

Brac University  
Department of Computer Science and Engineering  
Summer 2025

Name:  
Student ID:  
Section:

10 Marks  
25 Minutes

1. If  $\beta = 2$ , fraction is of 2 bits exponent is of 3 bits, What will be the non-negative lowest and highest number that can be generated using the Normalized form of floating point representation? N.B: Exponent range starts from 1.

exponent is of 3bits, So exponent possible  $2^3 = 8$  different values  
non-negative lowest =  $0.100 \times 2^1$   
highest =  $0.111 \times 2^8$

2. Derive the formula for machine epsilon ( $\epsilon_M$ ) for the Normalized form of floating point representation.

Same as lecture notes

3. For a system if  $\beta = 2$ ,  $m=4$  and  $e \in \{-3,3\}$  then how many non-negative number can be represented in a system following Denormalized form?

$$\text{Denormalized Form} = \pm 1.d_1d_2d_3d_4 \times 2^e$$

$$\therefore \text{Total non-negative values} = 2^4 \times 1 = 112$$

4. If  $x=31/16$ , find  $fl(x)$  where  $m=4$ ,  $e \in \{-3,3\}$  and the system follows Normalized Form of floating point representation. Also find the relative rounding error.

$$\begin{aligned} x = \frac{31}{16} &= \frac{16}{16} + \frac{8}{16} + \frac{4}{16} + \frac{2}{16} + \frac{1}{16} = 2^0 + 2^{-1} + 2^{-2} + 2^{-3} + 2^{-4} \\ &= 1.1111 = 0.11111 \times 2^1 \end{aligned}$$

So this number can be represented using our system.

$$\therefore fl(x) = x = \frac{31}{16}$$

$$\therefore \text{relative error} = \frac{\left| \frac{31}{16} - \frac{31}{16} \right|}{\left| \frac{31}{16} \right|} = 0.$$