dots & b_{1n} \\ & b_{22} & b_{23} & \dots & b_{2n} \\ & & b_{33} & \dots & b_{3n} \\ & & & \ddots & \vdots \\ & & & & b_{nn} \end{pmatrix}$$

dar. Die erste Givens-Rotation bei der Berechnung der QR-Zerlegung von  $\mathbf{H}^{(m)}$  wirkt auf die ersten beiden Spalten der Matrizen, die die Form

$$\begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \cdots & \cdots & a_{1n} \\ a_{21}^{(1)} & a_{22}^{(1)} & a_{23} & \cdots & a_{2n} \\ \gamma_1 & a_{32}^{(1)} & a_{33} & \cdots & a_{3n} \\ & & \ddots & \ddots & \vdots \\ & & & a_{n,n-1} & a_{nn} \end{pmatrix}, \quad \begin{pmatrix} b_{11}^{(1)} & b_{12}^{(1)} & \cdots & \cdots & b_{1n} \\ \delta_1 & b_{22}^{(1)} & b_{23} & \cdots & b_{2n} \\ & & b_{33} & \cdots & b_{3n} \\ & & & \ddots & \vdots \\ & & & & b_{nn} \end{pmatrix}$$

annehmen. Indem wir eine passende Givens-Rotation auf die beiden ersten Zeilen anwenden, können wir  $\delta_1$  eliminieren:

$$\begin{pmatrix} a_{11}^{(2)} & a_{12}^{(2)} & \cdots & \cdots & a_{1n}^{(2)} \\ a_{21}^{(2)} & a_{22}^{(2)} & a_{23} & \cdots & a_{2n}^{(2)} \\ \gamma_1 & a_{32}^{(1)} & a_{33} & \cdots & a_{3n} \\ & & \ddots & \ddots & \vdots \\ & & & a_{n,n-1} & a_{nn} \end{pmatrix}, \quad \begin{pmatrix} b_{11}^{(2)} & b_{12}^{(2)} & \cdots & \cdots & b_{1n}^{(2)} \\ & b_{22}^{(2)} & b_{23} & \cdots & b_{2n}^{(2)} \\ & & b_{33} & \cdots & b_{3n} \\ & & & \ddots & \vdots \\ & & & & b_{nn} \end{pmatrix}.$$

Den Eintrag  $\gamma_1$  eliminieren wir mit einer Givens-Rotation, die auf die zweite und dritte Zeile wirkt, so dass sich

$$\begin{pmatrix} a_{11}^{(2)} & a_{12}^{(2)} & \cdots & \cdots & a_{1n}^{(2)} \\ a_{21}^{(3)} & a_{22}^{(3)} & a_{23} & \cdots & a_{2n}^{(3)} \\ & a_{32}^{(3)} & a_{33} & \cdots & a_{3n}^{(3)} \\ & & \ddots & \ddots & \vdots \\ & & & a_{n,n-1} & a_{nn} \end{pmatrix}, \quad \begin{pmatrix} b_{11}^{(2)} & b_{12}^{(2)} & \cdots & \cdots & b_{1n}^{(2)} \\ & b_{22}^{(3)} & b_{23} & \cdots & b_{2n}^{(3)} \\ & \delta_2 & b_{33} & \cdots & b_{3n}^{(3)} \\ & & & \ddots & \vdots \\ & & & & b_{nn} \end{pmatrix}$$

ergibt. Wir beseitigen  $\delta_2$  mit einer Givens-Rotation, die auf die zweite und dritte Spalte wirkt und zu einem Eintrag  $\gamma_2$  in der zweiten Spalte der vierten Zeile der linken Matrix führt. Diesen Eintrag eliminieren wir mit einer Givens-Rotation, die auf die dritte und vierte Zeile wirkt und einen Eintrag  $\delta_3$  in die dritte Spalte der vierten Zeile der rechten Matrix erzeugt. Indem wir entsprechend fortfahren, können wir die störenden Einträge wieder „nach rechts unten aus den Matrizen heraus schieben“ und so die ursprüngliche Struktur der Matrizen wiederherstellen. Damit sind  $\widehat{\mathbf{A}}^{(m+1)}$  und  $\widehat{\mathbf{B}}^{(m+1)}$  gefunden. Da auf die erste Spalte der Matrizen lediglich die Givens-Rotation des ersten Schritts der QR-Zerlegung wirkt, lässt sich Satz 6.16 anwenden und folgern, dass die so berechneten Matrizen denen entsprechen, die im Rahmen der konventionellen QR-Iteration entstehen würden. Demzufolge ist zu erwarten, dass die Matrizen  $\mathbf{H}^{(m)}$  gegen obere Dreiecksform konvergieren werden, dass also die unteren Nebendiagonaleinträge gegen Null streben werden. Falls  $h_{k+1,k}^{(m)} = 0$  gilt, folgt aus (10.3)

