

Aplikace matematiky

Miroslav Šisler

Über die Konvergenz von Iterationsverfahren

Aplikace matematiky, Vol. 16 (1971), No. 1, 10–23

Persistent URL: <http://dml.cz/dmlcz/103323>

Terms of use:

© Institute of Mathematics AS CR, 1971

Institute of Mathematics of the Czech Academy of Sciences provides access to digitized documents strictly for personal use. Each copy of any part of this document must contain these *Terms of use*.



This document has been digitized, optimized for electronic delivery and stamped with digital signature within the project *DML-CZ: The Czech Digital Mathematics Library* <http://dml.cz>

ÜBER DIE KONVERGENZ VON ITERATIONSVERFAHREN

MIROSLAV ŠISLER

(Eingegangen am 6. März 1970)

Die Arbeit [2] befasste sich mit der Konvergenzbeschleunigung komplexer, durch die Formel

$$(1) \quad \mathbf{x}_{v+1} = \mathbf{P}_1^{-1} \mathbf{Q}_1 \mathbf{x}_v + \mathbf{P}_1^{-1} \mathbf{b}, \quad v = 0, 1, 2, \dots$$

definierter Iterationsverfahren, wo $\mathbf{A} = \mathbf{P}_1 - \mathbf{Q}_1$ eine solche Zerlegung der Matrix \mathbf{A} des Gleichungssystems $\mathbf{A}\mathbf{x} = \mathbf{b}$ ist, dass der Spektralradius $\varrho(\mathbf{P}_1^{-1} \mathbf{Q}_1)$ der Matrix $\mathbf{P}_1^{-1} \mathbf{Q}_1$ kleiner als 1 ist. Es wurde ein solcher Parameter k (im komplexen Gebiete) gesucht, für den das modifizierte, durch die Formel

$$(2) \quad \mathbf{x}_{v+1} = \mathbf{P}_k^{-1} \mathbf{Q}_k \mathbf{x}_v + \mathbf{P}_k^{-1} \mathbf{b}, \quad v = 0, 1, 2, \dots$$

definierte Iterationsverfahren am schnellsten konvergiert (d. h. für den der Spektralradius $\varrho(\mathbf{P}_k^{-1} \mathbf{Q}_k)$ möglichst klein ist); dabei ist $\mathbf{P}_k = k\mathbf{P}_1$, $\mathbf{Q}_k = (k-1)\mathbf{P}_1 + \mathbf{Q}_1$, $k \neq 0$.

Es ist gut bekannt, dass für eine nichtsymmetrische Matrix \mathbf{A} die Konvergenz mancher oft benutzter Iterationsverfahren (z. B. des Jacobischen Verfahrens, des Gauss-Seidelschen Verfahrens, Relaxationsverfahren usw.) nicht garantiert ist. In diesem Fall muss man manchmal das Gleichungssystem in die symmetrische Form überführen, was den Fehler vergrößert und die Berechnungen kompliziert. In diesem Artikel wird die Frage gelöst, ob und unter welchen Voraussetzungen man durch die geeignete Wahl des Parameters k ein konvergierendes Iterationsverfahren (2) bekommen kann, wenn das ursprüngliche Iterationsverfahren (1) nicht konvergiert. Der Artikel knüpft sich eng an die Arbeit [2] an und setzt die Kenntnis der Ergebnisse dieses Artikels voraus.

Man setzte voraus, dass die Matrix $\mathbf{P}_1^{-1} \mathbf{Q}_1$ n paarweise verschiedene Eigenwerte λ_i , $i = 1, \dots, n$ hat. In [2] wurde die Formel

$$(3) \quad \mu_i(k) = \frac{\lambda_i - 1}{k} + 1, \quad i = 1, \dots, n$$

angeführt, wo mit $\mu_i(k)$ die Eigenwerte der Matrix $\mathbf{P}_k^{-1} \mathbf{Q}_k$ bezeichnet sind (es ist also

$\lambda_i = \mu_i(1)$). Dabei bezeichnen $\mu_i(k)$ und λ_i die am gleichen Eigenvektor entsprechenden Eigenwerte (wie aus dem Artikel [2] bekannt ist, hängen die Eigenvektoren der Matrix $\mathbf{P}_k^{-1}\mathbf{Q}_k$ nicht von k ab).

Wie in [2] definiert man wieder die Funktionen g_i der komplexen Veränderlichen $k = k_1 + k_2i$, $k \neq 0$, durch die Formel

$$(4) \quad g_i(k) = |\mu_i(k)|^2 = \frac{|A_i|^2}{k_1^2 + k_2^2} + \frac{2(R_ik_1 + J_ik_2)}{k_1^2 + k_2^2} + 1, \quad i = 1, \dots, n,$$

wo $A_i = R_i + J_ii$, $\lambda_i = \lambda_i - 1$ ist.

Im Artikel [2] wurde vorausgesetzt, dass $\varrho(\mathbf{P}_1^{-1}\mathbf{Q}_1) < 1$ ist, d. h. dass die Zahlen A_i , $i = 1, \dots, n$ im Kreise mit dem Radius 1 und dem Mittelpunkt im Punkte -1 liegen. Wie aus der Bemerkung 1 des Artikels [2] folgt, sind gewisse Sätze aus der Arbeit [2] unabhängig von der Voraussetzung $\varrho(\mathbf{P}_1^{-1}\mathbf{Q}_1) < 1$. Es handelt sich um die Sätze 1, 2, 3, 5, 6, 7, 9, 11 und um die Bemerkung 3. Auf diese Sätze werden wir uns in dieser Arbeit berufen.

Die Eigenschaften der Funktionen g_i beschreibt der Hilfssatz 1 aus [2].

In der folgenden Betrachtung werden wir die Menge K solcher komplexer Zahlen k untersuchen, für die die Ungleichungen

$$g_i(k) < 1, \quad i = 1, \dots, n$$

gleichzeitig erfüllt sind. Da $\varrho(\mathbf{P}_k^{-1}\mathbf{Q}_k) = \max_{i=1, \dots, n} \sqrt{g_i(k)}$ ist, stellt diese Menge K die Menge aller Werte des Parameters k dar, für die das Iterationsverfahren (2) konvergiert. Nach Satz 9 des Artikels [2] (siehe auch [2], Hilfssatz 1, Behauptung c)) ist die Menge K durch den Durchschnitt der Halbebenen p_i

$$(5) \quad p_i \equiv 2R_ik_1 + 2J_ik_2 + |A_i|^2 < 0, \quad i = 1, \dots, n$$

gebildet (es ist klar, dass keine von diesen Halbebenen den Punkt $k = 0$ enthält). Man kann jetzt die Frage stellen, wann die Menge K nicht leer ist. Leicht kann man die folgende Behauptung beweisen:

1. Wenn $k' \in K$, ist, liegen in der Menge K auch alle Punkte der Hälfte der durch die Punkte $k = 0$ und $k = k'$ bestimmten Geraden, die ihren Ursprung in $k = k'$ hat und den Punkt $k = 0$ nicht enthält.

