AD-A089 636    WISCONSIN UNIV-MADISON MATHEMATICS RESEARCH CENTER    F/G 12/1
                REPRESENTATIONS OF INTERVALS AND OPTIMAL ERROR BOUNDS.(U)
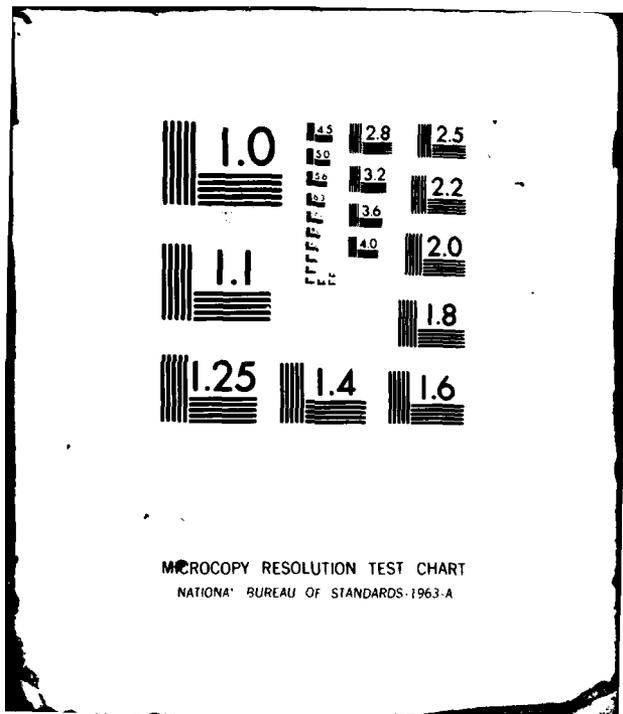                JUL 80    L B RALL                      DAA629-80-C-0041
UNCLASSIFIED    MRC-TSR-2098                                    NL

END
DATE
FILMED
10 80
DTIC

MICROCOPY RESOLUTION TEST CHART
NATIONAL BUREAU OF STANDARDS-1963-A

MRC Technical Summary Report #2098

REPRESENTATIONS OF INTERVALS
AND OPTIMAL ERROR BOUNDS

L. B. Rall

**Mathematics Research Center**

**University of Wisconsin—Madison**

**610 Walnut Street**

**Madison, Wisconsin 53706**

July 1980

(Received June 30, 1980)

**Approved for public release**
**Distribution unlimited**

80 9 24 043

UNIVERSITY OF WISCONSIN - MADISON
MATHEMATICS RESEARCH CENTER

REPRESENTATIONS OF INTERVALS AND OPTIMAL ERROR BOUNDS

L. B. Rall

Technical Summary Report #2098
July 1980

ABSTRACT

Calculations with interval arithmetic and interval extensions of real func-
tions are often used to obtain lower and upper bounds for what would be the
theoretical results of precise computations on exact data. A more traditional
way to represent approximate data and results contained in intervals is by means
of a representative point x in the interval and an associated measure of error
E. For example, the midpoint of an interval has absolute error as an approxi-
mation to other points of the interval which is bounded by half the width of
the interval. In general, the chosen point is said to be optimal with respect
to the measure of error if the maximum value of E over the interval is mini-
mized. A general method for determination of optimal points x* and minimal
error bounds E* is given. In particular, the harmonic mean of the endpoints
of positive intervals is shown to be optimal with respect to relative (or per-
centage) error, and the geometric mean plays the same role with respect to rela-
tive precision, which is a measure of error introduced recently by F. W. J. Olver
[SIAM J. Numer. Anal. 15 (1978), 368-393]. In addition, the optimal point and
error bound x*, E* may be regarded as alternative coordinates for representa-
tion of an interval. Explicit rules for interval arithmetic are given in terms
of the arithmetic, geometric, and harmonic means of the endpoints and the asso-
ciated optimal error bounds. Comparisons are also made between the results of
exact or rounded interval arithmetic and the a priori estimates of relative
precision presented in the cited paper by Olver, which show that the intervals
resulting from calculation with interval arithmetic can be expected to be
smaller than predicted.

## SIGNIFICANCE AND EXPLANATION

The numerical solution of practical problems almost always involves performing inaccurate calculations on inexactly known data. Thus, in addition to the answers produced by a computer, some indication of their reliability is required. Such an estimate of accuracy may be merely to justify the expense of the computation; however, in cases where the results have implications for human life and safety, as in the design of critical components of aircraft, a guarantee of reliability is imperative. One way to obtain numerical results with such assurance is by interval computation. If the data x are known to lie in an interval I, and one computes an output interval $J = F(I)$, where F is an interval transformation which is an interval extension of a real function f, then the output interval J will contain the values $y = f(x)$ of exact transformations of the data. Ordinarily, interval calculations are done in terms of lower and upper bounds; $I = [a,b]$ is transformed into $J = [c,d]$. The rules for interval calculation are customarily expressed for intervals in this standard form, and in many applications, lower and upper bounds for the exact results are satisfactory. Often, however, one prefers to think of a single number x and its associated error bound E, that is, the maximum error with which x approximates other points of the interval. In particular, x* is optimal with respect to E if E is minimized on I, the corresponding value E* being the optimal error bound obtainable. For example, $x* = m[a,b] = (a+b)/2$ is optimal with respect to absolute error, and $E* = (b-a)/2$. In other problems, one may wish to optimize different measures of error, such as percentage error or the recently introduced concept of relative precision. A general method for determination of x*,E* is given in this paper. For positive intervals, it turns out that the harmonic mean $x* = h[a,b] = 2ab/(a+b)$ is optimal with respect to percentage error, while the geometric mean $x* = g[a,b] = \sqrt{ab}$ of the endpoints is optimal with respect to relative precision. The numbers x*,E* may be used instead of the endpoints a,b as an alternate representation (or format) for the interval I. Rules for interval arithmetic in some of these alternative formats are given. Finally, it is shown that exact or rounded interval arithmetic produces intervals which are usually smaller than predicted by a priori error estimation, and in no case larger. The results presented are intended to increase the usefulness of interval methods of error estimation.

