

## あふれのない浮動小数点表示方式†

松井正一†† 伊理正夫††

従来の浮動小数点表示方式で問題となっていた指数部あふれを解決する新しい表現方式の提案を行う。新表現方式は、指数部と仮数部との境界を動的に変化させることにより、指数部あふれを防ぐとともに、普通の大きさの数に対してはより高い精度を確保することができる。また普通の数ではない数をいくつか考えることにより、“数の体系”を閉じたものとする。新しい表現方式が効果を発揮するような計算例も示す。

## 1. はじめに

数値計算に計算機を使う際には、実数（の近似値）を1語中に“浮動小数点”の形で表わすのが慣習になっている。初期の計算機では“固定小数点”が用いられたが、“桁取り（スケールリング）”の繁雑さのため、有効桁数の少々（指数部に割当てる分）の犠牲を払っても任意の大きさの数が同じ相対精度で取扱える浮動小数点方式に完全に軍配があがった。浮動小数点方式にも多くの変種があるが、最近では、IBM型の16進仮数部絶対値、指数部7ビット“げたばき”式の型が優勢である。経験が積まれるに従って、このIBM型に対する多くの不満が数値計算家の間に亢じている。最も切実な問題の一つに“あふれ（over-flow, under-flow）”がある。一松の解説<sup>6)</sup>にもこのような問題の指摘がある。一般に、Hamming<sup>4)</sup>も述べているように、計算が進むにつれて扱われる数の指数部の絶対値の範囲は広がるので、扱える数の範囲がたとえば $(16^{-2^k}, 16^{2^k}) \equiv (10^{-77}, 10^{77})$ のように、あらかじめ限られていると、それからはみ出す数が生じやすい。それを防ぐために計算途上で数の大きさを常に検査するのは「スケールリングの心配無用」という浮動小数点方式の最大の特長が失われてしまう。一松<sup>6)</sup>に紹介されているKahanの提案は一言でいえば、(i)1語中で指数部に割当てる部分を増大させ、(ii)“非数”を導入しあふれた後の面倒をもみることにある。

その他、記数法の底(=2, 4, 8, 16等)の良し悪しについての議論<sup>1), 2)</sup>(およびその参考文献)、通常のものとは異なる表示体系の提案<sup>3), 12)</sup>等々、浮動小数点方式の本質的な問題への関心が最近高まりつつあるよ

うにみえる。

可変語長にすれば問題はすべて解決するという議論もあるが<sup>7)</sup>、非常に特殊な目的を除けば、通常の数値計算では一つの数値の担うべき“情報量”には当然限度があり、計算時間、記憶場所の経済性の点からも、これは考察の外としてよいであろう。

本論文では、「固定語長内で実数を近似的に表現する」という大前提のもとで、一つの理想的な方式を提案し、その特徴を調べ、それを具体的に実現し、それによる数値計算例を示す。

本論文で提案する方式の基本的な考え方を、形式ばらずに述べると、次の2点となる。(i)与えられた1語中に数を表現する際、より重要な情報から順に必要なビット数を割当てる(指数部に必要なだけのビット数をまず確保した後に、残りを仮数部に用いる；指数部だけでも入り切らなくなったら、指数部を更に浮動小数点表示にして、“指数部の指数部”に優先的な割当てをする；以下同様なことを繰り返す)。(ii)数値演算の結果が常にプログラムで取扱い可能な形である(あふれによる割込みが生じたりしない)ために導入すべき必要にして十分な“非数”をも含んだ“数の体系”を確定させる。(この際四則演算に関する対称性を重視する。)

(i)についてはMorris<sup>10)</sup>に類似の着想が認められ、またKnuthの本<sup>9)</sup>の演習問題(第4章2.1節の17番)にそのような方向への示唆があるが、本論文ではその可能性を徹底的に追求してみる。(ii)は(i)とほとんど独立な考え方で、Kahan等が部分的に論じているが<sup>8)</sup>いくつかの妥当と考えられる前提のもとで必然的に一つの体系が確定されることを指摘する。

## 2. 基本的な考え方

## 2.1 指数部可変長方式とその一般化

実数 $x$ の浮動小数点表現は、底を $\beta$ (=2, 4, 10, 16

† An Over Flow/Under Flow-free Floating-Point Representation of Numbers by SHOUICHI MATSUI and MASAO IRI  
(Department of Mathematical Engineering and Instrumentation Physics, Faculty of Engineering, University of Tokyo).

†† 東京大学工学部計数工学科

等) とするとき,

<レベル 0>

$$x = (-1)^{s_0} F \times \beta^{-1} E_1 \quad (2.1)$$

の形をとる.  $s_0, s_1$  は 0 あるいは 1,  $E_1$  は非負整数で,  $F$  は正規化条件

$$1 > F \geq 1/\beta \quad (2.2)$$

あるいは

$$\beta > F \geq 1 \quad (2.3)$$

を満たす実数である。(通常は前者が多く用いられるが, ここでは 2.3 で述べる要求をも考えて後者を採用する.)  $x \neq 0$  に対して  $s_0, E_1, F$  は ( $E_1 \neq 0$  なら  $s_1$  も) 一意に定まる. 固定語長の 1 語中に  $E_1$  と  $F$  を表現する際いくらかの近似誤差が不可避であるとすればその誤差は  $F$  のみに負わせるのが明らかによい.(正規化条件の範囲で  $F$  が変化しても  $E_1$  が  $\pm 1$  変化するよりも,  $x$  に及ぼす影響は小さいから.)