Beweis. Die Behauptung folgt sofort aus der Tatsache, dass diese Eigenschaft jede von den Halbebenen (5) hat.

Nun beweisen wir den folgenden Satz:

2. Die Menge K ist genau dann nicht leer, wenn eine durch den Punkt $k = 0$ gehende Gerade p existiert, sodass alle Zahlen A_i in einer durch die Gerade p begrenzten offenen Halbebene liegen.

Beweis. I) Jede Gerade p , die durch den Punkt $k = 0$ geht (mit Ausnahme der Geraden $k_2 = 0$, d. h. der reellen Achse), kann man in der Form

$$(6) \quad k_1 - qk_2 = 0$$

schreiben, wobei q eine reelle Zahl ist.

a) Man setze voraus, dass die Gerade p die Form (6) hat und dass die Zahlen A_i in einer durch die Gerade p begrenzten offenen Halbebene liegen, dann haben die Zahlen

$$R_i - qJ_i, \quad i = 1, \dots, n$$

offensichtlich das gleiche Vorzeichen. Nun beweisen wir, dass auf der Geraden p' , die zu der Geraden p senkrecht ist und durch den Punkt $k = 0$ geht, immer der Punkt $k' \in K$ liegt. Die Gerade p' hat die Gleichung

$$(7) \quad qk_1 + k_2 = 0,$$

sodass $k_2 = -qk_1$ ist. Nach Einsetzung in (5) bekommt man die Ungleichungen

$$2R_ik_1 - 2J_iqk_1 + |A_i|^2 < 0, \quad i = 1, \dots, n$$

oder

$$(8) \quad k_1(R_i - qJ_i) < -\frac{|A_i|^2}{2}, \quad i = 1, \dots, n.$$

Legt man also $k_1 = k'_1$, wo k'_1 eine Zahl mit genügend grossem Absolutbetrag ist, deren Vorzeichen dem Vorzeichen der Zahlen $R_i - qJ_i$ entgegengesetzt ist, sind alle Ungleichungen erfüllt. Der Punkt $k' = k'_1 - iqk'_1$ liegt dann im Durchschnitt der Halbebenen p_i , d. h. es ist $k' \in K$ und also $K \neq \emptyset$.

b) Man setze voraus, dass die Gerade p die Form $k_2 = 0$ hat und dass die Zahlen A_i wieder in einer durch die Gerade p begrenzten offenen Halbebene liegen; dann haben die Zahlen J_i , $i = 1, \dots, n$ offensichtlich das gleiche Vorzeichen. Die zu der Geraden p senkrechte Gerade p' , die durch den Punkt $k = 0$ geht, hat die Gleichung $k_1 = 0$. Nach Einsetzung in (5) bekommt man die Ungleichungen

$$(9) \quad J_ik_2 < -\frac{|A_i|^2}{2}, \quad i = 1, \dots, n.$$

Wenn man wieder die Zahl $k_2 = k'_2$ mit genügend grossem Absolutbetrag, und mit einem dem Vorzeichen der Zahlen J_i entgegengesetzten Vorzeichen wählt, sind alle Ungleichungen (9) erfüllt. Auf der Geraden p' liegt wieder ein solcher Punkt $k' = ik'_2$, dass $k' \in K$ ist. Es ist also $K \neq \emptyset$, was zu beweisen war.

II) Entgegengesetzt beweisen wir jetzt, dass für $K \neq \emptyset$ eine durch den Punkt $k = 0$ gehende Gerade p existiert, sodass alle Zahlen A_i in einer durch die Gerade p begrenzten offenen Halbebene liegen.

Da $K \neq \emptyset$ ist, kann man den Punkt k' im Durchschnitt der Halbebenen (5) wählen.

a) Es sei $k' = k'_1 + ik'_2$ und es gelte $k'_2 \neq 0$. Legt man $q = k'_1/k'_2$, d. h. $k'_1 = qk'_2$, dann erfüllt der Punkt k' nach der Voraussetzung die Ungleichungen (5), sodass die Ungleichungen

$$2R_i q k'_2 + 2J_i k'_2 + |A_i|^2 < 0,$$

$$k'_2(qR_i + J_i) < -\frac{|A_i|^2}{2}, \quad i = 1, \dots, n$$

gelten. Daraus folgt, dass alle Ausdrücke $qR_i + J_i$ das gleiche Vorzeichen haben. Das heisst aber, dass alle Punkte A_i in einer durch die Gerade

$$p \equiv qk_1 + k_2 = 0$$

begrenzten Halbebene liegen.

b) Wenn $k'_2 = 0$ ist, folgen aus (5) die Ungleichungen

$$2R_i k'_1 + |A_i|^2 < 0,$$

d. h.

$$R_i k'_1 < -\frac{|A_i|^2}{2}.$$

Daraus folgt, dass die Zahlen R_i das gleiche Vorzeichen haben müssen, sodass alle Zahlen A_i in einer durch die imaginäre Achse begrenzten offenen Halbebene liegen.

Dadurch ist der Satz 2 bewiesen.

Bemerkung 1. Aus Satz 2 folgen jetzt folgende zwei Spezialfälle: Wenn alle imaginäre Teile der Zahlen λ_i des Spektrums der Matrix $\mathbf{P}_1^{-1}\mathbf{Q}_1$ das gleiche Vorzeichen haben, oder wenn alle Zahlen λ_i in einer offenen Halbebene liegen, die durch eine mit der imaginären Achse parallele und durch den Punkt 1 gehende Gerade begrenzt ist, dann existiert ein solcher Wert des Parameters k , dass die modifizierte Methode (2) konvergiert.

