# Consequences of a properly implemented computer arithmetic for periodicities of iterative methods

R. Klatte and Ch. Ullrich

Abstract - In ordered sets it is possible to show under certain assumptions two basic theorems concerning the cycle length of sequences of iterates generated by monotone operators. These results are applied to different iterative methods, where the conclusions are valid for the sequences of iterates produced by the numerical computations only, if the used computer arithmetic is properly implemented.

Index terms - Rounding invariant structures, cycle, weakly cyclic vector function, floating-point arithmetic.

## 1. Introduction and Summary

Because of the technical conditions the set R of the numbers representable in a computer is finite. If we execute an algorithm of the form $x_{n+1} = Fx_n$, $n = 0,1,2,\ldots$ , beginning with any element $x_0 \in R$ , it is clear, that the sequence $\{x_n\}$ of the iterates has only a finite number of different elements. Moreover, from an element $x_k$ the values of the m preceding iterates $x_{k-m}, x_{k-m+1}, \ldots, x_{k-1}$ return in the same order, i.e. the algorithm is ending in a cycle of the length m . After receiving a cycle it is useless to continue the iteration because we cannot obtain any further information. Therefore we can regard the end of the first cycle as the end of the numerical computation which easily can be recognized,if the cycle length is known.

For the studies of cycles it is usefull to know the mathematical structure of the arithmetic in R . Those algebraic and order properties, which should hold for any floating-point system, namely the usual rules for the neutral elements 0,1 and the minus sign as well as compatibility properties between the algebraic and the order structure well-known from the real number field $\mathbb{R}$ , have been summarized to the concept of the completely linearly ordered division ringoid $\{R,N,+,\cdot,/,\leq\}$ ([ 3 ]). In [ 5 ] it could be shown that R with a properly implemented floating-point arithmetic represents such a structure. In a similar way it is possible to describe computations with vectors $x \in V_n R$ and $n \times n$ matrices $A \in M_n R$ as well as the matrix vector multiplication $\cdot : M_n R \times V_n R \rightarrow V_n R$ by the structure of the ordered vectoid over an ordered division ringoid

([ 6 ]). The intent of this paper is to show the possibility of determining or estimating the length of cycles of iterative methods only by means of the algebraic and order properties of the above structures.

The following results are essentially based on two simple theorems for monotone operators in an ordered set $\{M, \leq\}$ : an isotone operator always generates a sequence of iterates with the cycle length 1, if two elements $x_{i_0}, x_{i_0+1}$ of the sequence are comparable with respect to $\leq$ and an antitone operator produces at most cycles of length 2, if three consecutive iterates are comparable.

It will be shown that for all non-negative resp. non-positive elements $A \in M_n R$ and for all $b \in V_n R$ the operator $T : V_n R \rightarrow V_n R$ with $Tx := Ax + b$ is an isotone resp. antitone operator. Therefore we can apply these theorems to T immediately. Without the assumption $A \geq 0$ resp. $A \leq 0$, it is nevertheless possible to get results for special matrices, which follow from general studies of weakly cyclic vector functions in $V_n R$ ([ 2 ]). So, for example we are able to determine the cycle length of iteration sequences produced analogously to the Jacobi resp. Gauss-Seidel method , if A is a weakly cyclic matrix ([ 7 ]) or a generalization of it. All results are illustrated by numerical examples.

## 2. Rounding Invariant Structures

In recent investigations it was shown that the algebraic and order properties of the mathematical space given on a computer allow to be described by the structure of the completely linearly ordered division ringoid ([ 3 ]). We shortly repeat this definition:

Definition 1: A non empty ordered set $\{R, \leq\}$ with two binary operations $+, \cdot$ is called an "ordered ringoid" $\{R,+,\cdot,\leq\}$, if the properties (D1) to (D6) and (OD1) to (OD3) hold:

(D1) $\bigwedge\limits_{a,b \in R} a + b = b + a$ .

(D2) $\bigvee\limits_{o \in R} \bigwedge\limits_{a \in R} a + o = a$ .

(D3) $\bigvee\limits_{e \in R} \bigwedge\limits_{a \in R} a \cdot e = e \cdot a = a$ .

(D4) $\bigwedge\limits_{a \in R} a \cdot o = o \cdot a = o$ .

(D5) There exists a uniquely defined element $x \neq e$

of $R$ such that with $-a := x \cdot a$ hold :

(a) $x \cdot x = e$ ,

(b) $\bigwedge\limits_{a,b \in R} -(a \cdot b) = (-a) \cdot b = a \cdot (-b)$ ,

(c) $\bigwedge\limits_{a,b \in R} -(a + b) = (-a) + (-b)$ .

(OD1) $\bigwedge\limits_{a,b,c \in R} (a \leq b \Rightarrow a + c \leq b + c)$ .

(OD2) $\bigwedge\limits_{a,b \in R} (a \leq b \Rightarrow -b \leq -a)$ .

(OD3) $\bigwedge\limits_{a,b,c \in R} (o \leq a \leq b \wedge c \geq o \Rightarrow a \cdot c \leq b \cdot c \wedge$
$$c \cdot a \leq c \cdot b) .$$

The elements $-e,o,e \in R$ are called the special elements of $\{R,+,\cdot,\leq\}$ . An ordered ringoid $\{R,+,\cdot,\leq\}$ with a further inner operation $/ : R \times R \backslash N \rightarrow R$ with $o \in N \subseteq R$ is called an "ordered division ringoid" $\{R,N,+,\cdot,/,\leq\}$ , if the properties (D6) to (D8) and (OD4) are fulfilled:

(D6) $\bigwedge\limits_{a \in R} a / e = a$ .

(D7) $\bigwedge\limits_{a \in R \backslash N} o / a = o$ .

(D8) $\bigwedge\limits_{a \in R} \bigwedge\limits_{b \in R \backslash N} -(a / b) = (-a)/b = a / (-b)$ .

(OD4) $\bigwedge\limits_{a,b,c \in R} (o <[1]) a \leq b \wedge c > o \Rightarrow o \leq a/c \leq b/c$
$$\wedge \quad o \leq c/b \leq c/a) .$$

An ordered ringoid $\{R,+,\cdot,\leq\}$ (resp. an ordered division ringoid $\{R,N,+,\cdot,/,\leq\}$) is called a "completely ordered ringoid" (resp. a "completely ordered division ringoid"),if $\{R, \leq\}$ is a complete lattice.