$|\log_{\beta}|x||$  が小さいうちは  $E_1$  の桁数は少ないから,  $E_1$  は少ないビット数で正確に表わせ, 1 語中の大きな部分を  $F$  に割当てることができる.  $|\log_{\beta}|x||$  が大きくなるにつれて,  $E_1$  の桁数が増すので, それに割当てる部分が増加し,  $F$  に割当てる部分は次第に減少する.

$E_1$  だけで 1 語を占有するくらいに  $|\log_{\beta}|x||$  が大きくなった後は, (2.1) の代わりに,

<レベル 1>

$$x = (-1)^{s_0} F \times \beta^{-1} E_1 \times \beta E_2 \quad (2.4)$$

という表示を採用する. このときには,  $F$  は“正規化条件を満たすある数”という情報しか有しない(これを表わすための場所は不用).  $E_1, E_2$  は (2.1) の表示をしたときの  $E_1$  (これを  $E$  と書く) が

$$E_1 \times \beta E_2 \leq E < (E_1 + 1) \times \beta E_2 \quad (2.5)$$

を満足し, かつ  $E_1, E_2$  が 1 語中に収まる範囲内で  $E_1$  の桁数なるべく大きくなるように選ぶ.

同様にして, 更に高次の“レベル”の表示に進むことができる. すなわち (2.4) において  $E_2$  が 1 語を占めるに至った後は,

<レベル 2>

$$x = (-1)^{s_0} F \times \beta^{-1} E_1 \times \beta E_2 \times \beta E_3 \quad (2.6)$$

という表示を採用する. ここで,  $F$  は“正規化条件を満たすある数”,  $E_1$  は“条件

$$1 \leq E_1 < \beta E_2 \times \beta E_3 - 1 \quad (2.7)$$

を満たすある数”,  $E_2, E_3$  は (2.1) の表示をしたときの  $E_1$  (これを  $E$  と記す) が

$$1 \times \beta E_2 \times \beta E_3 \leq E < 2 \times \beta (E_2 + 1) \times \beta E_3 - 1 \quad (2.8)$$

を満足し, かつ  $E_2$  と  $E_3$  が 1 語中に収まる範囲内で  $E_2$  の桁数なるべく大きくなるように選ぶ. レベル 3, 4, ... も同様に定める.

実際には,

(i) どのレベルで表示しているか,

(ii)  $F$  と  $E$  に割当てられている場所の境界 (レベル 0 の場合),  $E_i$  と  $E_{i+1}$  の境界 (レベル  $i \geq 1$ ) の場合) はどこか,

等のための情報も, 同じ語中に収めねばならない.

## 2.2 底の選択

浮動小数点表示の底  $\beta$  については, いろいろな観点から現在では  $\beta=2$  (2進法) を採用するのが最良であるとされている<sup>1), 2)</sup>. また, “正規化数の先頭ビットが 1” という事を利用してビット数の節約をはかれることも見逃せない長所である.

## 2.3 符号付き絶対値表示の採用

上記の 2.1, 2.2 の原則の長所を發揮させるためには, (2.1) の仮数部  $\pm F$ , 指数部  $\pm E$  は, “補数表示”, “けたばき表示” 等ではなく, 符号付き絶対値表示でなければならない. これを正規化条件 (2.3) と組み合わせることにより, 数の表示法にとって望ましい, 加法逆元, 乗法逆元に関する対称性<sup>11)</sup>も確保できる. (指数部の絶対値が最大の数では  $F=1$  であることに注意.)

## 2.4 丸め的方式

丸め的方式は特殊な目的 (たとえば区間演算) を除けば四捨五入 (2進なら 0 捨 1 入) が良いから, これを採用する. 他の丸め的方式でも本質的な差異は生じない. 以下では “0 捨 1 入” を前提とする.

## 2.5 “非数”の役割

固定長の 1 語で区別できる数の個数は有限であるが一方ではいくらかでも大きい (小さい) 数が現れるので, “ある数よりも大きい (小さい) 数” というようなものも 1 つの “数” として扱わざるを得ない. これはいわば “ $\pm\infty$ ” のようなものである. これが不可避であるならば, これを積極的に “数” の体系に組み入れて, 四則演算に関して閉じた, “あふれ” が起きても OS の厄介にならない体系を考えるのがよい. (このような方向への試みは Kahan<sup>8)</sup> にも見られるが, いくらかの曖昧性が残されている.) 四則演算に関する対称性を要請すれば, “非数” の体系は一意に確定する (4 を参照).

