

## An Interval Algorithm for Solving Systems of Linear Equations to Prespecified Accuracy

J. W. Demmel, New York, and F. Krückeberg, St. Augustin

Received July 27, 1983, revised June 29, 1984

### Abstract — Zusammenfassung

**An Interval Algorithm for Solving Systems of Linear Equations to Prespecified Accuracy.** We describe an interval arithmetic algorithm for solving a special class of simultaneous linear equations. This class includes but is not limited to systems  $Ax = b$  where  $A$  and  $b$  have integer entries. The algorithm uses fixed point arithmetic, and has two properties which distinguish it from earlier algorithms: given the absolute accuracy  $\varepsilon$  desired, the algorithm uses only as much precision as needed to achieve it, and the algorithm can adjust its own parameters to minimize computation time.

*AMS Subject Classifications:* 65G10 (primary), 65F05 (secondary).

*Key words:* Systems of linear equations, interval arithmetic, fixed point arithmetic.

**Ein Intervallalgorithmus für die Lösung von linearen Gleichungssystemen mit vorausgewählter Genauigkeit.** Wir beschreiben einen Intervallalgorithmus, der eine gewisse Klasse von linearen Gleichungssystemen löst. Diese Klasse enthält u. a. Systeme  $Ax = b$ , bei denen  $A$  und  $b$  ganzzahlige Komponenten haben. Dieser Algorithmus verwendet Festpunktarithmetik und unterscheidet sich von früheren Algorithmen wie folgt. Erstens: Bei Vorgabe der gewünschten absoluten Genauigkeit  $\varepsilon$  des Ergebnisses benötigt der Algorithmus nur so viel Zwischengenauigkeit wie notwendig, um die Fehlerschranke  $\varepsilon$  zu erreichen. Zweitens kann der Algorithmus selbststeuernd seine eigenen Parameter dynamisch ändern, um die Rechenzeit zu minimieren.

### 1. Introduction

In this paper we describe an interval arithmetic algorithm for solving a special class of systems of simultaneous linear equations. This class includes but is not limited to systems  $A\mathbf{x} = \mathbf{b}$  where  $A$  and  $\mathbf{b}$  have integer entries. (Capital italic letters denote matrices and lower case bold letters denote vectors.) The algorithm uses fixed point arithmetic where the precision may be chosen by the program. Our algorithm has two properties which distinguish it from previous interval linear system solvers:

- (1) Given the absolute accuracy  $\varepsilon$  desired in the solution, the algorithm uses only as much precision as needed to achieve it.
- (2) The algorithm can adjust its own parameters to minimize computation time.

Our research was motivated by ongoing work in Petri nets [1] at the Gesellschaft für Mathematik und Datenverarbeitung. Several decision problems in Petri nets can be reduced to deciding if a particular linear system  $A\mathbf{x} = \mathbf{b}$  with integer  $A$  and  $\mathbf{b}$  has a

nonnegative integer solution vector [2]. By solving  $A\mathbf{x}=\mathbf{b}$  with absolute accuracy  $\varepsilon < 1$ , we produce an interval vector  $\mathbf{x}$  such that  $A^{-1}\mathbf{b} \equiv \hat{\mathbf{x}} \in \mathbf{x}$  and the width of each component interval  $x_i$  is less than 1. Thus at most one vector  $\mathbf{x}_i$  of integer entries can lie within  $\mathbf{x}$ , and by testing to see if  $\mathbf{x}_i$  is nonnegative and satisfies  $A\mathbf{x}=\mathbf{b}$ , we may answer our decision problem.

Our assumptions and limitations in this paper are as follows:

- (1) Our underlying arithmetic delivers results of the operations add/subtract and multiply to within an *absolute* precision chosen by the program (the program will choose the precision once at the beginning of the computation and not change it). Thus addition/subtraction may be performed without rounding error if the precision of the result is no smaller than the precision of the operands. Fixed point arithmetic has these properties, and for our numerical tests we used a variable precision fixed point format (see section 5). This requirement implies that we can compute inner products to a given absolute accuracy.
- (2) We consider only problems  $A\mathbf{x}=\mathbf{b}$  where we can find an approximate (nonsingular) inverse matrix  $B$  of  $A$  with the following properties:
  - (2a) The matrix  $R \equiv I - BA$  and the vector  $\mathbf{c} \equiv B\mathbf{b}$  are exactly representable (i. e. without roundoff) in our number system.
  - (2b)  $\|R\|_\infty \equiv \max_i \sum_j |R_{ij}| = \text{max-row-sum norm} < 1$ .

Condition 2 a holds if  $A$  and  $\mathbf{b}$  have integer entries because if we round  $B$  to fit in our fixed point number system (and if integers are exactly representable) then we may then compute  $BA$ ,  $I - BA$ , and  $B\mathbf{b}$  without error by assumption (1). Condition 2 b is needed to guarantee that  $|R|$  is a contraction, that is  $|R|^n \rightarrow 0$  as  $n \rightarrow \infty$ . A condition similar to 2 b is required for convergence by virtually all interval arithmetic algorithms, as we will discuss below.

We do not care how  $B$  is computed. It may, for example, be computed using standard floating point library routines and then rounded to fit in our fixed point format.