An ordered ringoid $\{R,+,\cdot,\leq\}$ (resp. ordered division ringoid $\{R,\{o\},+,\cdot,/,\leq\}$) is called linearly ordered,if $\{R, \leq\}$ is a linearly ordered set (resp. and

(OD5) $\bigwedge\limits_{a \in R \backslash \{o\}} a / a = e$

holds). $\qquad\qquad\qquad\qquad\qquad$ □

Let $M_n R$ be the set of $n \times n$ - matrices with elements of the ordered ringoid $\{R,+,\cdot,\leq\}$ . With the usual addition and multiplication for matrices $A = (a_{ij})$ , $B = (b_{ij}) \in M_n R$

$\quad A + B := (a_{ij} + b_{ij})$

$$A \cdot B := (\sum_{j=1}^{n} a_{ij} \cdot b_{jk}) \quad [2]$$

and the order relation

$A \leq B : \iff (a_{ij} \leq b_{ij}$ for $i,j = 1(1)n)$ ,

$\{M_n R,+,\cdot,\leq\}$ is an ordered ringoid with the special elements

$$-E = \begin{bmatrix} -e & & 0 \\ & -e & \\ 0 & & \ddots \\ & & & -e \end{bmatrix}, O = \begin{bmatrix} o & \cdots\cdots & o \\ o & \cdots\cdots & o \\ o & \cdots\cdots & o \end{bmatrix}, E = \begin{bmatrix} e & & 0 \\ & e & \\ 0 & & \ddots \\ & & & e \end{bmatrix}.$$

Let further $V_n R$ be the set of n-tuples over an ordered ringoid $\{R,+,\cdot,\leq\}$ with the usual inner operation $+$

$\quad a + b := (a_i + b_i)$ , $a = (a_i)$ and $b = (b_i)$,i=1(1)n and the order relation

$\quad a \leq b : \iff (a_i \leq b_i$ for $i = 1(1)n)$ .

Additionally we can define an outer operation $\cdot$ : $M_n R \times V_n R \rightarrow V_n R$ by

$$A \cdot b := (\sum_{j=1}^{n} a_{ij} \cdot b_j), A = (a_{ij}), b = (b_j) ,$$
$$i,j=1(1)n .$$

To be able to describe the properties of these operations we need the following

<u>Definition 2</u>: Let $\{R,+,\cdot,\leq\}$ be an ordered ringoid with the special elements $\{-e,o,e\}$ and $\{V, \leq\}$ an ordered set with a commutative binary operation $+$ and the neutral element $o$ . $V$ is called an "ordered vectoid over $R$ " (or an "ordered R-vectoid"), if additionally an outer operation $\cdot$ : $R \times V \rightarrow V$ is defined so that the following properties hold:

(VD1) $\bigwedge\limits_{a \in R} \bigwedge\limits_{a \in V} (a \cdot o = o \wedge o \cdot a = o)$ .

(VD2) $\bigwedge\limits_{a \in V} e \cdot a = a$ .

(VD3) $\bigwedge\limits_{a \in R} \bigwedge\limits_{a \in V} (-e) \cdot (a \cdot a) = ((-e) \cdot a) \cdot a =$
$$= a \cdot ((-e) \cdot a) .$$

(VD4) $\bigwedge\limits_{a,b \in V} (-e) \cdot (a + b) = (-e) \cdot a + (-e) \cdot b$ .

(OV1) $\bigwedge\limits_{a,b,c \in V} (a \leq b \Rightarrow a + c \leq b + c)$ .

$(OV2)\ \bigwedge_{a,b\in V} (a \leq b \implies (-e)\cdot b \leq (-e)\cdot a)$ .

$(OV3)\ \bigwedge_{a\in R}\ \bigwedge_{a,b\in V} (o \leq a \wedge o \leq a \leq b \implies a\cdot a \leq a\cdot b)$. $\qquad\square$

**Theorem 1:** With the above defined order relation $\leq$ and the described operations, $V_nR$ is an ordered $M_nR$-vectoid $\{V_nR, M_nR, \leq\}$ and therefore also a R-vectoid.

**Proof:** Obviously $V_nR$ fulfills the assumptions with the element $o = (o_i)$, $o_i = o$ for $i=1(1)n$. The properties (VD1) and (VD2) follow immediately from the definition of the outer operation and from (D3), (D4). Let now be $A = (a_{ij}) \in M_nR$, $a = (a_i)$, $b = (b_i)$ and $c = (c_i) \in V_nR$ . Then we get:

$(VD3):\ (-E)\cdot(A\cdot a) =$

$$= \begin{bmatrix} -e & & 0 \\ & \ddots & \\ 0 & & -e \end{bmatrix} \cdot \begin{bmatrix} \sum\limits_{j=1}^{n} a_{1j}\cdot a_j \\ \vdots \\ \sum\limits_{j=1}^{n} a_{nj}\cdot a_j \end{bmatrix} = \begin{bmatrix} (-e)\cdot \sum\limits_{j=1}^{n} a_{1j}\cdot a_j \\ \vdots \\ (-e)\cdot \sum\limits_{j=1}^{n} a_{nj}\cdot a_j \end{bmatrix} =$$

$$\underset{(D5c)}{=} \begin{bmatrix} \sum\limits_{j=1}^{n} (-e)\cdot(a_{1j}\cdot a_j) \\ \vdots \\ \sum\limits_{j=1}^{n} (-e)\cdot(a_{nj}\cdot a_j) \end{bmatrix} \underset{(D5b)}{=} \begin{bmatrix} \sum\limits_{j=1}^{n} ((-e)\cdot a_{1j})\cdot a_j \\ \vdots \\ \sum\limits_{j=1}^{n} ((-e)\cdot a_{nj})\cdot a_j \end{bmatrix} =$$

$$\underset{(D5b)}{=} \begin{bmatrix} \sum\limits_{j=1}^{n} a_{1j}\cdot((-e)\cdot a_j) \\ \vdots \\ \sum\limits_{j=1}^{n} a_{nj}\cdot((-e)\cdot a_j) \end{bmatrix} = ((-E)\cdot A)\cdot a = A\cdot((-E)\cdot a).$$

$(VD4):\ (-E)\cdot(a + b) = ((-e)\cdot(a_i + b_i)) \underset{(D5c)}{=}$

$= ((-e)\cdot a_i + (-e)\cdot b_i) = (-E)\cdot a + (-E)\cdot b$ .

