

## 浮動小数点表現に代る実数値表現法

日立システム開発研 浜田 穂積

### まえがき

実数値の計算機内部における表現形式としての浮動小数点表現は、もちろんこれまで多大の貢献をしてきた。しかしながら現在の表現形式に対する不満がいくつか出てきた。その第一は、個々の計算機ごとにと決められた仕様が異なり、計算可能な数値の範囲、精度、丸めの方針などの不一致が計算結果を異なるものにする。あるいは、より効率を追求するため、内部表現のままデータの交換を行なおうとすると、仕様が完全に一致する必要があるのに、現状はそれから程遠い、などの不満がある。これを解決するものとして米国 U.C. Berkeley の W. Kahan 等による案をもとにした IEEE 標準<sup>1)2)</sup>案が作成された。

IEEE 標準案の特徴はまさに上述の不満を解消するためという点に集約されるが、その他にも二三の新しい着想がある。その一つは、丸めの方針をきちんと決めたこと、数でないもの（非数）の表現を導入すべきとしたことなどのほか、単精度では 8 ビット、倍精度では 11 ビットと指数部の長さを変

えたことである。ビット数の余裕のある倍精度で多少表現範囲が拡大されるが、それでもなおもう一つの不満であるオーバーフロー、アンダフローの発生という問題の根本的解決にはなっていない。

オーバーフロー、アンダフローについては、計算の途中で起こるとしても好ましいことではないので、これを防ぐ方法として松井・伊理による案<sup>3)</sup>が提出された。これは内部表現上で隣接する仮数部と指数部の境界を動かし得るものとし、指数部の長さを記録するフィールドを持つものである。この表現によると、代数方程式の解法の一つである Gräffe の方法を用いることができ、計算途中で指数の値がかなり大となって仮数部の精度が減少しても、最終結果は表現の本来持つ適切な精度が得られたと報告されている。

このようにいろいろ新しい表現法が提案されて、従来の形式についての不満が解消したかに見えたが、実は逆にもっと重要な点が満たされなくなってきた。というのは、データの形式がその長さごとに定義されるため、データの長さ、したがってその表す精度が異なるデータ間の相互変換が複雑になってきたことである。大型の電子計算機の場合は余裕があるから、32ビット、64ビットなど処理に好都合なデータ長の種類を少数選んで、その組合せを特別に扱えば何とか解決

できるが、マイコンコンピュータ、その他のデジタル応用機器やAD変換器等、条件の厳しいものの多種の結合を考えると、実数値の表現法を改めて考え直すなくてはならなくなってきた。

### 新しい実数値表現法の満たすべき条件

(4)5)

