

## Tuto 1 - Corrige'

## Exercise 2 :

(1) Represent 3.14 and -8.625 in IEEE 754 F.P.S.P.

$$3.14 = 3+14.$$

$$3 = (11)_2, \quad 0.14 = (0.001)_2, \quad \Rightarrow (3.14) = (11.001)_2$$

$$\text{Normalization} = 1.1001 \times 2^7.$$

$$e = 1 \implies 1 + 127 = 128 = (1000\ 0000)_2$$

$S = 0$  (positive).

SD  $3.14 = [0|10000000|100100\ldots]_2$

$$-8 - 62s = -8 + 62s \text{ , with } s=1 \text{ (Negative)}$$

$$8 = (1000)_2, \quad 0.625 = (0.101)_2 \rightarrow 8.625 = (1000.101)_2$$

Normalization:  $1.000101 \times 2^3$

$$c = 3 \rightarrow 3 + 127 = 130 = (10000010)_2$$

50	-8.625	=	<table border="1"> <tr> <td>1</td><td>10000010</td><td>0001010...6</td></tr> <tr> <td>8 bits</td><td>23 bits</td><td></td></tr> </table>	1	10000010	0001010...6	8 bits	23 bits	
1	10000010	0001010...6							
8 bits	23 bits								

(2) 35.5 in JEEE 754 D.P

$$35 = (100011)_2, \quad 0.5 = (0.1)_2 \Rightarrow 35.5 = (100011.1)_2$$

$$\text{Normalization} : 1.000111 \times 2^5. , S = 0$$

$$e = 5 \rightarrow 5 + 1023 - 1028 = (10000000100)$$

$$SB \quad 35.5 = \boxed{0 \mid 10000000\ 100 \mid 00011100 \dots 0}$$

11 bits .      52 bits .

(3) floating-point representation of A:

$$A = \begin{matrix} & \text{8 bits} \\ \begin{matrix} s & | & 100000 & 1000110011 & 000000000000000000000000 \\ e & | & & & M \end{matrix} \end{matrix}$$

$S = 0$  : positive

$$10000100 = (132)_{10} \Rightarrow e = 132 - 127 = 5$$

$$\text{Manhissa} = 1.0110011$$

$$so A = 1.0110011 \times 2^5 = 0.10110011 \times 2^6 = (44.95)_{10}$$

$$(4) a, k = 1.11 \times 2^e \quad B = 1.001 \times 2^f = 0.1001 \times 2^{f+1}$$

$A+B$ : we add the mantissas :

$$\begin{array}{r} 1.11 \\ + 0.1001 \\ \hline 10.0101 \end{array}$$

$$A + B = 10.0101 \times 2^2 = 1.00101 \times 2^3$$

$A + B$    $C = 10101010$

$$b. \quad A = 1.11 \times 2^4, \quad B = 1.0011 \times 2^3$$

We multiply mantissas, and we add exponents

$$\begin{array}{r}
 & 10011 \\
 \times & 1.11 \\
 \hline
 & 10011 \\
 + & 10011 \\
 \hline
 & 10011
 \end{array}$$

$$1 \times B = 10,000.101 \times 2^1 = 1.0000101x$$

$$A \times B = \boxed{0} \quad \boxed{1\ 000\ 00101} \quad \boxed{000\ 010100\ldots} \quad 0$$

8 bits                    93 bits

5) Interval of representable Numbers using IEEE 754 S.P.

We know that  $-128 \leq e \leq 127$  and  $1 \leq M \leq 2$  (1)

50

- 1 -

18

e

3

1

1

1

1

1

-128

1

1

1

19

### Exercise 3.

- 1) (a)  $6.7 \cdot 10^{-2}$ , (b)  $0.00173 : 3$ , (c)  $2.30 \times 10^{-3} : 3$   
 (d)  $3.0054 : 5$ , (e)  $0.0004000 : 4$ .

(2) (a)  $6.234$  using 2 S.D : 6.2.

(b)  $0.006928$  using 3 S.D :  $6.34 \times 10^{-3}$ .

(f)  $238.62$  using 3 S.D : 239 (with rounding).

(d)  $6.345$  using 2 S.D : 6.3.

