

Siegfried M. Rump

## References:

- /1/ Apostolatos, N., Kulisch, U., Krawczyk, R., Lortz, B., Nickel, K., Wippermann, H.-W.: The Algorithmic Language TRIPLEX ALGOL 60, Num. Math. 11, 175-180 (1968)
- /2/ Bohlander, G.: Floating-point Computation of functions with maximum accuracy. IEEE Trans. Comp. C-26, Nr. 7, 621-632 (1977)
- /3/ Bohlander, G., Kaucher, E., Klante, R., Kulisch, U., Miranker, W.L., Ulrich, Ch. und Wolff von Gudenberg, J.: FORTRAN for Contemporary numerical computation, Report RC 8348, IBM Thomas J. Watson Research Center 1980 and Computing 26, 277-314 (1981)
- /4/ Kulisch, U.: An axiomatic approach to rounded computations. TS Report No. 1020, Mathematics Research Center, University of Wisconsin, Madison, Wisconsin, 1969 und Numer. Math. 19, 1-17 (1971)
- /5/ Kulisch, U.: Interval arithmetic over completely ordered rings, TS Report No. 1105, Mathematics Research Center, University of Wisconsin, Madison, Wisconsin, 1970
- /6/ Kulisch, U.: Grundlagen des Numerischen Rechnens - Mathematische Begründung der Rechnerarithmetik. Reihe Informatik, Band 19, Wissenschaftsverlag des Bibliographischen Instituts Mannheim, 1976
- /7/ Kulisch, U., Miranker, W.L.: Computer Arithmetic in Theory and Practice, Academic Press, 1980
- /8/ Coonan, J. et al.: A proposed standard for floating-point arithmetic, SIGNUM newsletter, Oct. 1979
- /9/ INTEL 12.1586-001: The 8086 family user's manual, Numeric Supplement, July 1980
- /10/ Kulisch, U., Miranker, W.L. (Editors): A New Approach to Scientific Computation, Academic Press, 1983
- /11/ Kulisch, U., Miranker, W.L.: The Arithmetic of the Digital Computer, IBM Research Report RC 10580, 1984, to appear in SIAM Reviews
- /12/ High Accuracy Arithmetic, Subroutine Library, IBM Program Description and User's Guide, Program Number 5664-185, 1984
- /13/ Böhm, H.: Berechnung von Polynomnullstellen und Auswertung arithmetischer Ausdrücke mit garantierter maximaler Genauigkeit. Dissertation, Universität Karlsruhe 1984

Additional References are given in /7/, /10/ and /11/

Abstract. The computational results of traditional numerical algorithms on computers are usually good approximations to the solution of a given problem. However, no verification is provided for some bound on the maximum relative error of the approximation. As can be demonstrated by ill-conditioned examples, these approximations may be drastically wrong. The algorithms based on the inclusion theory (cf. [38]) do have an automatic verification process. Rather than approximations to the solution an inclusion of the solution is computed and the correctness of the bounds together with the existence and uniqueness of the solution within the bounds is automatically verified by the computer without any effort on the part of the user. The computing time of these algorithms is of the order of a comparable, standard floating-point algorithm (such as Gaussian elimination in case of general linear systems).

In the following some new results complementing the inclusion theory are given. One of the main results is that the inclusion sets need not to be convex. Therefore other types of inclusion sets such as torus-sectors can be used. Another main observation is that the new and old theorems can be proved without using fixed point theorems. Moreover improvements of existing theorems of the inclusion theory by means of weaker assumptions are presented.

Another fundamental observation is the following. It is well-known that a real iteration in  $\mathbb{R}^n$  with affine iteration function converges if and only if the spectral radius of the iteration matrix is less than one. It can be shown that a similar result holds for our inclusion algorithm: An inclusion will be achieved if and only if the spectral radius of the iteration matrix is less than one. This result is best possible.

It is demonstrated by means of theorems and examples that even for extremely ill-conditioned examples very sharp inclusions of the solution are computed. The inclusions are almost always of least significant bit accuracy, i.e. the left and right bounds of the inclusion are adjacent floating-point numbers.

$[A, B] \in \Pi T : \{x \in T \mid A \leq x \leq B\}$  for  $A, B \in T$  and  $A \leq B$ . Therefore  $\Pi T \subseteq PT$  where every element of  $\Pi T$  is nonempty, convex, closed and bounded.

By  $S$  we denote some finite set of real numbers (which could be regarded as the set of single precision floating-point numbers on some computer). We consider the set of  $n$ -tuples  $VS$  over  $S$  and the set of  $n$ -tuples  $MS$  over  $S$ . Similarly for some finite subset  $CS$  of  $\mathbb{R} \times \mathbb{R}$  we consider  $VCS$  and  $MCS$ . If  $U$  is some set out of  $S, VS, MS, CS, VCS, MCS$  and  $T$  is the corresponding set in  $\mathbb{R}, VR, MR, C, VC, MC$ , then intervals over  $U$  are defined by

$$[A, B] \in \Pi U : \{x \in T \mid A \leq x \leq B\} \text{ for } A, B \in U \text{ and } A \leq B.$$

For example, an interval  $[a, b] \in S$  consists of all real numbers  $x$  with  $a \leq x \leq b$ . Usually, interval operations  $\oplus, \otimes, \ominus, \ominus, \ominus, \ominus$  are defined by

$$A, B \in \Pi T : A \oplus B := \{x \in \Pi T \mid \exists a \in A, \exists b \in B, x = a + b\} \quad (3)$$

$$A, B \in \Pi U : A \otimes B := \{x \in \Pi U \mid \exists a \in A, \exists b \in B, x = a \cdot b\} \quad (4)$$

for  $T$  and  $U$  as above. Using (3) and (4) the operations  $\oplus$  are well defined (cf. [4], [8]). Sometimes the strict definition (4) is difficult to realize on computers, e.g. for complex division. Therefore we allow in our further discussions any isotone interval operation  $\otimes$ , i.e. any operation

$$\otimes : PT \times PT \rightarrow \Pi T \text{ with } A, B \in PT \rightarrow A \otimes B \subseteq A \oplus B \text{ respectively} \quad (5)$$

$$\otimes : PU \times PU \rightarrow \Pi U \text{ with } A, B \in PU \rightarrow A \otimes B \subseteq A \otimes B. \quad (6)$$

There will be no confusion caused by the fact that in  $\Pi T$  and  $\Pi U$  the same symbol  $\otimes$  is used. In (5) and (6) elements of the power set over  $T$  or  $U$  are allowed to define more general interval operations. As follows by (3) and (4),  $\otimes$  is best possible for interval arguments. We allow any interval rounding  $O : PT \rightarrow \Pi T$  resp.  $O : PU \rightarrow \Pi U$  satisfying

$$A \in PT : A \subseteq OA \text{ and } A \in PU : A \subseteq OA.$$

1. Introduction. Let  $\mathbb{R}$  be the set of real numbers,  $VR$  the set of real vectors with  $n$  components,  $MR$  the set of real square matrices with  $n$  rows,  $\mathbb{C}$  the set of complex numbers,  $VC$  the set of complex vectors with  $n$  components and  $MC$  the set of complex square matrices with  $n$  rows. In the following the letter  $n$  is reserved to denote the number of elements of a vector or the number of rows and columns of a square matrix. Vectors with other than  $n$  components are, for instance, denoted by  $V_{n+1, \mathbb{R}}$ , non-square matrices for example with  $l$  rows and  $m$  columns over the complex numbers by  $M_{l, m, \mathbb{C}}$ .  $I$  denotes the identity matrix.

The operations in the power set  $PT$  over  $T$  for  $T \in \{\mathbb{R}, VR, MR, \mathbb{C}, VC, MC\}$  are as usual defined by

$$A \in PT, B \in PT : A * B := \{a * b \mid a \in A, b \in B\} \quad (1)$$

for  $*$   $\in \{+, -, \cdot, /$  and well-known restrictions for  $/$ . Definition (1) applies for inner and outer operations. In case a set (an element of the power set) occurs more than once in a formula special care has to be taken. Consider as an example  $f : PV\mathbb{R} \rightarrow PV\mathbb{R}$  being defined by  $f(x) := Z + \mathbb{C} \cdot X$  for  $X, Z \in PV\mathbb{R}$  and  $\mathbb{C} \in PM\mathbb{R}$ . Then

$$f(f(X)) = Z_1 + \mathbb{C}_1 \cdot (Z_2 + \mathbb{C}_2 \cdot X) \text{ for } Z_1 = Z_2 = Z \text{ and } \mathbb{C}_1 = \mathbb{C}_2 = \mathbb{C}.$$

Moreover

$$\{z + \mathbb{C} \cdot (z + \mathbb{C} \cdot X) \mid z \in Z, \mathbb{C} \in \mathbb{C}, X \in X\} \subseteq f(f(X)) \quad (2)$$

where in general equality does not hold in (2).

The order relation  $\leq$  in  $\mathbb{R}$  is extended to  $VR$  and  $MR$  by

$$A, B \in VR : A \leq B : \{A_i \leq B_i \text{ for } 1 \leq i \leq n \text{ and } A, B \in MR : A \leq B : \{A_{ij} \leq B_{ij} \text{ for } 1 \leq i, j \leq n.\}$$

The order relation in  $\mathbb{C}$  is defined by

$$a + bi, c + di \in \mathbb{C} : a + bi \leq c + di : \{a \leq c \text{ and } b \leq d\}$$

and similarly in  $VC$  and  $MC$ .

The set of intervals  $\Pi T$  over  $T$  for  $T \in \{\mathbb{R}, VR, MR, \mathbb{C}, VC, MC\}$  is defined by

Of course, composed operations need not to be of high accuracy. Such highly accurate results for composed operations and even for whole algorithms for solving linear or nonlinear systems are delivered by the inclusion theory as described in the following.

2. Basic new results. In [4] one of the main steps to develop the inclusion theory was to apply Brouwer's Fixed Point Theorem to an affine mapping and to show contraction in order to obtain a zero from the fixed point. The following theorem gives Brouwer's Fixed Point Theorem for affine mappings, the elegant proof of which the author learnt at the University of Karlsruhe.

Theorem 1. Let  $X$  be a nonempty, compact, convex subset of the normed vector space  $E$ ,  $C: E \rightarrow E$  a linear transformation with continuous restriction  $C|_X$ ,  $z \in E$  and  $T(X) \subseteq X$  for the affine mapping  $Tx := z + Cx$ . Then there is a fixed point  $\bar{x} \in X$  of  $T: Tx = \bar{x}$ .

Proof. Let  $x^0 \in X$  and define  $x^m := \frac{1}{m+1} \cdot \{x^0 + Tx^0 + \dots + T^m x^0\}$ ,  $m \geq 1$ .

Then short computation shows  $Tx^m = \frac{1}{m+1} \cdot \{Tx^0 + T^2 x^0 + \dots + T^{m+1} x^0\}$

and by the convexity of  $X$  follows  $Tx^m \in X$ . Therefore

$$Tx^m - x^m = \frac{1}{m+1} \cdot (T^{m+1} x^0 - x^0)$$

and because  $X$  is bounded  $Tx^m - x^m \rightarrow 0$  for  $m \rightarrow \infty$ . Therefore some subsequence  $(x^{m_i})$  of  $(x^m)$  converges to some  $\bar{x} \in X$  with  $T\bar{x} = \bar{x}$ .  $\square$

To come from a fixed point of the affine mapping to a zero of the original problem (e.g. a system of linear or nonlinear equations) in [37] the contraction of the linear part is shown by assuming a proper inclusion  $TX \subseteq X$  rather than  $TX \subseteq X$ . Next we will show this contradiction for general compact subsets of  $\mathbb{R}^n$  without assuming convexity.

Lemma 2. Let  $Z \in PVR$ ,  $C \in PVR$  and  $\emptyset \neq X \in PVR$  compact. Then

- a)  $Y = \text{ch}(X)$  and  $Z + C \cdot X \subseteq Y \Rightarrow Z + C \cdot Y \subseteq Y$ .
- b)  $Y = X - X$  and  $Z + C \cdot X \subseteq X \Rightarrow C \cdot Y \subseteq Y$ .

Here  $0: PT \rightarrow IT$  resp.  $0: PU \rightarrow IU$  has to be regarded as an operator rather than a mapping. That means we do not necessarily require

$$0A = 0B \text{ for } A = B.$$

Instead the essential property for  $0$  is  $A \subseteq 0A$  for  $A \in PT$  resp.  $A \in PU$ . The reason for this is that in practical implementations of  $0$  it might happen that  $0A \not\subseteq 0B$  for  $A = B$  because of performance reasons. An overestimation of  $A$  may be allowed to save computing time, but this does not affect the following theorems.

