

Componentwise Error Estimates for Approximate Solutions of Nonlinear Equations

TETSURO YAMAMOTO*

Let $x^{(0)}$ be an approximate solution of a system of nonlinear equations. On the basis of the theory of pseudometric space due to Schröder, Collatz and others, a theorem is first proved which determines an existence region of a solution of the system. Next, the result is applied to derive a theorem which may be useful for finding a sharp error bound of each component of $x^{(0)}$. Further, after proving a uniqueness theorem, an algorithm is presented for finding the best uniqueness domain based upon it. Finally, the results are illustrated with a system of two nonlinear equations.

1. Introduction

There is abundant literature concerning the convergence of iterative methods for finding a solution of a nonlinear equation

$$f(x) = 0, \quad (1)$$

defined on a domain D in R^n . Many of the results developed there are applicable for estimating the error of an approximate solution $x^{(0)}$ which was obtained by some method. However, the usual norm estimates may not be suitable in general, when one wants to find a sharp error bound of each or specified component of $x^{(0)}$. In this case, it is desirable to apply the convergence theorems in pseudometric space which were obtained by Schröder [7], [8], [9], Collatz [2] and others.

They are formulated in a certain set $D_0 \subseteq D$ containing $x^{(0)}$ (cf., for example, Theorems 1-3 below). However, no criterion for the choice of D_0 has been given, while the estimates depend, in general, on the set D_0 . In fact, for finding a sharp error bound, the set D_0 should be chosen as small as possible. On the other hand, for finding a large uniqueness region, D_0 should be chosen large.

In this paper, after stating notation and definitions in §2, we first establish, in §3, a theorem giving a sharp existence region of a solution, which is independent of an initial choice of the set D_0 . It slightly improves Schröder's one. Next, it is shown how the result is applied in practical computation to find a sharp error bound of each component of $x^{(0)}$. Further, in §4, we prove a uniqueness theorem which somewhat generalizes a theorem of Kantorovich and Akilov [3] and describe a practical algorithm for finding the best uniqueness domain based upon the theorem. Finally, in §5, the results are illustrated with a system of two nonlinear equations.

2. Notation and Definitions

Let $x = (x_1, \dots, x_n)'$ and $y = (y_1, \dots, y_n)'$ be two column vectors of R^n . We write $x \geq y$ or $y \leq x$ to signify that $x_i \geq y_i$ for all i . Thus $x \geq 0$ means that all the elements of x are nonnegative. We put $v[x] = (|x_1|, \dots, |x_n|)'$ and $\rho(x, y) = v[x - y]$. Then R^n becomes a complete pseudometric space with a pseudodistance $\rho(x, y)$, where the limit of a sequence of vectors are naturally defined (see Collatz [2]). A symbol $x > y$ or $y < x$ means that $x \geq y$ but $x \neq y$. The same notation is used for two matrices $A = (a_{ij})$ and $B = (b_{ij})$ of the same type. For example, $A \geq B$ means that $a_{ij} \geq b_{ij}$ for all i, j , and we put $v[A] = (|a_{ij}|)$, etc.

The spectral radius of an $n \times n$ matrix A is denoted by $\sigma(A)$.

As a vector norm, we adopt the sum norm $\|x\| = |x_1| + \dots + |x_n|$. Observe that, if $x \geq 0$, $y \geq 0$ are two vectors, then

$$\|x + y\| = \|x\| + \|y\|. \quad (2)$$

The corresponding matrix norm is defined by

$$\|A\| = \max_{\|x\|=1} \|Ax\| = \max_j \sum_{i=1}^n |a_{ij}|.$$

If $C = (c_{ijk})$, $i, j, k = 1, 2, \dots, n$ is a bilinear operator (a tensor of the third order), then, analogously, we have

$$\|C\| = \max_{\|x\|=\|y\|=1} \|Cxy\| = \max_{j,k} \sum_{i=1}^n |c_{ijk}|.$$

For a nonnegative vector $d \geq 0$ and a positive number r , we put

$$\bar{U}(x^{(0)}, d) = \{x \in R^n | \rho(x, x^{(0)}) \leq d\},$$

$$S(x^{(0)}, r) = \{x \in R^n | \|x - x^{(0)}\| < r\},$$

and

$$\bar{S}(x^{(0)}, r) = \{x \in R^n | \|x - x^{(0)}\| \leq r\}.$$

$J(x)$ stands for the Jacobian matrix of $f(x) = (f_1(x), \dots, f_n(x))'$, provided that it exists on D .

Finally, we shall call an operator $G: D \subseteq R^n \rightarrow R^n$ a

*Department of Mathematics, Faculty of Science, Ehime University, Matsuyama 790, Japan.

P-contraction on a set $D_0 \subseteq D$ if there exists a matrix P with $P \geq 0$ and $\sigma(P) < 1$ such that, for every $x, y \in D$, $\rho(Gx, Gy) \leq P\rho(x, y)$.

3. Componentwise Error Estimates

Let $x^{(0)} = (x_1^{(0)}, \dots, x_n^{(0)})^t$ be an approximate solution of (1). In this section, we shall prove an existence theorem (Theorem 4) of a solution and gives a practical method (Theorem 5) for finding a sharp error bound of each component $x_i^{(0)} (1 \leq i \leq n)$.

We first refer to two important results which are special cases of the theorems originally obtained by Schröder.

