

Floating point correction

1-bit
Sign

11-bits
Exponent

52 bits
mantissa

Sign bit :- 0 or 1

0 - Positive
1 - negative

Exponent :- 11 bits for magnitude and sign are stored using a biased representation, with bias 1023. I.e

True Binary exponent

= Stored exponent - 1023

Mantissa :- 52 bits

Convert number to binary expansion; with digit next to decimal point to be 1 in the normalised representation. This gives 1-extra bit for free

Example :- 15. 8125

$$\left. \begin{array}{l} 15 = 2 \times 7 + 1 \\ 7 = 2 \times 3 + 1 \\ 3 = 2 \times 1 + 1 \\ 1 = 2 \times 0 + 1 \end{array} \right\} 15 = 1 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 = 111$$

$$\left. \begin{array}{l} 0.8125 \times 2 = 1. + 0.8125 \\ 0.625 \times 2 = 1 + 0.25 \\ 0.25 \times 2 = 0 + 0.5 \\ 0.5 \times 2 = 1 + 0 \end{array} \right\} 0.8125 = 1 \cdot 2^{-1} + 1 \cdot 2^{-2} + 0 \cdot 2^{-3} + 1 \cdot 2^{-4} = .1101$$

$$15.8125 = 1111.1101$$
$$= (1.111101 \times 2^3)$$

normalised

always 1

$$\text{Sign bit} \equiv 0$$

$$\text{stored exponent} = 3 + 1023 = 1026$$
$$= 100\ 0000\ 0010$$

$$\text{normalised mantissa} = 1111101 \underbrace{\text{45 zeros}}$$

$$15.8125 =$$

$$0100\ 0000\ 0010\ 1111\ 0100\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000$$

Other Conventions :-

Exponent 0 0 0 0 0 0 0 0 0 0 - reserved
 (! ! ! ! ! ! ! !)

Smallest positive number

$$0 0 0 0 0 0 0 0 0 1 \underbrace{0}_{s_2} 0 \quad | - 1023$$

$$= \left(1 + \sum_{j=1}^{s_2} 0 \cdot 2^j \right) \times 2$$

$$\approx 2.2 \times 10^{-308}$$

Largest positive number

$$0 \underbrace{(1 1 1 1 1 1 1 1 0)}_{s_2} 1 \quad | \quad (2046 - 1023)$$

$$= \left(1 + \sum_{j=1}^{s_2} 2^j \right) \times 2$$

$$\approx 1.8 \times 10^{308}$$

Special numbers

$$0 0 0 0 0 0 0 0 0 0 \underbrace{0}_{s_2} 0 \quad +0$$

$$1 0 0 0 0 0 0 0 0 0 \underbrace{0}_{s_2} 0 \quad -0$$

0 (1 1 1 1 1 1 1 1 0 s₂) +.inf

1 (1 1 1 1 1 1 1 1 0 s₂) -inf

0 or (1 1 1 1 1 1 1 1 0 s₁) NAN

Denormal numbers :- "Smaller than smallest"

- sacrifice on s₂-bit precision and
- enlarge the exponent storage
- thus going below .Machine \$double.xmin