$$0 = h_{k+1,k}^{(m)} = a_{k+1,k}^{(m)} / b_{k+1,k+1}^{(m)},$$

also muss  $a_{k+1,k}^{(m)} = 0$  gelten, wir hätten in diesem Fall also auch einen Nebendiagonaleintrag der Matrix  $\widehat{\mathbf{A}}^{(m)}$  eliminiert und damit einen Schritt in Richtung der gewünschten Dreiecksmatrix vollzogen.

### Deflation

Falls  $\hat{a}_{k+1,k}^{(m)} = 0$  für ein  $k \in \{1, \dots, n-1\}$  gelten sollte, können wir wie üblich eine Deflation durchführen und die Iteration mit kleineren Teilmatrizen fortführen.

Die Deflation kann uns aber auch dabei helfen, den Fall einer nicht invertierbaren Matrix  $\mathbf{B}$  zu behandeln: Falls  $\mathbf{B}$  nicht invertierbar ist, kann die Dreiecksmatrix  $\widehat{\mathbf{B}}$  ebenfalls nicht invertierbar sein. Bei einer Dreiecksmatrix ist das genau dann der Fall, wenn ein Diagonaleintrag gleich null ist, wir können diese Situation also sehr einfach erkennen.

Falls beispielsweise  $\hat{b}_{11}^{(m)} = 0$  gilt, sind wir in der Situation

$$\widehat{\mathbf{A}}^{(m)} = \begin{pmatrix} a_{11} & a_{12} & \dots & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ & a_{32} & a_{33} & \dots & a_{3n} \\ & & \ddots & \ddots & \vdots \\ & & & a_{n,n-1} & a_{nn} \end{pmatrix}, \quad \widehat{\mathbf{B}}^{(m)} = \begin{pmatrix} \mathbf{0} & b_{12} & \dots & \dots & b_{1n} \\ & b_{22} & b_{23} & \dots & b_{2n} \\ & & b_{33} & \dots & b_{3n} \\ & & & \ddots & \vdots \\ & & & & b_{nn} \end{pmatrix}$$

und können eine Givens-Rotation auf die ersten beiden Zeilen anwenden, um den Eintrag  $a_{21}$  zu eliminieren. Wir erhalten

$$\begin{pmatrix} a_{11}^{(1)} & a_{12}^{(1)} & \dots & \dots & a_{1n} \\ & a_{22}^{(1)} & a_{23} & \dots & a_{2n} \\ & a_{32}^{(1)} & a_{33} & \dots & a_{3n} \\ & & \ddots & \ddots & \vdots \\ & & & a_{n,n-1} & a_{nn} \end{pmatrix}, \quad \begin{pmatrix} \mathbf{0} & b_{12}^{(1)} & \dots & \dots & b_{1n} \\ & b_{22}^{(1)} & b_{23} & \dots & b_{2n} \\ & & b_{33} & \dots & b_{3n} \\ & & & \ddots & \vdots \\ & & & & b_{nn} \end{pmatrix}$$

Nun ist ein Nebendiagonalelement in der linken Matrix gleich null, so dass wir per Deflation zu den Teilmatrizen

$$\begin{pmatrix} a_{22}^{(1)} & a_{23} & \dots & a_{2n} \\ a_{32}^{(1)} & a_{33} & \dots & a_{3n} \\ & \ddots & \ddots & \vdots \\ & & a_{n,n-1} & a_{nn} \end{pmatrix}, \quad \begin{pmatrix} b_{22}^{(1)} & b_{23} & \dots & b_{2n} \\ b_{33} & \dots & b_{3n} \\ & \ddots & \vdots \\ & & b_{nn} \end{pmatrix}$$

übergehen können. Nicht invertierbare Matrizen können wir also nicht nur einfach bei unseren Berechnungen berücksichtigen, sie führen sogar dazu, dass wir besonders früh mit einer Deflation die Problemgröße reduzieren können.

