

On Smallest Isolated Gerschgorin Disks for Eigenvalues*

By

RICHARD S. VARGA**

* To ALSTON S. HOUSEHOLDER on his sixtieth birthday.

** The results of this paper were announced at the SIAM Matrix Symposium at Gatlinburg, Tenn., in April, 1964.

§ 1. Introduction

Let $A = (a_{i,j})$ be a fixed $n \times n$ complex matrix, and let $X(\mathbf{x}) = \text{diag}(x_1, x_2, \dots, x_n)$, where $x_i > 0$, $1 \leq i \leq n$. Upon forming the matrix $X^{-1}(\mathbf{x}) A X(\mathbf{x})$, the radii $A_i(\mathbf{x})$ of the Gerschgorin disks $|z - a_{i,i}| \leq A_i(\mathbf{x})$, determined from row sums of $X^{-1}(\mathbf{x}) A X(\mathbf{x})$, are given by

$$(1.1) \quad A_i(\mathbf{x}) = \left(\sum_{\substack{j=1 \\ j \neq i}}^n |a_{i,j}| x_j \right) / x_i, \quad 1 \leq i \leq n.$$

Further, let

$$(1.2) \quad d_{k,j} = |a_{k,k} - a_{j,j}|, \quad 1 \leq j, \quad k \leq n.$$

In analogy to [4], we make the following definition.

Definition 1. Let P_k be the set of all vectors $\mathbf{x} > \mathbf{0}$ such that

$$(1.3) \quad d_{k,j} - A_j(\mathbf{x}) - A_k(\mathbf{x}) \geq 0 \quad \text{for all } j \neq k.$$

If P_k is non-null, then the matrix A admits, under diagonal similarity transformations, an *isolated k -th Gerschgorin disk*.

Assuming there exists an $\mathbf{x}_0 \in P_k$ for which strict inequality is valid in (1.3) for all $j \neq k$, then it follows from GERSCHGORIN'S original paper [3] that there is *exactly* one eigenvalue of A in the disk $|z - a_{k,k}| \leq A_k(\mathbf{x}_0)$. In order to find improved eigenvalue bounds for this isolated eigenvalue, it is quite natural to think in terms of varying the vector $\mathbf{x}_0 \in P_k$ so as to decrease $A_k(\mathbf{x}_0)$. Indeed, this approach was already suggested by GERSCHGORIN [3], and has subsequently been employed numerically with success by TAUSSKY [8] and WILKINSON [12]. Yet, in terms of finding the *smallest* possible radius μ for this disk, obviously given by

$$(1.4) \quad \mu = \inf_{\mathbf{x} \in P_k} A_k(\mathbf{x}),$$

only recently did HENRICI [4] give a convergent algorithm for finding μ in a special tridiagonal case. The first object of this paper is to extend (Theorem 1) this convergent algorithm to the *general* case. We do this by using the theory of M -matrices [5].¹ Next, we also give a convergent non-linear Gauss-Seidel

¹ The next paper, by Professor JOHN TODD, gives an alternate matrix proof of this extension. The author wishes to express his thanks to Professor TODD for helpful discussions.

type algorithm (Theorem 2) which, for large order matrices, may be more useful in practical computations.

The results for the specific matrix A actually are valid for a set of matrices \mathring{Q}_A . Continuing the investigations of [10], we show (Theorem 3) that every point of an annulus $\tau \leq |z| \leq \mu$ is an isolated eigenvalue of some $n \times n$ matrix $B \in \mathring{Q}_A$. Finally, a numerical example is included.

§ 2. Some basic Lemmas

Since our goal is to obtain bounds for certain eigenvalues of the matrix A , we may first assume without loss of generality that A is *irreducible*². Next, if some set P_k is non-null, we may further assume, again without loss of generality, that $k=1$ so that the *first* Gerschgorin disk $|z - a_{1,1}| \leq A_1(\mathbf{x})$ can be isolated. This reduction to the case $k=1$ can obviously be accomplished by means of a similarity transformation by a permutation matrix applied to A . We henceforth assume that the set P_1 is non-null, and thus from Definition 1,

$$(2.1) \quad d_{1,j} - A_j(\mathbf{x}) - A_1(\mathbf{x}) \geq 0, \quad 2 \leq j \leq n,$$

for all $\mathbf{x} \in P_1$. Since $A_j(\beta \mathbf{x}) = A_j(\mathbf{x})$ for any $\beta > 0$ and any $1 \leq j \leq n$, we may assume, as a final normalization, that $x_1 = 1$ for any $\mathbf{x} \in P_1$. Note from (1.1) that $A_1(\mathbf{x})$ is then linear in the variables x_2, x_3, \dots, x_n .

Let $Q = (q_{i,j})$ now be a fixed real $n \times n$ matrix, whose entries are defined by

$$(2.2) \quad \begin{cases} q_{i,i} = d_{1,i} = |a_{1,1} - a_{i,i}|, & 1 \leq i \leq n, \\ q_{1,j} = |a_{1,j}|, & 2 \leq j \leq n; \quad q_{i,j} = -|a_{i,j}|, \quad i \neq j, \quad i \neq 1. \end{cases}$$

The matrix Q is, by construction, irreducible. For notational simplicity in what is to follow, we introduce the following partitioning of Q :

$$(2.3) \quad Q = \left[\begin{array}{c|c} 0 & \hat{\boldsymbol{\alpha}}^T \\ \hline -\hat{\boldsymbol{a}} & \tilde{Q} \end{array} \right],$$

where \tilde{Q} is an $(n-1) \times (n-1)$ principal submatrix of Q , $\hat{\boldsymbol{\alpha}}^T = (|a_{1,2}|, |a_{1,3}|, \dots, |a_{1,n}|)$ and $\hat{\boldsymbol{a}}^T = (|a_{2,1}|, |a_{3,1}|, \dots, |a_{n,1}|)$. Let I_{n-1} denote the $(n-1) \times (n-1)$ identity matrix, and let $\hat{\mathbf{y}}$ be the column vector with $n-1$ components obtained from \mathbf{y} such that $\hat{\mathbf{y}}^T = (y_2, \dots, y_n)$. Conversely, given $\hat{\mathbf{y}}$, let \mathbf{y} denote the unique column vector with n components determined by $\hat{\mathbf{y}}$ and $y_1 = 1$. With this notation, we have

