

SIMPLE BOUNDS FOR ZEROS OF SYSTEMS OF EQUATIONS

Arnold Neumaier

Institut für Angewandte Mathematik

Universität Freiburg

D-7800 Freiburg, W-Germany

**Abstract.** Let  $F: D(\subseteq \mathbb{R}^n) \rightarrow \mathbb{R}^n$  be a continuous function, and suppose that for some  $x_0 \in D$  and some nonsingular matrix  $A$  the vector  $\delta_0 := A^{-1} F(x_0)$  is "small". Assuming the existence of a nonnegative vector  $c \in \mathbb{R}^n$  such that

$$|F(x) - F(x_0) - A(x - x_0)| \leq \|\delta_0\| c$$

for all  $x$  in a suitable neighbourhood  $S$  of  $x_0$ , a simple condition is given which guarantees that  $S$  contains a zero  $\lambda$  of  $F$ . The resulting bounds are shown to be quite accurate. They contain as a special case the bounds obtainable from a theorem of Kantorovic. It is discussed how to compute the required vector  $c$ , and how to make efficient use of sparsity. The practical use of the bounds is demonstrated by an extensive example, the finite difference equations obtained by discretizing the minimal surface equation.

1. Introduction

Let  $D$  be a subset of  $\mathbb{R}^n$ , and let  $F: D \rightarrow \mathbb{R}^n$  be a bounded continuous function. We assume the following:

(I) A good approximation  $x_0 \in D$  for an isolated zero  $\lambda$  of  $F(x)$  is known.

(II) An approximation  $A$  for the Fréchet-derivative or a generalized slope operator near  $x_0$  is known.

(III)  $A$  is nonsingular, and the structure of  $A$  allows the practical solution of linear equations with coefficient matrix  $A$  by Gauss elimination.

With these assumptions, we ask for a bound on the magnitude of the error  $\lambda - x_0$ . The bound shall be such that we can take advantage of any special structure of  $A$  such as diagonal dominance or sparseness.

Assuming (I), (II), and (III), we have for  $x$  near  $x_0$

$$F(x) - F(x_0) \approx A(x - x_0), \tag{1}$$

hence, for  $x = \lambda$  with  $F(\lambda) = 0$ ,

$$\lambda - x_0 \approx -A^{-1} F(x_0).$$

Thus, if we put

$$\delta_0 := A^{-1} F(x_0), \tag{2}$$

we have in any norm  $\|\lambda - x_0\| \approx \|\delta_0\|$ , and we can expect an error estimate

$$\|\lambda - x_0\| \leq \kappa \|\delta_0\|$$

with a reasonably small factor  $\kappa > 1$ . To guarantee  $\lambda \in D$  we assume that

$$S := \{x \in \mathbb{R}^n \mid \|x - x_0\| \leq \kappa \|\delta_0\|\} \subseteq D, \tag{3}$$

and to make (1) precise we require that

$$|F(x) - F(x_0) - A(x - x_0)| \leq \|\delta_0\| c \quad \text{for all } x \in S. \tag{4}$$

Hence  $c$  is a nonnegative vector; absolute value and order is understood to be componentwise.

It turns out that (2), (3), and (4), together with a simple inequality are sufficient to show that  $S$  must contain at least one zero of  $F(x)$ . The precise statement and its proof is contained in Section 2. Section 3 discusses the practical computation of the quantities involved in the bound. In Section 4 we compare our bounds with those of the Kantorovic theorem, and in Section 5 we present some numerical examples.

The concepts used in this paper can be found, e.g. in the book by Ortega and Rheinboldt [12].

## 2. The main result

For the moment we shall forget about the heuristic introduction, and suppose that  $D$  is a subset of  $\mathbb{R}^n$ ,  $F: D \rightarrow \mathbb{R}^n$  is a continuous function,  $x_0 \in D$  is a vector,  $A$  is a nonsingular  $n \times n$ -matrix, and

$$\delta_0 := A^{-1} F(x_0). \quad (2)$$

We suppose further that there is a constant  $\kappa > 1$  such that

$$S := \{x \in \mathbb{R}^n \mid \|x - x_0\| \leq \kappa \|\delta_0\|\} \subseteq D, \quad (3)$$

and a nonnegative vector  $c \in \mathbb{R}^n$  such that, for some monotone norm  $\|\cdot\|$

$$\|F(x) - F(x_0) - A(x - x_0)\| \leq \|\delta_0\| c \quad \text{for all } x \in S. \quad (4)$$

Then the following theorem holds.

