

Günter Schumacher

Universität Karlsruhe  
 Institut für Angewandte Mathematik  
 Kaiserstr. 12, D-7500 Karlsruhe  
 West-Germany

**Abstract**

Interval arithmetic, computation with numbers which are affected by tolerances, always provides reliable results when applied in numerical algorithms. It guarantees that the exact result of an algorithm lies within the computed tolerance bounds. In ill-conditioned cases these bounds may be very wide and although the statements which have been recently introduced provide the possibility of successive refining the tolerance. Furthermore, the existence and uniqueness of the solution is proved (in a mathematical sense) by the computer.

These methods combine concepts of interval analysis with the computer arithmetic defined by Kulisch and Miranker. They are based on fixed point theorems and the tensor product is essential for their implementation.

**0. Introduction**

In [8] a mathematical approach has been laid out what we call "Computer Arithmetic". For example, the algebraic structure of a digital computer may be described in terms of "floating point numbers", and the process of truncating a real number to its machine representation may be seen to be a mapping from the real numbers to the set of floating-point numbers.

We will summarize some of the basic concepts:

- $\mathbb{R}$  ..... set of real numbers
- $\mathbb{S}$  ..... set of floating point numbers of a given digit length
- $V_n T$  ..... vector space of dimension  $n$  over a set  $T \in \{\mathbb{R}, \mathbb{S}\}$
- $M_n T$  ..... square matrices of  $n$  rows and components in  $T$
- $\mathbb{P}T$  ..... power set of  $T$  ( $\mathbb{P}\mathbb{R}, V_n \mathbb{R}, M_n \mathbb{R}, M_n \mathbb{S}$ )
- $IT$  ..... set of intervals  $A = [a, b]$  over an ordered set  $T \in \{\mathbb{R}, \mathbb{S}, V_n \mathbb{R}, M_n \mathbb{R}, M_n \mathbb{S}\}$ , i.e. the set of all  $[a, b]$  with  $a \leq x \leq b$  (For some special properties of interval arithmetic see [1],[6]).

Occasionally, we will denote a member of an interval  $A$  by  $\hat{A}$ .

For two intervals  $A, B \in T$  an arithmetic operation  $\circ \in \{+, -, \cdot, /\}$  is defined by

$$A \circ B := \{a \circ b \mid a \in A, b \in B\}$$

$\cup$  ..... convex union operation for two intervals  $A, B \in IS$ :

$$A \cup B := \{x \in S \mid \exists a \in A, b \in B : a \leq x \leq b \vee b \leq x \leq a\}$$

$d(\cdot)$  ..... diameter of an interval, i.e.

$$A = [a, b] \in IS \implies d(A) := b - a$$

$m(\cdot)$  ..... midpoint of an interval

(in our context  $m(A)$  may take any point out of an interval  $A \in IS$ )

$\diamond : \mathbb{P}\mathbb{R} \rightarrow IS$  ..... rounding from  $\mathbb{P}\mathbb{R}$  into  $IS$  with the property

$$\bigwedge_{A \in \mathbb{P}\mathbb{R}} \diamond A = \cap \{X \in IS \mid A \subset X\}$$

For vectors and matrices this definition applies componentwise

$\diamond, \diamond, \dots$  interval operations in  $IS, IV_n S$  and  $IM_n S$  defined for  $T \in \{\mathbb{S}, V_n \mathbb{S}, M_n \mathbb{S}\}$  by

$$\bigwedge_{A, B \in IT} \bigwedge_{\circ \in \{+, -, \cdot, /\}} A \diamond B := \diamond(A \circ B)$$

Throughout this paper we consider the following problem:

"Let  $f : D \rightarrow V_n \mathbb{R}$  be a continuously differentiable function,  $D \subset V_n \mathbb{R}$ . We seek an  $\hat{x} \in V_n \mathbb{R}$  with  $f(\hat{x}) = 0$ , or more practically, we ask for an inclusion  $X \in IV_n \mathbb{R}$  of  $\hat{x}$ , where the diameter  $d(X)$  is sufficiently small."