Lemma 1. $\mathbf{x} \in P_1$ if and only if $\hat{\mathbf{x}} > \hat{\mathbf{0}}$ and

$$(2.4) \quad (\tilde{Q} - A_1(\mathbf{x}) I_{n-1}) \hat{\mathbf{x}} \geq \hat{\mathbf{a}}.$$

Moreover, $\mathbf{x} \in P_1$ implies

$$(2.5) \quad Q\mathbf{x} \geq A_1(\mathbf{x})\mathbf{x}.$$

² An $n \times n$ matrix A is *irreducible* if there exists no $n \times n$ permutation matrix P such that $PAP^T = \begin{bmatrix} A_{1,1} & A_{1,2} \\ 0 & A_{2,2} \end{bmatrix}$, where $A_{1,1}$ and $A_{2,2}$ are square nonvoid submatrices. Otherwise, A is *reducible*.

Proof. Both parts of this lemma are a direct consequence of Definition 1 and our normalization in P_1 .

Based on a simple argument using the irreducibility of the matrix A and the inequalities of (2.1), we can show that there exist constants r_i, s_i such that

$$(2.6) \quad 0 < r_i \leq x_i \leq s_i, \quad 2 \leq i \leq n,$$

for any $\mathbf{x} \in P_1$. Thus, as $A_1(\mathbf{x})$ is linear in the components x_2, \dots, x_n , we see from (2.6) and (2.4) that the set of vectors P_1 is compact. Hence, there necessarily exist vectors \mathbf{z} and \mathbf{w} in P_1 such that

$$(2.7) \quad \mu \equiv \inf_{\mathbf{x} \in P_1} A_1(\mathbf{x}) = A_1(\mathbf{z}), \quad \text{and} \quad \sigma \equiv \sup_{\mathbf{x} \in P_1} A_1(\mathbf{x}) = A_1(\mathbf{w}).$$

From (2.5) of Lemma 1, we necessarily have that $Q\mathbf{z} \geq \mu\mathbf{z}$ and $Q\mathbf{w} \geq \sigma\mathbf{w}$. We shall actually show by means of the next lemmas that μ and σ are eigenvalues of Q . Moreover, we shall show that the eigenvectors \mathbf{z} and \mathbf{w} corresponding respectively to μ and σ are uniquely determined in P_1 . Finally, the method of proof we employ gives rise to an algorithm for determining $\mu = A_1(\mathbf{z})$ and \mathbf{z} which, in the tridiagonal case, reduces to the algorithm of [4].

The next lemma establishes a connection between our problem and the theory of M -matrices³ [5].

Lemma 2. For any real number $t \leq \sigma$, $(\tilde{Q} - tI_{n-1})$ is an M -matrix, and $(\tilde{Q} - tI_{n-1})^{-1} \hat{\mathbf{a}} > \hat{\mathbf{0}}$.

Proof. We consider first the special case $t = \sigma$. With $\sigma = A_1(\mathbf{w})$, $\mathbf{w} \in P_1$, let $\tilde{X} = \text{diag}(w_2, w_3, \dots, w_n)$, and let $\tilde{E} \equiv \tilde{X}^{-1}(\tilde{Q} - \sigma I_{n-1})\tilde{X} \equiv (e_{i,j})$, $1 \leq i, j \leq n-1$. It follows that the diagonal entries of \tilde{E} , given by $e_{i,i} = d_{1,i+1} - \sigma$, are positive real numbers from (2.1) since $\mathbf{w} \in P_1$, and the off-diagonal entries of \tilde{E} , given by $e_{i,j} = -|a_{i+1,j+1}|w_{j+1}/w_{i+1}$, are non-positive real numbers satisfying

$$(2.8) \quad \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |e_{i,j}| = A_{i+1}(\mathbf{w}) - (|a_{i+1,1}|/w_{i+1}), \quad 1 \leq i \leq n-1.$$

Thus, since $\mathbf{w} \in P_1$, we have from (2.1) and (2.8) that

$$(2.9) \quad e_{i,i} - \sum_{\substack{j=1 \\ j \neq i}}^{n-1} |e_{i,j}| \geq |a_{i+1,1}|/w_{i+1} \geq 0, \quad 1 \leq i \leq n-1,$$

which proves that \tilde{E} is *diagonally dominant* [9, p. 23]. Since A is irreducible, not all of the components of $\hat{\mathbf{a}}$ can vanish, so that for some l , $1 \leq l \leq n-1$,

$$(2.10) \quad e_{l,l} > \sum_{i \neq l} |e_{l,i}|.$$

If \tilde{Q} is irreducible, so is \tilde{E} , and thus, with (2.10), \tilde{E} is *irreducibly diagonally dominant* [9, p. 23]. Because of the sign pattern of the entries of \tilde{E} , it follows

³ A real $n \times n$ matrix $B = (b_{i,j})$ with $b_{i,j} \leq 0$ for all $i \neq j$ is defined to be an M -Matrix if B is non-singular, and $B^{-1} \geq 0$.

[9, p. 85] that \tilde{E} is an irreducible M -matrix, and as such $\tilde{E}^{-1} > 0$. This of course implies that $\tilde{Q} - \sigma I_{n-1}$ is an irreducible M -matrix with $(\tilde{Q} - \sigma I_{n-1})^{-1} > 0$. Again, since $\hat{\mathbf{a}}$ has at least one component positive, this positivity insures that $(\tilde{Q} - \sigma I_{n-1})^{-1} \hat{\mathbf{a}} > \hat{\mathbf{0}}$. But the above argument is actually valid for $t \leq \sigma$, and we have thus established the desired result in the case that \tilde{Q} is irreducible. If \tilde{Q} is reducible, the desired conclusions are established in a similar (but more tedious) way, the proofs now making use of the reduced normal form [2, p. 74; 9, p. 46] of a reducible matrix. For brevity, we have omitted the proof in the reducible case.