新しい実数値表現法を設定するにあたって、次の6条件を考慮した。これらのうち次の3条件を必須と考えた。

条件1：2進数としての指数、仮数から本表現へ、またその逆の変換が容易であること。このことは対数関数など、複雑な処理を用いなくて済むことを前提に置いている。

条件2：オーバフロー、アンダフローが事実上起らないこと。

条件3：表現仕様がデータの長さに依存せず、長さの異なるデータ間の相互変換が容易であること。

次の3条件は、必須というほどの必然性はないが、可能ならば満たしたいと考えた。

条件4：形式の取決め事項が不要あるいはなるべく少ないこと。

取決めのためのパラメータは少い程よく、ある一つの形式に自然に落ち着くようであれば望ましい。

条件5：固定長と考えたとき、すべてのパターンが異なる数値

に対応すること。従来の表現法には正規形、非正規形があり、異なるパターンが同じ値を表わす場合がある方式が多い。このようなパターンの無駄使いをさける。

条件6: 値の大小関係が、固定小数点と考えた場合のそれと一致すること。比較演算命令を固定小数点演算用の命令と共用できる。

### 新実数値表現法

表わそうとする数を  $x$  とする。これを

$$x = 2^e \cdot f \quad (1)$$

と、2つの数  $e$  と  $f$  で表現する。ここで値を一意的にするため  $e, f$  に次の条件を設ける。

$$e: \text{整数} \quad (2)$$

$$1 \leq f < 2 \quad \text{あるいは} \quad -2 \leq f < -1 \quad (3)$$

条件1により、この表現法は  $e, f$  のビットパターンから簡単な方法で得られるものとする必要がある。条件3を満たす変換法を得るために、より意味のある情報ほど左に記す。すなわち、数を大きく規定する情報ほど左に置き、長いデータへの変換は右に0のビットを補い、短いデータへの変換は左から必要な長さのビット数だけ切り出しとすることによって得られるものとする。このときビット0とビット1の意味づけ

を上手に選べば、条件6を満たすことが可能となる。具体的には次の順に置く。

- (i) 数  $x$  の符号 (1ビット)
- (ii) 指数部
- (iii) 仮数部

仮数の条件を(3)とすれば、 $x$ の符号にしたがって

$$\left. \begin{array}{l} f > 0 \text{ のとき } f = 01.f_1f_2 \dots f_j \dots \\ f < 0 \text{ のとき } f = 10.f_1f_2 \dots f_j \dots \end{array} \right\} \quad (4)$$

と表わせる。この表現は小数点の左に2ビットを有する、2進固定小数点(負数は2の補数)表現である。 $f_j$ についてはすべての組合せが起りうる。小数点の左の情報は符号ビットのみで定まる。したがって

符号部:  $x > 0$  のとき 0,  $x < 0$  のとき 1

仮数部:  $f_1f_2 \dots f_j \dots$

とすればよい。これを図示すると図1の通りである。

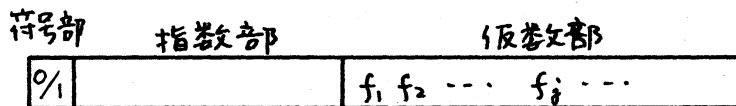


図1. データの構成

指数部は条件2を満たそうとするとき可変長にせざるを得ない。しかも指数の値に上限を設けたくない。また  $x = \pm 1$  となる  $e \sim 0$  のとき我々にとって重要な値であるから、なる

べく指数部を短くしたい。そこで自分自身で限界を知りうる表現として2次の可変長表現を採用する。

まず  $x > 0$  かつ  $e > 0$  の場合、 $e$  の値が2進  $i (> 0)$  桁で表わせる範囲は次の通りである。

$$2^{i-1} \leq e \leq 2^i - 1 \quad (5)$$

これに  $x$  の範囲も含めて  $x$  の範囲として表わせれば

$$2^{2^{i-1}} \leq x < 2^{2^i} \quad (6)$$

がある。次に  $e < 0$  のとき、同様に、対称性を考慮すれば

$$-2^{2^i} \leq x < -2^{2^{i-1}} \quad (7)$$

とすべきであり、これを  $e$  の範囲に翻訳すれば

$$-2^i \leq e \leq -2^{i-1} - 1 \quad (8)$$

となる。これは  $e$  を2の補数で表わすと解釈できる。このことを用いて、 $e$  の2進表現を以下の通りとしよう。

$$\left. \begin{array}{l} e > 0 \text{ のとき } 0 \cdots 0 1 e_{i-1} \cdots e_2 e_1 \\ e < 0 \text{ のとき } 1 \cdots 1 0 e_{i-1} \cdots e_2 e_1 \end{array} \right\} \quad (9)$$

(5), (8) には  $e = -1, 0$  の場合が含まれないが、これは  $i=0$  と考える。ここで定めた  $i$  と  $e_{i-1} \cdots e_2 e_1$  とから、 $i$  の識別ビット列を左に追加して指数部を次の通りに構成する。

$$\left. \begin{array}{l} e > 0 \text{ のとき } \underbrace{1 \cdots 1 0}_{i+1 \text{ 個}} e_{i-1} \cdots e_2 e_1 \\ e < 0 \text{ のとき } 0 \cdots 0 1 e_{i-1} \cdots e_2 e_1 \end{array} \right\} \quad (10)$$