Der aus den impliziten QR-Schritten und der Deflation entstehende Algorithmus ist unter dem Namen *QZ-Iteration* bekannt. Da ausschließlich unitäre Transformationen zum Einsatz kommen, ist er in der Praxis relativ unanfällig für Rundungsfehler und somit auch für eher schlecht konditionierte Eigenwertprobleme anwendbar.

## 10.2 Selbstadjungierte positiv definite verallgemeinerte Eigenwertprobleme

Wir wenden uns einem wichtigen Spezialfall unter den verallgemeinerten Eigenwertproblemen zu: Sei  $\mathbf{A} \in \mathbb{K}^{n \times n}$  selbstadjungiert, und sei  $\mathbf{B} \in \mathbb{K}^{n \times n}$  selbstadjungiert und positiv definit, es gelte also

$$0 < \langle \mathbf{B}\mathbf{x}, \mathbf{x} \rangle_2 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n \setminus \{\mathbf{0}\}.$$

Wir suchen nach Lösungen  $\lambda, \mathbf{x} \neq \mathbf{0}$  des verallgemeinerten Eigenwertproblems

$$\mathbf{A}\mathbf{x} = \lambda\mathbf{B}\mathbf{x}.$$

Unser Ziel besteht darin, dieses verallgemeinerte Eigenwertproblem auf ein gewöhnliches Eigenwertproblem zurückzuführen, ohne dabei die Selbstadjungiertheit der Matrix  $\mathbf{A}$  zu verlieren. Das Hilfsmittel der Wahl ist die *Cholesky-Zerlegung*

$$\mathbf{B} = \mathbf{L}\mathbf{L}^*$$

der Matrix  $\mathbf{B}$ , die für selbstadjungierte positiv definite Matrizen immer existiert und sich in der Praxis auch häufig stabil berechnen lässt. Wir erhalten

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \lambda\mathbf{B}\mathbf{x}, \\ \mathbf{A}\mathbf{x} &= \lambda\mathbf{L}\mathbf{L}^*\mathbf{x}, \\ \mathbf{L}\mathbf{L}^{-1}\mathbf{A}(\mathbf{L}^*)^{-1}\mathbf{L}^*\mathbf{x} &= \lambda\mathbf{L}\mathbf{L}^*\mathbf{x}, \end{aligned}$$

und durch Multiplikation mit  $\mathbf{L}^{-1}$  von links folgt

$$\mathbf{L}^{-1}\mathbf{A}(\mathbf{L}^*)^{-1}(\mathbf{L}^*\mathbf{x}) = \lambda\mathbf{L}^*\mathbf{x}.$$

Wir führen die Hilfsgrößen

$$\hat{\mathbf{x}} := \mathbf{L}^*\mathbf{x}, \quad \hat{\mathbf{A}} := \mathbf{L}^{-1}\mathbf{A}(\mathbf{L}^*)^{-1}$$

ein und erhalten das gewöhnliche Eigenwertproblem

$$\hat{\mathbf{A}}\hat{\mathbf{x}} = \lambda\hat{\mathbf{x}}. \quad (10.4)$$

Offenbar ist die Matrix  $\hat{\mathbf{A}}$  selbstadjungiert, es ist uns also gelungen, diese wichtige Eigenschaft zu erhalten.

Damit ist das Problem (10.4) unseren sämtlichen bisher untersuchten Verfahren zugänglich. Beispielsweise folgt aus Folgerung 3.54, dass eine unitäre Matrix  $\hat{\mathbf{Q}} \in \mathbb{K}^{n \times n}$  und eine reelle Diagonalmatrix  $\hat{\mathbf{D}} \in \mathbb{R}^{n \times n}$  mit