Corollary 1. If λ is any eigenvalue of \tilde{Q} , then $\text{Re } \lambda > \sigma$.

Proof. If B is an M -matrix, it is known [5] that any eigenvalue α of B satisfies $\text{Re } \alpha > 0$. The result then follows from Lemma 2.

For any $t \leq \sigma$, we define the vector \mathbf{y}_t from $(\tilde{Q} - tI_{n-1}) \hat{\mathbf{y}}_t \equiv \hat{\mathbf{a}}$, and it follows from Lemma 2 that $\hat{\mathbf{y}}_t > \hat{\mathbf{0}}$. Thus, as $A_1(\mathbf{y}_t) > 0$, we see that the function $g(t)$ defined by $g(t) \equiv A_1(\hat{\mathbf{y}}_t)$ is positive for all $t \leq \sigma$. We now show that $g(t)$ is strictly increasing and strictly convex upward.⁴

Lemma 3. For any $t < \sigma$ and for $\varepsilon > 0$ sufficiently small,

$$(2.11) \quad g(t + \varepsilon) = g(t) + \sum_{k=1}^{\infty} c_k \varepsilon^k, \quad \text{where } c_k > 0 \text{ for all } k \geq 1.$$

Proof. By definition, $(\tilde{Q} - (t + \varepsilon)I_{n-1}) \hat{\mathbf{y}}_{t+\varepsilon} = \hat{\mathbf{a}}$, and $(\tilde{Q} - tI_{n-1}) \hat{\mathbf{y}}_t = \hat{\mathbf{a}}$. If we let $B \equiv (\tilde{Q} - tI_{n-1})^{-1}$, then it follows that

$$(2.12) \quad \hat{\mathbf{y}}_{t+\varepsilon} = \{I - \varepsilon B\}^{-1} B \hat{\mathbf{a}} = \{I - \varepsilon B\}^{-1} \hat{\mathbf{y}}_t.$$

Thus, for $\varepsilon > 0$ sufficiently small,

$$(2.13) \quad \hat{\mathbf{y}}_{t+\varepsilon} = \hat{\mathbf{y}}_t + \sum_{k=1}^{\infty} \varepsilon^k B^{k+1} \hat{\mathbf{a}}.$$

From Lemma 2, we know that $B \geq 0$ and $B \hat{\mathbf{a}} > \hat{\mathbf{0}}$, from which it follows inductively that $B^{k+1} \hat{\mathbf{a}} > \hat{\mathbf{0}}$. Recalling that $A_1(\mathbf{x})$ is linear and that $A_1(\mathbf{x}) = \hat{\boldsymbol{\alpha}}^T \hat{\mathbf{x}}$, then the series expansion for $g(t + \varepsilon) = A_1(\hat{\mathbf{y}}_{t+\varepsilon})$ in (2.11) directly follows from (2.13) with

$$(2.14) \quad c_k \equiv \hat{\boldsymbol{\alpha}}^T B^{k+1} \hat{\mathbf{a}}, \quad k \geq 1.$$

Thus, the coefficients c_k are all positive, completing the proof.

Lemma 4. $\mathbf{y}_t \in P_1$ if and only if $t \geq g(t)$.

Proof. Since $\hat{\mathbf{y}}_t > \hat{\mathbf{0}}$ for all $t \leq \sigma$, it follows from Lemma 1 that $\mathbf{y} \in P_1$ if and only if $(\tilde{Q} - A_1(\mathbf{y})I_{n-1}) \hat{\mathbf{y}}_t \geq \hat{\mathbf{a}}$, or equivalently, $(\tilde{Q} - g(t)I_{n-1}) \hat{\mathbf{y}}_t \geq \hat{\mathbf{a}}$. Now,

$$(2.15) \quad (\tilde{Q} - g(t)I_{n-1}) \hat{\mathbf{y}}_t = \hat{\mathbf{a}} + (t - g(t)) \hat{\mathbf{y}}_t,$$

and it therefore follows that $(\tilde{Q} - g(t)I_{n-1}) \hat{\mathbf{y}}_t \geq \hat{\mathbf{a}}$ if and only if $t \geq g(t)$, which completes the proof.

We now show that there exist values of t such that $t \geq g(t)$.

⁴ More precisely, $g(t)$ is strictly absolutely monotonic [11].

Lemma 5. For any $\mathbf{x} \in P_1$, $t \geq g(t)$ where $A_1(\mathbf{x}) \equiv t$. Moreover, $(\tilde{Q} - tI_{n-1})\hat{\mathbf{x}} \geq \hat{\mathbf{a}}$ with strict inequality for some component if and only if $t > g(t)$.

Proof. Since $\mathbf{x} \in P_1$ and $A_1(\mathbf{x}) \equiv t$, then from Lemma 1 we know that $(\tilde{Q} - tI_{n-1})\hat{\mathbf{x}} \geq \hat{\mathbf{a}}$, and hence $\hat{\mathbf{x}} \geq \hat{\mathbf{y}}_t > \hat{\mathbf{0}}$, using Lemma 2. Thus, $t = A_1(\mathbf{x}) \geq A_1(\mathbf{y}_t) = g(t)$, which, by Lemma 4, assures us that $\mathbf{y}_t \in P_1$. Moreover, if equality is valid in all components, i.e., $(\tilde{Q} - tI_{n-1})\hat{\mathbf{x}} = \hat{\mathbf{a}}$, then surely $\mathbf{x} = \mathbf{y}_t$ and thus $t = A_1(\mathbf{x}) = A_1(\mathbf{y}_t) = g(t)$. To prove the remaining part of this Lemma, we assume for simplicity that \tilde{Q} is irreducible. From the proof of Lemma 2, we then have that $(\tilde{Q} - tI_{n-1})^{-1} > 0$. If $(\tilde{Q} - tI_{n-1})\hat{\mathbf{x}} \geq \hat{\mathbf{a}}$ with *strict* inequality for some component, the positivity of $(\tilde{Q} - tI_{n-1})^{-1}$ allows us to deduce that $\hat{\mathbf{x}} > \hat{\mathbf{y}}_t$, i.e., strict inequality in *all* components, and thus $t = A_1(\mathbf{x}) > A_1(\mathbf{y}_t) = g(t)$. The extension of these results to the case when \tilde{Q} is reducible is again based on the reduced normal form of \tilde{Q} , and is omitted.