## 3. 新浮動小数点方式の実現法

レベルを無限にとることは不可能であり、実際的にはどこかで打切らねばならない。本節ではレベル0の実現についてまず考察する。より高次のレベルについてもほぼ同じ方針で実現可能であろう。

3.1 レベル0の実現法

2進法で正規化条件(2.3)のもとで、(2.1)の形に数  $x(≠0)$  を表わせば、次の形のビット並びになる：

$$F=1.f_2f_3\dots f_m \quad (3.1)$$

$$E=1e_{n-1}e_{n-2}\dots e_2e_1 \text{ (あるいは } E=0) \quad (3.2)$$

そこで数  $x$  を表わす情報としては、次の  $L=m+n$  ビットが必要である：

$$s_0f_2f_3\dots f_me_{n-1}e_{n-2}\dots e_2e_1s_1$$

( $n=1$  のときは  $E=1$ ;  $n=0$  のときは  $E=0$  で  $s_1$  は不用)。  $L$  が固定されていれば、  $n$  は  $0, 1, \dots, L-1$  のいずれかの値をとる。仮数部  $F$  と指数部  $E$  との“境界”を示すこの値  $n$  のために、  $\lceil \log_2 L \rceil$  ビットがさらに必要となる。そこで、

$$L + \lceil \log_2 L \rceil \quad (3.3)$$

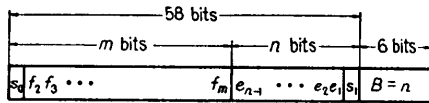
が1語長を超えないように  $L$  を選べばよい。数  $x$  に対して、このような表現が一意に定まることは明らかである。(なお、理想的な“0”は本体系では“非数”として扱う。4を参照。)

3.2 具体的な実現例

1語64ビットの場合を具体的に考える。(3.3)の条件から、  $L=58$ ,  $\lceil \log_2 L \rceil=6$  として図1のようなデータ構造を得る。

仮数部(の絶対値)の長さは58~1ビット(明示しない先頭ビットも含む)、指数部(符号も含む)の長さは0~57ビット( $B=n \neq 0$  のとき明示していない先頭ビットも含む)となり、  $|x|$  が約  $10^{-10^{16}}$  から  $10^{10^{16}}$  までの数  $x$  がレベル0で表現可能である。(表現誤差については5参照。)  $B$  の値のうち58~63は無意味となるが、それらは(i)非数、(ii)高次レベルの数、等を表わすのに使用できる。

このような(内部)表現を持つ(レベル0)の出力



仮数部  $(-1)^{s_0}F = (-1)^{s_0} \times (1.f_2f_3\dots f_m)_2$   
 指数部  $(-1)^{s_1}E = (-1)^{s_1} \times (1e_{n-1}\dots e_2e_1)_2$   
 $B=n$  は仮数部と指数部の境界を表わす。  
 ( $B=0, 1, \dots, 57$ ;  $B=0$  のとき  $E=0$ )

図1 新しい表現の浮動小数点数のデータ構造(レベル0)

Fig. 1 Data structure for a floating-point number in the proposed representation (level 0).

に際しては、印字幅のみを指定し、その内で指数部(10進整数)の場所をまず確保し、残りを(X.XX…の形に正規化された)仮数部に割当てればよい。(具体的には6の例を参照。)

数の入力の場合の表現は、従来許されている諸形式をすべて許容することとする。

4. 非数を含んだ数の体系

計算機で数値計算を行う際に扱う“数”は少なくとも四則演算に関して閉じた体系をなしている(演算が“合法的でない”ような被演算数の組があってもプログラムの正常な流れがそれによって妨げられない)ことが望ましい。そもそも“1語中に表現された数”は、数直線上の1点というよりは、ある区間(区間代数におけるようにそれを明示するかどうかは別として)を表わすと考えるべきである。2.3で触れた四則演算に関する対称性の観点からすると、我々がまず考えるべき実数は  $R^+UR^-$  である。ここで

$$R^+ = \{\text{正の実数}\}, R^- = \{\text{負の実数}\}. \quad (4.1)$$

3の表現法では  $R^+UR^-$  の一部分が表現可能である。考えるべき区間の“型”としては少なくとも次のもの(および  $\leq$  を  $<$  でおきかえたもの)が必要である。

- (i)  $\{x | a \leq x \leq b; a, b \in R^+\},$   
 $\{x | a \leq x \leq b; a, b \in R^-\};$
- (ii)  $\{x | a \leq x; a \in R^+\}, \{x | x \leq a; a \in R^-\};$
- (iii)  $\{x | 0 < x \leq a; a \in R^+\},$   
 $\{x | a \leq x < 0; a \in R^-\};$
- (iv)  $\{x | x \leq a \text{ or } b \leq x; a \in R^-, b \in R^+\};$
- (v)  $\{x | a \leq x \leq b; a \in R^-, b \in R^+\};$
- (vi)  $R^+, R^-;$
- (vii)  $R (= R^+UR^- \cup \{0\}).$

なぜならば、基本的な型(i)から出発して演算を繰返すと上記のすべての型が生起する可能性があるからである。これらを、それぞれ

- (i) +num, -num;
- (ii)  $+\infty, -\infty;$
- (iii) +0, -0;
- (iv)  $\infty;$
- (v) 0;
- (vi) +?, -?;
- (vii) ?

と略記し(図2)、型(i)以外のものを“非数”と呼ぶ。

3の表現方式で表現可能な数は  $\pm \text{num}$  の代表点であるとみなし、それらの間の演算は、代表点間の普通

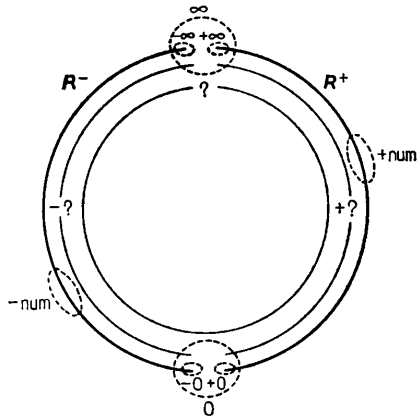


図 2 四則演算に関して閉じた数の体系の概念図

Fig. 2 Schematic diagram of the proposed number system including non-numbers.