Nun werden wir uns mit der Frage der Lage des Optimalparameters in der Menge K befassen (unter der Voraussetzung, dass diese Menge nicht leer ist). Als *Optimalparameter* wird wie in der Arbeit [2] eine solche komplexe Zahl k_0 bezeichnet, für die der Ausdruck $\max_i \sqrt{[g_i(k)]}$ sein Infimum bezüglich k in der Menge K erreicht (mit Rücksicht auf den Verlauf der Funktionen g_i in der Menge K ist das Infimum gleichzeitig ein Minimum). Unter dem *optimalen Spektralradius* verstehen wir dann die Zahl

$$(10) \quad \varrho(\mathbf{P}_{k_0}^{-1}\mathbf{Q}_{k_0}) = \min_{k \in K} \varrho(\mathbf{P}_k^{-1}\mathbf{Q}_k) = \min_k \max_i \sqrt{[g_i(k)]}.$$

In der Arbeit [2], Satz 2, wurde bewiesen, dass die Menge der Punkte k , für die $g_i(k) = g_j(k)$ ($i \neq j$) ist, eine Gerade

$$p_{ij} \equiv \alpha_{ij}k_1 + \beta_{ij}k_2 + \gamma_{ij} = 0$$

bildet, wo $\alpha_{ij} = R_i - R_j$, $\beta_{ij} = J_i - J_j$, $\gamma_{ij} = \frac{1}{2}(|A_i|^2 - |A_j|^2)$.

Da wir voraussetzen, dass $K \neq \emptyset$ ist, liegen nach Satz 2 die Zahlen A_i immer in einer offenen Halbebene, deren Grenzgerade den Punkt $k = 0$ enthält. Für beliebige A_i, A_j aus dieser Halbebene sind dann die Voraussetzungen des Satzes 3 aus [2] erfüllt. Aus diesem Satz folgt, dass der Punkt k^{ij} , in dem die Funktion g_i (bzw. g_j) auf der Geraden p_{ij} ihr Minimum annimmt, immer existiert und dass der Real- und Imaginärteil des Punktes k^{ij} durch die Formeln

$$(11) \quad k_1^{ij} = \frac{-(R_i + R_j) |A_i| |A_j| - R_j |A_i|^2 - R_i |A_j|^2}{2|A_i| |A_j| + 2(R_i R_j + J_i J_j)},$$

$$(12) \quad k_2^{ij} = \frac{-(J_i + J_j) |A_i| |A_j| - J_j |A_i|^2 - J_i |A_j|^2}{2|A_i| |A_j| + 2(R_i R_j + J_i J_j)}$$

gegeben ist.

Unter k^{ijk} werden wir ferner einen solchen Punkt verstehen, für den $g_i(k^{ijk}) = g_j(k^{ijk}) = g_k(k^{ijk})$ ist (in diesem Punkte schneiden sich also die Flächen $z = g_i(k)$, $z = g_j(k)$, $z = g_k(k)$). In [2] (Satz 5, Bemerkung 3) wurde gezeigt, dass der Punkt k^{ijk} genau dann existiert, wenn die Punkte A_i, A_j, A_k nicht kollinear sind. Dann schneiden sich die Geraden p_{ij}, p_{jk}, p_{ki} im Punkte k^{ijk} , wobei der Real- und Imaginärteil der Zahl k^{ijk} durch die Formeln

$$(13) \quad k_1^{ijk} = -\frac{1}{2} \frac{|A_i|^2 (J_j - J_k) + |A_j|^2 (J_k - J_i) + |A_k|^2 (J_i - J_j)}{R_i (J_j - J_k) + R_j (J_k - J_i) + R_k (J_i - J_j)},$$

$$(14) \quad k_2^{ijk} = -\frac{1}{2} \frac{|A_i|^2 (R_j - R_k) + |A_j|^2 (R_k - R_i) + |A_k|^2 (R_i - R_j)}{J_i (R_j - R_k) + J_j (R_k - R_i) + J_k (R_i - R_j)}$$

gegeben ist.

Für die Punkte k^{ij} gilt der folgende Satz:

3. Wenn die Zahlen A_i , $i = 1, \dots, n$ im Kreis $|A_i + a| < b$ liegen, wo $a \geq b > 0$ ist, dann liegen alle Zahlen k^{ij} im Kreis $|k^{ij} - a| < b$.

Beweis. Man setze voraus, dass für irgendeine von den Zahlen k^{ij} die Ungleichung $|k^{ij} - a|^2 \geq b^2$ gilt. Da $k^{ij} \neq 0$ ist für $i, j = 1, \dots, n$ (es kann nicht $k^{ij} = 0$ sein, da für $k = 0$ die Funktionen g_i, g_j keinen Sinn haben), kann man in diese Ungleichung $u^{ij} = A_i/k^{ij} + 1$ einsetzen. Nach einigen Umformungen bekommt man sukzessiv diese Ungleichungen:

$$(15) \quad |A_i - a(u^{ij} - 1)|^2 \geq b^2 |u^{ij} - 1|^2,$$

$$|A_i|^2 - a[A_i \bar{u}^{ij} + \bar{A}_i u^{ij}] + a(A_i + \bar{A}_i) \geq (b^2 - a^2) |u^{ij} - 1|^2.$$

Für die Real- und Imaginärteile der Zahlen u^{ij} gelten die Formeln

$$(16) \quad u_1^{ij} = \frac{|A_j|^2 - (R_i R_j + J_i J_j)}{(|A_i| + |A_j|) |A_j|},$$

$$(17) \quad u_2^{ij} = \frac{R_i J_j - R_j J_i}{(|A_i| + |A_j|) |A_j|}.$$

(Diese Formeln wurden in der Arbeit [2] im Beweis des Satzes 3 abgeleitet.)

Mit Hilfe (16), (17) bekommt man sofort nach einer Umformung

$$(18) \quad A_i \bar{u}^{ij} + \bar{A}_i u^{ij} = 2R_i u_1^{ij} + 2J_i u_2^{ij} = 2 \frac{R_i |A_j|^2 - R_j |A_i|^2}{(|A_i| + |A_j|) |A_j|}$$

und ferner

$$(19) \quad |u_{ij} - 1|^2 = (u_1^{ij} - 1)^2 + (u_2^{ij})^2 = \frac{2|A_i|^2 |A_j| + 2|A_i| (R_i R_j + J_i J_j)}{(|A_i| + |A_j|) |A_j|}.$$

Wenn man (18) und (19) in (15) einsetzt, bekommt man sukzessiv

$$(20) \quad |A_i|^2 - 2a \frac{R_i |A_j|^2 - R_j |A_i|^2}{(|A_i| + |A_j|) |A_j|} + 2aR_i \geq 2(b^2 - a^2) \frac{|A_i|^2 |A_j| + |A_i| (R_i R_j + J_i J_j)}{(|A_i| + |A_j|)^2 |A_j|},$$

$$|A_j| (|A_i|^2 + 2aR_i) + |A_i| (|A_j|^2 + 2aR_j) \geq 2(b^2 - a^2) \frac{|A_i| |A_j| + (R_i R_j + J_i J_j)}{|A_i| + |A_j|}.$$

Nach Voraussetzung ist $|A_i + a|^2 < b^2$, d. h. $|A_i|^2 + 2R_i a < b^2 - a^2 \leq 0$. Eine ähnliche Ungleichung gilt auch für A_j . Es gilt also die Ungleichung

$$(21) \quad |A_j| (|A_i|^2 + 2aR_i) + |A_i| (|A_j|^2 + 2aR_j) < (b^2 - a^2) (|A_i| + |A_j|).$$

Es sei $a = b$. Dann sind die rechten Seiten der Ungleichungen (20) und (21) gleich Null, was ein Widerspruch ist.