$(OV1):\ a \leq b \implies \bigwedge_{i=1(1)n} a_i \leq b_i \underset{(OD1)}{\implies}$

$\bigwedge_{i=1(1)n} a_i + c_i \leq b_i + c_i \implies a + c \leq b + c$ .

$(OV2):\ a \leq b \implies \bigwedge_{i=1(1)n} a_i \leq b_i \underset{(OD2)}{\implies}$

$\bigwedge_{i=1(1)n} (-e)\cdot b_i \leq (-e)\cdot a_i \implies (-E)\cdot b \leq (-E)\cdot a$ .

$(OV3):\ (o \leq A \wedge o \leq a \leq b) \implies \bigwedge_{i,j=1(1)n} (o \leq a_{ij} \wedge$

$\wedge o \leq a_j \leq b_j) \underset{(OD3)}{\implies} \bigwedge_{i,j=1(1)n} a_{ij}\cdot a_j \leq a_{ij}\cdot b_j \underset{(OD1)}{\implies}$

$\implies \bigwedge_{i=1(1)n} \sum\limits_{j=1}^{n} a_{ij}\cdot a_j \leq \sum\limits_{j=1}^{n} a_{ij}\cdot b_j) \implies$

$\implies A\cdot a \leq A\cdot b$ . $\quad\square$

In an ordered set $\{M, \leq\}$ an isotone (resp. antitone) operator $T : M \rightarrow M$ is defined by the property·

$$\bigwedge_{a,b\in M} (a \leq b \implies Ta \leq Tb)$$

$(resp.\ \bigwedge_{a,b\in M} (a \leq b \implies Ta \geq Tb)$ . As an immediate consequence of Theorem 1 and the property

$$\bigwedge_{a,b,c\in R} (a \leq b \wedge c \geq o \implies a\cdot c \leq b\cdot c \wedge c\cdot a \leq c\cdot b),$$

which holds for linearly ordered ringoids (see [3]), we get the

**Corollary:** Let $\{V_nR, M_nR, \leq\}$ be the ordered $M_nR$-vectoid over the linearly ordered ringoid $\{R,+,\cdot,\leq\}$ and the operator $T: V_nR \rightarrow V_nR$ defined by

$Tx := A\cdot x + b,\ A \in M_nR,\ b \in V_nR$ ,

with $A \geq 0$ (resp. $A \leq 0$). Then $T$ is an isotone (resp. antitone) operator. $\quad\square$

Furthermore, for isotone (resp. antitone) matrix operators we next show the following

**Theorem 2:** Let $\{R,+,\cdot,\leq\}$ be a linearly ordered ringoid with the special elements $\{-e,o,e\}$ . Then, an element $A \in M_nR$ with $A : V_nR \rightarrow V_nR$ is isotone (resp. antitone), iff $A \geq 0$ (resp. $A \leq 0$).

**Proof:** The first part of Theorem 2 is the statement of the above Corollary, if we choose $b = o$ . Conversely, we assume, that $A = (a_{ij})$ is not non-negative. Then there exists at least one negative entry $a_{kl}$ of $A$ . Now we choose the vectors $a = o$ and $b = (b_i)$ with $b_1 = e$ and $b_i = o$ for $i \neq 1$ and obtain $A\cdot a = o$ , $(A\cdot b)_k = a_{kl}\cdot e = a_{kl} < o$ , i.e. $A\cdot a \not\leq A\cdot b$ with $a < b$. Therefore $A$ ist not isotone, which completes the proof. $\quad\square$

## 3. Cycles in ordered sets

The numerical computation of an iterative method produces a sequence, in which beginning with a certain index a finite number of different elements repeat periodically. We precise this fact in the following

<u>Definition 3:</u> Let $M$ be a non empty set. A sequence $\{x_n\}$ in $M$ is called "cyclically ending", if the property (Z) holds:

$$(Z) \quad \bigvee_{\mu, n_0 \in \mathbb{N}} \bigwedge_{n \in \mathbb{N}} (n \geq n_0 \Rightarrow x_{n+\mu} = x_n) .$$

The number $m := \min\{\mu \in \mathbb{N} \mid (Z)\}$ is called the "length of the cycle" and the set $\{x_{n_0}, x_{n_0+1}, \ldots, x_{n_0+m-1}\}$ is called "cycle of length m" (see [1]). An iterative method

$$(IT) \quad x_{n+1} = \phi(x_n) , \quad x_0 \in M , \quad n = 0,1,2,\ldots ,$$

generated by an operator $\phi : M \longrightarrow M$ is called "cyclically ending in the cycle $Z(x_0)$", if the sequence $\{x_n\}$ is cyclically ending. □

In ordered sets we get for a sequence produced by an iterative method the following for the further investigations fundamental result:

<u>Theorem 3:</u> Let $\{M, \leq\}$ be an ordered set and $\phi : M \longrightarrow M$ an isotone operator. The sequence $\{x_n\}$ of iterates produced by $\phi$ be cyclically ending in a cycle of length $m$ . Then,

$$\bigwedge_{i,j \in \mathbb{N}} (x_i \leq x_j \Rightarrow m \mid |i-j|) .$$

<u>Proof:</u> Let $\{z_1, z_2, \ldots, z_s\}$ be the cycle produced by (IT). First we show, that all cycle elements are incomparable:

Assume $z_r < z_{r+d}$ with $d > 0$ . Since $\phi$ is isotone, we have $z_{r+d} = \phi^d z_r < \phi^{2d} z_r < \ldots < \phi^{sd} z_r = z_r$ , i.e. $z_r < z_r$ , a contradiction.

If the sequence $\{x_n\}$ ends cyclically, there exists a $r \in \mathbb{N}$ with $\phi^r(x_i), \phi^r(x_j) \in Z(x_0)$ . The assumption $x_i \leq x_j$ leads to $\phi^r x_i \leq \phi^r x_j$ and therefore to $\phi^r x_i = \phi^r x_j$ . Moreover, we get with $l := j - i$

$$z = \phi^r x_i = \phi^r(x_{i+l}) = \phi^{r+l}(x_i) = \phi^l(\phi^r x_i) = \phi^l z ,$$

which completes the proof. □

The assumption of Theorem 3 the sequence of iterates being cyclically ending means no restriction for our purposes, since the set of numbers representable on a computer is finite.