(10)のビット列を符号のない2進数と考えたとき、固定小数点数としての大小順は、 $e$ の大小順に一致する。 $x$ の大小順は  $x > 0$  のとき  $e$ の大小順に一致するから、条件6を満たす。

$x < 0$  のときは、 $x$ の大小順と  $e$ の大小順は逆であるから、(10)による指数の順を逆にする。これは(10)の指数部の1の補数をとったものとして定める。

### 新表現法の形式的定義

前節のような慣例的定義法では、条件3, 4, 5, 6が満たされているの否かが直観的に理解しにくい。そこでこれと同じ内容を以下に示す区間分割に基づく方法で形式的に定義する。数学的にはこちらを正式な定義とし、前節のものを従と考える。まず一般論として次の事項から規定する。

任意のビット列  $S$  は実数値の一つの区間に対応する。すなわち

$$S : \{ x \mid a \leq x < b \}$$

これを次の記法を用いて示す。なお  $S$  が正しく表わすと考えられる値は区間の下限  $a$  であることにする。

$$I(S) = [a, b) \quad (11)$$

$S$  の右にビット0を連結したものを  $S0$ 、ビット1を連結したものを  $S1$  と記す。このとき、 $S$  の形および  $a, b$  の値によ

で定まる第3の値  $c$  によつて区間  $I(S)$  が二分され、次の対応づけが行なわれる。

$$I(S_0) = [a, c), \quad I(S_1) = [c, b) \quad (12)$$

上述のことから、条件3, 4, 5, 6 が満たされることが理解できよう。

さて、区間の分割を次の4段階の手順で行なう。

### (I) 疎分割

$$I(1) = [-\infty, 0) \quad (13-1)$$

$$I(0) = [0, +\infty) \quad (13-2)$$

から出発する。これらを  $-1, 1$  で分割して2次を得る。

$$I(10) = [-\infty, -1) \quad (14-1)$$

$$I(11) = [-1, 0) \quad (14-2)$$

$$I(00) = [0, 1) \quad (14-3)$$

$$I(01) = [1, +\infty) \quad (14-4)$$

さらにこれらを  $-2, -0.5, 0.5, 2$  で分割して2次を得る。

$$I(100) = [-\infty, -2) \quad (15-1)$$

$$I(101) = [-2, -1) \quad (15-2)$$

$$I(110) = [-1, -0.5) \quad (15-3)$$

$$I(111) = [-0.5, 0) \quad (15-4)$$

$$I(000) = [0, 0.5) \quad (15-5)$$

$$I(001) = [0.5, 1) \quad (15-6)$$



$$I(010) = [1, 2) \quad (15-7)$$

$$I(011) = [2, +\infty) \quad (15-8)$$

(15)の8式のうち, ビット列の第2ビットから右にみて, 0あるいは1の連が止まるという(15-2), (15-3), (15-6), (15-7)については直ちに(IV)の等差分割に進むが, その他の4つの場合は次の二重指数分割を行なう.

### (II) 二重指数分割

ここで  $n \geq 2$  の場合について帰納的に次の分割を行なう.

$$I(10^n) = [-\infty, -2^{2^{n-2}}) \quad \text{を} \quad -2^{2^{n-1}} \quad \text{で}$$

$$I(11^n) = [-2^{-2^{n-2}}, 0) \quad \text{を} \quad -2^{-2^{n-1}} \quad \text{で}$$

$$I(00^n) = [0, 2^{-2^{n-2}}) \quad \text{を} \quad 2^{-2^{n-1}} \quad \text{で}$$

$$I(01^n) = [2^{2^{n-2}}, +\infty) \quad \text{を} \quad 2^{2^{n-1}} \quad \text{で}$$

これにより, 0あるいは1の連の止まるという次の次を得る.

