

浮動小数点数

荒田 実樹

2017 年 6 月 22 日

目次

- ① コンピューター上での数の表現
- ② IEEE754 による浮動小数点数の表現
- ③ 浮動小数点数の演算の注意点
- ④ まとめ

数の表現：(固定長) 整数

16 ビット符号なし整数

- $3 = [0003]_{16} = [0000\ 0000\ 0000\ 0011]_2$
- $2017 = [07E1]_{16} = [0000\ 0111\ 1110\ 0001]_2$
- 表せる最大の整数： $2^{16} - 1 = 65535 = [FFFF]_{16} = [1111\ 1111\ 1111\ 1111]_2$

16 ビット符号あり整数 (2 の補数表現)

- $-1 = [FFFF]_{16} = [1111\ 1111\ 1111\ 1111]_2$
- $-2 = [FFFE]_{16} = [1111\ 1111\ 1111\ 1110]_2$
- 表せる最小の整数： $-2^{15} = [8000]_{16} = [1000\ 0000\ 0000\ 0000]_2$
- 表せる最大の整数： $2^{15} - 1 = [7FFF]_{16} = [0111\ 1111\ 1111\ 1111]_2$

32 ビットや 64 ビットの整数が使われることが多い

数の表現：固定小数点数

小数点の位置を決め、小数点以下 N 桁を表せるようにする。

例：小数点以下（2進で）8 桁

- $0.5 = 128 \times 2^{-8} = [00.80]_{16} = [0000\ 0000.1000\ 0000]_2$
- $1.25 = 320 \times 2^{-8} = [01.40]_{16} = [0000\ 0001.0100\ 0000]_2$
- $0.1 = 25.6 \times 2^{-8} \overset{\text{丸め}}{\rightsquigarrow} 26 \times 2^{-8} = [00.1A]_{16} = [0000\ 0000.0001\ 1010]_2$

固定小数点数の演算

- 足し算、引き算は整数と同じように計算できる
- かけ算は、整数の掛け算とビットシフトを組み合わせる
- 整数演算を流用できる

数の表現：浮動小数点数

精度 p と底 b (2 か 10 のことが多い) を決め、数を

$$[X_0.X_1X_2\dots X_{p-1}]_b \times b^e \quad (X_0 \neq 0)$$

の形で表し、仮数部 X_i と指数部 e をそれぞれ保存する。

例：2 進小数の場合

$$[1.X_1X_2X_3\dots X_{p-1}]_2 \times 2^e$$

- 先頭の桁は必ず 1 にできる

浮動小数点数の演算

- 整数や固定小数点数よりも複雑

固定小数点数 vs 浮動小数点数

	固定小数点数	浮動小数点数
メリット	実装が簡単	表せる数の範囲が広い
デメリット	表せる数の範囲が狭い	実装が複雑

- コンピューター上で普通に計算する際は、浮動小数点数が使われることが多い
- 固定小数点数を採用している言語としては、 $\text{T}_{\text{E}}\text{X}$ などがある
 - $\text{T}_{\text{E}}\text{X}$ が作られたのは 1970 年代で、IEEE754 が登場したのは 1985 年
- 貧弱な CPU だと浮動小数点数がハードウェア的に実装されていないケースもある

余談：それ以外の数の表し方

- これまで紹介したのは、使用するビット数が決まっている表し方
 - 普通の用途ではこれで十分
 - 速度やメモリを犠牲にしても、数を「もっと広い範囲で」「もっと正確に」表したい用途では
 - 多倍長整数：メモリの許す限り何桁でも
 - 有理数：(多倍長) 整数の比で表す
 - 任意精度計算：精度 100 桁でも 1000 桁でも（無限精度ではない）
- などを使う

コストを払ってでも高精度計算が必要な状況の例

「OS 標準の電卓アプリで 10 進小数の計算が正しくできない！」みたいなクレームの対策

- 浮動小数点数の規格
 - 初版は IEEE754-1985、最新版は IEEE754-2008
 - うちの学内ネットワークからはタダで見れるみたい
- 同じ規格に従っていれば、言語やマシンが違ってても同じ計算結果になる
 - 例：パソコン vs スマートフォン
- データ交換の形式（どういう数にどういうビット列を対応させるか）を定める
- 2進小数 (binary) と 10進小数 (decimal) が定義されているが、この小話では以後 2進小数 (binary) を扱う

