

Rounding Error

Rounding

- * Computers can read upto certain bit or significant figures that why we need rounding.
- * If we want to round any numbers we need to draw number line.
→ Real number to F is known as rounding.
It is denoted by $f_1(x)$.

Example : (Picac - 1)

Given that $B=2$, $m=3$. The given number $(0.1000100)_2$ needs to rounded or we need to round this number.

- * For rounding any number at first we need to draw a number line to see where this number lies.

$$\begin{array}{c} + \\ \hline & (0.100)_2 \times 2^e & (0.101)_2 \times 2^e \\ \hline \end{array}$$

- * Now we need to find the middle point

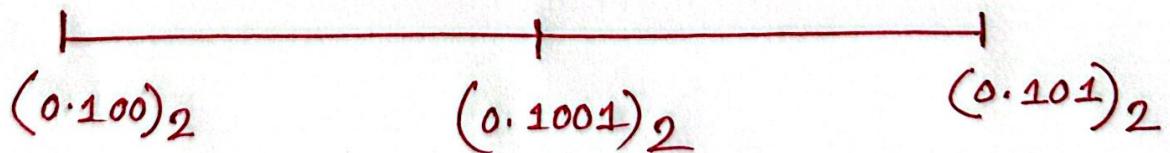
The process of finding the midpoint

$$\begin{aligned} & (0.100)_2 \\ &= 1 \times 2^{-1} \\ &= \frac{1}{2} \end{aligned}$$

$$\begin{aligned} & (0.101)_2 \\ &= 1 \times 2^{-1} + 1 \times 2^{-3} \\ &= \frac{1}{2} + \frac{1}{8} \\ &= \frac{1}{2} + \frac{1}{8} \\ &= \frac{4+1}{8} \\ &= 5/8 \end{aligned}$$

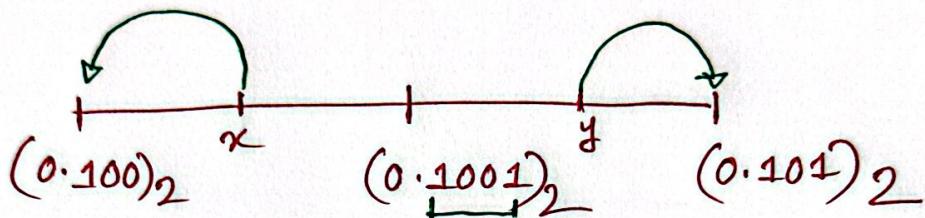
$$\begin{aligned} \frac{\left(\frac{1}{2} + \frac{5}{8}\right)}{2} &= \frac{9}{16} \\ &= \frac{8}{16} + \frac{1}{16} \\ &= \frac{1}{2} + \frac{1}{16} \\ &= 2^{-1} + 2^{-4} \\ &= (0.1001)_2 \end{aligned}$$

mid point



* How to round?

Ans:- 1) If a number lies on the left side of the number line then round the number to the left number.



Ex: $(0.\underbrace{1000100}_\text{left})_2 = (0.100)_2$

→ This number lies on the left side of the number line. That's why we rounded it to the left number $(0.100)_2$.

2) If a number lies on the right side of the number line then round the number to the right number.

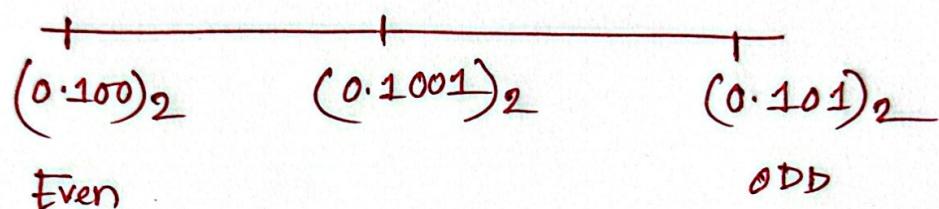
Ex: $(0.\underbrace{1001010001}_\text{right})_2 \approx (0.101)_2$

→ This number lies on the right side of the number line. That's why we round it to the right number $(0.101)_2$.

3) If the number lies exactly at middle or indicate middle point then it will get rounded to the even one in that range.

$$\text{Ex: } (0.1001)_2 = \underline{(0.100)}_2$$

This is the even number.



