

Floating Point Arithmetic

Fixed point Representation:

In our daily lives, we typically use fixed-point representation.

$$x = \pm (d_1 d_2 d_3 \dots d_{k-1} \cdot d_k \dots d_n)_{\beta}$$

$$\text{where } d_1 \dots d_n \in \{0, 1, \dots, \beta-1\}$$

Example:

$$x = + (11.1)_2$$

$$x = - (21.5)_{10}$$

Evaluating Fixed point numbers in base 10:

$$(11.1)_2 = 1 \times 2^1 + 1 \times 2^0 + 1 \times 2^{-1}$$

$$= 2 + 1 + \frac{1}{2}$$

$$= 3 + \frac{1}{2}$$

$$= \frac{7}{2}$$

$$= (3.5)_{10}$$

Floating Point Representation

Floating point numbers are the subset of all real numbers FCR . In computers, floating-point representation allows the system to handle a wide range of numbers, from tiny decimals to huge values - by breaking them down (organized) into a format that the computer can easily store and process.

$$F = \pm (0.d_1 d_2 d_3 \dots d_n)_{\beta} \times \beta^e$$

Significand / Fraction / mantissa

$e \rightarrow$ exponent
 $\beta \rightarrow$ Base

where, $\beta, d_i, e \in \mathbb{Z} \rightarrow$ integers

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

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

Examples: (Magic of Math)

[We have convert a given number into floating point representation]

$$\begin{aligned} * \quad & \pm (110.111)_2 \times 2^{\frac{-1}{2}} \\ &= \pm (11.0111)_2 \times 2^{\frac{1}{2}} \\ &= \pm (1.10111)_2 \times 2^{\frac{2}{2}} \\ &= \pm (0.110111)_2 \times 2^{\frac{3}{2}} \end{aligned} \quad \mid \quad \begin{aligned} * \quad & (1594.23)_{10} \times 10^{\frac{-1}{2}} \\ &= (159.423)_{10} \times 10^1 \\ &= (15.9423)_{10} \times 10^2 \\ &= (1.59423)_{10} \times 10^3 \\ &= (0.159423)_{10} \times 10^4 \end{aligned}$$

Conventions

Convention 1:

$$\pm (0.d_1 d_2 d_3 \dots d_m) \beta \times \beta^e$$

*** Hence $d_1 = 1$ and $d_1 \neq$ any (always) other value ***

d_1 always have to be 1.

Example: Given that, $\beta = 2$, $m = 3$, $e \in [-1, 2]$
 $e_{\min} \quad e_{\max}$

Highest possible Floating point number using convention:
is $(0.111)_2 \times 2^2$

Normalized Form : (Convention 2)

$$\pm (1.d_1d_2d_3 \dots d_m)_B \times B^e$$

* Hence d_1 can be 1 or any other value.

Example: Given that, $B=2$, $m=3$, $e \in [-1, 2]$
 \downarrow
 e_{\min} e_{\max}

Highest possible floating point number using Convention 2/
Normalized form is $(1.111)_2 \times 2^2$.

Denormalized Form : (Convention 3)

$$\pm (0.1d_1d_2d_3 \dots d_m)_B \times B^e$$

* Hence d_1 can be 1 or any other value.

Example: Given that, $B=2$, $m=3$, $e \in [-1, 2]$

Highest possible floating point number using
Convention 3/denormalized form is $(0.111)_2 \times 2^2$.

Total Combination

(How to find total Combination)

How many floating numbers can be possible

1. Given that, $\beta = 2$, $m = 3$, $e = [-1 \ 2]$

\downarrow
 e_{\min}

\downarrow
 e_{\max}

Using convention 1 : [where $d_1 = 1$ (always)]

For $e = -1$

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

0	0
0	1
1	0
1	1

[For $e = -1$, Hence 4 combinations possible]

$$\boxed{2^n = 2^2 = 4}$$

$e = [-1, 2]$ means " e " can be $-1, 0, 1, 2$.

For $e = 0$

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

0	0
0	1
1	0
1	1

[For $e = 0$, Hence 4 combinations possible]

Same thing goes for

$$\boxed{e = 1}, \boxed{e = 2}$$

4 combinations
possible

4 combinations
possible

\therefore total possible combinations $= 4 + 4 + 4 + 4 = 16$.