Concerning Theorem 3 we still give two

<u>Remarks:</u> 1. If two consecutive iterates are comparable, we get a cycle of length 1. Especially, in linearly ordered sets we always obtain this result.

$$2. \bigwedge_{y \in M} (x_i \leq y \leq x_j \Rightarrow Z(y) = Z(x_0)) .$$

For antitone operators we derive from Theorem 3 the following

<u>Corollary:</u> Let $\{M, \leq\}$ be an ordered set and $\phi : M \longrightarrow M$ an antitone operator. The sequence $\{x_n\}$ of iterates produced by $\phi$ be cyclically ending in a cycle of length $m$ . Then we have

$$\bigwedge_{i,j \in \mathbb{N}} ((x_i \leq x_j \wedge 2 \mid |i-j|) \Rightarrow m \mid |i-j|) .$$

The assumption of comparability of iterates with even-numbered difference of indices is fulfilled especially for generalized alternating sequences, i.e. for sequences with the property

$$\bigwedge_{i \in \mathbb{N}} \min(x_i, x_{i+d}) \leq x_{i+2d} \leq \max(x_i, x_{i+d})$$

resp.

$$\bigwedge_{i \in \mathbb{N}} \min(x_{i+d}, x_{i+2d}) \leq x_i \leq \max(x_{i+d}, x_{i+2d})$$

with $d \in \mathbb{N}$ .

## 4. Applications to iterative methods

First we consider iterative methods, where Theorem 3 and the corresponding corollary lead immediately to results concerning cycles. So, let be $\{R, +, \cdot, \leq\}$ a linearly ordered ringoid with the special elements $\{-e, o, e\}$ . We express the element $A = (a_{ij}) \in M_n R$ as the matrix sum $A = L + D + R$ with

$$L = \begin{pmatrix} o & \cdots & & o \\ a_{21} & o & \ddots & \vdots \\ \vdots & & \ddots & \vdots \\ a_{n1} & \cdots & a_{n,n-1} & o \end{pmatrix}, \quad D = \begin{pmatrix} a_{11} & & 0 \\ & \ddots & \\ 0 & & a_{nn} \end{pmatrix},$$

$$R = \begin{pmatrix} o & a_{12} & \cdots & a_{1n} \\ \vdots & & \ddots & \vdots \\ \vdots & & \ddots & a_{n-1,n} \\ o & \cdots & \cdots & o \end{pmatrix} \in M_n R .$$

Now we can formulate the following well-known iterative methods in $V_n R$ with arbitrary elements $x^{(0)}, b \in V_n R$ :

(a) $x^{(k+1)} = A x^{(k)} + b ,$      total-step method

(b) $x^{(k+1)} = Lx^{(k+1)} + \{(D + R)x^{(k)} + b\}$, single-step
$$\text{method}$$

(c) $x^{(k+1)} = (\omega L)x^{(k+1)} + \{[(e - \omega)E + \omega(D + R)]x^{(k)} +$
$$+ \omega \cdot b)\}, \; \omega \in R, \; k=0,1,2,\ldots \; . \; \text{successive}$$
$$\text{relaxation method}$$

By reason of monotony properties of the operators $A,L,D,R$, we are able to verify the isotony resp. antitony of the operators defined by (a), (b) and (c):

Theorem 4: Let $\{R,+,\cdot,\leq\}$ be a linearly ordered ringoid with the special elements $\{-e,o,e\}$. Then,

1. If $A : V_n R \rightarrow V_n R$ is an isotone resp. antitone matrix, then the operator defined by (a) is also isotone resp. antitone.

2. If $L$ is an isotone and $D + R$ an isotone (resp. antitone) matrix, then the operator defined by (b) is isotone (resp. antitone).

3. If $L$ is an isotone and $D + R$ an isotone (resp. antitone) matrix, then the operator defined by (c) is isotone for $o < \omega \leq e$ (resp. antitone for $\omega \geq e$).

Proof: 1. This is the statement of the corollary of Theorem 2.

2. With $L$ and $D + R$ isotone follows for the operator $Tx = L(Tx) + (D + R)x + b$:

(*) $x \leq y \Rightarrow (D + R)x + b \leq (D + R)y + b$ and therefore $(Tx)_1 \leq (Ty)_1$ . Now suppose

(**) $\bigwedge\limits_{i=1(1)k < n} (Tx)_i \leq (Ty)_i$ .

Then follows from (**) and $L$ isotone :
$(L(Tx))_{k+1} \leq (L(Ty))_{k+1}$
and with (*)
$(Tx)_{k+1} \leq (Ty)_{k+1}$ , i.e. $T$ is isotone.

The proof of the second assertion turns out analogously.

3. In a linearly ordered ringoid we have the property $\bigwedge\limits_{a \in R} a - a = o$ . Hence, we get

$o < \omega \leq e \underset{(OD2)}{\Rightarrow} -e \leq -\omega \leq o \underset{(OD1)}{\Rightarrow} o = e - e \leq e-\omega \leq e.$

Now it is clear, that $\omega L$, $\omega(D + R)$, $(e - \omega)E$ , $(e - \omega)E + \omega(D + R)$ are isotone operators, if $L$ and $D + R$ are isotone, and therefore 3. is reduced to 2.

Similarly, we verify the second assertion by applying

$e \leq \omega \underset{(OD2)}{\Rightarrow} -\omega \leq -e \Rightarrow e - \omega \underset{(OD1)}{\leq} e - e = o . \quad \square$

The operators considered in Theorem 4 fulfill the

assumptions of monotony in Theorem 3 and the corresponding corollary, i.e. we are able to determine the cycle length of the iteration sequences generated by the methods (a), (b) and (c), if comparable iterates occur. Consequently, the problem consists in starting the iteration with an element $x^{(o)}$, such that comparable iterates are produced with indices as small as possible. Because of this difficulty we will consider the question, whether the assumption of the comparability can be avoided for certain operators. To this, we generalize the concept of the weakly cyclic matrices of index $k$ to vector functions by

Definition 4: Let $\{R,+,\cdot,\leq\}$ be a linearly ordered ringoid. A vector function $F = (f_i): B \subseteq V_n R \rightarrow V_n R$ is called "weakly cyclic of index $k$" , if the set $N := \{1,2,\ldots,n\}$ is partitioned into $k$ disjoint subsets $N_1,\ldots,N_k$ with $\bigcup\limits_{i=1}^{k} N_i = N$ and the following property is valid:

$\bigwedge\limits_{i=1(1)k} \bigwedge\limits_{1 \in N_i} f_1 = f_1(x_{N_{P(i)}})$ with

$x_{N_i} := \{x_j | j \in N_i\}$ , where $P$ is a cyclic permutation of $\{1,2,\ldots,k\}$ . If $P$ is a permutation of $\{1,2,\ldots,k\}$ partitioned into $r$ cyclic permutations of $k_i$, $i = 1(1)r$, elements, the vector function $F$ is called "$(k_1,k_2,\ldots,k_r)$-weakly cyclic of index $k$".