Theorem 1 (Schröder [7], Ortega and Rheinboldt [4] and Collatz [2]). *Let G be a P -contraction on a closed set D_0 and*

$$GD_0 \subseteq D_0. \tag{3}$$

Then G has exactly one fixed point x^ in D_0 . If $x^{(0)} \in D_0$, the iterates $x^{(k+1)} = Gx^{(k)}$, $k = 0, 1, 2, \dots$, belong to D_0 and converge to x^* . The error estimate*

$$\rho(x^{(k)}, x^*) \leq (I - P)^{-1} P^k \rho(x^{(0)}, x^{(1)}), \tag{4}$$

$$k = 0, 1, 2, \dots,$$

holds. Further, if the condition (3) is replaced by

$$\bar{U} \equiv \bar{U}(x^{(1)}, (I - P)^{-1} P \rho(x^{(0)}, x^{(1)})) \subseteq D_0, \tag{5}$$

then $G\bar{U} \subseteq \bar{U}$ and the same conclusions hold in \bar{U} : the sequence $\{x^{(k)}\}$ belongs to \bar{U} and converges to the unique fixed point x^ in \bar{U} , and the estimate (4) holds.*

Theorem 2 (Schröder [8], [9] and Rheinboldt [6]). *Suppose that there exist a convex set $D_0 \subseteq D$ which contains $x^{(0)}$ and a symmetric bilinear operator $B(D_0)$ such that, for every x, y in D_0 ,*

$$v[J(x) - J(y)] \leq B(D_0)\rho(x, y). \tag{6}$$

Let A be a matrix and put $K = v[I - AJ(x^{(0)})]$ and $\varepsilon = v[Af(x^{(0)})]$. If the sequence

$$d^{(0)} = 0, \quad d^{(k+1)} = \frac{1}{2} v[A]B(D_0)d^{(k)2} + Kd^{(k)} + \varepsilon, \tag{7}$$

$$k = 0, 1, 2, \dots,$$

converges to a vector $d^ \geq 0$ and, if $\bar{U}(x^{(0)}, d^*) \subseteq D_0$, then (1) has exactly one solution x^* in $\bar{U}(x^{(0)}, d^*)$.*

On the other hand, the following theorem was proved by Urabe [10].

Theorem 3 *Given an equation (1) where $f(x)$ is continuously differentiable with respect to x in the set D of the x -space. Let $x^{(0)} \in D$ and suppose that there are an $n \times n$ matrix A , an $n \times n$ matrix $P \geq 0$ with the property $\sigma(P) < 1$, and an n -dimensional vector $\delta > 0$ such that*

- (i) $D_\delta = \bar{U}(x^{(0)}, \delta) \subseteq D$,
- (ii) $v[I - AJ(x)] \leq P$ for any $x \in D_\delta$,
- (iii) $(I - P)^{-1} v \leq \delta$,

where v is an n -dimensional vector satisfying $\varepsilon = v[Af(x^{(0)})] \leq v$. The equation (1) then possesses a unique solution x^ in D_δ , and it is valid that $\det[J(x^*)] \neq 0$ and $\rho(x^{(0)}, x^*) \leq (I - P)^{-1} v$.*

Obviously, this theorem follows from Theorem 1. In fact, if we define an operator G by $Gx = x - Af(x)$, then it follows from the condition (ii) of Theorem 3 that G is a P -contraction on D , since G has the first Fréchet derivative $I - AJ(x)$ at every x in D . Further, the conditions (i), (iii) imply (5) with $D_0 = D_\delta$. We thus obtain from Theorem 1

$$\rho(x^{(0)}, x^*) \leq (I - P)^{-1} \rho(x^{(0)}, x^{(1)}) \leq (I - P)^{-1} v.$$

The estimates obtained from these theorems depend on the choice of D_0 or D_δ . However, the theorems tell us nothing about how to choose them. For example, in Theorem 3, we should choose δ so small that the conditions (i), (ii) are satisfied. But, if δ is too small, then the condition (iii) may not hold.

To overcome this difficulty, we now set the following assumptions:

Assumption 1. There is an $n \times n$ matrix A such that the matrix $K = v[I - AJ(x^{(0)})]$ has the spectral radius which is smaller than unity.

Assumption 2. Given any convex set $D_0 \subseteq D$, there exists a symmetric tensor of the third order $H(D_0) = (h_{ijk}(D_0))$, $i, j, k = 1, 2, \dots, n$ such that $h_{ijk}(D_0) \geq 0$ for all i, j, k and

$$\rho(AJ(x), AJ(y)) \leq H(D_0)\rho(x, y), \tag{7}$$

holds for every x, y in D_0 .

Assumption 3. The tensor $H(D_0)$ is monotonically increasing and continuous. That is, if $D_0 \subseteq D'_0$, then $H(D_0) \leq H(D'_0)$ (i.e., $h_{ijk}(D_0) \leq h_{ijk}(D'_0)$ for all i, j, k) and $H(D_0) \rightarrow H(D'_0)$ as $D_0 \rightarrow D'_0$. Observe that Assumption 1 implies the nonsingularity of $J(x^{(0)})$ and that Assumptions 2 and 3 are satisfied if the second Fréchet derivative of $f(x)$ exists and is continuous in D .

