

Floating Point Representation :

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

fraction/mantissa/significand

β = Base, e = exponent.

The fraction bits are commonly referred to as 'mantissa'.

Convention 1 (Standard/General Form)

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

$d_1 = 1$ always

Convention 2 (Normalized Form)

$$\pm (1.\underline{d_1} d_2 d_3 \dots d_m)_{\beta} \times \beta^e$$

d_1 can be either 0 or 1.

Convention 3 (Denormalized Form)

$$\pm (0.1 \underline{d_1} d_2 d_3 \dots d_m)_{\beta} \times \beta^e$$

d_1 can be either 1 or 0.

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

Questions : ① Find largest possible/highest number for Convention 1, 2 and 3?

Convention 01 : $(0.111)_2 \times 2^2$

Convention 02 : $(0 \rightarrow 1.111)_2 \times 2^2$
(Normalized)

Convention 03 : $(0-1111)_2 \times 2^2$
(Denormalized)

- (ii) Find smallest possible (non-negative) number for Convention 01, 02 and 03?

Convention 01 : $(0.100)_2 \times 2^{-1}$

Convention 02 : $(1.000)_2 \times 2^{-1}$
(Normalized)

Convention 03 : $(0.1000)_2 \times 2^{-1}$
(Denormalized)

- (iii) Find smallest possible number for convention 01, 02 and 03? [Always remember to consider sign bit].

Convention 01 : $-(0.111)_2 \times 2^{+2}$

Convention 02 : $-(1.111)_2 \times 2^{+2}$
(Normalized)

Convention 03 : $-(0.1111)_2 \times 2^{+2}$
(Denormalized)

* Trick is to find the highest number and give negative sign.

IV Possible combinations of numbers for the three conventions?

$$e = \begin{cases} -1 \\ 0 \\ 1 \\ 2 \end{cases}$$

Convention 01 : 16 possible combinations.

$$(0 \cdot 1 \quad \quad)_2 \times 2^e$$

$d_1 d_2 d_3 \rightarrow 4$ possible combinations

00
01
10
11

Convention 02 : 32 possible combinations.

$$(1 \cdot \quad \quad)_2 \times 2^e$$

$d_1 d_2 d_3 \rightarrow 8$ possible combinations

000
001
010
011
101
110
100

Convention 03 : 32 possible combinations

$$(0 \cdot 1 \quad \quad \quad)_2 \times 2^e$$

$d_1 d_2 d_3 \rightarrow 8$ possible combinations

111
↑

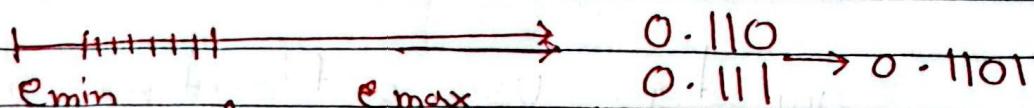
IEEE Standard (1985) for double precision
(64-bit arithmetic)

$\beta = 2$, 52 bits for the fraction \rightarrow precision

11 bits for the exponent \rightarrow (bigger range) $\rightarrow e_{\max}, e_{\min}$

1-bit for the sign

If we have more bits in our fractional part, then we will have more numbers nearby, so precision gets better.
in the number line



$(1 \cdot d_1 d_2 \dots d_{52})_2 \times 2^e$ We use binary numbers maximum value 1 in computer and 0 normalized form -
 $e = 2^{11} = 2048$ $[0, 2047]$ Form -
maximum value 2
 $e_{\min} = -1023$

largest possible number?

$$(1.11\dots1) \times 2^{2047}$$

smallest possible number?

$$(1.00\dots0) \times 2^{-1023}$$

We need smaller representation of number?

0.1

For this we need smaller values of e for which we need exponent biasing.

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

↳ exponent bias.

IEEE standard uses denormalized form

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

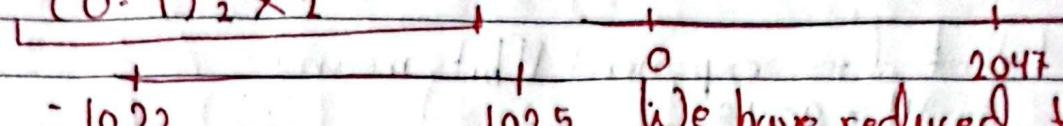
$e \in [0, 2047]$ - previous range
 b biasing $(e_i - 1022) \in [-1022, 1025]$

largest possible number

$$(0.11\ldots1)_2 \times 2^{1025}$$

smallest possible number

$$(0.1)_2 \times 2^{-1022}$$


We have reduced the range of largest possible number but we have incorporated very small range of numbers.