$$\hat{\mathbf{Q}}\hat{\mathbf{D}}\hat{\mathbf{Q}}^* = \hat{\mathbf{A}}$$

existieren. Indem wir die Transformation rückgängig machen, erhalten wir

$$\mathbf{Q} := (\mathbf{L}^*)^{-1}\hat{\mathbf{Q}}$$

und stellen fest, dass

$$\begin{aligned}\mathbf{Q}^* \mathbf{A} \mathbf{Q} &= \widehat{\mathbf{Q}}^* \mathbf{L}^{-1} \mathbf{A} (\mathbf{L}^*)^{-1} \widehat{\mathbf{Q}} = \widehat{\mathbf{Q}}^* \widehat{\mathbf{A}} \widehat{\mathbf{Q}} = \mathbf{D}, \\ \mathbf{Q}^* \mathbf{B} \mathbf{Q} &= \widehat{\mathbf{Q}}^* \mathbf{L}^{-1} \mathbf{L} \mathbf{L}^* (\mathbf{L}^*)^{-1} \widehat{\mathbf{Q}} = \widehat{\mathbf{Q}}^* \widehat{\mathbf{Q}} = \mathbf{I}\end{aligned}$$

gelten, dass also die Matrix  $\mathbf{Q}$  sowohl  $\mathbf{A}$  als auch  $\mathbf{B}$  auf Diagonalform transformiert.

Die Eigenwerte  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$  des transformierten Problems (und damit auch der verallgemeinerten Eigenwertproblems) lassen sich mit Hilfe des Rayleigh-Quotienten charakterisieren, der sich in der Form

$$\Lambda_{\widehat{\mathbf{A}}}(\widehat{\mathbf{x}}) = \frac{\langle \widehat{\mathbf{A}} \widehat{\mathbf{x}}, \widehat{\mathbf{x}} \rangle_2}{\langle \widehat{\mathbf{x}}, \widehat{\mathbf{x}} \rangle_2} = \frac{\langle \mathbf{L}^{-1} \mathbf{A} (\mathbf{L}^*)^{-1} \mathbf{L}^* \mathbf{x}, \mathbf{L}^* \mathbf{x} \rangle_2}{\langle \mathbf{L}^* \mathbf{x}, \mathbf{L}^* \mathbf{x} \rangle_2} = \frac{\langle \mathbf{A} \mathbf{x}, (\mathbf{L}^*)^{-1} \mathbf{L}^* \mathbf{x} \rangle_2}{\langle \mathbf{L} \mathbf{L}^* \mathbf{x}, \mathbf{x} \rangle_2} = \frac{\langle \mathbf{A} \mathbf{x}, \mathbf{x} \rangle_2}{\langle \mathbf{B} \mathbf{x}, \mathbf{x} \rangle_2}$$

darstellen lässt. Diese Eigenschaft ist vor allem für große verallgemeinerte Eigenwertproblem wichtig, da es mit ihrer Hilfe häufig möglich ist, auf die zeitaufwendige Berechnung der transformierten Matrix  $\widehat{\mathbf{A}}$  zu verzichten und direkt mit  $\mathbf{A}$  und  $\mathbf{B}$  zu arbeiten.

Mit kleinen Modifikationen lassen sich auch andere Verfahren durchführen, ohne explizit mit  $\widehat{\mathbf{A}}$  arbeiten zu müssen. Ein Beispiel ist der Lanczos-Algorithmus 8.11, der sich elegant formulieren lässt, wenn man das euklidische Skalarprodukt durch das zu der Matrix  $\mathbf{B}$  gehörende *Energie-Skalarprodukt* ersetzt, das durch

$$\langle \mathbf{x}, \mathbf{y} \rangle_B := \langle \mathbf{B} \mathbf{x}, \mathbf{y} \rangle_2 \quad \text{für alle } \mathbf{x}, \mathbf{y} \in \mathbb{K}^n$$

definiert ist. Es induziert die *Energienorm*

$$\|\mathbf{x}\|_B := \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle_B} \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n.$$