We restrict ourselves to the case where the components of  $f$  are arithmetic expression with operations  $+, -, \cdot, /$  and integer exponentiation. The algorithm to be introduced determines an inclusion  $X$  of the solution  $\hat{x}$  and performs an automatic verification of conditions such as the existence and uniqueness of a zero  $\hat{x}$  of  $f$  within  $X$ . These methods are called E-methods (E is the first letter of the three German words "Existenz" for existence, "Eindeutigkeit" for uniqueness, and "Einschließung" for inclusion).

<sup>†</sup>This paper is an extended version of [4].

For this purpose, the problem of finding a zero  $\hat{x}$  of  $f$  is transformed into a fixed point equation  $g(x) = x$  (see [7], [11], [12]). An inclusion of the solution of this fixed-point problem is computed iteratively in the space  $IV_n\mathbb{R}$  starting with an approximate solution  $\tilde{x} \in V_n\mathbb{R}$ . By use of residual correction techniques the diameter of the resulting interval is diminished more and more.

### 1. Theoretical foundations

We summarize some theorems from [7], [11], [12] concerning the existence and uniqueness of a solution of a system of equations. At first, we need the following

**Lemma 1:** (Schauder's fixed point theorem)

Let  $f : X \rightarrow V_n\mathbb{R}$  be a continuous function on  $X \subseteq V_n\mathbb{R}$ ,  $X$  nonempty, convex and compact. If  $f(X) \subseteq X$ , then, the equation  $f(x) = x$  has at least one solution  $\hat{x}$  in  $X$ .

**Proof:** see [5], e.g.  $\square$

From this lemma we have immediately

**Theorem 1:** Let  $f : X \rightarrow V_n\mathbb{R}$  be a continuous function on  $X \in IV_n\mathbb{R}$ , and  $F : PX \rightarrow PV_n\mathbb{R}$  with

$$\bigwedge_{A \in PX} \bigwedge_{x \in A} f(x) \in F(A). \quad (1.1)$$

If  $F(X) \subseteq X$ , then the equation  $f(x) = x$  has at least one solution  $\hat{x}$  in  $X$  and

$$\bigwedge_{k \geq 0} \hat{x} \in F^k(X),$$

where  $F^0(X) := X$ ,  $F^k(X) := F(F^{k-1}(X))$ .

**Proof:** From  $X \in IV_n\mathbb{R}$  we have:  $X \neq \emptyset$ ,  $X$  compact and convex. From the definition we have  $f(X) := \{f(x) \mid x \in X\} \subseteq F(X) \subseteq X$ . Therefore, Lemma 1 delivers the existence of a solution  $\hat{x} \in X$  of  $f(x) = x$ . The rest is shown by induction.  $\square$

$F$  may be arbitrarily chosen to satisfy (1.1). In practical applications  $F$  will be any interval extension of  $f$ .

Theorem 1 allows no conclusion about the uniqueness of the solution. For this purpose we define

$$\bigwedge_{A, B \in IV} A \subsetneq B : \Leftrightarrow A \subseteq B \wedge A \neq B$$

We also need the following two lemmata from [15]:

**Lemma 2:** Let  $A = ((a_{ij})) \in M_n\mathbb{R}$ ,  $x = (x_i) \in V_n\mathbb{R}$  and  $x > 0$ . Then for the spectral radius,

$$\max_{1 \leq i \leq n} \frac{\sum_{j=1}^n |a_{ij}| x_j}{x_i}$$

holds

**Lemma 3:** Let  $|B|$  and  $0 \leq |B| \leq A$ . Then,  $\rho(A)$ .

From this we

**Theorem 2:** Let  $C \in IM_n\mathbb{R}$ ,

$$X \in IV_n\mathbb{R}, \quad (1.2)$$

with  $f$  as defined componentwise.

Then

$$(C) < 1.$$

**Proof:** Observe that  $d(X) > 0$ . Furthermore, we have

$$d(Z + C \cdot X) = d(Z + C \cdot X) < d(X)$$

and

$$d(X).$$

From

$$d(C \cdot X) \geq |C| \cdot d(X) < d(X). \quad (1.3)$$

Let

$d(X)$ . Then by lemma 2

$$\frac{y_i}{d(X_i)} < \frac{d(X_i)}{d(X_i)} = 1. \quad (1.3)$$

Then

$\rho(C)$  we have

$$\rho(|C|) \leq \rho(|C|) < 1$$

by

lemma 3 twice.  $\square$

Now we can prove the following constructive theorem which is a constructive proof of the existence and uniqueness of a solution of a system of equations.

