

Fixed Point Numbers

- keeps the decimal point in fixed position
- faster, simpler calculations, but less precise
- no. are equally spaced.

Example

ample
 $(10.1)_2$, $-(10.1)_2$ etc. $(12.345)_{10}$ ↖ this point is fixed

fixed point representation : $x = \pm (d_1 d_2 \dots d_{k-1} \cdot d_k \dots d_n)_\beta$
where, $d_1 \dots d_n \in \{0, 1, \dots, \beta-1\}$

Converting $(10.1)_2$ to $(?)_{10}$

$$= 1 \times 2^1 + 0 \times 2^0 + 1 \times 2^{-1}$$
$$= (2.5)_{10}$$

Example → 64-bit integers

(signed) largest possible no. = $2^{63} - 1$
 \uparrow
 (MSB) 1 bit for sign $\rightarrow \begin{cases} 0 (+ve) \\ 1 (-ve) \end{cases}$ Range: -2^{63} to $2^{63} - 1$
 $= 2^{n-1} - 1$

Floating Point Numbers

- allows the decimal point to move
- more precision, wider range of values, more complex
- no. are not equally spaced

Both fixed point and ~~flat~~ floating points are how numbers are stored/represented in a computer.

Floating Point Representation

$$F = \underset{\substack{\uparrow \\ \text{sign}}}{\pm} (0 \cdot d_1 d_2 \dots d_n)_{\underset{\substack{\uparrow \\ \text{base}}}{\beta}} \times \beta^{\underset{\substack{\uparrow \\ \text{exponent}}}{e}}$$

significand/fraction/mantissa

where, $\beta, d_i, e \in \mathbb{Z}$ (set of integers)

$$0 \leq d_i \leq \beta - 1$$

$$e_{\min} \leq e \leq e_{\max}$$

Examples

$$\begin{aligned} 123.45 \\ &= 12.345 \times 10^1 \\ &= 1.2345 \times 10^2 \\ &= 0.12345 \times 10^3 \end{aligned}$$

$$\begin{aligned} 1001.11 \\ &= 0.100111 \times 2^4 \end{aligned}$$

base

Conventions

Convention 1 : $\pm (0 \cdot \boxed{d_1} d_2 \dots d_m)_{\beta} \times \beta^e$

$$d_1 = 1 \text{ (always)}$$

① Example: $\beta = 2, m = 3, e = e$

$m = 3$ means, after decimal point, there will be 3 digits.

$$\begin{aligned} \text{So, highest value} &= (0 \cdot d_1 d_2 d_3)_2 \times 2^e \\ &= (0 \cdot 1 d_2 d_3)_2 \times 2^e \text{ [by convention, } d_1 = 1 \text{ always]} \\ &= (0 \cdot 1 \underline{1} \underline{1})_2 \times 2^e \text{ (Ans)} \end{aligned}$$

② Example: $\beta = 2, m = 3, e_{\min} = -1, e_{\max} = 2$

$$\text{max. value} = +(0 \cdot 111)_2 \times 2^{2 \Rightarrow \max} = \boxed{\frac{7}{4}}$$

$$\text{min. value (non-negative)} = +(0 \cdot 100)_2 \times 2^{-1} = \boxed{\frac{1}{4}}$$

$$\text{min. value} = -(\text{max. value}) = \boxed{-\frac{7}{4}}$$

Calculation

$$(0.111)_2 \times 2^2$$

$$= (2^0 \times 0 + 2^{-1} \times 1 + 2^{-2} \times 1 + 2^{-3} \times 1) \times 2^2$$

$$= \left(\frac{7}{2}\right)_{10}$$

Convention : 2 \Rightarrow Normalized form

$$\pm (1.d_1 d_2 \dots d_m)_\beta \times \beta^e$$

Example: $\beta = 2, m = 3, e = [-1, 2]$

\Downarrow

$$-1, 0, 1, 2$$

Highest value = $(1.111)_2 \times 2^2$

$$= (1 \times 2^0 + 1 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3}) \times 2^2$$

$$= (7.5)_{10}$$

Note:

if $e = (0, 2]$

\Downarrow

1, 2

Convention : 3 \Rightarrow Denormalized form

$$\pm (0.1 d_1 d_2 \dots d_m)_\beta \times \beta^e$$

Example: $\beta = 2, m = 3, e = [-1, 2]$

Highest value = $(0.1 d_1 d_2 d_3)_2 \times 2^2$

$$= (0.1111)_2 \times 2^2$$

$$= (3.75)_{10}$$

Total combinations

Convention: 1

$$\beta = 2, m = 3, e = [-1, 2] = \{-1, 0, 1, 2\}$$

$\swarrow \quad \searrow$
 $e_{\min} \quad e_{\max}$

$$(0.1 \underline{d_2} \underline{d_3})_2 \times 2^{-1}$$

$$\begin{array}{cc} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{array} \left. \vphantom{\begin{array}{cc} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{array}} \right\} \begin{array}{l} 4 \text{ combinations.} \\ \text{for } e = -1 \end{array}$$