There is always a best rounding  $\diamond$  (cf. [4]) satisfying

$$A \in PT: \diamond A = \min\{B \in IT : A \subseteq B\} \text{ and}$$

$$A \in PU: \diamond A = \min\{B \in IU : A \subseteq B\}.$$

The rounding  $\diamond$  and the operations  $\diamond$  over  $IT$  for  $\ast \in \{+, -, \cdot, /$  and  $\cup \in \{S, VS, MS, CS, VCS, MCS\}$  can be effectively implemented on computers (cf. [23]).

For an interval  $X = [A, B]$  the lower and upper bound is defined by  $\inf(X) := A$  and  $\sup(X) := B$ , or, in short notation  $\bar{X} := A$  and  $\underline{X} := B$  (where is no confusion with the algebraic closure).

In the following we use beside the ordinary inclusion  $\subseteq$  three other types of inclusions:

$$A, B \in PVR: A \subseteq B: \Leftrightarrow A \subseteq B,$$

$$A, B \in ITVR: A \subseteq B: \Leftrightarrow A \subseteq B \wedge A_i \dagger B_i \text{ for all } i, 1 \leq i \leq n,$$

$$A, B \in IIVR: A \subseteq B: \Leftrightarrow A \subseteq B \wedge A_i \dagger B_i \text{ for some } i, 1 \leq i \leq n.$$

The definitions apply similarly over  $S$ . It is

$$A, B \in IIVR: A \subseteq B \Rightarrow A \subseteq B \Rightarrow A \subseteq B.$$

For any algorithm performing verification of the results on a computer a precisely defined computer arithmetic is mandatory. Preferably this computer arithmetic should deliver highly accurate or maximum accurate results. An arithmetic delivering always results of least significant bit accuracy even for matrix and vector operations has been given in [8].

$$\sum_{j=1}^n C_{ij} \cdot (\bar{y}_j - \underline{y}_j) < \bar{X}_i - \underline{X}_i \quad (11)$$

because one of both  $\leq$  can be replaced by  $<$ . By definition (8) follows  $C_{ij} \cdot (\bar{y}_j - \underline{y}_j) = |\sum_{j=1}^n C_{ij} \cdot (\bar{X}_j - \underline{X}_j)|$  and therefore demonstrating assertion b) using (11). Assertion a) follows by b).

ad c): Follows by the preceding proof for cases a) and b). □

Consider a system of linear equations  $Ax = b$ . For an approximate inverse  $R$  of  $A$  lemma 3 can be applied to the residual function  $f(x) := R \cdot b + (I - RA) \cdot x$ . If  $f(x) \in X$  for some  $x \in \Pi VR$ , then by Brouwer's Fixed Point Theorem there is a  $\bar{x} \in X$  with  $f(\bar{x}) = \bar{x} = Rb + (I - RA)\bar{x}$  and therefore  $b - A\bar{x} \in \ker R$ . If the non-singularity of  $R$  could be shown, then the existence of a zero  $\bar{x}$  of  $Ax - b = 0$  would be demonstrated. In the following we give criteria for the regularity of  $R$  and  $A$  and show the existence of a fixed point of  $f$  and a zero of  $Ax - b = 0$  without using Brouwer's Fixed Point Theorem.

**Lemma 4.** Let  $\mathcal{C} \in \mathbb{F}^{m \times n}$  and  $\emptyset \neq X \in \mathbb{F}^{m \times n}$  compact with  $0 \in X$ . If  $\mathcal{C} \cdot X \subseteq X$  (12)

then  $\forall C \in \mathcal{C} : \rho(C) < 1$ .

**Proof.** Let  $C \in \mathcal{C}$ . For  $y := X + iX$  holds  $\mathcal{C} \cdot y \subseteq y$  and  $0 \in y$ . Suppose  $C \cdot (0)$  and let  $\lambda \in \mathbb{C}$ ,  $x \in VC$  be an eigenvalue/eigenvector pair of  $C$ . Define

$$T \in \mathbb{F}^n \text{ by } T := (y \in \mathcal{C} \mid y \cdot x \in y).$$

$T$  is not empty because  $0 \in y$  and  $T$  is closed and bounded because  $y$  is compact. Therefore there is some  $y^* \in T$  with

$$|y^*| = \max_{y \in T} |y|.$$

Then by (12) we have  $y^* \neq 0$  and

$$C \cdot (y^* x) = (y^* \lambda) x \in \mathcal{C}.$$

Because  $y^* x \in y$  and the definition of  $y^*$  follows  $|y^*| > |y^* \lambda|$  and therefore  $|\lambda| < 1$ . □

With these preparations we get the following theorem.

**Remark.** All operations are the power set operations,  $ch(X)$  denotes the convex hull.

**Proof.** ad a): Let  $z \in Z$  and  $C \in \mathcal{C}$ . Then  $y \in Y \Leftrightarrow \exists x_1, x_2 \in X : y = x_1 + \lambda(x_2 - x_1)$  for  $0 \leq \lambda \leq 1$ . Then with  $x_3 := z + C \cdot x_1$  and  $x_4 := z + C \cdot x_2$  follows  $z + C \cdot y = x_3 + \lambda(x_4 - x_3) \in Y$  by  $x_3, x_4 \in Y$  and (7).

ad b): For  $z \in Z$ ,  $C \in \mathcal{C}$  is  $y \in X - X \Leftrightarrow \exists x_1, x_2 \in X : y = x_2 - x_1$ . Therefore  $C \cdot y = (z + C \cdot x_1) - (z + C \cdot x_2) \in Y$ . □

The proof in (40) of assertion b) in lemma 2 for convex, compact sets required Brouwer's Fixed Point Theorem by  $z + C \cdot \bar{x} = \bar{x} \Leftrightarrow C \cdot (X - \bar{x}) = z + C \cdot X + z - \bar{x} \subseteq X - \bar{x}$ . The assumption of proper inclusion  $\bar{x}$  can be sharpened for interval vectors. For this purpose we need the following lemma.

**Lemma 3.** Let  $z \in \mathbb{F}^n$ ,  $\mathcal{C} \in \mathbb{F}^{m \times n}$  and  $X \in \mathbb{F}^{m \times n}$ . Then for  $Y := X \diamond X$  holds

- a)  $0 \subseteq (z + \mathcal{C} \cdot X) \subseteq X \Rightarrow \diamond(\mathcal{C} \cdot Y) \subseteq Y$
- b)  $0 \subseteq (z + \mathcal{C} \cdot X) \subseteq X \Rightarrow \diamond(\mathcal{C} \cdot Y) \subseteq Y$
- c)  $0 \subseteq (z + \mathcal{C} \cdot X) \subseteq X \Rightarrow \diamond(\mathcal{C} \cdot Y) \subseteq Y$ .

**Proof.** ad a) and b): Suppose  $0 \subseteq (z + \mathcal{C} \cdot X) \subseteq X$  and  $z \in Z$ ,  $C \in \mathcal{C}$ . Obviously  $Y = [X - X, X - X]$  and  $(\mathcal{C} \cdot Y)_i = [\sum_{j=1}^n C_{ij} \cdot |X_j - X_j|, \sum_{j=1}^n C_{ij} \cdot |X_j - X_j|]$ . Therefore the proof is finished if  $\sum_{j=1}^n C_{ij} \cdot (X_j - X_j) < X_i - X_i$  for  $1 \leq i \leq n$  holds. Define

$$\bar{y}_j := \begin{cases} X_j & \text{if } C_{ij} \geq 0 \\ X_j & \text{otherwise} \end{cases} \quad \text{and} \quad \underline{y}_j := \begin{cases} X_j & \text{if } C_{ij} \geq 0 \\ X_j & \text{otherwise} \end{cases} \quad (8)$$

for some fixed  $1 \leq i \leq n$ . Then by assumption

$$X_i \leq z_i + \sum_{j=1}^n C_{ij} \bar{y}_j \leq X_i \quad \text{and} \quad X_i \leq z_i + \sum_{j=1}^n C_{ij} \underline{y}_j \leq X_i, \quad (9)$$

where either both left or both right  $\leq$  can be replaced by  $<$ . The left of the second inequality in (8) implies

$$-z_i - \sum_{j=1}^n C_{ij} \bar{y}_j \leq -X_i \quad (10)$$

and adding the right of the first inequality of (9) and (10) yields

$$C := \begin{Bmatrix} 1 & -0.25 \\ 2 & 3 \\ 3 & -1 \end{Bmatrix} \text{ with } \rho(C) = 0.5 \text{ and } \rho(|C|) = 1 + \sqrt{0.75} > 1.8.$$

Define  $X := \{\lambda \cdot \begin{pmatrix} 1 & 1 \\ 2 & 2 \end{pmatrix} + \mu \begin{pmatrix} 1 & 1 \\ 3 & 3 \end{pmatrix} \mid -1 \leq \lambda, \mu \leq 1\} \in \mathbb{P}VR$ .

Then  $C \cdot X = \{\lambda \cdot \begin{pmatrix} 0.5 \\ 1 \end{pmatrix} + \mu \begin{pmatrix} -0.5 \\ -3 \end{pmatrix} \mid -1 \leq \lambda, \mu \leq 1\} \subseteq X$ .

The vectors (1,2)' and (1,6)' are the eigenvectors of C with eigenvalues 0.5 and -0.5. It is an open problem whether something similar to  $\mathcal{C}$  can be defined for  $X \in \mathbb{P}VR$  to weaken the assumptions of Theorem 5. Theorem 5 has been proved for compact and convex sets in [37]. The assumptions of theorem 6 can still be weakened. For preparing this we need the following lemma.

**Lemma 7.** Let  $Z \in \mathbb{P}VR, C \in NR$  and  $X \in \mathbb{I}VR$ . If C is irreducible,  $d(X) > 0$  and  $0(Z + C \cdot X) \subseteq X$

then  $\rho(|C|) < 1$ .

**Proof.** By lemma 3, c),  $|C| \cdot Y \subseteq C \diamond Y \subseteq Y$  for  $Y = X \diamond X$ . The assertion follows by Perron-Frobenius Theory (cf. [45]).  $\square$

**Theorem 8.** Let  $Z \in \mathbb{P}VR, C \in \mathbb{P}NR$  and  $X \in \mathbb{I}VR$ . If  $0(Z + C \cdot X) \subseteq X$  using Einzelschrittverfahren,

$$(15)$$

i.e. for  $Y_i := (0Z + C \cdot (Y_1, \dots, Y_{i-1}, X_i, \dots, X_n))^T$ , holds  $Y_i \subseteq X_i$  for  $1 \leq i \leq n$ , then for every  $C \in \mathcal{C} : \rho(C) \leq \rho(|C|) < 1$ .

**Proof.** Let  $z \in Z$  and  $C \in \mathcal{C}$ , arbitrarily chosen. By Perron-Frobenius Theory (cf. [45], (2.5f)page 46) there is a permutation matrix P with

$$PCP^T = \begin{pmatrix} R_{11} & R_{12} & \dots & R_{1n} \\ & R_{22} & \dots & R_{2n} \\ & & \ddots & \\ & & & R_{nn} \end{pmatrix}$$

**Theorem 5.** Let  $Z \in \mathbb{P}VR, C \in \mathbb{P}NR$  and  $\emptyset \neq X \in \mathbb{P}VR$  compact. If

$$Z + C \cdot X \subseteq X \tag{13}$$

then  $\forall C \in \mathcal{C} : \rho(C) < 1$ .

**Proof.** Follows by lemma 2, b) and lemma 4.  $\square$

If the inclusion sets are intervals, then weaker conditions suffice to prove a stronger assertion than stated in theorem 5.

**Theorem 6.** Let  $Z \in \mathbb{P}VR, C \in \mathbb{P}NR$  and  $X \in \mathbb{I}VR$ . If

$$0(Z + C \cdot X) \subseteq X \tag{14}$$

then  $\forall C \in \mathcal{C} : \rho(C) \leq \rho(|C|) < 1$ .