**Theorem 3:** Let  $X \in IV_n\mathbb{R}$ ,  $f : X \cup \tilde{x} \rightarrow V_n\mathbb{R}$ ,  $f$  continuous. For arbitrary but fixed  $Y \in IV_n\mathbb{R}$  let  $J(Y) \in IM_n\mathbb{R}$  and assume that  $Y$  is a very set of  $n$  vectors  $y_1, \dots, y_n$ .

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_n}(y_1) \\ \vdots \\ \frac{\partial f_n}{\partial x_n}(y_n) \end{pmatrix} \in J(Y)$$

holds

$f : PX \rightarrow PV_n\mathbb{R}$  be defined by

$G(\tilde{x}) = \tilde{x} + (I - R \cdot J(\tilde{x} \cup Y)) \cdot (Y - \tilde{x})$  for  $\tilde{x} \in X$ , where  $R \in M_n\mathbb{R}$  is an arbitrary fixed matrix.

If  $G(Y) \subseteq X$  (1.4) holds, the equation  $f(x) = 0$  has one and only one solution  $\hat{x}$  in  $X$  and  $\hat{x} \in G^k(Y)$ .

**Proof:** Consider the continuous function  $g : X \cup X \rightarrow X$  defined by  $g(x) = -R \cdot f(x)$ . From the Banach contraction theorem we have that for  $Y \subseteq X$  and  $y \in Y$  there exists a unique fixed point  $\hat{x}$  in  $X$  such that  $g(\hat{x}) = \hat{x}$ . This implies  $(\hat{x} \cup Y) \cdot (X - \hat{x}) \subseteq X - \hat{x}$ . Hence,  $G(Y) \subseteq X - \hat{x}$  implies  $\rho(I - R \cdot M) < 1$  for all  $Y \subseteq X$ . When considering the eigenvalues of  $I - R \cdot M$  we see that  $R$  and  $M$  are contractions. With  $\hat{x} - R \cdot f(\hat{x}) = \hat{x}$ , it follows that  $f(\hat{x}) = 0$ . Let  $\hat{y} \in Y$  be a further fixed point. The generalized mean value theorem implies the existence of  $\hat{z} \in J(\hat{x} \cup X)$  with  $\hat{z}(\hat{y} - \hat{x}) = \hat{y} - \hat{x}$ . But  $\hat{z} \in J(\hat{x} \cup X)$  and hence  $\hat{y} = \hat{x}$ .  $\square$

In practical applications the matrix  $R$  will be an approximation of the inverse of the Jacobian  $J(\hat{x})$ . Theorem 3 has been used for application on computers. In this approach the goal of computer application is that for any  $G^k(Y) \subseteq X$  the set  $G^k(Y)$  satisfies  $G^k(Y) \subseteq G^k(Y)$ .

It is obvious that  $G(Y) \subseteq Y$  implies  $G(Y) \subseteq Y$ .

This fact allows us to assert that theorem 3 remains valid in the interval arithmetic approach. The definition of  $G$ , i.e.

$$G^*(Y) := \bigcap_{x \in V_n \mathbb{R}} \{x \mid f(x) = 0\} \quad (1.5)$$

$\diamond f$  denotes any function  $f : V_n \mathbb{R} \rightarrow IV_n \mathbb{S}$  with

$$\bigwedge_{x \in V_n \mathbb{R}} f(x) \subseteq f(x)$$

The definition of the diamond product

$$a \diamond b := \{a \cdot b\}$$

requires an exact summation algorithm as stated in [8]. With such a tool the critical residual term  $\diamond(I - R \cdot J(\hat{x} \cup Y))$  may be evaluated with only one rounding in each component.