Corollary 2. There exist two positive real numbers μ_1, σ_1 with $\mu_1 = g(\mu_1)$ and $\sigma_1 = g(\sigma_1)$ such that $t \geq g(t)$ if and only if $0 < \mu_1 \leq t \leq \sigma_1$.

Proof. Since, by Lemma 5, there exist positive values of t with $t \geq g(t)$, this result then follows directly from the strictly increasing and strictly convex property of $g(t)$, established in Lemma 3. We remark that if $\mu_1 < \sigma_1$, then $t > g(t)$ if and only if $\mu_1 < t < \sigma_1$.

Lemma 6. $\mu = \mu = \inf_{\mathbf{x} \in P_1} A_1(\mathbf{x})$, and $\sigma = \sigma = \sup_{\mathbf{x} \in P_1} A_1(\mathbf{x})$. Moreover, μ and σ are eigenvalues of Q , whose eigenvectors \mathbf{y}_μ and \mathbf{y}_σ are uniquely determined in P_1 .

Proof. Since $\mathbf{y}_{\mu_1} \in P_1$ by Lemma 4 and Corollary 2, it follows that $\mu_1 \geq \inf_{\mathbf{x} \in P_1} A_1(\mathbf{x}) = \mu$. On the other hand, there is a vector $\mathbf{z} \in P_1$ such that $A_1(\mathbf{z}) = \mu$. Thus, from Lemma 5 and Corollary 2, $\mu \geq \mu_1$, and we thus have $\mu = \mu_1$. Next, since $(\tilde{Q} - \mu I_{n-1})\hat{\mathbf{y}}_\mu = \hat{\mathbf{a}}$ by definition and $g(\mu) = \mu$, it is readily verified from (2.3) that $Q\mathbf{y}_\mu = \mu\mathbf{y}_\mu$, and as $\mathbf{y}_\mu > \mathbf{0}$, then \mathbf{y}_μ is an eigenvector of Q corresponding to the eigenvalue μ . To prove a somewhat stronger form of the final result, suppose that $Q\mathbf{x} = \mu\mathbf{x}$ where \mathbf{x} is any complex vector with $\mathbf{x} \neq \mathbf{0}$. If x_1 were zero, then μ would be an eigenvalue of \tilde{Q} , which contradicts the result of Corollary 1. Thus, $x_1 \neq 0$. Forming $(\hat{\mathbf{x}}/x_1)$, it follows from $Q\mathbf{x} = \mu\mathbf{x}$ that $(\tilde{Q} - \mu I_{n-1})(\hat{\mathbf{x}}/x_1) = \hat{\mathbf{a}}$, but as $(\tilde{Q} - \mu I_{n-1})$ is non-singular, then necessarily $(\hat{\mathbf{x}}/x_1) = \hat{\mathbf{y}}_\mu > \hat{\mathbf{0}}$. Thus, $\mathbf{x} = x_1\mathbf{y}_\mu$. Similar arguments prove analogous results for $\sigma = \sigma$.

To complete this section, we point out that the Gerschgorin circles for the irreducible matrix $X^{-1}(\mathbf{y}_\mu)QX(\mathbf{y}_\mu)$ all pass through μ , i.e.,

$$(2.16) \quad |q_{i,i} - \mu| = A_i(\mathbf{y}_\mu) \quad \text{for all } 1 \leq i \leq n,$$

which is related to a result of TAUSSKY [7]. From (2.16), it follows [10, Theorem 3] that μ is a boundary point of the minimal Gerschgorin set for Q . The same is true for σ .

§ 3. First Convergence Theorem

With the lemmas of the previous section, we now prove

Theorem 1. Let A be an irreducible $n \times n$ matrix which admits a first isolated Gerschgorin disk. Then, the smallest radius μ , under all diagonal similarity

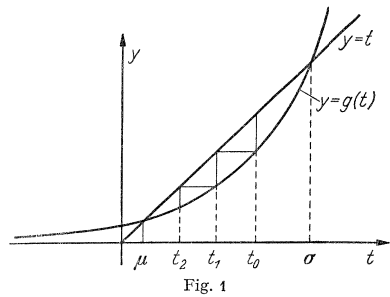
transformations, for this isolated disk is an eigenvalue of the $n \times n$ matrix Q , and its corresponding eigenvector \mathbf{y}_μ is uniquely determined in P_1 . If $\mathbf{x}_0 \in P_1$, and $Q\mathbf{x}_0 \geq A_1(\mathbf{x}_0)\mathbf{x}_0$, with strict inequality for at least one component, then the sequence of vectors $\{\mathbf{x}_i\}_{i=0}^\infty$ defined by

$$(3.1) \quad (\tilde{Q} - A_1(\mathbf{x}_i)I_{n-1}) \hat{\mathbf{x}}_{i+1} = \hat{\mathbf{a}}, \quad i \geq 0,$$

are all elements of P_1 with $\lim_{i \rightarrow \infty} \mathbf{x}_i = \mathbf{y}_\mu$, and the sequence $\{A_1(\mathbf{x}_i)\}_{i=0}^\infty$ is strictly decreasing with $\lim_{i \rightarrow \infty} A_1(\mathbf{x}_i) = \mu$.