### Theorem

If the vector

$$b := \|A^{-1}\|c \quad (5)$$

satisfies the condition

$$\|b\| \leq \kappa - 1 \quad (6)$$

then  $F(x)$  has at least one zero in  $S$ , and any such zero  $\lambda$  satisfies

$$(2 - \kappa) \|\delta_0\| \leq \|\lambda - x_0\| \leq \kappa \|\delta_0\|. \quad (7)$$

Proof. We show that the map  $\phi: D \rightarrow \mathbb{R}^n$  defined by

$$\phi(x) := x - A^{-1} F(x)$$

maps  $S$  into itself. Indeed, if  $x \in S$  then

$$\begin{aligned} |\phi(x) - x_0| &= |x - x_0 - A^{-1} F(x)| \\ &= |-A^{-1}(F(x) - F(x_0)) - A(x - x_0)| - A^{-1} F(x_0)| \\ &\leq |A^{-1}| \cdot \|\delta_0\| c + |A^{-1} F(x_0)| \\ &= \|\delta_0\| b + |\delta_0| \end{aligned}$$

by (4), (5), and (2), hence

$$\|\phi(x) - x_0\| \leq \|\delta_0\| \|b\| + \|\delta_0\| \leq \kappa \|\delta_0\|$$

since  $\|\cdot\|$  is monotone and by (6). Hence  $\phi(x) \in S$ . By Brouwer's fixed point theorem,  $S$  contains a fixed point  $\lambda$  of  $\phi$ , and clearly  $\lambda$  satisfies  $F(\lambda) = 0$ .

Now for any zero  $\lambda \in S$  of  $F(x)$ , the right hand inequality of (7) is obvious. To prove the left hand inequality, we insert  $x = \lambda$ ,  $F(\lambda) = 0$  into (4) and obtain

$$\|F(x_0) + A(\lambda - x_0)\| \leq \|\delta_0\| c,$$

hence

$$\begin{aligned} |\lambda - x_0| &= |\delta_0 - A^{-1}(F(x_0) + A(\lambda - x_0))| \\ &\geq |\delta_0| - |A^{-1}| \|F(x_0) + A(\lambda - x_0)\| \\ &\geq |\delta_0| - |A^{-1}| \cdot \|\delta_0\| c \\ &= |\delta_0| - \|\delta_0\| b, \end{aligned}$$

and therefore

$$\|\lambda - x_0\| \geq \|\delta_0\| - \|\delta_0\| \|b\| \geq (2 - \kappa) \|\delta_0\|. \quad (8)$$

This completes the proof.  $\square$

To analyze the quality of the bound and the influence of the various quantities involved we adopt again the heuristic attitude of Section 1. Two important questions arise immediately:

- (1) How restrictive is the condition (6) needed for the validity of (7)?
- (11) How good are the bounds (7)?

Obviously, (6) is satisfied if either  $c$  and hence  $b$  is sufficiently small, or if  $\kappa$  is sufficiently large.

Suppose first that  $\kappa$  is fixed. How small may we choose  $c$ ? To answer this, we assume that  $F(x)$  has a Fréchet-derivative  $F'(x_0)$  in  $x_0$ , and  $A = F'(x_0)$ . Then  $|F(x) - F(x_0) - A(x - x_0)| = o(\|x - x_0\|) = o(\|\delta_0\|)$  for  $x \in S$ , hence we may choose  $c$  arbitrarily small provided that  $\|\delta_0\|$  is small enough. If the second derivative exists, too, then we even have  $|F(x) - F(x_0) - A(x - x_0)| = o(\|\delta_0\|^2)$ , hence  $c = o(\|\delta_0\|)$ . If  $A$  is close to  $F'(x_0)$  then  $c$  has to be increased by a correction of order  $\|A - F'(x_0)\|$ . Thus we can satisfy (6) provided that  $\delta_0$  is sufficiently small (i.e.  $x_0$  is sufficiently close to a zero  $\lambda$ ) and  $A$  is sufficiently close to  $F'(x_0)$ .

Now we suppose that  $c$  is kept fixed. If  $F(x)$  is differentiable then  $|F(x) - F(x_0) - A(x - x_0)|$  will grow at least like a multiple of  $\|x - x_0\|$ , hence (4) implies that  $\|x - x_0\|$  may not increase indefinitely. Therefore, the bounding factor  $\kappa$  in (3) must be kept sufficiently small.

Thus, if we increase  $\kappa$ , the vector  $c$  will increase as well. Unless  $F(x)$  is linear or almost linear,  $c$  will increase much faster than  $\kappa$  (for example  $c = O(\kappa^2)$  if the second derivative exists). Therefore,  $\kappa$  has to be kept fairly small.

Of course there must be cases when no choice of  $c$  and  $\kappa$  will make (6) valid since  $F(x)$  may be nonzero in every set  $S$  of shape (3).