Theorem 5: Let $\{R, \leq\}$ be an ordered set, $F = (f_i):$ $V_n R \rightarrow V_n R$ a weakly cyclic vector function of index $k$ . Further $IS$ and $AN$ be disjoint subsets of $\{1,\ldots,k\}$ with $IS \cup AN = \{1,\ldots,k\}$ and $\bigwedge\limits_{i \in IS} \bigwedge\limits_{1 \in N_i} f_1$ isotone and $\bigwedge\limits_{i \in AN} \bigwedge\limits_{1 \in N_i} f_1$ antitone.

Then, $F^k: V_n R \rightarrow V_n R$ is an isotone operator, if $\#AN$ is even, and an antitone operator, if $\#AN$ is odd.

Proof: By means of the definition of the weakly cyclic vector function we get by $k$ applications of $F$ to an element $a \in V_n R$

$$F^k(a) = (\prod\limits_{1=o}^{k-1} (f_{N_{p^1(i)}})(a)) .$$

Therefore, $F^k$ is isotone, if $\#AN$ is even, and antitone, if $\#AN$ is odd. $\square$

In the next theorem we extend the results of Theorem 5 to $(k_1,\ldots,k_r)$-weakly cyclic vector functions of index $k$ .

Theorem 6: Let $\{R, \leq\}$ be an ordered set, $F = (f_i):$

$V_nR \to V_nR$ a $(k_1,...,k_r)$-weakly cyclic vector function of index $k$ and $p := 2 \cdot$ s.c.m. $(k_1,...,k_r)$. Let further IS and AN be disjoint subsets of $\{1,2,...,k\}$ with $IS \cup AN = \{1,2,...,k\}$ and $\bigwedge_{i \in IS} \bigwedge_{1 \in N_i} f_1$ isotone and $\bigwedge_{i \in AN} \bigwedge_{1 \in N_i} f_1$ antitone.

Then, $F^p$ is an isotone operator.

Proof: We consider the partitioning of $F$ into the $r$ weakly cyclic vector functions $F_i$ of index $k_i$, $i=1(1)r$. Applying Theorem 5, the $2k_i$-fold application of $F_i$ is isotone in any case. Therefore the s.c.m. $(2k_1,...,2k_r)$-fold application of $F$ is isotone, which yields immediately the statement. $\square$

To give an idea of the magnitude of the number 1 in Theorem 6, we give an estimation in the following

Corollary: Let be given the assumptions of Theorem 6. Then $F^p$ is isotone with

$$p \leqq 2 \cdot g(k) \text{ and } g(k) = \begin{cases} 3^{\frac{k}{3}} & k \equiv o \bmod 3 \\ 4 \cdot 3^{[\frac{k}{3}-1]} & k \equiv 1 \bmod 3 \\ 2 \cdot 3^{[\frac{k}{3}]} & k \equiv 2 \bmod 3 \end{cases}$$

Proof: The first part of the statement follows from Theorem 5 applied to the $r$ weakly cyclic vector functions $F_i$ of index $k_i$, $i = 1(1)r$. It is

$$m_{Max} \leqq \max_{\substack{r \\ \sum_{i=1}^{r} k_i = k, \, r \leqq k}} 2 \cdot \text{s.c.m.} (k_1, k_2, ..., k_r)$$

and hence,

$$m_{Max} \leqq 2 \cdot \max_{\substack{r \\ \sum_{i=1}^{r} k_i = k}} \prod_{i=1}^{r} k_i \, , \, r \leqq k \, .$$

For the determination of that additive partitioning $k_1,...,k_r$ of $k$, for which the product $\prod_{i=1}^{r} k_i$ becomes maximal, the following considerations are possible.

Let be $k = k_1 + k_2$, $k, k_1, k_2 \in \mathbb{N}$. Then,

$$k_1 + k_2 < k_1 \cdot k_2 \, , \, \text{if} \, k_1 \geqq 3 \, \text{and} \, k_2 \geqq 2 \, .$$

Hence, we get the estimation

$$5 \leqq k = k_1 + k_2 \leqq k_1 \cdot k_2 = (k_{11} + k_{12}) \cdot k_2 \leqq$$
$$\leqq k_{11} \cdot k_{12} \cdot k_2 \leqq \cdots \, .$$

The product on the right side of this inequality increases, if $k$ is further partitioned in terms not less 2. So we get a sum of the numbers 2 and 3.

Since $2 + 2 + 2 = 6 = 3 + 3$ and $2^3 < 3^2$, we get the maximal product by a partition into the numbers 2 and 3, where the number 2 does not occur more than twice. Now expression $g(k)$ follows immediately. $\square$

For $k = 5(1)12$ the function $g(k)$ assumes the following values:

| k | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|----|----|----|
| g(k) | 12 | 18 | 24 | 36 | 54 | 72 | 108 | 162 |

The results of Theorem 5 and Theorem 6 resp. of the corresponding Corollary are now applied to matrices

Theorem 7: Let $\{R,+,\cdot,\leqq\}$ be a linearly ordered ringoid and $A=(a_{ij}) \in M_nR$ a weakly cyclic matrix of index $n$. With $\#AN = \#\{a_{ij}|a_{ij} < o\}$ hold

1. If the iterative method defined by (a) is cyclically ending, then the cycle length $m$ divides $n$ resp. $2n$ for $\#AN$ even resp. $\#AN$ odd.

2. If $A$ is given in normal form and the iterative method (b) is cyclically ending, then the cycle has length 1 resp. 2 for $\#AN$ even resp. $\#AN$ odd.

