

ZAMM · Z. angew. Math. Mech. 68 (1988) 3, 181–184

ALEFELD, G.

## Berechenbare Fehlerschranken für ein Eigenpaar beim verallgemeinerten Eigenwertproblem

Herrn JOHANNES WEISSINGER zum 75. Geburtstag am 12. 5. 1988 gewidmet

*Für eine Näherung eines Eigenpaares des verallgemeinerten Eigenwertproblems  $Ax = \lambda Bx$  werden Schranken berechnet. Davon ausgehend wird ein Iterationsverfahren betrachtet, mit dessen Hilfe die Schranken verbessert werden können. Numerische Beispiele erläutern das Vorgehen. Alle Rundungsfehler werden dabei miteingerechnet.*

*For an approximate eigenpair of the generalized eigenvalue problem  $Ax = \lambda Bx$  we compute bounds. Using these bounds we consider an iteration method which improves the inclusions. The procedure is illustrated by some numerical examples. All rounding error are taken into account.*

Для приближенной собственной пары обобщенной задачи о собственных значениях  $Ax = \lambda Bx$  вычисляем грани. При помощи этих граней рассматриваем метод итераций улучшающий включения. Метод иллюстрируется численным примером. Все ошибки округления принимаются в расчет.

### 0. Einleitung

Wir gehen in dieser Arbeit von dem verallgemeinerten Matrixeigenwertproblem

$$Ax = \lambda Bx \quad (1)$$

mit quadratischen reellen  $(n, n)$ -Matrizen  $A$  und  $B$  aus. Wir setzen dabei stets voraus, daß  $B$  nichtsingulär ist. Dies ist unter praktischen Gesichtspunkten keine wesentliche Einschränkung, da bei Anwendungsbeispielen  $B$  gewöhnlich sogar symmetrisch und positiv definit ist.

Ausgehend von einer reellen Näherung  $\lambda$  für einen einfachen reellen Eigenwert und einem reellen Näherungsvektor  $x$  für den dazugehörigen Eigenvektor zweier gegebener reeller  $(n, n)$ -Matrizen  $A$  und  $B$  betrachten wir die Aufgabe, Fehlerschranken für diese Näherungen zu berechnen. Im Spezialfall  $B = I$  (Einheitsmatrix) wurde diese Aufgabenstellung in der Vergangenheit wiederholt betrachtet. Siehe dazu z. B. G. ALEFELD [1], S. RUMP [3], H. J. SYMM und J. WILKINSON [4] und T. YAMAMOTO [5].

In dieser Arbeit wird zunächst wie in [1] für das gewöhnliche Eigenwertproblem ein nichtlineares Gleichungssystem mit  $n$  Unbekannten aufgestellt, dessen exakte Auflösung mit der Bestimmung eines Eigenpaares für das verallgemeinerte Eigenwertproblem (1) äquivalent ist. Daran anschließend wird gezeigt, daß man mit hinreichend guten Näherungen  $\lambda$  und  $x$  den Eigenwert und die Komponenten des Eigenvektors in Schranken einschließen kann.

Beginnend mit dieser Einschließung kann man Folgen von Intervallvektoren berechnen, welche das Eigenpaar fortwährend einschließen und (theoretisch, d. h. bei Rundungsfehlerfreier Rechnung) gegen dieses konvergieren.

Wir erläutern das vorgeschlagene Verfahren abschließend an einigen Beispielen.

### 1. Berechnung von Schranken für ein Eigenpaar und deren iterative Verbesserung

Gegeben seien die reellen  $(n, n)$ -Matrizen  $A$  und  $B$ . Für einen reellen einfachen Eigenwert  $\lambda + \mu$  und einen dazugehörigen reellen Eigenvektor  $x + \tilde{y}$  der verallgemeinerten Matrixeigenwertaufgabe (1) seien die Näherungen  $\lambda$  und  $x$  bekannt, so daß also

$$A(x + \tilde{y}) = (\lambda + \mu) B(x + \tilde{y}) \quad (2)$$

gilt. Die Näherungen  $\lambda$  und  $x$  können z. B. mit einem der bekannten Näherungsverfahren berechnet worden sein. Siehe etwa [6].

Es sei mit  $x = (x_i)$

$$\|x\|_\infty = |x_s| > 0, \quad (3)$$

wobei wir, falls es mehrere Indizes gibt, für welche die Unendlichnorm angenommen wird, z. B.  $s$  als den kleinsten solchen Index wählen können.