To estimate the accuracy of the bound (7) we compare the lower and upper bound. If  $\kappa < 2$ , then (7) shows that the upper bound overestimates the true error by a factor of at most

$$\mu = \frac{\kappa}{2-\kappa} \quad \}$$

Thus for  $\kappa = 3/2$  the true error is caught within a factor of  $\mu = 3$ . If  $\kappa > 2$ , an overestimation factor of

$$\mu' = \frac{1}{1-\|\|b\|}$$

can be derived from (8) provided that  $\|b\| < 1$ . But in this case (6) is satisfied for some  $\kappa < 2$  as well. Hence from the point of view of accurate estimation, there is no incentive to use a factor  $\kappa > 2$ . On the other hand, there are cases when (6) can be satisfied with some  $\kappa > 2$  but not with  $\kappa \leq 2$  (see Section 4).

#### Remark

The argument leading to (8) can be adapted to give the component-wise error bound

$$\|x - x_1\| \leq \|\delta_0\| b \quad \text{for } x_1 := x_0 - \delta_0. \quad (9)$$

Thus, if the entries of  $b$  differ much in magnitude, it can be expected that some components of  $x_1$  are much more accurate than others. (This may be an indication that the problem is badly scaled).

#### 3. Practical aspects

In this section we discuss the choice of the parameters and the practical determination of the bound. For the sake of a simple presentation we assume that  $F(x)$  is (Fréchet-) differentiable.

3.1. Choice of  $x_0$  and  $A$ . For good bounds,  $x_0$  should be an approximation of a zero of  $F(x)$ . From (6) we see that this essentially means that  $\delta_0$  is kept small. Thus it is sensible to determine  $x_0$  from a quasi-Newton iteration.

$$z_0 \text{ suitable, } z_{i+1} := z_i - \eta_i A_i^{-1} F(z_i), \quad A_i \approx F'(z_i). \quad (10)$$

If  $A_m^{-1} F(z_m)$  is sufficiently small we put

$$x_0 := z_m, \quad A := A_m, \quad \delta_0 := A_m^{-1} F(z_m). \quad (11)$$

But it should be kept in mind that the bound (7) does not depend on this choice of  $x_0$  and  $A$ . For example,  $x_0$  might be the solution of a previously solved near by system  $\tilde{F}(x) = 0$ , or (in particular if  $F'(x_0)$  is an  $M$ -matrix),  $A$  might be the lower triangle of  $F'(x_0)$ .

3.2. Choice of  $\kappa$ . We already commented on the need of keeping  $\kappa$  small, preferably  $\kappa < 2$ . If  $x_0$  is indeed close to a zero of  $F(x)$ , and  $A = F'(x_0)$  then  $\kappa = 3/2$  works quite well since  $c = O(\|\delta_0\|)$ . But if  $F'(x_0)$  is ill-conditioned, or if  $A$  is only a crude (e.g. triangular) approximation of  $F'(x_0)$  then  $\kappa$  should be taken larger, e.g.  $\kappa = 10$ . If this does not work, a better approximation for  $F'(x_0)$  or a better approximation  $x_0$  of the zero is required.

3.3. Determination of  $c$ . There are various ways of finding a vector  $c$  satisfying (4).

Perhaps the most straightforward approach is the following. Suppose that the expression

$$F(x_0 + w) - F(x_0) - Aw \quad (12)$$

which comes from (4) by replacing  $x$  by  $x_0 + w$  can be rewritten into an analytically equivalent expression each term of which is of order  $O(w)$  or smaller. The resulting expression can be bounded by eliminating  $w$  with  $\|w\| \leq \kappa \|\delta_0\|$ , using inequalities like the triangle inequality. If this is possible, it gives very good values for  $c$ ; an example is given in Section 5.

Alternatively, the expression can be evaluated in interval arithmetic (see e.g. Moore [8]), using in place of  $w$  an interval vector containing all  $w$  with  $\|w\| \leq \kappa \|\delta_0\|$ . Interval arithmetic has the advantage that rounding errors are automatically taken into account (cf. 3.5), but at present most interval arithmetic implementations

are rather slow (cf., however, Moore [9]).

Another approach makes use of a technique of Alefeld [1], and can be applied if  $F(x)$  is a so-called polynomial operator. Alefeld shows how to find (with interval arithmetic) two matrices  $\bar{A}$  and  $\underline{A}$  such that for every  $x \in S$ ,

$$F(x) - F(x_0) = A_x(x - x_0) \quad (13)$$

for a suitable matrix  $A_x$  with

$$\bar{A} \leq A_x \leq \underline{A}; \quad (14)$$

then (4) amounts to finding a vector  $c$  such that

$$|(A_x - A)(x - x_0)| \leq \|c_0\|c \quad \text{for all } x \in S. \quad (15)$$

If a scaled  $l_\infty$ -norm is used,

$$\|x\| = \|x\|_u := \max_{1 \leq i \leq n} \frac{|x_i|}{|u_i|} = \min_{|x| \leq cu} |x| \leq cu, \quad (16)$$