If a number ends in 0, it is even

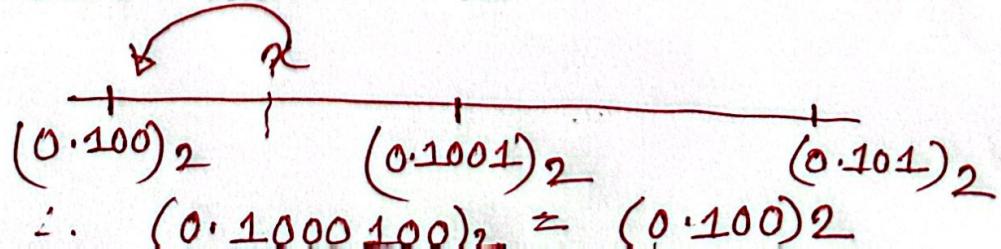
(0.100)₂ Here it ~~is~~ ends in 0 that's why it is an even number.

If a number ends in 1, it is an odd.

(0.101) Here it ends in 1 that's why it is an odd number

Our Example (previous) : Part 1

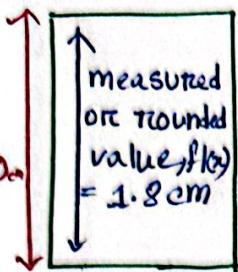
Given number $x = (0.1000100)_2$



Rounding Error

Case 1:

Actual value, $x = 2.0$

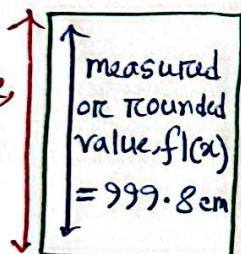


$$\begin{aligned} \text{Error}_{\text{rc}} &= |f(x) - x| \\ &= |1.8 - 2.0| \\ &= |-0.2| \\ &= 0.2 \end{aligned}$$

Annotations for the error calculation:
 - A green circle highlights the term $f(x)$ in the first equation.
 - A green bracket groups the terms $|f(x) - x|$. Arrows point from this bracket to the labels 'rounded value' (above), 'modulus' (below), and 'actual value' (below).
 - The final result '0.2' is underlined.

Case 2:

Actual value, $x = 1000$ cm



$$\begin{aligned} \text{Error}_{\text{rc}} &= |f(x) - x| \\ &= 0.2 \end{aligned}$$

From these two cases it is hard to understand the impact or significance just based on Error_{rc} .

Though the error_{rc} is same in both cases, the impact is different, and we understand this through

Scale invariant rounding error.

This is denoted by $S.$ (δ)

Formula

$$\therefore \text{Scale invariant rounding error, } S = \frac{|f_l(x) - x|}{|x|}$$

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

$$\Rightarrow S_x = f_l(x) - x$$

$$\begin{aligned}\Rightarrow f_l(x) &= S_x + x \\ &= x(1+S)\end{aligned}$$

However, we deal with maximum scale invariant rounding error, S_{\max} . This is also known as

Machine Epsilon, ϵ

$$S_{\max} = \epsilon$$

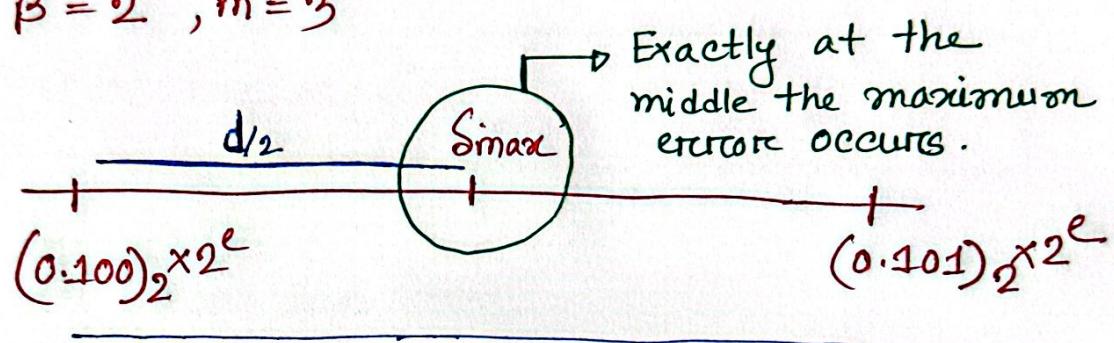
$$S_{\max} = \epsilon = \frac{|f_l(x) - x|_{\max}}{|x|_{\min}} \uparrow \downarrow$$

Formula of Machine Epsilon (ϵ) for different conventions.

