

IEEE floating point, based on IEEE Standard 754 *deprecated*

Single precision: 32 bits

*note
on floats*

S	exp	frac
1 bit	8 bits	23 bits

Double precision: 64 bits

S	exp	frac
1 bit	11 bits	52 bits

sign 符号位 exp: 指数 fraction: 小数
exponent fraction

Double precision 为 13 位
value = $(-1)^{\text{sign}} \times \text{fraction} \times 2^{\text{exponent}}$

$$\begin{aligned} \text{exponent} &= C - 1023 && \text{偏移量} \\ &\quad \uparrow \text{exponent} \\ &\quad \uparrow \text{bias} \end{aligned}$$

$$\text{fraction} = 1 + m && \text{尾数} \\ &\quad \uparrow \text{mantissa}$$

单精度浮点数偏移量为 127 bias = $2^{k-1} - 1$
 $k = \# \text{ of exp bits}$

$$C = \sum_{i=0}^{10} c_i 2^i$$

62 位到 52 位存的是 c , 共 11 位, 以表示 exponent, range [0, 2047]

所以 $C - 1023$ 的 range 为 $[-1023, 1024]$, 不考虑 $C=0$ 和 $C=2047$ 的情况 (特殊含义)

$C - 1023$ 的 range 为 $[-1022, 1023]$, 指数 exponent 的 range 为 $[2^{-1022}, 2^{1023}]$

$$m = \sum_{i=1}^{52} m_i 2^{-i}$$

51 位到 0 位存的是 m , 共 52 位, 以表示 fraction, range $[0, 1 - 2^{-52}]$

所以小数 fraction 的 range 为 $[1, 1 - 2^{-52}]$

当 $C=2046, m=1 - 2^{-52}$ 时, 双精度浮点数达到可以表示的最大正数 $\approx 2 \times 10^{308}$

当 $C=1, m=0$ 时, 最小正数 $\approx 2 \times 10^{-308}$

最小负数 $\approx -2 \times 10^{-308}$

最大负数 $\approx -2 \times 10^{-308}$

$C=0 \left\{ \begin{array}{l} m=0 \Rightarrow \pm 0 \\ m \neq 0 \text{ denormalized num} \end{array} \right.$

$C=2047 \left\{ \begin{array}{l} m=0 \Rightarrow \pm \infty \\ m \neq 0 \Rightarrow NaN \end{array} \right.$

IEEE floating point, based on IEEE Standard 754

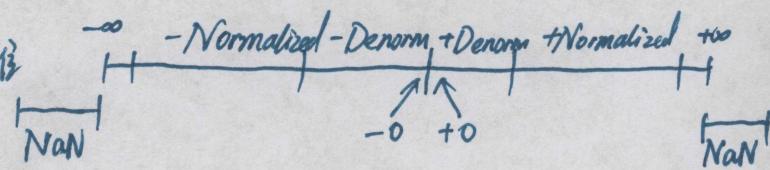
Single precision: 32 bits 有效位数: 7位

S	exp	frac
1 bit	8 bits	23 bits

Visualization: Float Point Encodings

Double precision: 64 bits 有效位数: 16位

S	exp	frac
1 bit	11 bits	52 bits



S: sign bit 符号位 exp: exponent 指数 frac: fraction 小数

$$\text{value} = (-1)^{\text{sign}} \times \text{fraction} \times 2^E$$

$$E = \begin{cases} \text{exp} - \text{bias} & \text{normalized} \\ 1 - \text{bias} & \text{denormalized} \end{cases}$$

bias = $2^{k-1} - 1$, k: # of exp bits
偏移量

$$\text{fraction} = \begin{cases} 1+m & \text{normalized} \\ 0+m & \text{denormalized} \end{cases}$$

mantissa 尾数 always have implied 1
no implied 1

example: double precision normalized values

63	62	52 51

63: sign bit

62-52: stores exp to encode E, 存储 exp 以表示 E ; $\text{exp} = \sum_{i=0}^{10} \text{exp}_i 2^i \in [0, 2047]$

去掉 exp=0 和 exp=2047 的情况 (denormalized values), $E = \text{exp} - \text{bias} \in [-1023, 1023]$

$$2^E \in [2^{-1023}, 2^{1023}]$$

51-0: stores frac to encode fraction, 存储 frac 以表示 fraction ; $\text{frac} = \sum_{i=1}^{52} \text{frac}_i 2^{-i} \in [0, 1 - 2^{-52}]$

when exp = 2046 ($E = 1023$), frac = $1 - 2^{-52}$ 时, 可表示最大正数 $\approx 1.8 \times 10^{38}$

when exp = 1 ($E = -1022$), frac = 0 时 可表示最小正数 $\approx 2.2 \times 10^{-308}$

2's complement: overflow in modulo way

denormalized floating point: overflow to infinity

denormalized: equal-spaced

normalized: once exp increment by 1, the spacing doubles
numbers with same exp field are equally spaced

Rounding: Nearest Even is the default method

less than a half: round down

more than a half: round up

in the middle: round to nearest even { if truncate is odd, round up and propagate
if truncate is even, round down
that is, just truncate

Rounding Binary Fractional Numbers:

if bits to rounding position start with 0 : round down

start with 1 { 100...0 $\frac{1}{2}$ > 0 halfway
1..0..1.. $\frac{1}{2}$ < 1 round up

floating point puzzles

int x = ...; $x == (\text{int})(\text{float})x$: false, we don't have enough frac bits in float to represent int, so we lose some bits by rounding

$x == (\text{int})(\text{double})x$: true, we have enough frac bits in double to represent int

float f = ...; $f == (\text{float})(\text{double})f$: true
neither d nor f is NaN

double d = ...; $d == (\text{double})(\text{float})d$: false $f = -(-f)$: true just toggling a bit

$2/3 == 2/3.0$: false

if $d < 0.0$ $d + 2 < 0.0$: true (even $2+d$ causes overflow, because overflow goes to negative infinity)

if $d > f$ $-f > -d$: true monotony

$d * d \geq 0.0$: false true (overflow goes to positive infinity)

$(d+f) - d == f$: false not associative