Further, we set

$$\varepsilon = v[Af(x^{(0)})], \quad e = (I - K)^{-1} \varepsilon,$$

$$C(d) = (I - K)^{-1} H(\bar{U}(x^{(0)}, d)),$$

$$C(e + \|e\|(1, \dots, 1)^t) = (c_{ijk}), \quad c_i = \max_{j,k} c_{ijk},$$

and

$$c = (c_1, \dots, c_n)^t.$$

Then, under the above notation and assumptions, we first have the following theorem.

Theorem 4 *Let $\bar{U}(x^{(0)}, 2\|e\|(1, \dots, 1)^t) \subseteq D$. If $2\|c\| \cdot \|e\| \leq 1$ or $2\|C(2\|e\|(1, \dots, 1)^t)\| \cdot \|e\| \leq 1$, then the monotonically increasing sequence*

$$\delta^{(0)} = 0, \quad \delta^{(k+1)} = \frac{1}{2} C(\delta^{(k)})\delta^{(k)2} + e, \quad k = 0, 1, 2, \dots, \tag{8}$$

converges to a vector δ^* and the equation (1) has a solution x^* in $\bar{U}(x^{(0)}, \delta^*)$. There is no other solution in $\bar{U}(x^{(0)}, \delta^*)$.

Proof. We first consider the case where $2\|c\| \cdot \|e\| \leq 1$. For every vector $u \geq 0$, we have $C(e + \|e\|(1, \dots, 1))u^2 \leq \|u\|^2 c$. But, as is easily verified, the equation

$$\frac{1}{2} \|u\|^2 c - u + e = 0, \quad u \geq 0,$$

has a nonnegative solution

$$\alpha = e + \frac{\|e\|^2}{1 - \|c\| \cdot \|e\| + \sqrt{1 - 2\|c\| \cdot \|e\|}} c \leq e + \|e\|(1, \dots, 1)^t.$$

Therefore, by Assumption 3 and by induction on k , it is shown that the sequence

$$u^{(0)} = 0, \quad u^{(k+1)} = \frac{1}{2} \|u^{(k)}\|^2 c + e, \quad k = 0, 1, 2, \dots,$$

satisfies

$$x^{(k+1)} - x^{(k)} = x^{(k)} - x^{(k-1)} - AJ(x^{(0)})(x^{(k)} - x^{(k-1)}) - A \int_0^1 \{J(x^{(k-1)} + t(x^{(k)} - x^{(k-1)})) - J(x^{(0)})\} (x^{(k)} - x^{(k-1)}) dt,$$

and, by induction on k ,

$$\begin{aligned} \rho(x^{(k+1)}, x^{(k)}) &\leq K\rho(x^{(k)}, x^{(k-1)}) + \frac{1}{2} H\{\rho(x^{(k)}, x^{(0)}) + \rho(x^{(k-1)}, x^{(0)})\} \rho(x^{(k)}, x^{(k-1)}) \\ &\leq K(v^{(k)} - v^{(k-1)}) + \frac{1}{2} H\{(v^{(k)} - v^{(0)}) + (v^{(k-1)} - v^{(0)})\} (v^{(k)} - v^{(k-1)}) \\ &= K(v^{(k)} - v^{(k-1)}) + \frac{1}{2} (Hv^{(k)2} - Hv^{(k-1)2}) = \psi(v^{(k)}) - \psi(v^{(k-1)}) = v^{(k+1)} - v^{(k)}, \end{aligned}$$

where we have used Assumption 2.

Further, if $v \leq \delta^*$, then $\psi(v) \leq \psi(\delta^*) = \delta^*$ since $C(\delta^*)\delta^{*2} - 2\delta^* + 2e = 0$. This shows that $v^{(k)} \leq \delta^* (k \geq 0)$ and the monotonically increasing sequence $\{v^{(k)}\}$ converges to a vector v^* . It is easy to see that $v^* = \delta^*$. (See Remark 1 below.) We thus conclude from (9) that $\{x^{(k)}\}$ converges to a vector x^* in $\bar{U}(x^{(0)}, \delta^*)$, which is a solution of (1). The uniqueness follows from the proof of Rheinboldt's theorem [6; Theorem 2.6]. In fact, we can show that if y^* is any solution in $\bar{U}(x^{(0)}, \delta^*)$, then

$$\rho(y^*, x^{(k)}) \leq \delta^* - v^{(k)}, \quad k = 0, 1, 2, \dots,$$

which implies that $y^* = x^*$. We should remark here that

$$\bar{U}(x^{(0)}, \delta^*) \subseteq \bar{U}(x^{(0)}, 2\|e\|(1, \dots, 1)^t) \subseteq D,$$

since, by (2),

$$\begin{aligned} \|\delta^*\| &\leq \|\alpha\| = \|e\| + \frac{\|e\|^2}{1 - \|c\| \cdot \|e\| + \sqrt{1 - 2\|c\| \cdot \|e\|}} \|c\| \\ &= \frac{2\|e\|}{1 + \sqrt{1 - 2\|c\| \cdot \|e\|}} \leq 2\|e\|. \end{aligned} \quad (10)$$

Next, we consider the case where $2\|C\|2\|e\|(1, \dots, 1)^t \cdot \|e\| \leq 1$. Let $C = C(2\|e\|(1, \dots, 1)^t)$ and

$$r = \frac{1 - \sqrt{1 - 2\|C\| \cdot \|e\|}}{\|C\|} = \frac{2\|e\|}{1 + \sqrt{1 - 2\|C\| \cdot \|e\|}}.$$

Then the sequence

$$\delta^{(k)} \leq u^{(k)} \leq \frac{1}{2} \|\alpha\|^2 c + e = \alpha.$$

Since the sequence $\{\delta^{(k)}\}$ is monotonically increasing, it converges to a vector $\delta^* \geq 0$. It is a solution of the equation $C(\delta)\delta^2 - 2\delta + 2e = 0$, since, by Assumption 3, $C(\delta^{(k)}) \rightarrow C(\delta^*)$ as $k \rightarrow \infty$.

