# Assignment 1

Name : Shihab Muhtasim

ID : 21301610

Sec : 3 (MAA)

course : CSE 330

## (a)

Given,
$$\beta = 2, \quad m = 4, \quad -4 \leq e \leq 2.$$

In this system,

① Lecture note form:

maximum number $= \pm(0.1111)_2 \times 2^2$

② Normalized form:

maximum number $= \pm(1.1111)_2 \times 2^2$

③ De normalized form:

maximum number $= \pm(0.11111)_2 \times 2^2$

## (b)

similarly for the system,

① Lecture note form:

minimum number $= +(0.1000)_2 \times 2^{-4}$

② Normalized form:

minimum number $= +(1.0000)_2 \times 2^{-4}$

③ Denormalized form:

minimum number $= +(0.10000)_2 \times 2^{-4}$

$\times 2^{-4}$

## (c)

Using eqr 1, for $e = -3$ the numbers generated are,

1) $(0.1000)_2 \times 2^{-3} = 0.5 \times 2^{-3} = 0.0625 = \frac{1}{16}$

2) $(0.1001)_2 \times 2^{-3} = \frac{9}{16} \times 2^{-3} = 0.0703 = \frac{9}{128}$

3) $(0.1010)_2 \times 2^{-3} = 0.625 \times 2^{-3} = 0.078 = \frac{5}{64}$

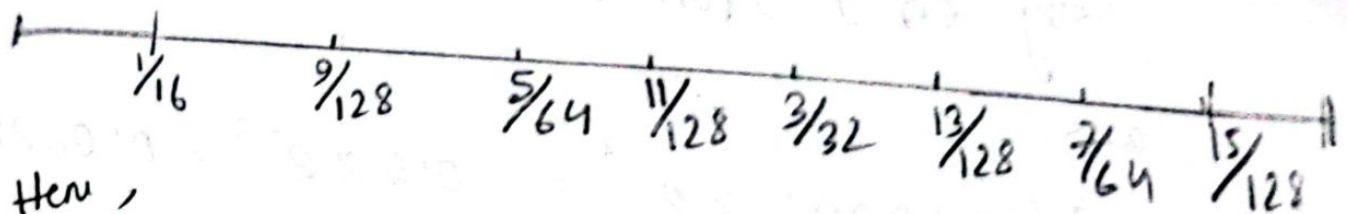4) $(0.1011)_2 \times 2^{-3} = 0.6875 \times 2^{-3} = 0.085 = \frac{11}{128}$

5) $(0.1100)_2 \times 2^{-3} = 0.75 \times 2^{-3} = 0.093 = \frac{3}{32}$

6) $(0.1101)_2 \times 2^{-3} = 0.8125 \times 2^{-3} = 0.101 = \frac{13}{128}$

7) $(0.1110)_2 \times 2^{-3} = 0.875 \times 2^{-3} = 0.109 = \frac{7}{64}$

8) $(0.1111)_2 \times 2^{-3} = 0.9375 \times 2^{-3} = 0.117 = \frac{15}{128}$

Number line is given below:



Hence,

The number line is equally spaced as the diffrence between all the numbers to its consecutive one is $\frac{1}{128}$.

However, when considering value of e to be lower / higher, it can differ.

## Ans to or no 2

### (a)

Given, $\beta = 2$, $m = 5$, $e_{min} = -2$, $e_{max} = 5$

In normalized form :

minimum number, $|x| = +(1.00000)_2 \times 2^{-2}$

$$= +(0.25)_{10}$$

In denormalized form :

minimum number, $|x| = +(0.100000)_2 \times 2^{-2}$

$$= +(0.125)_{10}$$

## (b)

Calculating machin epsilon :
for both normalized and denormalized,

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

Given, $\beta = 2$, $m = 5$,

$\therefore$ Machine epsilon for both normalized and denormalized will be,

$$\epsilon m = \frac{1}{2} \times 2^{-5} = \frac{1}{64} = 0.015625$$

## (c)

Maximum delta i's the machine epsilon value which will be for the normalized form ($e$ or 2) :

$$\epsilon m = \frac{1}{2} \times 2^{-5} = \frac{1}{64} = 0.015625$$

## Ans to qr 3

### (a)

i) $(2.23)_{10}$ to binary =

$(10.00111010111000010l)_2 \times 2^0$

Representation in normalized form :

$(1.00011010101110000101)_2 \times 2^1$

Since the system allows $m=3$ & finding the closest value :



$$\underset{(1.000)_2 \times 2^1}{\overset{(2)_{10}}{\vdash}} \qquad \underset{(1.000111\ldots)_2 \times 2^1}{\overset{2.23}{\mid}} \qquad \underset{(1.001)_2 \times 2^1}{\overset{(2.25)}{\mid}}$$

we can see that $(m+1)^{th}$ 's value is 1 : (1.001)

