

Floating point Arithmetics

Fixed Point Representation:

Regular numbers (সাধারণ সংখ্যা) that we use.

Example: Decimal $\rightarrow (10.215)_{10}$, Binary $\rightarrow (110.11)_2$

Here the point is fixed.

General form:

$$x = \pm \left(\underbrace{d_1 d_2 d_3 \dots d_{k-1}}_{\text{whole part}} \cdot \underbrace{d_k d_{k+1} \dots d_n}_{\text{fraction part}} \right)_{\beta}$$

Point

where,

d_i = digit

β = Base

But our main focus is on,

Floating Point Representation:

* Floating point number, F is subset of all Real number, \mathbb{R} .

$$\therefore F \subset \mathbb{R}$$

It follows a certain format. For example

$$(10.215)_{10} = (0.10215)_{10} \times 10^2$$

General form,

$$F = \pm \left(\underbrace{0 \cdot d_1 d_2 d_3 \dots d_m}_{\text{fraction/mantissa/significand}} \right)_{\beta} \times \beta^e$$

exponent
Base

$\left. \begin{array}{l} \text{an integer.} \\ \text{has a range defined} \\ \text{as,} \\ e_{\min}, e_{\max} \end{array} \right\}$

For example,
 $(11.101)_2$ this is a fixed point number. To convert it into a floating point number, we have to move the decimal point & adjust the power accordingly.

$$\begin{aligned}
 & (11.101)_2 \times 2^0 \\
 &= (1.1101)_2 \times 2^1 \\
 &= (0.11101)_2 \times 2^2 \quad \leftarrow \text{this is following the general form.}
 \end{aligned}$$

There are 3 conventions of floating point that we will discuss in this course.

Convention 1: (General/Standard form)

$$\pm(0.d_1d_2\dots d_m)_\beta \times \beta^e$$

where $d_1 \neq 0$

For binary

$$\pm(0.d_1d_2\dots d_m)_\beta \times \beta^e$$

where $d_1 = 1$

Note,
why $d_1 = 1$ or $d_1 \neq 0$?

In codes, if $d_1 = 0$,
then, $\pm(0.000)$ it
can produce error in
code).

Also, this can lead
unnecessary bit usage as
 0.001×2^2 .

→ These 2 bits are
unnecessary.

Convention 2: IEEE Normalized form

For binary,

$$\pm(0.1d_1d_2d_3\dots d_m)_\beta \times \beta^e$$

Here, d_1 can be either 0 or 1.

Convention 3: IEEE Denormalized form

For binary,

$$\pm(1.d_1d_2d_3\dots d_m)_\beta + \beta^e$$

Here, d_1 can be either 0 or 1.

Note, in previous lecture
notes & videos, conv 2
and conv 3 were
swapped. Follow this
lecture note!!!

Examples : ① Given, $\beta = 2$, $m = 3$ and a exponent e . What is the highest value using convention 1, 2, 3?

Solⁿ:

$$\text{conv1} ; \pm (0 \cdot d_1 \ d_2 \ d_3)_2 \times 2^e \text{ where } d_1 = 1,$$

$$\therefore \max = (0 \cdot 1 \underline{1} \ \underline{1})_2 \times 2^e$$

$$\text{conv2} ; \pm (0 \cdot 1 \ d_1 \ d_2 \ d_3)_2 \times 2^e$$

$$\therefore \max = (0 \cdot 1 \underline{\underline{1}} \ \underline{\underline{1}})_2 \times 2^e$$

$$\text{conv3} ; \pm (1 \cdot d_1 \ d_2 \ d_3)_2 \times 2^e$$

$$\therefore \max = (\underline{1} \ \underline{\underline{1}} \ \underline{\underline{1}})_2 \times 2^e$$

② Given, $\beta = 2$, $m = 3$, $e = [-1, 2]$, using conv 1 and conv 3 how many combination is possible? (non negative)

Solⁿ:

$$\text{Here, } e_{\min} = -1 \quad \therefore e = \{-1, 0, 1, 2\}$$

$$e_{\max} = 2$$

We know,

$$\text{conv1} ; (0 \cdot d_1 \ d_2 \ d_3)_2 \times 2^e \text{ where } d_1 = 1,$$

For each value of e ,

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

$$\begin{matrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{matrix} \left. \right\} 4 \text{ combinations}$$

We get 4 combination.