Da der Eigenvektor  $x + \tilde{y}$  zunächst nicht eindeutig festgelegt ist, können wir die  $s$ -te Komponente von  $\tilde{y} = (\tilde{y}_i)$  gleich Null setzen:

$$\tilde{y}_s = 0. \quad (4)$$

Die Gleichung (2) kann geschrieben werden als

$$(A - \lambda B)\tilde{y} - \mu Bx = (\lambda B - A)x + \mu B\tilde{y}. \quad (5)$$

Führen wir nun einen Vektor  $y = (y_i) \in \mathbb{R}^n$  ein durch

$$y_i = \begin{cases} \tilde{y}_i, & i \neq s, \\ \mu, & i = s, \end{cases} \quad (6)$$

so kann unter Beachtung von (4) die Gleichung (5) geschrieben werden als

$$Cy = r + B(y; \tilde{y}) \quad (7)$$

mit dem Residuenvektor

$$r = \lambda Bx - Ax \quad (8)$$

und der Matrix  $C$ , die aus der Matrix  $A - \lambda B$  dadurch entsteht, daß die  $s$ -te Spalte durch  $-Bx$  ersetzt wird und alle anderen Spalten beibehalten werden.

Die Umformung von (2) in (7) mit der Normierung (4) findet man bereits in [4].

Für nichtsinguläres  $B$  und hinreichend gute Näherungen  $\lambda$  und  $x$  ist die Matrix  $C$  aus (7) nichtsingulär.

Dies folgt aus Stetigkeitsgründen aus der Tatsache, daß für den Spezialfall  $B = I$  und ein exaktes Eigenpaar  $(\lambda, x)$  mit einem algebraischen einfachen Eigenwert  $\lambda$  die Matrix  $C$  nach Satz 1 in [1] nichtsingulär ist und daß außerdem (1) mit dem gewöhnlichen Eigenwertproblem  $B^{-1}Ax = \lambda x$  äquivalent ist.

Die Gleichung (7) ist der Ausgangspunkt für die Berechnung von Schranken für  $\tilde{y}$  und  $\mu$  und damit für das Eigenpaar  $(\lambda + \mu, x + \tilde{y})$ .

Wir setzen voraus, daß  $\lambda$  und  $x$  so gute Näherungen sind, daß  $C$  nichtsingulär ist.  $L$  sei eine Näherung für die Inverse von  $C$  oder die exakte Inverse. Damit kann die Gleichung (7) geschrieben werden als

$$y = Lr + (I - LC)y + L(B(y; \tilde{y})) \quad (9)$$

Wir bestimmen nun einen Intervallvektor  $[y] = ([y]_i)$ , für welchen

$$Lr + (I - LC)y + L(B(y; \tilde{y})) \in [y] \quad (10)$$

für alle  $y \in [y]$  gilt. Aufgrund des Brouwerschen Fixpunktsatzes besitzt die Gleichung (9) dann mindestens einen Fixpunkt  $y^*$  in  $[y]$ .

Zur Bestimmung eines Intervallvektors  $[y]$ , für welchen (10) gilt, machen wir wie in [1] den Ansatz

$$[y] = [-\beta, \beta] e \quad (11)$$

mit  $\beta > 0$  und  $e = (1, \dots, 1)^T \in \mathbb{R}^n$ .

Wir setzen

$$\rho = \|Lr\|_\infty, \quad \kappa = \|I - LC\|_\infty, \quad l = \| |L| \cdot |B| \|_\infty, \quad (12), (13), (14)$$

wobei der Index  $\infty$  die Unendlichvektor- bzw. Unendlichmatrixnorm bezeichnet und  $|L|$  und  $|B|$  aus  $L$  und  $B$  durch elementweise Betragsbildung entstehen. Dann gilt die folgende Aussage, die vollständig analog wie Satz 2 in [1] bewiesen werden kann.

Satz 1: Es seien  $\rho, \kappa, l$  nach (12)–(14) definiert. Es sei  $\kappa < 1$ ,  $(1 - \kappa)^2 - 4\rho l \geq 0$ , und

$$\beta_{1/2} = \frac{1 - \kappa \mp \sqrt{(1 - \kappa)^2 - 4\rho l}}{2l} \quad (15)$$

seien die (positiven) Nullstellen der quadratischen Gleichung

$$l\beta^2 + (\kappa - 1)\beta + \rho = 0. \quad (16)$$

Wird dann  $\beta \in [\beta_1, \beta_2]$  gewählt, so hat die Gleichung (9) mindestens eine Lösung  $y^*$  im Intervallvektor (11).

Eine Lösung der Gleichung (9) ist sicherlich dann eine Lösung von (7), wenn  $L$  nichtsingulär ist. Unter der Voraussetzung  $\kappa < 1$  von Satz 1 ist dies stets der Fall, da wegen  $\|I - LC\|_\infty = \kappa < 1$  die Inverse von  $I - (I - LC) = LC$  existiert.