$$I(10^n 1) = [-2^{2^{n-1}}, -2^{2^{n-2}}) \quad (16-1)$$

$$I(11^n 0) = [-2^{-2^{n-2}}, -2^{-2^{n-1}}) \quad (16-2)$$

$$I(00^n 1) = [2^{-2^{n-1}}, 2^{-2^{n-2}}) \quad (16-3)$$

$$I(01^n 0) = [2^{2^{n-2}}, 2^{2^{n-1}}) \quad (16-4)$$

### (III) 等比分割

(16)の4つの場合について, 次の分割を  $n-2$  回行なう.

$$c = \pm \sqrt{ab} \quad (\text{符号は } a \text{ のそれと同じ}) \quad (17)$$

で分割する. これにより区間の両端の比は2となる.

## (IV) 等差分割

最後に次の分割を任意回行なう。

$$c = (a+b)/2 \quad (18)$$

2分割する。

以上の分割法は、ビット列と区間との対応関係の生成法であって、所望の長さのビット列が得られれば、必ずしもすべての手順を完了しなくても途中で打ち切、よい。

上述の形式的定義によつて得られるビット列と、従来の表現における符号、指数部、仮数部との関係はほぼ次の通り。

- ・(13)によつて得られるビット列は符号ビットに対応する。
- ・(14), (15)と(II), (III)によつて新たに得られるビット列は指数部に対応する。
- ・(IV)によつて新たに得られるビット列は仮数部に対応する。

非数

この節ではデシマルの長さは  $n$  ビット固定として考える。

次の6つの極端な場合を考える。

$$(i) \quad 10 \dots 00 \rightarrow -\infty$$

$$(ii) \quad 10 \dots 01 \rightarrow -2^{2^{n-3}}$$

$$(iii) \quad 11 \dots 11 \rightarrow -2^{-2^{n-3}}$$

$$(iv) \quad 00 \dots 00 \rightarrow 0$$

$$(v) 00 \cdots 01 \rightarrow 2^{-2^{n-3}}$$

$$(vi) 01 \cdots 11 \rightarrow 2^{2^{n-3}}$$

0 は (v) において  $n \rightarrow \infty$  とした極限と考えられ、すべてが 0 の列で表現することは自然である。(13) の定義はこの点を考慮してなされたものである。

オーバーフロー、アンダフローを無くしたいが、ビット・パタンの種類は有限であるから、例えば絶対値の小さい数同志の乗算を行なえばもっと絶対値の小さな結果が得られる……と、いずれは表現不可能な値になろう。そこで、表現可能な最も絶対値小のパターンを極限的な値と解釈し、それより絶対値小にはなり得ないことにすればよい。これは結局のところアンダフローと同じであるが、その値が現実の計算ではまず絶対にといていっていほど起りにくいものであって、たとえ起ったとしてもそういう状態であることがわかればよいという程度のものであれば、アンダフローと考えなくてもよいだろう。ちなみに  $n=32$  としたときの、表現可能な絶対値最小の値は  $2^{-2^{29}}$  で、10進数で云えば小数点の後に  $1.61 \times 10^8$  個あまりの 0 を並べた後に初めて 0 でない数の出現するような値である。この場合の表現を、通常の解釈ではなく、特別に状態と考へ非数と呼ぶ。非数の概念は IEEE 標準案で提案され、松井・伊理方式でも採用されている。両方式とも数のパターンと

は切りはなしで独立に定義されているので、数の表現と連続しない。しかるに本表現は以下のように数の表現と連続して自然な解釈が可能である。

$$00 \cdots 01 \rightarrow +0, \quad 11 \cdots 11 \rightarrow -0$$

次に対称性を考慮して

$$01 \cdots 11 \rightarrow +\infty, \quad 10 \cdots 01 \rightarrow -\infty$$

と定める。こうすると  $10 \cdots 00$  は  $-\infty$  ではなく単なる  $\infty$  として扱う方が便利である。  $+\infty$  と  $-\infty$  はビット・パターンでは近い関係にあり、  $10 \cdots 00$  はその中間にあるからである。以上をまとめて次の通りに定義する。

$$10 \cdots 00 \rightarrow \infty$$

$$10 \cdots 01 \rightarrow -\infty$$

$$11 \cdots 11 \rightarrow -0$$

$$00 \cdots 00 \rightarrow 0$$

$$00 \cdots 01 \rightarrow +0$$

$$01 \cdots 11 \rightarrow +\infty$$

(19)