Proof: 1. From Theorem 5 follows, that $T^n$ with $Tx := Ax + b$ is an isotone operator for $\#AN$ even resp. an antitone operator for $\#AN$ odd, i.e. the component functions of $T^n$ are isotone resp. antitone. Therefore, the component sequences of the iterates $x^{(0)}, x^{(n)}, x^{(2n)}, ...$ resp. $x^{(0)}, x^{(2n)}, x^{(4n)}, ...$ are monotone. The application of Theorem 3 gives the desired result.

2. It is

$$L = \begin{bmatrix} o & & & \\ a_1 & \ddots & & 0 \\ & a_2 & & \\ 0 & & \ddots & \\ & & & a_{n-1} \, o \end{bmatrix}, \, D + R = \begin{bmatrix} o & \cdots & o & a_n \\ & & & o \\ & & \ddots & \vdots \\ 0 & & \ddots & \vdots \\ & & & o \end{bmatrix}$$

and therefore,

$$Tx = \begin{cases} a_n x_n + b_1 \\ a_1(a_n x_n + b_1) + b_2 \\ \vdots \\ a_{n-1}(a_{n-2}(...a_1(a_n x_n + b_1) + b_2) + ...) + b_n \end{cases} =:$$

$$=: (t_i(x_n)) \, .$$

Analogously, we get

$$T^2 x = (t_i(t_n(x_n))), \quad T^3 x = (t_i(t_n^2(x_n))), \ldots, \quad T^{k+1} x = (t_i(t_n^k(x_n))), \ldots .$$ Since $\{R,+,\cdot,\leq\}$ is a linearly ordered ringoid, we have $x_n \leq t_n x_n$ resp. $x_n \geq t_n x_n$. Hence, $t_n$ isotone yields

$t_n^{k-1} x_n \leq t_n^k x_n$ resp. $t_n^{k-1} x_n \geq t_n^k x_n$. If $t_n$ is antitone, we obtain from $x_n \leq t_n^2 x_n$ resp. $x_n \geq t_n^2 x_n$

the inequality $\quad t_n^{2(k-1)} x_n \leq t_n^{2k} x_n$ resp.

$t_n^{2(k-1)} x_n \geq t_n^{2k} x_n$. By Theorem 3 the sequence in each component $i$, $i=1(1)n$, has the cycle length 1, if the operator $t_n$ is isotone, and the cycle length 1 or 2, if $t_n$ is antitone. By the definition of $t_n$ it is clear, that for #AN even resp. odd $t_n$ is isotone resp. antitone. □

Remark: Theorem 6 gives us the hint that the determination of cycles is not only possible for weakly cyclic matrices. We consider a matrix $A$ which has only one entry not equal $o$ in each row and in each column. Then, by similar considerations the iterative method (a) generates a cycle of length $m$ with $m | 2 \cdot$ s.c.m. $(n_1, \ldots n_r)$, where $n_i$, $i = 1(1)r$ are the numbers of elements in the $r$ cyclic submatrices of $A$. In the following table we give for some values of $n$ the greatest attainable cycle length $m_{Max}$, the corresponding partition of $n$ and the estimation $m_E$ computed by $g(n)$ (see Corollary of Theorem 6).

| $n$ | $m_{Max}$ | partition of $n$ | $m_E$ |
|-----|-----------|------------------|-------|
| 2 | 4 | 2 | 4 |
| 3 | 6 | 3 | 6 |
| 4 | 8 | 4 | 8 |
| 5 | 12 | 2,3 | 12 |
| 6 | 12 | 6 | 18 |
| 7 | 24 | 3,4 | 24 |
| 8 | 30 | 3,5 | 36 |
| 9 | 40 | 4,5 | 54 |
| 10 | 60 | 2,3,5 | 72 |
| 11 | 60 | 1,2,3,5, | 108 |
| 12 | 120 | 3,4,5 | 162 |

After we have proved results concerning cycles produced by iterative methods for the solution of systems of equations, we finally want to include a theorem concerning iterative methods for the determination of zeros of real functions.

Theorem 8: Let $\{R,N,+,\cdot,/,\leq\}$ be a linearly ordered division ringoid with the special elements $\{-e,o,e\}$ and $f,g: R \rightarrow R$ mappings, where $f$ is isotone (resp. antitone) and $g$ is antitone (resp. isotone) in $J \subseteq R$ and $g > o$ (resp. $g < o$) in $J$. Furthermore, assume that the operator defined by $Tx := x - f(x)/g(x)$, $g(x) \notin N$, generates a sequence $\{x_n\} \subseteq J$, starting with $x_0 \in J$. Then, if the sequence $\{x_n\}$ is cyclically ending, the cycle length may be 1 or 2.

Proof: Let $\{z_1, \ldots, z_s\}$ be the cycle of the sequence $\{x_n\}$ with $s > 2$. Since $\{R, \leq\}$ is linearly ordered, with a permutation $P$ we can write

$$z_{P(1)} < z_{P(2)} < \cdots < z_{P(s)} .$$

From $f$ isotone and $g$ antitone and $g > o$ we obtain

$$f(z_{P(1)}) < f(z_{P(2)}) < \cdots < f(z_{P(s)})$$

and $\quad g(z_{P(1)}) > g(z_{P(2)}) > \cdots > g(z_{P(s)}) > o$.

With the property

$$\bigwedge_{a,b,c \in R} (a \leq b \wedge c \geq o \implies a/c \leq b/c) ,$$

which is an immediate consequence of (OD4), we get

$$f(z_{P(1)})/g(z_{P(1)}) < \cdots < f(z_{P(s)})/g(z_{P(s)}).$$

With (OD1) and (OD2) we have

$$z_{P(1)} - f(z_{P(1)})/g(z_{P(1)}) > \cdots > z_{P(s)} - f(z_{P(s)})/g(z_{P(s)})$$

and therefore,

$$z_{P(s)} = T z_{P(1)}$$

$$z_{P(1)} = T z_{P(1)} ,$$

i.e. $\{z_{P(1)}, z_{P(s)}\}$ forms a cycle of length 2, in contradiction to the assumption. The other case is proved analogously. □

