

# CSE330: Numerical Methods

Topic: Rounding, Rounding  
Error, Normalized and  
Denormalized System

Prepared by:

Saad Bin Sohan  
BRAC University

Email: [sohan.academics@gmail.com](mailto:sohan.academics@gmail.com)

GitHub: <https://github.com/saad-bin-sohan>

# IEEE standard

$$e_{\min} = 0$$

$$e_{\max} = 2047$$

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

denormalized

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

$$[-1022, 1025]$$

power 1025 અને માટેનો case છે, કમ્પ્યુટરને જાણ  
infinity હોવાનો  
save કરવો

$$\dots \times 2^{\overbrace{1025}} \rightarrow \infty \text{ (infinity)}$$

$$\dots \times 2^{-1022} \rightarrow \text{ZERO}$$

this is how computer represents zero

so,

highest possible positive number (except infinity),  $= (0.1 \underbrace{11 \dots 11}) \times 2^{1024}$   
52 1's  
one

lowest possible number (except zero),

$$= (0.1 \underbrace{00 \dots 00}) \times 2^{-1024}$$

52 0's zero

[same things go for negative numbers]

# Rounding

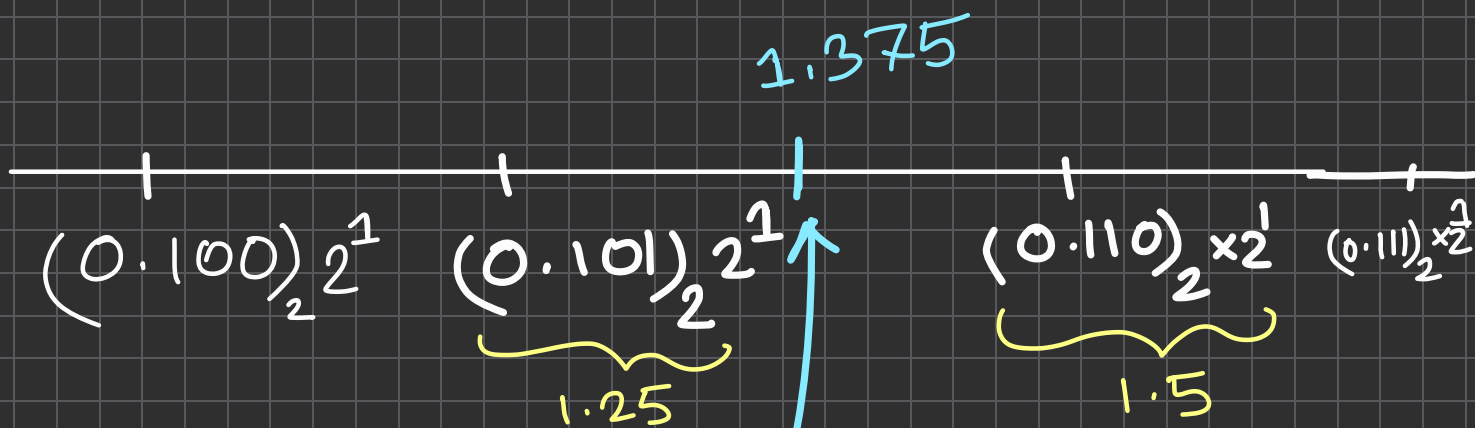
let,  $\beta = 2$

$e_{\min} = -1$

$m = 3$

$e_{\max} = 2$

constraint -  $\hookrightarrow$



$(0.1011)_2 \times 2^1$

exact middle point.  
so how to  
round up then?

given computer can  
store 3 bits

total  $m = 3$ , 3 bits  
mantissa & 4 bits.

**rule**: if the  
number is in

so  $\frac{1}{2}$  is the midpoint  
round up  $\frac{1}{2}$  represent  
 $10101$

the middle,  
then round it  
it off to the  
nearest  
(binary) even  
number

so, actual value,  $x = (0.1011)_2 \times 2^1$

round value,  $fl(x) = (0.110) \times 2^1$   
↓  
L

# Rounding Error

$$\text{Error} = \left| fl(x) - x \right| \rightarrow \text{modulus}$$

$$\underbrace{\text{relative rounding error}}_{\text{scale-invariant error}(\delta)} \delta = \frac{|fl(x) - x|}{|x|}$$

scale-invariant error( $\delta$ )

max value of  $\delta$  = Machine Epsilon ( $\epsilon_m$ )

[37 convention  $\epsilon_m < 10^{-37}$  value 2020]

Que Let,  $\beta = 2$   $e_{\min} = -1$   
 $m = 3$   $e_{\max} = 2$

convention - 1  $\hookrightarrow$   $e_m$   $\hookrightarrow$  value  $\frac{1}{2}$ ?