### 表現法に固有の精度に関する評価

多くの浮動小数点表現では、仮数部が数  $x$  の大きさによらず一定の長さをとるので、表現法に固有の精度は相対誤差がほぼ一定という形で述べられ、量として言うときは用いられ

る基数における仮数の桁数とする。仮数部の桁数の変りうる  
 松井・伊理方式とか本表現ではより細く調べこみる必要があ  
 る。ここでは誤差として、表現可能な最小単位の半分のもの  
 があるとして、相対誤差  $E_r$  と絶対誤差  $E_a$  について考察する。

64ビット・データの相対誤差  $E_r$  について、代表的な4種の  
 表現についてものを図2に示す。横軸は表現しようとする  
 値  $x$  の、縦軸は相対誤差  $E_r$  の、それぞれの絶対値の2を底と  
 する対数で示す。図2はかなり複雑に重なりあ、こ見にくい

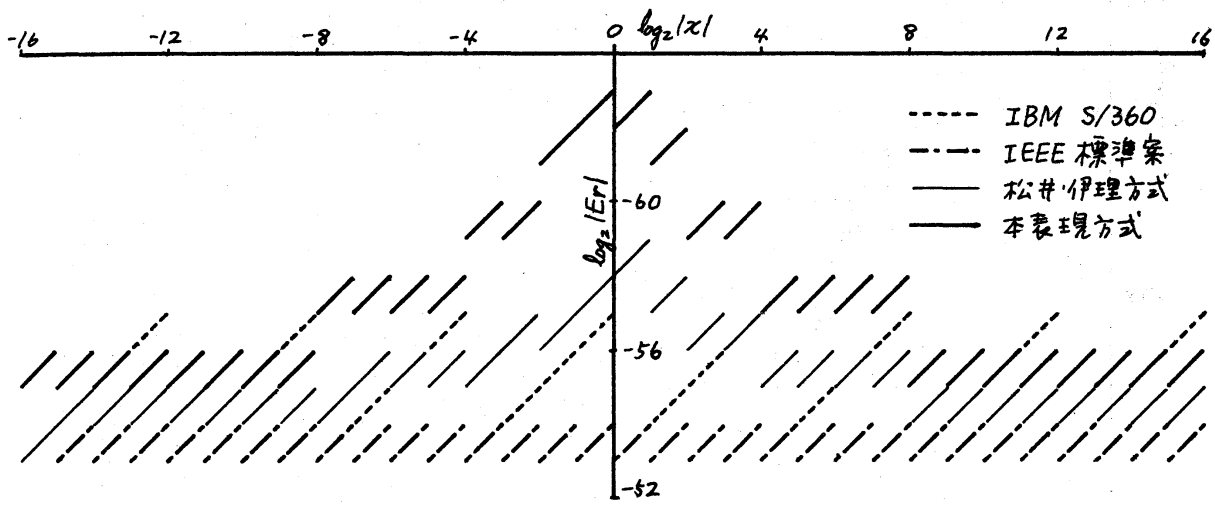


図2. 表現法固有の相対誤差

ので少し単純化する。同じ形が左右に並ぶ場合1本の横線で  
 結んだものと置き換えてみる。このような場合最も不利な条  
 件で考えねばならないので、各々の部分の下の縁を連ねなけ  
 ればならない。これを行なうとほぼ左右対称となる。そこで  
 右半分だけをとって、かつ横軸の尺度のさらに2を底とする  
 対数をとる。こうして得られたのが図3である。図4は同様

