

BRAC UNIVERSITY

CSE330

NUMERICAL METHODS

Assignment 1

Student Information

NAME: Md. Danial Islam ID: 20101534 SECTION: 10



Inspiring Excellence

Date: 04 February 2023

Answer to the question no 1

1

(a) given $\beta = 2$, $m = 4$, $e = [-3, 6]$

general form,

$$\Rightarrow \beta^{m-1} \times \text{count of exponent}$$

$$\Rightarrow 2^{4-3} \times 10$$

$$\Rightarrow 80 \text{ Possible numbers}$$

normalized form,

$$\Rightarrow \beta^m \times \text{count of exponent}$$

$$\Rightarrow 2^4 \times 10$$

$$\Rightarrow 160 \text{ Possible numbers}$$

denormalized form,

$$\Rightarrow \beta^m \times \text{count of exponent}$$

$$\Rightarrow 2^4 \times 10$$

$$= 160 \text{ (Possible numbers)}$$

(b) largest number possible,

$$\beta = 2$$

$$m = 4$$

$$e = [-3, 6]$$

general form,

$$\begin{aligned} & (0.1111)_2 \times \beta^{e_{\max}} \\ \Rightarrow & (0.1111)_2 \times 2^6 \Rightarrow \left(\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16} \right) \times 2^6 \\ & \Rightarrow 60 \end{aligned}$$

normalized form,

$$\begin{aligned} & (1.1111)_2 \times 2^6 \Rightarrow \left(1 + 2^{-1} + 2^{-2} + 2^{-3} + 2^{-4} \right) \times 2^6 \\ & \Rightarrow 124 \end{aligned}$$

denormalized form,

$$\begin{aligned} & (0.11111)_2 \times 2^6 \Rightarrow \left(2^{-1} + 2^{-2} + 2^{-3} + 2^{-4} + 2^{-5} \right) \times 2^6 \\ & \Rightarrow 62 \end{aligned}$$

(c) non-negative smallest number

$$\beta = 2$$

$$m = 4$$

$$e = [-3, 6]$$

general form,

$$(0.1000)_2 \times \beta^{e_{\min}}$$

$$\Rightarrow (0.1000)_2 \times 2^{-3}$$

$$\Rightarrow 2^{-1} \times 2^{-3} \Rightarrow \frac{1}{16}$$

normalized form,

$$(1.0000)_2 \times \beta^{e_{\min}}$$

$$\Rightarrow (1)_2 \times 2^{-3}$$

$$\Rightarrow \frac{1}{8}$$

de-normalized form,

$$(0.10000)_2 \times \beta^{e_{\min}}$$

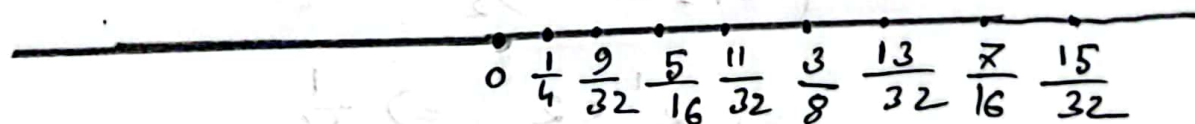
$$\Rightarrow 2^{-1} \times 2^{-3} \Rightarrow \frac{1}{16}$$

(d) using general form,

$$\beta = 2, m = 4, e = -1$$

then all the numbers are,

$$\begin{array}{lcl} (0.1000)_2 \times 2^{-1} \Rightarrow & \frac{1}{4} & \searrow \text{---} \frac{1}{32} \\ (0.1001)_2 \times 2^{-1} \Rightarrow & \frac{9}{32} & \searrow \frac{1}{32} \\ (0.1010)_2 \times 2^{-1} \Rightarrow & \frac{5}{16} & \searrow \frac{1}{32} \\ (0.1011)_2 \times 2^{-1} \Rightarrow & \frac{11}{32} & \searrow \frac{1}{32} \\ (0.1100)_2 \times 2^{-1} \Rightarrow & \frac{3}{8} & \searrow \frac{1}{32} \\ (0.1101)_2 \times 2^{-1} \Rightarrow & \frac{13}{32} & \searrow \frac{1}{32} \\ (0.1110)_2 \times 2^{-1} \Rightarrow & \frac{7}{16} & \searrow \frac{1}{32} \\ (0.1111)_2 \times 2^{-1} \Rightarrow & \frac{15}{32} & \searrow \frac{1}{32} \end{array}$$



here the difference between two numbers are $\frac{1}{32}$ and it's equal for every number of this series, so they are equally spaced.