Es sei $a > b > 0$. Wir werden beweisen, dass die durch die Verbindung der Ungleichungen (20) und (21) entstehende Ungleichung

$$(22) \quad (b^2 - a^2) (|A_i| + |A_j|) > 2(b^2 - a^2) \frac{|A_i| |A_j| + (R_i R_j + J_i J_j)}{|A_i| + |A_j|}$$

zu einem Widerspruch führt. Da $b^2 - a^2 < 0$ ist, bekommt man nach Umformungen der Ungleichung (22) sukzessiv die Ungleichungen

$$\begin{aligned} (|A_i| + |A_j|)^2 &< 2[|A_i| + |A_j| + (R_i R_j + J_i J_j)], \\ |A_i|^2 - 2(R_i R_j + J_i J_j) + |A_j|^2 &< 0, \\ |A_i - A_j|^2 &< 0, \end{aligned}$$

was ein Widerspruch ist. Dadurch ist der Satz 3 bewiesen.

Mit Hilfe des Satzes 3 beweist man sofort den allgemeineren Satz 4.

4. Wenn die Zahlen A_i , $i = 1, \dots, n$ im Kreise $|A_i + a| < b$ liegen, wo a eine komplexe Zahl und $|a| \geq b > 0$ ist, dann liegen die Zahlen k^{ij} im Kreise $|k^{ij} - a| < b$.

Beweis. Es sei φ eine reelle Zahl, für die die Zahl $\tilde{a} = ae^{i\varphi}$ reell und positiv ist. Man definiere ferner die Zahl $\tilde{\lambda}_i = A_i e^{i\varphi} = \tilde{R}_i + i\tilde{J}_i$ und die Zahl $\tilde{k}^{ij} = \tilde{k}_1^{ij} + i\tilde{k}_2^{ij}$, wo

$$(23) \quad \tilde{k}_1^{ij} = \frac{-(\tilde{R}_i + \tilde{R}_j) |\tilde{\lambda}_i| |\tilde{\lambda}_j| - \tilde{R}_j |\tilde{\lambda}_i|^2 - \tilde{R}_i |\tilde{\lambda}_j|^2}{2|\tilde{\lambda}_i| |\tilde{\lambda}_j| + 2(\tilde{R}_i \tilde{R}_j + \tilde{J}_i \tilde{J}_j)},$$

$$(24) \quad \tilde{k}_2^{ij} = \frac{-(\tilde{J}_i + \tilde{J}_j) |\tilde{\lambda}_i| |\tilde{\lambda}_j| - \tilde{J}_j |\tilde{\lambda}_i|^2 - \tilde{J}_i |\tilde{\lambda}_j|^2}{2|\tilde{\lambda}_i| |\tilde{\lambda}_j| + 2(\tilde{R}_i \tilde{R}_j + \tilde{J}_i \tilde{J}_j)}$$

ist. Nun beweisen wir, dass

$$(25) \quad \begin{aligned} \tilde{k}^{ij} &= \tilde{k}_1^{ij} + i\tilde{k}_2^{ij} = \\ &= (k_1^{ij} \cos \varphi - k_2^{ij} \sin \varphi) + i(k_2^{ij} \cos \varphi + k_1^{ij} \sin \varphi) = k^{ij} e^{i\varphi} \end{aligned}$$

gilt. Aus der Beziehung $\tilde{\lambda}_i = A_i e^{i\varphi}$ folgt sofort, dass

$$(26) \quad \tilde{R}_i = R_i \cos \varphi - J_i \sin \varphi, \quad \tilde{J}_i = J_i \cos \varphi + R_i \sin \varphi$$

ist, und ähnlicherweise für $\tilde{\lambda}_j$.

Nach Einsetzung von (26) in (23), (24) bekommt man mit Hilfe der Formeln (11), (12) nach einer Umformung die Formel (25).

Da nun $|\tilde{\lambda}_i + \tilde{a}| = |A_i e^{i\varphi} + a e^{i\varphi}| = |A_i + a| < b$ ist, wo \tilde{a} reell ist, $|\tilde{a}| = |a|$, $0 < \tilde{a} \leq b$, folgt aus Satz 3 die Ungleichung $|\tilde{k}^{ij} - \tilde{a}| < b$ und da $|\tilde{k}^{ij} - \tilde{a}| = |k^{ij} e^{i\varphi} - a e^{i\varphi}| = |k^{ij} - a|$ gilt, ist der Satz 4 bewiesen.

Ferner gilt der folgende Satz:

5. Wenn die Zahlen A_i , $i = 1, \dots, n$ im Kreise $|A_i + a| < b$ liegen, wo a eine komplexe Zahl ist, $|a| \geq b > 0$, dann liegt der Optimalparameter k_0 im Durchschnitt des Kreises $|k - a| < b$ und der Halbebenen

$$(5) \quad 2R_i k_1 + 2J_i k_2 + |A_i|^2 < 0, \quad i = 1, \dots, n.$$

Beweis. Aus [2], Satz 7 folgt, dass der Optimalparameter im kleinsten konvexen, die Punkte k^{ij} enthaltenden Vieleck liegt. Da nach Satz 4 diese Punkte im Kreis $|k - a| < b$ liegen, liegt in diesem Kreis auch der Optimalparameter. Der Optimalparameter muss ferner in der Menge K liegen, d. h. im Durchschnitt der Halbebenen (5). Dadurch ist der Satz 5 bewiesen.

