

# CSE 330

i) Standard form  $x = \pm (0.d_1d_2 \dots d_m) \times \beta^e$   
 $\underline{d_1 \neq 0}$

ii) Normalized

$$x = \pm (0.1d_1d_2 \dots d_m) \beta \times \beta^e$$

iii) Denormalized

$$x = \pm (1.d_1d_2 \dots d_m) \beta \times \beta^e$$

$$\textcircled{1} \quad \beta = \underline{2}, \underline{m=3} \quad \frac{-2}{e_{\min}} \leq e \leq \frac{1}{e_{\max}}$$

a) Find minimum and maximum number with and without negative exponent

i) Standard form  $\rightarrow (0.\underline{d_1d_2d_3}) \times 2^e$   
 w/o -ve

$$\text{Min} = (0.100)_2 \times 2^{-2}$$

$$\text{Max} = (0.111)_2 \times 2^1$$

$$\boxed{e = \{-2, -1, 0, 1\}}$$

with -ve

$$\text{Min} = - (0.111)_2 \times 2^1$$

$$\text{Max} = + (0.111)_2 \times 2^1$$

ii) Normalized  $\rightarrow (0.\underline{1d_1d_2d_3})_2 \times 2^e$

w/o -ve

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

$$1^{\text{Max}} = (0.1111)_2 \times 2^1$$

with -ve

$$\text{Min} = -(0.1111)_2 \times 2^1$$

$$\text{Max} = +(0.1111)_2 \times 2^1$$

$$\text{if i) Denormalized} \rightarrow (1.d_1d_2d_3)_2 \times 2^e$$

w/o -ve,

$$\text{Min} = (1.000)_2 \times 2^{-2}$$

$$\text{Max} = (1.111)_2 \times 2^1$$

w -ve.

$$\text{Min} = -(1.111)_2 \times 2^1$$

$$\text{Max} = (1.111)_2 \times 2^1$$

b) How many numbers can be represented in each form?

i) Standard  $\rightarrow (0.\underline{d_1d_2d_3})_2 \times 2^e \quad \{-2, -1, 0, 1\}$

$\downarrow d_1 \neq 0 \rightarrow d_1 = 1$

$(0.1\underline{00}) \times 2^{-2}$        $(0.1\underline{01}) \times 2^{-2}$        $(0.1\underline{10}) \times 2^{-2}$        $(0.1\underline{11}) \times 2^{-2}$

$\begin{array}{c} \square \quad \square \\ d_2 \quad d_3 \\ \frac{1}{2} \times 2 = 1 \end{array}$

$\Rightarrow 4 \times 4 = 16$

Free Support  
16x2

ii) Normalized form

$$(0.\underline{1d_1d_2d_3})_2 \times 2^e$$

$\begin{array}{c} \square \quad \square \quad \square \\ 0/1 \quad 0/1 \quad 0/1 \\ 2 \times 2 \times 2 \\ = 2^3 \end{array}$

$$2^3 \times 4 = 8 \times 4 = 32$$

-ve support

$\{0.1000\}$   
 $\{0.1001\}$   
 $\{0.1010\}$   
 $\{0.1011\}$   
 $\{0.1100\}$   
 $\{0.1101\}$   
 $\{0.1110\}$   
 $\{0.1111\}$

$$32 \times 2 = 64$$

iii)  $(1. \underline{d_1 d_2 d_3}) \times 2^e$

$\square \quad \square \quad \square = 2^3$

$2 \times 2 \times 2$   
 $0/1 \quad 0/1 \quad 0/1$

$2^3 \times 4 = 32$

-ve support  
 64

c) IEEE format (759) - 64 bit precision

$m = 52$ ,  $e = \underline{11 \text{ bits}}$  sign bit = 1 bit  
 bits

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

$e_{\min} = 0$

$e_{\max} = 2^{11} - 1 = \underline{2047}$

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

Smallest no  $\rightarrow (1. 0000 \dots 000) \times 2^0 = \underline{1}$

Largest no  $\rightarrow (1. 1111 \dots 111) \times 2^{2047} = \underline{1}$

$\sum$   $\rightarrow$   $\text{reg}$

Binarying

$$(1. d_1 d_2 \dots d_{52})_2 \times 2^{\frac{e-1023}{\text{Bias}}}$$

$$= (0. 1 d_1 d_2 \dots d_{52})_2 \times 2^{\frac{e-1023}{\text{Bias}}} \times 2^1$$