\therefore for 4 values of $e (-1, 0, 1, 2)$ we get, $4 + 4 + 4 + 4 = 16$ combination.

in conv1, $d_1 = 1$.
It is fixed. So, we
only find out
the number of
combination possible
using d_2, d_3, \dots

$$\text{conv 2: } (1 \cdot \underline{d_1} \underline{d_2} \underline{d_3})_2 \times 2^e$$

0	00
0	01
:	
1	10
1	11

} total $2^3 = 8$ combinations for each e.

∴ For 4 values of e (-1, 0, 1, 2), we get total $4 * 8 = 32$ values.

③ Given, $B = 2$, $m = 3$, $e = [-1, 2]$, using conv 1 & conv 2 what is the smallest possible non negative value & highest value?

Solⁿ:

$$e = \{-1, 0, 1, 2\}$$

We know,

$$\text{Conv1: } (0 \cdot \underline{d_1} \underline{d_2} \underline{d_3})_2 \times 2^e \text{ where } d_1 = 1.$$

$$\therefore \text{smallest possible non negative value} = (0 \cdot \underline{1} 00)_2 \times 2^{-1}$$

$$\therefore \text{highest value} = (0 \cdot 111)_2 \times 2^2$$

in conv 1, $d_1=1$
and $2^{e_{\min}}$
makes the
smallest value.
and $2^{e_{\max}}$
makes the
largest value

We know,

$$\text{conv2: } (0 \cdot 1 \underline{d_1} \underline{d_2} \underline{d_3})_2 \times 2^e$$

$$\therefore \text{smallest possible non negative value} = (0 \cdot 1 00)_2 \times 2^{-1}$$

$$\therefore \text{highest value} = (0 \cdot 111)_2 \times 2^2$$

④ Given, $B = 2$, $m = 3$, $e = [-1, 2]$, Find the lowest value in conv 1.

Solⁿ:

$$\text{Conv1: } \pm (0 \cdot \underline{d_1} \underline{d_2} \underline{d_3})_2 \times 2^e \text{ where } d_1 = 1.$$

$$\text{lowest value} = -(0 \cdot 111)_2 \times 2^2$$

Sometimes, the questions may ask for the decimal representation. In that case, you have to convert the number into decimal.

Recap: how to convert bases of a number

Number Conversion:

Type 1: Decimal to other bases

Type 2: Other Base to Decimal

Type 1

Decimal to other Base:

$(43.3125)_{10}$

Whole number | Fractions

- * Whole number \rightarrow Repeated Division by Base R.
- * Fraction number \rightarrow Repeated Multiplication by Base R.

Example:

$$(43.3125)_{10} = (?)_2$$

$$\begin{array}{r} 43 \\ \times 2 \quad \text{MSB} \\ \hline 21 \\ \times 2 \quad \rightarrow 1 \\ \hline 10 \\ \times 2 \quad \rightarrow 1 \\ \hline 5 \\ \times 2 \quad \rightarrow 0 \\ \hline 2 \\ \times 2 \quad \rightarrow 1 \\ \hline 1 \\ \times 2 \quad \rightarrow 0 \\ \hline 0 \\ \times 2 \quad \rightarrow 1 \end{array}$$

We stop at 0.

$$\therefore (43)_{10} = (101011)_2$$

We stop when it reaches 0 or the limit will begin

$$\begin{array}{r}
 3125 \times 2 = 0.625 \\
 \times 2 = 1.25 \\
 \times 2 = 0.5 \\
 \times 2 = 1.0
 \end{array}$$

$\begin{array}{r} 0 \\ \times 2 \\ \hline 1 \end{array}$	Carry ^{LSB}
$\begin{array}{r} 1 \\ \times 2 \\ \hline 0 \end{array}$	<u>0</u>
$\begin{array}{r} 0 \\ \times 2 \\ \hline 1 \end{array}$	<u>1</u>