where  $u = (u_1, \dots, u_n) > 0$ , then  $x \in S$  implies  $|x - x_0| \leq \kappa \|c_0\|u$ , whence  $|(A_x - A)(x - x_0)| \leq \kappa \|c_0\| \max(|\bar{A} - A|u, |\underline{A} - A|u)$ ; therefore (15) holds with

$$c := \kappa \cdot \max(|\bar{A} - A|u, |\underline{A} - A|u). \quad (17)$$

A similar, slightly more complicated argument shows that for the Euclidean norm

$$\|x\| = \|x\|_2 = \sqrt{x^T x} \quad (16')$$

we can satisfy (15) with

$$c := \kappa \cdot (v_1, \dots, v_n)^T, \quad (17')$$

where  $v_1$  is the maximum of the Euclidean norms of the  $i$ -th rows of  $\bar{A} - A$  and  $\underline{A} - A$ .

Slightly less accurate estimates are obtainable from lower and upper bounds of the derivative,

$$\bar{A} \leq F'(x) \leq \underline{A} \quad \text{for all } x \in S; \quad (18)$$

then (13) and (14) hold with

$$A_x = \int_0^1 F'(x_0 + t(x - x_0)) dt,$$

whence again (17) can be used to compute  $c$ .

A further possibility to get an expression for  $c$  is to use bounds for the norm of the second derivative  $F''(x)$ ; we do not recommend the use of  $F''(x)$  and hence give no details (but cf. Section 4).

### 3.4. Determination of $b$ . We assume that a triangular decomposition

$$A = LR \quad (19)$$

of  $A$  is already available (from the computation of  $\delta_0$ ). The straightforward way to compute  $b = |A^{-1}|c$  is to form  $A^{-1}$ , take absolute values, and multiply with  $c$ . The formation of  $A^{-1}$  at least triples the work needed in (19), but if  $A$  is sparse the work is much more (e.g. for tridiagonal  $A:O(n)$  for (19),  $O(n^2)$  for  $A^{-1}$ ). Thus we look for ways to avoid the formation of  $A^{-1}$ .

The simplest case is when  $A$  is an  $M$ -matrix; then  $A^{-1}$  is nonnegative, and  $b$  can be computed as the solution of the doubly triangular system  $LRb = c$ . (20)

In general, if  $A^{-1}$  has not a constant sign,  $b$  cannot be formed without explicit knowledge of  $A^{-1}$ . Fortunately we only need an upper bound for the norm of  $b$ , hence the computation of an upper bound  $\bar{b}$  of  $b$  is sufficient. It is shown in Neumaier [16] that such a bound is given by the solution of the system

$$\langle L \rangle \langle R \rangle \bar{b} = c, \quad (21)$$

where the Ostrowski operator  $\langle \cdot \rangle$  replaces a matrix  $M = (m_{ik})$  by

$$\langle M \rangle := (\bar{m}_{ik}) \quad \text{with } \bar{m}_{ik} = \begin{cases} |m_{ik}| & \text{if } i=k \\ -|m_{ik}| & \text{otherwise.} \end{cases} \quad (22)$$

If we write

$$\beta := \|b\|, \quad \bar{\beta} := \|\bar{b}\|$$

then  $\beta \leq \bar{\beta}$ , hence the condition

$$\bar{\beta} \leq \kappa - 1 \quad (23)$$

implies the hypothesis of our theorem.

If  $A$  is an  $M$ -matrix then the results of [6] imply that  $\langle L \rangle = L$  and  $\langle R \rangle = R$  whence  $\bar{\beta} = \beta$ . For other matrices, it is of critical importance to know how much  $\bar{\beta}$  overestimates  $\beta$ . Examples show that  $\bar{\beta}/\beta$  can grow exponentially with the dimension  $n$ , but for diagonally dominant matrices the situation is much better, see [10]. As an illustration, the following table gives the range of  $\log_2(\bar{\beta}/\beta)$  for certain dense  $20 \times 20$  matrices (band matrices behave slightly better). We chose constant entries for  $c$ , and randomly generated matrices  $B$  with unit diagonal and off-diagonal entries uniformly distributed in  $[-\max, \max]$ . For each value of  $\max \in \{5.0, 1.0, 0.25, 0.1\}$ , 40 such matrices were generated, and  $\bar{\beta}/\beta$  was calculated for  $A=B$ ,  $A=|B|$ , and  $A = \langle B \rangle$ , respectively.

max	5.0	1.0	0.25	0.1
A = B	15...21	12...18	2...4	≤ 1.5
A =  B	16...21	12...19	2...4	≤ 1.5
A = <B>	15...23	9...16	7...12	< 0.5

Table 1. Range of  $\log_2 (\bar{\beta}/\delta)$

An estimate for  $\beta$  which usually is smaller than  $\delta$  can be obtained by using  $\beta = \|b\| \leq \|A^{-1}\| \|c\|$  and replacing  $\|A^{-1}\|$  by a lower bound for  $\|A^{-1}\|$ . Several such lower bounds have been proposed, all based on  $\|A^{-1}\| \geq \|x\| / \|Ax\|$ , with  $x$  suitably chosen to maximize the lower bound (see [3], [4], [11]). While such an estimate  $\bar{\beta}$  cannot be used for bounding the error  $\|f - x_0\|$ , it is of value in assessing whether a computed  $\bar{\beta}$  is a realistic bound for  $\beta$ .