Exponent biasing is very important because we need numbers in the smaller range rather than in the larger range.

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

$$2^{-1022} \rightarrow +0 \rightarrow 1025 \text{ is reserved for infinity.}$$

$$\text{largest number: } (0.11\ldots1)_2 \times 2^{1024} \approx 1.798 \times 10^{308}$$

$$\text{smallest number: } (0.1)_2 \times 2^{-1022} \approx 2.225 \times 10^{-308}$$

underflow

$$\rightarrow 0$$

overflow

\rightarrow exception

underflow is mapped to 0.

overflow is exception.

Rounding Prior

When we round a number we need to draw a number line. We need to do rounding because computer can read upto certain bits or significant figures.

$m = 3$ (3 bits after the decimal place)
 (0.100)

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

$$\begin{array}{r} 5 \\ \times 2 \\ \hline \end{array}$$

(0.101),

$$= \frac{(1x_2^{-1} + 1x_2^{-2})}{(1 + 1 + 8 + 2)}$$

$$\frac{+1}{2} \quad \frac{+1}{8} \quad \frac{+1}{16}$$

$$= \frac{10}{16}, \quad \boxed{\begin{array}{|c|} \hline 5 \\ \hline 8 \\ \hline \end{array}}$$

$$\left(\frac{1}{2} + \frac{5}{8}\right) / 2 = \frac{(8 + (2 \times 5))}{16} \times 2$$

$$\frac{18}{16} \times 2 = \frac{18}{8} \rightarrow \boxed{\frac{9}{4}} = \frac{8}{16} + \frac{1}{16}$$

$$\begin{array}{r} 1 + 1 \\ \hline 2 \quad 2^4 \end{array}$$

2-1-10

12-42

If the last bit is 0 of a binary number then it is even and if the last bit is 1 of a binary number then it is odd.

Suppose a number is given $(0.1000100)_2$

H_2 

$$\rightarrow (0 \cdot 100), x (6 \cdot 100), y (0 \cdot 101).$$

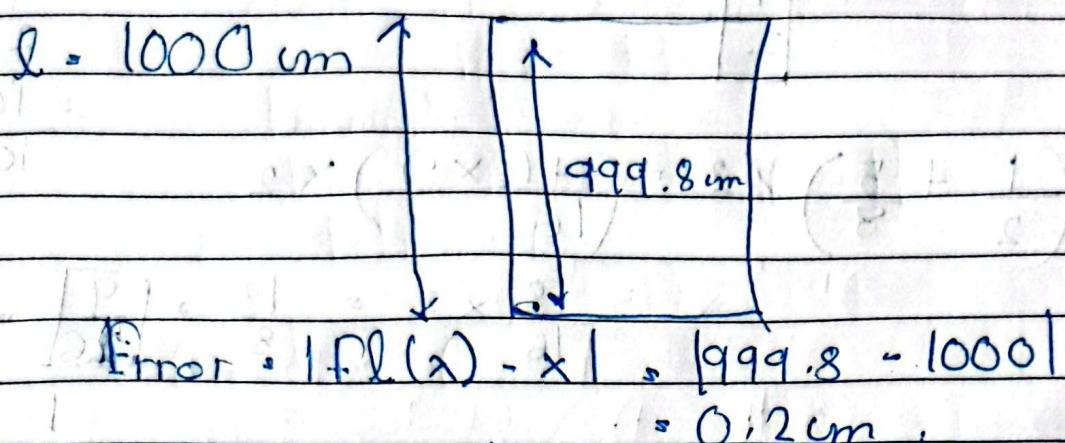
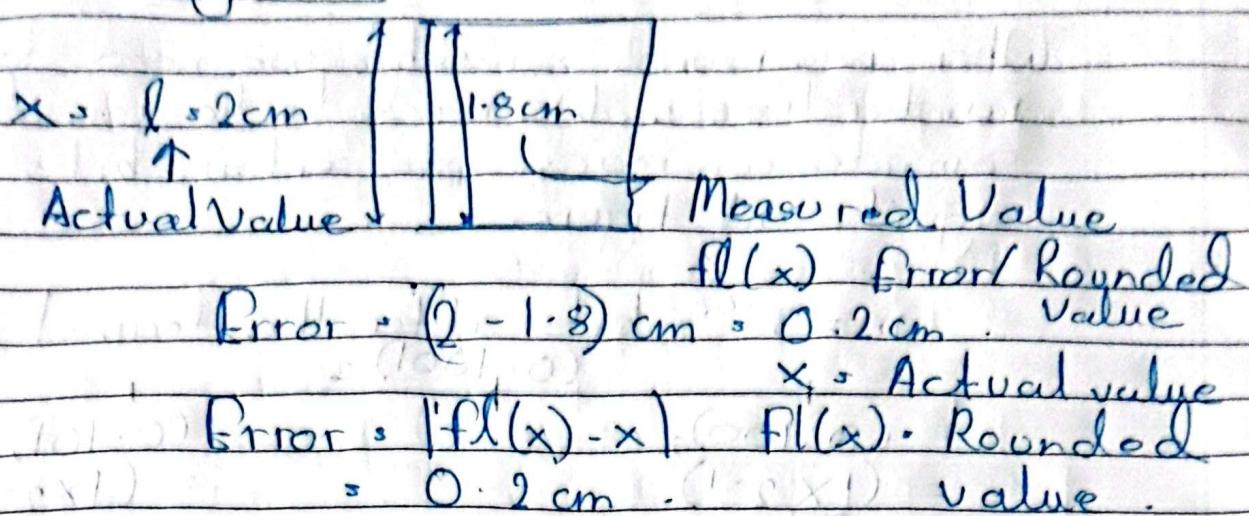
$$(0.1000100)_2 = (0.100)_2$$