Wir bezeichnen mit

$$\widehat{\mathbf{p}}^{(m)} := \mathbf{L}^* \mathbf{p}^{(m)}, \quad \widehat{\mathbf{q}}^{(m)} := \mathbf{L}^* \mathbf{q}^{(m)} \quad \text{für alle } m \in \mathbb{N}$$

die transformierten Vektoren und untersuchen die im Rahmen des Algorithmus 8.11 auftretenden Gleichungen. Die Berechnung von

$$\begin{aligned}\gamma &= \|\widehat{\mathbf{p}}^{(m)}\|_2 = \sqrt{\langle \widehat{\mathbf{p}}^{(m)}, \widehat{\mathbf{p}}^{(m)} \rangle_2} = \sqrt{\langle \mathbf{L}^* \mathbf{p}^{(m)}, \mathbf{L}^* \mathbf{p}^{(m)} \rangle_2} \\ &= \sqrt{\langle \mathbf{L} \mathbf{L}^* \mathbf{p}^{(m)}, \mathbf{p}^{(m)} \rangle_2} = \sqrt{\langle \mathbf{B} \mathbf{p}^{(m)}, \mathbf{p}^{(m)} \rangle_2} = \|\mathbf{p}^{(m)}\|_B\end{aligned}$$

lässt sich auf die Energienorm zurückführen, und ebenso die von

$$\beta_m = \langle \widehat{\mathbf{q}}^{(m)}, \widehat{\mathbf{p}}^{(m)} \rangle_2 = \langle \mathbf{L}^* \mathbf{q}^{(m)}, \mathbf{L}^* \mathbf{p}^{(m)} \rangle_2 = \langle \mathbf{L} \mathbf{L}^* \mathbf{q}^{(m)}, \mathbf{p}^{(m)} \rangle_2 = \langle \mathbf{q}^{(m)}, \mathbf{p}^{(m)} \rangle_B$$

auf das Energie-Skalarprodukt. Für die Berechnung von  $\widehat{\mathbf{p}}^{(m)}$  erhalten wir

$$\widehat{\mathbf{p}}^{(m)} = \widehat{\mathbf{A}} \widehat{\mathbf{q}}^{(m)} = \mathbf{L}^{-1} \mathbf{A} (\mathbf{L}^*)^{-1} \mathbf{L}^* \mathbf{q}^{(m)} = \mathbf{L}^{-1} \mathbf{A} \mathbf{q}^{(m)},$$

so dass sich

$$\mathbf{p}^{(m)} = (\mathbf{L}^*)^{-1} \widehat{\mathbf{p}}^{(m)} = (\mathbf{L}^*)^{-1} \mathbf{L}^{-1} \mathbf{A} \mathbf{q}^{(m)} = (\mathbf{L} \mathbf{L}^*)^{-1} \mathbf{A} \mathbf{q}^{(m)} = \mathbf{B}^{-1} \mathbf{A} \mathbf{q}^{(m)}$$

ergibt. Anders als im Fall des ursprünglichen Lanczos-Algorithmus müssen wir also für das verallgemeinerte Eigenwertproblem dazu in der Lage sein, Gleichungssysteme der Form

$$\mathbf{B}\mathbf{p}^{(m)} = \mathbf{A}\mathbf{q}^{(m)}$$

effizient zu lösen. Der resultierende verallgemeinerte Lanczos-Algorithmus nimmt mit diesen Modifikationen die folgende Form an:

**Algorithmus 10.1 (Verallgemeinerter Lanczos-Algorithmus)** *Seien  $\mathbf{A} \in \mathbb{K}^{n \times n}$  selbstadjungiert, sei  $\mathbf{B} \in \mathbb{K}^{n \times n}$  selbstadjungiert und positiv definit, und sei ein Startvektor  $\mathbf{q}^{(1)} \in \mathbb{K}^n$  mit  $\|\mathbf{q}^{(1)}\|_B = 1$  gegeben. Der folgende Algorithmus berechnet die bezüglich des zu  $\mathbf{B}$  gehörenden Energieskalarprodukts orthonormale Arnoldi-Basis  $\mathbf{q}^{(1)}, \dots, \mathbf{q}^{(m_0)}$  und die Matrizen  $\widehat{\mathbf{A}}^{(m)}$ . Der Algorithmus endet mit  $m = m_0$ .*

```

m ← 1;
a ← Aq(m);
Löse Bp = a;
γ ← √⟨a, p⟩2;
α1 ← ⟨q(1), a⟩2;   p ← p − α1q(1);
β1 ← √⟨a, p⟩2;
while βm ≥ εirγ do begin
    q(m+1) ← p/βm;
    m ← m + 1;
    a ← Aq(m);
    Löse Bp = a;
    γ ← √⟨a, p⟩2;
    αm ← ⟨q(m), a⟩2;   p ← p − β̄m−1q(m−1) − αmq(m);
    βm ← √⟨a, p⟩2
end
    
```