---

# REPRESENTATIONS OF INTERVALS AND OPTIMAL ERROR BOUNDS

## L. B. Rall

1. __Intervals and interval analysis__. Just as real and complex numbers are the basic units of real and complex analysis, respectively, the closed finite real intervals

(1.1) $$I = \{x \mid a \le x \le b\}, \quad -\infty < a \le b < +\infty,$$

are the basic units in the branch of mathematics known as __interval analysis__ [4], [5], which is the study of transformations of sets of the type (1.1), which will be called simply __intervals__ for brevity, into others. Geometrically, the set IR of all intervals (1.1) may be visualized as the closed halfplane in the a,b-plane which includes and is bounded below an on the right by the line $a = b$ (see Figure 1).



Figure 1.  IR as a subset of the a,b-plane.
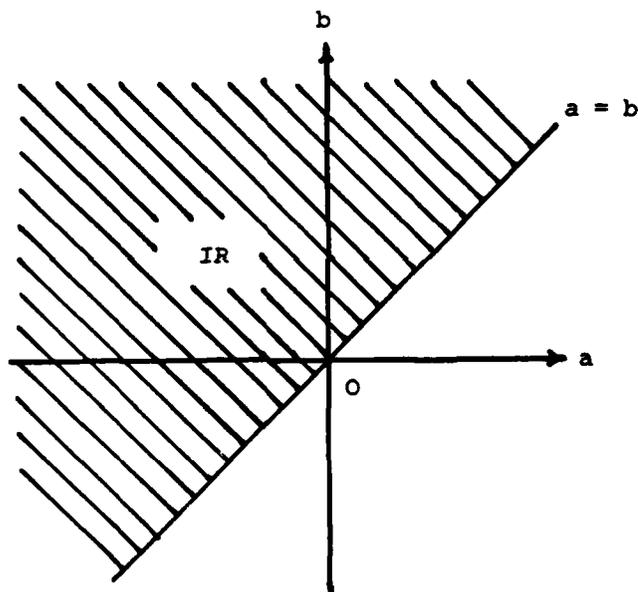
Thus, IR is two-dimensional. As in the case of complex numbers, which have representations in rectangular coordinates as $z = x + iy$, or in polar coordinates as $z = re^{i\theta}$, various representations are possible for intervals. Some of these representation, which are convenient in one situation or another for computation an error estimation, will be presented in this paper.

2.  <u>Interval extensions and error estimation</u>.  Interval analysis has important applications to numerical analysis in connection with error estimation.  Suppose that one wishes to calculate

(2.1)                    $y = f(x)$,

where $f: D \subset R \to R$ is a continuous real function on its domain of definition D.  It may happen that x is not known exactly, or is not representable exactly on a given computer; however, it is known that $x \in I$, where I is an interval which can be represented exactly.  The set

(2.2)                    $f(I) = \{y \mid y = f(x), x \in I\}$

will be an interval by the continuity of f, but it may not be possible to compute or represent $f(I)$ exactly.  For example, $f(x)$ may be defined by an infinite process (series, product, integral, etc.) so that a finite number of computational operations will give an approximation to $f(x)$ which is accurate up to some <u>truncation error</u>, and there will of course be <u>roundoff error</u> in the actual calculations.

Interval analysis deals with the problem of error estimation by the use of interval transformations $F: D \subset IR \to IR$ which are <u>computable</u> (or <u>rounded</u>) <u>interval extensions</u> of real functions f.  Thus, one actually computes the interval J defined by

(2.3)                    $J = F(I)$,

where F has the following properties:

(i)  <u>Representation</u>.  $F(I)$ is exactly representable for each representable interval I;

(ii)  <u>Extension</u>.  F is an extension of f in the sense that $f(I) \subset F(I)$ for each interval I contained in the domain of definition of f;

(iii)  <u>Inclusion monotonicity</u>.  F is <u>inclusion monotone</u>, which means that for intervals $I_1, I_2$ in the domain of definition of F, $I_1 \subset I_2 \Rightarrow F(I_1) \subset F(I_2)$.

The interval J may be taken to be the answer to the problem of evaluation of the real transformation (2.1) by inaccurate calculations on inexact data, as it contains the results of all exact values of f for all exact data points $x \in I$.

In traditional numerical analysis, it is customary to think in terms of the <u>error</u> (or more positively, <u>accuracy</u> or <u>precision</u>) of an approximate result, rather than use intervals directly.  Thus, one picks a representative point x from the interval, and assigns a measure of error to that point as a way of representing the interval.  The interval I is thought of as arising from errors

in the data, and the enlargement of $f(I)$ to $F(I)$ as being due to _transformation errors_ (i.e., truncation and roundoff errors, called _abbreviation errors_ by Olver [6]). If the representative point is picked in such a way that the maximum value of the error over the entire interval is minimized, then this point and the resulting error bound are said to be _optimal_. The representation of intervals by optimal points and error bounds will be illustrated by an example in the next section, following which the general theory will be discussed.

3. _An example of an optimal point and error bound._ A simple interval method for numerical integration [1] applied to

(3.1)
$$y = \int_0^2 \sqrt{1 + 4x}\ dx$$

yields lower and upper bounds for the value of the integral,

(3.2)
$$4.3306334 \le \int_0^2 \sqrt{1 + 4x}\ dx \le 4.3452908.$$

The _standard format_ for the representation of an interval (1.1) is, of course,

(3.3)
$$I = [a\ ,\ b]$$

in terms of its _lower_ endpoint

(3.4)
$$a = \ell(I) = \min\{x \mid x \in I\},$$

and _upper_ endpoint

(3.5)
$$b = u(I) = \max\{x \mid x \in I\}.$$

Thus, (3.2) may be expressed as

(3.6)
$$\int_0^2 \sqrt{1 + 4x}\ dx \ \epsilon\ [4.3306334\ ,\ 4.3452908].$$

Another way to write (3.2) is as

(3.7)
$$\int_0^2 \sqrt{1 + 4x}\ dx \simeq 4.3379621 \pm 0.0073287$$

in terms of the _midpoint_

(3.8)
$$m(I) = m[a\ ,\ b] = \frac{a + b}{2}$$

and the _halfwidth_ (or _radius_)

(3.9)
$$\alpha(I) = \alpha[a\ ,\ b] = \frac{b - a}{2}$$

of the interval $I$. This gives

(3.10)
$$I = m(I) \pm \alpha(I) = [m(I) \pm \alpha(I)]$$

as an alternative format for the representation of the interval $I$. In traditional terminology, (3.7) states that one may take $m(I) = 4.3379621$ as an

approximate value for the integral (3.1), with absolute error bounded by $\alpha(I)$
= 0.0073287. There are also other measures of error of interest in various
problems, such as percentage or relative error, for which the harmonic point
of the interval is optimal, or the newly introduced concept of relative pre-
cision [6], for which the geometric point plays the same rôle [7]. These lead
to alternative representations for intervals. The general situation will now
be discussed, after which some examples of useful interval formats and the cor-
responding rules of interval arithmetic will be given.