の意味での演算と考え、演算結果は、

- a. 3の方法で表現可能ならば(必要なら丸めを行った後)  $\pm num$ ,

b. 絶対値が大きすぎて表現不能ならば  $\pm\infty$ ,

c. 絶対値が小さすぎて表現不能ならば,

c1. 仮数部に符号の情報が残っていれば  $\pm 0$ ,

c2. 仮数部に符号の情報が残っていなければ 0,

であると約束する。その他の非数がどのようにして出現するかを、また型(i)~(vi)を考えれば四則演算に関して閉じた体系が(厳密には“最小の体系”)が得られることを、表1~4に示す。

数(および非数)  $x$  と  $y$  の大小関係は、上記 a, b, c および表2の減算の結果の約束を用いて、 $x-y$  の演算結果が+ (あるいは-) の符号を有するとき  $x > y$  ( $x < y$ ) であると定義する。 $x=y$  は  $x$  と  $y$  が同一のビットパターンで表現されているときのみに限ることにする。こうすれば、たとえば「 $-\infty < -num < -0 < +0 < +num < +\infty$ 」, 「 $-num < 0 < +num$ 」, 「 $x < y$ ,  $u < v$  なら  $x+u < y+v$ 」等は成立する。しかし、「 $x = y, u < v$  なら  $x+u < y+v$ 」等が成立するとは限ら

表 1 非数を含んだ数の体系に対する加算表

Table 1 Addition table for the number system including non-numbers.

op 2 \ op 1	+num	-num	+∞	-∞	+0	-0	+?	-?	∞	?	0
+num	+num +∞	±num ±0,0	+∞	-∞	+num	+num	+?	?	∞	?	+num
+num	±num ±0,0	-num -∞	+∞	-∞	-num	-num	?	-?	∞	?	-num
+∞	+∞	+∞	+∞	?	+∞	+∞	+∞	?	?	?	+∞
-∞	-∞	-∞	?	-∞	-∞	-∞	?	-∞	?	?	-∞
+0	+num	-num	+∞	-∞	+0	0	+?	?	∞	?	0
-0	+num	-num	+∞	-∞	0	-0	?	-?	∞	?	0
+?	+?	?	+∞	?	+?	?	+?	?	?	?	?
-?	?	-num	?	-∞	?	-?	?	-?	?	?	?
∞	∞	∞	?	?	∞	∞	?	?	?	?	∞
?	?	?	?	?	?	?	?	?	?	?	?
0	+num	-num	+∞	-∞	0	0	?	?	∞	?	0

表 2 非数を含んだ数の体系に対する減算表

Table 2 Subtraction table for the number system including non-numbers.

op 2 \ op 1	+num	-num	+∞	-∞	+0	-0	+?	-?	∞	?	0
+num	±num ±0,0	+num +∞	-∞	+∞	+num	+num	?	+?	∞	?	+num
-num	-num -∞	±num ±0,0	-∞	+∞	-num	-num	-?	?	∞	?	-num
+∞	+∞	+∞	?	+∞	+∞	+∞	+∞	?	?	?	-∞
-∞	-∞	-∞	-∞	?	-∞	-∞	-∞	?	?	?	0
+0	-num	+num	-∞	+∞	0	+0	?	+?	∞	?	0
-0	-num	+num	-∞	+∞	0	0	-?	?	∞	?	0
+?	?	+?	?	+∞	?	+?	?	+?	?	?	?
-?	-?	?	-∞	?	-?	?	-?	?	?	?	?
∞	∞	∞	?	?	∞	∞	?	?	?	?	∞
?	?	?	?	?	?	?	?	?	?	?	?
0	-num	+num	-∞	+∞	0	0	?	?	∞	?	0

表 3 非数を含んだ数の体系に対する乗算表  
Table 3 Multiplication table for the number system including non-numbers.

op 1 \ op 2	+num	-num	+∞	-∞	+0	-0	+?	-?	∞	?	0
+num	+num +∞, +0	-num -∞, -0	+∞	-∞	+0	-0	+?	-?	∞	?	0
-num	-num -∞, -0	+num +∞, +0	-∞	+∞	-0	+0	-?	+?	∞	?	0
+∞	+∞	-∞	+∞	-∞	+?	-?	+?	-?	∞	?	?
-∞	-∞	+∞	-∞	+∞	-?	+?	-?	+?	∞	?	?
+0	+0	-0	+?	-?	+0	-0	+?	-?	?	?	0
-0	-0	+0	-?	+?	-0	+0	-?	+?	?	?	0
+?	+?	-?	+?	-?	+?	-?	+?	-?	?	?	?
-?	-?	+?	-?	+?	-?	+?	-?	+?	?	?	?
∞	∞	∞	∞	∞	?	?	?	?	∞	?	?
?	?	?	?	?	?	?	?	?	?	?	?
0	0	0	?	?	0	0	?	?	?	?	0

表 4 非数を含んだ数の体系に対する除算表  
Table 4 Division table for the number system including non-numbers.

op 1 \ op 2	+num	-num	+∞	-∞	+0	-0	+?	-?	∞	?	0
+num	-num +∞, +0	-num -∞, -0	+0	-0	+∞	-∞	+?	-?	0	?	∞
+num	-num -∞, -0	+num +∞, +0	-0	+0	-∞	+∞	-?	+?	0	?	∞
+∞	+∞	-∞	+?	-?	+∞	-∞	+?	-?	?	?	∞
-∞	-∞	+∞	-?	+?	-∞	+∞	-?	+?	?	?	∞
+0	+0	-0	+0	-0	+?	-?	+?	-?	0	?	?
-0	-0	+0	-0	+0	-?	+?	-?	+?	0	?	?
+?	+?	-?	+?	-?	+?	-?	+?	-?	?	?	?
-?	-?	+?	-?	+?	-?	+?	-?	+?	?	?	?
∞	∞	∞	?	?	∞	∞	?	?	?	?	?
?	?	?	?	?	?	?	?	?	?	?	∞
0	0	0	0	0	?	?	?	?	0	?	?