Im Algorithmus wird bei der Berechnung von  $\beta_m$  ein kleiner Trick verwendet: Streng genommen müssten wir

$$\langle \mathbf{p}, \mathbf{p} \rangle_B = \langle \mathbf{B}\mathbf{p}, \mathbf{p} \rangle_2$$

berechnen und würden deshalb  $\mathbf{B}\mathbf{p}$  benötigen, also eine Multiplikation mit der Matrix  $\mathbf{B}$ . Glücklicherweise steht  $\mathbf{p}$  an dieser Stelle des Algorithmus nach Konstruktion senkrecht (bezüglich des Energie-Skalarprodukts) auf  $\mathbf{q}^{(m-1)}$  und  $\mathbf{q}^{(m)}$ , erfüllt also

$$\langle \mathbf{B}\mathbf{q}^{(m-1)}, \mathbf{p} \rangle_2 = \langle \mathbf{q}^{(m-1)}, \mathbf{p} \rangle_B = 0, \quad \langle \mathbf{B}\mathbf{q}^{(m)}, \mathbf{p} \rangle_2 = \langle \mathbf{q}^{(m)}, \mathbf{p} \rangle_B = 0,$$

so dass wir

$$\begin{aligned} \langle \mathbf{a}, \mathbf{p} \rangle_2 &= \langle \mathbf{A}\mathbf{q}^{(m)}, \mathbf{p} \rangle_2 = \langle \mathbf{B}\mathbf{B}^{-1}\mathbf{A}\mathbf{q}^{(m)}, \mathbf{p} \rangle_2 \\ &= \langle \mathbf{B}(\mathbf{B}^{-1}\mathbf{A}\mathbf{q}^{(m)} - \bar{\beta}_{m-1}\mathbf{q}^{(m-1)} - \alpha_m\mathbf{q}^{(m)}), \mathbf{p} \rangle_2 = \langle \mathbf{B}\mathbf{p}, \mathbf{p} \rangle_2 \end{aligned}$$

erhalten und so die zusätzliche Multiplikation vermeiden können. Der verallgemeinerte Lanczos-Algorithmus berechnet eine Tridiagonalmatrix, deren Eigenwerte die der Matrix  $\widehat{\mathbf{A}}$  approximieren, und damit auch die des verallgemeinerten Eigenwertproblems.

Im Fall der vorkonditionierten Eigenwertverfahren ist die Situation etwas besser, da sich die Inverse der Matrix  $\mathbf{B}$  in der Matrix des Vorkonditionierers unterbringen lässt: Wenn  $\widehat{\mathbf{N}}$  ein Vorkonditionierer für die Matrix  $\widehat{\mathbf{A}}$  ist, nimmt die Richardson-Iteration die Form

$$\widehat{\mathbf{x}}^{(m+1)} = \widehat{\mathbf{x}}^{(m)} - \theta \widehat{\mathbf{N}}(\widehat{\mathbf{A}}\widehat{\mathbf{x}}^{(m)} - \mu\widehat{\mathbf{x}}^{(m)}) \quad \text{für alle } m \in \mathbb{N}_0$$

an. Indem wir wie zuvor

$$\widehat{\mathbf{x}}^{(m)} = \mathbf{L}^*\mathbf{x}^{(m)}, \quad \mathbf{x}^{(m)} = (\mathbf{L}^*)^{-1}\widehat{\mathbf{x}}^{(m)} \quad \text{für alle } m \in \mathbb{N}_0$$