Now, to prove the existence of a solution in $\bar{U}(x^{(0)}, \delta^*)$, we apply the majorant principle due to Schröder [8], [9] and Rheinboldt [6]. For every vector v , let $H = H(\bar{U}(x^{(0)}, \delta^*))$ and

$$\psi(v) = \frac{1}{2} Hv^2 + Kv + e.$$

Then it is known [8], [9] that a sequence of vectors $v^{(k)} = \psi(v^{(k-1)}) (k \geq 1)$, $v^{(0)} = 0$, majorizes the sequence $x^{(k)} = x^{(k-1)} - Af(x^{(k-1)}) (k \geq 1)$:

$$\rho(x^{(k+1)}, x^{(k)}) \leq v^{(k+1)} - v^{(k)}, \quad k = 0, 1, 2, \dots \quad (9)$$

In fact, we have

$$w^{(0)} = 0, \quad w^{(k+1)} = \frac{1}{2} Cw^{(k)2} + e, \quad k = 0, 1, 2, \dots,$$

remains in $\bar{S}(0, r)$ and converges to a vector w^* , since, by induction on k ,

$$\|w^{(k+1)}\| \leq \frac{1}{2} \|C\| \cdot \|w^{(k)}\|^2 + \|e\| \leq \frac{1}{2} \|C\| r^2 + \|e\| = r.$$

This implies that

$$\delta^{(k)} \leq w^{(k)} \leq w^* \leq r(1, \dots, 1)^t \leq 2\|e\|(1, \dots, 1)^t.$$

Hence, $\{\delta^{(k)}\}$ again converges to δ^* . This means the unique existence of the solution x^* in $\bar{U}(x^{(0)}, \delta^*)$. Q.E.D.

Remark 1. The sequence $\{v^{(k)}\}$ converges to δ^* . For, if $v^* < \delta^*$, then $\|v^*\| < \|\delta^*\|$ and

$$\begin{aligned} \delta^* - v^* &= \psi(\delta^*) - \psi(v^*) \\ &= \frac{1}{2} H(\bar{U}(x^{(0)}, \delta^*))(\delta^{*2} - v^{*2}) + K(\delta^* - v^*), \end{aligned}$$

so that

$$\begin{aligned} \delta^* - v^* &= \frac{1}{2} (I - K)^{-1} H(\bar{U}(x^{(0)}, \delta^*))(\delta^{*2} - v^{*2}) \\ &= \frac{1}{2} C(\delta^*)(\delta^* + v^*)(\delta^* - v^*) \\ &\leq \frac{1}{2} \|\delta^* + v^*\| c(1, \dots, 1)(\delta^* - v^*) \\ &< \|\delta^*\| c(1, \dots, 1)(\delta^* - v^*), \end{aligned}$$

where we have used the facts that $\delta^* \leq e + \|e\|(1, \dots, 1)'$ and

$$C(\delta^*)u \leq C(e + \|e\|(1, \dots, 1)')u \\ \leq \|u\| \cdot \begin{pmatrix} c_1 \cdots c_1 \\ \cdots \\ c_n \cdots c_n \end{pmatrix} = \|u\|c(1, \dots, 1),$$

for every vector $u \geq 0$. Hence, by (10),

$$0 < \|\delta^* - v^*\| < \|\delta^*\| \cdot \|c\| \cdot \|\delta^* - v^*\| \\ \leq 2\|e\| \cdot \|c\| \cdot \|\delta^* - v^*\| \leq \|\delta^* - v^*\|,$$

which is a contradiction and we must have $v^* = \delta^*$.

Remark 2. If $2\|c\| \cdot \|e\| < 1$ or $2\|C\| \cdot \|e\| < 1$ where $C = C(2\|e\|(1, \dots, 1)')$, then the matrix $J(x^*)$ is nonsingular, that is, the solution x^* is simple. To prove this, let $L = I - AJ(x^{(0)})$. Then $\sigma(L) \leq \sigma(K) < 1$ so that $J(x^{(0)})$ is nonsingular and $J(x^{(0)})^{-1} = (I - L)^{-1}A$. Therefore, putting $M = I - J(x^{(0)})^{-1}J(x^*)$, we have

$$v[M] = v[(I - L)^{-1}(AJ(x^{(0)}) - AJ(x^*))] \\ \leq v[I + L + L^2 + \dots]H(\bar{U}(x^{(0)}, \delta^*))\rho(x^{(0)}, x^*) \\ \leq (I + K + K^2 + \dots)H(\bar{U}(x^{(0)}, \delta^*))\delta^* \\ = C(\delta^*)\delta^* \leq \|\delta^*\|c(1, \dots, 1),$$