Wir betrachten nun das Iterationsverfahren

$$\begin{cases} [y]^0 = [-\beta, \beta] e, \\ [y]^{k+1} = g([y]^k), \quad k = 0, 1, 2, \dots, \end{cases} \quad (17)$$

wobei

$$g([y]) = Lr + (I - LC)[y] + L(B[y], [\tilde{y}]). \quad (18)$$

Durch (17) wird eine Folge von Intervallvektoren  $\{[y]^k\}_{k=0}^\infty$  berechnet. Für diese gilt

Satz 2: Es sei  $\kappa < 1$ ,  $(\kappa - 1)^2 - 4\rho l > 0$ , und  $\beta_1, \beta_2$  seien durch (15) definiert. Genügt dann  $\beta$  in (17) der Ungleichung

$$\beta_1 \leq \beta < \frac{\beta_1 + \beta_2}{2},$$

so liefert (17) eine Folge von Intervallvektoren  $\{[y]^k\}_{k=0}^\infty$ , mit

$$y^* \in [y]^k, \quad k = 0, 1, 2, \dots, \quad (19)$$

und

$$\lim_{k \rightarrow \infty} [y]^k = y^*. \quad (20)$$

Dabei ist  $y^*$  die eindeutige Lösung von (9) in  $[y]^0$ .

Der Beweis kann vollständig analog wie der von Satz 3 in [1] geführt werden. Wir gehen daher nicht auf die Einzelheiten ein, bemerken aber noch, daß unter den Voraussetzungen von Satz 2 für die Iterierten sogar  $[y]^{k+1} \subseteq [y]^k$ ,  $k = 0, 1, 2, \dots$ , gilt.

## 2. Numerische Beispiele

Gegeben seien die beiden symmetrischen Matrizen

$$F = \begin{bmatrix} 10 & 2 & 3 & 1 & 1 \\ 2 & 12 & 1 & 2 & 1 \\ 3 & 1 & 11 & 1 & -1 \\ 1 & 2 & 1 & 9 & 1 \\ 1 & 1 & -1 & 1 & 15 \end{bmatrix} \quad \text{und} \quad G = \begin{bmatrix} 12 & 1 & -1 & 2 & 1 \\ 1 & 14 & 1 & -1 & 1 \\ -1 & 1 & 16 & -1 & 1 \\ 2 & -1 & -1 & 12 & -1 \\ 1 & 1 & 1 & -1 & 11 \end{bmatrix}$$

aus [6], Seite 312.

a) Wir wählen als  $\lambda$  von der in [6], Tabelle 3 (Seite 313) angegebenen Näherung für den kleinsten Eigenwert des verallgemeinerten Eigenwertproblems  $Fx = \lambda Gx$  sechs Dezimalstellen:

$$\lambda = 0.432787.$$

Als dazugehörige Näherung für den Eigenvektor wählen wir einen Vektor, dessen Komponenten mit Ausnahme der betragsgrößten ebenfalls mit den ersten sechs Ziffern der dort angegebenen Zahlen übereinstimmen. Die betragsgrößte Komponente — das ist die dritte Komponente — wird exakt übernommen:

$$x = \begin{bmatrix} 0.134591 \\ -0.612947 \times 10^{-1} \\ -0.157902562211 \\ 0.109466 \\ -0.414730 \times 10^{-1} \end{bmatrix}.$$

Mit diesen Werten und mit  $C = L^{-1}$  erhält man für  $\beta_1$  in Satz 1 den Wert

$$\beta_1 = 0.426040283320 \times 10^{-6}.$$

Nach zwei Schritten des Iterationsverfahrens (17) erhält man die folgenden Einschließungen für den Eigenwert  $\lambda + \mu$  und den Eigenvektor  $x + \tilde{y}$ :

$$\lambda + \mu \in [0.43278721101_7^6];$$

$$x + \tilde{y} = \begin{bmatrix} [0.13459057396_5^6] \\ [-0.61294722471_4^5 \times 10^{-1}] \\ [-0.157902562211] \\ [0.10946578772_4^3] \\ [-0.41473011796_5^6 \times 10^{-1}] \end{bmatrix}.$$

b) Wir betrachten das verallgemeinerte Eigenwertproblem

$$Gx = \lambda Fx.$$

Als Näherung für den kleinsten Eigenwert wählen wir die auf 6 Stellen abgeschnittene Näherung aus [6]:

$$\lambda = 0.231060 \times 10^1.$$

Als Näherung für den dazugehörigen Eigenvektor wählen wir wie in a) mit Ausnahme der betragsgrößten Komponente die auf 6 Stellen abgeschnittenen Näherungen aus [6]. Die betragsgrößte Komponente übernehmen wir exakt aus [6]:

$$x = \begin{bmatrix} -0.204587 \\ 0.931721 \times 10^{-1} \\ 0.240022507111 \\ -0.166395 \\ 0.630418 \times 10^{-1} \end{bmatrix}.$$