Convention 1

$$(0.d_1 d_2 \dots d_m)_B \times B^e$$

Let, $B = 2, m = 3$



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

$$= (0.\overset{1}{0}\overset{2}{0}\overset{3}{1})_2 \times 2^e$$

$$= 1 \times 2^{\textcircled{3}} \times 2^e$$

$$= 1 \times B^{-m} \times B^e$$

[We know,
Given that
 $m = 3$]

$$\begin{aligned} |f(x) - x|_{\max} &= \frac{d}{2} \\ &= \frac{B^{-m} \times B^e}{2} \end{aligned}$$

$$= \boxed{\frac{1}{2} B^{-m} \times B^e} \quad (\text{important})$$

$$\begin{aligned} |x|_{\min} &= (0.\overset{1}{1}\overset{0}{0})_2 \times 2^e \\ &= 1 \times 2^{-1} \times 2^e \\ &= \boxed{B^{-1} \times B^e} \quad (\text{important}) \end{aligned}$$

$$\begin{aligned}
 \therefore \text{Machine Epsilon}(\epsilon) &= \frac{|f(x) - x|_{\max}}{|x|_{\min}} \\
 &= \frac{\frac{1}{2} \beta^{-m} \times \beta^e}{\beta^{-1} \times \beta^e} \\
 &= \frac{1}{2} \beta^{1-m} \quad (\text{important})
 \end{aligned}$$

$$\therefore \text{Machine Epsilon}(\epsilon) = \delta_{\max} = \frac{1}{2} \beta^{1-m}$$

Convention 2: (Normalized Form)

$$(1.d_1d_2 \dots d_m)_{\beta} \times \beta^e$$

Let, $\beta = 2, m = 3$

$$\begin{array}{c}
 \overbrace{+ \dots +}^{d/2} \quad | \quad \delta_{\max} \\
 \hline
 (1.000)_2 \times 2^e \qquad \qquad \qquad (1.0001)_2 \times 2^e
 \end{array}$$

$$\begin{aligned}
 d &= (0.001)_2 \times 2^e \\
 &= 1 \times 2^{-3} \times 2^e \\
 &= 1 \times \beta^{-3} \times \beta^e \\
 &= 1 \times \beta^{-m} \times \beta^e
 \end{aligned}$$

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

$$\begin{aligned}
 |x|_{\min} &= (1.000)_2 \times 2^e \\
 &= 1 \times 2^0 \times 2^e \\
 &= 2^e \\
 &= \boxed{B^e} \quad (\text{important})
 \end{aligned}$$

$$\therefore \text{Machine Epsilon, } \epsilon = \delta_{\max} = \frac{\frac{1}{2} B^{-m} B^e}{B^e} \\
 = \boxed{\frac{1}{2} B^{-m}} \quad (\text{important})$$

$$\boxed{\therefore \text{Machine Epsilon } (\epsilon) = \delta_{\max} = \frac{1}{2} B^{-m}}$$

Convention 3: (Denormalized Format)

$$(0.1d_1d_2 \dots d_m)_B \times B^e$$

Let, $B=2, m=3$

$$\begin{array}{c}
 \overbrace{\qquad\qquad\qquad}^{d/2} \qquad \qquad \qquad \delta_{\max} \\
 \hline
 (0.1000)_2 \times 2^e \qquad \qquad \qquad (0.1001)_2 \times 2^e
 \end{array}$$

$$d = (0.001)_2 \times 2^{-4}$$

$$\begin{aligned}
 &= 1 \times 2^{-4} \times 2^e \\
 &= 1 \times 2^{-3-1} \times 2^e \\
 &= 1 \times B^{-m-1} \times B^e
 \end{aligned}$$

$$\begin{aligned}
 |f_l(x) - x|_{\max} &= \frac{d}{2} = \frac{1}{2} \beta^{-m-1} \times \beta^e \\
 &= \boxed{\frac{1}{2} \beta^{-(m+1)} \times \beta^e} \quad (\text{important})
 \end{aligned}$$

$$|x|_{\min} = (0.1000)_2 \times 2^e$$

$$= (0.\overset{-1}{1}000)_2 \times 2^e$$

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

$$= \boxed{\beta^{-1} \times \beta^e} \quad (\text{important})$$

$$\begin{aligned}
 \text{Machine Epsilon, } \epsilon &= S_{\max} = \frac{|f_l(x) - x|_{\max}}{|x|_{\min}} \\
 &= \frac{\frac{1}{2} \beta^{-m-1} \beta^e}{\beta^{-1} \beta^e} \\
 &= \boxed{\frac{1}{2} \beta^{-m}} \quad (\text{important})
 \end{aligned}$$

$$\therefore \text{Machine Epsilon, } \epsilon = S_{\max} = \frac{1}{2} \beta^{-m}$$

Note: Machine Epsilon, ϵ or S_{\max} doesn't affect with the value of exponent, e .