so that

$$\|M\| = \|v[M]\| \leq \|\delta^*\| \cdot \|c\| \leq 2\|e\| \cdot \|c\|.$$

Similarly we have $v[M] \leq C\delta^*$ and $\|M\| \leq \|C\| \cdot \|\delta^*\| \leq 2\|C\| \cdot \|e\|$. Consequently, if $2\|c\| \cdot \|e\| < 1$ or $2\|C\| \cdot \|e\| < 1$, then $\|M\| < 1$ and $J(x^*) = J(x^{(0)})(I - M)$ is nonsingular. But, if $2\|c\| \cdot \|e\| = 1$ or $2\|C\| \cdot \|e\| = 1$, then $J(x^*)$ may be singular. For example, let $n = 1, f(x) = x^2, x^{(0)} = a > 0$ and $A = (2x^{(0)})^{-1}$. Then a simple computation yields $c = C = H = a^{-1}, 2e = 2e = \delta^* = a$ and $2\|c\| \cdot \|e\| = 1$. Thus Theorem 4 asserts that there exists a unique solution x^* of (1) in the closed interval $[0, 2a]$. But, $x^* = 0$ and $J(x^*) = f'(0) = 0$.

Remark 3. It would be interesting to compare Theorem 4 with Theorems 2 and 3. If the convex set D_0 in Theorem 2 includes the region $\bar{U}(x^{(0)}, \delta^*)$, then, obviously, we have $\delta^* \leq d^*$. Next, as the region D_δ in Theorem 3, we take $D^* = \bar{U}(x^{(0)}, \delta^*)$, a sharp region whose a priori choice can not be expected in a practical computation. Then, for any x in $D_\delta = D^*$, we have

$$v[I - AJ(x)] = v[I - AJ(x^{(0)}) - A(J(x) - J(x^{(0)}))] \\ \leq K + H(D^*)\delta^*,$$

provided that (7) holds. (We may take $H(D^*) = v[A]B(D^*)$ if (6) holds.) So we can take $P = K + H(D^*)\delta^*$. Then we have

$$\varepsilon + P\delta^* - \delta^* = \frac{1}{2}H(D^*)\delta^* > 0,$$

which means that $(I - P)^{-1}\varepsilon > \delta^*$ if $\sigma(P) < 1$. Therefore, in order that Theorem 3 gives a sharper estimate than Theorem 4, δ should be chosen so small that $\delta < \delta^*$.

Remark 4. If the equation (1) is defined on a domain

D in C^n and $x^{(0)} \in C^n$, then our theorem holds by putting

$$\bar{U}(x^{(0)}, d) = \{x \in C^n | \rho(x, x^{(0)}) \leq d\}, \text{ etc.}$$

Corollary. Let $f(z) = 0$ be a single equation in C^1 or R^1 and apply the Newton method $z_{k+1} = z_k - f(z_k)/f'(z_k), k = 0, 1, \dots$. Set

$$\varepsilon_{k+1} = |f(z_{k+1})/f'(z_{k+1})|, \quad D_{k+1} = \{z | |z - z_{k+1}| \leq 2\varepsilon_{k+1}\}, \\ M_{k+1} = \sup_{z \in D_{k-1}} |f''(z)| \quad \text{and} \quad h_{k+1} = M_{k+1}/|f'(z_{k+1})|,$$

provided that z_{k+1} is defined, $f'(z_{k+1}) \neq 0$ and $f''(z)$ exists in D_{k+1} . If $2h_{k+1}\varepsilon_{k+1} \leq 1$, then a solution z^* exists and is unique in

$$D_{k+1}^* = \left\{ z \in C^1 | |z_{k+1} - z| \leq \frac{2\varepsilon_{k+1}}{1 + \sqrt{1 - 2h_{k+1}\varepsilon_{k+1}}} \right\}. \quad (11)$$

Remark 5. Let

$$\hat{D}_{k-1} = \{z | |z - z_k| \leq \varepsilon_{k-1}\} \quad \text{and} \quad \hat{M}_{k-1} = \sup_{z \in \hat{D}_{k-1}} |f''(z)|.$$

Then, Ostrowski's theorem [5; Theorem 7.2] asserts that, if $2\varepsilon_{k-1}\hat{M}_{k-1} \leq |f'(z_{k-1})|$, then

$$|z_{k+1} - z^*| \leq \frac{\hat{M}_{k-1}}{2|f'(z_k)|} \varepsilon_{k-1}^2, \quad (12) \\ 2\varepsilon_k |f'(z_k)| \leq \hat{M}_{k-1} \varepsilon_{k-1}^2, \\ \varepsilon_{k+1} \leq \frac{1}{2} \varepsilon_k \leq \frac{1}{4} \varepsilon_{k-1}, \quad \text{and} \quad 2\varepsilon_k M_{k-1} \leq |f'(z_k)|.$$

Obviously we have $D_{k+1} \subseteq \hat{D}_{k-1}$ and $M_{k+1} \leq \hat{M}_{k-1}$. Hence

$$\frac{2\varepsilon_{k+1}}{1 + \sqrt{1 - 2h_{k+1}\varepsilon_{k+1}}} \leq 2\varepsilon_{k+1} \leq \varepsilon_k \leq \frac{\hat{M}_{k-1}}{2|f'(z_k)|} \varepsilon_{k-1}^2.$$

That is, the estimate (11) is sharper than (12).

