

Fixed-point representation.

Example: 8 bits to store a real number.

$$(10111.011)_2$$

$$\underbrace{2^4 \ 2^3 \ 2^2 \ 2^1 \ 2^0}_{\text{integer part}} \quad \underbrace{2^{-1} \ 2^{-2} \ 2^{-3}}_{\text{fractional part}}$$

integer part fractional part.

More bits for the integer part \rightarrow increase range

fractional part \rightarrow precision

Scientific notation: $x = \pm r \times 10^m$, $r \in [1, 10]$.

floating number in the binary system:

$$x = \pm q \times 2^m$$

q is significand, normally a fractional value
in the range $[1.0, 2.0)$

m is the exponent.

Numerical Form: $x = \pm q \times 2^m = \pm \overbrace{b_0}^{\text{leading bit}} \underbrace{b_1 b_2 \dots b_n}_{\text{fractional}} \times 2^m$

$$b_i \in \{0, 1\}$$

Exponent range: $m \in [L, U]$.

Precision: $p = n + 1$.

$$0.40625 = 0.03125 \overset{2^{-5}}{x_1} + 0.0625 \overset{2^{-4}}{x_2} + 0.125 \overset{2^{-3}}{x_3}$$

$$x_1 = 3 \quad x_2 = 0 \quad x_3 = 1$$

$$(0.40625)_{10} = 2^{-2} + 2^{-3} + 2^{-5} \\ = (0.01101)_2$$

Normalized floating-point numbers.

FP number = $\langle \text{sign} \rangle 1. \langle \text{fraction field} \rangle \times 2^{(\langle \text{exponent field} \rangle - \text{bias})}$

$$x = \pm 1. b_1 b_2 b_3 \dots b_n \times 2^m = \pm 1. f \times 2^m$$

$$b_i \in \{0, 1\}$$

Exponent range: $m \in [L, U]$

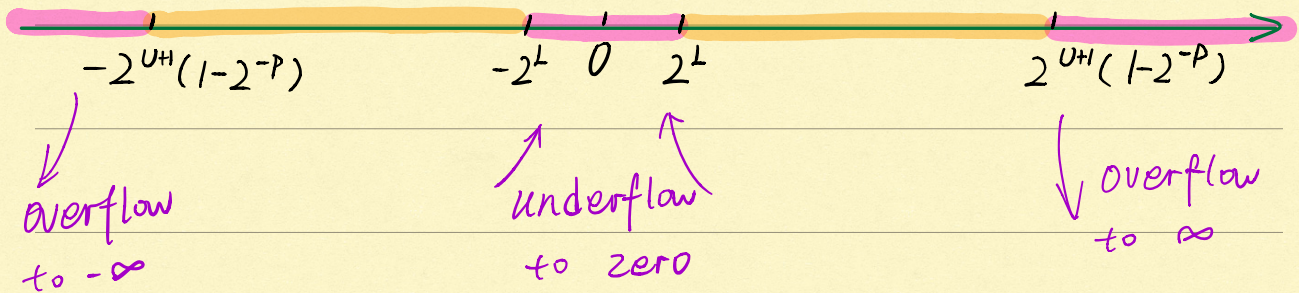
Precision: $p = n + 1$.

Smallest positive normalized FP number:

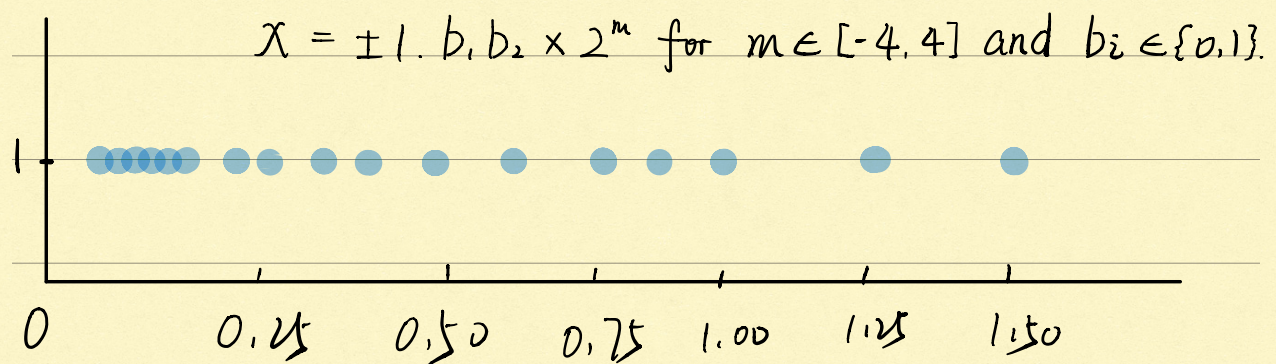
$$1.\underbrace{00\dots0}_n \times 2^L = 2^L \xrightarrow{\text{depends on}} \text{exponent}$$

Largest positive normalized FP number:

$$1.\underbrace{111\dots1}_n \times 2^U = 2^{U+1} (1 - 2^{-P}) \xrightarrow{\text{depends on}} \text{exponent precision.}$$



Machine Epsilon has relative precision. is defined as the distance (gap) between 1 and the next larger floating point number. $\epsilon_m = 1.25 - 1 = 0.25$



→ harder to be represented by x .

$$(1.00)_2 = 1.00 \times 2^0 = (1)_{10}$$

$$(10)_2 = 1.00 \times 2^1 = (2)_{10}$$

$$(11)_2 = 1.1 \times 2^1 = (3)_{10}$$

$$E_m = \underbrace{1.00 \dots 1}_n \times 2^0 - \underbrace{1.00 \dots 0}_n \times 2^0$$

$$= \underbrace{0.00 \dots 1}_n \times 2^0 = 2^{-n} \cdot 2^0 = 2^{-n}.$$

Subnormal floating point number:
(denormal)

bigger than 0, smaller than $1 \times 2^{\min.}$

$$0.1 \times 2^{\min.} = 2^{\min.-1}$$

$$0.01 \times 2^{\min.} \text{ --- } \text{---}$$

$$2^{\min.-n} \text{ (smallest).}$$