IEEE754 binary

全体のビット数によって、何種類か定められている

	binary16	binary32	binary64	binary128
全体のビット数	16	32	64	128
仮数部の精度 p	11	24	53	113
指数部の範囲	$[-14, 15]$	$[-126, 127]$	$[-1022, 1023]$	$[-16382, 16383]$
指数部のバイアス	15	127	1023	16383
指数部のビット数	5	8	11	15
仮数部のビット数	10	23	52	112

どれを使う？

- よく使われるのは binary64（倍精度）と binary32（単精度）
- 最近では、画像処理や機械学習等の用途で 16 ビット浮動小数点数（半精度）の需要があるらしい（質より量？）
- binary128（四倍精度）は 誰も使っていない 一般的とはいえない
- ~~※87には80ビットの...いや何でもない~~

IEEE754 における数の種類

IEEE754 における浮動小数点数は、以下の 5 つのどれかに該当する：

- 正規化数 (normal)
- ゼロ (zero)
- 非正規化数 (subnormal)
- 無限大 (infinity)
- 非数 (NaN; Not a Number)

binary64 の場合にそれぞれ見ていく

IEEE754 における数の種類：正規化数 (normal)

正規化数とは

$[1.X_1X_2\dots X_{52}]_2 \times 2^e$ ($-1022 \leq e \leq 1023$) の形で表せる数

例

$$1.25 \times 2^{1023} = [1.4]_{16} \times 2^{1023} = [1.01]_2 \times 2^{1023}$$

- 符号：正 (0)
- 指数部のビット列： $1023 + 1023 = 2046 = [7fe]_{16} = [111\ 1111\ 1110]_2$
 - 指数部をビット列として表す際には、バイアス（この場合 1023）を加える。
- 仮数部のビット列： $[4000\ 0000\ 0000\ 0]_{16} = [0100\ \underbrace{00\dots 00}_{48\text{bits}}]_2$

まとめると

$$1.25 \times 2^{1023} = \text{binary64}([\underbrace{0}_{\text{符号}} \underbrace{11111111110}_{\text{指数部}} \underbrace{00100000\dots 0000}_{\text{仮数部 (52 ビット)}}]_2)$$

IEEE754 における数の種類：正規化数 (normal)

正規化数で表せる範囲

最小の正の正規化数：

$$\begin{aligned}
 1 \times 2^{-1022} &= [1.\underbrace{000 \dots 000}_{52\text{bits}}]_2 \times 2^{-1022} \\
 &= \text{binary64}([\underbrace{0}_{\text{符号}} \underbrace{000000000001}_{\text{指数部}} \underbrace{00000000 \dots 0000}_{\text{仮数部 (52 ビット)}}]_2)
 \end{aligned}$$

最大の正の正規化数：

$$\begin{aligned}
 (2 - 2^{-52}) \times 2^{1023} &= [1.\underbrace{111 \dots 111}_{52\text{bits}}]_2 \times 2^{1023} \\
 &= \text{binary64}([\underbrace{0}_{\text{符号}} \underbrace{11111111110}_{\text{指数部}} \underbrace{11111111 \dots 1111}_{\text{仮数部 (52 ビット)}}]_2)
 \end{aligned}$$

- 13 / 31

IEEE754 における数の種類：ゼロ (zero)

浮動小数点数におけるゼロとは

計算結果がゼロだった、または絶対値が（非）正規化数で表現できないほど小さかったことを表す

ゼロは正規化数では表せないので、専用のビット列を使って表す：

$$\begin{aligned}
 +0 &= \text{binary}_{64}([\underbrace{0}_{\text{符号}} \underbrace{000000000000}_{\text{指数部}} \underbrace{00000000 \dots 0000}_{\text{仮数部 (52 ビット)}}]_2), \\
 -0 &= \text{binary}_{64}([\underbrace{1}_{\text{符号}} \underbrace{000000000000}_{\text{指数部}} \underbrace{00000000 \dots 0000}_{\text{仮数部 (52 ビット)}}]_2)
 \end{aligned}$$