ない。

5. 表現誤差

ある実数  $x$  に対応する代表点  $x^*$  は、仮数部の最後の桁の次の桁を丸めた数であるとされている。(我々のレベル 0 の表示の場合もこうする。レベル 1 以上の場合にも、仮数部および指数部に対して同様な考え方を適用できる。) このとき、“相対誤差”

$$e_{rep}(x) = \frac{|x^* - x|}{|x|} \tag{5.1}$$

を  $x$  の表現誤差と呼ぶ。また、

$$u(x) = \max_y \{ e_{rep}(y) \mid y \text{ は } x \text{ と同じ代表点 } x^* \text{ を持つ} \} \tag{5.2}$$

を丸めの単位と呼ぶ。β 進法で仮数部が  $l$  桁のとき、

$$u(x) = \frac{1}{2} \beta^{l-1} \tag{5.3}$$

である。(丸めは四捨五入。) なお、非数  $x$  については、 $e_{rep}(x)$  は定義されない(あるいは ∞ である)と

するのがよいであろう。(この意味でも、0 や ±0 は非数とみなすのが適当である。)

異なる表現方式の特長を表現誤差、丸めの単位を用いて比較すると興味深い。1 語長を 64 ビットと仮定し比較対象として以下の典型的なものを考える。

〈1〉 IBM 型: -16 進, 仮数部 56 ビット (丸めは“切捨て”), 指数部 7 ビット (“げたばき”式), 符号 1 ビット. (“倍精度”実数型.)

〈2〉 Kahan 型: -2 進, 仮数部 52 ビット (先頭ビットを隠すので実質的には 53 ビット; 丸めは“0 捨 1 入”), 指数部 11 ビット (“げたばき”式), 符号 1 ビット. (“倍精度数”<sup>8)</sup>.) (正規化すると under-flow するとき, それを遅らせるため “非正規化数”<sup>8)</sup> も考える.)

〈3〉 Morris 型: -Morris の原論文<sup>10)</sup>の精度を 1 語長が 64 ビットの場合に改作したもの. 2 進, “指数部長”を表わす部分 3 ビット (これが表わす数を  $G(=0, 1, \dots, 7)$  とする), 残り 61 ビットを仮数部, 指

数部に割当てる(丸めは0捨1入). 指数部(符号つき絶対値表示)の長さは

- (i) G+1 (Morris (i) 型)
- (ii) G+4 (Morris (ii) 型)

の2通りの解釈が提案されている.

<4> 本論文の方式: -3.2 参照. レベル0を主として考える(丸めは0捨1入).

正規化条件としては, <1>と<3>は(2.2)を, <2>と<4>とは(2.3)を用いる.

$|\log|x||$  が比較的小さい範囲での  $e_{rep}(x)$  を図3に, 広い範囲にわたっての  $u(x)$  を図4に示す.

これらの図からもわかるように, 本論文の方式は, “あふれ” が起こる “不連続点” がなく, “普通の大き

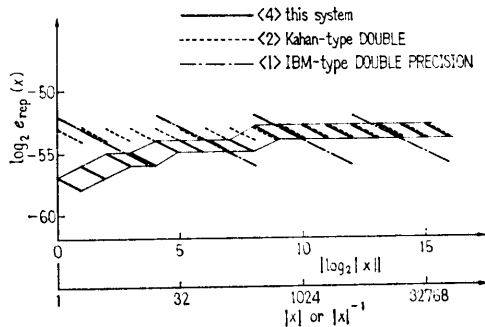


図3 各種表現方式による表現誤差  $e_{rep}(x)$  の比較  
Fig. 3 Representation error  $e_{rep}(x)$  of various systems.

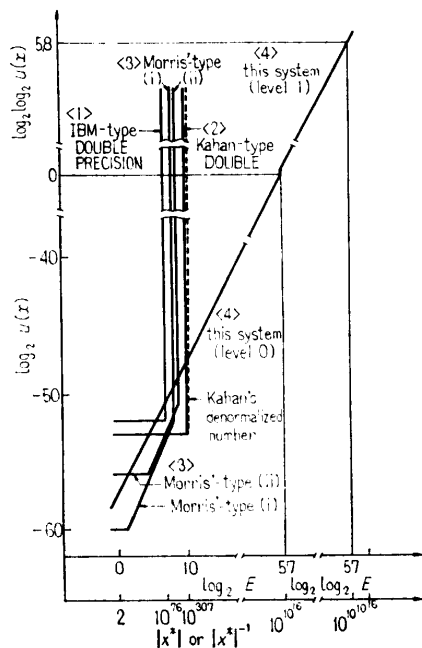


図4 各種表現方式による丸めの単位  $u(x)$  の比較  
Fig. 4 Units of round-off  $u(x)$  for various systems.