Next, a theorem can be formulated directly applicable on computers. For a better understanding of the following corollary it should be remarked that, in practical application, the inclusion of the absolute error  $\hat{x} - \tilde{x}$  of an approximate solution  $\tilde{x}$  leads to better results than an inclusion for  $\hat{x}$  itself. This is summarized in the following

**Corollary:** Let  $\tilde{x} \in V_n \mathbb{S}$ ,  $Z \in IV_n \mathbb{S}$ , and  $f \in C^1(\tilde{x} \diamond Z \cup \tilde{x})$ . For every arbitrary but fixed  $Y \subseteq (\tilde{x} \diamond Z) \cup \tilde{x}$ , let  $J(Y) \in IM_n \mathbb{S}$  and assume that for every set of  $n$  vectors  $y_1, \dots, y_n \in Y$ ,

$$\begin{pmatrix} \frac{\partial f_1}{\partial x_1}(y_1) & \dots & \frac{\partial f_1}{\partial x_n}(y_1) \\ \vdots & & \vdots \\ \frac{\partial f_n}{\partial x_1}(y_n) & \dots & \frac{\partial f_n}{\partial x_n}(y_n) \end{pmatrix} \in J(Y)$$

holds. Let  $H : \mathbb{P}Z \rightarrow IV_n \mathbb{S}$  be defined as

$$H(Y) := -R \diamond \diamond f(\tilde{x}) \diamond \diamond (I - R \cdot J(\tilde{x} \cup (\tilde{x} \diamond Y))) \diamond \diamond Y \quad (1.5)$$

for  $Y \in \mathbb{P}Z$ , where  $R \in M_n \mathbb{S}$  is an arbitrary but fixed matrix. If

$$H(Z) \subseteq Z, \quad (1.6)$$

then the equation  $f(x) = 0$  has one and only one solution  $\hat{x}$  in  $\tilde{x} \diamond Z$  and  $\hat{x} - \tilde{x} \in H^{(k)}(Z)$ ,  $k \geq 0$ .

## 2. Equation solving with error control

The corollary provides the second of the two basic steps in computing verified bounds for the solution of a system of equations:

- 1) an approximation step to determine a sufficiently good estimate  $\tilde{x}$ .
- 2) an inclusion step to determine the enclosure  $Z$  of the absolute error  $\hat{x} - \tilde{x}$  with respect to  $\tilde{x}$ .

In this chapter we will discuss some implementation aspects of these two points. Let us start with the inclusion step.

## 2.1 Implementation of the inclusion step

According to [11] we define an  $\epsilon$ -inflation for an interval  $A \in IS$  by

$$A^{\circ\epsilon} := \begin{cases} A + [-1, 1] \cdot \epsilon \cdot d(A) & \text{for } d(A) \neq 0 \\ A + [-\eta, +\eta] & \text{for } d(A) = 0 \end{cases}$$

Here  $\eta$  is the smallest positive floating-point number of the computer in use. For interval vectors this definition applies componentwise.  $\epsilon$ -inflation is indispensable to assure the convergence of the interval iteration. In practice, 0.1 turned out to be a good value for  $\epsilon$ .

Using the corollary, a verified inclusion may be obtained by performing the following iteration in the space  $IV_nS$ :

```

Y := 0 ; count := 0 ;
repeat
  count := count + 1 ;
  Z := Y  $\circ$   $\epsilon$  ;
  Y := -R  $\diamond$   $\diamond$  f(x)  $\diamond$ 
       $\diamond$  (I-R  $\cdot$  J(x  $\cup$  (x  $\diamond$  Z)))  $\diamond$  Z
until ( Y  $\subsetneq$  Z ) or ( count = count_max );

```

If the termination criterion  $Y \subsetneq Z$  is satisfied, then by corollary 1 it is verified that the exact solution  $\hat{x}$  lies in  $\tilde{x} \diamond Z$  and that no other solution lies between these bounds. (The integer variable count is used to prevent infinite looping in cases where the iteration is not convergent).

## 2.2 Performing the approximation step

There are a large number of different methods for obtaining an approximation for  $\hat{x}$  (see e.g. [9]). In the following any of those principles may be used. Here we describe only the "classical" Newton's method

$$x^{(k+1)} := x^{(k)} - (f'(x^{(k)}))^{-1} \cdot f(x^{(k)}) \quad (2.6)$$

since the problem of evaluating the function  $f$  at a point  $x^{(k)}$  with high accuracy is common to all methods. If it is not possible to compute the residual term  $f(x^{(k)})$  accurately, we have no measure how near we are to the exact solution.