## 5. Numerical examples

Let $\hat{R} := \{x \in \mathbb{R} \,|\, |x| \leq s\}$ be a bounded subset of the real number field. Then, a finite symmetric subset $\hat{R}_M$ of $\hat{R}$ with the property $-s, s \in \hat{R}_M$ is a symmetric screen of $\{\hat{R}, \leq\}$ (see [4]). Furthermore, $\{\hat{R}, \leq\}$ forms a completely, linearly ordered ringoid by using the real arithmetic. With a monotone, antisymmetric rounding

$$\Box : \hat{R} \to \hat{R}_M \ ,$$

i.e. an isotone, idempotent mapping from $\hat{R}$ to $\hat{R}_M$ with the property

$$\bigwedge_{a \in \hat{R}} \Box(-a) = -\Box a \ ,$$

and the inner operations $\boxdot$ defined by

$$\bigwedge_{a,b \in \hat{R}_M} a \boxdot b := \Box(a * b), \ * \in \{+, \cdot\} \ ,$$

$\hat{R}_M$ forms a completely, linearly ordered ringoid $\{\hat{R}_M, \boxplus, \boxdot, \leq\}$ (see [3]). By analogous considerations we obtain a completely, linearly ordered division ringoid $\{\hat{R}_M, \{o\}, \boxplus, \boxdot, \boxslash, \leq\}$

Thus, $\{\hat{R}_M, \boxplus, \boxdot, \leq\}$ resp. $\{\hat{R}_M, \{o\}, \boxplus, \boxdot, \boxslash, \leq\}$ and $\{V_n \hat{R}_M, M_n \hat{R}_M, \leq\}$ fulfill the assumptions of the theorems proved in the precedent sections.

Now we replace $\hat{R}_M$ by the set $\hat{R}_{m,b}$ of the normalized floating-point numbers $\hat{R}_{m,b} :=$ $\{x \in \mathbb{R} \,|\, x = * o.d_1 d_2 \ldots d_m \cdot b^e, \ * \in \{+, -\}$ , $d_i \in \{o, 1, \ldots, b-1\}, \ i = 1(1)m, \ d_1 \neq o,$ $e_{min} \leq e \leq e_{max}, \ e \in \mathbb{Z}\} \cup \{o\}$ with $m, b \in \mathbb{N}, \ b \geq 2$. Then, the set $\hat{R}$ is given by $\hat{R} :=$ $\{x \in \mathbb{R} \,|\, |x| \leq o.d_1 d_2 \ldots d_m \cdot b^{e_{max}}, d_i = b-1, \ i = 1(1)m\}$ . Furthermore, we apply the frequently used monotone, antisymmetric rounding

$$\Box a := \begin{cases} \nabla a \text{ if } a < (\nabla a + \triangle a)/2 \text{ or } a = (\nabla a + \triangle a)/2 \wedge \\ \qquad\qquad\qquad\qquad\qquad a \leq o \\ \triangle a \text{ if } a > (\nabla a + \triangle a)/2 \text{ or } a = (\nabla a + \triangle a)/2 \wedge \\ \qquad\qquad\qquad\qquad\qquad a > o \end{cases}$$

(rounding to the nearest number of $\hat{R}_{m,b}$), where $\nabla$ and $\triangle$ mean the uniquely defined monotone, directed roundings, ([5]).

The examples 2. and 4. has been performed on the digital computer EL X8 of the Computer Center of the University of Karlsruhe. Note, that the decimal numbers of $\hat{R}_{12,1o}$ used for the representation of the operators and cycles have been generated by conversion of binary numbers of $\hat{R}_{4o,2}$. We cannot expect, that the described operators produce the same cycles by computation in $\hat{R}_{12,1o}$ .

Examples:

1. Let be $n = 2$, $b = 1o$, $m = 2$. The operator

$$Tx = \begin{bmatrix} o & -o.78 \\ o.93 & o \end{bmatrix} \boxdot x \boxplus \begin{bmatrix} 1.2 \\ o.75 \end{bmatrix} \ ,$$

generates a sequence $\{x^{(n)}\}$ by $x^{(n+1)} :=$ $Tx^{(n)}, \ x^{(o)} \in V_2 \hat{R}_{2,1o}$. For different starting vectors $x^{(o)}$ the sequence is ending in different cycles:

$$Z_1 = \left\{ \begin{bmatrix} o.26 \\ 1.2 \end{bmatrix} , \begin{bmatrix} o.26 \\ o.99 \end{bmatrix} , \begin{bmatrix} o.43 \\ o.99 \end{bmatrix} , \begin{bmatrix} o.43 \\ 1.2 \end{bmatrix} \right\}$$

$$Z_2 = \left\{ \begin{bmatrix} o.34 \\ 1.2 \end{bmatrix} , \begin{bmatrix} o.26 \\ 1.1 \end{bmatrix} , \begin{bmatrix} o.34 \\ o.99 \end{bmatrix} , \begin{bmatrix} o.43 \\ 1.1 \end{bmatrix} \right\}$$

$$Z_3 = \left\{ \begin{bmatrix} o.34 \\ 1.1 \end{bmatrix} \right\} \quad .$$

According to Theorem 7(1) we get a maximal cycle length 4 and by Theorem 7 (2) the single-step method is ending in a cycle of length 1 or 2:

$$Z_3 \text{ and } Z_4 = \left\{ \begin{bmatrix} o.26 \\ o.99 \end{bmatrix} , \begin{bmatrix} o.43 \\ 1.2 \end{bmatrix} \right\} \quad .$$

Applying the monotone rounding towards zero we get beginning with $x^{(o)} = (o.26, 1.2)^T$ the cycle
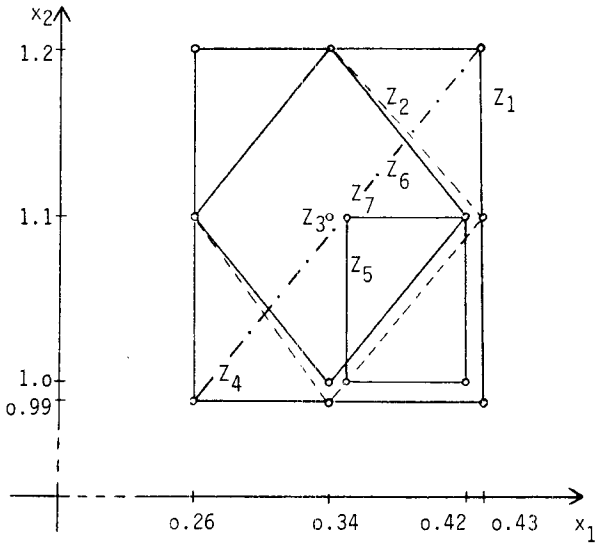
$$Z_5 = \left\{ \begin{bmatrix} o.42 \\ 1.1 \end{bmatrix} , \begin{bmatrix} o.35 \\ 1.1 \end{bmatrix} , \begin{bmatrix} o.35 \\ 1.o \end{bmatrix} , \begin{bmatrix} o.42 \\ 1.o \end{bmatrix} \right\}$$