さ”の数に対しては IBM 型, Kahan 型より誤差が小さく, 数が極端に大きくあるいは小さくなくても連続的に誤差が増加する形で対処している. このような点からみて, 本方式は各種の表現方式を比較検討する際の一つの基準となりうるであろう.

## 6. 新体系による数値計算例

### 6.1 実験プログラム

本論文で提案した方式を3で述べた形で実現するプログラム(サブプログラム群)を東京大学の大型計算機センターの HITAC-8800/8700 システムの FORTRAN および教育用計算機センターの Melcom-Cosmo 900 の FORTRAN で作成した(現在, 下記のものがある).

- 1° 四則演算用サブルーチン (計約 300 行),
- 2° 入出力用関数副プログラム (計約 200 行),
- 3° 型変換 (FORTRAN のデータと本体系間の変換を用う) 用関数幅プログラム (計約 200 行),
- 4° 平方根用関数副プログラム (約 60 行),
- 5° 補助サブルーチン群 (比較, 丸め, 仮数部と指数部の分離と結合等の処理) (計約 300 行).

なお現在の版では非数は “0” だけが特殊なビットパターンとして含まれている.

プログラムは実験用であり, 速度向上のための技巧を凝らしてないので, 速度は遅いが(通常の倍精度実数型演算に比較して約 600~1,000 倍程度の時間を要する (HITAC での比較)), この点は大幅に改良の余地がある. 特に演算回路を金物化した場合には, IBM 型や Kahan 型の 1 割増し程度の時間で実現できる見込みがある.

### 6.2 数値計算例

極端に大きな数, 小さな数が出現する典型的な数値計算の例として代数方程式の Graeffe 法<sup>13)</sup>による解を考えてみる. 簡単のため, 正の単純零点のみの多項式

$$\begin{aligned}
 P_1(x) &= x^8 - 11x^7 + 45.35x^6 - 88.55x^5 \\
 &\quad + 86.7524x^4 \\
 &\quad - 43.274x^3 + 10.984x^2 \\
 &\quad - 1.32x + 0.0576 \\
 &= (x-0.1)(x-0.2)(x-0.3)(x-0.4) \\
 &\quad \times (x-1)(x-2)(x-3)(x-4) \quad (6.1)
 \end{aligned}$$

および

$$\begin{aligned}
 P_2(x) &= x^4 - 10.43857593020613x^3 \\
 &\quad + 40.58740567587410x^2 \\
 &\quad - 69.60408570545396x
 \end{aligned}$$

$$+44.36715614906059 \\ = (x-2)(x-e)(x-\sqrt{7.4})(x-3) \quad (6.2)$$

$$(e=2.718281828\dots, \sqrt{7.4}=2.7202941017\dots)$$

を考える。一般に、

$P(x) = x^n - a_{n-1}x^{n-1} + \dots + (-1)^{n-1}a_1x + (-1)^na_0$   
 の零点を  $\alpha_i (i=1, \dots, n)$  とするとき、 $\alpha_i^{(r)} (i=1, \dots, n)$   
 を零点とする多項式

$$P^{(r)}(x) = x^n - a_{n-1}^{(r)}x^{n-1} + \dots + (-1)^na_0^{(r)} \quad (6.3)$$

は簡単な漸化式

$$a_k^{(0)} = a_k, \\ a_k^{(r+1)} = \sum_{i+j=2k} (-1)^{i+j} a_i^{(r)} a_j^{(r)} \quad (\nu=0, \dots) \quad (6.4)$$

で計算され、 $\alpha_i$  の近似値は

$$\alpha_i \approx (a_{i-1}^{(r)} / a_i^{(r)})^{1/2\nu} \quad (6.5)$$

で計算できる。

$P_1(x)$  に対する反復計算の結果を図5に、  
 $P_2(x)$  に対するものを図6に示す。

図5においては  $\nu=7$  くらいですべての零点  
 が十分な精度で得られている。(その後無駄な  
 反復を行っても精度が悪化していないことに注  
 目すべきである。) IBM 型のもでは  $\nu=6$  くら  
 いで破綻をきたす。

図6においては、近接零点があるため、必要  
 な反復回数が増加し、また最終的に得られる零  
 点の近似値の精度もいくらか劣る。しかし、本  
 方式では  $\nu=16$  くらいで、10進10桁の精度  
 で零点が得られている。IBM 型では  $\nu=6$  くら  
 いで破綻をきたし(そのときの近似零点の精  
 度は最悪のものではわずか2桁)、Kahan 型の  
 ものも  $\nu=7$  までしかもたない(そのときの  
 近似零点の精度は、最悪のものでは2桁)。

### 7. 今後の課題

本論文では、新しい浮動小数点方式の基本的  
 な着想を提案し、その特長を従来の諸方式と比  
 較考察し、それを実現する試作プログラムによ  
 る数値計算例をあげた。しかしこの方式につ  
 いては更に検討すべき課題が数多く残されてい  
 るので、それらについて逐次検討結果を発表し  
 てゆきたい。主なものは次の通りである。

(i) 演算の金物化: 本来この種の方式は  
 演算を金物化することによって始めてその特長  
 が発揮できる。そのための演算回路の設計試作  
 を行うこと。理論的には従来方式とほぼ同程度