注意

- ゼロにも符号がある！
- $+0$ と -0 は比較演算では同一視される

IEEE754 における数の種類：非正規化数 (subnormal)

非正規化数とは

「指数部のビット列が0で、なおかつ仮数部のビット列が0でない」ものを使って、0より大きく 1×2^{-1022} 未満の数をコードしたもの

$$[0.X_1X_2 \dots X_{52}]_2 \times 2^{-1022}$$

例

$$\begin{aligned} 1.25 \times 2^{-1024} &= [1.4]_{16} \times 2^{-1024} = [1.01]_2 \times 2^{-1024} = [0.0101]_2 \times 2^{-1022} \\ &= \text{binary64}([\underbrace{0}_{\text{符号}} \underbrace{000000000000}_{\text{指数部}} \underbrace{01010000 \dots 0000}_{\text{仮数部 (52 ビット)}}]_2) \end{aligned}$$

- 仮数部のビット列： $[5000 \ 0000 \ 0000 \ 0]_{16} = [0101 \ \underbrace{00 \dots 00}_{48\text{bits}}]_2$

IEEE754 における数の種類：非正規化数 (subnormal)

非正規化数で表せる範囲

最小の正の非正規化数：

$$\begin{aligned}
 1 \times 2^{-1074} &= 1 \times 2^{-1022-52} = [0.00 \dots 0001]_2 \times 2^{-1022} \\
 &= \text{binary64}([\underbrace{0}_{\text{符号}} \underbrace{000000000000}_{\text{指数部}} \underbrace{0000 \dots 0001}_{\text{仮数部 (52 ビット)}}]_2)
 \end{aligned}$$

最大の正の非正規化数：

$$\begin{aligned}
 (1 - 2^{52}) \times 2^{-1022} &= [0.111 \dots 111]_2 \times 2^{-1022} \\
 &= \text{binary64}([\underbrace{0}_{\text{符号}} \underbrace{000000000000}_{\text{指数部}} \underbrace{1111 \dots 1111}_{\text{仮数部 (52 ビット)}}]_2)
 \end{aligned}$$

Diagram illustrating the distribution of floating-point numbers on a number line, showing the range from 0 to 2^{-1021} .

The number line is divided into two main regions:

- Subnormal Region (Left):** The range from 0 to 2^{-1022} . The amplitude (幅) is 2^{-1074} (subnormal).
- Normal Region (Right):** The range from 2^{-1022} to 2^{-1021} . The amplitude (幅) is 2^{-1074} (normal).

Arrows indicate the constant amplitude for each region, showing that the spacing between numbers is constant within each region.

- 17 / 31

IEEE754 における数の種類：無限大 (infinity)

浮動小数点数における無限大とは

計算結果の絶対値が正規化数で表現できないほど大きかった (2^{1024} 以上)、あるいは、 $1/0$ や $\log 0$ を計算しようとしたことを表す

binary64 における無限大は

$$\begin{aligned}
 +\infty &= \text{binary64}([\underbrace{0}_{\text{符号}} \underbrace{1111111111}_{\text{指数部}} \underbrace{00000000 \dots 0000}_{\text{仮数部 (52 ビット)}}]_2), \\
 -\infty &= \text{binary64}([\underbrace{1}_{\text{符号}} \underbrace{1111111111}_{\text{指数部}} \underbrace{00000000 \dots 0000}_{\text{仮数部 (52 ビット)}}]_2)
 \end{aligned}$$

の 2 つ

$+0$ と -0 の逆数は、それぞれの符号の無限大になる

IEEE754 における数の種類：非数 (NaN; Not a Number)

NaN とは

計算結果が実数として ill-defined だったことを表す (例: $0/0$, $\sqrt{-1}$, $\infty - \infty$)

性質

- NaN が絡む演算の結果は NaN となる
 - 例: $NaN \times 0 = NaN$
- 比較演算では「自身と同一でない」と判断される
 - これを利用して計算結果が NaN かどうかを判断できる

ビットパターンの例

binary64($\underbrace{0}_{\text{符号}} \underbrace{1111111111}_{\text{指数部}} \underbrace{10000000 \dots 0000}_{\text{仮数部 (52 ビット)}}_2$)