4. Optimal points and error estimates. Precise definitions will now be made
of the concepts of measures of error and corresponding optimal points and er-
ror bounds in intervals.

Definition 4.1. A measure of error of the approximation of a real number
y by a real number x is a non-negative continuous function $E = E(x,y)$ which is
strictly monotone increasing with respect to the distance

(4.1)
$$d(x,y) = |x - y|$$

between x and y and such that

(4.2)
$$E(x,x) = 0.$$

Definition 4.2. The error of approximation of points of an interval I by
a real number x corresponding to a given measure of error E is

(4.3)
$$\epsilon(x) = \max_{y \in I}\{E(x,y)\}.$$

Definition 4.3. A point $x^* \in I$ is said to be optimal in I with respect to
the measure of error E if

(4.4)
$$\epsilon(x^*) = \min_{x \in I}\{\epsilon(x)\} = \min_{x \in I} \max_{y \in I}\{E(x,y)\}.$$

Although the following result is self-evident, it will be dignified with
a proof.

Remark 4.1. The optimal point $x^* \in I$ satisfying (4.4) exists, is unique,
and is the solution of the equation

(4.5)
$$E(x,a) = E(x,b).$$

Proof: If a = b, then $x^* = a = b$; otherwise, the function

(4.6)
$$h(x) = E(x,a) - E(x,b)$$

is continuous, and h(a) < 0, h(b) > 0 by hypothesis. Thus, an $x^* \in I$ exists
such that $h(x^*) = 0$, and thus $x^*$ satisfies equation (4.5). Now, suppose that
$y < x^*$. It follows that

(4.7)
$$E(y,b) > E(x^*,b) = E(x^*,a)$$

- 4 -

as E is strictly monotone increasing in $|x - b|$, and if $y > x^*$, then similarly

(4.8)  $\qquad$ $E(y,a) > E(x^*,a) = E(x^*,b)$

from which the uniqueness and optimality of $x^*$ follow.  QED.

Definition 4.4.  The value

(4.9)  $\qquad$ $E^* = E(x^*,b) = E(x^*,a)$

obtained for optimal $x^* \epsilon I$ is called the <u>optimal error estimate</u> on I for the measure of error E.

One object of this paper is to introduce $x^*, E^*$ as alternative coordinates for the representation of the interval I.  As the inverse functions

(4.10)  $\qquad$ $\ell(x,E) = \{y \mid E(x,y) = E, \; y \leq x\},$

and

(4.11)  $\qquad$ $u(x,E) = \{y \mid E(x,y) = E, \; y \geq x\},$

are single-valued and continuous by the hypotheses on the error function E, one has the equivalence of the standard representation of the interval I and the $x^*, E^*$ coordinates through the equations

(4.12)  $\qquad$ $a = \ell(x^*,E^*), \quad b = u(x^*,E^*),$

and the definitions of $x^*, E^*$ in terms of a,b by equations (4.5) and (4.9).

5.  <u>The standard format and interval arithmetic</u>.  Interval extensions of the arithmetic operations form the basis of what is called <u>interval arithmetic</u> [4], [5].  In the standard format (3.3) for representation of the intervals $I = [a,b]$ and $J = [c,d]$, the fundamental rules for the operations of interval arithmetic are:

(i)  <u>Addition</u>

(5.1)  $\qquad$ $I + J = [a,b] + [c,d] = [a + b , c + d];$

(ii)  <u>Subtraction</u>

(5.2)  $\qquad$ $I - J = [a,b] - [c,d] = [a - d , b - c];$

(iii)  <u>Multiplication</u>

(5.3)  $\qquad$ $I \cdot J = [a,b] \cdot [c,d] = [\min\{ac,ad,bc,bd\} , \max\{ac,ad,bc,bd\}];$

(iv)  <u>Reciprocation</u> is defined only for intervals J such that $0 \notin J$,

(5.4)  $\qquad$ $J^{-1} = [c , d]^{-1} = \left[\frac{1}{d} , \frac{1}{c}\right]$ if $cd > 0$.

Remark 5.1.  Arithmetic operations between real numbers k and intervals are defined by identification of real numbers with <u>degenerate</u> intervals; one writes

(5.5) $k = [k, k]$

and uses the rules for interval arithmetic; for example, if $k \geq 0$, then one has $k \cdot [a, b] = [ka, kb]$.

Remark 5.2. <u>Division</u> is defined as a compound operation; if $0 \notin J$, then

(5.6) $I/J = I \cdot J^{-1}$ if $c \cdot d = \ell(J) \cdot u(J) > 0$.

Definition 5.1. An interval I is <u>positive</u> if $\ell(I) > 0$, negative if $u(I) < 0$, and <u>zero</u> if equal to $0 = [0, 0]$; I is <u>positive</u> (<u>negative</u>) <u>semidefinite</u> if nonzero, $\ell(I) \cdot u(I) = 0$ and $u(I) > 0$ ($\ell(I) < 0$); finally, I is <u>indefinite</u> if $\ell(I) \cdot u(I) < 0$.

Remark 5.3. By analysis of the signs of the intervals I,J, it is possible to calculate the product $I \cdot J$ in certain cases by shortcut methods. However, even with a programmable pocket calculator, it is easy to compute all four products and sort out the smallest and largest. It is also possible to derive explicit (but clumsy) formulas for the product of two intervals, based on the identities

(5.7) $\min\{x,y\} = \dfrac{x + y - |x - y|}{2}$, $\max\{x,y\} = \dfrac{x + y + |x - y|}{2}$.

The derivation of such formulas for the product of intervals is left as an exercise for the reader.

Before leaving the familiar standard format for an interval I, it will be noted that this representation provides upper and lower bounds for the points of I, which is the information required in many practical situations. A biologist, for example, may be interested in the range of water temperatures in which a certain species of fish can survive. Another illustration is the "weight and balance" calculation performed by a pilot before flying to determine if the moment of his airplane about its center of gravity along its longitudinal axis is within safe limits for its gross weight. Many other examples of the direct use of intervals may be found, and of problems to which, "The interval is the answer."

