

EIN NEUES APPROXIMATIONSVERFAHREN ZUR SPEKTRALZERLEGUNG VON MATRIZEN

Von

T. FREY

Lehrstuhl für Mathematik, Fakultät für Elektrotechnik
Technische Universität, Budapest

(Eingegangen am 3. September 1971)

1. Die Dreiecksiteration [1] bzw. einige neuere Varianten dieses Verfahrens, wie z. B. der L — R - resp. der Q — R -Algorithmus [2, 3] sind die am meisten angewandten Verfahren für die Spektralzerlegung von reellen, quadratischen Matrizen. Alle diese Verfahren haben jedoch zwei wesentliche Nachteile: einerseits besteht Konvergenz (dann und) nur dann, wenn alle Eigenwerte reell sind und verschiedene Modulen haben, andererseits ist die Konvergenz langsam, nämlich linear. Wir geben eine solche Weiterentwicklung der Grundidee der Dreiecksiteration an, durch die alle genannten Nachteile behoben werden: die Konvergenz ist quadratisch, und besteht auch dann, wenn komplexe oder mehrfache Eigenwerte auftreten. Einen Nachteil, u. zw. von algorithmischem Charakter hat jedoch das hier angegebene Verfahren: es arbeitet durch eine Halbierung der Dimensionszahl, und deshalb müssen mehrere Iterationsprozesse einander nachgeschaltet werden. Die Grundidee könnte man auch so anwenden, daß nur ein Iterationsprozeß erforderlich sei, jedoch wären dann die Rechenformeln zu kompliziert; es werden einige Hinweise auch auf diese Möglichkeit angegeben.

2. Es ist nicht schwer einzusehen (s. Satz 1), daß man zu jeder reellen quadratischen Matrix A , die in der Form

$$\begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}$$

partitioniert ist, reelle obere bzw. untere Blockdreiecksmatrizen S, T

$$S = \begin{pmatrix} S_{11} & S_{12} \\ \mathbf{0} & S_{22} \end{pmatrix}; \quad T = \begin{pmatrix} E_{11} & \mathbf{0} \\ H & E_{22} \end{pmatrix}$$

mit

$$A = TST^{-1} \tag{1}$$

angeben kann, wo A_{11}, S_{11} bzw. A_{22}, S_{22} quadratische Blöcke sind, und mindestens eines dieser Paare eine gerade Dimensionszahl besitzt (E_{11} bzw. E_{22} sind

Einheitsmatrizen mit derselben Dimensionszahl wie \mathbf{A}_{11} , \mathbf{S}_{11} bzw. wie \mathbf{A}_{22} , \mathbf{S}_{22}). Es ist nun bekannt, daß durch die Spektralzerlegung von \mathbf{S}_{11} bzw. von \mathbf{S}_{22} sogleich auch die Spektralzerlegung von \mathbf{S} angegeben wird. Die Angabe von \mathbf{S} und \mathbf{T} bedeutet also, daß wir uns mit der Spektralzerlegung von \mathbf{S}_{11} bzw. \mathbf{S}_{22} beschäftigen, und somit die ursprüngliche Dimensionszahl ungefähr halbieren können. Ist nun \mathbf{S}_{11} bzw. \mathbf{S}_{22} ein- oder zweidimensional, so kann die Spektralzerlegung durch elementare Operationen erfolgen. Die Hauptfrage ist also, die Blockdreiecksmatrizen \mathbf{S} und \mathbf{T} auszusuchen. Dazu kann einerseits die Verallgemeinerung der Dreiecksiteration angewandt werden

$$\mathbf{A}\mathbf{T}_{n-1} = \mathbf{T}_n \mathbf{S}_n, \quad (2)$$

u. zw. bei beliebiger Wahl von \mathbf{T}_0 . Es soll gezeigt werden (s. Satz 2), daß die verallgemeinerte Dreiecksiteration (2) konvergent ist, falls \mathbf{A} mindestens zwei Eigenwerte mit verschiedenen Modulen besitzt und die Partizionierungen von \mathbf{A} , \mathbf{T} und \mathbf{S} der Regel entsprechen, und daß mindestens einer der Blöcke \mathbf{A}_{11} und \mathbf{A}_{22} eine gerade Dimensionszahl hat. Jedoch ist die Konvergenz des Prozesses (2) nur linear. Es läßt sich jedoch auch eine quadratische Konvergenz erreichen — u. zw. auch dann, wenn alle Eigenwerte von \mathbf{A} einen gleichen Modulus haben —, falls eine hinreichend gute Approximation \mathbf{T}_n von \mathbf{T} bekannt ist (s. Satz 3). Ist nämlich \mathbf{T}_n bekannt, so wird dadurch die Gleichung

$$\mathbf{A}\mathbf{T}_n = \mathbf{T}^{(n)} \mathbf{S}^{(n)} \quad (3)$$

befriedigende obere bzw. untere Blockdreiecksmatrix $\mathbf{S}^{(n)}$ bzw. $\mathbf{T}^{(n)}$ bestimmt. Wäre nun $\mathbf{T}_n = \mathbf{T}^{(n)}$, so hätten wir schon die Lösung von (1) in der Hand. Das ist natürlich im allgemeinen nicht der Fall, doch ist \mathbf{T}_n eine gute Annäherung von \mathbf{T} , dann gibt eine geeignete geringe Variation von \mathbf{T}_n — und dadurch von $\mathbf{T}^{(n)}$ und $\mathbf{S}^{(n)}$ — die gesuchte Lösung. Es soll diese Variation von \mathbf{T}_n , $\mathbf{S}^{(n)}$ und $\mathbf{T}^{(n)}$ in der Form

$$\Delta\mathbf{T}_n = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \delta & \mathbf{0} \end{pmatrix}; \quad \Delta\mathbf{T}^{(n)} = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \delta^{(n)} & \mathbf{0} \end{pmatrix}; \quad \Delta\mathbf{S}^{(n)} = \begin{pmatrix} \Delta_{11} & \Delta_{12} \\ \mathbf{0} & \Delta_{22} \end{pmatrix} \quad (4)$$

gesucht werden. Diese Variationen sollen einerseits die Gleichung

$$\mathbf{T}_n + \Delta\mathbf{T}_n = \mathbf{T}^{(n)} + \Delta\mathbf{T}^{(n)},$$

d. h.