LSB

$$\therefore (3125)_{10} = (0.101)_2$$

$$\therefore (43.3125)_{10} = (101011.0101)_2$$

LSB = Least Significant Bit
MSB = Most Significant Bit

$$*(34.215)_{10} = (?)_5$$

$$\begin{array}{r} 34 \\ \hline 5 | \quad 6 \\ \quad \quad \quad \rightarrow 4 \\ \hline 1 \\ \hline 5 | \quad 1 \\ \quad \quad \quad \rightarrow 1 \\ \hline 0 \end{array}$$

$$(34)_{10} = (114)_5$$

$$\begin{array}{r} 215 \times 5 = 1075 & 1 \\ 075 \times 5 = 375 & 0 \\ 375 \times 5 = 1875 & 1 \\ 875 \times 5 = 4375 & 4 \\ 375 \times 5 = 1875 & 1 \\ 875 \times 5 = 4375 & 4 \\ \vdots & \vdots \\ \vdots & \vdots \\ \vdots & \vdots \end{array}$$

$$(\cdot 215)_{10} = (101414\ldots)_5$$

$$\therefore (34.215)_{10} = (114.101414\ldots)_5$$

Type 2

Any Base to Decimal.

$$(1101.101)_2 = (?)_{10}$$

$$\begin{array}{ccccccccc} & & & & & & & & \\ & \xleftarrow{\text{Positions}} & & \xrightarrow{\text{---}} & & & & & \\ \dots & 3 & 2 & 1 & 0 & -1 & -2 & -3 & \dots \\ & 1 & 1 & 0 & 1 & \cdot & 1 & 0 & 1 \\ & \downarrow & \downarrow & \downarrow & \downarrow & & \downarrow & & \downarrow \\ = 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} & & & & & & & & \end{array}$$

Base (Position)

$$= (13.625)_{10}$$

$$*(572.6)_8 = (?)_{10}$$

$$\begin{array}{r} 2 \ 1 \ 0 \ -1 \\ \hline 5 \ 7 \ 2 \cdot 6 \end{array}$$

$$= 5 \times 8^2 + 7 \times 8^1 + 2 \times 8^0 + 6 \times 8^{-1}$$

$$= (378.75)_{10}$$

⑤ Given, $B = 2$, $m = 4$, $e = [-1, 3]$, Find the highest & lowest value in conv 1 in decimal.

Soln:

$$\text{Conv1} : \pm (0 \cdot d_1 d_2 d_3 d_4)_2 \times 2^e \quad \text{where } d_1 = 1,$$

$$\begin{aligned}\text{highest value} &= (0 \cdot \underset{-1-2-3-4}{111})_2 \times 2^3 \\ &= (1 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3} + 1 \times 2^{-4}) \times 2^3 \\ &= (\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16}) \times 2^3 \\ &= \frac{15}{2} \\ &= 7.5\end{aligned}$$

$$\text{lowest value} = -7.5$$

(Ans)

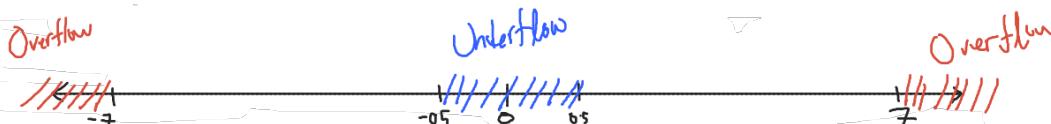
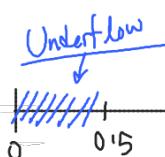
Think, how we directly found the lowest value.

Now, if we put the highest and lowest value for

$$B = 2, m = 3 \quad Q = [0, 3]$$

$$\text{highest v} = (0 \cdot 111)_2 \times 2^3 = (2^{-1} + 2^{-2} + 2^{-3}) \times 2^3 = 7$$

$$\text{lowest (non-ve) v} = (0 \cdot 100)_2 \times 2^0 = 2^{-1} \times 2^0 = 0.5$$



Floating point Representation Limits:

- Range are finite and limited. Causes Overflow & Underflow
- Possible numbers are not equally spaced.

Now, what does this mean?

Possible numbers are not equally spaced.