Answer to the question no 2

(a) given $\beta = 2$, $m = 3$, $e_{exp} = [0, 15]$ as 4 bit
as, e_{max} & e_{min} reserved then $e = [1, 14]$
 \therefore new $e_{max} = 14$, $e_{min} = 1$

\therefore ~~smallest~~

\therefore ~~normalized~~ form,

$$\text{smallest} = (1.000)_2 \times 2^1$$

$$\text{maximum} = (1.111)_2 \times 2^{14}$$

denormalized form,

$$\text{smallest} = (0.1000)_2 \times 2^2$$

$$\text{maximum} = (0.1111)_2 \times 2^{14}$$

— 0 —

(b) Machine epsilon, ϵ_m for normalized form

$$\Rightarrow \frac{1}{2} \beta^{-m} \Rightarrow \frac{1}{2} 2^{-3}$$

$$\Rightarrow \frac{1}{16} \text{ (Ans)}$$

(c) For max delta error value for General form

$$\begin{aligned}
 E_m &= \frac{1}{2} \beta^{1-m} \\
 &= \frac{1}{2} 2^{1-3} \\
 &= \frac{1}{2} \times 2^{-2} = \frac{1}{8} \text{ (Ans)}
 \end{aligned}$$

Answer to the question no 3

$\beta = 2, m = 3, \text{ ~~exp = 4 bit~~ } e_{\min} = -1$

$e_{\max} = 2$

Normalized form,

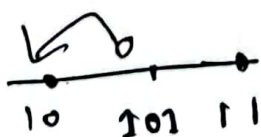
$x = 6.25$

$\therefore (6.25)_{10} \Rightarrow \textcircled{1100}. (6)_{10} = 110$

$(.25)_{10} = (.01)_2$

$(6.25)_{10} = (110.01)_2 \times 2^0 \Rightarrow (1.1001)_2 \times 2^2$

$\Rightarrow (1.10)_2 \times 2^2$



Answer to the question no 3

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

normalized form,

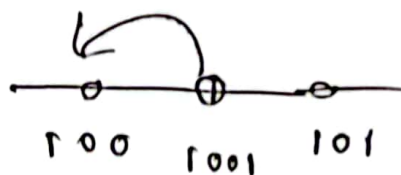
$$x = 6.25 = \frac{25}{4} = \frac{1+8+16}{4} = 2^{-2} + 2^1 + 2^2$$

$$= \text{~~(6.01)}_2~~$$

$$= (110.01)_2$$

$$= (1.1001)_2 \times 2^2$$

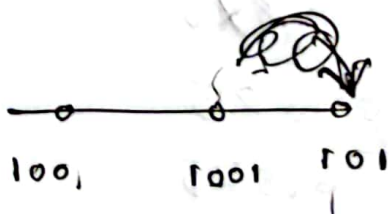
$$= (1.100)_2 \times 2^2 \quad (\text{Ans})$$



$$x = 6.875 = (110.111)_2 \times 2^0$$

$$= (1.10111)_2 \times 2^2$$

$$= (1.101)_2 \times 2^2 \quad (\text{Ans})$$



$$\begin{aligned}
 (b) \quad f_1(x) &= f_1(6.25) \\
 &= (1.100)_2 \times 2^2 \\
 &= (110.0)_2 \times 2^0 \\
 &= 6
 \end{aligned}$$

$$\begin{aligned}
 \delta_1 &= \frac{|x - f_1(x)|}{x} = \frac{|6.25 - 6|}{6.25} \\
 &= 0.04
 \end{aligned}$$

$$\begin{aligned}
 f_1(x) &= f_1(6.875) \\
 &= (1.101)_2 \times 2^2 \\
 &= (110.1)_2 \times 2^0 \\
 &= 6.5
 \end{aligned}$$

$$\begin{aligned}
 \therefore \delta_2 &= \frac{|x - f_1(x)|}{x} = \frac{|6.875 - 6.5|}{6.875} \\
 &= \frac{0.375}{6.875} \\
 &= 0.0545
 \end{aligned}$$

(c) Norm (b).

$$\begin{aligned} (6.25)_{10} &= (\overline{1.1001})_2 \times 2^2 \\ &= (110.01)_2 \times 2^0 \end{aligned}$$

$$\text{and } (6.875)_{10} = (110.111)_2 \times 2^0$$

to convert denormalized form

the floating point will be in this form

$$(0.1d_1d_2d_3 \dots d_n) \times \beta^e$$

$$\begin{aligned} (\overline{6.25})_{10} &= (110.01)_2 \times 2^0 \\ \Rightarrow (0.11001)_2 \times 2^3 \end{aligned}$$

$$\begin{aligned} & \& (110.111)_2 \times 2^0 \\ \Rightarrow (0.110111)_2 \times 2^3 \end{aligned}$$

So, if we try to convert it to denormalized form ~~it will be~~ the exponent will become 3 but an $e_{max} = 2$ is not possible.

(d)

Machine epsilon, ϵ_m

Machine rate,

$$\begin{aligned}\frac{1}{2} \beta^{1-m} &= \frac{1}{2} \times 2^{1-3} \\ &= \frac{1}{2} \times 2^{-2} \\ &= \frac{1}{8} \text{ (Ans)}\end{aligned}$$

Normalized and denormalized,

$$\begin{aligned}\frac{1}{2} \beta^{-m} &= \frac{1}{2} \times 2^{-3} \\ &= \frac{1}{16} \text{ (Ans)}\end{aligned}$$