Bemerkung 2. Eine unmittelbare Folgerung des Satzes 5 ist der Satz 4 des Artikels [2], wenn man $a = b = 1$ legt. (In dem Artikel [2] untersuchen wir den Fall, wenn $|A_i + 1| < 1$ für alle i gilt.)

Aus den Sätzen 2 und 5 folgt dieser Satz:

6. Wenn die Zahlen A_i in irgendeinem den Punkt $k = 0$ nicht enthaltenden Kreis K_0 liegen, dann ist die Menge K nicht leer und der Optimalparameter liegt im Durchschnitt dieser Menge und des zu dem Kreise K_0 bezüglich des Koordinatensprungs symmetrisch liegenden Kreises.

Folgender Satz ist eine Verallgemeinerung des Satzes 8 aus dem Artikel [2].

7. Es gelte für eine gewisse komplexe Zahl k_1 die Ungleichung $\varrho(\mathbf{P}_{k_1}^{-1}\mathbf{Q}_{k_1}) < 1$. Dann liegt der Parameter k_0 , für welchen der Spektralradius $\varrho(\mathbf{P}_k^{-1}\mathbf{Q}_k)$ sein Minimum annimmt, innerhalb des Kreises $|k - k_1| < |k_1|$.

Beweis. Man bezeichne $\mathbf{P}_{k_1} = \tilde{\mathbf{P}}_1$, $\mathbf{Q}_{k_1} = \tilde{\mathbf{Q}}_1$. Es ist also $\varrho(\tilde{\mathbf{P}}_1^{-1}\tilde{\mathbf{Q}}_1) < 1$. Man definiere jetzt die Matrizen

$$\tilde{\mathbf{P}}_m = m\tilde{\mathbf{P}}_1, \quad \tilde{\mathbf{Q}}_m = (m - 1)\tilde{\mathbf{P}}_1 + \tilde{\mathbf{Q}}_1.$$

Wenn man $mk_1 = k$ legt, dann gelten die Beziehungen

$$\begin{aligned} \tilde{\mathbf{P}}_m &= m\tilde{\mathbf{P}}_1 = m\mathbf{P}_{k_1} = mk_1\mathbf{P}_1 = k\mathbf{P}_1 = \mathbf{P}_k, \\ \tilde{\mathbf{Q}}_m &= (m - 1)\tilde{\mathbf{P}}_1 + \tilde{\mathbf{Q}}_1 = (m - 1)\mathbf{P}_{k_1} + \mathbf{Q}_{k_1} = \\ &= (m - 1)k_1\mathbf{P}_1 + (k_1 - 1)\mathbf{P}_1 + \mathbf{Q}_1 = (mk_1 - 1)\mathbf{P}_1 + \mathbf{Q}_1 = \\ &= (k - 1)\mathbf{P}_1 + \mathbf{Q}_1 = \mathbf{Q}_k. \end{aligned}$$

Aus dem Satz 8 des Artikels [2] wissen wir, dass der Spektralradius $\varrho(\tilde{\mathbf{P}}_m^{-1}\tilde{\mathbf{Q}}_m)$ sein Minimum im Kreis $|m - 1| < 1$ annimmt. Da nun $\varrho(\tilde{\mathbf{P}}_m^{-1}\tilde{\mathbf{Q}}_m) = \varrho(\mathbf{P}_k^{-1}\mathbf{Q}_k)$ für $m = k/k_1$ ist, nimmt der Spektralradius $\varrho(\mathbf{P}_k^{-1}\mathbf{Q}_k)$ sein Minimum im Kreis $|k/k_1 - 1| < 1$, d. h. im Kreis $|k - k_1| < |k_1|$, an, was zu beweisen war.

Bemerkung 3. Wenn $\varrho(\mathbf{P}_1^{-1}\mathbf{Q}_1) < 1$ und $\varrho(\mathbf{P}_{k_1}^{-1}\mathbf{Q}_{k_1}) < 1$ für $|k_1| < 1$ ist, gibt der Satz 7 ein kleineres Gebiet für den Optimalparameter an, als der Satz 8 aus dem Artikel [2].

Aus dem Satz 6 des Artikels [2] wissen wir, dass der Optimalparameter k_0 immer irgendetwas von den Elementen der Menge M_2 oder M_3 gleich ist, d. h. $k_0 \in M_2 \cup M_3$. Die Menge M_2 enthält dabei die Zahlen k^{ij} , für die die Beziehungen

$$(27) \quad g_i(k^{ij}) = g_j(k^{ij}) \geq g_l(k^{ij}), \quad l = 1, \dots, n, \quad l \neq i, \quad l \neq j$$

gelten. Die Menge M_3 enthält ferner nur die Zahlen k^{ijk} , für die die Beziehungen

$$(28) \quad g_i(k^{ijk}) = g_j(k^{ijk}) = g_k(k^{ijk}) \geq g_l(k^{ijk}), \quad l = 1, \dots, n, \quad l \neq i, \quad l \neq j, \quad l \neq k$$

gelten. Der folgende Satz charakterisiert geometrisch die Elemente der Mengen M_2 und M_3 .

8. a) Es ist $k^{ij} \in M_2$ genau dann, wenn k^{ij} die Ungleichungen

$$(29) \quad \begin{aligned} \alpha_{il}k_1^{ij} + \beta_{il}k_2^{ij} + \gamma_{il} &\geq 0, \\ \alpha_{jl}k_1^{ij} + \beta_{jl}k_2^{ij} + \gamma_{jl} &\geq 0 \end{aligned}$$

erfüllt, wo $l = 1, \dots, n$, $l \neq i$, $l \neq j$ ist.

b) Es ist $k^{ijk} \in M_3$ genau dann, wenn k^{ijk} die Ungleichungen

$$(30) \quad \begin{aligned} \alpha_{il}k_1^{ijk} + \beta_{il}k_2^{ijk} + \gamma_{il} &\geq 0, \\ \alpha_{jl}k_1^{ijk} + \beta_{jl}k_2^{ijk} + \gamma_{jl} &\geq 0, \\ \alpha_{kl}k_1^{ijk} + \beta_{kl}k_2^{ijk} + \gamma_{kl} &\geq 0 \end{aligned}$$

erfüllt.

