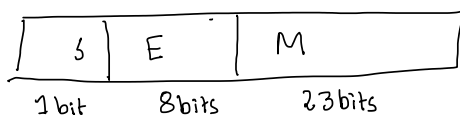1st Ques)  Floating Point Representation:

There are various architectures (ARM, INTEL) and each Architecture has

their own representation.

To avoid ambiguity, All these processor builders, (building Companies) should
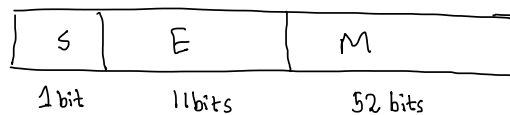
Follow Standerds.

One of the famous Standerds for Floating point Representation is IEEE 754.

SINGLE PRECISION:

- It is a 32-bit Representation

| S | E | M |
|---|---|---|
| 1 bit | 8 bits | 23 bits |

DOUBLE PRECISION:

- It is a 64-bit Representation

| S | E | M |
|---|---|---|
| 1 bit | 11 bits | 52 bits |

STANDERD CONVERSION:    ↙SP    ↙DP

BY Seeing this expression we Can understand that we Can even store Smaller numbers.

$$(-1)^S \times (1 \cdot M) \times B^{(E-127)} \text{ or } (E-1023)$$

where,

( to decimal )
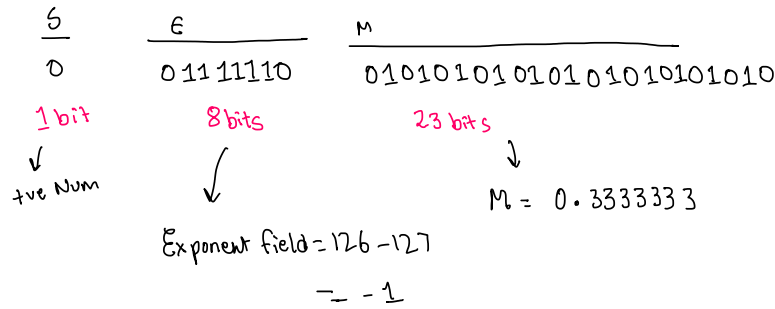
S → Sign

← B → Base   For Binary = 2

M → Mantissa

Mantissa Plays an important role in defining Precision of a Number in Floating Point Representation.

Eg:  $\frac{2}{3} = 0.66666666 \ldots$  ( here 0.6666..67 is more Precise than 0.66 )
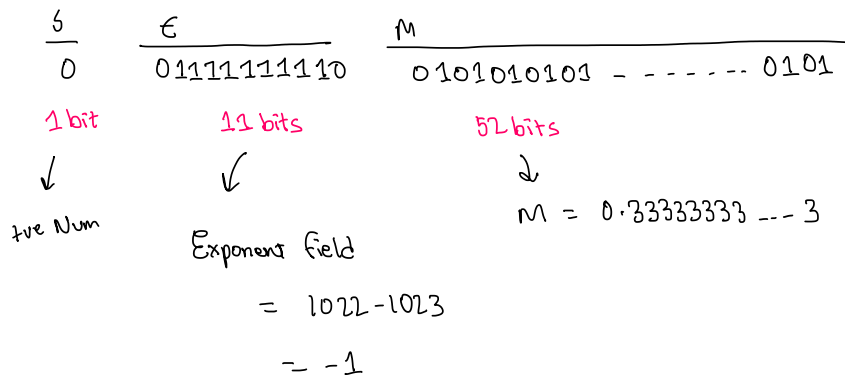
This above example Can have more precision if number of Mantissa bits are more. Single Precision has 23 bits fractional part and double Precision has 52 bits fractional parts. More number of fractional bits will result in avoiding loss of information of the number. So in double precision, precision is mattered more.
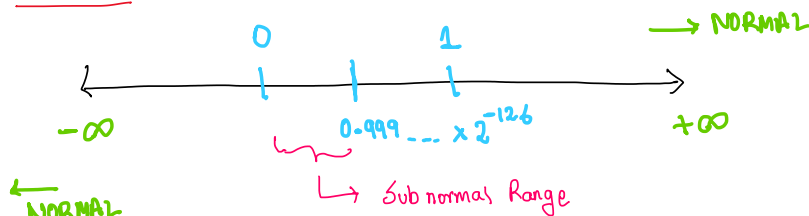
Eg: $\frac{2}{3} = 0.6666 --- 67$

SINGLE PRECISION:

| $\frac{S}{0}$ | $\frac{E}{01111110}$ | $\frac{M}{01010101010101010101010}$ |
|---|---|---|

1 bit → +ve Num

8 bits →

23 bits → $M = 0.3333333$

Exponent field = 126 - 127
= -1

DOUBLE PRECISION: (MORE PRECISION)

| $\frac{S}{0}$ | $\frac{E}{01111111110}$ | $\frac{M}{0101010101 -------- 0101}$ |
|---|---|---|

1 bit → +ve Num

11 bits → Exponent field
= 1022 - 1023
= -1

52 bits → $M = 0.33333333 --- 3$

2 Ques) NUMBER LINE:



NORMAL NUMBERS:

Format of Normal Numbers: $(-1)^S \times (1.M) \times 2^{E-127}$
Here the bit before the decimal point is always 1.

SUBNORMAL NUMBERS:

In Short, Numbers Smaller than NORMAL NUMBERS ARE
SUBNORMAL NUMBERS.

Format of Sub Normal Numbers: $(-1)^S \times (0.M) \times 2^{E-127}$

Here the bit before the decimal point is 0.

• From the Number line

SUBNORMAL NUMBER < Smallest Normal Number.

3 Ques) According to the IEEE 754 vv Standerds, The five rounding methods
Are:

- Round to nearest even
- Round to nearest away
- Round Up (Round to infinity)
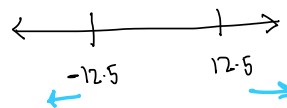- Round down (Round to -infinity)
- Round to Zero.

1. <u>Round to nearest even</u>: Rounded to the nearest possible value s.t
if the number is the middle, moves to the even least
Significant bit

Eg: $1.5 \rightarrow 1.0$, $-1.5 \rightarrow -2$

2. <u>Round to nearest away</u>: Rounded to the nearest possible value
s.t it is rounded away from zero. (nearest value above the num)
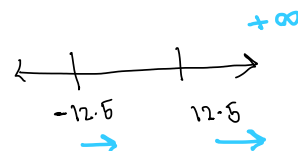
Eg: $12.5 \rightarrow 13$

$-12.5 \rightarrow -13$



3. <u>Round up</u>: Rounded to the number larger than itself i.e number
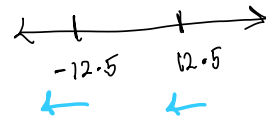is rounded towards $+\infty$.

Eg: $12.5 \rightarrow 13.0$

$-12.5 \rightarrow -12.0$



4. <u>Round down</u>: Rounded to the number smaller than itself
i.e number is rounded towards $-\infty$

Eg:    12.5 $\longrightarrow$ 12.0

      -12.5 $\longrightarrow$ -13.0

$-\infty$

$\longleftarrow$ | -12.5     | 12.5 $\longrightarrow$

5. <u>Round to zero</u>: The number is rounded to zero.

Eg:    12.5 $\longrightarrow$ 12.0

      -12.5 $\longrightarrow$ -12.0