Same goes for $e = 0, e = 1, e = 2$

$\Downarrow \quad \Downarrow \quad \Downarrow$
 $4 \quad 4 \quad 4$

total = $4 + 4 + 4 + 4 = 16$ combinations / floating numbers that can be represented.

Convention: 2

$$1 \cdot \frac{d_1}{0/1} \frac{d_2}{0/1} \frac{d_3}{0/1} = 2^3 = 8 \text{ combinations. for each exponent.}$$

$$\text{So, total} = 8 \times 4 = 32$$

Convention: 3

$$0.1 \underline{d_1} \underline{d_2} \underline{d_3}$$

$$\text{total} = 8 \times 4 = 32$$

Floating Point numbers are not equally spaced

$$e = [-1, 1] \quad \beta = 2, m = 3 \rightarrow \text{convention 1}$$

$$e = -1$$

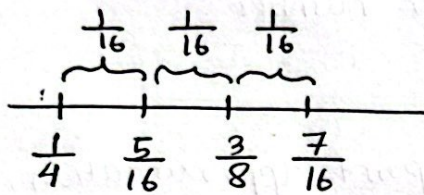
Smallest non-negative number combinations

$$1^{\text{st}} \text{ smallest} = (0.100) \times 2^{-1} = \frac{1}{4}$$

$$2^{\text{nd}} \quad \quad = (0.101) \times 2^{-1} = \frac{5}{16}$$

$$3^{\text{rd}} \quad \quad = (0.110) \times 2^{-1} = \frac{3}{8}$$

$$4^{\text{th}} \quad \quad = (0.111) \times 2^{-1} = \frac{7}{16}$$



□ when exponent is constant, numbers are equally spaced.

$$e = 0$$

$$(0.100) \times 2^0 = \frac{1}{2} \quad] \text{ diff.} = \frac{1}{8}$$

$$(0.101) \times 2^0 = \frac{5}{8} \quad] \text{ diff.} = \frac{1}{8}$$

$$(0.110) \times 2^0 = \frac{3}{4} \quad] \text{ diff.} = \frac{1}{8}$$

$$(0.111) \times 2^0 = \frac{7}{8} \quad] \text{ diff.} = \frac{1}{8}$$

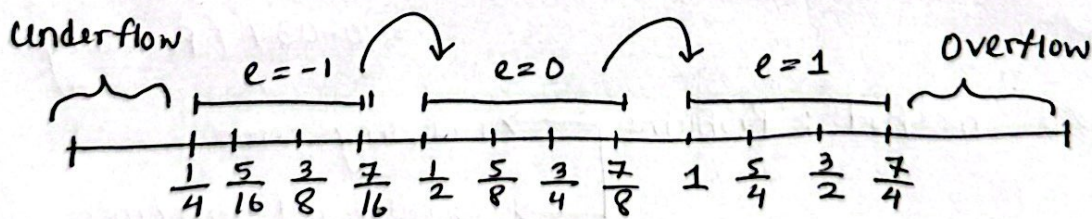
$$e = 1$$

$$(0.100) \times 2^1 = 1 \quad] \text{ diff.} = \frac{1}{4}$$

$$(0.101) \times 2^1 = \frac{5}{4} \quad] \text{ diff.} = \frac{1}{4}$$

$$(0.110) \times 2^1 = \frac{3}{2} \quad] \text{ diff.} = \frac{1}{4}$$

$$(0.111) \times 2^1 = \frac{7}{4} \quad] \text{ diff.} = \frac{1}{4}$$



□ with change in e , the gap/interval between numbers changes.

For convention 1 & denormalised form, we don't have zero in our number system.