$$\mathbf{H}_n + \delta = \mathbf{H}^{(n)} + \delta^{(n)}, \quad (5)$$

andererseits die Gleichung

$$\mathbf{A} \cdot (\mathbf{T}_n + \Delta\mathbf{T}_n) = (\mathbf{T}^{(n)} + \Delta\mathbf{T}^{(n)}) \cdot (\mathbf{S}^{(n)} + \Delta\mathbf{S}^{(n)}),$$

d. h. die Gleichungen

$$\mathbf{A}_{11} + \mathbf{A}_{12} \mathbf{H}_n + \mathbf{A}_{12} \delta = \mathbf{S}_{11}^{(n)} + \Delta_{11}, \quad (6)$$

$$\mathbf{A}_{12} = \mathbf{S}_{12}^{(n)} + \Delta_{12}, \quad (7)$$

$$\mathbf{A}_{21} + \mathbf{A}_{22} \mathbf{H}_n + \mathbf{A}_{22} \delta = \mathbf{H}^{(n)} \cdot \mathbf{S}_{11}^{(n)} + \mathbf{H}^{(n)} \Delta_{11} + \delta^{(n)} \cdot \mathbf{S}_{11}^{(n)} + \delta^{(n)} \cdot \Delta_{11}, \quad (8)$$

$$\mathbf{A}_{22} = \mathbf{H}^{(n)} \mathbf{S}_{12}^{(n)} + \mathbf{H}^{(n)} \cdot \Delta_{12} + \delta^{(n)} \mathbf{S}_{12}^{(n)} + \delta^{(n)} \cdot \Delta_{12} + \mathbf{S}_{22} + \Delta_{22} \quad (9)$$

befriedigen. Das Gleichungssystem (5)–(9) ist in geschlossener Form nicht auflösbar. Wird jedoch das Glied zweiter Ordnung $\delta^{(n)} \Delta_{11}$ in (8) vernachlässigt, so erhält man die folgende, näherungsweise Lösung des betrachteten Systems:

$$\Delta_{11} = \mathbf{A}_{12} - \mathbf{S}_{12}^{(n)} \quad (10)$$

$$\delta^{(n)} = \mathbf{H}_n - \mathbf{H}^{(n)} + \delta; \quad (11)$$

$$\Delta_{11} = \mathbf{A}_{11} + \mathbf{A}_{12} \mathbf{H}_n + \mathbf{A}_{12} \delta - \mathbf{S}_{11}^{(n)}; \quad (12)$$

$$\mathbf{A}_{22} \delta - \mathbf{H}^{(n)} \cdot \Delta_{11} - \delta^{(n)} \mathbf{S}_{11}^{(n)} \cong \mathbf{H}^{(n)} \cdot \mathbf{S}_{11}^{(n)} - \mathbf{A}_{21} - \mathbf{A}_{22} \mathbf{H}_n,$$

folglich, nach (11)–(12)

$$\begin{aligned} (\mathbf{A}_{22} - \mathbf{H}^{(n)} \mathbf{A}_{11}) \delta - \delta \mathbf{S}_{11}^{(n)} \cong \mathbf{H}^{(n)} \mathbf{S}_{11}^{(n)} - \mathbf{A}_{21} - \mathbf{A}_{22} \mathbf{H}_n + \\ + \mathbf{H}^{(n)} \mathbf{A}_{11} + \mathbf{H}^{(n)} \mathbf{A}_{12} \mathbf{H}_n - \mathbf{H}^{(n)} \mathbf{S}_{11}^{(n)} + \mathbf{H}_n \mathbf{S}_{11}^{(n)} - \mathbf{H}^{(n)} \mathbf{S}_{11}^{(n)}. \end{aligned} \quad (13)$$

und schließlich

$$\Delta_{22} = -\mathbf{H}^{(n)} \mathbf{S}_{11}^{(n)} - \mathbf{H}^{(n)} \Delta_{12} + \mathbf{A}_{22} - \delta^{(n)} \mathbf{S}_{12}^{(n)} - \delta^{(n)} \cdot \Delta_{12} - \mathbf{S}_{22}. \quad (14)$$

Den Annäherungswert von δ — gewonnen aus der angenäherten Gleichung (13) — in (11) eingesetzt erhält man die Annäherung von $\delta^{(n)}$, und dadurch eine verbesserte Annäherung von \mathbf{T} . Das Newtonsche Verfahren (11)–(14) sichert nun eine quadratische Konvergenz. (Es sei hier bemerkt, daß mit Hilfe dieser ersten Annäherung in (8) auch das quadratische Glied abgeschätzt und somit eine zweite Verbesserung in den Unbekannten erreicht wird.)

3. In diesem Abschnitt sollen die oben benutzten Sätze in exakter Weise formuliert und bewiesen werden.

Satz I. *Genügt die Matrix \mathbf{A} bzw. ihre Partitionierung den in 2 angegebenen Voraussetzungen, so bestehen die Relation (1) und in der angegebenen Darstellung partionierbare Blockdreiecksmatrizen \mathbf{S} und \mathbf{T} .*