**Proof.** By lemma 3, b) we have  $C \cdot Y \subseteq C \diamond (C \cdot Y) \subseteq Y$  for  $Y = X \diamond X$ . If the i-th component of X is  $[X_i, \bar{X}_i]$ , then  $Y_i = [X_i - \bar{X}_i, \bar{X}_i - X_i]$ . Therefore  $Y_i = -Y_i$  for  $1 \leq i \leq n$  and  $|C| \cdot Y \subseteq C \cdot Y \subseteq Y$  and the assertion follows by lemma 4 and Perron-Frobenius Theory (cf. [45]).  $\square$

Given a general set of vectors  $x \in \mathbb{P}VR$ , for  $Y = X - X$  holds of course  $Y = -Y$ , but something similar to  $Y_i = -Y_i$  is, in general, not true. In theorem 6 with the weaker assumption  $\mathcal{C}$  instead of  $\mathcal{E}$  in theorem 5, the stronger assertion holds that the spectral radius of the absolute value of every matrix  $C \in \mathcal{C}$  is less than one. This follows from the special structure of the elements of  $\mathbb{I}VR$ . Therefore the fact, that operations in  $\mathbb{I}VS$  are easy to implement and are fast, has to be paid by the impossibility to show  $\rho(C) < 1$  where  $\rho(|C|) \geq 1$  using (14). With spending more computing time for operations in other structures this deficiency can be eliminated. However, in practical cases the matrices C are of the form  $I - RA$  where  $R \in A^{-1}$ . In our experience never a case occurred where  $\rho(I - RA) < 1$  and  $\rho(|I - RA|) \geq 1$ . On the other hand there are cases where  $\|C\| \geq 1$  for the common norms but  $X^{k+1} := Z \diamond C \diamond X^k \subseteq X^k$  holds true for  $k = 2$  or  $k = 3$ . Later some examples are shown demonstrating this fact.

With the assumptions of theorem 5  $\rho(|C|) < 1$  for all  $C \in \mathcal{C}$  is, in general, not true. A counterexample is:

where  $R_{ii} = 0$  or  $R_{ij}$  irreducible for  $1 \leq i \leq n$ . Let  $1 \leq v \leq n$  be fixed but arbitrary. If  $R_{vv} = 0$  then  $\rho(R_{vv}) < 1$ ; suppose  $R_{vv} \neq 0$  and let  $P = (1, 2, \dots, n)^T = (\sigma(1), \sigma(2), \dots, \sigma(n))^T$ ,  $\tau(\sigma(j)) := j$  for  $1 \leq j \leq n$ . Let  $R_{vv} \in M_{\mathbb{R}}^n$ ,  $k \leq n$  with indices  $1, \dots, k+1$ . The indices of  $A$  belonging to  $R_{vv}$  are  $\tau(\ell) := \alpha_1, \dots, \tau(k+1) := \alpha_k, \alpha_m := \max(\alpha_1, \dots, \alpha_k)$ . Because of the irreducibility of  $R_{vv}$ , there is a path from  $\alpha_i$  to  $\alpha_m$  for every  $1 \leq i \leq k$  (cf. [48]). Therefore there is a  $j$  with  $1 \leq j \leq k$  and  $C_{\alpha_j, \alpha_m} \neq 0$ . For arbitrary  $\bar{y}_i, \bar{y}_j \in Y_i, 1 \leq i \leq j-1$  and  $\bar{x}_i, \bar{x}_j \in X_i, 1 \leq i \leq n$  and arbitrary  $i$  with  $1 \leq i \leq n$ , holds

$$\begin{aligned} & (C \cdot (\bar{y}_1 - \bar{y}_1, \dots, \bar{y}_{i-1} - \bar{y}_{i-1}, \bar{y}_i - \bar{y}_i, \dots, \bar{y}_n - \bar{y}_n)^T)_i = \\ & (z + C \cdot (\bar{y}_1, \dots, \bar{y}_{i-1}, \bar{x}_1, \dots, \bar{x}_n)^T)_i = (z + C \cdot (\bar{y}_1, \dots, \bar{y}_{i-1}, \bar{x}_2, \dots, \bar{x}_n)^T)_i = Y_i - X_i \\ & \text{and therefore } (C \cdot (\bar{y}_1 - \bar{y}_1, \dots, \bar{y}_{i-1} - \bar{y}_{i-1}, \bar{y}_i - \bar{y}_i, \dots, \bar{y}_n - \bar{y}_n)^T)_i = Y_i - X_i \\ & = Y_i \otimes Y_i - X_i - X_i. \end{aligned}$$

That means for  $U := X - X$

$C \otimes U \notin U$  using Einzelschrittverfahren. (16)

Let  $M := (\alpha_1, \dots, \alpha_k)$  and define

$$C_{if}^* := \begin{cases} C_{ik} & \text{if } i \in M, i \in M \\ 0 & \text{otherwise} \end{cases} \quad \text{and } U_i^* := \begin{cases} U_i & \text{if } i \in M \\ 0 & \text{otherwise} \end{cases}$$

Because  $0 \in U_i, 1 \leq i \leq n$  we have

$$(C^* \cdot U^*)_i = (C \cdot (U_1 - X_1, \dots, U_m - X_m, \dots, U_m - X_m, \dots, U_n - X_n)^T)_i = (C \cdot (U_1 - X_1, \dots, U_m - X_m, \dots, U_m - X_m, \dots, U_n - X_n)^T)_i$$

because  $C_{\alpha_j, \alpha_m} \neq 0$  and the definition of  $U_m$ . Define  $C_M \in M_{\mathbb{R}}^m$  resp.  $U_M \in V_{\mathbb{R}}^m$  to be the matrix resp. vector out of  $C$  resp.  $U$  with indices in  $M$ . Then  $C_M \cdot U_M \in U_M$ . But  $C_M$  is the permuted  $R_{vv}$  so that by lemma 7 the spectral radius  $\rho(C_M)$  is less than one. Because  $v$  was arbitrary with  $1 \leq v \leq n$  the theorem is proved.  $\square$

Theorem 8 provides a sharp test for  $\rho(|C|) < 1$ . As compared to the general case of compact inclusion sets  $X \in \text{EPVR}$ , for  $X \in \text{IIVR}$  the weaker inclusion  $\S$  and the Einzelschrittverfahren modus suffices.

Furthermore, as will be shown later, inclusion using Einzelschrittverfahren saves computing time and memory. It is an open problem whether something similar can be defined and be achieved for the general case  $X \in \text{EPVR}$  compact.

Theorem 9. Let  $z \in \text{IVR}, C \in M_{\mathbb{R}}^n$  and  $X \in \text{IIVR}$ ,  $d(X) > 0$ . If

$$\diamond(z + C \cdot X) \supseteq X \quad (17)$$

then  $\rho(|C|) \geq 1$ .

Proof. Let  $m(X)$  be the midpoint of  $X$  and  $Y := X - m(X)$ . Then  $Y \in \text{IIVR}$ ,  $0 \in Y$  and

$$\begin{aligned} \diamond(z + C \cdot X) &= z + C \cdot X + C \cdot m(X) + Y = z + C \cdot m(X) + C \cdot Y \supseteq X \\ &= C \cdot Y \supseteq Y + m(X) = z - C \cdot m(X). \end{aligned}$$

Because  $Y \in \text{IIVR}$  and  $Y = -Y$  this implies (abbreviating  $v := m(X) - z - C \cdot m(X)$ )

$$-|C| \cdot |Y| \leq -|Y| + v \leq |Y| + v \leq |C| \cdot |Y|$$

and therefore

$$|C| \cdot |Y| \geq |Y|. \quad (18)$$

Observing  $|Y| > 0$  the proof is finished by applying exercise 2, p.47 in [48].  $\square$

If an iteration  $X^{k+1} := \diamond(z + C \cdot X^k)$  is performed, theorem 9 gives a stopping criterion because  $X^{k+1} \subseteq X^k$  would imply  $\rho(|C|) < 1$  for every  $C \in \mathbb{C}$ . Theorem 9 holds true because of the special shape of interval vectors. However, (17) does not imply  $\rho(C) > 1$  because of interval dependencies. A counterexample is:

$$C := \begin{pmatrix} 1 & -1 \\ 1 & -1 \end{pmatrix} \quad \text{for } z = 0 \text{ and } X := \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mid -1 \leq x_1, x_2 \leq 1.$$

Then  $z + C \cdot X = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mid -2 \leq x_1 = x_2 \leq 2$ , but  $z + C \cdot X = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \mid -2 \leq x_1, x_2 \leq 2 \supseteq X$  with  $\rho(C) = 0$ , where  $\rho(|C|) = 2$ .

Applying iteration schemes to  $f(x) = Ax - b$  the aim is to demonstrate the existence (and maybe uniqueness) of a zero of  $f(x)$  within a certain domain. A preparation for this is the following lemma.

Lemma 10. Let  $z \in VR$ ,  $C \in MR$  and  $\emptyset * X \in PVR$  compact. If  $\rho(C) < 1$  and

$$z + C \cdot X \subseteq X \tag{19}$$

then there is one and only one  $\hat{x} \in VR$  with  $z + C \cdot \hat{x} = \hat{x}$ . It is  $\hat{x} = (I-C)^{-1} \cdot z \in X$ .

Proof. Define  $f: VR \rightarrow VR$  by  $f(x) := z + C \cdot x$ , then for  $m \in \mathbb{N}$

$$f^m(x) = \sum_{i=0}^{m-1} C^i \cdot z + C^m \cdot x \text{ and } (I-C) \cdot f^m(x) = z + C^m \cdot (x - z)$$

By (19)  $f^m(x) \in X$  for  $x \in X$  and therefore  $(I-C) \cdot f^m(x) + z$  for every  $x \in X$ ,  $m \rightarrow \infty$  because  $C$  is convergent. But  $I-C$  is invertible and  $X$  closed and bounded, therefore  $f^m(x) + (I-C)^{-1} \cdot z$  for  $m \rightarrow \infty$ , i.e.  $\hat{x} := (I-C)^{-1} \cdot z \in X$ . Now

$$f(\hat{x}) - \hat{x} = z + C \cdot (I-C)^{-1} \cdot z - (I-C)^{-1} \cdot z = 0 \tag{20}$$

Suppose  $z + C \cdot \hat{y} = \hat{y}$ . Then

$$C \cdot (\hat{x} - \hat{y}) = (z + C \cdot \hat{x}) - (z + C \cdot \hat{y}) = \hat{x} - \hat{y} \tag{21}$$

implying  $\hat{x} - \hat{y} = 0$  because  $\rho(C) < 1$ .

With the previous theorems 5 and 8 we have the tool to verify  $\rho(C) < 1$  by means of conditions (13) and (18), which can be implemented on computers. This leads to the following.

Theorem 11. Let  $Z \in PVR$ ,  $C \in PVR$  and  $\emptyset * X \in PVR$  compact. If

$$z + C \cdot X \subseteq X \tag{22}$$

then for every  $z \in Z$  and for every  $C \in \mathcal{C}$  holds  $\rho(C) < 1$  and there is one and only one  $\hat{x} \in VR$  with  $z + C \cdot \hat{x} = \hat{x}$ . It is  $\hat{x} = (I-C)^{-1} \cdot z \in \hat{X}$ .

Proof. Let  $z \in Z$  and  $C \in \mathcal{C}$  fixed but arbitrary. Then (22) and theorem 5 yields  $\rho(C) < 1$  and (19) and by lemma 10 the existence of a  $\hat{x} \in X$  with  $z + C \cdot \hat{x} = \hat{x}$ . Moreover  $\hat{x} = z + C \cdot \hat{x} \in z + C \cdot X \subseteq X$  by (22).  $\square$

Theorem 12. Let  $Z \in PVR$ ,  $C \in PMR$  and  $X \in IIMR$ . If

$$O(Z + C \cdot X) \subseteq X \text{ using Einzelschrittverfahren,} \tag{23}$$

i.e. for  $Y_i := \{O(Z + C \cdot (Y_1, \dots, Y_{i-1}, X_1, \dots, X_n)^T)\}_i$  holds  $Y_i \subseteq X_i$  for  $1 \leq i \leq n$ , then for every  $z \in Z$  and for every  $C \in \mathcal{C}$  holds  $\rho(C) < 1$  and there is one and only one fixed point  $\hat{x} \in VR$  with  $z + C \cdot \hat{x} = \hat{x}$ . It is  $\hat{x} = (I-C)^{-1} \cdot z \in Y$  and  $Z + C \cdot Y \subseteq Y$ .

Proof. Let  $z \in Z$  and  $C \in \mathcal{C}$  fixed but arbitrary. Then (23) and theorem 8 implies  $\rho(C) < 1$ . Moreover

$$\{z + C \cdot (Y_1, \dots, Y_n)^T\}_i \subseteq \{z + C \cdot (Y_1, \dots, Y_{i-1}, X_1, \dots, X_n)^T\}_i = Y_i \text{ for } 1 \leq i \leq n$$

implying  $z + C \cdot Y \subseteq Y$ . Lemma 10 concludes the proof.  $\square$

The next theorem gives another improvement of the inclusion assumption (13) in theorem 5, combined with an iteration scheme.