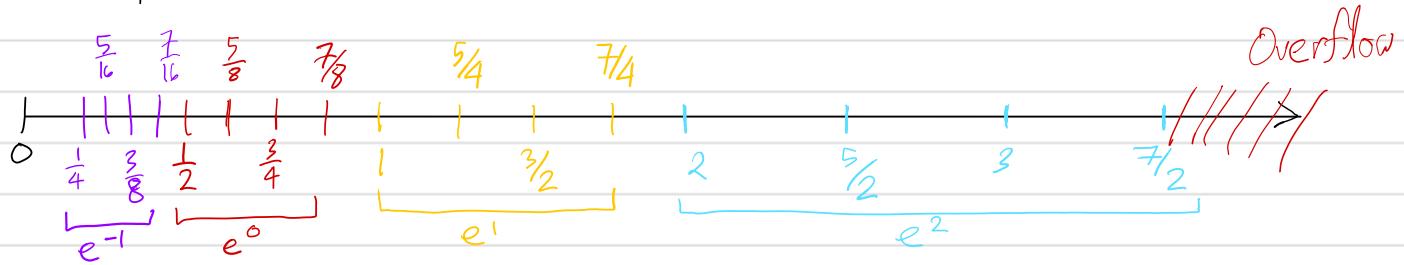
Let, $B = 2$, $m = 3$, $e = [-1, 2]$

Possible numbers are,

$$\begin{array}{lll} (0.100)_2 \times 2^{-1} = \frac{1}{4} & (0.100)_2 \times 2^0 = \frac{1}{2} & (0.100)_2 \times 2^1 = 1 \\ (0.101)_2 \times 2^{-1} = \frac{5}{16} & (0.101)_2 \times 2^0 = \frac{5}{8} & (0.101)_2 \times 2^1 = \frac{5}{4} \\ (0.110)_2 \times 2^{-1} = \frac{3}{8} & (0.110)_2 \times 2^0 = \frac{3}{4} & (0.110)_2 \times 2^1 = \frac{3}{2} \\ (0.111)_2 \times 2^{-1} = \frac{7}{16} & (0.111)_2 \times 2^0 = \frac{7}{8} & (0.111)_2 \times 2^1 = \frac{7}{4} \end{array}$$

$$\begin{array}{l} (0.100)_2 \times 2^2 = 2 \\ (0.101)_2 \times 2^2 = \frac{5}{2} \\ (0.110)_2 \times 2^2 = 3 \\ (0.111)_2 \times 2^2 = \frac{7}{2} \end{array}$$

If we put this values in a number line,



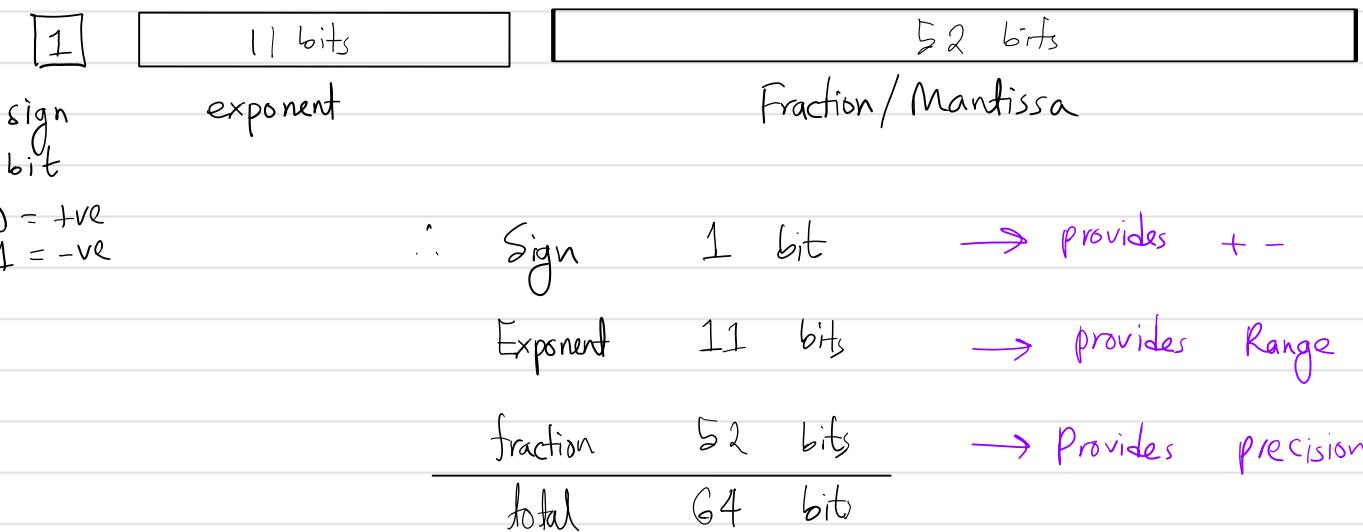
Observation:

① If the value of e stays the same, then the number are equally spaced.

② the lower the value of e , the narrower is the gap. So, more precise the value get.

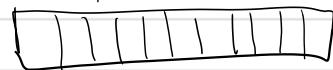
IEEE Standard (1985) For double Precision (64 bits)

To ensure every system are following a standard.