の速度が得られる見込みがある。

(ii) 新表現方式に適した丸め誤差理論の確立: 一  
 本論文では数の表現に直接関連した誤差のみを論じた  
 が、「表現誤差がほぼ一定」という従来の誤差理論<sup>13)</sup>  
 の仮定の成立しない本方式における、四則その他の演  
 算に伴う誤差の拡大、累積に関する理論を確立するこ  
 と。従来の理論が“非数”の場合をも含めてどのよう  
 に変更されるかを調べることは重要であろう。

(iii) 高次レベル、非数の実際的な取扱い: 非数  
 やレベル1以上の入出力をどうするか、どのようにし  
 て実現するか、またどのくらいのレベルまで考えるべ  
 きか、等々について、理論・実際面から検討の余地が  
 ある。

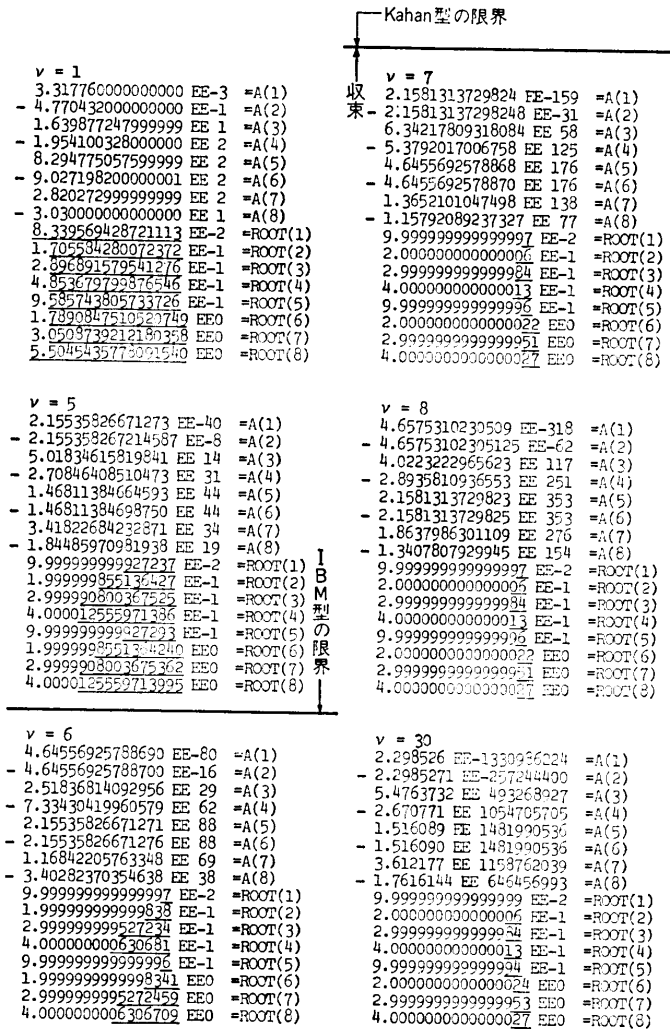


図5  $P_1(x)$  の零点を求める Graeffe 反復  
 Fig. 5 Iteration by the Graeffe method for  $P_1(x)$ .

		← Kahan型の限界	
<p><math>v = 1</math></p> <p>1.968444544755123 EE 3 =A(1)</p> <p>- 1.243233216278532 EE 3 =A(2)</p> <p>2.829367444181849 EE 2 =A(3)</p> <p>- 2.778905609893071 EE 1 =A(4)</p> <p>1.2583031670279898 EEO =ROOT(1)</p> <p>2.0961947379335457 EEO =ROOT(2)</p> <p>3.1908602092232556 EEO =ROOT(3)</p> <p>5.2715326138544105 EEO =ROOT(4)</p>	<p><math>v = 5</math></p> <p>5.08125451548211 EE 52 =A(1)</p> <p>- 1.18320169534957 EE 43 =A(2)</p> <p>3.02539330106157 EE 29 =A(3)</p> <p>- 2.01284256209763 EE 15 =A(4)</p> <p>1.9999931360575824 EEO =ROOT(1)</p> <p>2.6592468559773725 EEO =ROOT(2)</p> <p>2.7735140870227926 EEO =ROOT(3)</p> <p>3.0077660744101009 EEO =ROOT(4)</p>	<p><math>v = 7</math></p> <p>6.6662837510202 EE 210 =A(1)</p> <p>- 1.9590447225817 EE 172 =A(2)</p> <p>9.6227697791204 EE 116 =A(3)</p> <p>- 1.17902661948570 EE 61 =A(4)</p> <p>1.999999999999518 EEO =ROOT(1)</p> <p>2.7045783330773281 EEO =ROOT(2)</p> <p>2.7340770775614731 EEO =ROOT(3)</p> <p>3.0000001622446730 EEO =ROOT(4)</p>	<p><math>v = 16</math></p> <p>6.7375701284 EE 107941 =A(1)</p> <p>- 3.36284975669 EE 88213 =A(2)</p> <p>4.01365489311 EE 59751 =A(3)</p> <p>- 4.15479228694 EE 31268 =A(4)</p> <p>1.999999999999518 EEO =ROOT(1)</p> <p>2.7182818285951764 EEO =ROOT(2)</p> <p>2.7202941016100686 EEO =ROOT(3)</p> <p>3.0000000000009399 EEO =ROOT(4)</p>
		↑ 収束	
		I B M 型 の 限 界	
<p><math>v = 6</math></p> <p>2.5819147451107 EE 105 =A(1)</p> <p>- 1.39965879601065 EE 86 =A(2)</p> <p>4.38981732472593 EE 58 =A(3)</p> <p>- 3.44645651957942 EE 30 =A(4)</p> <p>1.999999999999518 EEO =ROOT(1)</p> <p>2.6899448172248597 EEO =ROOT(2)</p> <p>2.7487913512249591 EEO =ROOT(3)</p> <p>3.0001740432911321 EEO =ROOT(4)</p>			<p><math>v = 30</math></p> <p>1.577772 EE 1768518918 =A(1)</p> <p>- 3.759240 EE 1445290421 =A(2)</p> <p>2.5108124 EE 978970272 =A(3)</p> <p>- 2.0511859 EE 512305046 =A(4)</p> <p>1.999999999999518 EEO =ROOT(1)</p> <p>2.7182818285951764 EEO =ROOT(2)</p> <p>2.7202941016100686 EEO =ROOT(3)</p> <p>3.0000000000009399 EEO =ROOT(4)</p>