Theorem 13. Let  $Z \in PVR$ ,  $C \in PMR$  and  $\emptyset * X \in PVR$ , all  $Z, C$  and  $X$  being compact. Define  $f: PVR \rightarrow PVR$  by  $f(Y) := Z + C \cdot Y$  for  $Y \in PVR$ . If

$$f^{k+m}(X) \subseteq f^k(X) \text{ for some } k, m \in \mathbb{N}, k \geq 0, m \geq 1 \tag{24}$$

then for every  $z \in Z$  and for every  $C \in \mathcal{C}$  holds  $\rho(C) < 1$  and there is one and only one  $\hat{x} \in VR$  with  $z + C \cdot \hat{x} = \hat{x}$ . It is  $\hat{x} = (I-C)^{-1} \cdot z \in f^k(X)$ .

Proof. Let  $z \in Z$  and  $C \in \mathcal{C}$  fixed but arbitrary. Define  $g: VR \rightarrow VR$  by  $g(x) := z + C \cdot x$  and let  $Y := f^k(X)$ .  $Y$  is compact because  $Z, C$  and  $X$  are, so with

$$g^m(Y) = \sum_{i=0}^{m-1} C^i \cdot z + C^m \cdot Y \subseteq f^m(Y) \subseteq Y, \tag{25}$$

theorem 5 implies  $\rho(C) < 1$  and therefore  $\rho(C) < 1$ , and  $g^m(\hat{x}) = \hat{x}$  for  $\hat{x} := (I-C)^{-1} \cdot ( \sum_{i=0}^{m-1} C^i \cdot z )$ . On the other hand

$$\left( \sum_{i=0}^{m-1} C^i \right) (I-C) = I - C^m, \tag{26}$$

and by  $\rho(C) < 1$  and  $\rho(C^m) < 1$  follows  $\hat{x} = (I-C)^{-1} \cdot z \in Y$ . Now (20) and (21) complete the proof.  $\square$

When aiming to compute an inclusion of the solution of a given problem (e.g. a system of linear or nonlinear equations), an iteration scheme is very useful. Since the problem of computing an inclusion is later reduced to verify contraction for an affine mapping, let us consider (24) more closely. Consider the one-dimensional, con-

tracting affine function

$$f(x) := 3 - 0.5 \cdot x. \tag{27}$$

For  $X := [1.9, 2.1]$  is  $f(X) = [1.95, 2.05] \subseteq X$ . Therefore  $f$  has one and only one fixed point by theorem 11 and  $\hat{x} = (1 - (-0.5))^{-1} \cdot 3 = 2$ . By the formula used in (25) is

$$f^m(x) = \sum_{i=0}^{m-1} (-0.5)^i \cdot 3 + (-0.5)^m \cdot x = \frac{1 - (-0.5)^m}{1 - (-0.5)} \cdot 3 + (-0.5)^m \cdot x = 2 + (-0.5)^m \cdot (x-2).$$

For  $X^0 := [1.9, 2.3]$  and  $X^m := f^m(X)$  is

$$X^1 = [1.85, 2.05], X^2 = [1.975, 2.075], X^3 = [1.9625, 2.0125], \dots$$

It is easy to show  $X^{m+1} \subseteq X^m$  for every  $0 \leq m \in \mathbb{N}$ . Therefore the contraction cannot be established by theorem 11. However,  $X^0 \subseteq X^0$  and theorem 13 can be applied. For  $X^0 := [3, 5]$  is

$$X^1 = \{0.5, 1.5\}, X^2 = \{2.25, 2.75\} \text{ and } X^m = 2 + (-0.5)^m \cdot [1, 3]. \tag{28}$$

Therefore  $2 \notin X^m$  for every  $0 \leq m \in \mathbb{N}$  and by theorem 13, (24) can never be satisfied. Therefore, although the function  $f$  defined by (27) is contracting, no inclusion is possible using any of the theorems so far (for the Einzelschrittverfahren as well since the function is one-dimensional).

Next we seek for methods providing inclusions even under extreme circumstances. The next theorems allows an inclusion for the example above.

Theorem 14. Let  $Z \in \mathbb{P}VR$ ,  $C \in \mathbb{P}MR$  and  $\emptyset \neq X \in \mathbb{P}VR$ , all  $Z, C$  and  $X$  being compact. Define  $f : \mathbb{P}VR \rightarrow \mathbb{P}VR$  by  $f(V) := Z + C \cdot V$  for  $V \in \mathbb{P}VR$ . If

$$f^{m+1}(X) \subseteq \bigcup_{i=0}^m f^i(X) \text{ for some } 0 \leq m \in \mathbb{N}, \tag{29}$$

then for every  $z \in Z$  and for every  $C \in \mathbb{C}$  holds  $\rho(C) < 1$  and there is one and only one  $\hat{x} \in VR$  with  $z + C \cdot \hat{x} = \hat{x}$ . It is  $\hat{x} = (I - C)^{-1} \cdot z \in \bigcup_{i=0}^m f^i(X)$ .

Proof. For  $U, V \in \mathbb{P}VR$  is

$$f(U \cup V) = f(U) \cup f(V).$$

For  $Y := \bigcup_{i=0}^m f^i(X)$  is  $f^{m+1}(X) \subseteq Y$  and by induction for  $k > 1$

$$f^{m+k+1}(X) \subseteq f(f^{m+k}(X)) \subseteq f(Y) \subseteq f^{m+1}(X) \cup \bigcup_{i=0}^m f^i(X) \subseteq Y. \tag{30}$$

Therefore

$$f^{m+1}(Y) = \bigcup_{i=0}^m f^{m+1}(f^i(X)) = \bigcup_{i=0}^m f^{m+i+1}(X) \subseteq Y.$$

$Y$  is compact because  $Z, C$  and  $X$  are and applying theorem 13 completes the proof.  $\square$

Corollary 15. Let  $Z \in \mathbb{P}VR$ ,  $C \in \mathbb{P}MR$  and  $\emptyset \neq X \in \mathbb{P}VR$ , all  $Z, C$  and  $X$  being compact. Define  $f : \mathbb{P}VR \rightarrow \mathbb{P}VR$  by  $f(V) := Z + C \cdot V$  for  $V \in \mathbb{P}VR$ . If  $(\bigcup_{i=0}^m f^i(X)) \subseteq \bigcup_{i=0}^m f^i(X)$  for some  $0 \leq m \in \mathbb{N}$ ,

then for every  $z \in Z$  and every  $C \in \mathbb{C}$  holds  $\rho(C) < 1$  and there is one and only one  $\hat{x} \in VR$  with  $z + C \cdot \hat{x} = \hat{x}$ . It is  $\hat{x} = (I - C)^{-1} \cdot z \in \bigcup_{i=0}^m f^i(X)$ .

Proof. Follows by replacing  $U$  by  $\bigcup_{i=0}^m f^i(X)$  in the proof of theorem 14 or, by lemma 2, a) and theorem 13.  $\square$

Note, that corollary 15 holds true because  $f$  is an affine function. Let  $A, B \in \mathbb{P}VR$ , then

$$\begin{aligned} f(A \cup B) &= \{f(x) \mid x = a + \lambda(b-a) \text{ for } a \in A, b \in B, \lambda \in \mathbb{R} \text{ with } 0 < \lambda \leq 1\} \\ &= \{z + C \cdot (a + \lambda(b-a)) \mid a \in A, b \in B, \lambda \in \mathbb{R} \text{ with } 0 < \lambda \leq 1\} \\ &= \{z + C \cdot a + \lambda \cdot (z + C \cdot b) - (z + C \cdot a) \mid a \in A, b \in B, \lambda \in \mathbb{R} \text{ with } 0 < \lambda \leq 1\} \\ &= \{z + C \cdot a \mid a \in A\} \cup \{z + C \cdot b \mid b \in B\} \\ &= f(A) \cup f(B). \end{aligned}$$

If the function  $f$  is not affine, there need not be a relation between  $f(A \cup B)$  and  $f(A) \cup f(B)$ . Consider as examples:

- 1.)  $f : \mathbb{R} \rightarrow \mathbb{R}$  with  $f(x) = x$  and  $A := \{1\}, B := \{-1\}$ . Then  $f(A \cup B) = \{f(x) \mid x \in \{-1, 1\}\} = \{0, 1\}$  and  $f(A) \cup f(B) = \{1\} \cup \{-1\}$  with  $f(A \cup B) \not\subseteq f(A) \cup f(B)$ .



Proof. Short computation yields for  $x \in \text{VFR}$

$$\tilde{X} + R(b - A\tilde{X}) + (I - RA) \cdot (x - \tilde{X}) = R \cdot b + (I - RA) \cdot x. \quad (32)$$

Therefore  $R \cdot b + (I - RA) \cdot X \subseteq X$  and by theorem 11  $\rho(I - RA) < 1$  proving the non-singularity of  $R$  and  $A$ . The fixed point  $\tilde{x}$  of the function  $R \cdot b + (I - RA) \cdot X$  satisfies  $\tilde{x} \in X$  and  $\tilde{x} = \{I - (I - RA)\}^{-1} \cdot R \cdot b = A^{-1} \cdot b$  by theorem 11.  $\square$

Obviously (31) can be replaced by  $R \cdot b + (I - RA) \cdot X \subseteq X$ , which is, for appropriate  $X$ , also verifiable on computers.

One significant improvement of the quality of inclusions is not to include the solution itself but the difference to an approximate solution (cf. [37] and [38]). As shown by (32) the approximation  $\tilde{x}$  to  $x$  does not play any role in (31). If, instead the inclusion shall satisfy  $\tilde{x} - \tilde{x} \in X$ , then the linear system  $Ax = b - A\tilde{x}$  has to be solved because

$$A(\tilde{x} - \tilde{x}) = A\tilde{x} - A\tilde{x} = b - A\tilde{x}.$$

This leads to the well-known residue iteration scheme. The corresponding inclusion theorem is the following.

Theorem 17. Let  $A, R \in \text{MR}$ ,  $b, \tilde{x} \in \text{VFR}$ . If for some compact  $\emptyset + X \in \text{FP VFR}$

$$R \cdot (b - A\tilde{x}) + (I - RA) \cdot X \subseteq X, \quad (33)$$

then the matrices  $A$  and  $R$  are not singular and the uniquely determined solution  $\hat{x} := A^{-1} \cdot b$  of  $Ax = b$  satisfies  $\hat{x} \in \tilde{x} + X$ . If

$$R \cdot (b - A\tilde{x}) + (I - RA) \cdot X \cap X = \emptyset, \quad (34)$$

then there is no solution of  $Ax = b$  in  $\tilde{x} + X$ .

Proof. The first assertion follows by applying theorem 16 replacing  $b$  by  $b - A\tilde{x}$  and regarding (32). Suppose  $A\hat{x} = b - A\tilde{x}$ . Then

$$R \cdot (b - A\tilde{x}) + (I - RA) \cdot (\hat{x} - \tilde{x}) = \hat{x} - \tilde{x} \quad (35)$$

would contradict (34) if  $\hat{x} - \tilde{x} \in X$ .  $\square$

2.)  $f: \mathbb{R}^2 \rightarrow \mathbb{R}^2$  with  $f(x, y) = (x^2, y^3)^T$  and  $A = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ ,  $B = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$ . Then

$$f(A \cup B) = \{f(x, y) \mid 0 \leq x = y \leq 1\} = \left\{ \begin{pmatrix} x^2 \\ y^3 \end{pmatrix} \mid 0 \leq x \leq 1 \right\} \text{ and}$$

$$f(A) \cup f(B) = \left\{ \begin{pmatrix} 0 \\ 0 \end{pmatrix} \right\} \cup \left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\} = \left\{ \begin{pmatrix} x \\ x \end{pmatrix} \mid 0 \leq x \leq 1 \right\} \text{ with}$$

$$f(A) \cup f(B) \not\subseteq f(A \cup B) \text{ and } f(A \cup B) \not\subseteq f(A) \cup f(B). \text{ In fact}$$

$$(f(A) \cup f(B)) \cap f(A \cup B) = \left\{ \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right\} = f(A) \cup f(B) \text{ which is the bare minimum the intersection must contain.}$$

For the function  $f$  defined in (27) and  $X := \{3, 5\}$  is by (28)

$$f^2(X) = [2.25, 2.75] \cup X \cup f(X) = \{3, 5\} \cup [0.5, 1.5] = [0.5, 5].$$

This implies  $f$  to be contractive and gives a (poor) inclusion  $[0.5, 5]$  for the fixed point  $\hat{x} = 2$  of  $f$ . In the following chapter theorems are given for improving the quality of an inclusion.

3. Linear systems. Our next aim are theorems for inclusions of the solution of systems of linear equations of those theorems should be verifiable on computers such that inclusions sets can be calculated. In [17] and [39] a number of those theorems are given. Next those theorems will be improved. The inclusion formula of the next theorem occurred in [21]. There, Krawczyk supposed  $\rho(I - RA) < 1$  a priori. Using the proper inclusion  $\underline{e}$  we have an easy verification scheme for the iteration matrix to be convergent.