3.5. Control of rounding errors. Due to rounding errors, the quantities  $\delta_0$ ,  $c$ ,  $b$ , and  $\beta$  will usually be calculated only approximately. Then the bounds (7) and (9) are valid only approximately, and if  $\|b\|$  is close to  $\kappa^{-1}$ , the test (6) on which (7) depends may be spoiled by rounding errors. Nevertheless, if  $A$  is not ill-conditioned,  $\|b\|$  and  $\|c_0\|$  will have the right order of magnitude, whence

$$\|b\| \ll \kappa^{-1} \quad (6a)$$

is still an indication that (7) holds. If reliable bounds are required we may use rounded interval arithmetic (Moore [8]) which automatically takes account of the rounding errors made by working with the intervals of uncertainty of each variable. The solution of linear interval equations is easiest if  $A$  is diagonally dominant or at least an  $H$ -matrix, for then the interval version of Gauss elimination can be carried out without problems (cf. Alefeld [1]) moreover the system (21) can be programmed in the equivalent but simpler form

$$[-\bar{b}, \bar{b}] := \text{solution of } LRz = f - c, c]. \quad (21a)$$

#### 4. The Kantorovic bound

In this section we compare our bounds with the a priori bound which is part of the famous theorem of Kantorovic. This theorem gives conditions which guarantee that Newton's iteration converges to a zero of a function  $F(x)$ . As a byproduct, the following bounds for a zero are obtained (cf. e.g. Ortega and Rheinboldt [12]).

##### Proposition (Kantorovic)

Let  $D$  be a subset of  $\mathbb{R}^n$ , and suppose that the function  $F: D \rightarrow \mathbb{R}^n$  has in  $D$  a continuous Fréchet derivative  $F'(x)$ . Assume that for some  $x_0 \in D$ ,  $F'(x_0)$  is nonsingular, and

$$\|F'(x_0)^{-1} F(x_0)\| \leq \alpha < 1, \quad (24)$$

$$\|F'(x_0)^{-1}\| \leq \beta,$$

$$\|F'(x) - F'(y)\| \leq \gamma \|x - y\| \quad \text{for } x, y \in S_0,$$

where  $S_0$  is a suitable convex subset of  $D$ . If

$$\gamma \beta := 2\alpha \beta \gamma \leq 1, \quad (25)$$

and

$$S := \{x \in \mathbb{R}^n \mid \|x - x_0\| \leq r := \frac{2\alpha}{1 + \sqrt{1 - \gamma \beta}}\} \subseteq S_0 \quad (26)$$

then  $F(x)$  has a unique zero in  $S$ .

REMARKS. 1. There is also an affine invariant form of the proposition, see e.g. Kiel [7]. But in practical applications this simply amounts to applying the preceding to  $BF(x)$  in place of  $F(x)$ , with a suitable preconditioning matrix  $B$ .

2. Usually, the constant  $\gamma$  is determined as a bound of the norm of the second derivative  $F''(x)$  in  $S_0$ .

3.  $\gamma$  usually depends on the choice of  $S_0$ . Since in (26),  $\alpha < r \leq 2\alpha$ , it is sensible to choose

$$S_0 := \{x \in \mathbb{R}^n \mid \|x - x_0\| \leq \kappa \alpha\}. \quad (27)$$

with a parameter  $\kappa$  satisfying

$$1 < \kappa \leq 2; \quad (28)$$

of course, we must assume that  $S_0 \subseteq D$  (note that Rall [14], Theorem 24.1, recommends to use (27) with  $\kappa = 2$ ). With this choice of  $S_0$ ,

(25) and (26) are equivalent with  

$$h := 2\alpha\beta\gamma \leq 4(k-1)/\kappa^2. \quad (29)$$

4. The relation

$$\|F(x) - F(x_0) - F'(x_0)(x - x_0)\| \leq \frac{1}{2}\gamma \|x - x_0\|^2 \quad \text{for } x, y \in S_0 \quad (30)$$

is a consequence of (24); see [12].