6. <u>The midpoint-halfwidth format and absolute error</u>. The simplest measure of the error of approximation of y by x is the <u>absolute error</u>

(6.1) $E(x,y) = |x - y|$,

which is simply the distance between x and y. For this error function, equation (4.5) becomes

(6.2) $x - a = b - x$,

which is easily solved for the optimal point and error estimate

(6.3) $x^* = (a + b)/2 = m[a, b]$,

(6.4) $\qquad E^* = (b - a)/2 = \alpha[a, b] = \frac{1}{2} w[a, b],$

respectively, where the <u>width</u> $w(I) = w[a, b]$ of the interval I is of course

(6.5) $\qquad w(I) = w[a, b] = b - a.$

The representation of the interval I in <u>midpoint-halfwidth</u> format (or <u>A-format</u>) is then

(6.6) $\qquad I = [m(I) \pm \alpha(I)],$

where the comma in the standard format has been replaced by the $\pm$ sign. (One can also use the notation (6.6) without the square brackets, if no confusion is entailed.) As has been shown above (and has been known for a long time), this representation is optimal with respect to absolute error. The transformation from (6.6) back to the standard format is by means of the evident relationships

(6.7) $\qquad a = \ell(I) = m(I) - \alpha(I), \quad b = u(I) = m(I) + \alpha(I).$

The basic rules for interval arithmetic in this format are:

   (i)   <u>Addition</u>

(6.8) $\qquad m(I + J) = m(I) + m(J), \quad \alpha(I + J) = \alpha(I) + \alpha(J);$

   (ii)  <u>Subtraction</u>

(6.9) $\qquad m(I - J) = m(I) - m(J), \quad \alpha(I - J) = \alpha(I) + \alpha(J);$

   (iii) <u>Multiplication</u>

$$m(I \cdot J) = m(I) \cdot m(J) + \frac{1}{2}\{ |\alpha(I) \cdot m(J) + A(I) \cdot \alpha(J)| -$$
$$- |\alpha(I) \cdot m(J) - A(I) \cdot \alpha(J)| \},$$

(6.10)

$$\alpha(I \cdot J) = \mu(I) \cdot \alpha(J) + \frac{1}{2}\{ |\alpha(I) \cdot m(J) + A(I) \cdot \alpha(J)| +$$
$$+ |\alpha(I) \cdot m(J) - A(I) \cdot \alpha(J)| \},$$

where

$$A(I) = \frac{1}{2}( |m(I) + \alpha(I)| - |m(I) - \alpha(I)| ),$$

(6.11)

$$\mu(I) = \frac{1}{2}( |m(I) + \alpha(I)| + |m(I) - \alpha(I)| ).$$

   (iv) <u>Reciprocation</u> is defined only for intervals J such that $0 \notin J$,

$$m(J^{-1}) = m(J)/(m(J)^2 - \alpha(J)^2),$$

(6.12)

$$\alpha(J^{-1}) = \alpha(J)/(m(J)^2 - \alpha(J)^2), \text{ if } m(J)^2 > \alpha(J)^2.$$

Remark 6.1. There are several alternative expressions to (6.10)-(6.11)

for the product of two intervals in midpoint-halfwidth format, which the reader
may wish to derive.

Remark 6.2. The condition $m(J)^2 > \alpha(J)^2$ in (6.12) enforces $|m(J)| > \alpha(J)$,
and thus the interval J will be positive or negative.

Remark 6.3. The midpoint $m[a , b]$ is also called the <u>average</u>, or <u>arithmetic
mean</u> of the numbers a,b.

Finally, the set IR of intervals may be represented in the $m,\alpha$ coordinate
system simply as the closed upper halfplane of the $m,\alpha$-plane (see Figure 2).
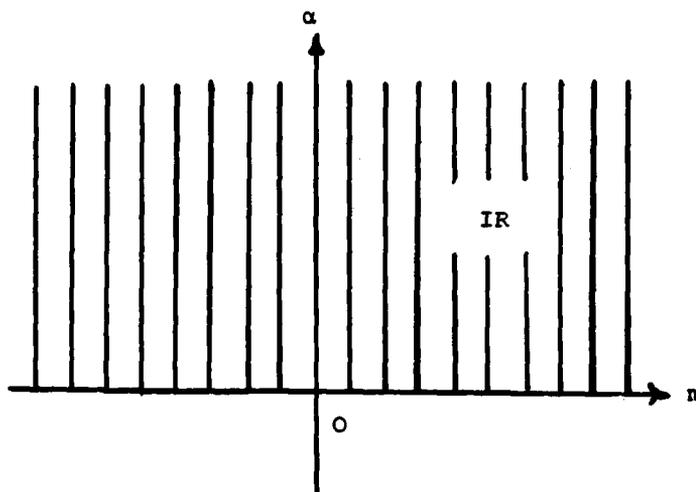


Figure 2. IR as a subset of the $m,\alpha$-plane.

## 7. <u>The harmonic point-relative width format.  Relative and percentage error</u>.

A measure of error which is useful in many circumstances is the <u>relative</u> <u>error</u>
of approximation of a real number $y \neq 0$ by a real number x of the same sign,
which is defined by

$$(7.1) \qquad E(x,y) = \left| \frac{x - y}{y} \right| .$$

For simplicity, it will be assumed that the numbers being considered are posi-
tive; the same results hold, with appropriate changes of sign, for intervals
of negative numbers.  For relative error, equation (4.5) becomes

$$(7.2) \qquad \frac{x - a}{a} = \frac{b - x}{b} ,$$

which has as its solution

$$(7.3) \qquad x^* = \frac{2ab}{a + b} = h[a , b] = h(I),$$

the <u>harmonic</u> <u>point</u> of I (also called the <u>harmonic</u> <u>mean</u> of a,b), with the cor-
responding optimal error estimate

$$(7.4) \qquad E^* = \frac{b - a}{a + b} = r[a , b] = r(I),$$

- 8 -

the _relative_ _width_ of I. This gives the corresponding _harmonic_ _point-relative_
_width_ format (or _R-format_) for the representation of the interval I:

(7.5) $\qquad I = [h(I) \ R \ r(I)].$

The transformation from this format back to the standard format is given by

(7.6) $\qquad a = \ell(I) = h(I)/(1 + r(I)), \quad b = u(I) = h(I)/(1 - r(I)).$

For positive intervals I,J, the basic rules of interval arithmetic in the
R-format are:

    (i)   _Addition_

$$h(I + J) = \frac{(h(I) + h(J))^2 - (h(I)r(J) + h(J)r(I))^2}{h(I) + h(J) - (h(I)r(J)^2 + h(J)r(I)^2)} \ ,$$

(7.7)

$$r(I + J) = \frac{h(I)r(I) + h(J)r(J) - (h(I)r(I)r(J)^2 + h(J)r(J)r(I)^2)}{h(I) + h(J) - (h(I)r(J)^2 + h(J)r(I)^2)} \ ;$$

    (ii)   _Subtraction_ is defined only if I - J is positive,

$$h(I - J) = \frac{(h(I) - h(J))^2 - (h(I)r(J) + h(J)r(I))^2}{h(I) - h(J) - (h(I)r(J)^2 + h(J)r(I)^2)} \ ,$$

(7.8)

$$r(I - J) = \frac{h(I)r(I) + h(J)r(J) - (h(I)r(I)r(J)^2 + h(J)r(J)r(I)^2}{h(I) - h(J) - (h(I)r(J)^2 + h(J)r(I)^2)} \ ,$$

    if $h(I) - h(J) > h(I)r(J) + h(J)r(I)$;

    (iii)   _Multiplication_

(7.9) $\qquad h(I \cdot J) = \dfrac{h(I)h(J)}{1 + r(I)r(J)} \ , \quad r(I \cdot J) = \dfrac{r(I) + r(J)}{1 + r(I)r(J)} \ ;$

    (iv)   _Reciprocation_

(7.10) $\qquad h(J^{-1}) = \dfrac{1 - r(J)^2}{h(J)} \ , \quad r(J^{-1}) = r(J).$

Remark 7.1. The condition on subtraction in (7.8) insures that d < a in
the standard format I = [a , b], J = [c , d].

Remark 7.2. The formulas (7.7)-(7.8) for addition and subtraction in this
format are unwieldy; a transformation to and from the standard format is prob-
ably preferable for these operations.

Example 7.1. The interval expression (3.6) may be written

(7.11) $\qquad \displaystyle\int_{0}^{2}\sqrt{1 + 4x} \ dx \ \epsilon \ [4.3379497 \ R \ 0.001689438]$

in the R-format. This means that one may take x* = h(I) = 4.3379497 as the ap-
proximate value of the integral (3.1), with relative error bounded by E* = r(I)
= 0.001689438, which can be symbolized by

(7.12) $\qquad \int_0^2 \sqrt{1 + 4x}\ dx \simeq 4.3379497\ R\ 0.001689438.$

Remark 7.3. In conversion from [a , b] to [h(I) R r(I)] format, r(I) is rounded upward, if necessary, so that [a , b] ⊂ [h(I) R r(I)]. This <u>directed rounding</u> has been employed in (7.11) and the numerical examples cited below.

Remark 7.4. The <u>percentage error</u> is 100 times the relative error (7.1), so that x* = h(I) is also optimal for percentage error, with

(7.13) $\qquad$ E* = p(I) = 100·r(I)

being the corresponding optimal error bound.

Thus, one may write (7.12) as

(7.14) $\qquad \int_0^2 \sqrt{1 + 4x}\ dx \simeq 4.3379497 \pm 0.1689438\%,$

showing that the percentage error in taking h(I) = 4.3379497 as an approximation to y in this example is less than 0.17%. The <u>percentage format</u> (or <u>P-format</u>) for representation of positive intervals is thus

(7.15) $\qquad$ I = h(I) ± p(I)% = [h(I) % p(I)].

8. <u>The geometric point-ratio format. Relative precision and approximate relative precision</u>. The <u>geometric point</u> g(I) = g[a , b] = $\sqrt{ab}$ and the <u>ratio</u> ρ(I) = ρ[a , b] = $\sqrt{b/a}$ may also be used to represent a positive interval I = [a , b]. This coordinate system will be called the <u>geometric point-ratio format</u> (or <u>G-format</u>), and symbolized by

(8.1) $\qquad$ I = [g(I) $*$ ρ(I)].

The sign $*$ is a combination of the × and ÷ signs, symbolizing the transformation

(8.2) $\qquad$ a = ℓ(I) = g(I)/ρ(I), b = u(I) = g(I)·ρ(I)

from the format (8.1) into the standard format, in the same way that ± is used in the midpoint-halfwidth format. The geometric point g(I) = g[a,b] = $\sqrt{ab}$ is also called the <u>geometric mean</u> of the positive numbers a,b.

The G-format is related to error estimation through the <u>maximum ratio function</u> M defined by

(8.3) $\qquad$ $M(x,y) = \max \left\{ \dfrac{x}{y},\ \dfrac{y}{x} \right\} = \dfrac{x^2 + y^2 + |x^2 - y^2|}{2xy}$

for x > 0, y > 0. On the positive interval I = [a , b], the minimum of the maximum value of M is attained for x* satisfying

(8.4) $\qquad$ $\dfrac{x}{a} = \dfrac{b}{x}$ ,

that is, $x^* = g(I)$. This point will be optimal for each error function

$$(8.5) \qquad E(x,y) = f(M(x,y)),$$

where f is continuous, strictly monotone increasing, and $f(1) = 0$. As

$$(8.6) \qquad M(x^*,a) = M(x^*,b) = \rho(I),$$

one has

$$(8.7) \qquad E^* = f(\rho(I))$$

as the optimal error estimate for the error function (8.5). In particular, the <u>relative</u> <u>precision</u>

$$(8.8) \qquad r.p.(I) = \ln \rho(I)$$

and <u>approximate</u> <u>relative</u> <u>precision</u>

$$(8.9) \qquad a.r.p.(I) = \rho(I) - 1$$

defined by Olver [6] fit into this category, and lead to the respective formats

$$(8.10) \qquad I = [g(I) \; r.p. \; \ln \rho(I)] \text{ and } I = [g(I) \; a.r.p. \; \rho(I) - 1]$$

for the representation of intervals.