Dabei hat die Ungleichung

$$(31) \quad \alpha_{rs}k_1 + \beta_{rs}k_2 + \gamma_{rs} \geq 0$$

diese geometrische Bedeutung: der Punkt $k = k_1 + k_2i$ liegt in jener abgeschlossenen, durch die Gerade p_{rs} begrenzten Halbebene, die den Koordinatenursprung enthält (bzw. nicht enthält), wenn $|A_r| > |A_s|$ (bzw. $|A_r| < |A_s|$) ist. Falls $|A_r| = |A_s|$ ist, sagt die Ungleichung (31), dass der Punkt $k = k_1 + ik_2$ in jener abgeschlossenen durch die Gerade p_{rs} begrenzten Halbebene liegt (die Gerade p_{rs} geht in diesem Falle durch den Punkt $k = 0$), die den Punkt A_r enthält.

Beweis. a) Der Punkt k^{ij} liegt in der Menge M_2 , wenn (27) gilt. Aus den Ungleichungen $g_i(k^{ij}) \geq g_i(k^{ij})$, $g_j(k^{ij}) \geq g_j(k^{ij})$ und aus der Beziehung (4) folgen nach einer Umformung die Ungleichungen (29). Ähnlicherweise bekommt man auch die Ungleichungen (30). Es ist ferner klar, dass für $|A_r| > |A_s|$ $\gamma_{rs} > 0$ ist und dass in der Halbebene (31) der Punkt $k = 0$ liegt. Wenn im Gegenteil $|A_r| < |A_s|$ ist, ist $\gamma_{rs} < 0$ und der Punkt $k = 0$ liegt nicht in der Halbebene (31). Wenn $|A_r| = |A_s|$ ist, ist $\gamma_{rs} = 0$ und nach Einsetzung in die Ungleichung (31) für $\gamma_{rs} = 0$, $k_1 = R_r$, $k_2 = J_r$ bekommt man mit Rücksicht auf die Gleichung $R_r^2 + J_r^2 = R_s^2 + J_s^2$ sukzessiv die Ungleichungen

$$\begin{aligned} (R_r - R_s)R_r + (J_r - J_s)J_r &\geq 0, \\ \frac{1}{2}R_r^2 + \frac{1}{2}J_r^2 + \frac{1}{2}R_s^2 + \frac{1}{2}J_s^2 - R_rR_s - J_rJ_s &\geq 0, \\ (R_r - R_s)^2 + (J_r - J_s)^2 &\geq 0. \end{aligned}$$

Die letzte Ungleichung ist offensichtlich erfüllt, sodass der Punkt A_r in der Halbebene (31) liegt. Dadurch ist der Satz 8 bewiesen.

Nun werden wir noch einige geometrische Konstruktionen beschreiben.

I. Die Konstruktion der Halbebene

$$(5) \quad p_i \equiv 2R_i k_1 + 2J_i k_2 + |A_i|^2 < 0.$$

Die Grenzgerade dieser Halbebene ist offensichtlich die auf die Richtung A_i senkrechte Gerade, die durch den Punkt $-A_i/2$ geht. Offensichtlich enthält diese Halbebene den Koordinatenursprung nicht. Daraus folgt ihre Konstruktion.

II. Die Konstruktion der Geraden p_{ij} . Die Gerade p_{ij} hat die Gleichung

$$\alpha_{ij} k_1 + \beta_{ij} k_2 + \gamma_{ij} = 0,$$

d. h.

$$(R_i - R_j) k_1 + (J_i - J_j) k_2 + \frac{1}{2}(|A_i|^2 - |A_j|^2) = 0.$$

Es ist klar, dass diese Gerade senkrecht auf die Richtung $A_i - A_j$ ist und durch den Punkt $A = -\frac{1}{2}(A_i + A_j)$ geht.

III. Die Konstruktion des Punktes k^{ijk} . Der Punkt k^{ijk} existiert, wie wir schon wissen, nur im Falle, wenn die Punkte A_i, A_j, A_k nicht auf einer Geraden liegen. Dann schneiden sich die Geraden p_{ij}, p_{jk}, p_{ki} genau im Punkte k^{ijk} .

IV. Die Konstruktion des Punktes k^{ij} auf der Geraden p_{ij} . Der Punkt k^{ij} ist der Schnittpunkt der Geraden p_{ij} und der Kreislinie K_{ij} , die durch die Punkte $0, -A_j, -A_i$ geht. Der Punkt k^{ij} ist derjenige von den beiden Schnittpunkten der Kreislinie K_{ij} und der Geraden p_{ij} , der auf dem Segment $-A_j, -A_i$ liegt, das den Punkt $k = 0$ nicht enthält (siehe Abb. 1).

Beweis. a) Es sei $|A_i| < |A_j|$. Definiert man in der komplexen Ebene eine Kreisinvolution durch die Beziehung

$$(32) \quad u = \frac{A_i}{k} + 1$$

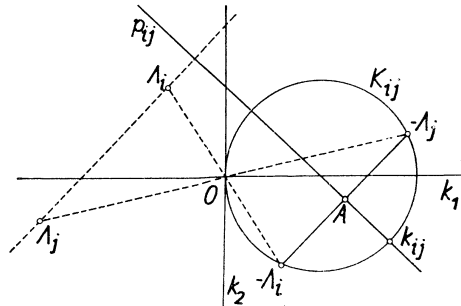


Abb. 1.

(siehe Beweis des Satzes 3 in [2]), dann entspricht der Geraden p_{ij} mit der Gleichung

$$(33) \quad p_{ij} \equiv \alpha_{ij} k_1 + \beta_{ij} k_2 + \gamma_{ij} = 0$$

eine Kreislinie mit der Gleichung

$$(34) \quad (u_1 - d_1)^2 + (u_2 - d_2)^2 - (d_1 + 1)^2 - d_2^2 = 0,$$

wo

$$(35) \quad d_1 = \frac{\alpha_{ij}R_i + \beta_{ij}J_i}{2\gamma_{ij}} - 1, \quad d_2 = \frac{\beta_{ij}R_i - \alpha_{ij}J_i}{2\gamma_{ij}}$$

ist. Wie wir schon aus dem Artikel [2] wissen, entspricht der Minimalpunkt k^{ij} der Funktion g_i auf der Geraden p_{ij} in der Kreisinverson (32) dem Punkte u^{ij} der Kreislinie (34), der dem Koordinatenursprung am nächsten liegt, d. h. dem Schnittpunkte dieser Kreislinie mit der durch den Mittelpunkt der Kreislinie (34) und durch den Nullpunkt (d. h. durch die Punkte $[0, 0]$, $[-d_1, d_2]$) gehenden Geraden. Diese Gerade hat die folgende Gleichung:

$$(36) \quad d_2u_1 + d_1u_2 = 0.$$

In der Kreisinverson (32) entspricht nun dieser Geraden die Kreislinie K_{ij} mit der Gleichung

$$(37) \quad K_{ij} \equiv \left(k_1 + \frac{d_2R_i + d_1J_i}{2d_2}\right)^2 + \left(k_2 - \frac{d_1R_i - d_2J_i}{2d_2}\right)^2 - \left(\frac{d_2R_i + d_1J_i}{2d_2}\right)^2 - \left(\frac{d_1R_i - d_2J_i}{2d_2}\right)^2 = 0.$$