As stated above, the components  $f_i$  of  $f$  are considered as arithmetic expressions. The computation of  $f_i(x)$  is usually done on a computer step by step, i.e. using an algorithm which performs one operation in the formula after the other. The results of the operations are stored in some intermediate variables  $z_1, \dots, z_n$ , say.

Example: The

$$z_1 := (x_3 + x_4) / (x_5 - x_6)$$

will

$$z_2 := \dots$$

$$z_3 := \dots$$

$$z_4 := \dots$$

$$z_5 := \dots$$

$$z_6 := \dots$$

using the results  $z_1, \dots, z_5$ .

We will state a way to arrive at a highly accurate approximation of the function values:

The computation of the function values may be considered as a successive evaluation of a system of nonlinear equations for  $z_1, \dots, z_n$ . This system of equations has a triangular form. Hence, for the computation of  $z_n = f_n(x)$  with a guaranteed accuracy we consider the following special type of system:

$$\begin{aligned} g_1(z_1) &= 0 \\ g_2(z_1, z_2) &= 0 \\ &\vdots \\ g_n(z_1, \dots, z_n) &= 0 \end{aligned}$$

(In [2] and [3] a nonlinear system was transformed into a system of linear equations by using term rewriting techniques. The resulting system also has triangular form and can be solved with the principles of linear algebra.)

We assume  $g_k$  differentiable in terms of  $z_1, \dots, z_k$ . Let  $\hat{z}_1, \dots, \hat{z}_k$  be the approximations and  $\tilde{z}_1, \dots, \tilde{z}_k$  the exact solutions.

We define

$$\Delta z_j := \hat{z}_j - \tilde{z}_j, \quad j = 1, \dots, n$$

and seek an increment  $\Delta z_j$  of the correction

$$\Delta z_j \in \Delta, \quad j = 1, \dots, n$$

Applying the mean value theorem to  $g_k$  with respect to the last variable  $z_k$  in the interval  $\hat{z}_k \cup \tilde{z}_k$  yields

$$g_k(\hat{z}_1, \dots, \hat{z}_k) = g_k(\hat{z}_1, \dots, \hat{z}_k)$$

Continuing this procedure results in

$$g_k(\hat{z}_1, \dots, \hat{z}_k) = g_k(\hat{z}_1, \dots, \hat{z}_k) + \sum_{j=1}^k \frac{\partial g_k}{\partial z_j}(\hat{z}_1, \dots, \hat{z}_k) \Delta z_j$$

Since  $\hat{z}_1, \dots, \hat{z}_k$  is an exact solution we have

$$g_k(\hat{z}_1, \dots, \hat{z}_k) = 0$$

$$\Delta z_k = \left[ -g_k(\hat{z}_1, \dots, \hat{z}_k) - \sum_{j=1}^{k-1} \frac{\partial g_k}{\partial z_j}(\hat{z}_1, \dots, \hat{z}_k) \Delta z_j \right] / \frac{\partial g_k}{\partial z_k}(\hat{z}_1, \dots, \hat{z}_k)$$

under the assumption that

Via induction we show

$$\Delta z_k \in \Delta Z_k := [-g_k(\hat{z}_1, \dots, \hat{z}_k) - \sum_{j=1}^{k-1} \frac{\partial g_k}{\partial z_j}(\hat{z}_1, \dots, \hat{z}_k) \Delta z_j, \frac{\partial g_k}{\partial z_k}(\hat{z}_1, \dots, \hat{z}_k) \Delta z_k]$$

$$- \sum_{j=1}^{k-1} \frac{\partial g_k}{\partial z_j}(\hat{z}_1, \dots, \hat{z}_k) \Delta z_j$$

We must also ensure that the denominator does not contain zero. This follows from the condition from the induction theorem that

$$\frac{\partial g_k}{\partial z_k}(\hat{z}_1, \dots, \hat{z}_k) \neq 0$$

to guarantee that the denominator is not zero in some neighborhood of  $(\hat{z}_1, \dots, \hat{z}_k)$ .