einsetzen, erhalten wir

$$\begin{aligned} \mathbf{x}^{(m+1)} &= (\mathbf{L}^*)^{-1}\widehat{\mathbf{x}}^{(m)} - \theta(\mathbf{L}^*)^{-1}\widehat{\mathbf{N}}(\mathbf{L}^{-1}\mathbf{A}(\mathbf{L}^*)^{-1}\widehat{\mathbf{x}}^{(m)} - \mu\widehat{\mathbf{x}}^{(m)}) \\ &= (\mathbf{L}^*)^{-1}\widehat{\mathbf{x}}^{(m)} - \theta(\mathbf{L}^*)^{-1}\widehat{\mathbf{N}}\mathbf{L}^{-1}(\mathbf{A}(\mathbf{L}^*)^{-1}\widehat{\mathbf{x}}^{(m)} - \mu\mathbf{L}\widehat{\mathbf{x}}^{(m)}) \\ &= \mathbf{x}^{(m)} - \theta(\mathbf{L}^*)^{-1}\widehat{\mathbf{N}}\mathbf{L}^{-1}(\mathbf{A}\mathbf{x}^{(m)} - \mu\mathbf{B}\mathbf{x}^{(m)}) \quad \text{für alle } m \in \mathbb{N}_0. \end{aligned}$$

Wir definieren

$$\mathbf{N} := (\mathbf{L}^*)^{-1}\widehat{\mathbf{N}}\mathbf{L}^{-1}, \quad \widehat{\mathbf{N}} = \mathbf{L}^*\mathbf{N}\mathbf{L}$$

und erhalten

$$\mathbf{x}^{(m+1)} = \mathbf{x}^{(m)} - \theta\mathbf{N}(\mathbf{A}\mathbf{x}^{(m)} - \mu\mathbf{B}\mathbf{x}^{(m)}) \quad \text{für alle } m \in \mathbb{N}_0.$$

Für die Untersuchung der Konvergenz ist wichtig, dass

$$\widehat{\mathbf{N}}^{-1} = \mathbf{L}^{-1}\mathbf{N}^{-1}(\mathbf{L}^*)^{-1}, \quad \widehat{\mathbf{A}} = \mathbf{L}^{-1}\mathbf{A}(\mathbf{L}^*)^{-1}$$

gelten, so dass wir die Gleichungen

$$\begin{aligned} \langle \widehat{\mathbf{N}}^{-1}\widehat{\mathbf{x}}, \widehat{\mathbf{x}} \rangle_2 &= \langle \mathbf{L}^{-1}\mathbf{N}^{-1}(\mathbf{L}^*)^{-1}\mathbf{L}^*\mathbf{x}, \mathbf{L}^*\mathbf{x} \rangle_2 = \langle \mathbf{N}^{-1}\mathbf{x}, \mathbf{x} \rangle_2, \\ \langle \widehat{\mathbf{A}}\widehat{\mathbf{x}}, \widehat{\mathbf{x}} \rangle_2 &= \langle \mathbf{L}^{-1}\mathbf{A}(\mathbf{L}^*)^{-1}\mathbf{L}^*\mathbf{x}, \mathbf{L}^*\mathbf{x} \rangle_2 = \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle_2 \quad \text{für alle } \mathbf{x} \in \mathbb{K}^n, \widehat{\mathbf{x}} = \mathbf{L}^*\mathbf{x} \end{aligned}$$

erhalten. Damit gilt die für die Konvergenz des vorkonditionierten Verfahrens wichtige Bedingung

$$(1 - \gamma)\langle \mathbf{N}^{-1}\mathbf{x}, \mathbf{x} \rangle \leq \langle \mathbf{A}\mathbf{x}, \mathbf{x} \rangle \leq (1 + \gamma)\langle \mathbf{N}^{-1}\mathbf{x}, \mathbf{x} \rangle \quad \text{für alle } \mathbf{x} \in \mathbb{R}^n$$

genau dann, wenn

$$(1 - \gamma)\langle \widehat{\mathbf{N}}^{-1}\widehat{\mathbf{x}}, \widehat{\mathbf{x}} \rangle_2 \leq \langle \widehat{\mathbf{A}}\widehat{\mathbf{x}}, \widehat{\mathbf{x}} \rangle_2 \leq (1 + \gamma)\langle \widehat{\mathbf{N}}^{-1}\widehat{\mathbf{x}}, \widehat{\mathbf{x}} \rangle_2 \quad \text{für alle } \widehat{\mathbf{x}} \in \mathbb{R}^n$$

erfüllt ist. Falls  $\widehat{\mathbf{N}}$  also ein guter Vorkonditionierer für  $\widehat{\mathbf{A}}$  ist, ist  $\mathbf{N}$  auch ein guter Vorkonditionierer für  $\mathbf{A}$ .