Nach Einsetzung in die Gleichung (6) und mit Hilfe der Beziehungen (35) stellt man sofort fest, dass auf der Kreislinie K_{ij} die Punkte 0 , $-A_i$, $-A_j$ liegen. Daraus folgt die Konstruktion der Kreislinie K_{ij} (siehe Abb. 1). (Der Mittelpunkt dieser Kreislinie liegt immer auf der Geraden p_{ij} , da diese Gerade durch den Punkt $-\frac{1}{2}(A_i + A_j)$ geht und senkrecht zur Richtung $A_i - A_j$ ist.)

Es ist klar, dass der Punkt k^{ij} nun ein Schnittpunkt der Geraden p_{ij} und der Kreislinie K_{ij} ist, denn der Punkt k^{ij} entspricht in der Kreisinverson (32) dem Schnittpunkt u^{ij} der Kreislinie (34) und der Geraden (36).

Man kann jetzt beweisen, dass der Punkt k^{ij} jener von der Schnittpunkten der Gerade p_{ij} und der Kreislinie K_{ij} ist, der auf dem den Punkt 0 nicht enthaltenden Segment liegt.

Wir untersuchen jetzt die Punkte 0 , u^{ij} , z^{ij} , wo $z^{ij} = \left[\frac{|A_i| + |A_j|}{|A_j|}\right] u^{ij}$ ist. Diese Punkte liegen alle auf der Geraden (36) in der angeführten Reihenfolge, denn es ist $\left(\frac{|A_i| + |A_j|}{|A_j|}\right) > 1$. Man kann leicht feststellen, dass dem Punkte z^{ij} in der Kreisinverson (32) der Punkt $-A_j$ entspricht. Für den Punkt u^{ij} gelten nämlich die Beziehungen

$$u_1^{ij} = \frac{|A_j|^2 - (R_iR_j + J_iJ_j)}{|A_j|(|A_i| + |A_j|)}, \quad u_2^{ij} = \frac{R_iJ_j - J_iR_j}{|A_j|(|A_i| + |A_j|)}$$

(siehe die Formeln (16), (17)). Aus der Beziehung (32) nach Einsetzung der Zahl z^{ij} für u und mit Hilfe der oben angeführten Beziehungen für u bekommt man nach einer Umformung die Gleichung $-A_j = A_i(z^{ij} - 1)$, was wir beweisen sollten.

Den Punkten $\infty, 0, u^{ij}, z^{ij}, \infty$, die in der angeführten Reihenfolge auf der Geraden (5) liegen, entsprechen dann in der Kreisinverson (32) Punkte der Kreislinie K_{ij} in folgender Reihenfolge: $0, -A_i, k^{ij}, -A_j, 0$. Hiervon folgt, dass der Punkt k^{ij} jener von den Schnittpunkten der Geraden p_{ij} und der Kreislinie K_{ij} ist, der auf dem den Punkt 0 nicht enthaltenden Segment $-A_i, -A_j$ liegt.

b) Es sei jetzt $|A_i| = |A_j|$. Die Gerade p_{ij} hat dann die Gleichung

$$(R_i - R_j) k_1 + (J_i - J_j) k_2 = 0.$$

Daraus folgt ihre Konstruktion. Wenn man diese Gerade in der Form

$$(33') \quad p_{ij} \equiv Ak + \bar{A}\bar{k} = 0$$

schreibt, wo $A = \beta_{ij} + i\alpha_{ij}$ ist, dann bekommt man mit Hilfe der Beziehung (32) und nach einer Umformung die Gerade

$$(34') \quad -\bar{A}\bar{A}_i u + AA_i \bar{u} - AA_i + \bar{A}\bar{A}_i = 0,$$

die offensichtlich durch den Punkt $u = 1$ geht. Die durch den Nullpunkt gehende und zu der Geraden (34') senkrechte Gerade schneidet diese Gerade im Punkte u^{ij} (dieser Punkt entspricht in der Kreisinverson dem Minimalpunkt k^{ij} der Funktionen g_i, g_j auf der Geraden p_{ij}). Es handelt sich um die Gerade

$$(36') \quad AA_i \bar{u} + \bar{A}\bar{A}_i u = 0.$$

Dieser Geraden entspricht in der Kreisinverson allgemein wieder die Kreislinie mit der Gleichung

$$(37') \quad K_{ij} \equiv k\bar{k}(AA_i + \bar{A}\bar{A}_i) + |A_i|^2 (Ak + \bar{A}\bar{k}) = 0.$$

Diese Kreislinie schneidet also die Gerade p_{ij} in Punkten 0, k^{ij} . Man kann leicht beweisen, dass auf der Kreislinie K_{ij} die Punkte 0, $-A_i, -A_j$ liegen. Für die Punkte 0, $-A_i$ folgt dies sofort nach der Einsetzung in (37). Wenn $k = -A_j$ ist, bekommt man nach einer Umformung die Gleichung

$$|A_j|^2 (AA_i + \bar{A}\bar{A}_i) - |A_i|^2 (AA_j + \bar{A}\bar{A}_j) = 0.$$

Da $|A_i| = |A_j| \neq 0$ ist, bekommt man sukzessiv

$$AA_i + \bar{A}\bar{A}_i = AA_j + \bar{A}\bar{A}_j,$$

$$Re A(A_i - A_j) = 0,$$

$$\beta_{ij}\alpha_{ij} - \alpha_{ij}\beta_{ij} = 0.$$

Das ist aber eine gültige Beziehung.

Die Konstruktion ist also ein Spezialfall der Konstruktion für $|A_i| \neq |A_j|$.

Es sei noch bemerkt, dass für $A_j = -A_i$ die Kreislinie (37') in die durch die Punkte $A_i, 0, -A_i$ gehende Gerade übergeht. (In diesem Falle ist nämlich $AA_i + \bar{A}\bar{A}_i = 0$.) Diese Kreislinie fließt aber mit der Gerade (33') zusammen, sodass kein Schnittpunkt k^{ij} existiert.