Example : (Practice Problem)

Given that, $\beta = 2$, $m = 3$, $x = 5/8$, $y = 7/8$

Find - $f_1(x * y)$.

$$\Rightarrow x * y = \frac{5}{8} * \frac{7}{8}$$

$$= \frac{35}{64}$$

Now convert $\frac{35}{64}$ into binary.

$$\begin{aligned} & \frac{32}{64} + \frac{2}{64} + \frac{1}{64} \\ &= \frac{1}{2} + \frac{1}{32} + \frac{1}{64} \\ &= 2^{-1} + \frac{1}{2^5} + \frac{1}{2^6} \\ &= 2^{-1} + 2^{-5} + 2^{-6} \end{aligned}$$

\therefore The binary value = $(0.\underline{\hspace{2cm}}100011)_2$
We need to round the value because $m = 3$ [given]

$(0.100)_2$ $(0.1001)_2$ $(0.101)_2$

$$\begin{aligned} \text{So, } (0.100011)_2 &= (0.100)_2 \\ &= \frac{1}{2} \end{aligned}$$

P.T.O

$$\text{Actual} = \frac{35}{64}$$

$$\text{Rounded} = \frac{1}{2}$$

So rounding error occurs because

$$\frac{35}{64} \neq \frac{1}{2}$$

Loss of Significance

We know,

$$f_{\ell}(x) \neq x$$

$$f_{\ell}(y) \neq y$$

Because,

$$f_{\ell}(x) = x(1 + \delta_1)$$

$$f_{\ell}(y) = y(1 + \delta_2)$$

Now if we want to calculate $x \pm y$

$$x \pm y = f_{\ell}(x) \pm f_{\ell}(y)$$

$$= x(1 + \delta_1) + y(1 + \delta_2)$$

$$= (x \pm y) \left(1 + \underbrace{\frac{x\delta_1 \pm y\delta_2}{x \pm y}}_1 \right)$$

Scale invariant Error.

Now if we want to calculate $x-y$

For scale invariant Error, we will get

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

Note: Now, if x and y are more closer values ($x \approx y$), scale invariant error will be significantly high.

Again, if $x=y$ then denominator will be 0 and scale invariant rounding error will be infinite (∞).

This is called loss of significance.

This can occur if 1) we subtract two closer numbers, 2) Denominator will be very small or 0.
3) Scale invariant errors will be very high.

Example :

$$x^2 - 56x + 1 = 0$$

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

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

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

Lets assume computer can calculate upto 4 significant figures.

$$\sqrt{783} = 27.98$$

$$x_1 = 28 + 27.98 = 55.98$$

$$x_2 = 28 - 27.98 = 0.02000$$

Loss of significance

Because 28 and 27.98 are very close to each other.

$$0.01786 \neq 0.02000 \quad [\text{Loss of significance}]$$

→ denominators very small. That's why rounding errors etc scale invariant errors very high.

* Solution of the previous problem that occurs

$$x^2 - 56x + 1$$

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

$\alpha, \beta \rightarrow$ roots

$$\alpha\beta = 1$$

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

[From the
previous
problem]

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

$$= \frac{1}{55.98}$$

$$= 0.01786 \quad (\text{same as original } x_2)$$

∴ Hence target is to avoid subtraction.

Example: Find the average of 5.01 & 5.02.

$$\text{Actual} = \frac{(5.01 + 5.02)}{2} = 5.015$$

Let assume computer can read upto 3 significant figure.

$$\begin{aligned} \text{so, fl} \left(\frac{5.01 + 5.02}{2} \right) &= 5.015 \neq 5 \\ &= \text{fl} \left(\frac{10.03}{2} \right) \\ &= \text{fl} \left(\frac{10.0}{2} \right) \rightarrow \text{take 3 significant figures.} \\ &= 5 \end{aligned}$$

$5.015 \neq 5$
∴ Truncating occurs