

# 浮動小数点数 (実数) の表現

コンピュータシステム 第3回講義

<http://www.cuc.ac.jp/~miyata/classes/system/>

## 本日の内容

- 少数点のある2進数を10進数に変換
- 少数点のある10進数を2進数に変換
- 浮動小数点数
- 浮動小数点数の内部形式
- 浮動小数点数と誤差

2

## 小数がある2進数⇒10進数

3

- 小数点のある2進数  $b_n \cdots b_1 b_0 . a_1 a_2 a_3 \cdots$  は  
 $2^n b_n + \cdots + 2^1 b_1 + 2^0 b_0 + 2^{-1} a_1 + 2^{-2} a_2 + 2^{-3} a_3 + \cdots$   
を表している

- 例: 2進数 101.011 は  
 $1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3}$   
 $= 1 \times 4 + 0 \times 2 + 1 \times 1 + 0 \times \frac{1}{2} + 1 \times \frac{1}{4} + 1 \times \frac{1}{8}$   
 $= 4 + 1 + 0.25 + 0.125$   
 $= 5.375$

- 練習: 2進数 11.1011 を10進数にしなさい

## 小数がある10進数⇒2進数

4

- 小数点のある10進数から2進数への変換方法
  - 整数部と小数部に分ける (整数部はふつうに2進数に変換)
  - 小数部に2を掛けて整数部 (0か1) を取り出す
  - 小数部が0になるか、適当な精度まで求めたら終了。

- 例: 10進数の「3.7」は  
3 → 「11」  
 $0.7 \times 2 \rightarrow \underline{1}.4$   
 $0.4 \times 2 \rightarrow \underline{0}.8$   
 $0.8 \times 2 \rightarrow \underline{1}.6$   
 $0.6 \times 2 \rightarrow \underline{1}.2$   
 $0.2 \times 2 \rightarrow \underline{0}.4$   
以下、0110を繰り返す

よって、11.101100110011001100110...

- 練習: 10進数の「1.1」を2進数で表せ。ただし小数点以下11桁目以降を切り捨てて、小数点以下10桁まで求めなさい。

## 浮動小数点数 (floating point number)

5

- 限られた情報量 (ビット数) で広い範囲の数値を表現する方法

$$\pm m \times R^e$$

の形で実数を表現する。 $m$ を仮数部 (mantissa),  $e$ を指数部 (exponent) と呼ぶ。また  $R$ は基数 (radix)と呼ばれ, コンピュータでは  $R = 2$  が使われる (人間は一般的に  $R = 10$  を使う)。

例: 基数  $R = 10$  とする

$$1\ 2\ 3 = 1\ 2.3 \times 10^1 = 1.2\ 3 \times 10^2 = 0.1\ 2\ 3 \times 10^3$$

- 0以外の数は指数を調節して  $1 \leq m < R$  となるように表現することができる。これを正規化表現という。

上の例では 「 $1.2\ 3 \times 10^2$ 」が正規化表現

## 浮動小数点数の例 (10進数)

6

- 簡単のため基数 10 (10進数) で説明する。
- 仮数部3桁, 指数部を符号付で1桁 ( $-9 \sim 9$ ) として正規化表現だけを許すものとする。

- 光の速さ (m/s) :

$$2\ 9\ 9\ 7\ 9\ 2\ 4\ 5\ 8 \rightarrow 3.0\ 0 \times 10^8 \quad \text{※四捨五入した}$$

- ロト6の1等当選確率

$$0.0\ 0\ 0\ 0\ 0\ 0\ 1\ 6\ 4\ 0\ 2\ 9\ 7\ 7 \rightarrow 1.6\ 4 \times 10^{-7}$$

- 10進4桁 (+符号情報2bit) の情報量で幅広い数を扱うことができる

- 最大  $9.9\ 9 \times 10^9 = 9,990,000,000$
- 正の数で最小  $1.0\ 0 \times 10^{-9} = 0.000000001$
- 負の数で最大  $-1.0\ 0 \times 10^{-9} = -0.000000001$
- 最小  $-9.9\ 9 \times 10^9 = -9,990,000,000$

## 浮動小数点数どうしの演算

7

- 簡単のため仮数部2桁とする
- 乗算 (仮数部どうしを乗算, 指数部は加算)
  - $8.4 \times 10^5 \times 2.1 \times 10^2$   
 $(8.4 \times 2.1) \times 10^{5+2}$   
 $= 17.64 \times 10^7$   
 $\rightarrow 1.8 \times 10^8$
- 除算 (仮数部どうしを除算して指数部は減算)
- 加算 (指数部をそろえてから計算)
  - $8.4 \times 10^5 + 2.1 \times 10^2$   
 $8.4 \times 10^5 + 0.0021 \times 10^5$   
 $= 8.4021 \times 10^5$   
 $\rightarrow 8.4 \times 10^5$
- 減算 (指数部をそろえてから計算)

## 浮動小数点数の例 (2進数)

8

- 次の2進数を正規化し仮数部4桁の浮動小数点数にせよ。ただし仮数部の小数点以下第4桁は0捨1入せよ。

- $(1\ 1\ 0\ 1\ 0\ 0\ 1\ 0)_2$

$$\rightarrow (1.1\ 0\ 1)_2 \times 2^7$$

- $(0.0\ 0\ 1\ 1\ 0\ 0\ 1\ 1\ 0\ 0)_2$

$$\rightarrow (1.1\ 0\ 1)_2 \times 2^{-3}$$

- ※ 正規化すると先頭は“1.”で始まることに注意 (ただし0は例外!)