and applying the monotone rounding away from zero we get beginning with $x^{(o)} = (o.27, 1.2)^T$ the cycle

$$Z_6 = \left\{ \begin{bmatrix} o.26 \\ 1.1 \end{bmatrix} , \begin{bmatrix} o.34 \\ 1.o \end{bmatrix} , \begin{bmatrix} o.42 \\ 1.1 \end{bmatrix} , \begin{bmatrix} o.34 \\ 1.2 \end{bmatrix} \right\}$$

and the cycle $Z_7 = \left\{ \begin{bmatrix} o.34 \\ 1.1 \end{bmatrix} \right\}$ .

We illustrate the cycles $Z_1$ to $Z_7$ graphically:

with $a_1 = -0.7551928744688$    $b_1 = 0.6305774756001$

$a_2 = 0.1032844094188$    $b_2 = 0.5755482108107$

$a_3 = 0.6376119369179$    $b_3 = 0.265077o197379$

$a_4 = 0.4747402850753$    $b_4 = 0.1190302626437$

$a_5 = 0.8622977826562$    $b_5 = -0.03283479505717$

With the elements

$x_1 = 0.1459944151620$    $y_1 = 0.1459944151626$

$x_2 = 0.5906271577596$    $y_2 = 0.5906271577587$

$x_3 = 0.6416679457934$    $y_3 = 0.6416679457925$

$x_4 = 0.07339623131165$    $y_4 = 0.07339623131184$

$x_5 = -0.0961242o257263$    $y_5 = -0.0961242o257272$

we can describe the generated cycle in the following way:

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix}, \begin{bmatrix} x_1 \\ y_2 \\ x_3 \\ y_4 \\ x_5 \end{bmatrix}, \begin{bmatrix} x_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ x_4 \\ y_5 \end{bmatrix}, \begin{bmatrix} y_1 \\ x_2 \\ y_3 \\ x_4 \\ x_5 \end{bmatrix}, \begin{bmatrix} y_1 \\ x_2 \\ x_3 \\ y_4 \\ y_5 \end{bmatrix}, \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ y_4 \\ y_5 \end{bmatrix}, \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ y_5 \end{bmatrix}, \begin{bmatrix} x_1 \\ y_2 \\ x_3 \\ x_4 \\ x_5 \end{bmatrix}, \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ x_5 \end{bmatrix},$$

$$\begin{bmatrix} y_1 \\ x_2 \\ y_3 \\ y_4 \\ y_5 \end{bmatrix}, \begin{bmatrix} y_1 \\ x_2 \\ y_3 \\ x_4 \\ y_5 \end{bmatrix} .$$

2. Another example of Theorem 7 (1) is given by the operator

$$Tx := \begin{bmatrix} 0 & 0 & -0.6638554373967 \\ -0.6686745454153 & 0 & 0 \\ 0 & -0.6795139955248 & 0 \end{bmatrix} \boxdot x \boxplus$$

$$\boxplus \begin{bmatrix} 4.235750979760 \\ -2.577679425907 \\ 0.7134761463822 \end{bmatrix}$$

in $V_3\hat{R}_{13,10}$, which produces e.g. the following cycle

$$\begin{bmatrix} 1.996957891017 \\ -3.912994335897 \\ 3.372410562035 \end{bmatrix}, \begin{bmatrix} 1.996957891017 \\ -3.912994335897 \\ 3.372410562031 \end{bmatrix}, \begin{bmatrix} 1.996957891020 \\ -3.912994335897 \\ 3.372410562031 \end{bmatrix},$$

$$\begin{bmatrix} 1.996957891020 \\ -3.912994335900 \\ 3.372410562031 \end{bmatrix}, \begin{bmatrix} 1.996957891020 \\ -3.912994335900 \\ 3.372410562035 \end{bmatrix}, \begin{bmatrix} 1.996957891017 \\ -3.912994335900 \\ 3.372410562035 \end{bmatrix},$$

$$\begin{bmatrix} 1.996957891017 \\ -3.912994335897 \\ 3.372410562035 \end{bmatrix}, \ldots .$$

3. Finally, we give an operator in $V_5\hat{R}_{40,2}$, which generates according to the remark of Theorem 7 a cycle of maximum length 12:

$$Tx := \begin{bmatrix} 0 & 0 & a_1 & 0 & 0 \\ a_2 & 0 & 0 & 0 & 0 \\ 0 & a_3 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & a_4 \\ 0 & 0 & 0 & a_5 & 0 \end{bmatrix} \boxdot x \boxplus \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \\ b_5 \end{bmatrix}$$

References:

[1] COLLATZ, L.: Funktionalanalysis und numerische Mathematik, Springer Verlag, Berlin 1964

[2] KLATTE, R.: Zyklisches Enden bei Iterationsverfahren. Dr.-Dissertation, Universität Karlsruhe, 1975

[3] KULISCH, U.: Rounding Invariant Structures. Mathematics Research Center, University of Wisconsin, Technical Summary Report # 1103, September 1970, 1-47

[4] KULISCH, U.: An Axiomatic Approach to Rounded Computations. Numerische Mathematik 18, 1-17 (1971)

[5] KULISCH, U.: Implementation and Formalization of floating-point arithmetic. To appear in Computing

[6] ULLRICH, Chr.: Rundungsinvariante Strukturen mit äußeren Verknüpfungen. Dr.-Dissertation , Universität Karlsruhe , 1972

[7] VARGA, R.: Matrix Iterative Analysis, Prentice Hall, Inc., Englewood Cliffs, New Jersey 1962

Address:

Dr. R. Klatte, Dr. Chr. Ullrich
Institut für Angewandte Mathematik
Universität Karlsruhe

D-75 Karlsruhe
Kaiserstr. 12
West Germany