Theorem 16. Let  $A, R \in \text{MR}$ ,  $b, \tilde{x} \in \text{VFR}$ . If for some compact  $\emptyset + X \in \text{FP VFR}$

$$\tilde{X} + R(b - A\tilde{X}) + (I - RA) \cdot (X - \tilde{X}) \subseteq X, \quad (31)$$

then the matrices  $A$  and  $R$  are not singular and the uniquely determined solution  $\hat{x} := A^{-1} \cdot b$  of  $Ax = b$  satisfies  $\hat{x} \in \tilde{x} + X$ .

Theorem 17 can be applied on computers if operations extending the power set operations, possibly for special inclusion sets, are available. One possibility is to restrict the inclusion sets to interval vectors, i.e. hyperrectangles or, in the complex number space, to n-dimensional balls or torus-sectors. For the application on computers we formulate theorem 17 for interval vectors over some subset G of R.  $0: PVR \rightarrow II VG$  is some isotone rounding, i.e. satisfying  $X \in PVR \Rightarrow X \subseteq OX$ .

Theorem 18. Let  $G \subseteq R$  be some subset of R,  $A, R \in MG$ ,  $b, \tilde{X} \in VG$  and  $X \in II VG$ . If

$$O(R \cdot (b - AX) + (I - RA) \cdot X) \subseteq X \text{ using Einzelschrittverfahren, (36)}$$

i.e. for  $Y_i := \{O(R(b - AX) + (I - RA) \cdot (Y_1, \dots, Y_{i-1}, X_i, \dots, X_n))^T\}_i$  holds  $Y_i \subseteq X_i$  for  $i = 1, \dots, n$ , then the matrices A and R are not singular and the uniquely determined solution  $\tilde{X} := A^{-1} \cdot b$  of  $AX = b$  satisfies  $\tilde{X} \subseteq X + Y$ .

Proof. By theorem 12 is  $\rho(I - RA) < 1$  and therefore A and R regular, and  $\{(I - (I - RA))^{-1} \cdot R \cdot (b - AX) = A^{-1} \cdot (b - AX) = A^{-1} \cdot b - \tilde{X} \in Y$ .  $\square$

Assumption (36) can be verified on computers. If G is a set of floating-point numbers and operations  $\theta: II G \rightarrow II G$  satisfying  $A, B \in II G: A \cdot B \subseteq A \theta B$  for  $\theta \in \{+, \cdot, /, \setminus\}$  are available, then

$$R \theta (b \theta A \theta \tilde{X}) \subseteq (I \theta R \theta A) \theta X \subseteq X \text{ using Einzelschrittverfahren (37)}$$

implies (36) and therefore the assertions of theorem 18 hold true. If a precise scalar product is available, then

$$R \theta O(b - Ax) \subseteq O(I - RA) \theta X \subseteq X \text{ using Einzelschrittverfahren (38)}$$

can be used. Obviously, (38) implies (36) but is much sharper than (37) because the critical parts, the residuals, are enclosed with maximum accuracy.

It occurs frequently, that either the input data of a numerical problem is not exactly representable in a given floating-point screen or, that input data is afflicted with tolerances. In the latter case, one is interested in an inclusion of the set of all possible solutions for all possible combinations of input data. In the first case, the

input data could be enclosed in the input range of the immediate predecessor and successor in the floating-point grid and the resulting problem with input afflicted with tolerances be solved. For input data with tolerances consider the following theorem.

Theorem 19. Let  $A \in PMR$ ,  $b \in PVR$ ,  $R \in MR$  and  $\tilde{X} \in VR$ . If for some compact  $\theta + X \in PVR$

$$R \cdot (b - AX) + (I - RA) \cdot X \subseteq X, \quad (39)$$

then for every  $A \in A$  and for every  $b \in b$  the following is true: the matrices A and R are not singular and the uniquely determined solution  $\tilde{X} := A^{-1} \cdot b$  of  $AX = b$  satisfies  $\tilde{X} \subseteq X$ .

Proof. Follows by applying theorem 17 to every  $A \in A$ ,  $b \in b$ .  $\square$

For immediate applications on computers the power set operations can be replaced by corresponding isotone operations, which are executable on computers. The input data A and b, which are afflicted with tolerances, can be replaced by computer representable sets. Usually, one would choose R and X to be point data. An example of such a theorem is the following.

Theorem 20. Let  $G \subseteq R$  be some subset of R,  $A \in II MG$ ,  $b \in II VG$ ,  $R \in MG$  and  $X \in VG$ . If for some  $X \in II VG$

$$R \theta O(b - AX) \subseteq O(I - RA) \theta X \subseteq X \text{ using Einzelschrittverfahren,}$$

then for every  $A \in MR$  and  $b \in VR$  with  $A \in A$  and  $b \in b$  the following is true: the matrices A and R are not singular and the uniquely determined solution  $\tilde{X} := A^{-1} \cdot b$  of  $AX = b$  satisfies  $\tilde{X} \subseteq X$ .

Proof. Follows by applying theorem 18 to every  $A \in A$ ,  $b \in b$ .  $\square$

Note, that the assertion of theorem 20 applies to every real matrix  $A \in MR$  and every real vector  $b \in VR$  with  $A \in A$ ,  $b \in b$  and not only to those with elements in G. Moreover, a computer verification holds true without any restriction, e.g. with respect to a given floating-point precision.

In [38] an algorithm is given for computing an inclusion of the solution of a system of linear equations. After computing an approxi-

mate inverse  $R$  of  $A$  a residual correction method is applied yielding some approximate solution  $\tilde{X}$ . After computing an inclusion for  $\mathcal{C} := O(I-RA)$  an iteration is started. To illustrate the differences between the different iteration schemes for obtaining an inclusion consider the following example. Let

$$A := \begin{pmatrix} 1 & 2 \\ 2 & 3 \end{pmatrix}; A^{-1} = \begin{pmatrix} -3 & 2 \\ 2 & -1 \end{pmatrix}; R := \begin{pmatrix} -2.85 & 2.1 \\ 1.85 & -1.1 \end{pmatrix}. \text{ Then}$$

$$C := \begin{pmatrix} -0.35 & -0.6 \\ 0.35 & 0.6 \end{pmatrix} \text{ and } \rho(|C|) = 0.95; \quad b := \begin{pmatrix} 1 \\ 0 \end{pmatrix}. \tag{40}$$

The following computational results are obtained on a IBMS/370 processor in single precision, i.e. 6 hexadecimal digits or approximately 7 decimal digits. We apply three different iteration schemes according to theorems 16,17 and 18:

- 1)  $Z := \tilde{X} \oplus R \oplus O(b-A\tilde{X}); X^0 := Z; X^{k+1} := Z \oplus C \oplus (X^k \oplus \tilde{X});$
- 2)  $Z := R \oplus b; X^0 := Z; X^{k+1} := Z \oplus C \oplus X^k;$
- 3)  $Z := R \oplus O(b-A\tilde{X}); X^0 := Z; X^{k+1} := Z \oplus C \oplus X^k;$

Then  $X^{k+1} \subseteq X^k$  proves the non-singularity of  $A$  and  $R$  and  $\tilde{X} \in X$  resp.  $\tilde{X} \in X$  in cases 1) and 2) resp. 3),  $\tilde{X} := A^{-1} \cdot b$  according to theorems 16,17 and 18.

In the following table from left to right the number  $i$  of the iteration scheme, the number of iterations  $k$  and the maximum relative error  $\epsilon$  for the two inclusions (defined by  $\max_{x_1, x_2 \in X} |x_1 - x_2|$  for an interval  $X$ ) are displayed.

$i$	$k$	$\epsilon$
1	no inclusion	
2	10	$1.5 \cdot 10^{-5}$
3	1	$6.0 \cdot 10^{-8}$

For the first iteration, no inclusion is obtained because the iterates  $X^k$  become cyclic. The third scheme enclosing the difference  $\tilde{X} - \tilde{X}$  is superior with respect to the number of iterations necessary and with respect to the quality of the inclusion. Nevertheless it is

possible that  $\rho(|C|) < 1$  but neither of the iteration schemes (41) obtain an inclusion. Consider

$$f(x) := 3 + 0.5 \cdot x. \tag{42}$$

Then by the formula used in (25)

$$f^m(x) = \int_{i=0}^{m-1} 0.5^i \cdot 3 + 0.5^m \cdot x = \frac{1-0.5^m}{1-0.5} \cdot 3 + 0.5^m \cdot x = 6 + 0.5^m \cdot (x-6).$$

For  $X^0 := [4, 5]$  is

$$X^1 = [5, 5.5]; X^2 = [5.5, 5.75]; X^3 = [5.75, 5.875] \text{ etc.}$$

It is easy to see that  $X^{k+m} \subseteq X^k$  and  $X^{k+m} \subseteq_{i=0}^{m-1} X^{k+i}$  for every  $k, m \in \mathbb{N}$ .

For the purpose of obtaining an inclusion even in these cases, in [37] the  $\epsilon$ -inflation has been introduced. An iteration with  $\epsilon$ -inflation is

$$\begin{aligned} &\text{repeat } Y := X \oplus \epsilon; X := Z + C \cdot Y \\ &\text{until } X \subseteq Y \end{aligned} \tag{43}$$

for some starting interval  $X$  and the property  $X \in \text{PVR} : X \subseteq X \oplus \epsilon$ . One could, e.g., define  $X \oplus \epsilon := X + [-0.5, +0.5]$ . Then with  $X := [4, 5]$  for the example in (42)

$$Y = [3.5, 5.5]; X = [4.75, 5.75]; Y = [4.25, 6.25]; X = [5.125, 6.125]$$

with an inclusion  $[5.125, 6.125]$  and all the verifications. When including the difference to an approximate solution rather than the solution itself this inclusion will be sharpened.

What is observed in the example can be generalized. It can be shown, that an iteration (43) terminates if and only if  $\rho(C) < 1$ . It is clear from the previous theorems, that if (43) terminates then  $\rho(C) < 1$ . The other direction is demonstrated in the following.

Lemma 21. Let  $Z \in \text{PVR}, C \in \text{HR}$  and  $E_i \in \text{PVR}, i \geq 1$ . Suppose  $Z$  and the  $E_i$  to be bounded,  $E_{i+1} \subseteq E_i$  and  $U_\epsilon(O) \subseteq E_i$  for all  $i \geq 1$  and some fixed  $0 < \epsilon \in \mathbb{R}$ . Define  $f: \text{PVR} \rightarrow \text{PVR}$  by

$$X \in \text{PVR} : f(X) := Z + C \cdot X. \tag{44}$$

For some  $X^0 \in \text{PVR}$  bounded define

Proof. (1)  $\Rightarrow$  (2) follows by theorem 5 and the fact, that  $X^m$  is compact and nonempty. (2)  $\Rightarrow$  (1) follows by lemma 21.  $\square$

Theorem 22 shows that whenever in inclusion is possible, namely if the iteration matrix is convergent, an inclusion will be achieved using an iteration scheme (43) with the  $\varepsilon$ -inflation. This equivalence is a best possible result, it is the equivalent to the same condition for a real iteration

$$X^{k+1} := X^k + R \cdot (b - AX^k)$$

for the linear system  $Ax = b$ ,  $R \in A^{-1}$ . Similar conditions for  $\rho(C) < 1$  in case  $X^0 \in \text{IVR}$  are given in [39].

In practice, the  $\varepsilon$ -inflation would be dependent on the iterate. For more details see [38, 39].