Beweis: Es sei \mathbf{A} in der Jordanschen Normalform

$$\mathbf{A} = [s_1, s_2, \dots, s_n] \cdot \mathbf{J} \cdot \begin{bmatrix} z_1^* \\ z_2^* \\ \vdots \\ z_n^* \end{bmatrix} \quad (15)$$

dargestellt, u. zw. in einer Reihenfolge der rechtsseitigen Hauptvektoren, daß, einerseits, in der Hauptdiagonalen von $[s_1, s_2, \dots, s_n]$ alle Elemente von 0 verschieden seien, und daß andererseits, die zu konjugierten Eigenwerten gehörenden und somit konjugierten Hauptvektoren nebeneinander stehen. Im ersten Schritt werden aus der Darstellung (15) die komplexen Zahlen eliminiert. (Die Paritätsbedingungen, die für die Dimensionszahlen von \mathbf{A}_{11} und \mathbf{A}_{22} vorgeschrieben wurden, hängen mit diesem Schritt zusammen.) Sind z. B. s_i und s_{i+1} konjugiert, so sind

$$\hat{s}_i = \frac{1}{2}(s_i + s_{i+1}) \quad \text{und} \quad \hat{s}_{i+1} = \frac{1}{2j}(s_i - s_{i+1}) \quad (16)$$

schon reell und linear unabhängig. Betrachten wir nun einerseits das Produkt

$$[s_1, s_2, \dots, s_i, s_{i+1}, \dots, s_n] \cdot \begin{bmatrix} e_1^* \\ e_2^* \\ \vdots \\ \frac{1}{2}(e_i^* + e_{i+1}^*) \\ \frac{1}{2j}(e_i^* - e_{i+1}^*) \\ \vdots \\ e_n^* \end{bmatrix}, \quad (17)$$

andererseits das Produkt

$$\begin{bmatrix} e_1^* \\ e_2^* \\ \vdots \\ \frac{1}{2}(e_i^* + e_{i+1}^*) \\ \frac{1}{2j}(e_i^* - e_{i+1}^*) \\ \vdots \\ e_n^* \end{bmatrix}^{-1} \cdot \mathbf{J} = \mathbf{H}_{i,i+1}^{-1} \cdot \mathbf{J}. \quad (18)$$

Es ist so gleich zu sehen, daß (17) die Matrix

$$[s_1, s_2, \dots, s_{i-1}, \hat{s}_i, \hat{s}_{i+1}, \dots, s_n]$$

angibt, und daß $\mathbf{H}_{i,i+1}^{-1} \cdot \mathbf{J}$ überall mit \mathbf{J} übereinstimmt, abgesehen vom zwei-dimensionalen Block in den i - und $i+1$ -ten Spalten und Zeilen, der von

$$\begin{bmatrix} \lambda & 0 \\ 0 & \bar{\lambda} \end{bmatrix} \quad (19)$$

in

$$\begin{bmatrix} \lambda & \bar{\lambda} \\ j\bar{\lambda} & -j\lambda \end{bmatrix}$$

übergeht. Es folgt dann, daß in der Darstellung

$$\mathbf{A} = \{[s_1, s_2, \dots, s_n] \cdot \mathbf{H}_{i,i+1}\} \cdot \{\mathbf{H}_{i,i+1}^{-1} \mathbf{J} \mathbf{H}_{i,i+1}\} \cdot \left\{ \mathbf{H}_{i,i-1}^{-1} \cdot \begin{bmatrix} \hat{z}_1^* \\ \hat{z}_2^* \\ \vdots \\ \hat{z}_n^* \end{bmatrix} \right\} \quad (20)$$

der erste Faktor in der i -ten und in der $i + 1$ -ten Spalte nur reelle Zahlen und in der Hauptdiagonalen von 0 verschiedene Elemente besitzt, während der zweite Faktor mit \mathbf{J} übereinstimmt, abgesehen vom in (19) genannten Block, der in

$$\begin{bmatrix} \frac{\lambda + \bar{\lambda}}{2} & \frac{\lambda - \bar{\lambda}}{2j} \\ -\frac{\lambda - \bar{\lambda}}{2j} & \frac{\lambda + \bar{\lambda}}{2} \end{bmatrix} \quad (21)$$

übergeht. Wird also für alle konjugierten Paare die angegebene Transformation benützt, erhält man schließlich die Darstellung

$$[\hat{s}_1, \hat{s}_2, \dots, \hat{s}_n] \cdot \mathbf{J} \cdot \begin{bmatrix} \hat{z}_1^* \\ \hat{z}_2^* \\ \vdots \\ \hat{z}_n^* \end{bmatrix}, \quad (22)$$

wo bereits alle Elemente reell sind, und in $[\hat{s}_1, \hat{s}_2, \dots, \hat{s}_n]$ die Hauptdiagonale nur von 0 verschiedene Zahlen enthält. Um nun in der i -ten Zeile alle Elemente — bzw. vom $(k + 1)$ -ten zum n -ten Element (falls $i > k$ ist; k ist die Dimensionszahl von \mathbf{A}_{11}) — zu annullieren (vom i -ten Element abgesehen), ist von links mit der Matrix

$$\mathbf{H}_i^{(1)} = \begin{bmatrix} -\frac{s_{1i}}{s_{ii}} e_i^* + e_1^* \\ s_{ii} \\ \vdots \\ -\frac{s_{i-1,i}}{s_{ii}} e_i^* + e_{i-1}^* \\ s_{ii} \\ e_i^* \\ \vdots \\ -\frac{s_{i,i}}{s_{ii}} e_i^* + e_k^* \\ s_{ii} \\ \vdots \\ -\frac{s_{n,i}}{s_{ii}} e_i^* + e_n^* \\ s_{ii} \end{bmatrix} \quad \text{bzw.} \quad \mathbf{H}_i^{(2)} = \begin{bmatrix} e_1^* \\ \vdots \\ e_k^* \\ -\frac{s_{n-1,i}}{s_{ii}} e_i^* + e_{k+1}^* \\ s_{ii} \\ \vdots \\ e_i^* \\ \vdots \\ -\frac{s_{n,i}}{s_{ii}} e_i^* + e_n^* \\ s_{ii} \end{bmatrix}$$

zu multiplizieren; in der Darstellung

$$\mathbf{A} = \{[\hat{s}_1, \dots, \hat{s}_n] \cdot \mathbf{H}_i\} \cdot \{\mathbf{H}_i^{-1} \cdot \hat{\mathbf{J}} \cdot \mathbf{H}_i\} \cdot \left\{ \mathbf{H}_i^{-1} \cdot \begin{bmatrix} \hat{x}_1^* \\ \hat{x}_2^* \\ \vdots \\ \hat{x}_n^* \end{bmatrix} \right\}$$