Now assume that  $\|\cdot\|_u$  is a scaled  $l_\infty$ -norm. Then, with  $S_0$  as in (27), we find from (30) that

$$|F(x) - F(x_0) - F'(x_0)(x - x_0)|_u \leq \frac{1}{2}\gamma \kappa^2 \alpha^2 u. \quad (31)$$

Hence we are in the situation of Section 2 with

$$A = F'(x_0), \quad \|\delta_0\| = \alpha, \quad c = \frac{1}{2}\alpha\gamma\kappa^2 u,$$

and we find  $\|b\| = \|A^{-1}c\| = \frac{1}{2}\alpha\gamma\kappa^2 \|A^{-1}\|_u \leq \frac{1}{2}\alpha\gamma\kappa^2 \beta = \frac{1}{2}h\kappa^2$ .

We see that (6) and (29) are equivalent inequalities. Hence, for scaled  $l_\infty$ -norms, the a priori Kantorovic bound is equivalent to our theorem, when  $c$  is computed from the relation (30).

Now we show by an example (in dimension  $n=1$ ) that our theorem may produce a bound when Kantorovic's bound does not work. Consider

$$F(x) = x^3 + 12x + 12, \quad D = \mathbb{R}$$

$$x_0 = 0, \quad A = F'(x_0) = 12.$$

Then Kantorovic's constants are

$$\alpha = 1, \quad \beta = \frac{1}{12}, \quad \gamma = \sup_{x, y \in S_0} 3|x + y| \geq 6 \sup_{x \in S_0} |x|.$$

Condition (25) implies that  $\gamma \leq 6$ , hence  $\sup_{x \in S_0} |x| \leq 1 = \alpha$ , whence

$$S_0 \subseteq \{x \in \mathbb{R} \mid \|x - x_0\| \leq \alpha\},$$

and (26) cannot be satisfied. Therefore, no bound is obtained.

On the other hand, our theorem, applied with  $\kappa = 3/2$  gives

$$\delta_0 = A^{-1} F(x_0) = 1,$$

$$S = \{x \in \mathbb{R} \mid |x| \leq 3/2\},$$

$$|F(x) - F(x_0) - A(x - x_0)| = |x^3| \leq 27/8,$$

whence

$$c = 27/8, \quad b = \|A^{-1}\|c = 9/32 \leq 1/2 = \kappa - 1.$$

Therefore,  $F(x)$  has a zero  $\delta$  with  $|\delta| \leq 3/2$ . Note that (9) gives the sharper estimate  $|\delta+1| \leq 9/32$ , whence  $-1.3 < \delta < -0.7$ . In fact, we have  $\delta \approx -0.932441$ .

From the relationship to Kantorovic's bound it again appears that it is sensible to choose  $\kappa \leq 2$ . But there are cases when a value of  $\kappa > 2$  is useful. For example, choose  $t > 0$  and put

$$F(x) = \begin{cases} e^x - 1 & \text{if } x \leq t, \\ e^t - 1 + e^t(x - t) & \text{if } x > t. \end{cases} \quad (32)$$

Then  $F(x)$  has a continuous derivative

$$F'(x) = e^{\min(x, t)}.$$

The unique zero of  $F(x)$  is  $\delta = 0$ . Hence if we take

$$x_0 := t, \quad A := F'(x_0) = e^t,$$

so that

$$\delta_0 = 1 - e^{-t} (> 0),$$

the set (3) contains the zero only if  $\kappa \|\delta_0\| > t$ , or

$$\kappa \geq t/(1 - e^{-t}). \quad (33)$$

Now suppose we choose  $\kappa$  according to (33), so that

$$r := \kappa \delta_0 \geq t.$$

Since

$$F(x_0 + w) - F(x_0) - Aw = \begin{cases} e^t(e^w - w - 1) & \text{if } w < 0, \\ 0 & \text{if } w \geq 0. \end{cases}$$

is nonnegative and decreasing, it takes its maximal absolute value for the smallest admissible  $w$ . Therefore, (4), (5) hold with

$$c = e^t(e^{-r} + r - 1)/\delta_0,$$

$$b = (e^{-r} + r - 1)/\delta_0 = \kappa - (1 - e^{-r})/(1 - e^{-t}) \leq \kappa - 1.$$

By our theorem, the bound

$$|\delta - x_0| \leq \kappa \delta_0 \quad (34)$$

is valid for every  $\kappa$  satisfying (33).

Remarks. 1. If  $t \geq 2$  then (33) implies  $\kappa > 2$ , whence no value of  $\kappa \leq 2$  can be used to bound the zero.

2. If  $x_0$  is "close" to the zero  $\delta$  then  $t$  is small, and we may take  $\kappa$  small as well; note that  $t/(1 - e^{-t}) = 1 + \frac{1}{2}t + O(t^2)$ .