IEEE754 における数の種類：非数 (NaN; Not a Number)

余談：NaN の応用



- 仮数部におよそ 51 ビット分の情報を持てる
 - 仮数部は 52 ビットあるが、先頭ビットは NaN の種類 (quiet/signaling) を表すのに使われる
 - 仮数部が完全に 0 であってはいけない
- NaN tagging / NaN trick (スクリプト言語処理系の実装に使われるテクニック)
 - 一つの 64 ビット値に、スクリプト言語における値を保持できる (普通はデータの種類を表すのに数ビット、実際のデータを表すのに 64 ビット必要)
 - LuaJIT が発祥 (のはず; 2009 年ごろ) で、その後 JavaScript の処理系などでも採用されているらしい

IEEE754 における数の種類

まとめ

指数部のビット列	仮数部	種類	値
$[000\ 0000\ 0000]_2$	0	ゼロ	± 0
$[000\ 0000\ 0000]_2$	$\neq 0$	非正規化数	$\pm 0.[\text{仮数部}] \times 2^{-1022}$
$[000\ 0000\ 0001]_2, \dots$ $\dots, [111\ 1111\ 1110]_2$		正規化数	$\pm 1.[\text{仮数部}] \times 2^{(\text{指数部のビット列})-1023}$
$[111\ 1111\ 1111]_2$	0	無限大	$\pm \infty$
$[111\ 1111\ 1111]_2$	$\neq 0$	NaN	

演算と丸め方向

計算結果を正確に表せない場合にどうするか？

例

0.1 は 2 進法だと循環小数になる：

$$0.1 = \frac{1}{10} = [1.9999\cdots]_{16} \times 2^{-4} = [1.1001\ 1001\ 1001\cdots]_2 \times 2^{-4}$$

計算結果を正確に表せない場合は、近い値へ丸める：

- 最近接丸め： $0.1 \rightsquigarrow [1.1001\ \cdots\ 1001\ 1010]_2 \times 2^{-4}$
- 負の無限大方向： $0.1 \rightsquigarrow [1.1001\ \cdots\ 1001\ 1001]_2 \times 2^{-4}$
- 正の無限大方向： $0.1 \rightsquigarrow [1.1001\ \cdots\ 1001\ 1010]_2 \times 2^{-4}$
- ゼロ方向： $0.1 \rightsquigarrow [1.1001\ \cdots\ 1001\ 1001]_2 \times 2^{-4}$

丸め方向の利用：区間演算と精度保証

- 丸めが発生すると、正確な計算はできない
- それでも、丸め方向を上手く制御すると、計算結果の上界や下界を与えることはできる
- 詳しくは「区間演算」や「精度保証」で調べて

(2進) 浮動小数点数の罫：10進小数との兼ね合い

最近接丸めで $(0.1 + 0.1) + 0.1$ を計算してみよう

まず、 $0.1 + 0.1$ は

$$\begin{array}{rcl}
 0.1 \overset{\text{丸め}}{\rightsquigarrow} & [1.1001 \ 1001 \ \dots \ 1001 \ 1010]_2 \times 2^{-4} \\
 +) \ 0.1 \overset{\text{丸め}}{\rightsquigarrow} & [1.1001 \ 1001 \ \dots \ 1001 \ 1010]_2 \times 2^{-4} \\
 \hline
 & [11.0011 \ 0011 \ \dots \ 0011 \ 0100]_2 \times 2^{-4} \\
 & \overset{\text{丸め}}{\rightsquigarrow} [11.0011 \ 0011 \ \dots \ 0011 \ 010]_2 \times 2^{-4}
 \end{array}$$

なので、 $(0.1 + 0.1) + 0.1$ は

$$\begin{array}{rcl}
 & [11.0011 \ 0011 \ \dots \ 0011 \ 010]_2 \times 2^{-4} \\
 +) \ 0.1 \overset{\text{丸め}}{\rightsquigarrow} & [1.1001 \ 1001 \ \dots \ 1001 \ 1010]_2 \times 2^{-4} \\
 \hline
 & [100.1100 \ 1100 \ \dots \ 1100 \ 1110]_2 \times 2^{-4} \\
 & \overset{\text{丸め}}{\rightsquigarrow} [1.0011 \ 0011 \ \dots \ 0011 \ 0100]_2 \times 2^{-2}
 \end{array}$$