Proof. The first part of this theorem is a restatement of Lemma 6. Next, for any $\mathbf{r} \in P_1$, it is readily verified that $Q\mathbf{r} \geq A_1(\mathbf{r})\mathbf{r}$ with strict inequality for at least one component if and only if $(\tilde{Q} - A_1(\mathbf{r})I_{n-1})\hat{\mathbf{r}} \geq \hat{\mathbf{a}}$, with strict inequality for at least one component. Thus, with Lemma 5, the hypothesis tells us that $t_0 > g(t_0)$ where $t_0 \equiv A_1(\mathbf{x}_0)$. Next, from (3.1), it follows that $\mathbf{x}_1 \equiv \mathbf{y}_{t_0}$, and we now ask if $Q\mathbf{x}_1 \geq A_1(\mathbf{x}_1)\mathbf{x}_1$, with strict inequality for some component. Since $A_1(\mathbf{x}_1) = g(t_0)$, it follows from (2.15) that

$$(3.2) \quad \begin{aligned} &(\tilde{Q} - A_1(\mathbf{x}_1)I_{n-1})\hat{\mathbf{x}}_1 \\ &= (\tilde{Q} - g(t_0)I_{n-1})\hat{\mathbf{y}}_{t_0} \\ &= \hat{\mathbf{a}} + (t_0 - g(t_0))\hat{\mathbf{y}}_{t_0} > \hat{\mathbf{a}}, \end{aligned}$$



so that inequality now holds for all $n - 1$ components. This means that the above argument applied to \mathbf{x}_0 can be applied to \mathbf{x}_1 , and thus $t_1 = A_1(\mathbf{x}_1) > g(t_1)$. Inductively, with $t_i \equiv A_1(\mathbf{x}_i)$ and $\mathbf{x}_{i+1} = \mathbf{y}_{t_i}$, it follows that

$$(3.3) \quad t_0 > g(t_0) = t_1 > g(t_1) = t_2 > \dots$$

so that the sequence $\{A_1(\mathbf{x}_i)\}_{i=0}^\infty$ is strictly decreasing. Similarly, the components of the vector sequence $\{\mathbf{x}_i\}_{i=0}^\infty$ are nonincreasing. If $\lim_{i \rightarrow \infty} t_i = \delta$, then clearly $\delta \geq \mu$, but since $t_i > g(t_i) = t_{i+1} > t_{i+2}$, it necessarily follows that $\delta = g(\delta)$, which from Corollary 2 implies that $\delta = \mu$, completing the proof.

The iterative procedure of (3.1) has a rather simple geometrical interpretation. Starting with $\mathbf{x}_0 \in P_1$ such that $\mu < A_1(\mathbf{x}_0) < \sigma$, it follows that the Gerschgorin disks $|q_{k,k} - z| \leq A_k(\mathbf{x}_0)$, $k > 1$, do not all touch the first (isolated) disk $|z| \leq A_1(\mathbf{x}_0)$. The essence of (3.1) is that we now determine a positive vector \mathbf{x}_1 such that the new Gerschgorin disks $|q_{k,k} - z| \leq A_k(\mathbf{x}_1)$, $k > 1$, are now all tangent to the old first disk $|z| \leq A_1(\mathbf{x}_0)$. This is done by decreasing the components x_j of \mathbf{x}_0 , $j > 1$, in such a way that the radii $A_j(\mathbf{x})$ all increase for $j > 1$. But then, the first disk has smaller radius, viz. $A_1(\mathbf{x}_1)$, and the process can be repeated.

The iteration of (3.1) to find μ , the smallest radius of the first isolated Gerschgorin disk, is actually the method of successive substitutions, which is represented by the staircase curve lying between $y=t$ and $y=g(t)$ in Fig. 1. We also remark that monotone convergence of the $A_1(\mathbf{x}_i)$ to μ is had for any real initial value $A_1(\mathbf{x}_0)$, with $A_1(\mathbf{x}_0) < \sigma$.

The iterative procedure of (3.4) requires at each step the solution of a system of $(n - 1)$ linear equations in $(n - 1)$ unknowns. In [4], HENRICI considered the special case where \tilde{Q} is essentially tridiagonal, since such matrices arose naturally in the problems of finding the zeros of a polynomial. For this case, the algorithm of [4] for solving the matrix equation (3.4) simply reduces to Gaussian elimination which is known to be efficient for such matrix problems.

It is natural to ask about the *rate* of convergence of the $A_1(\mathbf{x}_i)$ to μ . First, it is quite clear that the rate of convergence for the iterative method of (3.4) must be *linear*. In fact, writing $\mathbf{x}_i = \mathbf{y}_\mu + \boldsymbol{\epsilon}_i$, then $A_1(\mathbf{x}_i) = \mu + A_1(\boldsymbol{\epsilon}_i)$, $i \geq 0$, and since $(\tilde{Q} - A_1(\mathbf{x}_i)I_{n-1})\hat{\mathbf{x}}_{i+1} = \hat{\mathbf{a}} = (\tilde{Q} - \mu I_{n-1})\mathbf{y}_\mu$, it follows that

$$(3.4) \quad \frac{A_1(\boldsymbol{\epsilon}_{i+1})}{A_1(\boldsymbol{\epsilon}_i)} = A_1\{(\tilde{Q} - A_1(\mathbf{x}_i)I_{n-1})^{-1}\hat{\mathbf{y}}_\mu\} \sim A_1\{(\tilde{Q} - \mu I_{n-1})^{-2}\hat{\mathbf{a}}\}, \quad i \rightarrow \infty.$$

It is also natural to apply other iterative methods to this problem of finding the smallest real zero, μ , of $f(t) \equiv t - g(t)$. As an example, consider *Newton's method*. Given any real number $t < \sigma$, then previous definitions give us

$$(3.5) \quad (\tilde{Q} - tI_{n-1})\hat{\mathbf{y}}_t = \hat{\mathbf{a}} \quad \text{and} \quad A_1(\mathbf{y}_t) = g(t).$$