に32ビット・データの場合である。松井・伊理方式については、原論文では32ビット・データを定義しているが、容易に類推可能のためそれによった。

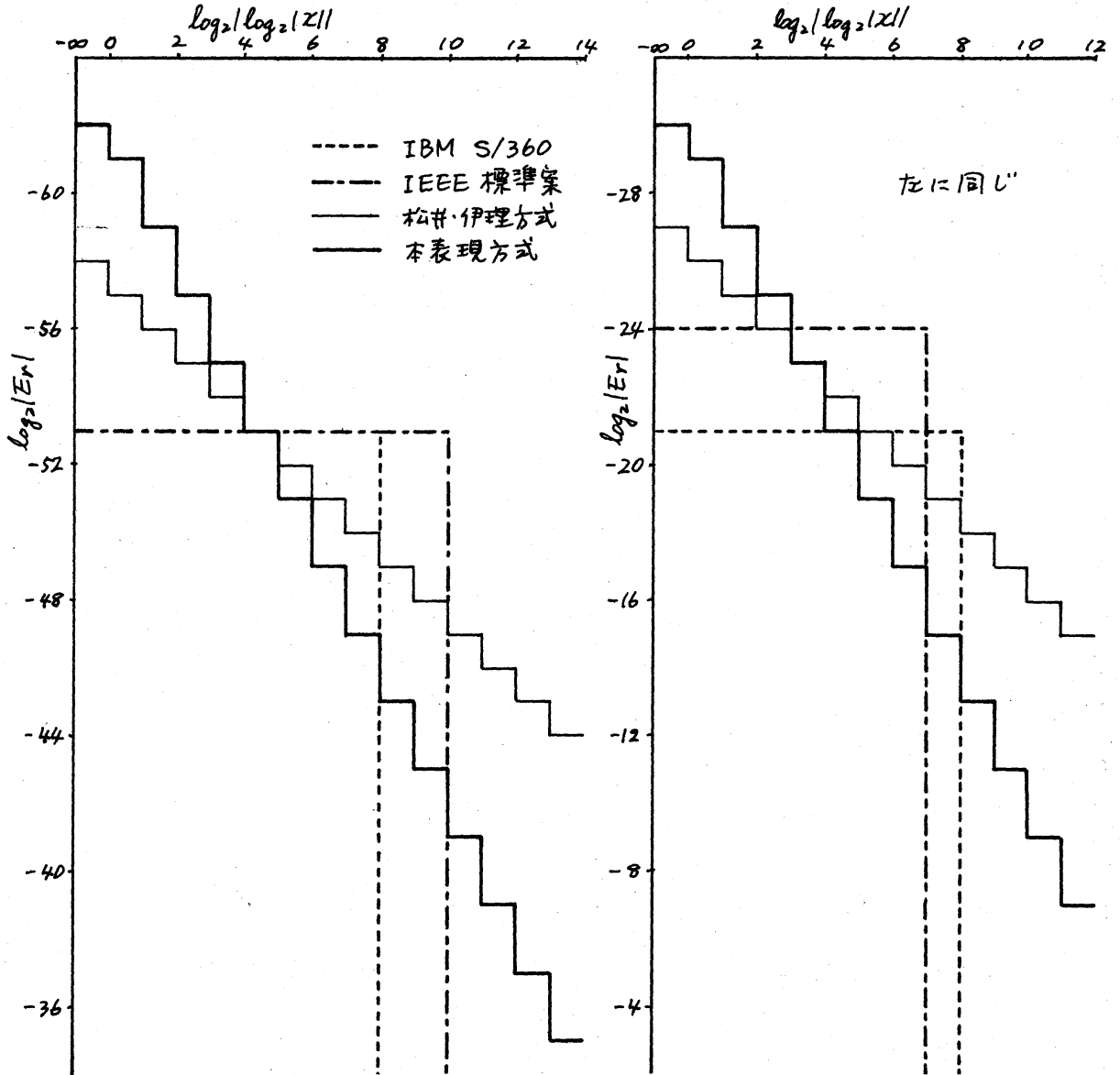


図3. 表現範囲と相対誤差

(64ビット・データ)