となる

(2 進) 浮動小数点数の罣：10 進小数との兼ね合い

一方、

$$0.3 = 3/10 = [1.3333\cdots]_{16} \times 2^{-2} = [1.0011\ 0011\ 0011\cdots]_2 \times 2^{-2}$$

$$\quad \quad \quad \overset{\text{丸め}}{\rightsquigarrow} [1.0011\ 0011\ \cdots\ 0011\ 0011]_2 \times 2^{-2}$$

なので、binary64 では $(0.1 + 0.1) + 0.1 \neq 0.3$ となる

解決策

10 進小数を正確に扱いたいなら 2 進の浮動小数点数ではなくて 10 進の（浮動）小数点数を使え

浮動小数点数の罫：演算の結合法則

$1 + 2^{-53} + 2^{-53}$ を (binary64 における最近接丸めで) 計算してみよう

$(1 + 2^{-53}) + 2^{-53}$ の場合：

$$(1 + 2^{-53}) + 2^{-53}$$

$1 + (2^{-53} + 2^{-53})$ の場合：

$$1 + (2^{-53} + 2^{-53})$$

浮動小数点数の罫：演算の結合法則

$1 + 2^{-53} + 2^{-53}$ を (binary64 における最近接丸めで) 計算してみよう

$(1 + 2^{-53}) + 2^{-53}$ の場合：

$$(1 + 2^{-53}) + 2^{-53} \xrightarrow{\text{丸め}} 1 + 2^{-53}$$

$1 + (2^{-53} + 2^{-53})$ の場合：

$$1 + (2^{-53} + 2^{-53}) = 1 + 2^{-52}$$

浮動小数点数の罫：演算の結合法則

$1 + 2^{-53} + 2^{-53}$ を (binary64 における最近接丸めで) 計算してみよう

$(1 + 2^{-53}) + 2^{-53}$ の場合：

$$(1 + 2^{-53}) + 2^{-53} \xrightarrow{\text{丸め}} 1 + 2^{-53} \xrightarrow{\text{丸め}} 1$$

$1 + (2^{-53} + 2^{-53})$ の場合：

$$1 + (2^{-53} + 2^{-53}) = 1 + 2^{-52} \quad (\text{これは binary64 で正確に表せる})$$

なので、 $(1 + 2^{-53}) + 2^{-53} \neq 1 + (2^{-53} + 2^{-53})$ となる (結合法則が成り立たない！)

浮動小数点数の罠：指数部のオーバーフロー・アンダーフロー

$\text{hypot}(x, y) = \sqrt{x^2 + y^2}$ 関数を次のように素朴に実装したとする：

```
double naive_hypot(double x, double y) {  
    return sqrt(x * x + y * y);  
}
```

(注：hypot は直角三角形の斜辺 (hypotenuse) の略)

問題点：指数部のオーバーフロー・アンダーフロー

$x = 3 \times 2^{1000}$, $y = 4 \times 2^{1000}$ に対して $\text{hypot}(x, y) = 5 \times 2^{1000}$ となるべきだが…？

浮動小数点数の罠：指数部のオーバーフロー・アンダーフロー

$\text{hypot}(x, y) = \sqrt{x^2 + y^2}$ 関数を次のように素朴に実装したとする：

```
double naive_hypot(double x, double y) {
    return sqrt(x * x + y * y);
}
```

(注：hypot は直角三角形の斜辺 (hypotenuse) の略)

問題点：指数部のオーバーフロー・アンダーフロー

$x = 3 \times 2^{1000}$, $y = 4 \times 2^{1000}$ に対して $\text{hypot}(x, y) = 5 \times 2^{1000}$ となるべきだが…？

$$x^2 = 9 \times 2^{2000} \quad , \quad y^2 = 16 \times 2^{2000}$$

浮動小数点数の罠：指数部のオーバーフロー・アンダーフロー

$\text{hypot}(x, y) = \sqrt{x^2 + y^2}$ 関数を次のように素朴に実装したとする：

```
double naive_hypot(double x, double y) {
    return sqrt(x * x + y * y);
}
```