3. The example is unusual in that every large  $\kappa$  gives a bound. This is due to the fact that the function  $F(x)$  is asymptotically linear for large  $x$ .

4. If  $\kappa$  is chosen optimally (equality in (33)) then the bound (34) gives the error exactly (equality in the bound). Although this does not generalize to arbitrary functions  $F(x)$ , it indicates the quality of the estimate.

#### 5. Numerical example

As an illustration, we consider the nonlinear bounding value problem ("minimal surface equation")

$$(1 + u_y^2) u_{xx} - 2u_x u_{xy} + (1 + u_x^2) u_{yy} = 0 \quad \text{in } D, \quad (35)$$

$$u(x, y) = g(x, y) \quad \text{on the boundary } \partial D.$$

We look at the problem when  $D$  is the square,

$$D = \{ (x, y) \mid 0 \leq x \leq 1, 0 \leq y \leq 1 \},$$

replace the square by a grid of  $(N+1)^2$  points,

$$\left( \frac{l}{N}, \frac{k}{N} \right), \quad l, k \in \{0, 1, \dots, N\},$$

and approximate the function  $u$  and its derivative in these points as follows:

$$u\left(\frac{l}{N}, \frac{k}{N}\right) \approx v_{lk},$$

$$u_x\left(\frac{l}{N}, \frac{k}{N}\right) \approx \frac{N}{2} \Delta_x v_{lk}, \quad u_y\left(\frac{l}{N}, \frac{k}{N}\right) \approx \frac{N}{2} \Delta_y v_{lk},$$

$$u_{xx}\left(\frac{l}{N}, \frac{k}{N}\right) \approx N^2 \Delta_{xx} v_{lk}, \quad u_{yy}\left(\frac{l}{N}, \frac{k}{N}\right) \approx N^2 \Delta_{yy} v_{lk},$$

$$u_{xy}\left(\frac{l}{N}, \frac{k}{N}\right) \approx \frac{N^2}{4} \Delta_{xy} v_{lk},$$

where

$$\Delta_x v_{lk} = v_{l+1, k} - v_{l-1, k}, \quad \Delta_y v_{lk} = v_{l, k+1} - v_{l, k-1},$$

$$\Delta_{xx} v_{lk} = v_{l+1, k} - 2v_{l, k} + v_{l-1, k}, \quad \Delta_{yy} v_{lk} = v_{l, k+1} - 2v_{l, k} + v_{l, k-1},$$

$$\Delta_{xy} v_{lk} = v_{l+1, k+1} - v_{l+1, k-1} - v_{l-1, k+1} + v_{l-1, k-1}.$$

The boundary conditions give

$$v_{lk} = g\left(\frac{l}{N}, \frac{k}{N}\right) \quad \text{if } l \in \{0, N\} \quad \text{or } k \in \{0, N\},$$

and the remaining  $v_{lk}$  are arranged as an  $(N-1)^2$ -dimensional block vector

$$v = (v_{11}, \dots, v_{N-1, N-1})^T, \quad \text{with } v_{11} = (v_{11}, \dots, v_{1, N-1})^T.$$

At a fixed interior point  $\left(\frac{l}{N}, \frac{k}{N}\right)$ , the differential equation is replaced by the discrete equation

$$0 = F_{lk}(v) = (\alpha + \delta_x^2) \delta_{xx} - \frac{1}{2} \delta_x \delta_{xy} + (\alpha + \delta_x^2) \delta_{yy}, \quad (36)_{lk}$$

where

$$\alpha = 4/N^2,$$

$$\delta_x = \Delta_x v_{lk}, \quad \delta_y = \Delta_y v_{lk}, \quad \text{etc.}$$

For a small vector  $w$  with

$$\|w\|_\infty \leq w, \quad (37)$$

$$w_{lk} = 0 \quad \text{if } l \in \{0, N\} \quad \text{or } k \in \{0, N\}, \quad (38)$$

the quantities

$$\epsilon_x = \Delta_x v_{lk}, \quad \epsilon_y = \Delta_y v_{lk}, \quad \text{etc.}$$

are bounded by

$$|\epsilon_x| \leq 2w, \quad |\epsilon_y| \leq 2w,$$

$$|\epsilon_{xx}| \leq 4w, \quad |\epsilon_{xy}| \leq 4w, \quad |\epsilon_{yy}| \leq 4w.$$

Now

$$F_{lk}(v+w) = (\alpha + (\delta_x + \epsilon_x)^2) (\delta_{xx} + \epsilon_{xx}) - \frac{1}{2} (\delta_x + \epsilon_x) (\delta_{xy} + \epsilon_{xy}) (\delta_{xy} + \epsilon_{xy}) + (\alpha + (\delta_x + \epsilon_x)^2) (\delta_{yy} + \epsilon_{yy}).$$