Example 8.1. In the formats (8.1) and (8.10), the interval expression (3.6) may be written

$$(8.11) \qquad \int_0^2 \sqrt{1 + 4x} \; dx \; \epsilon \; [4.3379559 \; * \; 1.0016909],$$

or as

$$(8.12) \qquad \int_0^2 \sqrt{1 + 4x} \; dx \simeq 4.3379559 \; r.p. \; 0.001689438,$$

or

$$(8.13) \qquad \int_0^2 \sqrt{1 + 4x} \; dx \simeq 4.3379559 \; a.r.p. \; 0.0016909.$$

Directed rounding has been used in the conversion of (3.6) into the above formats.

As conversion between the r.p. and a.r.p. formats (8.10) and the G-format (8.1) is immediate, the rules of interval arithmetic will be stated for the G-format.

(i) <u>Addition</u>

$$g(I + J)^2 = (g(I) + g(J))^2 - \frac{g(I)g(J)}{\rho(I)\rho(J)}(\rho(I) - \rho(J))^2,$$

$$(8.14)$$

$$\rho(I + J)^2 = \rho(I)\rho(J)\frac{g(I)\rho(I) + g(J)\rho(J)}{g(I)\rho(J) + g(J)\rho(I)} ;$$

(ii)  <u>Subtraction</u> is defined only if I − J is positive,

$$g(I - J)^2 = (g(I) - g(J))^2 - \frac{g(I)g(J)}{\rho(I)\rho(J)}(\rho(I)\rho(J) - 1)^2,$$

(8.15)

$$\rho(I - J)^2 = \frac{\rho(I)}{\rho(J)}(\frac{g(I)\rho(I)\rho(J) - g(J)}{g(I) - g(J)\rho(I)\rho(J)}) \text{ if } g(I) - g(J)\rho(I)\rho(J) > 0;$$
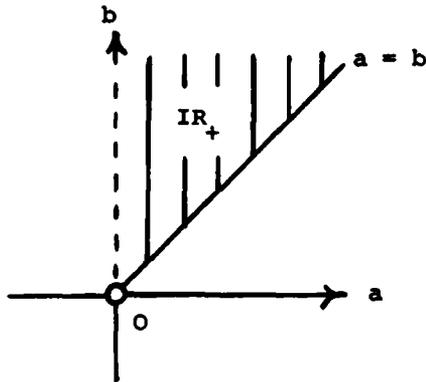
(iii)  <u>Multiplication</u>

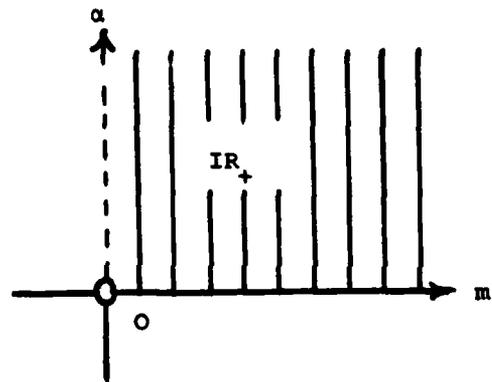(8.16)      $g(I \cdot J) = g(I) \cdot g(J), \quad \rho(I \cdot J) = \rho(I) \cdot \rho(J);$

(iv)  <u>Reciprocation</u>

(8.17)      $g(J^{-1}) = 1/g(J) = g(J)^{-1}, \quad \rho(J^{-1}) = \rho(J).$
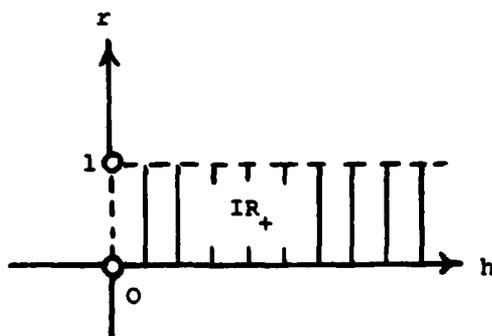
The set $IR_+$ of positive intervals is shown in Figure 3 in the [a , b], [m ± α], [h R r], and [g * ρ] coordinate systems.
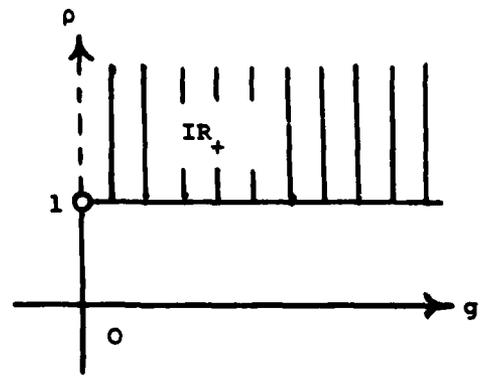


(a)  The a,b-plane.          (b)  The m,α-plane.

(c)  The h,r-plane.          (d)  The g,ρ-plane.

Figure 3.  The set $IR_+$ of positive intervals.

Remark 8.1.  The transformations from the G-format into the relative precision (r.p.) and approximate relative precision (a.r.p.) formats may be made

by using the formulas

$$(8.18) \qquad \rho(I) = e^{r.p.(I)},$$

and

$$(8.19) \qquad \rho(I) = 1 + a.r.p.(I),$$

which follow directly from (8.8) and (8.9), respectively. From (8.14)-(8.17), the rules of interval arithmetic in the r.p. format are obtained from the following formulas for the relative precision, the geometric point being given in the same way.

(i) <u>Addition</u>

$$(8.20) \qquad r.p.(I + J) = \frac{r.p.(I) + r.p.(J)}{2} + \frac{1}{2}\ln\{\frac{g(I)\rho(I) + g(J)\rho(J)}{g(I)\rho(J) + g(J)\rho(I)}\};$$

(ii) <u>Subtraction</u> is defined only if $I - J$ is positive,

$$r.p.(I - J) = \frac{r.p.(I) - r.p.(J)}{2} + \frac{1}{2}\ln\{\frac{g(I)\rho(I)\rho(J) - g(J)}{g(I) - g(J)\rho(I)\rho(J)}\},$$

(8.21)

$$\text{if } g(I) - g(J)\rho(I)\rho(J) > 0;$$

(iii) <u>Multiplication</u>

$$(8.22) \qquad r.p.(I \cdot J) = r.p.(I) + r.p.(J);$$

(iv) <u>Reciprocation</u>

$$(8.23) \qquad r.p.(J^{-1}) = r.p.(J).$$

Remark 8.2. The [g * ρ] format is related to the [m ± α] format for positive intervals by means of logarithms. If I is positive, then the <u>logarithm</u> of I is defined to be the interval

$$(8.24) \qquad \ln(I) = [\ln\{\ell(I)\}, \ln\{u(I)\}].$$

Thus,

$$(8.25) \qquad \ln\{g(I)\} = \frac{1}{2}(\ln\{\ell(I)\} + \ln\{u(I)\}) = m(\ln(I)),$$

and

$$(8.26) \qquad \ln\{\rho(I)\} = r.p.(I) = \frac{1}{2}(\ln\{u(I)\} - \ln\{\ell(I)\}) = \alpha(\ln(I)),$$

as noted by Olver [6].