Under Assumptions 1-3, we now obtain the following theorem which may be useful when one wants to find a sharp error bound of each or specified component of $x^{(0)}$.

Theorem 5 Let $\bar{U}(x^{(0)}, 2\|e\|(1, \dots, 1)') \subseteq D$.

(i) If $2\|c\| \cdot \|e\| \leq 1$, then there exists a solution of (1) in $\bar{U}(x^{(0)}, \alpha)$, where

$$\alpha = e + \frac{\|e\|^2}{1 - \|c\|\|e\| + \sqrt{1 - 2\|c\| \cdot \|e\|}} c.$$

(ii) Let $\{\delta^{(k)}\}$ be the iterated sequence of vectors defined in (8) and set $\xi^{(k)} = 2(\delta^{(k+1)} - \delta^{(k)})$ and $\eta^{(k)} = \delta^{(k)} + \xi^{(k)}$, provided that $2\|c\| \cdot \|e\| \leq 1$ or $2\|C(2\|e\|(1, \dots, 1)')\| \cdot \|e\| \leq 1$. If the conditions

$$\eta^{(k)} \leq 2\|e\|(1, \dots, 1)',$$

and

$$\{C(\eta^{(k)}) - C(\delta^{(k)})\}\delta^{(k)2} + 2C(\eta^{(k)})\delta^{(k+1)}\xi^{(k)} \leq \xi^{(k)},$$

are satisfied for some $k \geq 0$, then there exists a solution of (1) in $\bar{U}(x^{(0)}, \eta^{(k)})$.

Proof. It remains to prove (ii). Assume now that $\delta^{(l)} \leq \eta^{(k)}$ for some $l > k$. Then we have

$$\begin{aligned} \delta^{(l+1)} &= \frac{1}{2} C(\delta^{(l)}) \delta^{(l)2} + e \leq \frac{1}{2} C(\eta^{(k)}) \eta^{(k)2} + e \\ &= \frac{1}{2} C(\eta^{(k)}) (\delta^{(k)} + \xi^{(k)})^2 + e \\ &= \frac{1}{2} \{C(\eta^{(k)}) - C(\delta^{(k)})\} \delta^{(k)2} + C(\eta^{(k)}) \left(\delta^{(k)} + \frac{1}{2} \xi^{(k)} \right) \xi^{(k)} \\ &\quad + \frac{1}{2} C(\delta^{(k)}) \delta^{(k)2} + e \\ &= \frac{1}{2} \{C(\eta^{(k)}) - C(\delta^{(k)})\} \delta^{(k)2} + C(\eta^{(k)}) \delta^{(k+1)} \xi^{(k)} + \delta^{(k+1)} \\ &\leq \frac{1}{2} \xi^{(k)} + \delta^{(k+1)} = \eta^{(k)}. \end{aligned}$$

Therefore, by induction on l , we have $\delta^{(l)} \leq \eta^{(k)}$ for every $l > k$. Letting $l \rightarrow \infty$, we now obtain $\delta^* \leq \eta^{(k)} \leq 2\|e\|(1, \dots, 1)^t$, which means the existence of a solution in $\bar{U}(x^{(0)}, \eta^{(k)})$. Q.E.D.

4. Determination of a Uniqueness Domain

Given an approximate solution $x^{(0)}$ of (1), it is also of importance to find a uniqueness domain as large as possible. In this section, taking account that the parallelopete $\bar{U}(x^{(0)}, d)$ with a vector $d \geq 0$ is included in the corresponding ball $\bar{S}(x^{(0)}, \|d\|)$, we shall present a method for finding a uniqueness domain $S(x^{(0)}, r)$ as large as possible. We keep the assumptions of §3.

We first prove the following theorem which slightly generalizes a theorem of Kantorovich and Akilov [3], [6; Corollary 3.3].

Theorem 6. Let $D_0 = \bar{S}(x^{(0)}, r) \subseteq D$ for some $r > 0$ and set $\hat{C}(r) = (I - K)^{-1} H(D_0)$. If $2\|\hat{C}(r)\| \cdot \|e\| < 1$, then (1) has at most one solution in $S(x^{(0)}, r^{**}) \cap D_0$, where

$$r^{**} = \frac{1 + \sqrt{1 - 2\|\hat{C}(r)\| \cdot \|e\|}}{\|\hat{C}(r)\|}.$$

Proof. For simplicity, let $H = H(D_0)$ and $Gx = x - Af(x)$. Then, for every x, y in D_0 ,

$$\begin{aligned} Gx - Gy &= x - y - AJ(x^{(0)})(x - y) \\ &\quad - A \int_0^1 \{J(y + t(x - y)) - J(x^{(0)})\} (x - y) dt, \end{aligned}$$

and, by (7),

$$\rho(Gx, Gy) \leq K\rho(x, y) + \frac{1}{2} H\{\rho(x, x^{(0)}) + \rho(y, x^{(0)})\} \rho(x, y). \tag{13}$$