Mit diesen Werten und mit  $C = L^{-1}$  erhält man für  $\beta_1$  in Satz 1

$$\beta_1 = 0.432157544139 \times 10^{-5}.$$

Nach zwei Iterationsschritten mit dem Verfahren (17) erhält man die folgenden Einschließungen

$$\lambda + \mu \in [0.23106043213_5^4 \times 10^1];$$

$$x + \tilde{y} \in \begin{bmatrix} [-0.20458671818_4^5] \\ [0.93172097743_6^5 \times 10^{-1}] \\ [0.240022507111] \\ [-0.1663953544_7^8] \\ [0.63041765310_7^6 \times 10^{-1}] \end{bmatrix}.$$

In beiden Beispielen stimmen nach zwei Iterationsschritten mit Ausnahme der letzten Stelle alle Stellen in den Mantissen überein.

Die Beispiele wurden auf einem IBM-Personalcomputer unter Verwendung der Programmiersprache PASCAL-SC gerechnet. Es stehen dabei intern 12 Dezimalstellen in der Mantisse zur Verfügung.

Wir erwähnen abschließend, daß (18) für  $L = C^{-1}$  in der Form

$$g([y]) = C^{-1}\{r + B([y], [\tilde{y}])\}$$

mit

$$r = \lambda Bx - Ax$$

geschrieben werden kann. Wenn wir hinreichend gute Näherungen  $(\lambda, x)$  für ein Eigenpaar besitzen, so wird bei der Berechnung von  $r$  in einem Gleitpunktzahlensystem Auslöschung eintreten.  $r$  ist daher möglichst genau zu berechnen. Dies kann mit dem sogenannten „genauen Skalarprodukt“ (siehe dazu [2]) folgendermaßen erfolgen:

Wir nehmen an, daß die Näherungen  $\lambda$  und  $x$  in dem Gleitpunktsystem exakt darstellbar sind (was keine Einschränkung für eine Näherung bedeutet). Bezeichnet  $gl(\lambda x_i)$  das berechnete Gleitpunktprodukt von  $\lambda x_i$  und bildet man mit den beiden zweidimensionalen Vektoren

$$\begin{pmatrix} \lambda \\ gl(\lambda x_i) \end{pmatrix} \quad \text{und} \quad \begin{pmatrix} x_i \\ -1 \end{pmatrix}$$

das genaue Skalarprodukt, das wir etwa mit  $ex_i$  bezeichnen, so gilt (exakt!)

$$\lambda x_i = gl(\lambda x_i) + ex_i.$$

Eine ähnliche Darstellung kann man für das Produkt  $[y]_s [\tilde{y}]$  angeben. Als Konsequenz ergibt sich, daß die Komponenten des in geschweiften Klammern stehenden Vektors

$$r + B([y]_s [\tilde{y}])$$

jeweils mit einem einzigen Skalarprodukt (aus Vektoren genügend großer Länge) berechnet werden können. Hierzu bietet sich die Verwendung des genauen Skalarprodukts an. Die Beispiele wurden unter Verwendung dieser von M. NEHER angegebenen Technik gerechnet.

### Literatur

- 1 ALEFELD, G.: Berechenbare Fehlerschranken für ein Eigenpaar unter Einschluß von Rundungsfehlern bei Verwendung des genauen Skalarprodukts. Z. angew. Math. Mech. **67** (1987) 3, 145–152.
- 2 KULISCH, U.; MIRANKER, W. (Eds.): A new approach to scientific computation. Academic Press 1983.
- 3 RUMP, S.: Solving algebraic problems with high accuracy. In [2], pp. 53–120.
- 4 SYMM, H. J.; WILKINSON, J. H.: Realistic error bounds for a simple eigenvalue and its associate eigenvector. Numer. Math. **35** (1980), 113–126.
- 5 YAMAMOTO, T.: Error bounds for computed eigenvalues and eigenvectors. Numer. Math. **34** (1980), 189–199.
- 6 WILKINSON, J. H.; REINSCH, C.: Handbook for automatic computation. Volume 2: Linear algebra. Springer Verlag, Berlin 1971.

Eingegangen am 25. Mai 1987

*Anschrift:* Prof. Dr. G. ALEFELD, Institut für Angewandte Mathematik, Universität Karlsruhe, Kaiserstraße 12, D-7500 Karlsruhe 1, BRD