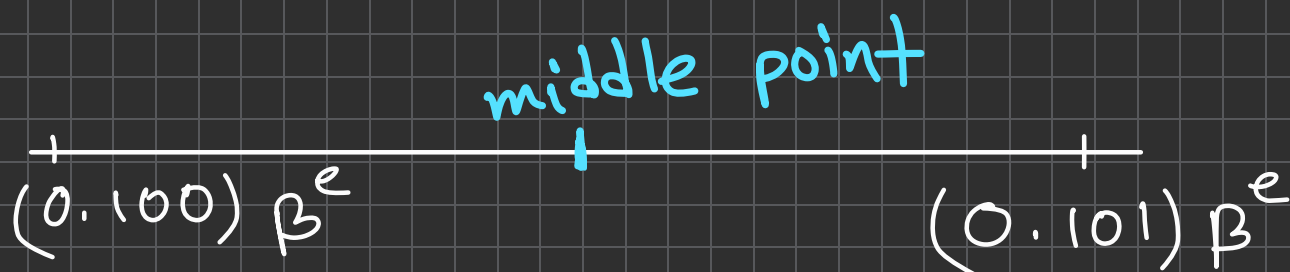
Sol<sup>n</sup>:  $e_m \in \mathbb{Z}^m$  rel. rounding  $\hookrightarrow$  max value.

so,  $|f_l(x) - x|$  has to be maximum  
and  $|x|$  has to be minimum.

we will find

the  $|f_l(x) - x|$  to be maximum, when  
actual given number  $\frac{1}{2}$   $\hookrightarrow$   $\frac{2^{m-1}}{2^m} = \frac{1}{2}$   
representable number  $\frac{1}{2}$  middle  $\hookrightarrow$   $\frac{1}{2}$ .

so, for any  $e$ ,



Maximum difference,

$$= \frac{1}{2} \left[ (0.101) \beta^e - (0.100) \beta^e \right]$$

$$= \frac{1}{2} (0.001) \beta^e$$

$$= \frac{1}{2} \times 2^{-3} \beta^e$$

$$= \frac{1}{2} \times \beta^{-3} \beta^e$$

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

$$\text{min of } x = (0.100) \beta^e$$

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



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

and now,

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

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

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

→ for convention-1 system

[ N.B:  $E_m$  is always decimal or floating point or  $2^n$  or  $2^m$  else - 1 marks gone.

IIII: 0.625 ✓ correct

$\frac{1}{2} \times 2^{-3}$  ✗ wrong

for both normalized and denormalised system,

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

arithmetic operation

Que

let,

$$\beta = 2$$

$$x = \frac{5}{8}$$

$$m = 3$$

$$y = \frac{7}{8}$$

$$e_{\min} = -1$$

$$e_{\max} = 2$$

$$x * y = ?$$

sol<sup>n</sup>:

$$x * y = f_l(x) * f_l(y)$$

$$x = \frac{5}{8} = (0.\underbrace{101})_2 \times 2^0 = f_l(x)$$

$$y = \frac{7}{8} = (0.\underbrace{111})_2 \times 2^0 = f_l(y)$$

m limit 10 2 (5) 3 2 x, y 10

rounding not needed

so,

$$x * y = fl(x) * fl(y)$$

$$= (0.101)_2 \times 2^0 \times (0.111)_2 \times 2^0$$

$$= (0.100011)_2 \times 2^0$$

0.100:011  
m=3  
division  
common sense  
to rounding in  
left number is

270, because of  
270

$$\text{so, } fl(xy) = \underbrace{(0.100)_2 \times 2^0}$$

that's how the  
computer will save it,  
as a rounded value, instead of actual  
numbers

(NB:

if  $(m+1)$ th digit is zero, then round it to previous number

elif that's 1, round it up to the next number)

$$\left[ \text{error} \Rightarrow \frac{|f(x) - x|}{|x|} \right]$$