$$(0 \cdot 1001010001), - (0 \cdot 101),$$

$$(0 \cdot \bar{1}00)_2 = (0 \cdot 100)_2$$

Since it is exactly the middle value it will get rounded to the nearest integer.

Rounding Error



It is a bit difficult to interpret the impact using only Error = $|f(l(x)) - x|$ so we need to do Error = $\frac{|f(l(x)) - x|}{|x|}$.

Scale invariant rounding error δ :

$$|f(l(x)) - x|$$

We deal with maximum scale invariant rounding error, Machine epsilon, ϵ

$$\cancel{\max} (0.01 \cdot 0) = 1(001001 \cdot 0)$$

$$(101 \cdot 0) = 1(001001001 \cdot 0)$$

$$1(001 \cdot 0) = 1(001 \cdot 0)$$

off by one digit places in floating point

s_{\max}/e

Scale invariant Rounding Error is $|f(x)| \times |s_{\max}|$

Convention 01: $(0.d_1d_2d_3)_{\beta} \times \beta^e$

$$(0.100)_2 \times 2^e \quad (0.101)_2 \times 2^e$$

$$d = (0.001)_2 \times 2^e$$

$$m=3, \quad (0.001)_2 \times 2^e$$

$$(1 \times 2^{-3}) \times 2^e.$$

$$\beta^{-m} \beta^e \quad (1 \times 2^{-m}) \times 2^e.$$

$$|f(x)| \times |s_{\max}| = \frac{1}{2} \beta^{-m} \beta^e$$

Convention (1)

$$(0.100)_2 \times 2^e$$

$$1 \times 2^{-1} \beta^e$$

$$[1 \times 1_m = \beta^{-1} \beta^e]$$

$$\epsilon = \frac{1}{2} \beta^{-m} \beta^e = \frac{1}{2} \beta^{1-m}$$

Two types of questions: $m=?$, $\beta=?$
(Given)

Calculate Machine Epsilon value for
mantissa convention 11

$$\text{Answer: } \epsilon = \frac{1}{2} \beta^{1-m}$$

Minimum value of x for convention 01?

$$1 \times 1_{\min} = \beta^{-1} \beta^e$$

There is no exponent e in Machine Epsilon

So the value of machine epsilon won't be affected by the value of exponent for convention 1.

Convention 02 (Normalized Form) $m=3$

$$|f(x) - x|_{\max} = \frac{1}{2} \beta^{-m} \beta e$$

$$|x|_{\min} = (1.000)_2 \times 2^e$$

$$|x|_{\min} = (1.0)_2 \times 2^e$$

$$= (1 \times 2^0) \times 2^e$$

$$|x|_{\min} = 2^e \times 2^0$$

$$|x|_{\min} = \beta^0 \beta e$$

$$|x|_{\min} = \beta e$$

$$\boxed{\beta^0 \cdot 1}$$

$$\epsilon = \frac{1}{2} \beta^{-m} (\beta e)$$

$$(\beta e)^{1-2} = |x|$$

$$\boxed{\epsilon = \frac{1}{2} \beta^{-m}}$$

Convention 03 (De-normalized Form)