図 6  $P_4(x)$  の零点を求める Graeffe 反復  
Fig. 6 Iteration by the Graeffe method for  $P_4(x)$ .

8. おわりに

本研究に関しては、非公式研究会などにおいて、電気通信大学教授森口繁一先生、慶応大学教授相磯秀夫先生ほか多くの方々から貴重なご助言を賜り、また研究室の方々からも有用な討論をいただいた。それらに基づいて内容の多くの部分を改良することができた。ここに記して、感謝の意を表する次第である。

なお、本研究は文部省科学研究費の援助による部分を含んでいる。

参 考 文 献

- 1) Brent, R. P.: On the Precision Attainable with Various Floating-Point Number System. IEEE Trans. on Computers, Vol. C-22, No. 6, pp. 601-607 (1973).
- 2) Cody, W. J., Jr.: Static and Dynamic Numerical Characteristics of Floating-Point Arithmetic. IEEE Trans. on Computers, Vol. C-22, No. 6, pp. 598-601 (1973).
- 3) Edger, A. D. and Lee, S. C.: FOCUS: Microcomputer Number System, Comm. ACM, Vol. 22, No. 3, pp. 166-177 (1979).
- 4) Hamming, R. W.: On the Distribution of Numbers. The Bell System Technical Journal, Vol. 40, No. 8, pp. 1609-1625 (1970).
- 5) Herner, E. C. R. and Horspool, R. N. S.: A New Representation of the Rational Numbers

- for Fast Easy Arithmetic. SIAM J. on Computing, Vol. 8, No. 2, pp. 124-134 (1979).
- 6) 一松 信: 新標準浮動小数点体系の提案, 情報処理. Vol. 20, No. 9, pp. 793-797 (1979).
- 7) Ida, T. and Goto, E.: Over-flow Free and Variable Precision Computing in FLATS. Journal of Information Processing, Vol. 1, No. 3, pp. 140-142 (1978).
- 8) Kahan, W. and Palmer, J.: On a Proposed Floating-Point Standard, ACM SIGNUM Newsletter, Special Issue, pp. 13-21 (Oct. 1979).
- 9) Knuth, D. E.: The Art of Computer Programming, Vol. 2: Seminumerical Algorithms. Addison-Wesley, Reading, Massachusetts (1969).
- 10) Morris, R.: Tapered Floating Point: A New Floating-Point Representation. IEEE Trans. on Computers, Vol. C-20, No. 6, pp. 1678-1679 (1973).
- 11) Reinsh, C. H.: Principle and Preference for Computer Arithmetic. ACM SIGNUM Newsletter, Vol. 14, No. 1, pp. 12-27 (1979).
- 12) Swartzlander, E. E., Jr., and Alexopoulos, A. G.: The Sign/Logarithm Number System. IEEE Trans. on Computers, Vol. C-24, No. 12, pp. 1238-1242 (1975).
- 13) Wilkinson, J. H.: Rounding Errors in Algebraic Processes. Prentice-Hall, Englewood Cliffs (1963).

(昭和 55 年 3 月 14 日受付)

(昭和 55 年 4 月 22 日採録)



## 訂 正

Vol. 21, No. 4, pp. 306~313 掲載の松井, 伊理の論文「あふれのない浮動小数点表示方式」中の以下の誤りを訂正いたします.

i) p. 309 右下から 2, 3 行目

誤: 「 $x < y, u < v$  なら  $x+u < y+v$ 」 正: 「 $x < y, y < z$  なら  $x < z$ 」

ii) p. 310 左上から 8 行目の式 (5.1)

誤:  $e_{rep}(x) = \frac{|x^* - x|}{|x|}$  正:  $e_{rep}(x) = \max_y \left\{ \frac{|x^* - y|}{|y|} \mid y \text{ は } x \text{ と同じ代表点を持つ} \right\}$

iii) p. 310 左下から 5 行目の式 (5.2)

誤:  $u(x) = \max_y \{e_{rep}(y) \mid y \text{ は } x \text{ と同じ代表点 } x^* \text{ を持つ}\}$

正:  $u(x) = \max_y \{e_{rep}(y) \mid y \text{ は } x \text{ と同じ指数を持つ}\}$