9. <u>Choice of format for intervals</u>. In various applications, one of the formats presented above may be preferable to the others. For error estimation in positive intervals, the classical inequality between the arithmetic, geometric, and harmonic means of two positive numbers a,b is

$$(9.1) \qquad h[a, b] \le g[a, b] \le m[a, b],$$

with equality if and only if $a = b$ [2]. Thus, the choice of $x^* = g(I)$ to

- 13 -

optimize relative precision may be viewed as a compromise between minimization of absolute and relative errors over I. In addition, the choice of the geometric point has important theoretical and computational implications, as discussed by Olver [6].

For expression of interval arithmetic, the standard and the midpoint-halfwidth representations are simpler for addition and subtraction, while the harmonic point-relative width and geometric point-ratio formats are simpler for multiplication and reciprocation (and hence division). It should be pointed out that for _positive_ intervals,

$$(9.2) \qquad I \cdot J = [ac, bd]$$

in standard format, and for the A-format, it follows from (6.11) that $A(I) = \alpha(I)$, $\mu(I) = m(I)$, so that (6.10) simplifies to

$$m(I \cdot J) = m(I) \cdot m(J) + \alpha(I) \cdot \alpha(J),$$

$$(9.3)$$

$$\alpha(I \cdot J) = m(I) \cdot \alpha(J) + \alpha(I) \cdot m(J).$$

These formulas compare favorably in simplicity to the corresponding rules in the R- and G-formats.

Rather than insisting on one format or another throughout a computation, it may be expedient to transform from one to another during the calculation, or to represent input data or output results. Relations such as

$$(9.4) \qquad m(J^{-1}) = 1/h(J) = h(J)^{-1}, \quad \alpha(J^{-1}) = \alpha(J)/g(J)^2,$$

which are among the many charming relationships between the arithmetic, geometric, and harmonic means, indicate the possiblities of coordinate transformations in interval calculations.

10. A priori _error estimates and excess width_. At any stage in an interval computation, one has an output interval $J = F(I)$ for which the representative point $x^*(J)$ and the optimal error estimate (or bound) $E^*(J)$ for a given measure of error E may be found by expressing J in $x^*, E^*$ coordinates. This is called a _posteriori_ error estimation, as the calculations called for by the transformation F have been performed. Another type of error analysis is concerned with the estimation of $x^*(J)$, $E^*(J)$ (usually $E^*(J)$ in particular), given $x^*(I)$, $E^*(I)$ for the input interval I, before the transformation F is performed to obtain J. This is called a _priori_ error estimation. From the standpoint of interval analysis, a _priori_ error analysis gives an estimated interval K such that $J \subset K$, and hence

$$(10.1) \qquad E^*(K) \geq E^*(J),$$

with the inequality being strict if the _excess width_ of K over its subinterval

is positive, where excess width is defined to be

(10.2)        $d(K,J) = \max\{\ell(J) - \ell(K), u(K) - u(J)\}$   for $J \subset K$,

which is simply the distance between K and J in the ordinary metric topology for intervals [4], [5].

A priori error estimation may be viewed as a type of approximate interval arithmetic, in which the results of various operations are expressed by simple formulas which yield intervals containing the results which would be actually computed.  The highly interesting methods of Olver [6] may be interpreted in this light.  For example ([6], (2.2)), the expression

(10.3)          $a \simeq \tilde{a}; \; ap(\alpha)$,

which means that $\tilde{a}$ is an approximation to a of absolute precision $\alpha$ in the terminology of [6], may be expressed in the notation of the present paper as

(10.4)          $a \in [\tilde{a} \pm \alpha] = I$,

i.e., $\tilde{a} = m(I)$, $\alpha = \alpha(I)$.  Olver gives ([6], (2.4)) the formula

(10.5)          $ab \simeq \widetilde{ab}; \; ap(|\tilde{a}|\beta + |\tilde{b}|\alpha + \alpha\beta)$,

which in interval notation defines the interval K with

                 $m(K) = m(I)m(J)$,

(10.7)

                 $\alpha(K) = |m(I)|\alpha(J) + |m(J)|\alpha(I) + \alpha(I)\alpha(J)$,

where $J = [\tilde{b} \pm \beta]$ contains the point b.  Comparison of (10.7) with (9.3) shows that (10.5) represents an approximation K to I·J which, for I,J positive, has excess width

(10.8)          $d(K, I\cdot J) = 2\alpha(I)\alpha(J)$,

which is positive if I,J are each nondegenerate.  It may be noted that $u(K) = u(I\cdot J)$, i.e., formula (10.7) gives the correct upper endpoint of I·J for positive I,J.

Similarly, Olver's definition ([6], (3.2)) of the relative precision of approximation of a by $\tilde{a}$,

(10.9)          $a \simeq \tilde{a}; \; rp(\alpha)$,

means

(10.10)         $a \in [\tilde{a} * e^{rp(\alpha)}] = I$

in the notation of the present paper, that is, $\tilde{a} = g(I)$, $rp(\alpha) = \ln \rho(I)$.  For relative precision of addition ([6], (3.5)), Olver gives the formula

$$(10.11) \qquad a + b \simeq \tilde{a} + \tilde{b}; \quad rp(\ln\{\frac{\tilde{a}e^{\alpha} + \tilde{b}e^{\beta}}{\tilde{a} + \tilde{b}}\}).$$

In the notation of §8, this defines an interval K with

$$(10.12) \qquad g(K) = g(I) + g(J), \quad \rho(K) = \frac{g(I)\rho(I) + g(J)\rho(J)}{g(I) + g(J)}.$$

In standard format, the interval I + J is

$$(10.13) \qquad I + J = [\frac{g(I)}{\rho(I)} + \frac{g(J)}{\rho(J)}, g(I)\rho(I) + g(J)\rho(J)],$$

which shows that $u(K) \simeq u(I + J)$, and K has excess width

$$(10.14) \qquad d(K, I + J) = \frac{g(I)g(J)(\rho(I) - \rho(J))^2}{\rho(I)\rho(J)(g(I)\rho(I) + g(J)\rho(J))},$$

which will be positive if $\rho(I) \neq \rho(J)$.

