

IEEE浮動小数点数のハードウェアによる精度の違いについて

福井 義成
(株) 東芝, 総合情報システム部

分散処理を行なう場合, 浮動小数点数の表現方法が異なると色々と問題が起きる. さらに, 浮動小数点数の表現方法がIEEE形式である計算機でも, システムによって結果が異なっている. IEEEの浮動小数点数の表現は規格で決っているが, 規格で規定している以外の部分のインプリメント方法がシステムによって異なり, 計算結果が違ふことがある. 各種のワークステーションやスーパーコンピュータをネットワークで結合し, 分散処理を行なう場合, データ形式の互換性が重要となる. ここでは, 現存するのシステムではIEEE形式の浮動小数点であっても結果が異なることに対する注意を喚起したい.

Hardware Dependent Precision Difference of
IEEE Floating Point Numbers

Yoshinari Fukui

Total Information & Systems Division
TOSHIBA CORPORATION
580-1, Horikawa-cho, Saiwai-ku, Kawasaki 210, JAPAN

In distributed computing, problem arises when different floating point number formats are used on each computers. Even when IEEE standard floating point number is used in different systems, nonidentical results may be derived. The reason for these differences is the different implementations of floating point number format not specified in IEEE standard. Data format compatibility is important in distributed computing over network connected supercomputers and/or workstations. And cautions should be taken since IEEE standard floating point number format may have nonidentical values over different systems.

1. はじめに

ここでは、“IEEE浮動小数点数”の計算精度について述べる。筆者自身、ネットワークを使用した分散処理を行なう場合、浮動小数点数の表現方法について問題意識を持っていたところ、「IEEEの浮動小数点数の計算機でも計算結果が異なるものがある。」という話を聞き、試したところ、本当にシステムによって結果が異なっていた。IEEEの浮動小数点数の表現は規格(754, 854)で決っているが、規格で規定している以外の部分のインプリメント方法がシステムによって異なり、計算結果が違ふことがある。各種のワークステーションやスーパーコンピュータをネットワークで結合し、分散処理を行なう場合、データ形式の互換性が重要となる。

2. データ形式について

各種のワークステーションやスーパーコンピュータをネットワークで結合し、分散処理を行なう場合、相互にデータを交換する都合上、データ形式の互換性が重要である。通常、データ交換の対象となるものは以下の3つである。

- (1) 文字(数字, 英字, 漢字)
- (2) 整数
- (3) 浮動小数点数

データ形式が異なる場合の変換の方式には

- (1) 文字形式を介する
- (2) バイナリで変換を行なう

等がある。「文字形式を介する」方法は、FORTRANの書式付きの入出力などを利用するものである。整数や浮動小数点数を文字形式に変換し、その後、ASCII-EBCDIC変換等を行えばよい。この方法は非常に簡便であるが、以下のような問題点がある。

- (1) 桁数をどの程度にするか(桁あふれ等)
- (2) 浮動小数点数の2進(16進) - 10進変換の誤差(桁数に依存)
- (3) 処理時間がかかる(場合によっては本来の計算時間以上かかる)

「バイナリで変換を行なう」方法はCRAYなどで実現されており、FORTRANの書式なしの入出ルーチンの中で文字、整数、浮動小数点数の変換を自動的に行なってくれるものである。変数の属性はコンパイラが処理し、属性にあった入出ルーチンへの呼び出しを生成してくれる。

整数の場合、16ビット、32ビット、64ビットの整数などがある。また、負の値の

表現方法には2の補数表示，絶対値表示，1の補数表示などがあるが，現在では2の補数表示のみを考えればよいであろう．そのため，異なる表現同士の変換といっても，ビット長の変換だけでよいであろう．長い整数から短い整数に変換するためには，オーバーフローに注意する必要があるが，長い整数の下の部分を取り出せばよい．逆に，短い整数から長い整数に変換するためには，符号ビットの拡張をし，長さを合わせればよい．

浮動小数点数の表現には

IEEE方式
IBM方式
CRAY方式
VAX方式
ACOS6方式 (GE635)