# Index

- Adjungierte, 23
- Ähnliche Matrizen, 20
- Ähnlichkeitstransformation, 20
- Arnoldi-Basis, 172
  - Algorithmus, 174
- Bauer-Fike-Störungssatz, 84
- Begleitmatrix, 16
- Bisektion
  - Algorithmus, 149
- Block-Gradientenverfahren, 200
- Cauchy-Schwarz-Ungleichung, 23
- Charakteristisches Polynom, 16
  - Algorithmus, 149
- Deflation, 123
- Diagonalisierbarkeit
  - komplex, 41
  - reell, 38
- Eigenraum, 15
- Eigenvektor, 15
- Eigenwert, 15
  - dominant, 71
- Eigenwert-Zweigitterverfahren, 207
- Frobenius-Norm, 42
- Gerschgorin-Kreise, 158
- gestörte Matrix
  - Eigenwerte, 84
- Glättungsverfahren, 203
- Gradient, 168
- Gradientenverfahren, 190
  - Block-Variante, 200
  - vorkonditioniert, 194
- Gramsche Matrix, 27
- Hauptachsentransformation, 38
- Hessenberg-Form, 124
  - Algorithmus, 127
- Hessenberg-Matrix, 124
- Hessenbergmatrix
  - irreduzibel, 132
- Householder-Spiegelung, 31
- invarianter Unterraum, 28
- Inverse Iteration, 91
  - mit Rayleigh-Shift, 96
  - mit Shift, 93
- Isometrische Matrix, 30
- Jacobi-Iteration, 66
- Jordan-Normalform, 47
- kanonische Einheitsvektoren, 21
- Kongruenztransformation, 159
- Krylow-Raum, 169
- Kutta-Schukowski-Transformation, 177
- Lanczos-Algorithmus, 175
  - für verallgemeinerte Eigenwertprobleme, 215
- LOBPCG, 201
- Lokal optimales vorkonditioniertes Block-cg-Verfahren, 201
- Lokal optimales vorkonditioniertes cg-Verfahren, 198
- LOPCG, 198
- Matrix
  - Hessenberg, 124
  - irreduzibel, 52
  - isometrisch, 30
  - normal, 40
  - reduzibel, 52

## INDEX

- selbstadjungiert, 30
- unitär, 30
- Metrische Äquivalenz, 40
- Norm
  - euklidisch, 23
- Normale Matrix, 40
- Orthogonale Iteration, 112
- orthogonale Projektion, 100
- Perron-Frobenius-Theorie, 49
- PINVIT, 192
- positiv definit, 26
- positiv semidefinit, 26
- Prolongation, 203
- QR-Iteration, 120
- QR-Zerlegung, 32
- QZ-Iteration, 212
- Rayleigh-Iteration, 96
- Rayleigh-Quotient, 36
  - verallgemeinert, 110
- Residuum, 81
- Richardson-Iteration, 185
- Satz von Courant-Fischer, 36
- Schukowski-Transformation, 177
- Schur-Zerlegung, 34
- Schwachbesetzte Matrix, 166
- Selbstadjungierte Matrix, 30
- Shift, 92
- Singuläres Mehrgitterverfahren, 208
- Singulärwertzerlegung, 139
- Skalarprodukt
  - euklidisch, 22
- Spektrallücke, 85
- Spektralnorm, 24
- Spektralradius, 41
- Spektrum, 18
- Sturmsche Kette, 152
  - Algorithmus, 157
- Sylvester-Gleichung, 45
- Tschebyscheff-Polynom, 177
- Unitäre Matrix, 30
- Vektoriteration, 75
  - mit Abbruchkriterium, 78
- Vielfachheit, 18
- Vorkonditioniertes Gradientenverfahren, 194
- Wilkinson-Shift, 128
- Winkel, 70

## Literaturverzeichnis

- [1] G. Frobenius. Ueber Matrizen aus nicht negativen Elementen. *Sitzungsber. Königl. Preuss. Akad. Wiss.*, pages 456–477, 1921.
- [2] G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, London, 1996.
- [3] O. Perron. Zur Theorie der Matrices. *Mathematische Annalen*, 64(2):248–263, 1907.
- [4] H. Wielandt. Unzerlegbare, nicht negative Matrizen. *Mathematische Zeitschrift*, 52:642–648, 1950.