It is necessary to check the following:

- To use the induction theorem to compute bounds for  $\Delta z_k$  we must ensure that the denominator does not depend on  $\Delta z_k$  (this is always true for all kinds of functions arising from arithmetic expressions and the example above).
- The denominator  $\frac{\partial g_k}{\partial z_k}(\hat{z}_1, \dots, \hat{z}_k)$  must be of

$$g_k(\hat{z}_1, \hat{z}_2, \dots, \hat{z}_k) + \Delta z_k$$

$$\text{with } \zeta_k \in \hat{z}_k \cup \tilde{z}_k$$

$$g_k(\hat{z}_1, \hat{z}_2, \dots, \hat{z}_k) + \Delta z_k$$

$$g_k(\hat{z}_1, \hat{z}_2, \dots, \hat{z}_k) + \Delta z_k$$

act solution we have

$$g_k(\hat{z}_1, \dots, \hat{z}_k) = 0$$

$$g_k(\hat{z}_1, \dots, \hat{z}_k) \neq 0$$

show

$$g_k(\hat{z}_1, \dots, \hat{z}_k) + \Delta z_k$$

$$g_k(\hat{z}_1, \dots, \hat{z}_k) + \Delta z_k$$

denominator does not contain zero. This follows from the induction theorem that

$$\frac{\partial g_k}{\partial z_k}(\hat{z}_1, \dots, \hat{z}_k) \neq 0$$

to guarantee that the denominator is not zero in some neighborhood of  $(\hat{z}_1, \dots, \hat{z}_k)$ .

check the following:

- To use the induction theorem to compute bounds for  $\Delta z_k$  we must ensure that the denominator does not depend on  $\Delta z_k$  (this is always true for all kinds of functions arising from arithmetic expressions and the example above).
- The denominator  $\frac{\partial g_k}{\partial z_k}(\hat{z}_1, \dots, \hat{z}_k)$  must be of

high accuracy to guarantee a sufficiently sharp inclusion. This is also possible for systems arising from arithmetic expressions by applying the exact scalar product as developed in [8]. For this purpose a "division equation"  $z_k = z_1 / z_j$  has to be transformed into  $z_k \cdot z_j = z_1$ .

In most cases it is not necessary to compute the value of the partial derivatives with the same effort as the function value since the success of (2.6) mainly depends on the accuracy of  $f(x^{(k)})$ . Nevertheless, it is also possible to determine the partial derivatives with comparable accuracy. (For the technique of computing the partial derivatives of a function we refer to [10]).

With the methods for linear systems from [7], each Newton step may be performed with high accuracy since each step is the solution of a system of linear equations

$$f'(x^{(k)}) \cdot (x^{(k)} - x^{(k+1)}) = f(x^{(k)})$$

The methods from [7] are specializations of the method described here, and so this way of performing a Newton step may be considered as a recursive call of the procedure in this paper. The possibility to do that is important in the cases where the Jacobian is ill-conditioned.

### 2.3 An algorithm for the solution of a system of nonlinear equation with maximum accuracy

We use the following notations:  $\hat{x}$  denotes the exact solution of the system of nonlinear equations  $f(x) = 0$ ,  $b$  is the base of the floating-point system in use,  $t$  the length of the mantissa.

The approximation of  $\hat{x}$  in the  $k$ -th step is  $\sum_{i=0}^k \tilde{x}^{(i)}$

where  $\tilde{x}^{(i)} \in S$ ,  $i=1, \dots, n$ . This form of our approximation  $\tilde{x}$  is usually called a "staggered correction form" of  $\hat{x}$  since all values  $\tilde{x}^{(i)}$  are actually stored and the summation will be done when we have finished. To this approximation we compute an inclusion  $Z$  of the next correction

$$x^{(k)} := \hat{x} - \sum_{i=0}^{k-1} \tilde{x}^{(i)}$$

In step (N3) of this algorithm, least bit accuracy is forced. This may be weakened if  $t$  is replaced by some positive integer  $t_0 < t$ . In case one of the components of the exact solution is zero, some refined termination criterion may replace the one stated here. This should be done since the computation of an inclusion  $[-\eta, \eta]$  of zero is rather expensive.