hat also die i -te Zeile von $[\hat{s}_1, \dots, \hat{s}_n] \cdot \mathbf{H}_i$ die Form, die der Form von \mathbf{T} entspricht. Multipliziert man mit allen \mathbf{H}_i , so ergibt sich die gewünschte Darstellung für \mathbf{T} . Da aber der dritte Faktor die Reziproke des ersten ist, und \mathbf{T}^{-1} notwendigerweise die angegebene Form hat, entspricht auch der dritte Faktor der angegebenen Form. Endlich gewährleistet die Struktur von $\hat{\mathbf{J}}$, daß der zweite Faktor

$$\left\{ \left(\prod_{i=1}^n \mathbf{H}_i^{-1} \right) \cdot \hat{\mathbf{J}} \cdot \left(\prod_{i=1}^n \mathbf{H}_i \right) \right\}$$

eine obere Blockdreiecksmatrix sei; \mathbf{H}_i^{-1} bedeutet ja (von links) eine Reihenelementoperation, bildet aber nur in der i -ten Reihe neue Elemente, und ist $i > k$, so bildet sie neue Elemente nur mit Hilfe solcher Reihen, die einen größeren Index als k haben; daraus folgt, daß der zweite Faktor eine obere Blockdreiecksmatrix ist, und damit ist unser Satz vollkommen bewiesen.

Um nun Satz 2 in exakter Weise zu formulieren, betrachten wir erst den Beweis von Satz 1. Wie wir gesehen haben, wird durch eine Jordansche Normalform dann, und nur dann, die Darstellung von \mathbf{A} in der gewünschten Form ermöglicht, wenn die rechtsseitigen Hauptvektoren in $[x_1, x_2, \dots, x_n]$ eine Reihenfolge haben, in welcher einerseits die Konjugierten nebeneinander stehen, andererseits wo in der Hauptdiagonalen keine O -Elemente vorkommen. Alle möglichen Reihenfolgen, die die obigen Anforderungen erfüllen, werden in der folgenden, eine Partitionierung-zulassenden Reihenfolge genannt.

Satz 2. *In der Blockdreiecksiteration (2) konvergieren — u. zw. linear — die Matrizen \mathbf{T}_n und \mathbf{S}_n gegen die Gleichung (1) genügende Grenzmatrizen \mathbf{T} und \mathbf{S} , falls es für \mathbf{A} eine solche, die Partitionierung-zulassende Reihenfolge der Hauptvektoren gibt, in welcher der größte Modulus der ersten k -Eigenwerte ungleich dem größten Modulus der letzten $(n - k)$ -Eigenwerte ist.*

Beweis: Betrachtet man die Iteration (2), so erhält man nacheinander die Relationen:

$$\begin{aligned} \mathbf{A}\mathbf{T}_0 &= \mathbf{T}_1 \mathbf{S}_1 ; \\ \mathbf{A}^2 \mathbf{T}_0 &= \mathbf{A}(\mathbf{A}\mathbf{T}_0) = \mathbf{A}(\mathbf{T}_1 \mathbf{S}_1) = \mathbf{A}\mathbf{T}_1 \mathbf{S}_1 = \mathbf{T}_2 \mathbf{S}_2 \mathbf{S}_1 ; \\ \mathbf{A}^3 \mathbf{T}_0 &= \mathbf{A}(\mathbf{A}^2 \mathbf{T}_0) = \mathbf{A}(\mathbf{T}_2 \mathbf{S}_2 \mathbf{S}_1) = \mathbf{T}_3 \mathbf{S}_3 \mathbf{S}_2 \mathbf{S}_1 ; \dots \\ \mathbf{A}^n \mathbf{T}_0 &= \mathbf{A}(\mathbf{A}^{n-1} \mathbf{T}_0) = \mathbf{A}(\mathbf{T}_{n-1} \mathbf{S}_{n-1} \mathbf{S}_{n-2} \dots \mathbf{S}_1) = \\ &= \mathbf{T}_n \mathbf{S}_n \mathbf{S}_{n-1} \dots \mathbf{S}_1 . \end{aligned} \tag{23}$$

Da nun das Produkt aus diesen oberen Blockdreiecksmatrizen selbst eine obere Dreiecksmatrix ist, folgt aus (23), daß

$$\mathbf{A}^n = \mathbf{T}_n \mathbf{S}^{(n)} \mathbf{T}_0^{-1} \quad (24)$$

gilt. Die Partitionierung von \mathbf{A}^n entsprechend (24) verläuft nun ebenso, wie es beim Beweis des Satzes 1 gezeigt wurde (hier ist im allgemeinen, $q_0^{-1} \neq q_n^{-1}$), jedoch ist \mathbf{T}_0 vorgeschrieben und \mathbf{T}_0^{-1} hat eine gleiche Struktur wie \mathbf{T}_0 ; deshalb muß der Grundgedanke des gezeigten Algorithmus so modifiziert werden, daß rechts und links verschiedene Spaltenoperationen angewendet werden.