$\therefore$ Rounding up $= (1.001)_2 \times 2^1$

(11) $(2.2018)_{10} = (10.0011000111010100100...)_2 \times 2^0$

$\times 2^1$

Normalized : $(1.000110011101010010011)_2 \times 2^1$

Since $m = 3$



$(2)_{10}$      $(2.2018)_{10}$     $(2.25)_{10}$

$(1.000)_2 \times 2^1$   $(1.00011...)_2 \times 2^1$   $(1.001)_2 \times 2^1$

Since $(m+1)$th val is $1$ :

∴ Rounded value $= (1.0001)_2 \times 2^1$

---

## (b)

$^{1)}$ for $(2.23)_{10}$ :

Rounded value , $fl(x) = (2.25)_{10} = (1.001)_2 \times 2^1$

actual value , $x = (2.23)_{10} = (1.00011101011110000 101)_2 \times 2^1$

∴ Rounding error $= |fl(x) - x| = |$
$= |(1.001)_2 \times 2^1 - (1.0001110101110000101)_2 \times 2^1|$
$= (0.0000001010001111)_2 \times 2^1 \times 2^1)_2$

$^{11)}$ for $(2.2018)_{10}$ : $= (1.001)_2 \times 2$

Rounded val , $fl(x) = (2.25)_{10}$
actual val , $x = (2.2018)_{10} = (1.00011001110101000100011)_2 \times 2^1$

$\searrow$ Rounding error $= |fl(x) - x| = |$
$=$

(i) for ... (L...

q

## (C)

... this system and denormalized form,

$max = (0.1111)_2 \times 2^2 = (3.75)_{10}$

$min = (0.1000)_2 \times 2^{-2} = (2) \cdot_{10} (0.125)_{10}$

∴ The numbers are within the range. So it is representable.

(1) $(2.23)_{10}^{\circ}$ :

$(2.23)_{10} = (10.001110101111000101)_2 \times 2^0$

In denormalized form:

$(0.100011101011110000101)_2 \times 2^2$

since, m = 3 and (m+1)th val is 1 :

```
        2            2.23              2.25
        |             |                 |
  (0.1000)₂×2²  (0.1000111.)₂×2²  (0.1001)₂×2²
```

∴ Rounded = $(0.1001)_2 \times 2^2$

(11) $(2 \cdot 2018)_{10} = (10 \cdot 0011001110101010010011)_2 \times 2^0$

Denormalized $= (0 \cdot 100011001110101010010011)_2 \times 2^2$

Since $m = 3$ and $(m+1)^{th}$ val is 1 :

$$\overset{2}{\underset{\underset{(0\cdot 1000)_2 \times 2^2}{\big|}}{\vdash}}\underline{\hspace{2cm}}\overset{2 \cdot 2018}{\underset{\underset{(0 \cdot 100011..)_2 \times 2^2}{\big|}}{\big|}}\underline{\hspace{2cm}}\overset{2 \cdot 25}{\underset{\underset{(0 \cdot 1001)_2 \times 2^2}{\big|}}{\big|}}$$

∴ Rounded $= (0 \cdot 1001)_2 \times 2^2$

## (b)

(i) for $(2.23)_{10}$ :

Rounded value, $fl(x) = (1.001)_2 \times 2^1$

$\qquad\qquad\qquad\quad = (2.25)_{10}$

Actual value, $x = (2.23)_{10}$

$\therefore$ Rounding error $= |fl(x) - x|$

$\qquad\qquad\qquad = (2.25 - 2.23)$

$\qquad\qquad\qquad = (0.02)_{10}$

$\qquad\qquad\qquad = (0.0000010100001111011)_2$

$\qquad\qquad\qquad = (1.01000111011)_2 \times 2^0 \times 2^{-6}$

$\qquad\qquad\qquad = (1.010)_2 \times 2^{-6}$

(ii) for $(2.2018)_{10}$ :

Rounded value, $fl(x) = (1.001)_2 \times 2^1$

$\qquad\qquad\qquad\quad = (2.25)_{10}$

Actual value, $x = (2.2018)_{10}$

Rounding error $= |fl(x) - x|$

Rounding error $= |2.25 - 2.2018|$

$= (0.0482)_{10}$

$= (0.0000110001010110110101)_2$

$= (1.10001010110110101)_2 \times 2^{-5} \times 2^0$

$= (1.100)_2 \times 2^{-5}$