The benefits of these assumptions are the following:

- (1) The condition number of the problem ( $\|A\|_\infty \|A^{-1}\|_\infty$ ) does not determine the limiting accuracy of the algorithm. The condition number will determine if we can find an approximate inverse  $B$  so that  $R = I - BA$  has norm less than 1, but as long as we can find such a  $B$ , we can compute  $A^{-1}\mathbf{b}$  to as much absolute accuracy  $\varepsilon$  as desired (limited, of course, by the accuracy available in the underlying arithmetic). This is in contrast to other interval arithmetic algorithms using a fixed, problem independent amount of precision in which the accuracy achievable degrades as the condition number grows [5, 6].
- (2) We may decide ahead of time (that is after computing  $B$ ,  $R$ , and  $\mathbf{c}$  but before starting the algorithm proper) exactly how much precision we need to use to achieve the desired accuracy  $\varepsilon$ . In fact, there are some parameters in the

algorithm which may be chosen to minimize the precision needed as a function of  $\|R\|_\infty$  and  $\varepsilon$ . By tuning the algorithm to the problem this way we can save time and possibly memory if a variable width fixed point format is used.

Let us put this algorithm into historical perspective. Almost all interval linear system solvers, including ours, convert the original problem  $A\mathbf{x}=\mathbf{b}$  into an iteration of the form

$$\mathbf{x}_{n+1} = R\mathbf{x}_n + \mathbf{c} \quad (1)$$

for some matrix  $R$  and vector  $\mathbf{c}$ , where  $|R|$  is a contraction. In our case  $R = I - BA$  and  $\mathbf{c} = B\mathbf{b}$  for some nonsingular  $B \approx A^{-1}$ , where  $A\mathbf{x}=\mathbf{b}$  is the equation we are trying to solve. The algorithm automatically verifies the nonsingularity of  $A$  and  $B$ .

A general linear system solver does not assume that  $R$  and  $\mathbf{c}$  can be computed exactly, as we do. Thus, in general, they are intervals. The width of  $\mathbf{c}$  is a lower bound on the width of all  $\mathbf{x}_{n+1}$ . Typically, the width of the intervals in  $R$  will depend on the condition number of  $A$  [6]. Thus, the width of  $\mathbf{x}_{n+1}$  will depend on the condition number of the problem, and, it turns out, of the size of the solution  $\mathbf{x}$  itself. Wongwises [5] shows that naive use of (1) does indeed produce solution intervals whose width is proportional to the condition number of  $A$ , a performance limitation shared by Gaussian elimination without any iterative improvement at all. Cleverer use of (1) can deliver the answer to an accuracy equivalent to roundoff error in the largest components of the solution in most cases, but there is still a decline in accuracy as the condition number of  $A$  gets too large [6, 8]. This inescapable dependence of accuracy on condition number is reflected in our algorithm by requiring higher precision to be able to deliver a desired accuracy  $\varepsilon$  if  $\|R\|_\infty$  is close to 1, but by assuming  $R$  and  $\mathbf{c}$  to be exact, we eliminate the complicated dependence of the achievable accuracy on the size the answer (which we do not, after all, know ahead of time) and on the width of the intervals in  $R$  and  $\mathbf{c}$ . In fact, the proof of our algorithm exploits the ease of determining the width of  $\mathbf{x}_{n+1}$  from  $\mathbf{x}_n$ .

We also note that almost every iterative interval algorithm has assumed that the initial interval vector  $\mathbf{x}_0$  contains the solution, and then iterates as follows:

$$\mathbf{x}_{n+1} = (R\mathbf{x}_n + \mathbf{c}) \cap \mathbf{x}_n. \quad (2)$$

$\mathbf{x}_0$  is often determined by using a very coarse approximate interval inverse matrix guaranteed to contain the actual inverse of  $A$ , and then multiplying this matrix by  $\mathbf{b}$  [5, 6].

Our algorithm, in contrast, makes no assumption about the accuracy of  $\mathbf{x}_0$ , but only about its width. Our algorithm works by taking an  $\mathbf{x}_n$  of width less than  $\varepsilon$ , artificially expanding it by adding a carefully chosen constant to the right endpoints and subtracting it from the left endpoints of each component interval, and then contracting this expanded interval using (1) several times to generate an  $\mathbf{x}_{n+1}$  also of width less than  $\varepsilon$ . This artificial expansion combined with the contraction (1) guaranteed that the solution  $\hat{\mathbf{x}}$  eventually lies in some  $\mathbf{x}_n$ . The algorithm terminates when  $\mathbf{x}_{n+1}$  lies within the expanded version of  $\mathbf{x}_n$ . Choosing the parameters such as precision, how much to artificially expand  $\mathbf{x}_n$ , and how many times to repeat (1) to

guarantee both  $\text{span}(\mathbf{x}_{n+1}) < \varepsilon$  and that the algorithm terminates precisely when it has found the solution constitutes the proof of the algorithm, given in Section 3.

Finally, we compare this work to that of Rump [9]. Our goal is to provide an answer to user specified accuracy while minimizing the work as a function of the required accuracy. We assume the input satisfies the constraints in (1), (2 a) and (2 b) above, and that our underlying arithmetic is variable precision. Rump's goal, on the other hand, is to deliver the answer to the same accuracy as the input (single precision floating point) while making no special assumptions about the input. His underlying arithmetic, in contrast to ours, has two precisions, single and "dot precision", the latter of which guarantees correctly rounded dot products of any two single precision vectors. (In the course of Rump's algorithm, he may compute some or all of the components of the solution vector to much higher than single precision, and his algorithm could be modified to make this information available on request. The format of this information would be rather different than ours: each component either as the sum of successively smaller single precision floating point numbers, or in a "dot precision" variable which can only participate in certain add, subtract, and rounding operations.) We both assume that the problem is not so ill-conditioned that our  $R$  matrices are too large for convergence. The different organizations of our algorithms reflect the different underlying arithmetics we work with, and so in practice the choice of algorithm would not surprisingly depend on what the machine supports most efficiently, as well as whether constraints (1), (2 a) and (2 b) were satisfied. We also note that the expansion technique of the last paragraph was discovered independently by Rump [9]. He expands his intervals by a certain fraction of their size, where he determined the fraction empirically as the one which gave best convergence [4]. We must choose the amount by which we expand subject to mathematical constraints given below; this careful choice of expansion leads to the guaranteed performance described in the theorem below.