Starting with De normalized form: (Conv 3),

$$(1. d_1 d_2 \dots d_{52}) \times 2^e \quad \text{and} \quad e = 11 \text{ bits} = 2^{11} \text{ combinations} = 2048$$



$$\therefore e_{\min} = (0000000000)_2 = 0$$

$$\therefore e_{\max} = (1111111111)_2 = 2047$$

$$\therefore e = [0, 2047]$$

$$\therefore \text{largest number} = (1.11\dots1) \times 2^{2047} = \text{too large value}$$

$$\therefore \text{Smallest } \dots = (1.0\dots0) \times 2^0 = 1 = \text{Not too small.}$$

(non negative)

Problem: We can't represent 0.1, 0.11, 0.0001

We have unnecessarily big number & not too small number.

To solve this, we use Exponent Biasing,

Exponent biasing means scaling the middle value of $[0, 2047]$;

which is 1023 in such a way that it becomes 0.

Previously,



With exponent biasing,



So, we do this by subtracting 1023 from e.

$$\therefore (1 \cdot d_1 d_2 d_3 \dots d_{52}) \times 2^{e-1023}$$

this value is called bias.

This subtraction is called exponent biasing.

In case of Normalized form,

We will focus in
this form.

$$F = \pm (0 \cdot 1 d_1 d_2 \dots d_{52}) \times 2^{e-1023}$$

because,

$$(1 \cdot d_1 d_2 \dots) \times 2^{1023}$$
$$(0 \cdot 1 d_1 d_2 \dots) \times 2^{1022}$$

$$e \in [0, 2047]$$

Previously, $e \in [0, 2047]$

$$(e-1023) \in [-1023, 1024]$$

So, largest possible number,

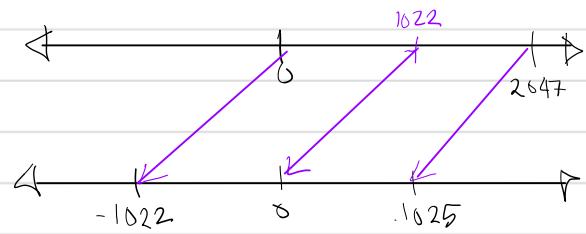
$$(0.111\dots 11) \times 2^{1025}$$

smallest possible number,

$$(0.100\dots 0) \times 2^{-1022}$$

= a very tiny value close to 0.

In normalized form, exponent is shifted like this,



* But what about absolute 0 and ∞ ?

These are special case, We have reserved exponents for them.

$$0 = 2^{e_{\min}} = 2^{-1022}$$
$$\infty = 2^{e_{\max}} = 2^{1025}$$

Finally,

$$\text{the highest possible value (except } \infty) = (0.111\ldots 1)_2 \times 2^{1024}$$
$$\approx 1.798 \times 10^{308}$$

$$\text{the lowest possible non negative value (except } 0) = (0.100\ldots 0)_2 \times 2^{-1021}$$
$$\approx 2.225 \times 10^{-308}$$

Note: You don't have to memorize these decimal values. Just see, how we calculate the value.

Example: What is the lowest value possible (non negative) value if bias is 512 in IEEE Normalized form 64 bit?

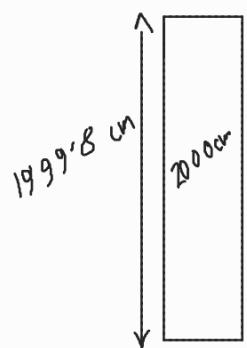
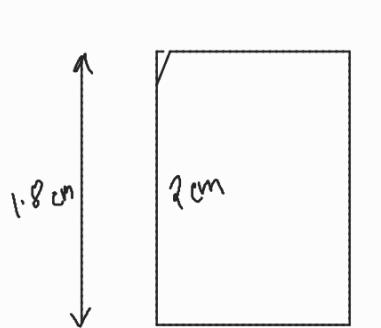
$$F = \pm (0.1d_1 d_2 \dots d_{52}) \times 2^{e - \text{bias}}$$
$$= 0.1 \times 2^{0-512}$$
$$= 0.1 \times 2^{-512}$$

(Ans)