2. Given that, $B=2$, $m=3$, $e=[-1, 2]$

* Using normalized form / Convention 2

$$(1. \underline{d_1} \underline{d_2} \underline{d_3})_2 \times 2^{-1}$$

0 0 0

0 0 1

0 1 0

0 1 1

1 0 0

1 0 1

1 1 0

1 1 1

[For e = -1, Hence 8 combinations possible]

$$2^n = 2^3 = 8$$

Same thing goes for

$$e=0, e=1, e=2$$

8 combination

8 combination

8 combination

$$\therefore \text{Total possible Combinations} = 8 + 8 + 8 + 8 \\ = 32$$

3. Given that $B = 2$, $m = 3$, $e = [-1, 2]$

* Using normalized Form / convention 3

$$(0 \cdot 1 \underline{d_1} \underline{d_2} \underline{d_3})_2 \times 2^{-1}$$

0 0 0

0 0 1

0 1 0

0 1 1

1 0 0

1 0 1

1 1 0

1 1 1

[For $e = -1$, Hence 8 combinations]

possible

$$\boxed{2^3 = 8}$$

Same thing goes for $\underline{e=0}$, $\underline{e=1}$, $\underline{e=2}$

8 combinations

8 combinations

8 combinations

$$\begin{aligned}\therefore \text{Total possible Combinations} &= 8 + 8 + 8 + 8 \\ &= 4 \times 8 \\ &= 32\end{aligned}$$

* Practice Problem

[Find the smallest and largest values]

1. Given that $B=2$, $m=3$, $e=[-1, 2]$. Find the largest or maximum value or Highest Value.

[Using convention 1]

$$\begin{aligned}
 \text{Highest value} &= +(0 \cdot d_1 d_2 d_3)_2 \times 2^e \\
 &= +(0 \cdot 1 \underline{d_2} \underline{d_3})_2 \times 2^e \\
 &= +(0 \cdot \boxed{1} \ 1 \ 1)_2 \times 2^e \quad 2 \rightarrow e_{\max} \\
 &\quad d_1=1 \\
 &[\text{As we follow convention 1 here}] \\
 &= (+1/2)_{10}
 \end{aligned}$$

2. Given that $B=2$, $m=3$, $e=[-1, 2]$, Find the non-negative smallest or minimum value.

[Using Convention 1]

$$\text{Smallest value [non-negative]} = + (0 \cdot \boxed{1} 0 0)_2 \times 2^{\frac{-1}{2}} = \frac{1}{4} \quad \downarrow e_{\min}$$

3. Given that, $B=2$, $m=3$, $e=[-1, \frac{1}{2}]$ [as we follow convention]

Find the negative minimum or smallest value.

$$\begin{aligned}
 \text{Negative minimum or smallest value} &= - \text{Highest value} \\
 &= -(0 \cdot 1 1 1)_2 \times 2^{\frac{1}{2}} \\
 &= - \frac{7}{2}
 \end{aligned}$$

* * Important * * *

* IF it mentions only the smallest number, you need to include both the smallest non-negative and negative value.

Note

Binary to Decimal

$$\begin{aligned}1. \quad & + (0.\overset{-1}{1}\overset{-2}{1}\overset{-3}{1})_2 \times 2^2 \\& = (1 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3}) \times 2^2 \\& = \left(\frac{1}{2} + \frac{1}{2^2} + \frac{1}{2^3} \right) \times 2^2 \\& = \left(\frac{1}{2} + \frac{1}{4} + \frac{1}{8} \right) \times 2^2 \\& = \frac{4+2+1}{8} \times 4 \\& = \left(\frac{7}{2} \right)_{10} = 7/2\end{aligned}$$

$$\begin{aligned}2. \quad & (0.\overset{-1}{1}\overset{-2}{0}\overset{-3}{0})_2 \times 2^{-1} \\& = (1 \times 2^{-1} + 0 \times 2^{-2} + 0 \times 2^{-3}) \times \frac{1}{2} \\& = \frac{1}{2} \times \frac{1}{2} \\& = \frac{1}{4}\end{aligned}$$

Practice Problem

Given $\beta = 2, m = 3, e = [-1, 2]$. Find the highest value using Normalized Form.

$$\pm (1 \cdot d_1 d_2 d_3 \dots d_m) \beta \times B^e$$

$$\text{Highest value} = + (1 \cdot 1 1 1) \underbrace{2}_{} \times \underbrace{2^e}_{} \quad (e=2)$$

$$= (1 \cdot 1 1 1) \underbrace{2}_{} \times \underbrace{2^2}_{} \quad (e=2)$$

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

$$= \left(1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8}\right) \times 2^2$$

$$= \left(1 + \frac{4+2+1}{8}\right) \times 2^2$$

$$= \frac{15}{8} \times 4$$

$$= \left(\frac{15}{2}\right)_{10}$$

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

Given $\beta=2$, $m=3$, $e = [-1, 2]$. Find the highest value using denormalized form.

$$\pm (0.1d_1 d_2 d_3 \dots d_m) \beta \times \beta^e$$

$$\begin{aligned}
 \text{Highest value} &= (0.1111)_2 \times 2^2 \\
 &= (1 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3} + 1 \times 2^{-4}) \times 2^2 \\
 &= \left(\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} \right) \times 4 \\
 &= \frac{8+4+2+1}{16} \times 4 \\
 &= \frac{15}{4} \\
 &= \left(\frac{15}{4} \right)_{10} \\
 &= (3.75)_{10}
 \end{aligned}$$

Floating Point numbers are not equally spaced.