(注：hypot は直角三角形の斜辺 (hypotenuse) の略)

問題点：指数部のオーバーフロー・アンダーフロー

$x = 3 \times 2^{1000}$, $y = 4 \times 2^{1000}$ に対して $\text{hypot}(x, y) = 5 \times 2^{1000}$ となるべきだが…？

$x^2 = 9 \times 2^{2000} \xrightarrow{\text{丸め}} +\infty$, $y^2 = 16 \times 2^{2000} \xrightarrow{\text{丸め}} +\infty$ なので、

$$\text{naive_hypot}(x, y) = \text{sqrt}((+\infty) + (+\infty)) = +\infty$$

となる (不適切！)

浮動小数点数の罣：指数部のオーバーフロー・アンダーフロー

解決策

- 計算前に x, y の指数部を適切な範囲に収める
(この場合なら $m = 1002$ として $2^m \cdot \sqrt{(x \cdot 2^{-m})^2 + (y \cdot 2^{-m})^2}$ と計算させる)
- または、割り算を使って

$$\begin{cases} |y| \cdot \sqrt{1 + (x/y)^2} & (|x| \leq |y|) \\ |x| \cdot \sqrt{1 + (y/x)^2} & (|y| \leq |x|) \end{cases}$$

と計算する

複素数の除算（逆数）にも似たような罣がある

浮動小数点数の罣：情報落ち・桁落ち

$\sinh x = \frac{\exp(x) - \exp(-x)}{2}$ を次のように素朴に実装したとする：

```
double naive_sinh(double x) {  
    return (exp(x) - exp(-x)) / 2.0;  
}
```

問題点：0 の近くでの挙動が不適切

$|x| \ll 1$ の場合 $\sinh x \approx x$ となるべきだが…？

浮動小数点数の罫：情報落ち・桁落ち

$\sinh x = \frac{\exp(x) - \exp(-x)}{2}$ を次のように素朴に実装したとする：

```
double naive_sinh(double x) {
    return (exp(x) - exp(-x)) / 2.0;
}
```

問題点：0 の近くでの挙動が不適切

$|x| \ll 1$ の場合 $\sinh x \approx x$ となるべきだが…？
 $x = 2^{-100}$ の場合

$$\exp(2^{-100}) = 1 + 2^{-100} + \dots, \quad \exp(-2^{-100}) = 1 - 2^{-100} + \dots$$

浮動小数点数の罫：情報落ち・桁落ち

$\sinh x = \frac{\exp(x) - \exp(-x)}{2}$ を次のように素朴に実装したとする：

```
double naive_sinh(double x) {
    return (exp(x) - exp(-x)) / 2.0;
}
```

問題点：0 の近くでの挙動が不適切

$|x| \ll 1$ の場合 $\sinh x \approx x$ となるべきだが…？
 $x = 2^{-100}$ の場合

$$\exp(2^{-100}) = 1 + 2^{-100} + \dots \overset{\text{丸め}}{\rightsquigarrow} 1, \quad \exp(-2^{-100}) = 1 - 2^{-100} + \dots \overset{\text{丸め}}{\rightsquigarrow} 1$$

なので、 $\text{naive_sinh}(2^{-100}) = (1 - 1)/2 = 0$ となる

浮動小数点数の罣：情報落ち・桁落ち

解決策

- $\text{expm1}(x) = \exp(x) - 1 = x + x^2/2 + \dots$ を導入
- ```
double sinh(double x) {
 return (expm1(x) - expm1(-x)) / 2.0;
}
```

と定義すれば良い

- $$\frac{\text{expm1}(2^{-100}) - \text{expm1}(-2^{-100})}{2.0} = \frac{2^{-100} - (-2^{-100})}{2.0} = 2^{-100}$$

# TL;DR

- コンピューターで計算をさせるときは浮動小数点数の性質を知っておこう
- 浮動小数点数の長所も短所も把握した上で、上手く付き合っていこう
- 「0.1 を 10 回足しても 1.0 にならない、これはバグだ」と大騒ぎするような大人にはなるな