$$= (0. 1 d_1 d_2 \dots d_{52})_2 \times 2^{\frac{e-1022}{\text{Bias}}}$$

$$e_{\min} \rightarrow 0 \quad e_{\max} = \underline{2047}$$

$$\text{Minimum} = (0. 1000 \dots 0_{52}) \times 2^{\frac{-1022}{\text{Bias}}}$$

$$\text{Maximum} = (0. 1111 \dots 1_{52}) \times 2^{\frac{1025}{\text{Bias}}}$$

Theoretical

By convention

$\approx 0$

$\approx \text{inf}$

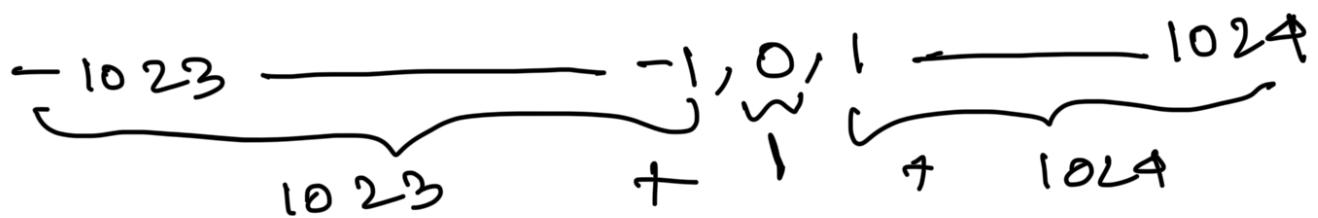
Actual

$$\text{minimum} = (0. 1000 \dots 0_{52}) \times 2^{-1021}$$

$$\text{maximum} = (0. 1111 \dots 1_{52}) \times 2^{1024}$$

Bias  $\rightarrow$  why  $\boxed{1023}$ ?

exponent, e range  $\frac{0 \text{ to } 2047}{2048}$



$$\begin{aligned}
 \text{Bias} &= \frac{2^e}{2} - 1 \\
 &= \frac{2^{e-1}}{2} - 1 \\
 &= 2^{e-1} - 1
 \end{aligned}
 \quad = 2048$$

②  $(6.25)_{10}$

Convert into **Standard floating point representation** where  $\beta=2, m=3$ ,

$$-1 \leq e \leq 3$$

$$(0.\underline{d_1d_2d_3}) \times 2^e$$

$$\begin{aligned}
 (6.25)_{10} &= (110.01)_2 \\
 &= (0.\underline{110}|\underline{01})_2 \times 2^3
 \end{aligned}$$

$$\begin{aligned}
 (0.\underline{1d_1d_2d_3}) \\
 = (0.\underline{1100})_2 \times 2^3
 \end{aligned}$$

Rounding?

$$0.\underbrace{110}_{m=3} |\boxed{0}|_1$$

i) If  $m+1$  th bit = 0, just truncate

ii) If  $m+1$ th bit = 1

ii-a)  $0.110 \underline{\underline{1}}$        $0.111 \underline{1}$

Round to the nearest even no

$0.110$

ii-b)  $0.110 \underline{\underline{101}} \dots$

Original number

$0.11 \boxed{1}$  increase  $m$ th bit by 1

$\underline{(6.25)_{10}}$

After rounding  $\rightarrow$

$$\underline{(0.110)_2} \times 2^3 =$$

$$= (2^{-1} \times 1 + 2^{-2} \times 1) \times 2^3$$

$$= \left(\frac{1}{2} + \frac{1}{4}\right) \times 2^3 = \frac{3}{4} \times 8 = \underline{16}$$

$$R.E = \sqrt{6.25 - 6} = \sqrt{31(x) - x} \\ = 0.25$$

③ If  $x = \frac{3}{8}, y = \frac{5}{8}$  find  $f_1(x \times y)$   
where  $\underline{m = 4}$  find rounding error if any.

$$x = \frac{1}{8}$$

$$= 0.375$$

$$= \underline{\underline{(0.011)_2}}$$

$$0.375$$

$$\begin{array}{r} \times 2 \\ \hline 0.75 \end{array}$$

$$\begin{array}{r} \times 2 \\ \hline 1.5 \end{array}$$

$$\begin{array}{r} \times 2 \\ \hline 1.0 \end{array}$$

$$y = \frac{1}{8}$$

$$= \frac{0.625}{101_2}$$

$$= \underline{\underline{(0.101)_2 \times 2^0}}$$

$$= \underline{\underline{(0.11)_2 \times 2^{-1}}}$$

$$f_1(x) = (0.11)_2 \times 2^{-1}$$

$$f_1(y) = (0.101)_2 \times 2^0$$

Standard form  $(0. \underline{d_1 d_2 d_3 d_4}) \times 2^{\underline{e}}$   
 $\boxed{d_1 \neq 0}$

$$\boxed{x} \times \boxed{y} = f_1(x \times y)$$

$$= f_1(x) \times \underline{f_1(y)}$$

$$= \underline{(0.11)_2 \times 2^{-1}} \times \underline{(0.101)_2 \times 2^0}$$

$$= \frac{3}{4} \times \frac{1}{2} \times \frac{5}{8}$$

$$= \frac{15}{64} = 0.234375$$

$$= (0.\underline{\underline{001111}})_2$$

$$= \underline{\underline{(0.1111)_2 \times 2^{-2}}}$$

$$\text{Rounding} = \frac{15}{16} \times 2^{-2}$$

$$f1(x \otimes y) = \frac{15}{64}$$

$$R.E = \left| \frac{15}{64} - \frac{15}{64} \right| = \left| f1(x \otimes y) - f1(x \otimes y) \right| \\ = 0$$