$$(0.1d_1d_2d_3)_\beta \times \beta e$$

An extra bit 1 is added in this

$$(0.1d_1d_2d_3)_2 \times 2^e$$

$$|f(x) - x|_{\max} = \frac{1}{2} \beta^{-m} \beta e$$

$$= \frac{1}{2} \beta^{m+1} \beta e$$

$$\frac{1}{2} \beta^{-(m+1)} \beta e$$

$$|x|_{\min} = \frac{1}{2} \beta^{-m-1} \beta e$$

$\xrightarrow{\quad \beta^{-1} \beta e \quad}$

$$\frac{1}{2} \beta^{-m}$$

ϵ

$$(0.1)_2 \times 2^e$$

$$1 \times 2^{-1} \times 2^e$$

$$\beta^{-1} \times \beta e$$

Question: Base and mantissa will be given, calculate Machine Epsilon.

Part-5 $\beta = 2, m = 3$, Calculate the value of $f(x * y)$.
We need to round the output.

$$x = \frac{5}{8}, y = \frac{7}{8}$$

① Convert to binary

$$\frac{7}{8} = \frac{6}{8} + \frac{1}{8}$$

$$\frac{5}{8} = \frac{4}{8} + \frac{1}{8}$$

$$= \frac{1}{2} + \frac{1}{2^3} \quad (1/2 + 1/8) \times \frac{1}{2} + \frac{1}{4} + \frac{1}{8}$$

$$= 2^{-1} + 2^{-3}$$

$$= (0.101)_2$$

$$= (0.111)_2$$

$$\text{actual } (x * y) = (0.101)_2 * (0.111)_2$$

$$\text{value} = \frac{5}{8} \times \frac{7}{8} = \frac{35}{64} = \frac{32}{64} + \frac{2}{64} + \frac{1}{64}$$

We need to round the value since $m=3$.

$$\text{value} = 2^{-1} + 2^{-3} + 2^{-5} + 2^{-6}$$

$$\Rightarrow 0.100011_2$$

$$0.1001_2 = 2^{-1} + 2^{-5} + 2^{-6}$$

$$(0.101)_2 = (0.100011)_2$$

Value In the left range $(0.100011)_2 = (0.100)_2$

$$\text{rounded value} = \frac{1}{2}$$

$\frac{3.5}{6.4} \neq \frac{1}{2}$. rounding errors

Loss of significance x will be rounded to $f_l(x)$
 y " " " " " $f_l(y)$

$$f_l(x) = x + \delta_1 x$$

\uparrow \uparrow

$f_l(x) \neq x$ (may not be equal)

underrounded actual rounding error

$f_l(y) \neq y$ (n)

value value is being incorporated.

$$f_l(y) = y + \delta_2 y$$

$x \neq y$

$$f_l(x) + f_l(y)$$

$$x - y$$

$$= x + \delta_1 x + y + \delta_2 y$$

$$\cancel{x\delta_1 - y\delta_2}$$

$$x - y \leftarrow \text{approximately } x(1 + \delta_1) + y(1 + \delta_2)$$

$$x - y \leftarrow (x + y)(1 + \delta_1 + \delta_2)$$

rounded value will
be at infinity

Scale invariant error becomes scale invariant
infinitely large when we subtract error.

two numbers. This phenomena is loss of significance.
When I subtract two numbers, and both the numbers
are closer then in case of rounding or scale invariant
error this denominator will be very small or
approximately close to zero, the scale invariant
error will be very high.

Example-1 :

$$x^2 - 5.6x + 1 = 0 \quad \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$x_1 = 28 + \sqrt{783} = 55.98$$

$$x_2 = 28 - \sqrt{783} = 0.01786$$

\rightarrow 4sf - computer can only process.

$$\sqrt{783} = 27.98$$

$$x_1 = 28 + 27.98 = 55.98$$

$$x_2 = 28 - 27.98 = 0.02000$$

$0.01786 \neq 0.02000 \rightarrow$ loss of significance.

As x_2 numbers are very small, denominator will be very small so rounding error will be very high.

Solution: We ~~never~~ will not subtract the numbers.

$$x^2 - (\alpha + \beta)x + \alpha\beta$$

$$\alpha' = x_1$$

$$\beta' = \frac{1}{x_1} = x_2$$

$$5.01 \quad 5.02$$

$$\text{average} = \underline{[5.015]} \rightarrow \text{3sf}$$

$$\text{fl}(\frac{5.01 + 5.02}{2}) = \text{fl}(\frac{10.03}{2})$$

$$= \text{fl}(\frac{10.0}{2})$$

5. Rounding error