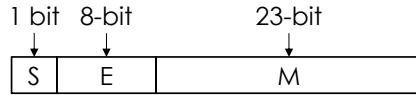
- 練習: 次の2進数を正規化し仮数部4桁の浮動小数点数にせよ。ただし仮数部の小数点以下第4桁は0捨1入せよ。

- $(1\ 1\ 1\ 0\ 1\ 0\ 0)_2$

- $(0.0\ 0\ 0\ 1\ 0\ 1\ 0\ 1)_2$

# 浮動小数点数の内部形式

## IEEE754 単精度浮動小数点 (32 bit)



- S: 仮数部の符号 (0: 正, 1: 負)
- E: 指数部
- M: 仮数部 (の少数部)

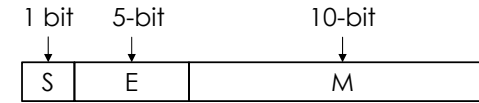
- 正規数として表現できる範囲: およそ  $\pm 1.0 \times 10^{-38} \sim \pm 1.0 \times 10^{38}$
- 精度: 10進数で7桁程度

## IEEE754 倍精度浮動小数点 (64 bit)



- 正規数として表現できる範囲: およそ  $\pm 1.0 \times 10^{-308} \sim \pm 1.0 \times 10^{308}$
- 精度: 10進数で16桁程度

# IEEE754 半精度浮動小数点数 (16bit)



$(-1)^S \times (1 + M) \times 2^{E-15}$  (正規化数)  
 $(-1)^S \times (0 + M) \times 2^{-14}$  (非正規化数)

- S: 仮数部の符号 (0: 正, 1: 負)
- E: 指数部
  - 符号なし整数 ( $0 \leq E \leq 31$ ), オフセット表現 (オフセット15)
  - $E=0$  のときは0または非正規化数,  $E=255$  のときは,  $\pm$ 無限大または非数 (NaN)
- M: 仮数の少数部
  - $E=0$  のとき, 正規化されていない仮数の少数部
  - $1 \leq E \leq 30$  のとき, 正規化された仮数の少数部
  - $E=31$  のとき, 無限大 ( $M=0$ ), NaN ( $M \neq 0$ )
- 特殊なハードウェアでしか使われていない (あまり一般的ではないが, 試験問題向きなので取り上げた)

# 補足

- 2進数の場合 (0以外の数は) 正規化すると必ず"1."で始まるので仮数部には小数部 ("1.XXXXX"のXXXXXの部分) だけを格納する。

- 0は例外扱い。+0と-0がある。指数部0の場合非正規数。
- 無限大や非数 (NaN) がある

- 単精度浮動小数点数の仮数部は実質的に11ビットの精度

- 10進数で3桁程度の精度

$$\log_{10} 2^{11} = 11 \log_{10} 2 = 11 \times 0.301 = 3.311$$

- 正規化数として表現できる正の数の範囲は  $6 \times 10^{-5} \sim 6 \times 10^4$  程度
- 指数部を2の補数表現にせずに, オフセット表現 (15足した値) にするのは, 数値の大小の比較を簡単にするため。

# 半精度浮動小数点数の例

S	E	M	
0	10000	0100000000000000	$(1.01)_2 \times 2^{16-15} = (10.1)_2 = 2.5$
1	01110	0100000000000000	$(-1.01)_2 \times 2^{14-15} = (-0.101)_2 = -0.625$
0	00000	0000000000000000	+0
1	00000	0000000000000000	-0
0	00000	0100000000000000	$(0.01)_2 \times 2^{-14} = 2^{-16}$ (非正規数)
0	11111	0000000000000000	$+\infty$
0	11111	0100000000000000	NaN (非数)

# 例題

- 「-10.75」をIEEE754半精度浮動小数点形式で表せ

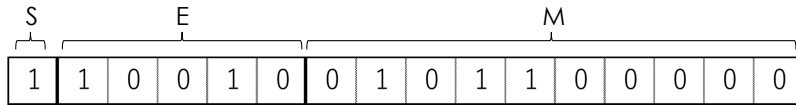
M: 10.75を2進数に変換して正規化

$$(10.75)_{10} = (1010.11)_2 = (1.01011)_2 \times 2^3$$

E: 指数の3にオフセット15を足すと18

$$(18)_{10} = (10010)_2$$

S: 符号部は1



# 浮動小数点数と誤差 (1)

- 丸め誤差
  - 仮数部の桁が有限なので、仮数部に入りきらない桁は四捨五入や切り上げ、切り捨てされることになる。このときに生じる誤差。
- オーバーフロー (overflow, 桁あふれ)
  - 計算結果 (の絶対値) が大きすぎて、正規数として表現できる範囲を超えること
- アンダーフロー (underflow)
  - 計算結果 (の絶対値) が小さすぎて、正規数として表現できる範囲を超えること

# 浮動小数点数と誤差 (2)

- 情報落ち誤差
  - (絶対値の) 大きな数と小さな数の加減算を行ったときに、小さいほうの値が演算結果に反映されないことによって生じる誤差。
  - スライド7参照
  - データの総和を出すような場合、小さいものから足すことで回避できる場合もある
- 桁落ち
  - (絶対値の) 近い2数で加減算を行ったときに精度が失われること

$$\begin{array}{r}
 1.1110110101 \\
 -) 1.1110110001 \\
 \hline
 0.0000000100
 \end{array}$$

11ビットの情報が3ビットに減少!