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

For convention 1

* Smallest non-negative number combinations.

For $e=-1$,

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

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

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

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

For $e=0$

$$(0.100) \times 2^0 = \frac{1}{2}$$

$$(0.101) \times 2^0 = \frac{5}{8}$$

$$(0.110) \times 2^0 = \frac{3}{4}$$

$$(0.111) \times 2^0 = \frac{7}{8}$$

For $e=1$,

$$(0.100) \times 2^1 = 1$$

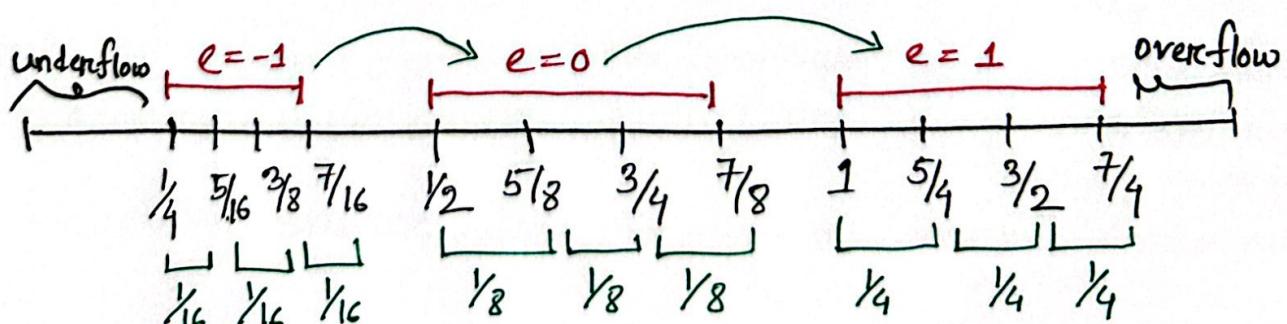
$$(0.101) \times 2^1 = \frac{5}{4}$$

$$(0.110) \times 2^1 = \frac{3}{2}$$

$$(0.111) \times 2^1 = \frac{7}{4}$$

*** important ***

When exponent e is constant, numbers are equally spaced. However, when exponent e changes, the interval or gap between numbers change.



IEEE Standard for double precision (64 bit)

$$\beta = 2$$

Fraction/mantissa = 52 bits

Exponent = 11 bits

Sign = 1 bits

Using Normalize Form

$$\pm (1.d_1 d_2 d_3 \dots d_{52})_2 \times 2^e$$

For $e = 11$ bits, there will be total $2^{11} = 2048$.

range of e $\boxed{[0, 2047]}$

\downarrow \downarrow
 e_{\min} e_{\max}

Largest possible number = $+ (1.111\dots1)_2 \times 2^{2047}$

Smallest possible number (non-negative) = $+ (1.000\dots0)_2 \times 2^0$
= 1 [1 is not enough]

the enough
small numbers
In this way
FP can't store
smaller values]

* To get a small number, the exponent value must be small.

* We obtain smaller values through exponent biasing.

Formula:

$$\begin{aligned} & 2^{k-1} - 1 \\ & = 2^{11-1} - 1 \\ & = 2^{10} - 1 \\ & = 1024 - 1 = 1023 \end{aligned}$$

$$\therefore (0.1d_1d_2 \dots d_{52})_2 \times 2^{e-1023} \xrightarrow{\text{exponent biasing}}$$

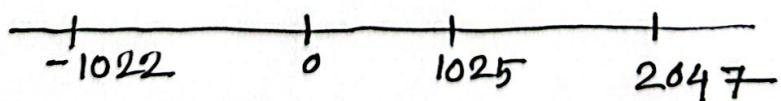
For denormalized form

$$(0.1d_1d_2 \dots d_{52})_2 \times 2^{-1022}$$

$$e-1022 \in [-1022, 1025]$$

$$\text{Largest possible number} = (0.11 \dots 1)_2 \times 2^{1025} \approx \infty$$

$$\text{Smallest possible number} = (0.10 \dots 0)_2 \times 2^{-1022} \approx 0$$



✳ Reduce the range for largest possible number so that we can get the smaller values.

$2^{\frac{1025}{2}} \rightarrow \pm \infty$ [highest power is used to store ∞]

$2^{-\frac{1022}{2}} \rightarrow \pm 0$ [used to store 0]

\therefore Highest value = $(0.111\dots1)_2 \times 2^{1024}$ [if we put 1025 it will denote ∞]
 $\approx 1.798 \times 10^{308}$

\therefore Smallest value = $(0.100\dots0)_2 \times 2^{-1022}$
 $\approx 2.225 \times 10^{-308}$