In Section 2 we define our notation, in Section 3 we state our algorithm and prove it has the stated properties, in Section 4 we show how to choose certain parameters of the algorithm to minimize the computation time, and in Section 5 we present numerical results.

## 2. Notation

$\mathbf{F}_d$  denotes the set of fixed point numbers  $\{nd : n \text{ is an integer and } d \leq 1 \text{ is the distance between adjacent numbers}\}$  over which we perform our computations. Typically  $d$  will be 1 over a power of the radix (e.g.  $1/2^j$  or  $1/10^j$ ). We assume that  $1/d$  is an integer so that the integers are exactly representable in  $\mathbf{F}_d$ .  $\mathbf{I}_d$  denotes the set of intervals over  $\mathbf{F}_d$ . Scalars in  $\mathbf{F}_d$  and  $\mathbf{I}_d$  will be represented by lower case italic letters.  $w \in \mathbf{I}_d$  can be represented  $w = [\underline{w}, \bar{w}]$  where  $\underline{w}, \bar{w} \in \mathbf{F}_d$  and  $\underline{w} \leq \bar{w}$ . The span of an interval  $w$  is  $\text{span}(w) = \bar{w} - \underline{w}$ .  $\mathbf{F}_d^{n,m}$ ,  $\mathbf{F}_d^{n,m}$ ,  $\mathbf{I}_d^n$ , and  $\mathbf{I}_d^{n,m}$  denote  $n$ -vectors and  $n$  by  $m$  matrices over  $\mathbf{F}_d$  and  $\mathbf{I}_d$ . Vectors will be written as lower case bold letters, and matrices as capital italic letters. If  $\mathbf{w} \in \mathbf{I}_d^n$ , then  $\text{span}(\mathbf{w}) = \max_i \text{span}(w_i)$ , and similarly for  $\text{span}(R)$ ,  $R \in \mathbf{I}_d^{n,m}$ . All

quantities will be assumed to consist of intervals, unless we wish to emphasize that a certain variable is a point interval, i.e. an interval of span zero. Such variables will be written with a point over them, like  $\dot{R}$  and  $\dot{\mathbf{c}}$ .  $\dot{\mathbf{x}}$  will denote the solution of  $\mathbf{x} = \dot{R}\mathbf{x} + \dot{\mathbf{c}}$ .

$\| \mathbf{x} \|_\infty$  will denote the infinity norm of the vector  $\mathbf{x}$ :

$$\| \mathbf{x} \|_\infty \equiv \max_i |x_i|$$

and  $\| R \|_\infty$  will denote the matrix norm induced by this vector norm, the max-row-sum norm:

$$\| R \|_\infty \equiv \sup_{\mathbf{x} \neq 0} \frac{\| R \mathbf{x} \|_\infty}{\| \mathbf{x} \|_\infty} = \max_i \sum_j |R_{ij}|.$$

The definitions of addition, subtraction, and multiplication of intervals can be found in the literature [12]. Since we are interested in controlling the amount of precision used in our calculations, we use a notation for interval arithmetic which explicitly displays the precision:  $\text{int}(\text{'expression'}, d)$  denotes the result of evaluating the 'expression' in interval arithmetic over  $\mathbf{I}_d$ . In our application 'expression' will only contain additions, subtractions, and multiplications, so since we are assuming unbounded range in  $\mathbf{F}_d$ , the value of 'expression' is always well defined. 'Expression' may contain scalars, vectors, and matrices. If 'expression' is a single variable, the operation  $\text{int}(\text{'expression'}, d)$  only involves rounding outward.

We will need one more piece of notation to denote the artificial expansion of an interval mentioned above. If  $w = [w, \bar{w}] \in \mathbf{I}_d$ , then

$$\text{roundout}(w, d, r) \equiv [w - rd, \bar{w} + rd] \in \mathbf{I}_d.$$

$r$  is the integer number of interpoint distances  $d$  to round out.

### 3. The Algorithm

The algorithm is as follows:

- (3.1) Input  $\varepsilon, \hat{\mathbf{c}}, \hat{R}, \mathbf{x}_0$  ( $\text{span}(\mathbf{x}_0) \leq \varepsilon$  and  $n$  is the dimension of  $R$ )
- (3.2) Compute parameters  $m, t, d$ , and  $e$  according to the theorem
- (3.3)  $i := 0$
- (3.4) repeat
- (3.5)      $i := i + 1$
- (3.6)      $\mathbf{w} := \text{roundout}(\mathbf{x}_{i-1}, d, t)$
- (3.7)      $\mathbf{x}_i := \mathbf{w}$
- (3.8)     for  $j := 1$  to  $m$  do
- (3.9)          $\mathbf{x}_i := \text{int}(\text{int}(\hat{R} \mathbf{x}_i + \hat{\mathbf{c}}, de), d)$
- (3.10) until  $\mathbf{x}_i \subseteq \mathbf{w}$
- (3.11) Output  $\mathbf{x}_i$

The subscript  $i$  is only needed for stating the theorem below;  $\mathbf{x}_i$  may be written over  $\mathbf{x}_{i-1}$ . The expression in (3.9) indicates that  $\hat{R} \mathbf{x}_i + \hat{\mathbf{c}}$  is to be computed in  $\mathbf{I}_{de}$  ( $e \leq 1$ ) and the result rounded out to fit in  $\mathbf{I}_d$ .