Others

$$\begin{cases} x + 2y = 0 \\ 2x - \pi y = 1 \end{cases} \text{ we can solve this simultaneous equation easily (direct method)}$$

Problem: $\pi \rightarrow$ transcendental number

= computer has finite memory

= hence, it uses a floating-point approximation, which incurs rounding error.

$$\begin{aligned} \text{Example: } \pi &= 3.14159 \dots \\ &\downarrow \\ \pi &= 3.142 \end{aligned} \quad \left. \vphantom{\begin{aligned} \pi &= 3.14159 \dots \\ \pi &= 3.142 \end{aligned}} \right\} \text{floating point approximation.}$$

$$\begin{aligned} \text{Example: } 2.5763 \\ 2.58 \end{aligned} \quad \left. \vphantom{\begin{aligned} 2.5763 \\ 2.58 \end{aligned}} \right\} \text{FP approximation}$$

$$\begin{aligned} \text{difference between the numbers} &= 2.58 - 2.5763 \\ &= 0.0037 \text{ (Rounding error)} \end{aligned}$$

direct method : Problems $\begin{cases} \rightarrow \text{rounding error} \\ \rightarrow \text{computational time/memory} \end{cases}$

Example : $\sin(1.2)$ [Iterative method]

\hookrightarrow adding extra terms to improve approximation.

$$\sin(1.2) = 1.2 - \frac{(1.2)^3}{3!} + \frac{(1.2)^5}{5!} - \frac{(1.2)^7}{7!} + \dots \text{ (infinite series)}$$

Even if a computer could do real arithmetic, there would be an error due to stopping the iterative process at some finite point. This is called truncation error

$$\sin(1.2) = \underbrace{1.2 - \frac{(1.2)^3}{3!} + \frac{(1.2)^5}{5!}}_{\text{approximate}} - \underbrace{\frac{(1.2)^7}{7!} + \dots}_{\text{truncation error}}$$

□ for limited resources, we are taking first 3 terms, eliminating the rest. The error caused by eliminating those terms are known as truncation error. //

IEEE standard for double precision (1985)

64 bit Architecture / Arithmetic

$\beta = 2$, 52 bits \rightarrow fraction \rightarrow precision
 11 bits \rightarrow exponent \rightarrow range
 1 bit \rightarrow sign (0 means +ve, 1 means -ve)

normalised form = $(1.d_1d_2\dots d_{52}) \times 2^e$

sign	exponent	fraction
1 bit	11 bits	52 bits

For exponent, 2^{11} combinations = 2048 combinations.

possible values = $[0, 2047]$ $\therefore 2047 = 2^{11} - 1$

\Downarrow
total 2048 values

$e_{\min} = 0$
 $e_{\max} = 2047$ } However, this excludes negative values of e (ie. e^{-1} etc)

$$\text{largest possible non-negative no.} = (1.11\dots1)_2 \times 2^{2047}$$

$$\text{smallest possible non-negative no.} = (1.00\dots0)_2 \times 2^0$$

$= 1$ [not enough small,
cannot express
smaller values]

Therefore, we do **exponent biasing**

$$2048/2 = 1024$$

1024 values to left

1024 values to right.

$$\Rightarrow 1023 \mid 0 \mid 1024$$

2048 numbers.

exponent biasing
(done to represent smaller
numbers < 1)

$$\text{new, } e_{\min} = -1023$$

$$e_{\max} = 1024$$

$$\left. \begin{array}{l} e_{\min} = -1023 \\ e_{\max} = 1024 \end{array} \right\} \begin{array}{l} 2^{e-1023} \rightarrow \\ \downarrow \quad \downarrow \\ 0-1023 \quad 2047-1023 \\ = -1023 \quad = 1024 \end{array}$$

Now,

$$(1.d_1 d_2 \dots d_{52})_2 \times 2^{-1023}$$

$$= (0.1 d_1 d_2 \dots d_{52})_2 \times 2^{-1022}$$

$$(1.d_1 d_2 \dots d_{52})_2 \times 2^{1024}$$

$$= (0.1 d_1 d_2 \dots d_{52})_2 \times 2^{1025}$$

$$\text{Before, } e \in [0, 2047]$$

$$\text{After, } e \in [-1022, 1025]$$

Now, after exponent biasing:

$$\text{highest no.} = (0.11111 \dots 1)_2 \times 2^{1025} \approx \infty$$

$$\text{lowest no.} = (0.100 \dots 0)_2 \times 2^{-1022} \approx 0$$

Important

In IEEE standard, 2 bits from exponent is used/reserved to store ∞ and 0.

Now,

$$\begin{aligned} \text{Highest possible num (except infinity)} &= (0.11 \dots 1)_2 \times 2^{1024} \\ &\approx 1.798 \times 10^{308} \end{aligned}$$

$$\begin{aligned} \text{Lowest possible num (except 0)} &= (0.10 \dots 0)_2 \times 2^{-1021} \\ &\approx 2.225 \times 10^{-308} \end{aligned}$$

For exponent,

$$2^{1025} \rightarrow \pm \infty$$

$$2^{-1022} \rightarrow \pm 0$$