Beispiel. Setzt man z.B. voraus, dass die Matrix $P_1^{-1}Q_1$ die folgenden Eigenwerte besitzt:

$$\lambda_1 = -0,4 + 1,4i, \quad \lambda_2 = -2,5 + 1,8i, \quad \lambda_3 = -2,9 - 1,2i.$$

Dann ist $\varrho(P_1^{-1}Q_1) = 3,14$, was bedeutet, dass das Iterationsverfahren (1) nicht konvergiert. Die Zahlen

$$A_1 = -1,4 + 1,4i, \quad A_2 = -3,5 + 1,8i, \quad A_3 = -3,9 - 1,2i$$

liegen offensichtlich im Kreise mit der Gleichung

$$|A + a| < b,$$

wo $a = 2,72 - 0,17i$, $b = 1,8$ ist. Laut Satz 4 liegt der Optimalparameter im Durchschnitt des Kreises

$$|k - a| < b$$

und des Gebietes K , das diesem Falle der Durchschnitt der folgenden den Koordinatensprung nicht enthaltenden Halbebenen ist:

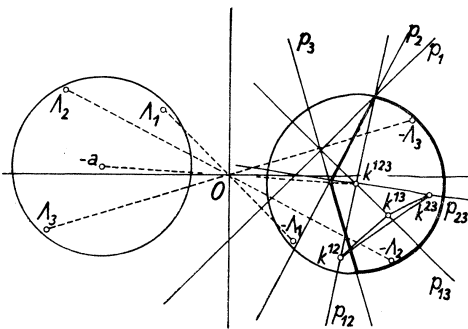


Abb. 2.

$$\begin{aligned} p_1 &\equiv -2,8k_1 + 2,8k_2 + 1,98 < 0, \\ p_2 &\equiv -7k_1 + 3,6k_2 + 3,94 < 0, \\ p_3 &\equiv -7,8k_1 - 2,4k_2 + 4,08 < 0. \end{aligned}$$

In Abb. 2 ist dieses Gebiet durch eine starke Linie gekennzeichnet. Man sieht, dass in diesem Gebiet die Zahlen

$$k^{23} = 4,3 - 0,38i,$$

$$k^{13} = 3,43 - 0,85i,$$

$$k^{123} = 2,8 - 0,16i$$

liegen, während die Zahl $k^{12} = 2,37 - 1,72i$ zu diesem Gebiet nicht gehört. Nach Satz 6 des Artikles [2] ist der Optimalparameter irgendeiner von der Zahlen k^{12} , k^{23} , k^{13} , k^{123} gleich. Mit Rücksicht auf die Tatsache, dass die Zahl k^{123} nicht im kleinsten die Punkte k^{12} , k^{23} , k^{13} enthaltendem konvexen Vieleck liegt, kann nach Satz 7 des Artikles [2] der Punkt k^{123} nicht als Optimalparameter in Betracht kommen. Da ferner $k^{12} \notin K$ ist, kommen nur die Punkte k^{23} , k^{13} in Betracht. Nach Einsetzung dieser Zahlen in die Funktionen g_1, g_2, g_3 (oder mit Hilfe des Satzes 8)

stellt man leicht fest, dass der Optimalparameter der Punkt k^{13} ist und dass der optimale Spektralradius der Zahl

$$\varrho(\mathbf{P}_{k^{13}}^{-1}\mathbf{Q}_{k^{13}}) = \sqrt{[g_1(k^{13})]} = 0,696$$

gleich ist. Dass bedeutet aber, dass das Iterationsverfahren für den Parameter k^{13} gut konvergiert.

In Abb. 2 ist die Konstruktion der Punkte k^{12} , k^{13} , k^{23} , k^{123} und der den Optimalparameter enthaltenden Gebiete angedeutet.

Literaturverzeichnis

- [1] *Isaacson, E., Keller, H. B.:* Analysis of Numerical Methods, John Wiley & Sons, Inc., New York, London, Sydney, 1966.
 [2] *Šisler, M.:* Über die Konvergenzbeschleunigung komplexer Iterationsverfahren, Aplikace matematiky 15, 156–166 (1970).

Souhrn

O KONVERGENCI ITERAČNÍCH PROCESŮ

MIROSLAV ŠISLER

Práce se zabývá otázkami konvergence iteračních procesů pro řešení soustavy m lineárních rovnic $\mathbf{Ax} = \mathbf{b}$. Předpokládá se, že je dán iterační proces definovaný vzorcem

$$(1) \quad \mathbf{x}_{v+1} = \mathbf{P}_1^{-1}\mathbf{Q}_1\mathbf{x}_v + \mathbf{P}_1^{-1}\mathbf{b}, \quad v = 0, 1, 2, \dots,$$

kde $\mathbf{A} = \mathbf{P}_1 - \mathbf{Q}_1$ je jistý rozklad matice \mathbf{A} . Přitom je lhostejné, zda tento iterační proces konverguje či nikoliv (tj. uvažují se i takové rozklady $\mathbf{A} = \mathbf{P}_1 - \mathbf{Q}_1$, že spektrální poloměr matice $\mathbf{P}_1^{-1}\mathbf{Q}_1$ je větší nebo roven 1).

V práci se vyšetřuje modifikovaný iterační proces definovaný vzorcem

$$(2) \quad \mathbf{x}_{v+1} = \mathbf{P}_k^{-1}\mathbf{Q}_k\mathbf{x}_v + \mathbf{P}_k^{-1}\mathbf{b}, \quad v = 0, 1, 2, \dots,$$

kde $\mathbf{P}_k = k\mathbf{P}_1$, $\mathbf{Q}_k = (k - 1)\mathbf{P}_1 + \mathbf{Q}_1$ a k je nenulový komplexní parametr (původní proces je tedy speciálním případem pro $k = 1$).

V článku jsou dokázány nutné a postačující podmínky pro to, aby existovala neprázdná množina K parametrů k , pro něž iterační proces (2) konverguje (věta 2). Dále je zkoumána otázka polohy optimálního parametru v komplexní rovině (pokud $K \neq \emptyset$), tj. parametru k , pro něž je spektrální poloměr matice $\mathbf{P}_k^{-1}\mathbf{Q}_k$ minimální (věty 5, 6 a 7).

V článku je uvedena i geometrická konstrukce optimálního parametru. Dosažené výsledky jsou demonstrovány na numerickém příkladě.

Anschrift des Verfassers: RNDr. Miroslav Šisler, CSc, Matematický ústav ČSAV v Praze, Žitná 25, Praha 1.