The parameters  $m, t, d$ , and  $e$  have the following meanings:

$d$  determines the precision ( $\mathbf{F}_d$ ) in which we will represent our data ( $\hat{R}, \hat{\mathbf{c}}$ , and  $\mathbf{x}$ ),

$e$  is the extra precision used to compute the inner products in (3.9) ( $e \leq 1$  means  $\mathbf{F}_{de}$  is at least as precise as  $\mathbf{F}_d$ ; if  $d$  is a power of the radix (e.g.  $(1/2)^j$  or  $(1/10)^j$  for some integer  $j$ ) then we may also take  $e$  to be a power of the radix,  $t$  (integer) determines by how many units of  $d$  we round  $\mathbf{x}_{i-1}$  out in (3.6), and  $m$  (integer) is the iteration count in (3.8).

The following theorem describes how the algorithm works.

**Theorem:** *Given an arbitrary  $\varepsilon > 0$ , an arbitrary starting vector  $\mathbf{x}_0$  (such that  $\text{span}(\mathbf{x}_0) < \varepsilon$ ),  $R$  (such that  $r \equiv \|\hat{R}\|_\infty < 1$ ), and  $\hat{\mathbf{c}}$ , it is possible to choose the parameters  $d, e, t$ , and  $m$  so that the algorithm generates a sequence of interval vectors  $\{\mathbf{x}_i\}$  with the following properties:*

- (4.1)  $\text{span}(\mathbf{x}_i) < \varepsilon$  for all  $i$ . Thus, as soon as we know  $\hat{\mathbf{x}} \in \mathbf{x}_i$ ,  $\mathbf{x}_i$  is our answer.
- (4.2)  $\hat{\mathbf{x}} \in \mathbf{x}_i$  must occur for some finite  $i$ . Thus, the algorithm must eventually find  $\hat{\mathbf{x}}$ .
- (4.3)  $\hat{\mathbf{x}} \in \mathbf{x}_i$  implies the algorithm terminates on the next iteration. Combined with property (2) above, this property means the algorithm is finite.
- (4.4) If the algorithm terminates, the final  $\mathbf{x}_i$  must contain  $\hat{\mathbf{x}}$ .

Further, if  $\hat{R} = I - BA$ , and  $r = \|\hat{R}\|_\infty < 1$ , then both  $A$  and  $B$  are nonsingular, guaranteeing that  $\hat{\mathbf{x}}$  is the unique solution of  $A\mathbf{x} = \mathbf{b}$ .

In short, it is possible to choose the parameters to compute an arbitrarily narrow interval vector  $\mathbf{x}_i$  containing the solution  $\hat{\mathbf{x}}$ . In addition, we will see how to choose the parameters to minimize the cost of computation in a problem dependent way.

*Proof:* More succinctly, the four properties above are:

- (4.1')  $\text{span}(\mathbf{x}_i) < \varepsilon \rightarrow \text{span}(\mathbf{x}_{i+1}) < \varepsilon$ ,
- (4.2')  $\hat{\mathbf{x}} \in \mathbf{x}_i$  occurs for some finite  $i$ ,
- (4.3')  $\hat{\mathbf{x}} \in \mathbf{x}_i \rightarrow \mathbf{x}_{i+1} \subseteq \mathbf{w}$ , and
- (4.4')  $\mathbf{x}_i \subseteq \mathbf{w} \rightarrow \hat{\mathbf{x}} \in \mathbf{x}_i$ .

We will prove the four properties in the order (4.4'), (4.2'), (4.1'), and (4.3'). (4.4') is an application of Brouwer's fixed point theorem [13] standard in interval analysis. To prove (4.2'), we use the contraction mapping theorem [14]. Let  $\hat{\mathbf{x}}_0$  be any point in  $\mathbf{x}_0$ , and let  $\hat{\mathbf{x}}_i$  be the point image of  $\hat{\mathbf{x}}_{i-1}$  under  $\mathbf{F}$ , where  $\mathbf{F}$  was defined by (3.6–3.9). Then since

$$\hat{\mathbf{x}}_i = R^m \hat{\mathbf{x}}_{i-1} + \hat{\mathbf{c}}' \quad (5)$$

for some constant vector  $\hat{\mathbf{c}}'$  we have

$$\|\hat{\mathbf{x}}_i - \hat{\mathbf{x}}\|_\infty \leq r^m \|\hat{\mathbf{x}}_{i-1} - \hat{\mathbf{x}}\|_\infty \leq r^{mi} \|\hat{\mathbf{x}}_0 - \hat{\mathbf{x}}\|_\infty. \quad (6)$$

Since we use interval arithmetic, the point vector  $\hat{\mathbf{x}}_i$  is a member of the interval vector  $\mathbf{x}_i$ .