Rounding Error

$$x = \sim \quad f_l(x) =$$

$$\begin{aligned} \text{Error} &= |(\text{Actual value} - \text{Rounded value})| \\ &= |x - f_l(x)| \end{aligned}$$



$$(2 - 1.8) = 0.2 \checkmark$$

$$(2000 - 1999.8) = 0.2 \checkmark$$

$$\text{Error} = \frac{|\text{Actual} - \text{Rounded}|}{\text{Actual}}$$

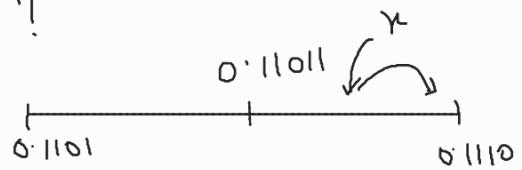
$$\frac{0.2}{2} = \underline{\underline{0.1}}$$

$$\frac{0.2}{2000} = \underline{\underline{0.0001}}$$

Scale Invariant Rounding Error,

$$S = \frac{|x - f_l(x)|}{|x|}$$

Let, $x = (0.11011101)_2$, $m = 4$, $S = ?$



$$fl(x) = (0.1110)_2$$

$$S = \frac{|(0.11011101)_2 - (0.1110)_2|}{|(0.11011101)_2|} =$$

Machine Epsilon (ϵ)

Maximum Scale Invariant Rounding Error;

$$\hookrightarrow S \uparrow = \frac{|x - fl(x)|}{|x|} \uparrow$$

$$\epsilon = S_{\max} = \frac{|x - fl(x)|_{\max}}{|x|_{\min}}$$

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

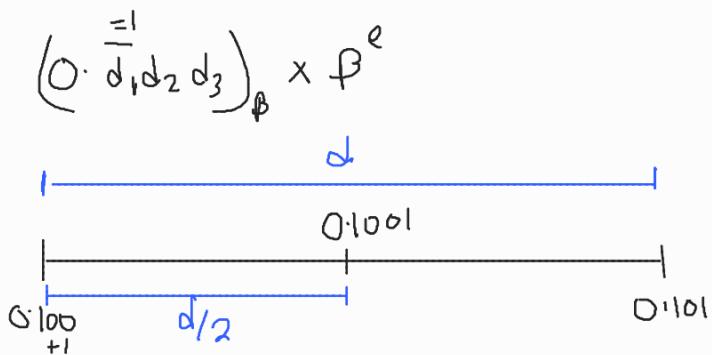
Conv 2: $\epsilon = \frac{1}{2} \beta^{-m}$

Conv 3: $\epsilon = \frac{1}{2} \beta^{-m}$

Conv 1:

$$m=3$$

$$\epsilon = S_{\max} = \frac{|x - f(x)|_{\max}}{|x|_{\min}}$$



$$d = (0.100)_2 \times 2^e \quad d/2 = \frac{1}{2} \beta^{-m} \times \beta^e$$

$$= 2^3 \times 2^e \quad \therefore |x - f(x)|_{\max} = \frac{1}{2} \beta^{-m} \times \beta^e$$

$$= \beta^{-m} \times \beta^e$$

$$\begin{aligned}
 |x|_{\min} &= (0.100)_2 \times 2^e \\
 &= (0.100)_2 \times 2^e \\
 &= 2^{-1} \times 2^e \\
 &= \beta^{-1} \times \beta^e
 \end{aligned}$$

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

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

Conv 3:

$$m=3$$

$$(1.d_1d_2d_3) \times 2^e$$

$$\epsilon = S_{\max} = \frac{|x - f(x)|_{\max}}{|x|_{\min}}$$

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

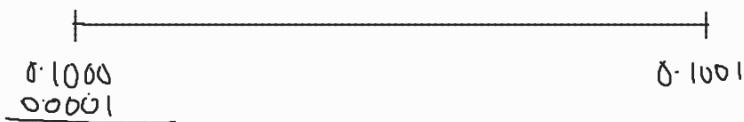
$$\begin{aligned} |x|_{\min} &= (1.000)_2 \times 2^e \\ &= 2^0 \times 2^e \\ &= 2^e = \beta^e \end{aligned}$$

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

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

Convention 2: m=3

$$(0.1d_1 d_2 d_3)_2 \times 2^e$$



$$\begin{aligned} d &= (0.0001)_{\beta} \times \beta^e \\ &= 2^{-4} \times \beta^e \\ &= \beta^{-m-1} \times \beta^e \end{aligned}$$

$$\begin{aligned} |x|_{\min} &= (0.1000)_2 \times 2^e \\ &= 2^{-1} \times 2^e \\ &= \beta^{-1} \times \beta^e \end{aligned}$$