Betrachten wir also eine Jordansche Normalform von \mathbf{A}_n , wo einerseits, die Reihenfolge der Hauptvektoren zulässig, andererseits, der größte Modulus der Eigenwerte in der ersten k -Zeile ungleich dem der letzten $(n - k)$ -Zeile ist. Wird die n -te Potenz des größten Modulus der Eigenwerte von \mathbf{A} in \mathbf{J}^n ausgeklammert, bleiben an einigen Stellen der Hauptdiagonalen von $\check{\mathbf{J}}^n$ komplexe Zahlen mit dem Modulus 1, an anderen Stellen aber Elemente, die eine Größenordnung $O(q^n)$ haben (mit $0 \leq q < 1$); u. zw. bleiben Elemente mit dem Modulus 1 nur entweder in der ersten k - oder aber in den letzten $n - k$ -Zeilen. Wendet man nun die in Zusammenhang mit dem Beweis des Satzes 1 erklärten Transformationen an, um $[x_1, x_2, \dots, x_n]$ in die Form \mathbf{T}_n ,

$\begin{bmatrix} z_1^* \\ \vdots \\ z_n^* \end{bmatrix}$ aber in die Form \mathbf{T}_0^{-1} zu bringen, so transformiert sich $\check{\mathbf{J}}^n$ in eine Form

— u. zw. in die \mathbf{S}^n darstellende Form —, wo gewisse Elemente die Größenordnung $1 + O(q^n)$, andere aber die Größenordnung $O(q^n)$ haben. Von dem ausgeklammerten Faktor abgesehen, konvergiert $\mathbf{S}^{(n)}$ zu einer Grenzmatrix, und deswegen gilt dasselbe für \mathbf{T}_n . Wenn aber \mathbf{T}_n konvergiert, so folgt aus (2), daß auch \mathbf{S}_n konvergent ist, w. z. B. w.

Es sei hier bemerkt, daß falls \mathbf{A} überhaupt Eigenwerte mit verschiedenen Modulen hat, aber bei der gewählten Partitionierung keine Reihenfolge existiert, die die Forderung des Satzes 2 erfüllt, so kann entweder die Partitionierung von \mathbf{A} gewechselt oder eine Ähnlichkeitstransformation mit Hilfe von Permutationsmatrixenpaaren angewendet werden, um eine Vertauschung der Komponentenfolge der Hauptvektoren zu erzwingen, wodurch auch eine Reihenfolgenänderung der zulässigen Reihenfolgen dargestellt ist bzw. auch neue zulässige Reihenfolgen stattfinden.

Satz 3. *Der Newtonsche Iterationsprozeß (3)—(4)—(10)—(11)—(12)—(13)—(14) konvergiert quadratisch für eine beliebige Matrix \mathbf{A} und für beliebige Partitionen mit entsprechenden Dimensionszahlen, falls \mathbf{T}_n eine genügend gute Approximation von \mathbf{T} ist. Nehmen wir die Korrekturen in (8) in Betracht, die sich so ergeben, erhält man einen Iterationsprozeß mit kubischer Konvergenz.*

Beweis: Im Satz 1 wurde bewiesen, daß bei geeigneten Dimensionszahlen die gewünschte Darstellung (1) von \mathbf{A} immer existiert. Die Korrektur $\mathbf{T}^{(n)}$

bzw. die dadurch definierten Korrekturen S_n und T_n werden in (5)—(6)—(7) und in (9) ganz genau in Betracht gezogen und hier erhält man die gesuchte Lösung durch einfache Umordnung der Gleichungen. Nur in der Gleichung (8) haben wir ein quadratisches Glied, das in erster Approximation außer acht gelassen wird. Für die Unbekannte δ ergibt sich hier eine lineare Gleichung der Form

$$\mathbf{B}\delta + \delta\mathbf{C} = \mathbf{D}, \quad (25)$$

wo \mathbf{B} , \mathbf{C} und \mathbf{D} bekannte Matrizen sind. Hätten wir in (8) auch das quadratische Glied in Betracht gezogen, u. zw. mit seinem genauen (und hinreichend kleinen) Wert, so würde statt (25) die Gleichung

$$\mathbf{B}\delta + \delta\mathbf{C} = \mathbf{D} - \delta^{(n)} S_{11}^{(n)} \quad (25')$$

die genaue Lösung δ angeben. \mathbf{A} hat nun mehrere Darstellungen der Form (1), die sich in der Anordnung der rechtsseitigen Hauptvektoren unterscheiden (s. den Nachweis des Satzes 1); eben deshalb haben die möglichen \mathbf{T} - und \mathbf{S} -Matrizen, die die Darstellung (1) geben — einen Abstand (in einer gewissen Matrixnorm gemessen), der von 0 streng abgegrenzt ist. Ist also $\mathbf{T}^{(n)}$ eine Approximation eines möglichen \mathbf{T} mit einem viel kleineren Abstand von diesem \mathbf{T} als die untere Grenze der Abstände der möglichen \mathbf{T} - bzw. \mathbf{S} -Matrizen (u. zw. mit einem so geringen Abstand, daß durch (3) ein solches Paar \mathbf{T}_n, S_n definiert wird, dessen Glieder auch viel kleinere Abstände von \mathbf{T} bzw. von \mathbf{S} als diese untere Grenze aufweisen) so definiert (25') eindeutig die Lösung δ . Das kann aber nur dann der Fall sein, wenn die Gleichung (25) nicht entartet ist, also auch eine eindeutige Lösung besitzt. (Über die Gleichung (25) s. Teil 4.) Daraus folgt aber, daß die Lösung von (25) nur quadratisch von der genauen Lösung von (25') abweichen kann, w. z. B. w. Nehmen wir noch die quadratische Approximation von $\delta^{(n)}$ und $S_{11}^{(n)}$ in (25') in Betracht, so ergibt sich eine kubische Approximation von δ , wie es behauptet wurde.

4. Wir müssen uns noch mit der Lösung der linearen Gleichung (25) beschäftigen. Besitzen die quadratischen Matrizen \mathbf{B} , \mathbf{C} und \mathbf{D} die Dimensionszahl n , so ergibt sich das folgende Gleichungssystem, wenn wir die Spalten der Matrizen δ bzw. \mathbf{D} in die (Spaltenvektor) Elemente der Hypervektoren x bzw. g umschreiben:

$$\mathbf{F} \cdot x = g, \quad (26)$$

wo

$$x = \begin{bmatrix} \delta_1 \\ \delta_2 \\ \vdots \\ \delta_n \end{bmatrix}; \quad g = \begin{bmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \end{bmatrix};$$

$$F = \begin{bmatrix} B + c_{11} E & c_{21} E & c_{31} E & \dots & c_{n1} E \\ c_{12} E & B + c_{22} E & c_{32} E & \dots & c_{n2} E \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ c_{1n} E & c_{2n} E & c_{3n} E & \dots & B + c_{nn} E \end{bmatrix}. \quad (27)$$

Die Gleichung (26) hat nun — wie es aus (27) sichtbar ist — einen speziellen Charakter, der jedoch die direkte Lösung nicht erleichtert: bei einer Gauß- oder Gauß—Jordan-Elimination sind $O(n^3)$ Operationen durchzuführen. Es scheint also viel günstiger, die Lösung von (25) direkt, durch ein Iterationsverfahren zu suchen, da die Ersetzung einer Approximation von δ in (25) nur $O(n^3)$ Operationen erfordert.