図4. 表現範囲と相対誤差

(32ビット・データ)

次に64ビット・データの場合の絶対誤差を図5に示す。な

おこれには松井・伊理方式に代えて、64ビットで $-2 \leq x < 2$ の範囲を表わす固定小数点表現の場合を追加して記してある。

これによると、本表現が固定小数点表現と比べて不利であるのは3箇所だけで、しかもいずれにおいても1ビットの不利

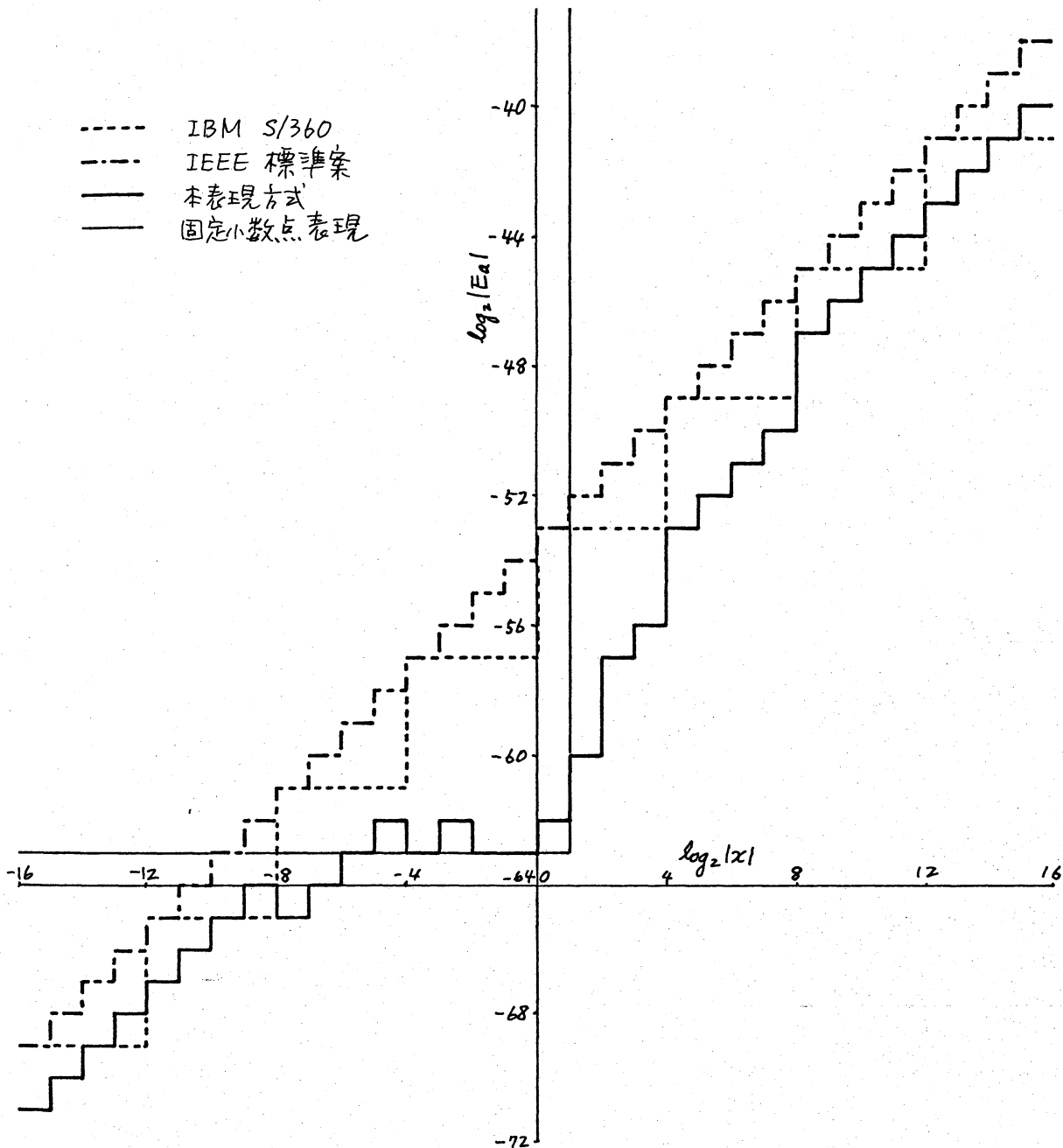


図5. 表現法固有の絶対誤差 (64ビット)

にすぎず、なおかつ浮動小数点表現の持つ長所、すなわち大きい数を表現でき、小さい数はより細かく表現可能という長所を併せ持つことがわかる。このことは別の観点から見るとき、ビット数の小さなデータもかなりの表現能力を持つことを示している。

以上の考察により、本表現は汎用大型計算機での実数値計算のみならず、ミニコンピュータ、マイクロコンピュータにも使え、また広義のアナログ・デジタル変換におけるデジタル側表現にも適する。これによりすべてのデジタル機器における実数値表現が一つの表現法で統一できるから、これら機器を有機的に結合して、自由な計算処理が可能となる。

### あとがき

現行の浮動小数点表現法に対して指摘されている欠点をほぼすべて取除くことのできる新しい実数値表現法を考察した。本表現法は最初に設定した6条件をすべて満たすことができた。これらを含めて、次の8項目が本表現法の特長である。

- (i) 2進数としての指数と仮数との間で容易に相互変換が可能である。
- (ii) オーバフロー、アンダフローが事実上発生しない。
- (iii) 表現の仕様がデータの長さによらず一様で、そのため



長さの異なる系を相互に簡単に結合できる。

- (iv) 形式の取決め事項が不要である。データの特性を定めるパラメータはデータの長さのみである。
- (v) データの長さを固定するとき、すべてのパターンは異なる数値に対応し、無駄がない。