Short technique  $\rightarrow$

Decimal  $\rightarrow$   $\frac{675}{10 \rightarrow 10^1} = 67.5$

$$\frac{675}{10^3 \rightarrow 1000} = 0.675$$

$$\frac{675}{10^4 \rightarrow 10000} = 0.0675$$

$$6.75 \times 10^1 = 67.5$$

$$\frac{6.75 \times 10^2}{6.75 \times 10^3} = \frac{675}{675} = 1$$

$$\begin{array}{rcl} \frac{3}{\boxed{8}} & = & \frac{3}{2^3} \quad \frac{(11)_2}{\boxed{2^3}} \quad \frac{(11)_2}{2^2} \\ & & = (0.11)_2 \\ & = & (0.011)_2 \rightarrow (0.11)_2 \times 2^{-1} \end{array}$$

$$\frac{1}{8} = \frac{1}{2^3}$$

$$= (0.101)_2$$

$$\frac{15}{64} = \frac{(1111)_2}{2^6} \quad \frac{(1111)_2}{2^1} = (11.1)_2$$

$$= (0.100111)_2$$

④  $\beta = 2, m = 9, -100 \leq e \leq 100$   
 Using IEEE normalized form -

a) Find machine epsilon  $\epsilon_M$

$$\epsilon = \frac{|f(x) - x|}{|x|}$$

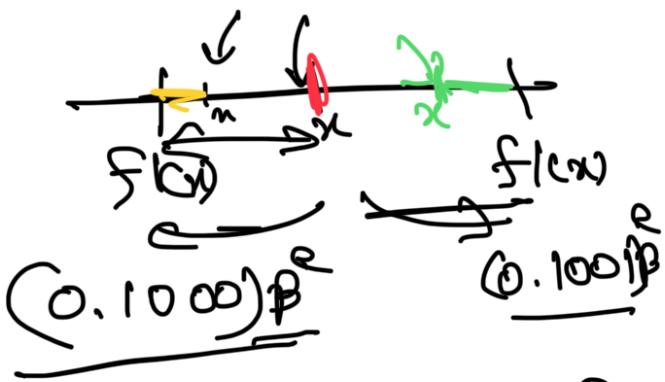
$$\frac{\epsilon_{\max}}{\epsilon_M} = \frac{R.E_{\max}}{|x|_{\min}} = \frac{\frac{1}{2} \beta^{1-m}}{\frac{1}{2} \beta^{1-4}} = \frac{\frac{1}{2} \beta^{-3}}{\frac{1}{2} \beta^{-3}} = \frac{1}{16}$$

$$(R.E)_{\max} = \frac{\frac{1}{2} \beta^{-m}}{\beta^{-1}} \boxed{\beta^e}$$

$$|x|_{\min} = \frac{\beta^{-1}}{\beta^{-m}} \boxed{\beta^e} = 2^{-1} \times 2^{-100}$$

$$= 2^{-101}$$

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



$$\underline{(0.1001) \beta^e} - \underline{(0.1000) \beta^e}$$

$$= \underline{\underline{(0.0001) \beta^e}} - \underline{\underline{\beta^{-1} \beta^e}} \rightarrow \underline{\frac{1}{2} \beta^{-1} \beta^e}$$

$$\textcircled{5} \quad x^2 - 12x + 5 \leq 0$$

Compute the roots of the eqn keeping

Four significant figures  $\Rightarrow$

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

$$= \frac{12 \pm \sqrt{144 - 20}}{2}$$

$$= \frac{12 \pm 2\sqrt{31}}{2} = 6 + \sqrt{31} = x_1$$

$$6 - \sqrt{31} = x_2$$

1 2 3 4

$$x_1 = 11.5677 = 11.57$$

Non-zero 5.f

$$x_2 = 0.43223 \underset{1234}{=} \boxed{0.4322}$$

Toy computer  $\rightarrow$

first calculate  $\sqrt{31}$

$$= \boxed{5.567} \times$$

$x_1 = 6 + 5.567$

$$= \underline{11.567} = \underline{11.57}$$

$x_2 = \underline{6 - 5.567} = 0 \boxed{.4330}$

Loss of significance

~~0.4320~~

$$x^2 - (\underline{\alpha} + \underline{\beta})x + \underline{\alpha\beta} = 0$$

$$\alpha\beta = 5$$

$$\alpha = \underline{11.57} \quad , \quad \beta = \frac{5}{\alpha}$$

$$= \frac{5}{11.57}$$

$$= 0.43215$$

$\boxed{= 0.4322}$