We extract from this the part linear in the epsilons and obtain a formula for the derivative:

$$F'(v)w_{lk} = (\alpha + \delta_x^2) \epsilon_{xx} - \frac{1}{2} \delta_x \delta_{xy} \epsilon_{xy} + (\alpha + \delta_x^2) \epsilon_{yy} + (2\delta_x \delta_{xy} - \frac{1}{2} \delta_x \delta_{xy}) \epsilon_x + (2\delta_y \delta_{xy} - \frac{1}{2} \delta_x \delta_{xy}) \epsilon_y. \quad (39)_{lk}$$

If we take differences we find

$$F_{lk}(v+w) - F_{lk}(v) - [F'(v)w]_{lk} = \epsilon_y^2 (\delta_{xx} + \epsilon_{xx}) - \frac{1}{2} \epsilon_x \delta_{xy} (\delta_{xy} + \epsilon_{xy}) + \epsilon_x^2 (\delta_{yy} + \epsilon_{yy}) + 2\delta_x \delta_{xy} \epsilon_x - \frac{1}{2} (\delta_x \delta_{xy} + \epsilon_x \delta_{xy}) \epsilon_y + 2\delta_x \delta_{xy} \epsilon_y.$$

hence

$$|F(v+w) - F(v) - F'(v)w|_{lk} \leq 4w^2 (\epsilon_{1k} + 10w), \quad (40)_{lk}$$

where

$$\epsilon_{1k} = |\delta_{xx}| + \frac{1}{2} |\delta_{xy}| + |\delta_{yy}| + 5|\delta_x| + 5|\delta_y|. \quad (41)_{lk}$$





Acknowledgment I want to thank Prof. Collatz for his critical remarks, and Prof. Törnig for suggesting the minimal surface equation as an interesting application.

#### References

1. G. Alefeld, Über die Durchführbarkeit des Gaußschen Algorithmus bei Gleichungen mit Intervallen als Koeffizienten, Computing, Suppl. 1 (1977), 15-19.
2. G. Alefeld, Bounding the slope of polynomial operators and some applications, Computing 26 (1981), 227-237.
3. A.L. Cline, C.B. Moler, G.W. Stewart, and J.H. Wilkinson, An estimate for the condition number of a matrix, SIAM J.Numer.Anal. 16 (1979), 368-375.
4. J.J. Dongarra, J.R. Bunch, C.B. Moler, and G.W. Stewart, LINPACK Users Guide. SIAM, Philadelphia 1979.
5. J. Douglas, A method of numerical solution of the problem of Plateau, Annals of Math. (2) 29 (1928), 180-188.
6. M. Fiedler and V. Ptak, On matrices with nonpositive off-diagonal elements and positive principal minors, Czech. Math.J. 12 (1962), 382-400.
7. G. Miel, An updated version of the Kantorovich theorem for Newton's method, Computing, to appear.
8. R.E. Moore, Methods and Applications of Interval Analysis, SIAM Publications, Philadelphia 1979.
9. R.E. Moore, New results on nonlinear systems. In: Interval Mathematics 1980 (Ed. K. Nickel), Acad. Press, New York - London 1980, pp. 165-180.
10. A. Neumaier, Hybrid norms, the Ostrowski operator, and bounds for solutions of linear equations, to appear.
11. D.P. O'Leary, Estimating matrix condition numbers, SIAM J. Sci.Stat.Comput. 1 (1980), 205-209.
12. J.M. Ortega and W.C. Rheinboldt, Iterative solution of non-linear equations in several variables. Acad. Press, New York - London 1970.
13. A.M. Ostrowski, Über die Determinanten mit überwiegender Hauptdiagonale, Comm:Math.Helv. 10 (1937), 69-96.
14. L.B. Rail, Computational solution of nonlinear operator equations. Wiley, New York - London 1969.

has error parameters

$$\bar{\beta} = 2.6193, \quad \omega = 0.0742.$$

In both cases we have  $\bar{\beta} > \kappa - 1$  although the true maximal error is  $\ll \omega$ . This is partially explained by the fact that the Newton iteration converges very slowly: For the first 25 iterations starting from Douglas' approximation, the convergence is only linear, with factor 0.7. The first iterate with  $\bar{\beta} < \kappa - 1$  is the 16th, with

$$\bar{\beta} = 0.2808, \quad \omega = 0.0082$$

(this allows checking the validity of the other two bounds). Thus we have here a very nonlinear problem.

All examples were calculated on a UNIVAC 1110. The time for one Newton iteration was 0.13 sec (without error information) resp. 0.16 sec (with error information)

Use was made of (21) to compute  $\bar{\beta}$ , although the matrices A were neither M-matrices nor diagonally dominant. Perhaps this accounts for the fact that in the examples  $\omega$  was a bound even when  $\bar{\beta}$  was too large.