Thus, since we assumed  $r < 1$ , some point in  $\mathbf{x}_i$  is eventually closer to the solution  $\hat{\mathbf{x}}$  than  $t d$  (here we use the fact that  $t$  is an integer  $\geq 1$ ). Then, the next time through the main loop (3.4)–(3.10),  $\mathbf{x}_i$  will be rounded out far enough so that  $\mathbf{w}$  contains  $\hat{\mathbf{x}}$ . Since  $\hat{\mathbf{x}} \in \mathbf{F}(\mathbf{w}) \subseteq \mathbf{x}_i$  ( $\hat{\mathbf{x}}$  being a fixed point of  $\mathbf{F}$ ), (4.2') is satisfied. It is easy to see that the  $i$  for

which (4.2') is true is bounded by

$$i \leq \frac{\log(td/\|\dot{\mathbf{x}}_0 - \dot{\mathbf{x}}\|_\infty)}{m \log r} + 2. \tag{7}$$

We will now derive two inequalities in the parameters  $d, e, t,$  and  $m$  that are sufficient conditions for the validity of (4.1') and (4.3'). By solving these inequalities simultaneously for the parameters, we will prove the theorem. Since the two inequalities do not determine the four parameters uniquely, we can choose the parameters to minimize the cost of computation. The first inequality is in  $\varepsilon, d, e, t$  and  $m,$  and is a sufficient condition for (4.1') to be true. Since the coordinates of  $\mathbf{x}_i$  are all multiples of  $d,$  we introduce an integer variable  $l$  which satisfies

$$ld \leq \varepsilon \tag{8}$$

and replace (4.1') by

$$\text{span}(\mathbf{x}_i) < ld \rightarrow \text{span}(\mathbf{x}_{i+1}) < ld.$$

Our inequality will be in terms of  $l$  rather than  $d.$  We take  $l$  to be an integer variable since  $\text{span}(\mathbf{x}_i)$  must be an integer multiple of  $d.$

Now  $\text{span}(\mathbf{x}_i) < ld$  means

$$\text{span}(\mathbf{w}) < (l + 2t)d \tag{9}$$

by the definition of roundout used in (3.6).

Using the standard [12] formulas for spans of sums and products of intervals, we see in (3.9) that

$$\text{span}(\text{int}(R\mathbf{w} + \mathbf{c}, de)) \leq r \text{span}(\mathbf{w}) + 2nde. \tag{10}$$

Therefore, after  $m$  iterations of (3.9) we get

$$\text{span}(\mathbf{x}_{i+1}) \leq r^m \text{span}(\mathbf{w}) + (1 + r + \dots + r^{m-1}) \cdot 2d(ne + 1). \tag{11}$$

Substituting  $(l + 2t)d$  for  $\text{span}(\mathbf{w})$  in the R. H. S. of (11) and requiring this quantity to be less than  $ld$  (which is a sufficient condition for (4.1')) yields (after some manipulation) our first inequality for  $l, e, t$  and  $m.$

$$l > \frac{2r^m}{1-r^m} t + \frac{2(ne+1)}{1-r}. \tag{12}$$

We will derive our second inequality from (4.3').  $\dot{\mathbf{x}} \in \mathbf{x}_i$  means that the distance from  $\dot{\mathbf{x}}$  to the edge of  $\mathbf{w}$  is at least  $td$  and no more than  $(l+t)d$  (see (3.6)). We seek a relationship among  $l, d, e, t,$  and  $m$  that guarantees that every point of  $\mathbf{x}_{i+1}$  is no farther from  $\dot{\mathbf{x}}$  than  $td,$  because this will imply that  $\mathbf{x}_{i+1} \subseteq \mathbf{w}$  as desired. We expect to be able to find such a relationship because the contractive property of  $\mathbf{F}$  implies all points in  $\mathbf{w}$  will approach  $\dot{\mathbf{x}}$  under the action of  $\mathbf{F}.$

Thus, if  $\dot{\mathbf{x}} \in \mathbf{x}_i$  and  $\mathbf{w} = \text{roundout}(\mathbf{x}_i, d, t)$  we have

$$\sup_{\dot{\mathbf{y}} \in \mathbf{w}} \|\dot{\mathbf{y}} - \dot{\mathbf{x}}\|_\infty \leq (l+t)d. \tag{13}$$

After every iteration of (3.9) we have from an analysis similar to the one leading to (12)

$$\sup_{\dot{\mathbf{y}} \in \mathbf{x}_{t+1} \text{ after (3.9)}} \|\dot{\mathbf{y}} - \dot{\mathbf{x}}\|_{\infty} \leq r \cdot \sup_{\dot{\mathbf{y}} \in \mathbf{x}_{t+1} \text{ before (3.9)}} \|\dot{\mathbf{y}} - \dot{\mathbf{x}}\|_{\infty} + nde + d, \quad (14)$$

so after  $m$  iterations of (3.9) we get

$$\sup_{\dot{\mathbf{y}} \in \mathbf{x}_{t+1}} \|\dot{\mathbf{y}} - \dot{\mathbf{x}}\|_{\infty} \leq r^m(l+t)d + \left(\frac{1-r^m}{1-r}\right)d(ne+1). \quad (15)$$

We require that the R.H.S. of (15) be no larger than  $td$ :

$$l \leq \frac{1-r^m}{r^m}t - \frac{1-r^m}{r^m(1-r)}(ne+1). \quad (16)$$

(16) is our second inequality relating  $l$ ,  $e$ ,  $t$ , and  $m$ .