Jedoch scheint es nicht leicht, die Gleichung (26) in eine Form

$$\delta = f(\delta, B, C, D) \quad (28)$$

so umzuschreiben, daß das Iterationsverfahren

$$\delta_{n-1} = f(\delta_n, B, C, D) \quad (29)$$

auch bei günstig gewähltem Anfangswert δ_0 dem Fixpunkt von f zustrebe, da die verallgemeinerte Fréchet-Derivierte von f im allgemeinen eine Norm über 1 hat. Die verallgemeinerte Regula-falsi-Methode ist natürlich anwendbar, jedoch ist für diese Methode eine gute Anfangsapproximation notwendig. Um eine solche Anfangsapproximation zu erhalten, wird das Schrödersche Verfahren verallgemeinert, das für lineare Gleichungen in der Regel bessere Konvergenz aufweist als das Iterationsverfahren (s. z. B. [4—6]). Betrachten wir also den Halbordnungs-Banachraum R , wo die Halbordnung eine Netzstruktur bildet, ferner den Operator $T : R \rightarrow R$, der in Form

$$T = T_1 + T_2 + T_3 + T_4$$

zerlegt werden kann, mit isotonen T_1 und T_3 bzw. mit antitonen T_2 und T_4 . Es seien ferner vorausgesetzt:

a) Es gibt einen konvexen Bereich $D \subseteq R$, wo $S_1 = T_1 + T_2$ und $S_2 = T_3 + T_4$ jedes Intervall in einer kompakten Menge durchführt;

b) $E - S_1$ ist von monoton wachsender Art;

c) es gibt ein Intervall $\langle w_0, z_0 \rangle \subseteq D$, und zu jedem inneren Intervall $\langle w_i, z_i \rangle \subseteq \langle w_0, z_0 \rangle$ ein Intervallpaar $\langle u_{10}(i), v_{10}(i) \rangle \subseteq D$ bzw. $\langle u_{20}(i), v_{20}(i) \rangle \subseteq D$, mit

$$u_{11}(i) = T_1 u_{10}(i) + T_2 v_{10}(i) + T_3 w_i + T_4 z_i \geq u_{10}(i); \quad (30)$$

$$v_{11}(i) = T_1 v_{10}(i) + T_2 u_{10}(i) + T_3 w_i + T_4 z_i \leq v_{10}(i); \quad (31)$$

$$u_{21}(i) = T_1 u_{20}(i) + T_2 v_{20}(i) + T_3 z_i + T_4 w_i \geq u_{20}(i); \quad (32)$$

$$v_{21}(i) = T_1 v_{20}(i) + T_2 u_{20}(i) + T_3 z_i + T_4 w_i \leq v_{20}(i). \quad (33)$$

d) Es gibt ein Indexpaar n_1 bzw. n_2 so, daß für die Folgen

$$u_{1,n+1} = T_1 u_{1n} + T_2 v_{1n} + T_3 w_0 + T_4 z_0; \quad (34)$$

$$v_{1,n+1} = T_1 v_{1n} + T_2 u_{1n} + T_3 w_0 + T_4 z_0; \quad (35)$$

$$u_{2,n+1} = T_1 u_{2n} + T_2 v_{2n} + T_3 z_0 + T_4 w_0; \quad (36)$$

$$v_{2,n+1} = T_1 v_{2n} + T_2 u_{2n} + T_3 z_0 + T_4 w_0, \quad (37)$$

$$u_{1,n_1} \text{ und } v_{1,n_1} \in \langle w_0, z_0 \rangle; u_{2,n_2} \text{ und } v_{2,n_2} \in \langle w_0, z_0 \rangle \quad (38)$$

feststehen.

Satz 4. Neben den obigen Voraussetzungen besitzt der Operator T mindestens einen Fixpunkt in $\langle w_0, z_0 \rangle$. Ist $E - S_1$ eindeutig umkehrbar und von streng monotoner Art, und sind T_3 und T_4 paarweise streng monoton (d. h. daß für $w < z$: $T_3 w + T_4 z < T_3 z + T_4 w$ gelte), und kann für kein $w < z$, $w \in D$ $z \in D$ das Gleichungspaar

$$w = (E - S_1)^{-1} (T_3 w + T_4 z); \quad z = (E - S_1)^{-1} (T_3 z + T_4 w)$$

feststehen, so besitzt T einen und nur einen Fixpunkt in $\langle w_0, z_0 \rangle$, den man mit Hilfe eines Verfahrens des Typs (34)—(37) approximieren kann.

Beweis: Es ist bekannt (s. z. B. [7]), daß unter den Voraussetzungen (30), (31) die Punkte, definiert durch (34)—(35), ineinander eingeschachtelte Intervallen bilden, d. h.

$$\langle u_{1,n+1}, v_{1,n+1} \rangle \subseteq \langle u_{1n}, v_{1n} \rangle.$$

Laut des Schröderschen Fixpunktsatzes besitzt dann der Operator $S_1 + T_3 w_0 + T_4 z_0$ wegen a) mindestens einen Fixpunkt in $\langle u_{1n}, v_{1n} \rangle$, und da (38) auch erfüllt ist, auch mindestens einen Fixpunkt in $\langle w_0, z_0 \rangle$. Ebenso kann gezeigt werden, daß auch der Operator $S_1 + T_3 z_0 + T_4 w_0$ mindestens einen Fixpunkt in $\langle w_0, z_0 \rangle$ besitzt. Betrachten wir einen Fixpunkt $w_1 \in \langle w_0, z_0 \rangle$ des ersten, und einen Fixpunkt $z_1 \in \langle w_0, z_0 \rangle$ des zweiten Operators.