Theorem 22 does not necessarily hold true for sets of matrices. In fact it may happen that for a convex set of matrices  $C$  every  $C \in C$  is convergent but  $C_1 \cdot C_2 \in C$  is not. An example is

$$C := \{A + \sigma(B-A) \mid 0 \leq \sigma \leq 1\} \text{ for } A := \begin{pmatrix} 0.5 & 0.27 \\ 0.92 & 0.5 \end{pmatrix}, B := \begin{pmatrix} 0.17 & 0.6 \\ 0.94 & 0.25 \end{pmatrix}$$

Then  $\rho(A) \approx 0.9984$ ,  $\rho(B) \approx 0.9621$  and  $\max_{0 \leq \sigma < 1} \rho(A + \sigma(B-A)) \approx 0.999617$ , where the maximum is achieved for  $\sigma \approx 0.1266$ . On the hand

$$A \cdot B = \begin{pmatrix} 0.3388 & 0.3675 \\ 0.6284 & 0.677 \end{pmatrix} \text{ and } \rho(A \cdot B) \approx 1.0166.$$

Note, that every  $A + \sigma(B-A)$ ,  $0 \leq \sigma \leq 1$  is positive. The condition  $\rho(C_1 \cdot C_2) < 1$  for  $C_1, C_2 \in C$  is necessary because

$$Z + C \cdot X \subseteq X \text{ implies } Z + C \cdot Z + C \cdot X \subseteq Z + C \cdot X \subseteq X$$

for  $Z \in \text{FVR}$ ,  $X \in \text{FVR}$ ,  $X$  compact and  $C \in \text{FVR}$ . Then

$$(Z + C \cdot Z) + (C_1 \cdot C_2) \cdot X \subseteq Z + C \cdot Z + C \cdot X \subseteq X \text{ for every } C_1, C_2 \in C$$

implying  $\rho(C_1 \cdot C_2) < 1$  by theorem 5.

$$X^{k+1} := f(X^k) + E_{k+1} \text{ for } 0 \leq k \in \mathbb{N}. \tag{45}$$

If  $\rho(C) < 1$ , then there is some  $m \in \mathbb{N}$  with

$$f(X^m) \subseteq X^m. \tag{46}$$

Proof. First we prove by induction

$$X^m = \bigcup_{i=0}^{m-1} C^i \cdot (Z + E_{m-1}) + C^m \cdot X^0 \text{ for } 0 \leq m \in \mathbb{N}. \tag{47}$$

(47) is true for  $m = 0$ . Supposing (47) to hold for  $m \in \mathbb{N}$  yields by definition (44) and (45)

$$X^{m+1} = Z + E_{m+1} + C \cdot \bigcup_{i=0}^{m-1} C^i \cdot (Z + E_{m-1}) + C^{m+1} \cdot X^0 + C^{m+1} \cdot X^0. \tag{48}$$

Because  $C$  is contracting there is an  $m \in \mathbb{N}$  satisfying

$$C^m \cdot (Z + E_1) + C^{m+1} \cdot X^0 - C^m \cdot X^0 \subseteq E_{m+1}$$

because  $U_\varepsilon(0) \subseteq E_i$ ,  $i \geq 1$ . This implies

$$Z + \bigcup_{i=1}^{m-1} C^i \cdot (Z + E_{m-1}) + C^m \cdot (Z + E_1) + C^{m+1} \cdot X^0 \subseteq Z + E_{m+1} + \bigcup_{i=1}^{m-1} C^i \cdot (Z + E_{m-1}) + C^m \cdot X^0$$

(for  $A, B, C \in \text{FVR}$ :  $A - B \subseteq C$  implies  $A \subseteq B + C$ , but the contrary is not true).

Therefore

$$Z + C \cdot \left( \bigcup_{i=0}^{m-1} C^i \cdot (Z + E_{m-1}) + C^m \cdot X^0 \right) \subseteq \bigcup_{i=0}^{m-1} C^i \cdot (Z + E_{m-1}) + C^m \cdot X^0.$$

Regarding (47) and  $E_{i+1} \subseteq E_i$  proofs the lemma.  $\square$

Theorem 22. Let  $Z \in \text{FVR}$ ,  $C \in \text{FVR}$  and  $E_i \in \text{FVR}$  for  $i \geq 1$ , all  $Z$  and  $E_i$  being compact. Suppose  $E_{i+1} \subseteq E_i$  and  $U_\varepsilon(0) \subseteq E_i$  for all  $i \geq 1$  and some fixed  $0 < \varepsilon \in \mathbb{R}$ . Define  $f : \text{FVR} \rightarrow \text{FVR}$  by

$$X \in \text{FVR}: f(X) := Z + C \cdot X.$$

For some compact  $X^0 \in \text{FVR}$  define

$$X^{k+1} := f(X^k) + E_{k+1} \text{ for } 0 \leq k \in \mathbb{N}.$$

Then the following is equivalent:

- (1)  $f(X^m) \subseteq X^m$  for some  $m \in \mathbb{N}$
- (2)  $\rho(C) < 1$ .

Finally we note the generalization of theorem 14 using  $\epsilon$ -inflation. This is not clear, because there the inclusion in the convex union of all previous iterates is assumed.

**Theorem 23.** Let  $Z \in \mathbb{P}VR$ ,  $C \in \mathbb{P}MR$  and  $\emptyset \neq X^0 \in \mathbb{P}VR$ , all  $Z, C$  and  $X^0$  being compact. Define  $f : \mathbb{P}VR \rightarrow \mathbb{P}VR$  by  $f(V) := Z + C \cdot V$  and let  $g_i : \mathbb{P}VR \rightarrow \mathbb{P}VR$  with  $V \in \mathbb{P}VR \rightarrow g_i(V) \supseteq V$  for  $i \geq 0$ . Define

$$X^{k+1} := g_{k+1}(f(X^k)) \text{ for } 0 \leq k \in \mathbb{N}.$$

If then

$$f(X^m) \subseteq \bigcup_{i=0}^m X^i \text{ for some } 0 \leq m \in \mathbb{N}, \tag{49}$$

then for every  $z \in Z$  and for every  $C \in \mathbb{C}$  holds  $\rho(C) < 1$  and there is one and only one  $\bar{X} \in \mathbb{P}VR$  with  $z + C \cdot \bar{X} = \bar{X}$ . It is  $\bar{X} = (I-C)^{-1} z$ .

Proof. Let  $Y := \bigcup_{i=0}^m X^i$ . We first proof by induction

$$f^k(X^{m-k+1}) \subseteq Y \text{ for } 1 \leq k \leq m+1.$$

This is true for  $k=1$  by assumption (49). Then

$$f^{k+1}(X^{m-k}) = f^k(f(X^{m-k})) \subseteq f^k(\bigcup_{i=0}^{m-k+1} f(X^{m-k-i})) = f^k(X^{m-k+1}) \subseteq Y. \tag{50}$$

Furthermore

$$f(Y) = \bigcup_{i=0}^m f(X^i) \subseteq \bigcup_{i=0}^{m-1} f(X^i) \cup \bigcup_{i=0}^{m-1} g_{i+1}(f(X^i)) = f(X^m) \cup \bigcup_{i=1}^m X^i \subseteq Y, \tag{51}$$

implying  $f^i(Y) \subseteq Y$  for  $0 \leq i \in \mathbb{N}$  (but not yet  $f^k(Y) \subseteq Y$  for some  $k \in \mathbb{N}$ ).

But

$$f^{m+1}(Y) = \bigcup_{i=0}^m f^{m+1}(X^i) = \bigcup_{i=0}^m \bigcup_{k=1}^{m+1-i} f^k(f(X^{m-k+1})) \subseteq \bigcup_{k=1}^{m+1} f^{m+1-k}(Y) \subseteq Y$$

by (50) and (51).  $Y$  is compact because  $Z, C$  and  $X^0$  are. Therefore application of theorem 13 completes the proof.  $\square$

Especially interesting for practical applications is the following corollary.

**Corollary 24.** Let  $Z \in \mathbb{P}VR$ ,  $C \in \mathbb{P}MR$  and  $\emptyset \neq X^0 \in \mathbb{P}VR$ , all  $Z, C$  and  $X^0$  being compact. Define  $f : \mathbb{P}VR \rightarrow \mathbb{P}VR$  by  $f(Y) := Z + C \cdot Y$  and let  $g_i : \mathbb{P}VR \rightarrow \mathbb{P}VR$  with  $V \in \mathbb{P}VR \rightarrow g_i(V) \supseteq V$  for  $i \geq 0$ . Define

$$X^{k+1} := g_{k+1}(f(X^k)) \text{ for } 0 \leq k \in \mathbb{N}.$$

If then

$$f(X^m) \subseteq \bigcup_{i=0}^m X^i \text{ for some } 0 \leq m \in \mathbb{N}, \tag{52}$$

then for every  $z \in Z$  and for every  $C \in \mathbb{C}$  holds  $\rho(C) < 1$  and there is one and only one  $\bar{X} \in \mathbb{P}VR$  with  $z + C \cdot \bar{X} = \bar{X}$ . It is  $\bar{X} = (I-C)^{-1} z$ .

Proof. Follows by replacing  $U$  by  $\bar{U}$  in the proof of theorem 23 because  $f$  is affine (see the remarks after corollary 15).  $\square$

The application of theorem 23 and corollary 24 to intervals of vectors are like shown before. Corollary 24 allows inclusions for iteration matrices  $C$  with  $\rho(C) < 1$  but  $\rho(|C|) \geq 1$ . An example is

$$C := \begin{bmatrix} -10 & -9.4 \\ 10.1 & 9.5 \end{bmatrix} \text{ with } \rho(C) = 0.6 \text{ and } \rho(|C|) \approx 19.5.$$

For  $f(X) := C \cdot X$  and  $X^0 := (1, -1)^T$  is  $f(X^0) = (-0.9, 0.6)^T$ . With an  $\epsilon$ -inflation let  $X^1 := [-0.61, -0.59], [0.59, 0.61]^T$ . Then

$$f(X^1) = \begin{bmatrix} 0.166 & 0.354 \\ -0.164 & 0.354 \end{bmatrix} \cup \begin{bmatrix} 0.366 & 0.366 \\ -0.366 & 0.366 \end{bmatrix} \cup \begin{bmatrix} 0.554 & 0.554 \\ -0.556 & 0.556 \end{bmatrix} \subseteq X^0 \cup X^1 \tag{53}$$

proving  $\rho(C) < 1$ . As follows by theorem 6,  $f(X^1) \subseteq X^0 \cup X^1$  cannot be satisfied. In practice, especially for larger  $n$ , condition (52) is difficult to verify. To compute the convex hull of the  $X^0, \dots, X^m$  is rather time consuming. On the other hand,  $\bigcup_{i=0}^m X^i$  resp.  $\bigcup_{i=0}^m X^i$  can be computed very fast. Therefore one might wish to generalize corollary 24 in the way, that (52) could be replaced by

$$f(X^m) \subseteq \bigcup_{i=0}^m X^i \text{ for some } m \in \mathbb{N}. \tag{54}$$

Unfortunately this is not true. Consider the following example:

$$C := \begin{bmatrix} 0.26 & 0.76 \\ 0.76 & 0.26 \end{bmatrix} \text{ with } \rho(C) = 1.02.$$

For  $X^{k+1} := C \cdot X^k$  and  $X^0 := (1, 0)^T$  is  $X^1 = (0.26, 0.76)^T$  and  $X^2 = (0.6452, 0.3952)^T$  with

$$X^2 \subseteq \begin{bmatrix} 0.26, 1 \\ 0, 0.76 \end{bmatrix}.$$

The example shows, that (54) does not imply  $\rho(C) < 1$ . However, the contraction is very probable if (54) is satisfied as turned out in many practical examples. Therefore, in an algorithm, condition (54) could be tested and, for a final verification,

$$f(Y) \stackrel{\Delta}{=} Y \text{ with } Y := \bigwedge_{i=0}^{\infty} \left( \bigcup_{j=0}^i X^j \right)$$

has to be checked. This method is very effective and inexpensive.

4. Nonlinear systems of equations. In this chapter theorems are developed for the inclusion of the solution of general systems of nonlinear equations, the assumptions of which, again, are verifiable on computers. For this purpose we first linearize a given function  $f: VR \rightarrow VR$ ,  $f \in C^1$  locally. For  $f': VR \rightarrow VR$  being the Jacobian of  $f$  holds

$$f'_i(x) = f'_i(\tilde{x}) + f'_i(\tilde{x} + \theta_1(x - \tilde{x})) \cdot (x - \tilde{x}) \text{ for } 1 \leq i \leq n, 0 < \theta_1 < 1, \quad (55)$$

where  $f = (f_1, \dots, f_n)^T$  and  $f'_i$  is the  $i$ -th row of the Jacobian  $f'$ . Defining  $f': VR \rightarrow VR$  by

$$f'_{ij}(x) := \{f'_{ij}(x) \mid x \in X\} \text{ for } x \in PVR \quad (56)$$

yields

$$f(x) \in f(\tilde{x}) + f'(\tilde{x} \cup x) \cdot (x - \tilde{x}) \text{ for } \tilde{x}, x \in VR. \quad (57)$$

Applying (57) to theorem 11 yields the following.

Theorem 25. Let  $f: VR \rightarrow VR$  with  $f \in C^1$ ,  $\tilde{x} \in VR$ ,  $R \in MR$  and  $\emptyset \neq X \in PVR$  compact be given. If with  $f': PVR \rightarrow PVR$  defined in (56)

$$\tilde{x} - R \cdot f(\tilde{x}) + (I - R \cdot f'(\tilde{x} \cup X)) \cdot (X - \tilde{x}) \stackrel{\Delta}{=} X, \quad (58)$$

Then  $X$  and every  $M \in MR$  with  $M \in f'(\tilde{x} \cup X)$  is not singular and there is one and only one  $x \in X$  with  $f(x) = 0$ . It is  $x \in \tilde{x}$ .