Define now the vector \mathbf{w}_t and the scalar $h(t)$ by

$$(3.6) \quad (\tilde{Q} - tI_{n-1})\hat{\mathbf{w}}_t = \hat{\mathbf{y}}_t \quad \text{and} \quad A_1(\mathbf{w}_t) = h(t).$$

Clearly, $\hat{\mathbf{w}}_t = (\tilde{Q} - tI_{n-1})^{-2}\hat{\mathbf{a}}$, and thus from (2.13), we see that $h(t) = g'(t)$. Hence, Newton's method is then defined by

$$(3.7) \quad t_{i+1} \equiv t_i - \frac{t_i - g(t_i)}{1 - h(t_i)}, \quad \text{when } h(t_i) \neq 1.$$

The point here is that the derivative $g'(t)$, which is necessary in the formulation of Newton's method, can be directly calculated by solving an additional system (3.6) of $(n - 1)$ equations in $(n - 1)$ unknowns. Though the rate of convergence of this method is quadratic, convergence of these iterates t_i to μ cannot be guaranteed for all initial values $t_0 < \sigma$, as is clear from Fig. 1. On the other hand, we remark that the iterates t_i of Newton's method (3.7) converge monotonically for the initial value $t_0 = 0$, but the associated vectors \mathbf{y}_i are no longer elements of P_1 .

§4. Second Convergence Theorem

The particular algorithm given in Theorem 1 required the solution of a system of $(n - 1)$ linear equations in $(n - 1)$ unknowns at each iteration. For large n , this algorithm may not be practical unless the matrix \tilde{Q} has some special structure (e.g. tridiagonal). To derive another convergent algorithm, we recall from Lemma 1 that

$$(4.1) \quad \tilde{Q}\hat{\mathbf{x}} \geq \hat{\mathbf{a}} + A_1(\mathbf{x})\hat{\mathbf{x}} \quad \text{for any } \mathbf{x} \in P_1.$$

As in Theorem 1, we assume that we are given an initial $\mathbf{x}^{(0)} \in P_1$ such that strict inequality is valid for at least one component in (4.1). We may, without loss of generality, assume that strict inequality in (4.1) occurs in the second vector

component. This implies, upon rewriting, that

$$(4.2) \quad |a_{1,2}|(x_2^{(0)})^2 + \left\{ \sum_{l>2} |a_{1,l}| x_l^{(0)} - q_{2,2} \right\} (x_2^{(0)}) + \sum_{l \neq 2} |a_{2,l}| x_l^{(0)} < 0,$$

which is a (possibly degenerate) *quadratic* inequality in $x_2^{(0)}$. Since A is irreducible and $\mathbf{x}^{(0)} > \mathbf{0}$, the last term of (4.2) cannot vanish and is thus positive. Hence, if we define $x_2^{(1)}$ as the least positive root of

$$(4.3) \quad |a_{1,2}|(x_2^{(1)})^2 + \left\{ \sum_{l>2} |a_{1,l}| x_l^{(0)} - q_{2,2} \right\} (x_2^{(1)}) + \sum_{l \neq 2} |a_{2,l}| x_l^{(0)} = 0,$$

it is clear that $0 < x_2^{(1)} < x_2^{(0)}$. We now show that this new vector $\mathbf{x}^{(1)}$, obtained by reducing the second component $x_2^{(0)}$ to $x_2^{(1)}$ is still an element of the set P_1 . The remaining inequalities of (4.1) can be written for $j > 2$ as

$$(4.4) \quad \left\{ |a_{1,2}| x_j^{(0)} + |a_{j,2}| x_2^{(0)} \right\} x_j^{(0)} + \left\{ (x_j^{(0)}) \left(\sum_{l>2} |a_{1,l}| x_l^{(0)} \right) + \sum_{\substack{l \neq j \\ l \neq 2}} |a_{j,l}| x_l^{(0)} - q_{j,j} x_j^{(0)} \right\} \leq 0,$$

which are *linear* inequalities in $x_2^{(0)}$ for all $j > 2$. Because the coefficients of $x_2^{(0)}$ in (4.4) are nonnegative for $j > 2$, it is clear that replacing $x_2^{(0)}$ by a smaller quantity leaves all these inequalities unchanged. Thus, we have that $\mathbf{x}^{(1)} \in P_1$. Moreover, since $\mathbf{x}^{(1)} \leq \mathbf{x}^{(0)}$, then $A_1(\mathbf{x}^{(1)}) \leq A_1(\mathbf{x}^{(0)})$.

We now point out that since A is irreducible, at least one coefficient of $x_2^{(0)}$ in (4.4) is positive for some $j > 2$. Thus, while the new vector $\mathbf{x}^{(1)}$ by definition gives equality in (4.3), there is some $j > 2$ for which strict inequality is now valid in (4.4) for this vector. This means that the above process can be *continued*. In fact, there are several natural ways in which the iteration can be continued. One can *cyclically* improve the components x_j in succession, $2 \leq j \leq n$, by solving at each step a quadratic equation similar to (4.3). This corresponds to a *non-linear Gauss-Seidel iterative method* [cf. 1]. One can also improve the components x_j by a *free-steering method* [cf. 6], where one makes the added assumption that each component x_j , $2 \leq j \leq n$, is improved infinitely often. For simplicity, we consider only the cyclic non-linear Gauss-Seidel iterative method, although the basic conclusions are valid also for free-steering methods. Starting with strict inequality in (4.2) for the second component x_2 , let $\mathbf{x}^{(1)}$ now denote the vector obtained after having improved the components x_j in succession, $2 \leq j \leq n$. Using the irreducibility of A , it follows that $A_1(\mathbf{x}^{(1)}) < A_1(\mathbf{x}^{(0)})$, and in general $A_1(\mathbf{x}^{(n+1)}) < A_1(\mathbf{x}^{(n)})$. Moreover, since each $\mathbf{x}^{(n)} \in P_1$, and $\mathbf{x}^{(n+1)} \leq \mathbf{x}^{(n)}$, it is not difficult to verify that $\lim_{n \rightarrow \infty} \mathbf{x}^{(n)} = \mathbf{y}_\mu$ and $\lim_{n \rightarrow \infty} A_1(\mathbf{x}^{(n)}) = \mu$, which gives us