**Algorithm:**

```

k := 0 ; { correction counter }
(N1) Solve  $f(\sum_{i=0}^{k-1} x^{(i)} + x^{(k)}) = 0$ 
      approximately for the unknown  $x^{(k)}$ 
      using Newton's method (2.6) with
      formula evaluation technique; in ill-
      conditioned cases perform each step
      with high accuracy using the linear
      techniques from [7]; call the approx-
      imation  $\tilde{x}^{(k)}$ ;

(N2) Compute an approximation R of
       $f'(\sum_{i=0}^k \tilde{x}^{(i)})^{-1}$ ;
      j := 0 ; { iteration counter }
      Z := Y ; X :=  $\diamond(\sum_{i=0}^k \tilde{x}^{(i)})$  ;
      repeat
        j := j+1;
        Z := Y  $\circ$   $\epsilon$ ;
        Y := W  $\diamond$   $\diamond(I-R \cdot J(X \cup (X \diamond Z))) \diamond Z$  ;
        success := (Y  $\subseteq$  Z) ;
      until success or (j = 10);

(N3) if success
      then dm := d(X  $\diamond$  Z) ;
           if dm < |X  $\diamond$  Z|  $\cdot$  b-t+1 goto (N4)
      else
        k := k+1 ;
        if k > 10
          then write('No solution
                    achieved'); stop
          else goto (N1)

(N4) Result:  $\hat{x} \in X \diamond Z$ 
  
```

**3. Examples of application**

The following two examples demonstrate the typical behaviour of the equation solver presented above. They will show the need of an exact scalar product and, furthermore, a formula evaluation technique using this scalar product.

**Example 1:**

We define the following system of equations for an unknown vector  $x \in V_nR$ :

$$f(x) := \alpha(H \cdot x + \psi(x)) = 0$$

where  $H \in M_nR$  is the Hilbert matrix of degree n, i.e.  $H = ((h_{ij}))$ ,  $h_{ij} := 1/(i+j-1)$ ,  $i, j = 1, \dots, n$ .  $\psi : V_nR \rightarrow V_nR$  is defined by

$$\psi(x) := (\psi_1, \dots, \psi_n) = x_1 + \dots + x_i, \epsilon_i > 0, i=1, \dots, n$$

The whole procedure is scaled by the lowest common multiple of the components of H to guarantee an exact scalar product. The problem has exactly one solution, for increasing n and small  $\epsilon_i$ , the solution becomes small. The characteristic polynomial is computed on a 68000 microprocessor using a 13 digit decimal arithmetic. The values of the  $\epsilon_i$  were chosen as an example for all other values, where the result for the component  $x_i$  is 3 we have the number of Newton steps. Column 4 answers the question whether the Newton approximation is performed with high accuracy. Column 5 displays the number of inclusion steps. The initial value for x was always the vector (1, 1, ..., 1).

n	initial value	error	appr. steps	high acc.?	incl. steps
5	-5.000	E-12	2	no	2
6	6.000	E-11	3	no	2
7	-6.000	E-12	3	no	2
8	7.000	E-12	4	no	2
9	-1.770	E-11	8	yes	2
10	-3.770	E-12	13	yes	2
11	-1.534	E-04	27	yes	2
12	-2.167	E-03	34	yes	2
13	8.000	E-01	27	yes	1
14	4.000	E-06	95	yes	1
15	1.000	E-06	35	yes	2

**Example 2:**

The following system arises from control engineering circuits and is known as the so-called "control equation":

$$x_0 = \alpha$$

$$x_i = 4x_{i-1} - 1 - x_{i-1}^2, i = 1, \dots, n.$$

The values of the components are theoretically determined by an explicit formula. Nevertheless, the best arithmetic will

approximate results for  $x_n$  without exact scalar products. In table 3.2b when we apply the method of evaluation by simple recursion, the correct digits are underlined. For  $n=10$ , as done on an IBM 4361, the values are obtained by using the ACRITH package. The results of table 3.2 are compared with the values of column 2 of table 3.2a.