Therefore, if x^* is a solution in $S(x^{(0)}, r^{**}) \cap D_0$, then

$$\begin{aligned} \rho(x^{(0)}, x^*) &\leq \rho(x^{(0)}, Gx^{(0)}) + \rho(Gx^{(0)}, Gx^*) \\ &\leq \varepsilon + K\rho(x^{(0)}, x^*) + \frac{1}{2} H\rho(x^{(0)}, x^*)^2, \end{aligned}$$

so that

$$\begin{aligned} \rho(x^{(0)}, x^*) &\leq (I - K)^{-1} \varepsilon + \frac{1}{2} (I - K)^{-1} H\rho(x^{(0)}, x^*)^2 \\ &= e + \frac{1}{2} \hat{C}(r) \rho(x^{(0)}, x^*)^2. \end{aligned}$$

In particular, we have

$$\|x^{(0)} - x^*\| \leq \|e\| + \frac{1}{2} \|\hat{C}(r)\| \cdot \|x^{(0)} - x^*\|^2.$$

This gives us

$$\|x^{(0)} - x^*\| \leq r^* = \frac{1 - \sqrt{1 - 2\|\hat{C}(r)\| \cdot \|e\|}}{\|\hat{C}(r)\|},$$

since, by assumption, $\|x^{(0)} - x^*\| < r^*$.

Now, to prove the uniqueness, let \hat{x} be another solution in $S(x^{(0)}, r^{**}) \cap D_0$. Then, it follows from (13) that $\rho(\hat{x}, x^*) = \rho(G\hat{x}, Gx^*)$

$$\leq K\rho(\hat{x}, x^*) + \frac{1}{2} H\{\rho(\hat{x}, x^{(0)}) + \rho(x^*, x^{(0)})\} \rho(\hat{x}, x^*),$$

or

$$\rho(\hat{x}, x^*) \leq \frac{1}{2} (I - K)^{-1} H\{\rho(\hat{x}, x^{(0)}) + \rho(x^*, x^{(0)})\} \rho(\hat{x}, x^*).$$

Since $\|\hat{x} - x^{(0)}\| \leq r^*$, we thus obtain

$$\begin{aligned} \|\hat{x} - x^*\| &\leq \|\hat{C}(r)\| r^* \|\hat{x} - x^*\| \\ &= \{1 - \sqrt{1 - 2\|\hat{C}(r)\| \cdot \|e\|}\} \|\hat{x} - x^*\|, \end{aligned}$$

which means $\hat{x} = x^*$ since $2\|\hat{C}(r)\| \cdot \|e\| < 1$. Consequently, (1) has at most one solution in $S(x^{(0)}, r^{**}) \cap D_0$. Q.E.D.

Corollary. Let $\bar{U}(x^{(0)}, 2\|e\|(1, \dots, 1)^t) \subseteq D$ and $2\|C(2\|e\|(1, \dots, 1)^t)\| \cdot \|e\| < 1$. Then the solution of (1) (which exists in $\bar{U}(x^{(0)}, \delta^*)$ where δ^* is the limit of the sequence (8)) is unique in $\bar{S}(x^{(0)}, 2\|e\|)$.

Proof. Take $r = 2\|e\|$ in Theorem 6, by noting that $\|\hat{C}(r)\| \leq \|C(r(1, \dots, 1)^t)\|$ and $r^{**} > 2\|e\|$. (The solution of (1) does exist in $\bar{U}(x^{(0)}, \delta^*)$ by Theorem 4.) Q.E.D.

The best uniqueness domain based upon Theorem 6 is obtained, if the radius r of the closed ball D_0 is chosen so that

$$r = \frac{1 + \sqrt{1 - 2\|\hat{C}(r)\| \cdot \|e\|}}{\|\hat{C}(r)\|}.$$

Such an r can be found by the following procedures, provided that the domain D is sufficiently large.

Procedure 1. Let $2p_0 = 2\|\hat{C}(2\|e\|)\| \cdot \|e\| < 1$ and set $r_0 = 2\|e\|$ and

$$s_0 = \frac{1 + \sqrt{1 - 2p_0}}{\|\hat{C}(2\|e\|)\|}.$$

Procedure 2. Assume that $S(x^{(0)}, s_0) \subseteq D$ and define the bounded sequences $\{r_i\}, \{s_i\} (i = 1, 2, \dots)$ as follows:

If $2p_i = 2\|\hat{C}(s_{i-1})\| \cdot \|e\| < 1$, then set

$$\begin{aligned} \omega_i &= \frac{1 + \sqrt{1 - 2p_i}}{\|\hat{C}(s_{i-1})\|}, \\ r_i &= \max \{r_{i-1}, \min \{s_{i-1}, \omega_i\}\}, \\ s_i &= \frac{1}{2} \{r_i + \max \{s_{i-1}, \omega_i\}\}. \end{aligned}$$

If $2p_i \geq 1$, then set $r_i = r_{i-1}$ and $s_i = \frac{1}{2}(r_{i-1} + s_{i-1})$.

Then, the solution is unique in $S(x^{(0)}, r_i)$ and the sequence $\{r_i\}$ is monotonically increasing and converges to a positive number r , the desired one. Note that $r_i \leq s_i \leq s_0$ for each i .

For automatic computation, however, it may be convenient to replace $\hat{C}(s_i)$ by $C(s_i(1, \dots, 1)^t)$, although the result obtained becomes rather rough.