等がある．精度についても32ビット，64ビット，128ビットなどがある．表現方法が異なる計算機の間でデータ交換を行なう場合，何等かの変換が必要となる．

複数の計算機の間で浮動小数点データを交換する場合，互いに同じ形式であれば問題ないが，異なる場合，

- (1) 変換誤差
- (2) オーバーフロー／アンダーフロー

などの問題が生じる．

3. IEEE浮動小数点数のハードウェアによる精度の違い

システム（ハードウェア，コンパイラなど）によって，同じIEEE形式の浮動小数点数であっても，計算結果が異なることがある．同じ浮動小数点数の表現であっても，大規模な計算を行なった場合，コンパイラが生成する演算順序が異なり，計算結果が微妙に（？）異なってしまうことは珍しいことではない．特に，ベクトル計算機を使用する時には，ベクトル化のためよくあることである．しかし，異なるシステムのIEEE形式の場合は，非常に単純な演算でも結果が異なることがある．最初に試したプログラムを下記の示す（テストはすべてFORTRANを使用して行なった．アセンブラ等は使用していない）．この例は $n=3, 30, 300, \dots$ とし， $1/n$ を n 回加算するプログラムである．色々なシステムで計算を行なったところ，結果は2つに分類できた．（計算結果2でSの値が1.0になっているにもかかわらず，誤差が出ているのはSの2進10進変換が10進17桁目を四捨五入しているためである．）

```

プログラム 1 .
program ex01
implicit real*8 (a-h,o-z)
read(5,*) n, loop
do 300 l=1,loop
  a=1.0d0/n
  s=0.0d0
  do 100 i=1,n
    s=s+a
100  continue
  d=s-1.0d0
  write(6,6100) l, n, s, d
6100  format(1x,i3,i10,1pd27.19,d27.19)
      n=n*10
300  continue
end

```

計算結果 1 .

loop	n	S の値	誤差
1	3	1.000000000000000000d+00	0.000000000000000000d+00
2	30	9.999999999999999000d-01	-1.110223024625157000d-16
3	300	9.999999999999961000d-01	-3.885780586188048000d-15
4	3000	9.99999999999564000d-01	-4.363176486776865000d-14
5	30000	9.99999999998996000d-01	-1.003641614261142000d-13
6	300000	1.00000000004593000d+00	4.593436742084123000d-12
7	3000000	1.00000000060662000d+00	6.06621419763087000d-11

計算結果 2 .

loop	n	S の値	誤差
1	3	1.000000000000000000d+00	-5.551115123125782700d-17
2	30	1.000000000000000000d+00	-1.387778780781445700d-17
3	300	1.000000000000000000d+00	6.418476861114186200d-17
4	3000	1.000000000000000000d+00	-3.339342691255353700d-17
5	30000	1.0000000000004000d+00	3.527993869267831400d-16
6	300000	9.999999999998000d-01	-1.976229509897953200d-15
7	3000000	1.0000000000038000d+00	3.797916842129822600d-14

計算結果 1 と 2 では n= 3000000 の場合、10進で約3桁も異なっている。演算の性質上

中間結果を保存するレジスタの長さが影響していると予想し、次のプログラムでのテストを行なった。このプログラムは1に 2^{-n} を加算し、 2^{-n} が小さくなり、積み残しが起きる状況をテストする。sub1では $1 + 2^{-n} + 2^{-n}$ を計算し、sub2では $2^{-n} + 2^{-n} + 1$ を計算している。この順序通りに計算すれば、sub1よりはsub2のほうが有限桁計算の影響が出にくいと考えられる。

```

プログラム 2 .
program ex02
read(5,*) m, n
call sub1(m,n)
call sub2(m,n)
end
subroutine sub1(m,n)
implicit real*8 (a-h,o-z)
dimension c(2)
a=1.0d0
b=1.0d0
do 100 i=1,m
100 b=b*0.5d0
do 300 i=m,n
do 200 j=1,2
200 c(j)=b
d=a+c(1)+c(2)
write(6,6100) i, d, d
6100 format(1x,i3,2x,z16)
300 b=b*0.5d0
end
subroutine sub2(m,n)
↓
d= c(1)+c(2)+a
↓
end

```

計算結果 3 .

```

1 + 2-n + 2-nの結果
m    dの内部表現
50  3ff0000000000008
51  3ff0000000000004
52  3ff0000000000002
53  3ff0000000000000 ←
54  3ff0000000000000

```

```

2-n + 2-n + 1の結果
m    dの内部表現
50  3ff0000000000008
51  3ff0000000000004
52  3ff0000000000002
53  3ff0000000000001 ←
54  3ff0000000000000

```

計算結果 4 .

```

1 + 2-n + 2-nの結果
m    dの内部表現
50  3ff0000000000008
51  3ff0000000000004
52  3ff0000000000002
53  3ff0000000000001 ←
54  3ff0000000000000

```

```

2-n + 2-n + 1の結果
m    dの内部表現
50  3ff0000000000008
51  3ff0000000000004
52  3ff0000000000002
53  3ff0000000000001 ←
54  3ff0000000000000

```

色々なシステムで計算を行なったところ、結果はやはり、計算結果3と4の2つに分類できた。計算結果3では、IEEE形式の浮動小数点数の倍精度の仮数部の長さは52ビットであるため、 2^{-53} のところでは変化が起きている。計算結果3の $1 + 2^{-n} + 2^{-n}$ の場合、 $1 + 2^{-n}$ の演算で積み残しが発生していると考えられる。 $2^{-n} + 2^{-n} + 1$ の場合は $2^{-n} + 2^{-n}$ を先に演算し、結果的に $1 + 2^{-n+1}$ と等価になるため、 2^{-53} のところでは積み残しは起きていない。計算結果4では 2^{-53} のところでは両者とも積み残しは起きていないので、中間結果を保存するレジスタが長いと想像できる。