Now we have to solve the inequalities (12) and (16) simultaneously in  $l$ ,  $e$ ,  $t$ , and  $m$  (note that  $d$  does not appear; we will deal with it later). We will only show here that a solution does exist, deferring to the next section a discussion of finding the best solution. We may choose any  $e$  such that  $0 < e \leq 1$ . Considering (12) and (16) only as inequalities in  $l$  and  $t$ , we see they are both linear, and so both determine half planes in the  $t, l$  plane. It is easy to see that (12) and (16) can only have a common, positive solution in  $t$  and  $l$  if  $(1-r^m)/(r^m) > (2r^m)/(1-r^m)$ , or

$$m > \frac{\log(\sqrt{2}-1)}{\log r}. \quad (17)$$

Having chosen  $m$  subject to this last constraint, the region of common solution of (12) and (16) is a sector in the  $t, l$  plane, and so must contain points with integer ( $t, l$ ) coordinates, any of which are candidate solutions for (12) and (16). The smallest possible values of  $t$  and  $l$  are the (not necessarily integer) coordinates of the apex of the sector, yielding the lower bounds:

$$t > (ne+1) \frac{1-r^m}{1-r} \frac{1+r^m}{1-2r^m-r^{2m}}, \quad (18)$$

and

$$l > (ne+1) \frac{1-r^m}{1-r} \frac{2}{1-2r^m-r^{2m}}.$$

So far we have shown how to choose  $e$ ,  $t$ ,  $l$  and  $m$  in order to simultaneously satisfy (12) and (16). It remains to choose  $d$ .  $d$  may be any number satisfying  $ld < \varepsilon$ , or  $d < \varepsilon/l$ . The lower bound on  $l$  in (18) translates into an upper bound on  $d$ :

$$d < \frac{\varepsilon(1-2r^m-r^{2m})}{2(ne+1)} \frac{1-r}{1-r^m}. \quad (19)$$

This completes the proof of the theorem. Q.E.D.



#### 4. Choosing the Parameters to Minimize Cost

In the proof of the theorem we showed that subject to certain inequality constraints, we could choose the parameters  $m$ ,  $t$ ,  $d$ , and  $e$  to make the algorithm behave as claimed. It became clear in the discussion that the constraints left some freedom in the choice of these parameters. In this section we will exploit that freedom and show how the parameters may be chosen to minimize the cost of the algorithm. In particular, we will see that the naive choice of parameters suggested by the proof of the theorem results in a cost function with poles at a countable number of values of  $r$ , so that some care really must be exercised in choosing parameters.

The cost is proportional to the product of the following four factors:

- (20.1)  $i_M$  = the number of iterations of the main loop (3.4–3.10) of the algorithm. A bound for  $i_M$  is  $i+1$ , where  $i$  is bounded in (7).
- (20.2)  $m$  = the number of iterations of the inner loop (3.8–3.9). A lower bound for  $m$  is given in (17).
- (20.3) The cost of a multiplication in the inner loop (3.9). This cost is a function, mult, of the number of places,  $p$ , used in the computation.  $p$  is proportional to  $\log(1/de)$ , and mult can grow as slowly as  $p \log p \log \log p$  [15] or as quickly as  $p^2$  depending on the implementation of multiplication.
- (20.4) The number of multiplies in the inner loop,  $2n^2$ .

Thus, we model the cost as follows:

$$\text{Cost} = 2n^2 \times \left( \frac{\log(\|\dot{\mathbf{x}}_0 - \dot{\mathbf{x}}\|_\infty / td)}{\log 1/r} + 3m \right) \times \text{mult} \left( \log \left( \frac{1}{de} \right) \right). \quad (21)$$

We do not include the cost of computing  $R$  and  $\mathbf{c}$ , which is  $n^3$  (fixed point) multiplies no more expensive than those in the inner loop above, as well as  $O(n^3)$  (floating point) multiplies to compute the approximate inverse  $B$ . The cost in (21) above may be less or more than the cost of obtaining  $R$  and  $\mathbf{c}$ , depending on the initial error  $\|\dot{\mathbf{x}}_0 - \dot{\mathbf{x}}\|_\infty$  and the desired precision  $\varepsilon$ . Here we address only the problem of minimizing the cost in (21).

To minimize (21) exactly, we would need to know the initial error  $\|\dot{\mathbf{x}}_0 - \dot{\mathbf{x}}\|_\infty$ , the function mult, and the set of discrete values to which  $d$  is restricted. Since we do not know all these things in general, we just show that the following algorithm makes a reasonable choice of  $m$ ,  $t$ ,  $d$ , and  $e$ :

- (22.1)  $e := 1$  (i.e. no extra precision),
- (22.2) choose the smallest  $m$  such that  $r^m < 1/3$ ,
- (22.3)  $d :=$  largest value less than the upper bound in (24) (recall  $d$  must be a negative power of the radix),
- (22.4)  $l :=$  largest integer such that  $l < \varepsilon/d$ ,
- (22.5) given  $l$ , find the largest integer  $t$  satisfying both (12) and (16),
- (22.6) if no such  $t$  exists, or if  $t/l < \frac{1}{2} \frac{1-r^m}{2r^m}$ , decrease  $d$  to the next smaller value and goto (22.4).

With this choice of parameters, we will see that the cost is

$$\text{Cost} \approx 2n^2 \times \left( \frac{\log(\|\dot{\mathbf{x}}_0 - \dot{\mathbf{x}}\|_\infty/\varepsilon) + K_1}{\log\left(\frac{1}{r}\right)} \right) \times \text{mult} \left[ \log\left(\frac{1}{\varepsilon}\right) + \log\left(\frac{1}{1-r}\right) + K_2 \right] \tag{23}$$

for modest constants  $K_1$  and  $K_2$ . The cost goes to  $\infty$  as either the desired precision  $\varepsilon$  goes to zero or the residual norm  $r$  goes to 1, as expected.