Theorem 2. Let A be an irreducible $n \times n$ matrix which admits a first isolated Gerschgorin disk, and let $\{\mathbf{x}^{(n)}\}_{n=0}^\infty$ be the sequence of vector iterates of the cyclic non-linear Gauss-Seidel iterative method, where $Q\mathbf{x}^{(0)} \geq A_1(\mathbf{x}^{(0)})\mathbf{x}^{(0)}$ with strict inequality for at least one component. Then, all the vectors $\mathbf{x}^{(n)}$ are elements of P_1 with $\lim_{n \rightarrow \infty} \mathbf{x}^{(n)} = \mathbf{y}_\mu$, and $\lim_{n \rightarrow \infty} A_1(\mathbf{x}^{(n)}) = \mu$.

This nonlinear iterative method of Theorem 2 can also be given a simple geometric interpretation. Starting with $\mathbf{x}_0 \in P_1$ such that $\mu < A_1(\mathbf{x}_0) < \sigma$, we now

cyclically treat each component x_j separately, $2 \leq j \leq n$. Evidently, decreasing x_j , $2 \leq j \leq n$, decreases $A_1(\mathbf{x})$, but increases $A_j(\mathbf{x})$. Thus, the iterative method of Theorem 2 in essence selects a new component x_j such that the disks $|z| \leq A_1(\mathbf{x})$ and $|q_{k,k} - z| \leq A_j(\mathbf{x})$ are now tangent. Since the radii $A_k(\mathbf{x})$ of the remaining disks $k \neq j$, are made smaller in the process, we can then apply this procedure to other components.

As a final comment in this section, we point out that the procedure of WILKINSON [12] in essence is the first step of the above Gauss-Seidel procedure for very diagonally dominant matrices Q . He shows in these cases that an excellent approximation to the solution $x_2^{(1)}$ of (4.3) is obtained simply by dropping the first term of (5.3) and solving the resultant linear equation for $x_2^{(1)}$. For such very diagonally dominant matrices, it is also clear that standard linear methods, such as the Gauss-Seidel method, should be rapidly convergent when used in conjunction with the iterative procedure of (3.1). In fact, just one iteration is often sufficient to give good estimates of μ in such cases [12]. In any event, since the vector iterates of these convergent algorithms are elements of the set P_1 , rigorous bounds for the first isolated eigenvalue can be obtained at any step in the iteration.

§ 5. The Set of Matrices $\mathring{\Omega}_A$

In this section, we again assume that the $n \times n$ matrix $A = (a_{i,j})$ is an irreducible matrix which admits a first isolated Gerschgorin disk. We further assume, without loss of generality, that the entry $a_{1,1}$ of A is zero. Otherwise, we could consider the new matrix $A - a_{1,1}I$. In analogy with [10], let $\mathring{\Omega}_A$ be the set of all $n \times n$ matrices $B = (b_{i,j})$ such that $|b_{i,j}| = |a_{i,j}|$ for all $1 \leq i, j \leq n$. Then, the $n \times n$ matrix Q of (2.2), derived from A , is the same for all $B \in \mathring{\Omega}_A$. Hence, each matrix $B \in \mathring{\Omega}_A$ has an eigenvalue in the disk $0 \leq |z| \leq \mu$, where $\mu = \inf_{\mathbf{x} \in P_1} A_1(\mathbf{x})$. Let $s(\mathring{\Omega}_A)$ be the set of all complex numbers z with $0 \leq |z| \leq \mu$ such that z is an eigenvalue of some $B \in \mathring{\Omega}_A$. In other words, $s(\mathring{\Omega}_A)$ is the spectrum of the set $\mathring{\Omega}_A$ restricted to the disk $0 \leq |z| \leq \mu$. Our object now is to determine precisely $s(\mathring{\Omega}_A)$.

Lemma 7. $s(\mathring{\Omega}_A)$ is a closed, bounded, and connected set. Moreover, if $z \in s(\mathring{\Omega}_A)$, then $z \exp i\varphi \in s(\mathring{\Omega}_A)$ for any real φ .

Proof. That $s(\mathring{\Omega}_A)$ is closed and bounded is obvious from the definitions above. To show that $s(\mathring{\Omega}_A)$ is connected, let $\lambda_0, \lambda_1 \in s(\mathring{\Omega}_A)$, where $B_l \mathbf{z}_l = \lambda_l \mathbf{z}_l$, $B_l \in \mathring{\Omega}_A$, $l=0, 1$. Writing $B_l = (|a_{j,k}| \exp i\vartheta_{j,k}^{(l)})$, then the $n \times n$ matrix T_ν , defined by

$$(5.1) \quad T_\nu = (|a_{j,k}| \exp i\{(1-\nu)\vartheta_{j,k}^{(0)} + \nu\vartheta_{j,k}^{(1)}\}), \quad 0 \leq \nu \leq 1,$$

is surely in the set $\mathring{\Omega}_A$, and $T_0 = B_0$, $T_1 = B_1$. Since the entries of T_ν vary continuously with ν , then T_ν possesses an eigenvalue $t(\nu)$ with $|t(\nu)| \leq \mu$ which also varies continuously with ν . Since $t(0) = \lambda_0$ and $t(1) = \lambda_1$, then $s(\mathring{\Omega}_A)$ is connected.