n	relative error
10	0.9
20	0.9
30	0.9
40	0.6
50	0.7
60	0.9
70	0.9
80	0.9
90	0.9
100	0.9

n	relative error
10	0.9
20	0.9
30	0.9
40	0.6
50	0.7
60	0.9
70	0.9
80	0.9
90	0.9
100	0.9

places per step. The results for  $x_n$  without exact scalar products are shown in table 3.2b when the method of evaluation by simple recursion is used. The correct digits are underlined. For  $n=10$ , as done on an IBM 4361, the values are obtained by using the ACRITH package. The results of table 3.2 are compared with the values of column 2 of table 3.2a.

n	relative error	evaluation by interval arithmetic
10	0.9	0.93608886 <sup>12</sup> <sub>08</sub>
20	0.9	0.9484 <sup>8</sup> <sub>4</sub>
30	0.9	43.3
40	0.6	-31.2
50	0.7	"overflow"
60	0.9	----
70	0.9	----
80	0.9	----
90	0.9	----
100	0.9	----

n	relative error	near system evaluation with exact scalar product
10	0.9	0.933608886103744 <sup>5</sup> <sub>4</sub>
20	0.9	0.94841664586818 <sup>9</sup> <sub>8</sub>
30	0.9	0.949706664490856 <sup>5</sup> <sub>4</sub>
40	0.6	0.61487509460460 <sup>2</sup> <sub>1</sub>
50	0.7	0.772273942734024 <sup>8</sup> <sub>1</sub>
60	0.9	0.331008960225076 <sup>1</sup> <sub>0</sub>
70	0.9	0.33101258545544 <sup>8</sup> <sub>6</sub>
80	0.9	0.59730842574404 <sup>3</sup> <sub>2</sub>
90	0.9	0.210001469055217 <sup>8</sup> <sub>7</sub>
100	0.9	0.58608531833946 <sup>1</sup> <sub>0</sub>

References:

- [ 1 ] Alefeld, G., Herzberger, J.: Einführung in die Intervallrechnung, Bibl. Inst., Mannheim-Wien-Zürich (1974)
- [ 2 ] Böhm, H.: Berechnung von Polynomnullstellen und Auswertung arithmetischer Ausdrücke mit garantierter maximaler Genauigkeit, Doktor-Dissertation, Universität Karlsruhe (1983)
- [ 3 ] Böhm, H., Rump, S.M.: Least Significant Bit Evaluation of Arithmetic Expressions in Single-Precision, Computing 30, 189-199 (1983)
- [ 4 ] Böhm, H., Rump, S.M., Schumacher, G.: E-Methods for Nonlinear Problems, to appear in: "Computer Arithmetic: Scientific Computation and Programming Languages", Teubner Verlag, Stuttgart (1987)
- [ 5 ] Heuser, H.: Funktionalanalysis, 2nd ed., Teubner Verlag, Stuttgart (1986)
- [ 6 ] Kaucher, E.: Interval Analysis in the Extended Interval Space IIR, Computing Suppl. 2, 33-49 (1980)
- [ 7 ] Kaucher, E., Rump, S.M.: E-Methods for Fixed Point Equations  $f(x)=x$ , Computing 28, 31-42 (1982)
- [ 8 ] Kulisch, U., Miranker, W.: Computer Arithmetic in Theory and Practice, Academic Press, New York (1981)
- [ 9 ] Ortega, J.M., Rheinboldt, W.C.: Iterative Solution of Nonlinear Equations in Several Variables, Academic Press, New York, 1970
- [ 10 ] Rall, L.B.: Automatic Differentiation, Techniques and Applications, Lecture Notes in Computer Science, No 12, Springer, Berlin (1981)
- [ 11 ] Rump, S.M.: Solving nonlinear systems with least significant bit accuracy, Computing 29, 183-200 (1982)
- [ 12 ] Rump, S.M.: Solving algebraic problems with high accuracy in: "A new Approach to Scientific Computation", p. 51-120, Academic Press, New York (1983)
- [ 13 ] Stark, J.: Maximal genaue Auswertung arithmetischer Ausdrücke durch Lösen nichtlinearer Gleichungssysteme, Diplomarbeit at the Institut für Angewandte Mathematik, Universität Karlsruhe (1985)
- [ 14 ] Varga, R.S.: Matrix Iterative Analysis, Prentice Hall, Englewood Cliffs (1962)