The error analysis given by Olver in [6] also provides formulas to predict the results obtained with various types of rounding. As adapted to directed rounding used in interval arithmetic, excess width in the predicted intervals will also be observed as above. This does not detract from the usefulness of the a priori methods given in [6]; it simply means that one can expect the results obtained by the use of properly rounded interval arithmetic to be better than predicted.

11. Historical remarks. The optimality of the midpoint of an interval with respect to absolute error has been known for a long time. The same relationship between the harmonic point and the relative (or percentage) error was learned from George Pólya in a seminar he conducted at the University of California, Berkely, in 1956. On the basis of these results, the A-, R-, and P-formats were made available as optional alternatives to the standard format in the programs [1], [3], for interval numerical integration and the solution of nonlinear systems of equations, respectively. The results given above on the significance of the geometric point for relative precision had to wait, of course, for the fundamental paper by Olver [6].

12. Acknowlegments. This paper was written during a leave at the Institute for Applied Mathematics, University of Freiburg, Germany. The gracious host of the Author was Prof. Dr. Karl Nickel, founder of the Interval Library, a unique and effective research tool.

# References

1. Julia H. Gray and L. B. Rall. INTE: A UNIVAC 1108/1110 program for numerical integration with rigorous error estimation, Mathematics Research Center, MRC Tech. Rep. 1428, University of Wisconsin-Madison, October, 1975.

2. G. H. Hardy, J. E. Littlewood, and G. Pólya. Inequalities, 2nd Ed., Cambridge University Press, Cambridge, England, 1952.

3. Dennis Kuba and L. B. Rall. A UNIVAC 1108 program for obtaining rigorous error estimates for approximate solutions of systems of equations, Mathematics Research Center, MRC Tech. Rep. 1168, University of Wisconsin-Madison, January, 1972.

4. R. E. Moore. Interval Analysis, Prentice-Hall, Englewood Cliffs, N. J., 1966.

5. R. E. Moore. Methods and Applications of Interval Analysis, SIAM Publications, Philadelphia, Pa., 1979.

6. F. W. J. Olver. A new approach to error arithmetic, SIAM J. Numer. Anal. 15 (1978), 368-393.

7. L. B. Rall. Applications of software for automatic differentiation in numerical computation, Computing, Suppl. 2 (1980), 141-156.

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS<br>BEFORE COMPLETING FORM |
|---|---|---|
| 1. REPORT NUMBER<br>2098 | 2. GOVT ACCESSION NO.<br>AD-A089 636 | 3. RECIPIENT'S CATALOG NUMBER |
| 4. TITLE (and Subtitle)<br><br>REPRESENTATIONS OF INTERVALS AND OPTIMAL<br>ERROR BOUNDS | | 5. TYPE OF REPORT & PERIOD COVERED<br>Summary Report - no specific<br>reporting period |
| | | 6. PERFORMING ORG. REPORT NUMBER |
| 7. AUTHOR(s)<br><br>L. B. Rall | | 8. CONTRACT OR GRANT NUMBER(s)<br><br>DAAG29-80-C-0041<br>511-15849 |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS<br>Mathematics Research Center, University of<br>610 Walnut Street                         Wisconsin<br>Madison, Wisconsin 53706 | | 10. PROGRAM ELEMENT, PROJECT, TASK<br>AREA & WORK UNIT NUMBERS<br>Work Unit Number 3 -<br>Numerical Analysis<br>and Computer Science |
| 11. CONTROLLING OFFICE NAME AND ADDRESS<br>U. S. Army Research Office<br>P.O. Box 12211<br>Research Triangle Park, North Carolina 27709 | | 12. REPORT DATE<br>July 1980 |
| | | 13. NUMBER OF PAGES<br>17 |
| 14. MONITORING AGENCY NAME & ADDRESS(if different from Controlling Office) | | 15. SECURITY CLASS. (of this report)<br><br>UNCLASSIFIED |
| | | 15a. DECLASSIFICATION/DOWNGRADING<br>SCHEDULE |

16. DISTRIBUTION STATEMENT (of this Report)

Approved for public release; distribution unlimited.

17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)

18. SUPPLEMENTARY NOTES

19. KEY WORDS (Continue on reverse side if necessary and identify by block number)

Optimal error bounds, Interval Arithmetic, Absolute error, Relative error,
Percentage error, Arithmetic, geometric and harmonic means, Excess width

20. ABSTRACT (Continue on reverse side if necessary and identify by block number)

Calculations with interval arithmetic and interval extensions of real
functions are often used to obtain lower and upper bounds for what would be
the theoretical results of precise computations on exact data. A more tradi-
tional way to represent approximate data and results contained in intervals
is by means of a representative point  x  in the interval and an associated
measure of error  E.  For example, the midpoint of an interval has absolute
error as an approximation to other points of the interval which is bounded
(continued)

DD $_{1 \text{ JAN } 73}^{\text{FORM}}$ 1473     EDITION OF 1 NOV 65 IS OBSOLETE

UNCLASSIFIED

ABSTRACT (cont.)

by half the width of the interval. In general, the chosen point is said to be
optimal with respect to the measure of error if the maximum value of  E  over the
interval is minimized. A general method for determination of optimal points  x*
and minimal error bounds  E*  is given. In particular, the harmonic mean of the
endpoints of positive intervals is shown to be optimal with respect to relative
(or percentage) error, and the geometric mean plays the same role with respect to
relative precision, which is a measure of error introduced recently by F. W. J.
Olver [SIAM J. Numer. Anal. 15 (1978), 368-393]. In addition, the optimal point
and error bound  x*, E*  may be regarded as alternative coordinates for represen-
tation of an interval. Explicit rules for interval arithmetic are given in terms
of the arithmetic, geometric, and harmonic means of the endpoints and the asso-
ciated optimal error bounds. Comparisons are also made between the results of
exact or rounded interval arithmetic and the a priori estimates of relative pre-
cision presented in the cited paper by Olver, which show that the intervals
resulting from calculation with interval arithmetic can be expected to be smaller
than predicted.