以上の例では、計算結果は2つに分類できた。この2つの結果を出す環境は、今回のテスト(FORTRANを使用)では、以下のようであった。

(1) 誤差が小さいケース

- ・ $80 * 87 + MS - FORTRAN$
- ・ $68881 + F77$ で最適化を行なった場合

(2) 誤差が大きいケース

- ・ $68881 + F77$ で最適化を行わない場合
- ・ テストしたその他の環境

今回はFORTRAN環境で行なった。この問題はレジスタの扱いによって大きく影響が出るため、アセンブラ等を使用した場合は今回の結果と異なる可能性もある。 $80 * 87$ の場合は常に、拡張精度の演算を行なっているようである。 68881 の場合、最適化を行なわないと加算のところでは `f a d d` という命令を生成しており、最適化を行なうと `f a d d x` という命令を生成している。`f a d d` は倍精度演算、`f a d d x` は拡張精度演算を行なう。

4. 長い中間レジスタを使用することの是非

今回の計算例では、長い中間レジスタを使用することは良い結果を与えている。また、内積計算の場合、長い中間レジスタを使用することによって非常に良い計算結果を得ることができることは良く知られている。しかし、長い中間レジスタを使用することが困った結果を与えることもある。次の例(プログラム3)はACOS6の計算例である。ACOS6は単精度の浮動小数点数は36ビット(仮数部27ビット)であるが、演算結果は仮数部71ビットのレジスタに残る。

$$ZBZ02=ZB*ZB-(Z-Z0)**2$$

の計算では $ZB*ZB$ の結果は一度メモリーに格納されるため、27ビットの精度になるが、 $(Z-Z0)**2$ の結果は71ビットのレジスタに残っており、これと27ビットの数の差を取るため、ゼロにならない。ACOS6の場合は、浮動小数点レジスタが1つしか

いため、メモリーへの格納が多くなり、この現象が起きやすい。ZBZ02 の値をそのまま条件分岐に使用すると、意図しない方向への分岐が発生することがある（ただし、浮動小数点演算で、このようなことが発生する可能性のある条件分岐を行なうこと自体問題であるが、計算機のなかでこのような処理が行なわれていることを知っているユーザばかりではないであろう）。

プログラム 3.

```
Z=0
ZB=0.18996827E03
ZO=0.18996827E03
ZBZB=ZB*ZB
ZOZO=(Z-ZO)**2
WRITE(6,1000) "ZBZB      =",ZBZB
WRITE(6,1000) "(Z-ZO)**2=",ZOZO
ZBZ01=ZBZB-ZOZO
ZBZ02=ZB*ZB-(Z-ZO)**2
WRITE(6,1000) "ZBZ01    =",ZBZ01
WRITE(6,1000) "ZBZ02    =",ZBZ02
STOP
1000 FORMAT(" ",A10,E16.8)
END
```

計算結果 5.

```
ZBZB      = 0.36087943E 05
(Z-ZO)**2= 0.36087943E 05
ZBZ01     = 0.
ZBZ02     = 0.24581095E-06
```

5. おわりに

ここでテストしたどのシステムも 32 / 64 ビットの形式では IEEE の規格に違反しているわけではないが、コンパイラまで含めた環境で、拡張倍精度の使用方法に違いがあるようである。

もともと浮動小数点数の表現が異なるを意識しているシステム間で分散処理を行なう場合、浮動小数点数の違いに注意する人は多いであろうが、どれも“同じ” IEEE 形式であるシステム間で分散処理を行なう場合、浮動小数点数の違いに注意する人はあまり多くないかもしれない。使用者の立場からは、精度がよいほうが望ましいが、それ以前の問題として、結果が同じであってほしいと望みたい（このことは、拡張倍精度の扱いかたと演算順序の両方に関係しますが）。ここでは、少なくとも現在のシステムでは IEEE 形式の浮動小数点であっても結果が異なることに対する注意を喚起できればよいと考えます。

数値計算屋の立場からは、拡張倍精度をうまく利用できるほうが、望ましいと考えていますが、演算の高速性を追求することとは矛盾するのでしょうか？ 長いレジスタの利用は演算順序にも大きく影響されます。最近の高度に発達した(?)最適化 FORTRAN コンパイラを使用しているとプログラムを書く人間の意志が尊重されない（考えてい

る計算順序が尊重されない) ことがあり、歯がゆい思いをしています。やって欲しい部分だけの最適化と長いレジスタの活用がうまく制御できるようになってほしいと望んでいます。ユーザの立場からは、個々の演算が問題ではなく、計算全体の精度/速度が重要だからです。

参考文献

- [1] 別冊インターフェース, 数値演算プロセッサ, CQ出版社, 1987.
- [2] 数値計算ガイド, 日本サン・マイクロシステムズ(株)発行マニュアル, 1991.