Next, suppose $z \in s(\dot{\Omega}_A)$. Then, there is a $B \in \dot{\Omega}_A$ such that $Bx = zx$, $x \neq 0$. But clearly, $\exp(i\varphi)B$ is also in the set $\dot{\Omega}_A$, and thus $z \exp i\varphi \in s(\dot{\Omega}_A)$, which completes the proof.

With this lemma, we now prove

Theorem 3. There exists a nonnegative real number τ with $\tau < \mu$ such that $s(\dot{\Omega}_A) = \{z \mid \tau \leq |z| \leq \mu\}$.

Proof. First, we note that the $n \times n$ matrix Q of (2.2) is an element of the set $\dot{\Omega}_A$. Thus, from Lemmas 6 and 7, $\mu \exp i\varphi \in s(\dot{\Omega}_A)$ for any real φ . It is now clear that $s(\dot{\Omega}_A) = \{z \mid \tau \leq |z| \leq \mu\}$ for some $\tau \geq 0$, and it remains to show that $\tau < \mu$. Equivalently, we shall show that there is some $z \in s(\dot{\Omega}_A)$ with $|z| < \mu$. Starting with the $n \times n$ matrix Q of (2.2), change the sign of the last diagonal entry of Q , thereby forming the new $n \times n$ matrix Q_1 . Clearly, $Q_1 \in \dot{\Omega}_A$. One can then show that the minimal Gerschgorin set (defined in [10]) for Q_1 , is disconnected, with one connected component lying in the open disk $|z| < \mu$. From this, it follows [10] that Q_1 has one eigenvalue which satisfies $|z| < \mu$, which completes the proof.

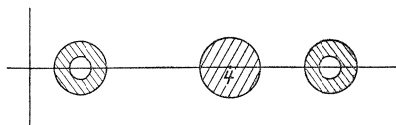


Fig. 2

We remark that the nonnegative quantity τ of Theorem 3 can be positive in some cases and thus $s(\dot{\Omega}_A)$ is a proper annulus. This is illustrated in § 6.

§ 6. Examples

To illustrate the preceding results for a particularly simple example, consider the following matrix

$$(6.1) \quad A = \begin{bmatrix} 1 & i/2 & i/2 \\ 1/2 & 4 & i/2 \\ 1/2 & 1/2 & 6 \end{bmatrix}.$$

It is readily verified from Definition 1 that the vector $\xi = (1, 1, 1)^T$ is simultaneously in the sets P_1, P_2 , and P_3 , so that all the eigenvalues of A can be isolated by positive diagonal similarity transformations. Using the results of the previous sections, it can be shown that there is exactly one eigenvalue of A in each of the following annuli:

$$(6.2) \quad \begin{cases} 0.01584 \leq |z - 1| \leq 0.1608, \\ 0 \leq |z - 4| \leq 0.3139, \\ 0.04593 \leq |z - 6| \leq 0.2301, \end{cases}$$

which is schematically indicated in Fig. 2.

The actual eigenvalues of the matrix A of (6.1) are:

$$\lambda_1 = 0.9897 - 0.1243i; \quad \lambda_2 = 4.0121 - 0.06423i, \quad \text{and} \quad \lambda_3 = 5.9982 + 0.1885i.$$

We remark that the iterative method of Theorem 1 is a rapidly convergent method for finding the exterior radii of the annuli of (6.2), for the initial iterate

unity. For the annulus with center $z=1$, the first four iterates $A_1(x_i)$ of Theorem 1 are 1, 0.2258, 0.1645, and 0.1610. For the cyclic non-linear Gauss-Seidel iterative method of Theorem 2, the corresponding first four iterates $A_1(x_i)$ are 1, 0.2957, 0.1644, and 0.1609.

Finally, we briefly mention that there is a 6×6 example given⁵ in [3], which isolates each of two complex eigenvalues in a disk of radius 0.000287. Because of the increased generality of Theorem 1 of § 3, it is possible to *decrease* this upper bound to 0.0001242. We omit the numerical details.

References

- [1] BERS, L.: On mildly nonlinear partial difference equations of elliptic type. J. Research Nat. Bureau of Standards **51**, 229–236 (1953).
- [2] GANTMACHER, F. R.: The Theory of Matrices, Vol. Two. New York: Chelsea Publishing Co. 1959.
- [3] GERSCHGORIN, S.: Über die Abgrenzung der Eigenwerte einer Matrix. Izv. Akad. Nauk SSSR Ser. Mat. **7**, 749–754 (1931).
- [4] HENRICI, P.: Bounds for eigenvalues of certain tridiagonal matrices. J. Soc. Indust. Appl. Math. **11**, 281–290 (1963).
- [5] OSTROWSKI, A. M.: Über die Determinanten mit überwiegender Hauptdiagonale. Comment. Math. Helv. **10**, 69–96 (1937).
- [6] SCHECHTER, S.: Iteration methods for nonlinear problems. Trans. Amer. Math. Soc. **104**, 179–189 (1963).
- [7] TAUSKY, O.: Bounds for characteristic roots of matrices. Duke Math. J. **15**, 1043–1044 (1948).
- [8] — A method for obtaining bounds for characteristic roots of matrices with application to flutter calculations. Aeronautical Research Council of Great Britain, Report **10**, 508 (1947).
- [9] VARGA, R. S.: Matrix Iterative Analysis. Englewood Cliffs, New Jersey: Prentice-Hall Inc. 1962.
- [10] — Minimal Gerschgorin Sets (to appear in the Pacific J. of Math.).
- [11] WIDDER, D.: The Laplace Transform. Princeton: Princeton University Press 1946.
- [12] WILKINSON, J. H.: Rigorous error bounds for computed eigensystems. The Computer J. **4**, 230–241 (1961).

Computing Center, Case Institute of Technology
10900 Euclid Avenue, Cleveland, Ohio 44106, USA

(Received May 15, 1964)

⁵ The numerical results of [3] are corrected in the J. Soc. Indust. Appl. Math. **12** (1964); pp. 497.