- (vi) 値の大小関係が、固定小数点数と考えた場合のそれと一致する。
- (vii) 固定小数点表現と比べ、絶対誤差はたかだか1ビットの不利でしかない。
- (viii) データの延長で、任意の実数値を限りなく近似可能である。

本表現が数値解析的な実際の応用においてどのような結果をもたらすかについては、さらに研究を待たねばならない。これについては残念ながらここで言及できない。今後に残された問題点と考える。

なお文献[4]で規定した表現法と、文献[5]で規定した表現法との間には細部で相異がある。これは文献[4]を投稿した後で、計算機の論理回路構成を容易にするために変更した方がよいと考えたので、その結果を文献[5]に記したものである。結果的には理解の点でも後者が優れていた。そこで、本報告書も文献[5]に規定したものによる。

参考文献

- [1] J. Coonen, et al: A Proposed Standard for Binary Floating Point Arithmetic, ACM SIGNUM Newsletter Special Issue. (Oct. 1979), pp. 4-12.
- [2] D. Stevenson, et al: A Proposed Standard for Binary Floating-Point Arithmetic, Draft 8.0 of IEEE, COMPUTER (Mar. 1981), pp. 51-62.
- [3] 松井正一・伊理正夫: あぶれのない浮動小数点表示方式, 情報処理学会論文誌, Vol. 21, No.4, (1980年7月), pp. 306-313.
- [4] 浜田穂積: =重指数分割に基づくデータ長独立実数値表現法, 情報処理学会論文誌, Vol. 22, No.6, (1981年11月), pp. 521-526.
- [5] 浜田穂積: 浮動小数点表現に代る実数値表現法, 情報処理学会第23回全国大会予稿集 (1981年10月), pp. 151-152.