Proof. Condition (58) implies

$$x - R \cdot (f(\tilde{x}) + f'(\tilde{x} \cup X) \cdot (x - \tilde{x})) \stackrel{\Delta}{=} X \text{ for every } x \in X.$$

By (57) and the definition of  $f'$  this yields

$$\{x - R \cdot f(x) \mid x \in X\} \stackrel{\Delta}{=} X. \quad (59)$$

By Brouwer's Fixed Point Theorem this implies the existence of an  $\tilde{x} \in X$  with  $R \cdot f(\tilde{x}) = 0$ , where (59) shows  $\tilde{x} \in \tilde{x}$ . By theorem 11 the matrix  $R$  and every matrix  $M \in f'(\tilde{x} \cup X)$  is not singular implying  $f(\tilde{x}) = 0$ . Suppose  $f(y) = 0$  for  $y \in X$ . Then by the definition of  $f'$  and by (58) there is a  $M \in f'(\tilde{x} \cup X)$  with

$$f(\tilde{x}) = f(y) + M(\tilde{x} - y)$$

implying  $M \cdot (\tilde{x} - y) = 0$  and therefore  $\tilde{x} = y$ . □

It has been noted in the literature (cf. [5], [14]) that the set of Jacobians defined in (56) can be weakened. Consider the following lemma.

Lemma 26. Let  $f: VR \rightarrow VR$  with  $f \in C^1$  and  $X \in VR$ . Then for every  $x \in VR$  there are  $\theta_j \in R$ ,  $1 \leq j \leq n$  with  $0 < \theta_j < 1$  and

$$f(x) = f(\tilde{x}) + \sum_{j=1}^n \frac{\partial f}{\partial x_j}(\tilde{x}_1, \dots, \tilde{x}_{j-1}, \tilde{x}_j + \theta_j(x_j - \tilde{x}_j), x_{j+1}, \dots, x_n) \cdot (x_j - \tilde{x}_j). \quad (60)$$

Proof. Follows by straightforward calculation. □

Let  $f: VR \rightarrow VR$  with  $f \in C^1$  and  $f = (f_1, \dots, f_n)^T$ . Then for  $\tilde{x} \in VR$  and the definition of  $f': PVR \rightarrow PVR$  for  $X \in VR$  by

$$f'_{ij}(x) := \left( \frac{\partial f_i}{\partial x_j}(\tilde{x}_1, \dots, \tilde{x}_{j-1}, \tilde{x}_j \cup x_j, x_{j+1}, \dots, x_n) \mid x_k \in X_k \text{ for } j \leq k \leq n \right) \quad (61)$$

follows by lemma 26:

$$f(x) \in f(\tilde{x}) + f'(X) \cdot (x - \tilde{x}). \quad (61)$$

Note, that definition (61) depends on  $\tilde{x}$ . Formula (62) implies that  $f'$  defined by (56) can be replaced by definition (61) in theorem 25. However, the uniqueness of the zero  $\tilde{x}$  of  $f$  in  $X$  cannot be proved like in theorem 25. There, for some  $\tilde{x}, y \in X$  with  $f(\tilde{x}) = f(y) = 0$  the existence of a matrix  $M \in f'(X)$  was assumed with  $f(\tilde{x}) = f(y) + M \cdot (\tilde{x} - y)$ . This need not be true when using (61) instead of (56).

For our purposes, especially when applying Einzelschrittverfahren, another formulation of lemma 26 gives better results. Consider the following lemma.

Lemma 27. Let  $f: VR \rightarrow VR$  with  $f \in C^1$  and  $\tilde{x} \in VR$ . Then for every  $x \in VR$  there are  $\theta_{ij} \in R$ ,  $1 \leq i, j \leq n$  with  $0 < \theta_{ij} < 1$  and

and therefore

$$(x - R \cdot f(\tilde{x} + x) | x \in X) \subseteq X.$$

Therefore regarding theorem 6 the proof of theorem 25 can be applied.  $\square$

Like in the case of systems of linear equations the inclusion of the difference of an approximate solution to the correct solution yields much sharper results than an inclusion of the solution itself. Next the application of the Einzelschrittverfahren is presented combined with the technique of including the residue with respect to an approximate solution.

Theorem 29. Let  $f : VR \rightarrow VR$  with  $f \in C^1$ ,  $\tilde{x} \in VR$ ,  $R \in M(R)$  and  $x \in \Pi VR$  be given. Let  $f' : \Pi VR \rightarrow \Pi M(R)$  be defined by

$$f'_{ij}(V) := \frac{\partial f_i}{\partial x_j}(x_1, \dots, x_{j-1}, \tilde{x}_j, x_{j+1}, \dots, \tilde{x}_n) | x_i \in V_i \text{ for } 1 \leq i < j, \tilde{x}_j \in \tilde{x}_j \cup V_j \quad (66)$$

for  $V \in \Pi VR$  and  $1 \leq i, j \leq n$ . Define recursively  $Y \in \Pi VR$  by

$$Y_i := 0(-R \cdot f(\tilde{x}) + (I - R \cdot f'(\tilde{x} + Z))) \cdot Z_i \text{ for } 1 \leq i \leq n \quad (67)$$

where  $Z_i := (Y_1, \dots, Y_{i-1}, \tilde{x}_i, \dots, \tilde{x}_n)^T \in \Pi VR$ . If then

$$Y_i \subseteq X_i \text{ for } 1 \leq i \leq n, \quad (68)$$

then the matrix  $R$  and every matrix  $M \in M(R)$  with  $M \in f'(\tilde{x} + Y)$  is not singular and there is an  $\tilde{x} \in \tilde{x} + Y$  with  $f(\tilde{x}) = 0$ . Moreover with

$$W^0 := Y; W^{k+1} := 0(-R \cdot f(\tilde{x}) + (I - R \cdot f'(\tilde{x} + W^k))) \cdot W^k$$

holds

$$\tilde{x} \in \tilde{x}_0 \cup W^k.$$

Proof. By (67), the definition of the  $Z_i$  and (68) holds

$$0(-R \cdot f(\tilde{x}) + (I - R \cdot f'(\tilde{x} + Y))) \cdot X \subseteq X \text{ using Einzelschrittverfahren (69)}$$

as described in theorem 18 with the rounding 0 defined in chapter 1. Therefore by theorem 12 the matrix  $R$  and every matrix  $M \in M(R)$  with  $M \in f'(\tilde{x} + Y)$  is not singular. Furthermore by (67), (69) and theorem 12

$$f(x) = f(\tilde{x}) + M \cdot (x - \tilde{x}) \text{ with } M \in M(R) \text{ and} \quad (63)$$

$$M_{ij} := \frac{\partial f_i}{\partial x_j}(x_1, x_2, \dots, x_{j-1}, \tilde{x}_j + \theta_{ij}(x_j - \tilde{x}_j), \tilde{x}_{j+1}, \dots, \tilde{x}_n).$$

Proof. Let  $i$  with  $1 \leq i \leq n$  fixed but arbitrary. Then by the -dimensional Mean-Value Theorem

$$f_i(x_1, x_2, \dots, x_n) = f_i(x_1, \dots, x_{n-1}, \tilde{x}_n) + \frac{\partial f_i}{\partial x_n}(x_1, \dots, x_{n-1}, \tilde{x}_n + \theta_{in}(x_n - \tilde{x}_n)) \cdot (x_n - \tilde{x}_n),$$

$$f_i(x_1, \dots, x_{n-1}, \tilde{x}_n) = f_i(x_1, \dots, x_{n-2}, \tilde{x}_{n-1} + \theta_{i(n-1)} \cdot (x_{n-1} - \tilde{x}_{n-1}), \tilde{x}_n), \dots,$$

$$f_i(x_1, x_2, \dots, \tilde{x}_n) = f_i(x_1, \dots, \tilde{x}_n) + \frac{\partial f_i}{\partial x_1}(x_1 + \theta_{i1}(x_1 - \tilde{x}_1), \tilde{x}_2, \dots, \tilde{x}_n) \cdot (x_1 - \tilde{x}_1)$$

and therefore

$$f_i(x) = f_i(\tilde{x}) + \sum_{j=1}^n \frac{\partial f_i}{\partial x_j}(x_1, \dots, x_{j-1}, \tilde{x}_j + \theta_{ij}(x_j - \tilde{x}_j), \tilde{x}_{j+1}, \dots, \tilde{x}_n) \cdot (x_j - \tilde{x}_j)$$

demonstrating the lemma.  $\square$

Next the usage of the Einzelschrittverfahren for systems of nonlinear equations will be introduced. Before doing this the technique of computing an inclusion of the difference between an approximate solution and the true solution (like for linear systems theorem 18) will be introduced using definition (63) for  $f'$ .

Theorem 28. Let  $f : VR \rightarrow VR$  with  $f \in C^1$ ,  $\tilde{x} \in VR$ ,  $R \in M(R)$  and  $X \in \Pi VR$  be given. Let  $f' : \Pi VR \rightarrow \Pi M(R)$  be defined by

$$f'_{ij}(V) := \frac{\partial f_i}{\partial x_j}(x_1, \dots, x_{j-1}, \tilde{x}_j, x_{j+1}, \dots, \tilde{x}_n) | x_i \in V_i \text{ for } 1 \leq i < j, \tilde{x}_j \in \tilde{x}_j \cup V_j \quad (64)$$

for  $V \in \Pi VR$  and  $1 \leq i, j \leq n$ . If then

$$-R \cdot f(\tilde{x}) + (I - R \cdot f'(\tilde{x} + X)) \cdot X \subseteq X, \quad (65)$$

then  $R$  and every matrix  $M \in M(R)$  with  $M \in f'(\tilde{x} + X)$  is not singular and there is an  $\tilde{x} \in \tilde{x} + X$  with  $f(\tilde{x}) = 0$ .

Proof. By the definition (64) of  $f'$ , lemma 27 and (65) follows

$$(x - R \cdot (f(\tilde{x}) + f'(x + X) \cdot x) | x \in X) \subseteq X$$

$$\{x - R \cdot (f(\bar{x}) + f'(\bar{x} + Y) \cdot X) \mid x \in Y\} \subseteq Y$$

and therefore by definition (66) and lemma 27

$$\{x - R \cdot f(\bar{x} + x) \mid x \in Y\} \subseteq Y. \quad (70)$$

By Brouwer's Fixed Theorem there is an  $\hat{y} \in Y$  with  $\hat{y} - R \cdot f(\bar{x} + \hat{y}) = \hat{y}$  implying  $f(\bar{x} + \hat{y}) = 0$  by the non-singularity of  $R$ . With  $\hat{x} := \bar{x} + \hat{y}$  the proof is complete observing (69) and (70).  $\square$

5. Implementation and examples. Following some implementation hints and numerical examples will be given. Implementation details will be given for systems of linear equations; they apply for systems of nonlinear equations as well.

To calculate an inclusion of a linear system  $Ax = b$  first an approximate inverse  $R$  is required according to theorem 18. Note, that  $R$  can be replaced by LU from an LU-decomposition. The success of the algorithm, i.e. whether an inclusion will be computed or not, depends highly on  $R$ . In fact, as theorem 22 shows, an inclusion will be found if and only if  $I - RA$  is convergent. However, it is important to note, that the user does not have to know a priori, whether  $R$  or  $A$  is not singular, he does not have to know in advance, whether the spectral radius of  $I - RA$  is less than one. This will be demonstrated by the algorithm automatically a posteriori.

In theorem 18, also an approximate solution  $\bar{x}$  is required. The number of iterations (3. in (41)) necessary depends, of course, on the quality of  $\bar{x}$ . However, again no additional information on  $\bar{x}$  such as  $\|\bar{x} - \bar{x}\|$  is required.

Suppose,  $R$  and  $\bar{x}$  are given (e.g. computed using some traditional method, cf. [11], [42], [43]), then an algorithm based on theorem 17, on the third iteration in (41) and the  $\varepsilon$ -inflation would be the following:

```
X := 0(b - A\bar{x}); Z := R \cdot X; C := 0(I - R \cdot A); X := Z; k := 0;
repeat k := k + 1; Y := X \oplus e; X := Z \oplus C \oplus Y until X \subseteq Y or k > 15;
```

Algorithm 1. Traditional way.