The rationale behind this algorithm is as follows. From (21) we see that the cost decreases as the precision decreases (*de* increases) and the amount of round out in (3.6) increases (*td* increases). Approximating  $d$  by  $\varepsilon/l$  (see (8)) we see that the precision *de* is

$$de = \frac{\varepsilon}{2} \cdot \frac{e}{ne + 1} \cdot \frac{1-r}{1-r^m} \cdot (1 - 2r^m - r^{2m}). \tag{24}$$

The term depending on the extra precision  $e$  is maximized at  $e=1$  (no extra precision). This justifies (22.1).

The  $1-r$  factor reflects the inescapable effect of condition number on precision, since the larger the condition number of  $A$ , the closer to one  $r$  is likely to be, and so the more precision needed.

The interesting factor is  $1 - 2r^m - r^{2m}$ . If we chose  $m$  to be the smallest integer greater than or equal to its lower bound  $\log(\sqrt{2}-1)/\log r$  (see (17)), then as a function of  $r$ ,  $1 - 2r^m - r^{2m}$  would have a countable number of zeroes ( $\{(\sqrt{2}-1)^{1/j}\}_{j=1, \infty}$ ) where the required precision goes to infinity! We avoid this strange behavior by simply increasing  $m$  (by at most a factor of 1.25) if  $r^m$  is too close to  $\sqrt{2}-1$ , guaranteeing that  $r^m < 1/3$  and so  $1 - 2r^m - r^{2m} > .22$ . This justifies (22.2).

The next four steps of the algorithm maximize the roundout  $td \approx \varepsilon t/l$ . The maximum value of  $t/l$  is approached as both  $t$  and  $l$  approach  $+\infty : (1-r^m)/2r^m$  (the reciprocal of the slope of the lower boundary; see (12)). Therefore, in line (22.4) we pick the largest possible  $l$  subject to  $d < \varepsilon/l$ , and in (22.5) we pick the largest  $t$  for that  $l$ . Such a  $t$  may not exist if the solution sector for (12) and (16) is unfortunately located; even if it does exist  $t/l$  may be far from its maximum value  $(1-r^m)/2r^m$ . In either case we decrease  $d$  in (22.6) (by at least a factor of 2 in binary arithmetic, in general by a factor of the radix) and recompute  $l$  and  $t$ . The constraint  $r^m < 1/3$  and the looseness of the test  $t/l < (1-r^m)/(4r^m)$  guarantees that  $t$  and  $l$  will be recomputed by (22.4–22.5) at most once. This justifies the rest of the algorithm. We may now substitute the values of  $m$ ,  $t$ ,  $d$ , and  $e$  computed by this algorithm into our cost function (21) to get the result (23).

We mention the case when  $r \ll 1$  (for which  $m=1$ ) which is likely when  $A$  is not too ill-conditioned. In this case we may pick  $t$  very large while keeping  $l$  small since the lower boundary of the solution sector is almost horizontal (slope =  $2r/(1-r)$ ). If we can pick  $t$  large enough so that  $dt$  is larger than the initial error  $\|\dot{\mathbf{x}}_0 - \dot{\mathbf{x}}\|_\infty$  then the first  $\mathbf{w}$  will contain the solution  $\dot{\mathbf{x}}$ , and our algorithm will terminate after at most two

iterations of the outer loop. If our initial guess  $\hat{x}_0$  is good enough (and many good interval algorithms start with a reasonably accurate solution obtained cheaply using noninterval arithmetic), then this quick termination is likely.

### 5. Numerical Results

To test our algorithm we tried to invert the Hilbert matrix scaled to have integer entries. Specifically, we multiplied  $H_{ij} = 1/(i+j-1)$  by the least common multiple of  $1, \dots, 2n-1$  ( $n$  is the dimension of the matrix). The Hilbert matrix is known to be very ill-conditioned and so is a good test of our algorithm.

The algorithm was written in FORTRAN and executed on an IBM 370/158. The arithmetic used was also written in FORTRAN (and accessed via subroutine calls). It used a fixed point format, with 80 decimal places before and 80 places after the decimal point. Decimal digits were available in blocks of 4, making the radix effectively  $10^4$ .

The approximate inverse  $B$  was computed using the NAG [16] scientific subroutine library. The double precision versions of F01ADF and F01ACF were used. F01ADF computes the approximate inverse of a symmetric, positive-definite matrix, such as the Hilbert matrix, and F01ACF computes an accurate inverse using iterative refinement.

In all cases,  $\varepsilon$  was  $10^{-15} \|H^{-1}\|_\infty$  (which means that the largest element in  $H^{-1}$  was computed to a relative accuracy of at least  $10^{-15}$ ),  $e$  was taken to be  $10^{-4}$ ,  $m$  turn out to be 1 (since  $r$  was  $< 1/3$  in all cases), and  $\hat{x}_0$  was taken from the approximate inverse supplied by the NAGLIB routines. Our results are shown in Table 1. Column 1 (labeled  $n$ ) gives the dimension of the Hilbert matrix, column 2 (labeled Routine) indicates which NAGLIB routine was used to obtain  $B$  (C for F01ACF and D for F01ADF), column 3 (labeled Cond) gives the condition number  $\|H\|_\infty \|H^{-1}\|_\infty$  of the Hilbert matrix, columns 4, 6, 7, and 8 give the values of the parameters  $r$ ,  $d$ ,  $t$ , and  $l$ , column 5 gives  $\varepsilon = 10^{-15} \|H^{-1}\|_\infty$ , and column 9 (labeled  $i_M$ ) gives the maximum number of iterations of the algorithms main loop used by any of the columns of the inverse.