Da weiterhin $E - S_1$ von monoton wachsender Art ist, und $T_3 w_0 + T_4 z_0 \leq T_3 z_0 + T_4 w_0$ gilt, steht auch $w_0 \leq w_1 \leq z_1 \leq z_0$ fest.

Laut (30)—(33) kann also der obige Gedankengang auch für die Operatoren $S_1 + T_3 w_1 + T_4 z_1$ bzw. $S_1 + T_3 z_1 + T_4 w_1$ verfolgt werden. Es bleibt nur die Frage offen, ob man auch jetzt ein Indexpaar $n_1(1)$, $n_2(1)$ mit der Eigenschaft

$$\langle u_{1,n_1(1)}(1), v_{1,n_1(1)}(1) \rangle \in \langle w_1, z_1 \rangle; \quad \langle u_{2,n_2(1)}(1), v_{2,n_2(1)}(1) \rangle \subseteq \langle w_1, z_1 \rangle$$

finden kann. Da aber beide Operatoren gewiß einen Fixpunkt in $\langle u_{1,n}(1), v_{1,n}(1) \rangle$ bzw. in $\langle u_{2,n}(1), v_{2,n}(1) \rangle$ haben, und $T_3 w_1 + T_4 z_1 \geq T_3 w_0 + T_4 z_0$ bzw. $T_3 z_1 + T_4 w_1 \leq T_3 z_0 + T_4 w_0$ auch feststeht, ferner der Operator $E - S_1$ von monotoner Art ist, so folgt daß sowohl $S_1 + T_3 w_1 + T_4 z_1$ einen Fixpunkt $w_2 \geq w_1$ als auch $S_1 + T_3 z_1 + T_4 w_1$ einen Fixpunkt $z_2 \leq z_1$ haben, ferner, daß auch $w_2 \leq z_2$ feststeht. Die Konstruktion kann also unbedingt fortgesetzt, und mit vollständiger Induktion bewiesen werden, daß sowohl der Operator $S_1 + T_3 w_n + T_4 z_n$ einen Fixpunkt $z_{n-1} \leq z_n$ besitzen, und $w_{n+1} \leq z_{n+1}$ auch feststeht. Es soll nun gezeigt werden, daß

$$T(\langle w_n, z_n \rangle) \subseteq \langle w_n, z_n \rangle \tag{38}$$

gültig ist. Es sei $x \in \langle w_n, z_n \rangle$. Da T_3 isoton, und T_4 antiton ist, folgt, daß $S_1 x + T_3 w_n + T_4 z_n \leq S_1 x + T_3 x + T_4 x = Tx \leq S_1 x + T_3 z_n + T_4 w_n$. Da ferner $S_1 + T_3 w_n + T_4 z_n$ einen Fixpunkt $w_{n+1} \in \langle w_n, z_n \rangle$ und $S_1 + T_3 z_n + T_4 w_n$ einen Fixpunkt $z_{n+1} \in \langle w_n, z_n \rangle$ besitzen, folgt, daß $w_{n+1} \leq Tx \leq S_1 x + T_3 w_n + T_4 z_n \leq Tx \leq S_1 x + T_3 z_n + T_4 w_n \leq T_1 z_n + T_2 w_n + T_3 z_n + T_4 w_n \leq z_{n+1}$, folglich ist

$$T(\langle w_n, z_n \rangle) \subseteq \langle w_{n+1}, z_{n+1} \rangle, \tag{39}$$

also steht auch (38) umso mehr fest. Damit haben wir den ersten Teil unseres Satzes bewiesen.

Der zweite Teil des Satzes folgt fast unmittelbar; da $(E - S_1)$ eindeutig umkehrbar und von streng monotoner Art ist, kann eine Relation der Form

$$\lim_{k \rightarrow \infty} u_{1k}^{(n)}(i) = u_1^{(n)}(i) < \lim_{k \rightarrow \infty} v_{1k}^{(n)}(i) = v_1^{(n)}(i) \tag{40}$$

bzw.

$$\lim_{k \rightarrow \infty} u_{2k}^{(n)}(i) = u_2^{(n)}(i) < \lim_{k \rightarrow \infty} v_{2k}^{(n)}(i) = v_2^{(n)}(i) \tag{41}$$

nicht richtig sein, da

$$u_1^{(n)}(i) = (E - S_1)^{-1} (T_3 z_n + T_4 z_n) = v_1^{(n)}(i) \tag{42}$$

bzw.

$$u_2^{(n)}(i) = (E - S_1)^{-1} (T_3 z_n + T_4 w_n) = v_2^{(n)}(i) \tag{43}$$

gelten. Die Punkte w_n, z_n sind also eindeutig definiert. Es kann aber auch

$$\lim_{n \rightarrow \infty} w_n = w < \lim_{n \rightarrow \infty} z_n = z \tag{44}$$

nicht gelten, da dann

$$w = (E - S_1)^{-1} (T_3 w + T_4 z)$$

und

$$z = (E - S_1)^{-1} (T_3 z + T_4 w)$$

auch gültig wären in Gegensatz zu unseren Voraussetzungen. Es folgt also, daß in diesem Falle der eindeutig definierte Fixpunkt von T folgendermaßen gefunden werden kann:

$$\begin{aligned} \text{Die Folge} \quad u_{1,n+1}^{(k)} &= T_1 u_{1n}^{(k)} + T_2 v_{1n}^{(k)} + T_3 w_k + T_4 z_k \\ v_{1,n+1}^{(k)} &= T_1 v_{1n}^{(k)} + T_2 u_{1n}^{(k)} + T_3 w_k + T_4 z_k \end{aligned} \quad (45)$$