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

$$\begin{aligned} \epsilon &= \frac{|x - fl(x)|_{\max}}{|x|_{\min}} \\ &= \frac{\frac{1}{2} \beta^{-m-1} \times \beta^e}{\beta^{-1} \times \beta^e} \\ &= \frac{1}{2} \beta^{-m-1-(-1)} \\ &= \frac{1}{2} \beta^{-m} \end{aligned}$$

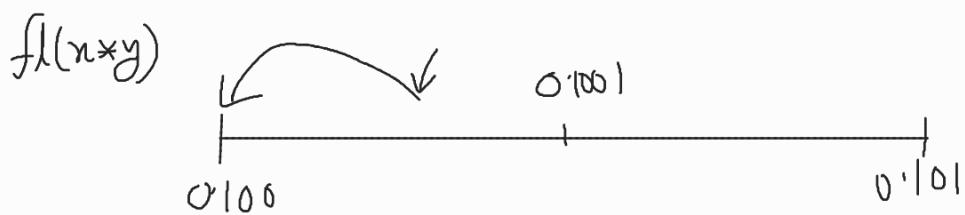
Example: $B=2$, $m=3$ calculate rounding error for $(x*y)$

$$x = \frac{5}{8} = \frac{1}{8} + \frac{4}{8} = \frac{1}{2^3} + \frac{1}{2^1} = 2^{-1} + 2^{-3} = (0.\overset{-1}{1}\overset{-3}{0}1)_2$$

$$y = \frac{7}{8} = \frac{4}{8} + \frac{2}{8} + \frac{1}{8} = \frac{1}{2^1} + \frac{1}{2^2} + \frac{1}{2^3} = 2^{-1} + 2^{-2} + 2^{-3} = (0.111)_2$$

$$x*y = \frac{5}{8} * \frac{7}{8} = \boxed{\frac{35}{64}}$$

$$\begin{aligned} x*y &= \frac{32}{64} + \frac{2}{64} + \frac{1}{64} = \frac{1}{2} + \frac{1}{32} + \frac{1}{64} = \frac{1}{2^1} + \frac{1}{2^5} + \frac{1}{2^6} \\ &= 2^{-1} + 2^{-5} + 2^{-6} = (0.\overset{-1}{1}\overset{-5}{0}\overset{-6}{0}011)_2 \end{aligned}$$



$$fl(x*y) = (0.\overset{-1}{1}0\overset{-1}{0})_2 = 2^{-1} = \boxed{\frac{1}{2}}$$

$$\begin{aligned} \text{rounding error, } \delta &= \frac{|(x*y) - fl(x*y)|}{|x*y|} \\ &= \frac{\left| \frac{35}{64} - \frac{1}{2} \right|}{\left| \frac{1}{2} \right|} = 0.09375 \end{aligned}$$

Ans

Loss of significance:

Significant figure!

5 s.f

19.215 | 105

= 19.215

$$0.0031256712 = 0.0031256 \\ = 3.1256 \times 10^{-2}$$

[Ex:

$$x^2 - 56x + 1 = 0, x_1 = ? \quad x_2 = ?$$

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

$$\frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

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

$$\sqrt{783} = 27.98$$

4 s.f

$$x_1 = 28 + 27.98 = 55.98 \checkmark$$

Solution

$$x_2 = 28 - 27.98 = 0.02$$

$$x^2 + (\alpha + \beta)x + \alpha\beta = 0$$

$$\alpha = x_1$$

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

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

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

↓ In the system

$$x \pm y \Rightarrow f(x) \pm f(y)$$

$$\Rightarrow x + \delta_1 x \pm y + \delta_2 y$$

$$\Rightarrow x(1 + \delta_1) \pm y(1 + \delta_2)$$

$$\Rightarrow (x \pm y) \pm x\delta_1 \pm y\delta_2$$

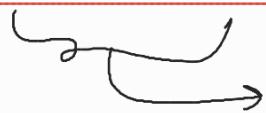
$$\Rightarrow (x \pm y) \left(1 + \frac{x\delta_1 \pm y\delta_2}{x \pm y} \right)$$

Scale invariant
error

$$x-y$$

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

$x \approx y$, $(x-y) \approx 0$.



error

significantly increases
when x closer to y .