$0 : PVR + I \cdot VG$  resp.  $0 : PMR + I \cdot MG$  is any isotone rounding, i.e.  $A \subseteq 0A$ . As has been pointed out before, it is of utmost importance to compute the residuals  $b - Ax$  and  $I - RA$

with one rounding. Preferably, a precise scalar product with maximum accuracy introducing only one rounding error is used (cf. [7], [8]). In this case the rounding  $0$  is best possible and therefore equal to  $\diamond$ .  $X, Y$  and  $Z$  are interval vectors and  $C$  is an interval matrix; therefore the additional storage needed is  $5n + 2n^2 + 0(1)$ . Next, the Einzeilschrittverfahren can be used according to theorem 18 yielding the following algorithm.

```
X := 0(b - A\bar{x}); Z := R \cdot X; C := 0(I - R \cdot A); X := Z; k := 0;
repeat k := k + 1; X := X \oplus e; incl := true;
  for i := 1 to n do {Q := Z_i \oplus C_i \cdot X; if Q \not\subseteq X_i then incl := false;
    X_i := A}
until incl or k > 15;
```

Algorithm 2. Einzelschrittverfahren.

Here  $C_i$  denotes the  $i$ -th row of  $C$ ,  $X_i$  the  $i$ -th component of  $X$ .  $Q$  is an interval and the additional storage needed is  $4n + 2n^2 + 0(1)$ .

In order to avoid the additional  $O(n^2)$  storage consider the following. When using interval vectors as subsets of  $VR$  it has been shown in theorem 12 that

$$Z + C \cdot X \subseteq X \quad \text{for } X, Z \in II \cdot VR; C \in II \cdot MR \quad (71)$$

already implies  $\rho(|C|) < 1$  for every matrix  $C \in \mathcal{E}$ . Therefore, instead of (71),

$$|Z| + |C| \cdot |X| < |X| \quad (72)$$

can be checked. Necessary and sufficient conditions for an iteration using (72) to stop have been given in [39].

A significant amount of storage can be saved because only the absolute value of  $X, Z$  and  $C$  are needed. Therefore the additional storage needed essentially halves to  $3n + n^2 + 0(1)$ . Moreover, now the matrix  $C$ , which is a point matrix, can be computed in the storage of  $R$ .

This is possible because  $C := I - RA$  can be computed rowwise, storing the intermediate (row) result somewhere. After the first row of  $C$  has been computed, the first row of  $R$  is no longer needed. This leads to the following algorithm.



[x,y] := 0(b - Ax); z := Δ(|R·[x,y]|);  
 for i := 1 to n do (x := R<sub>i</sub>; R<sub>i</sub> := Δ(|I<sub>i</sub> - x·A|)); x := z; k := 0;  
 repeat k := k + 1; x := x Δ ε; incl := true;  
 for i := 1 to n do (q := z<sub>i</sub> Δ R<sub>i</sub> Δ x; if q ≥ x<sub>i</sub> then incl := false;  
 x<sub>i</sub> := q)  
 until incl or k > 15;

Algorithm 3. Checking (72) with Einzelschrittverfahren

Algorithms 1 and 2 imply  $R \in \tilde{X} \otimes X$  ( $R := A^{-1} \cdot b$ ) whereas algorithm 3 implies  $R \in \tilde{X} \otimes [-x, x]$ . The validity of algorithm 3 follows by

$$\begin{aligned} |z_i + \epsilon| \cdot |x| < |x| &\Rightarrow -|x| < -|z_i - \epsilon| \cdot |x| \leq |z_i + \epsilon| \cdot |x| < |x| \Rightarrow \\ &\Rightarrow (-|z_i|, |z_i|) + (-|\epsilon|, |\epsilon|) \cdot (-|x|, |x|) \subseteq [-|x|, |x|] \Rightarrow \\ &\Rightarrow z + \epsilon \cdot (-|x|, |x|) \subseteq [-|x|, |x|]. \end{aligned}$$

All three algorithms verify the non-singularity of R and A and therefore the unique solvability of the linear system Ax = b. The additional storage required for algorithms 1 to 3 is:

algorithm	1	2	3
storage	$2n^2 + 6n + O(1)$	$2n^2 + 4n + O(1)$	$3n + O(1)$

For algorithm 3 the computing reduces as well.

The non-singularity of R and A could also be demonstrated by showing  $\|I - RA\| < 1$  for some norm  $\|\cdot\|$ :  $MR \rightarrow R$ . On the computer,  $\|\Delta(\|I - RA\|)\| < 1$  would be verified. In the following some numerical examples are shown where  $\|I - RA\| \geq 1$  for a number of norms but nevertheless an algorithm based on the inclusion theory does verify the contraction and delivers sharp bounds for the solution of the linear system.

The following examples are calculated on an IBM S/370 in single precision (~ 7.5 decimal digits) and double precision (~ 16.5 decimal digits). The approximate inverse X is computed using Gauss-Jordan algorithm. As examples consider (n is the number of rows)

Hilbert\* - matrices  $H_{ij}^* := \frac{1}{\text{lcm}(1, 2, \dots, 2n-1)}$

Pascal - matrices  $P_{ij} := \binom{i+j}{j}$

Pascal\* - matrices  $P_{ij}^* := \binom{i+j-1}{j}$

Zielke - matrices  $Z_{ij} := \frac{\binom{n+i-1}{i-1} \cdot n \cdot \binom{n-1}{n-i}}{i+j-1}$  (cf. [50])

$S_{ij}(q) := 1 - q \cdot r_{ij}$  where  $r_{ij} \in [0, 1]$  randomly.

All matrices except the last have integer entries. Where the precision in use does not suffice to store an entry precisely, the entry rounded to nearest is used.

In the following tables the matrix type and the dimension n is listed. In the column  $\|I - RA\|$  the minimum value of sum, maximum and Frobenius-norm is listed. k is the number of interval iterations and "verified digits" is the minimum number of decimal digits coinciding of the bounds for the components of the solution using algorithm 3 listed above. In all cases a linear system with the depicted matrix and right hand side (1, ..., 1)T is solved. First the single precision results are displayed.

matrix	n	$\ I - RA\ $	k	verified digits
H	7	1.7	2	7.5
P	8	1.2	1	7.5
	9	38	2	7.3
P*	9	3.5	2	7.4
Z	7	1.6	3	7.3
$S(10^{-5})$	25	0.74	2	7.4

Table 1. Single precision results (~ 7.5 decimal digits)

As can be seen the number of verified decimal digits is widely independent on the condition of the matrix. This is due to the technique not to compute an inclusion of the solution itself but of the difference of the true solution X to an approximate solution X. Furthermore, the convergence of I - RA can be shown even if norm estimates fail to show  $\rho(I - RA) < 1$ . In the next table there are more extreme examples in this respect.

matrix	n	$\ I-RA\ $	k	verified digits
P	20	11	1	16.5
	22	210	1	16.5
	24	670	1	16.5
	26	280000	1	16.5
	$S(10^{-3})$	50	0.02	1
	100	0.03	1	16.3
	200	0.49	2	16.4

Table 2. Double precision results (~16.5 decimal digits)

The Pascal  $26 \times 26$  matrix shows the extreme example where the spectral radius of  $I-RA$  is estimated by 280000 where in fact the new methods show  $\rho(I-RA) < 1$ . Also for larger dimensions the inclusions behave very stable.

5. Conclusion. In the preceding chapters the theoretical foundations of the inclusion theory were extended and new proofs were given without using Brouwer's Fixed Point Theorem. Some hints on the implementation were also provided. In [38] the theoretical background and corresponding algorithms are given for other numerical problems such as linear systems with band matrix, symmetric matrix or sparse matrix, for over- and underdetermined linear systems, evaluation and zeros of polynomials, algebraic eigenvalue problems, linear, quadratic and convex programming problems, evaluation of arithmetic expressions and others. Algorithms corresponding to a number of those problems are implemented in the IBM Program Product ACRITH, which is available since March 1984 and with a second Release since early 1985 (cf[51]).

The key property of the new algorithms is that the verification of the validity of the result is performed automatically by the computer without any effort on the part of the user. The verification includes the existence and uniqueness of a solution within the computed bounds. The input data may be real point or interval data as well as complex point or interval data. Especially if the data is afflicted with tolerances the verification process is of great help. In this case it is

verified that any problem within the tolerances is solvable and the solution of any of the (infinitely many) problems within the tolerances is enclosed within the calculated inclusion interval.

The computing time is of the order of a comparable floating-point algorithm (e.g. Gaussian elimination in case of general linear systems with full matrix), the latter, of course, without the verification of the result.

The computed bounds are of high accuracy, i.e. the difference of the left and right bound of the inclusion of every component is of the order of the relative rounding error unit. By our experience, very often the inclusions are of least significant bit accuracy, i.e. the left and right bound of the inclusion of every component are adjacent floating-point numbers.

#### 7. References

- [1] Abbott, J.P., Brent, R.P. (1975). Fast Local Convergence with Single and Multistep Methods for Nonlinear Equations, *Austr. Math. Soc.* 19 (series B), 173-199.
- [2] Alefeld, G., Intervallrechnung Über den Komplexen Zahlen und einige Anwendungen. Dissertation, Universität Karlsruhe, 1968.
- [3] Alefeld, G. and Herzberger, J., "Einführung in die Intervallrechnung". Reihe Informatik, 12. Wissenschaftsverlag des Bibliographischen Instituts Mannheim, 1974.
- [4] Alefeld, G. and Herzberger, J., "Introduction to Interval Analysis", Academic Press, New York (1982).
- [5] Alefeld, G. (1979). Intervallanalytische Methoden bei nicht-linearen Gleichungen. In "Jahrbuch Überblicke Mathematik 1979", B.I. Verlag, Zürich.
- [6] Bauer, F.L. and Samelson, K. Optimale Rechengenauigkeit bei Rechenanlagen mit gleitendem Komma, *Z. Angew. Math. Phys.* 4, 312-316 (1953).
- [7] Bohlander, G., Floating-point computation of functions with maximum accuracy. *IEEE Trans. Comput.* C-26, No. 7, 621-632 (1977).

- [35] Rall, L.B. (1981). Mean value and Taylor forms in interval analysis, *SIAM J. Math. Anal.* 14, No. 2 (1983).
- [36] Reinsch, Ch., Die Behandlung von Rundungsfehlern in der Numerischen Analysis, "Jahrbuch Oberblicke Mathematik 1979", Wissenschaftsverlag des Bibliographischen Instituts Mannheim, 43-62 (1979).
- [37] Rump, S.M. (1980). Kleine Fehlerschranken bei Matrixproblemen, Dissertation, Universität Karlsruhe.
- [38] Rump, S.M. (1983). Solving Algebraic Problems with High Accuracy, Habilitationsschrift, in Kullisch/Miranker: A New Approach to Scientific Computation, Academic Press, New York.
- [39] Rump, S.M. (1982). Solving Non-linear Systems with Least Significant Bit Accuracy, *Computing* 29, 183-200.
- [40] Rump, S.M. (1984). Solution of Linear and Nonlinear Algebraic Problems with Sharp, Guaranteed Bounds, *Computing Suppl.* 5, 147-168.
- [41] Rump, S.M. and Kaucher, E., Small bounds for the solution of systems of linear equations, *Computing Suppl.* 2, 157-164 (1980).
- [42] Stoer, J. (1972). Einführung in die Numerische Mathematik I. Heidelberg Taschenbücher, Band 105, Springer-Verlag, Berlin-Heidelberg-New York.
- [43] Stoer, J., Bulirsch, R. (1973). Einführung in die Numerische Mathematik II. Heidelberg Taschenbücher, Band 114, Springer-Verlag, Berlin-Heidelberg-New York.
- [44] Ullrich, Ch., Zur Konstruktion komplexer Kreisarithmetiken *Computing Suppl.* 1, 135-150 (1977).
- [45] Varga, R.S. (1962). *Matrix Iterative Analysis*. Prentice-Hall, Englewood Cliffs, New Jersey.
- [46] Walter, W. (1970). *Differential and Integral Inequalities*. Berlin-Heidelberg-New York: Springer.

- [47] Wilkinson, J.H., "Rounding Errors in Algebraic Processes", Prentice-Hall, Englewood Cliffs, New Jersey, 1963.
- [48] Wongwiss, P. Experimentelle Untersuchungen zur numerischen Auflösung von linearen Gleichungssystemen mit Fehlerfassung, Interner Bericht 75/1, Institut für Praktische Mathematik, Universität Karlsruhe.
- [49] Yohe, J.M., Roundings in floating-point arithmetic, *IEEE Trans. Comput.* C.12 No. 6, 577-586 (1973).
- [50] Zielke, R., *Algol-Katalog Matrizenrechnung*, Oldenburg Verlag, München, Wien (1972).
- [51] ACRITH High-Accuracy Arithmetic Subroutine Library: General Information Manual, IBM Publications, GC33-6163, (1985).

Acknowledgement: The author wants to thank his students of the summer lecture 1985 for several helpful comments.

Address of the author:

Priv.-Doz. Dr. Siegfried M. Rump  
IBM Development and Research  
Schönaicher-Straße 220  
D-7030 Böblingen  
Federal Republic of Germany