Table 1. Numerical results

$n$	Routine	Cond	$r$	$\varepsilon$	$d$	$t$	$l$	$i_M$
10	D	$3.5 \cdot 10^{13}$	.0084	$5.2 \cdot 10^{-11}$	$10^{-12}$	3	2	5
11	C	$1.2 \cdot 10^{15}$	.023	$1.8 \cdot 10^{-9}$	$10^{-12}$	3	2	2
12	C	$4.1 \cdot 10^{16}$	.23	$2.5 \cdot 10^{-9}$	$10^{-12}$	5	3	3
13	C	$1.3 \cdot 10^{18}$	20.1	$2.5 \cdot 10^{-8}$	—	—	—	—

All the inequalities and inclusions predicted in the theorem were automatically tested and verified. The values of  $t$  and  $l$  were not chosen according to algorithm 27 but as the smallest solutions of (14) and (18); our concern was not to minimize time but to verify the conclusions of the theorem. F01ADF (approximate inverse) could

not be used for  $n > 10$  because the  $B$  it produced was so inaccurate that  $r$  was greater than 1. For  $n = 13$   $r$  was greater than 1 even with the more accurate library routine, so we could not use the algorithm. The number of iterations  $i_M$  did not grow with  $n$  because the initial approximation  $\hat{x}_0$  came from the NAGLIB inverse  $B$  and so was rather accurate to start with. Indeed, most good interval algorithms wisely attempt to attain as accurate a solution as possible using noninterval arithmetic, requiring only a few iterations of the relatively expensive interval arithmetic at the end to refine the result, and ours is no exception.

### Acknowledgement

The authors acknowledge the kind assistance of several of Prof. Krückeberg's students who implemented the fixed point arithmetic and assisted in programming the algorithm. The authors also acknowledge the financial support of the Gesellschaft für Mathematik und Datenverarbeitung and the U.S. Department of Energy, Contract DE-AM03-76SF00034, Project Agreement DE-AS03-79ER10358. This work was performed while the first author was visiting the GMD in St. Augustin and while the Computer Science Division, University of California, Berkeley, CA 94720.

### References

- [1] Genrich, H. J., Lautenbach, K., Thiagarajan, P. S.: Elements of general net theory. In: Net Theory and Applications (Lecture Notes in Computer Science, Vol. 84), pp. 121 – 164. Berlin-Heidelberg-New York: Springer 1980.
- [2] Memmi, G., Roucairol, G.: Linear algebra in net theory. In: Net Theory and Applications (Lecture Notes in Computer Science, Vol. 84), pp. 213 – 224. Berlin-Heidelberg-New York: Springer 1980.
- [3] Wilkinson, J. H., Reinsch, C.: Linear algebra (Handbook for Automatic Computation, Vol. 2), pp. 41 – 44. Berlin-Heidelberg-New York: Springer 1971.
- [4] Rump, S. M.: Private communication, 1982.
- [5] Wongwises, P.: Experimentelle Untersuchungen zur Numerischen Auflösung von linearen Gleichungssystemen mit Fehlererfassung. In: Interval Mathematics (Lecture Notes in Computer Science, Vol. 29) (Nickel, K., ed.), pp. 316 – 325. Berlin-Heidelberg-New York: Springer 1975.
- [6] Krier, N., Spelluci, P.: Untersuchungen der Grenzgenauigkeit von Algorithmen zur Auflösung linearer Gleichungssysteme mit Fehlererfassung. In: Interval Mathematics (Lecture Notes in Computer Science, Vol. 29) (Nickel, K., ed.), pp. 288 – 297. Berlin-Heidelberg-New York: Springer 1975.
- [7] Krawczyk, R.: Newton-Algorithmen zur Bestimmung von Nullstellen mit Fehlerschranken. *Computing* 4, 187 – 201 (1969).
- [8] Wilkinson, J. H.: Rounding Errors in Algebraic Processes. Englewood Cliffs, N. J.: Prentice-Hall 1963.
- [9] Rump, S. M.: Kleine Fehlerschranken bei Matrixproblemen. Dissertation, Universität Karlsruhe, 1980.
- [10] Bareiss, E. H.: Sylvester's identity and multistep integer-preserving Gaussian elimination. *Math. Comp.* 22, 565 – 578 (1968).
- [11] Howell, J. A., Gregory, R. T.: An algorithm for solving linear algebraic equations using residue arithmetic. *BIT* 9, 200 – 234, 324 – 337 (1969).
- [12] Moore, R. E.: Methods and applications of interval analysis. Philadelphia: SIAM 1979.
- [13] Dunford, N., Schwartz, J. T.: Linear Operators, Part I, p. 453. New York: Interscience 1957.
- [14] Apostol, T.: Mathematical Analysis, 2nd ed., p. 92. Reading, Mass.: Addison-Wesley 1974.

- [15] Aho, A. V., Hopcroft, J. E., Ullman, J. D.: The Design and Analysis of Computer Algorithms, pp. 270–274. Reading, Mass.: Addison-Wesley 1974.
- [16] NAG Reference Manual, Oxford, Numerical Algorithms Group, Ltd., 1976.

J. Demmel  
Department of Computer Science  
Courant Institute  
New York University  
251 Mercer Street  
New York, NY 10012  
U.S.A.

Prof. F. Krückeberg  
Gesellschaft für Mathematik  
und Datenverarbeitung  
Schloss Birlinghoven  
D-2500 St. Augustin  
Federal Republic of Germany