bzw. die Folge

$$\begin{aligned} u_{2,n+1}^{(k)} &= T_1 u_{2n}^{(k)} + T_2 v_{2n}^{(k)} + T_3 z_k + T_4 w_k \\ v_{2,n+1}^{(k)} &= T_1 v_{2n}^{(k)} + T_2 u_{2n}^{(k)} + T_3 z_k + T_4 w_k \end{aligned} \quad (46)$$

sollen bis zu einem Index n_1 bzw. n_2 fortgesetzt werden, für die bereits

$$u_{1,n_1}^{(k)} > w_k \quad \text{bzw.} \quad v_{2,n_2}^{(k)} < z_k \quad (47)$$

feststehen; wir können dann

$$w_{k+1} = u_{1,n_1}^{(k)} \cap v_{2,n_2}^{(k)}; \quad z_{k+1} = u_{1,n_1}^{(k)} \cup v_{2,n_2}^{(k)} \quad (48)$$

wählen, wo \cap bzw. \cup die Bildung der unteren bzw. der oberen Grenze bedeuten.

Nun stehen für die betrachtete lineare Gleichung $B\delta + \delta C - D = 0$ alle Voraussetzungen des Satzes 4 fest, falls wir sie in eine Form

$$\delta = \alpha(\delta + \varepsilon(B\delta + \delta C - D)) + (1 - \alpha)(\delta + \varepsilon(B\delta + \delta C - D))$$

umschreiben, mit $\alpha \sim 3/4$ und genügend kleinem $|\varepsilon|$.

(Es sei erwähnt, daß für $\alpha \sim 1/2$ auch Elemente w und $z < w$ angegeben werden können, die die Gleichungen

$$w = (E - S_1)^{-1} (T_3 w + T_4 z); \quad z = (E - S_1)^{-1} (T_3 z + T_4 w)$$

gleichzeitig befriedigen.)

Es sei noch bemerkt, daß das Verfahren (54)—(57) nicht leicht durchzuführen ist und nicht schnell konvergiert; eben deswegen wenden wir es nur für eine grobe Eingrenzung der Lösung von $B\delta + \delta C - D = 0$ an, und nur nachfolgend bedienen wir uns der verallgemeinerten Regula-falsi-Methode. Es soll aber auch über diese Methode eine Bemerkung hinzugefügt werden. Die Regula-falsi-Methode wird in der Regel in Banachräumen angewandt, die gleichzeitig einen kommutativen Ring bilden und in denen eine Invertier-

menge existiert, jedoch ist der Raum der n -dimensionalen quadratischen Matrizen kein kommutativer Ring, und deshalb wendet man die Methode nicht in der im Buche von Collatz angegebenen Form an, sondern in der Form

$$u_{n+1} = u_n - (u_{n-1} - u_n) \cdot (Tu_{n-1} - Tu_n)^{-1} Tu_n. \quad (49)$$

(Es sei bemerkt, daß hier der »Differenzenquotient« $(Tu_{n-1} - Tu_n) \cdot (u_{n-1} - u_n)^{-1}$ keine Approximation der Ableitung ist, da letztere ein »Array« mit vier Indexen ist; die Formel ist jedoch für die linearen Gleichungen $AX = B$ bzw. $X \cdot A = B$ noch genau; für den allgemeineren Fall $BX + XC = D$ ist das aber schon nicht der Fall.)

5. Wie schon erwähnt wurde, kann die hier angegebene Methode auch so formuliert werden, daß sie Formen der L — R -Methode näher steht. Man kann nämlich die Matrizen A , T und S »stärker« partitionieren, z. B. in eine Form

$$\begin{pmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{pmatrix} \cdot \begin{pmatrix} E & O & O \\ T_{21} & E & O \\ T_{31} & T_{32} & E \end{pmatrix} = \begin{pmatrix} E & O & O \\ T_{21} & E & O \\ T_{31} & T_{32} & TE \end{pmatrix} \cdot \begin{pmatrix} S_{11} & S_{12} & S_{13} \\ O & S_{22} & S_{23} \\ O & O & S_{33} \end{pmatrix}.$$

und um so mehr auch in höchstens zweidimensionale Blöcke; jedoch werden dann die Gleichungen viel komplizierter.

Zusammenfassung

Eine Verallgemeinerung der L — R -Methode wird angegeben, die einerseits eine quadratische bzw. kubische Konvergenz sichert, andererseits auch dann konvergiert, wenn die Eigenwerte komplex oder mehrfach sind, oder einige den gleichen Modulus haben. Verfasser beschäftigt sich auch mit der iterativen Lösung der Matrixgleichung $BX + XC = D$, die ein Hilfsmittel der angegebenen Methode ist.

Literatur

1. BAUER, F. L.: Beiträge zum Danilewski-Verfahren. Ber. Internat. Math. Koll. Nov. 1955. (1957)
2. RUTISHAUER, H.: Une méthode pour la détermination des valeurs propres d'une matrice. C.R. Ac. Sci. Paris 240 (1955).
3. FRANCIS, J. G. F.: The QR Transformation. I—II. Computing J. 4. (1961)—4 (1962).
4. SCHRÖDER, J.: Anwendung von Fixpunktsätzen bei der numerischen Behandlung nicht-linearer Gleichungen in halbgeordneten Räumen. Arch. Rat. Mech. Anal. 4. (1960).
5. ALBRECHT, J.: Fehlerschranken und Konvergenzbeschleunigung bei einer monotonen oder alternierenden Iterationsfolge. Num. Math. (1962).
6. ALBRECHT, J.: Zur Fehlerabschätzung beim Gesamt- und Einzelschrittverfahren für lineare Gleichungssysteme. Z. angew. Math. Mech. 43. (1963).
7. COLLATZ, L.: Funktionalanalysis und numerische Mathematik. Springer, Berlin—Göttingen—Heidelberg, 1964.

Prof. Dr. Tamás FREY, Budapest XI., Egry József u. 18—20, Ungarn