5. A Numerical Example

To illustrate our results, we take up the system

$$\begin{aligned} f_1 &= 3x_1^2x_2 + x_2^3 - 1 = 0, \\ f_2 &= x_1^4 + x_1x_2^3 - 1 = 0, \end{aligned} \tag{14}$$

which was considered by Kantorovich and Akilov [3]. By using the Newton method starting from an initial value (0.98, 0.32), they obtained an approximate solution $x^{(0)t} = (0.991189, 0.327382)$ and concluded that there is a solution x^* such that

$$\begin{aligned} 0.991173 \leq x_1^* \leq 0.991205, \\ 0.327366 \leq x_2^* \leq 0.327398. \end{aligned} \tag{15}$$

We shall apply Theorem 5 to (14) and $x^{(0)}$, by choosing $A = J(x^{(0)})^{-1}$. The second Fréchet derivative of $f(x)$ at x is

$$\begin{aligned} J'(x) &= \begin{pmatrix} f_{111} & f_{112} & f_{121} & f_{122} \\ f_{211} & f_{212} & f_{221} & f_{222} \end{pmatrix} \\ &= \begin{pmatrix} 6x_2 & 6x_1 & 6x_1 & 6x_2 \\ 12x_1^2 & 3x_2^2 & 3x_2^2 & 6x_1x_2 \end{pmatrix}, \end{aligned}$$

where

$$f_{ijk} = \frac{\partial^2 f_i(x)}{\partial x_j \partial x_k}.$$

Therefore, $B = B(\bar{U}(x^{(0)}, d))$ defined in Theorem 2 is obtained if we replace each $x_i^{(0)}$ in the expression of $J'(x^{(0)})$ by $|x_i^{(0)}| + d_i$. Set $H(\bar{U}(x^{(0)}, d)) = v[A]B$. Then, it is easily verified that $2\|c\| \cdot \|e\| = 0.21 \dots E - 4 < 1$ and

$$\alpha^t = (0.5215503 \dots E - 6, 0.1331679 \dots E - 5).$$

Thus, by (i) of Theorem 5, we can assert that the solution x^* does exist in $\bar{U}(x^{(0)}, \alpha)$ and

$$\begin{aligned} 0.9911884 \dots \leq x_1^* \leq 0.9911895 \dots, \\ 0.3273806 \dots \leq x_2^* \leq 0.3273833 \dots \end{aligned} \tag{16}$$

This estimate is sharper than (15). Further, if we execute the iteration (8) with a stopping criterion $\delta^{(k+1)} - \delta^{(k)} \leq$

$10^{-13}(1, \dots, 1)^t$, then this condition is satisfied at $k=2$, as well as the conditions in (ii) of Theorem 5, and we obtain

$$\eta^{(2)t} = (0.5215459 \dots E - 6, 0.1331677 \dots E - 5),$$

which slightly improves the estimate (16). Observe that, if $\|c\| \cdot \|e\|$ is sufficiently small, then

$$e \leq \delta^* \leq \alpha \leq e + \|e\|^2 c = e.$$

Therefore, in this case, there is little difference between α and $\eta^{(k)}$. The result of the modified procedure described at the end of the previous section is also shown in Table 1, which shows that the solution is unique in $S(x^{(0)}, r_i)$ for each i .

The computation in this example was done on a FACOM 230-28 computer of Ehime University, with double precision arithmetic chopping to hexadecimal 14-digit numbers. The values of r_i and s_i in Table 1 are respectively chopped and rounded to six digits.

Table 1 The radii of the uniqueness domain $S(x^{(0)}, r_i)$.

i	r_i	s_i
0	0.370643E-05	0.350150E+00
1	0.188274E+00	0.269212E+00
2	0.213309E+00	0.241261E+00
3	0.223187E+00	0.232224E+00
4	0.226535E+00	0.229380E+00
5	0.227606E+00	0.228493E+00

Acknowledgement

The author wishes to express his thanks to Professor Lothar Collatz of Hamburg University for his many helpful comments and to Dr. M. T. Noda of Ehime University for discussions during the preparation of this paper. Thanks also go to the referees for their constructive comments.

References

1. BOHL, E., MONOTONIE, Lösbarkeit und Numerik bei Operatorgleichungen. New York: Springer (1974).
2. COLLATZ, L. Functional analysis and numerical mathematics. New York: Academic Press (1966).
3. KANTOROVICH, L. V. and AKILOV, G. P. Functional analysis in normed spaces. Oxford: Pergamon Press (1964).
4. ORTEGA, J. M. and RHEINBOLDT, W. C. Iterative solution of nonlinear equations in several variables. New York: Academic Press (1970).
5. OSTROWSKI, A. M. Solution of equations in Euclidean and Banach spaces. New York: Academic Press (1973).
6. RHEINBOLDT, W. C. A unified convergence theory for a class of iterative processes. SIAM J. Numer. Anal. 5 (1968) 42-63.
7. SCHRÖDER, J. Das Iterationsverfahren bei allgemeinerem Abstandsbegriff. Math. Z. 66 (1956) 111-116.
8. SCHRÖDER, J. Nichtlineare Majoranten beim Verfahren der schrittweisen Näherung. Arch. Math. 7 (1956) 471-484.
9. SCHRÖDER, J. Über das Newtonsche Verfahren. Arch. Rat. Mech. Anal. 1 (1957) 154-180.
10. URABE, M. A posteriori component-wise error estimation of approximate solutions to nonlinear equations. Lect. Notes in Computer Sci. 29 (Springer) (1975) 99-117.

(Received July 10, 1978: revised June 20, 1979)