3)  $\epsilon_{\text{machine}} = 10^{-6}, -20 < \epsilon < 20.$

$$(a) 1 + 10^{-6} = (0.1 + 0.000001) \times 10^7 = 0.1 \times 10^7 = 10^6 \text{ (this is (e))}$$

$$(b) 1 + 10^{-4} = (0.1 + 0.00001) \times 10^5 = 0.1 \times 10001 \times 10^5 = 10001.$$

$$(c) 1 + 10^{-6} = 1 + 0.000001 = 1. \text{ (this is (g))}$$

$$(d) 10^3 + 10^6 = 1000000$$

$$(e) \frac{10^6}{10^{-12}} = 10^{6+12} = 10^{18} \quad (18 < 20).$$

$$(f) 10^{-9} \times 10^{-16} = 10^{-9-16} = 10^{-25} = 0 \quad (-25 < -20).$$

$$(g) 10^{28} + 10^4 = \text{Inf} \quad (28 > 20, 10^{28} = \text{Inf}).$$

$$(h) \frac{10^5}{10^{-21}} = 10^{5+21} = 10^{26} = \text{Inf}.$$

$$(i) \sqrt{10^0 - 10} \Rightarrow 10^0 - 10 = 10^0, \sqrt{10^0} = 10^0.$$

$$(j) \ln(10^{-25}), -25 < -20 \text{ the } 10^{-25} = 0, \text{ so } \ln(0) = \text{NaN}.$$

### Exercise 4:

- (1) Absolute true error and Percent Relative error when approximating  $\pi$  with (we take 5 sign. digits in the results).

$$(a) 3 : e_A = |\pi - 3| = 0.14159, e_R \% = \frac{0.14159}{\pi} \times 100 = 4.519\%$$

$$(b) 3.14 : e_A = |\pi - 3.14| = 0.0015926, e_R \% = \frac{e_A}{\pi} \times 100 = 0.050\%$$

$$(c) \frac{22}{7} : e_A = \left| \pi - \frac{22}{7} \right| = 0.0012645, e_R \% = 0.0443\%.$$

We observe that  $\frac{22}{7}$  is the best approximation because it gives the smallest error.

$$(2) A = 4\pi r^2. \text{ the impact means } e_A = \Delta A$$

(b) Five approximations used to calculate  $A$  in a machine:

① the earth is supposed Spherical, Impact: Difficult to find.

② the value of  $r$  is based on measures, so it's approximated

$$\text{Impact: } \Delta A = 2\pi^2 r \Delta r.$$

③ the value of  $\pi$  is not exact, Impact:  $\Delta A = 4r^2 \Delta \pi$

④ the numbers  $A$ ,  $\pi$ ,  $r$  are stored with errors (it depends on the machine properties) Impact: Difficult to find.

⑤ floating-point Arithmetic errors (Rounding and chopping).

Impact: Difficult to find.

(b) the radius is known with precision 2% means that the relative percent error is 2% and this means percent relative error  $e_R = \frac{\Delta r}{r} \times 100 = 2$ .

we find  $e_A = \Delta r$  (the relative error in the radius).

$$\Delta r = \frac{2r}{100}$$

and we want to find the relative percent error in  $A = \frac{\Delta A}{A}$   
we know that:

$$\Delta A = 4\pi^2 r \Delta r \Rightarrow 4\pi^2 r \cdot \frac{2r}{100} = \frac{4\pi r^2}{25}$$

$$\text{so } \frac{\Delta A}{A} \times 100 = \frac{4\pi r^2}{25} \times \frac{1}{4\pi r^2} \times 100 = 4\%$$

### Exercise(s)

$$\begin{aligned} (1) 122 \times (333 \times 695) &= (0.122 \times 10^3 \times (0.333 \times 10^3 + 0.695 \times 10^3)) \\ &= (0.122 \times 10^3) \times (0.028 \times 10^3) \quad \cancel{+ 0.03 \times 10^3} \\ &= 0.122 \times 10^3 \times 0.103 \times 10^4 \\ &= 0.012566 \times 10^7 = \boxed{0.126 \times 10^6}. \end{aligned}$$

$$\begin{aligned} (122 \times 333) + (122 \times 695) &= (0.122 \times 10^3 \times 0.333 \times 10^3) + (0.122 \times 10^3 \times 0.695 \times 10^3) \\ &= (0.040626 \times 10^6) + (0.08479 \times 10^6) \\ &= 0.0406 \times 10^6 + 0.0848 \times 10^6 \\ &= (0.406 \times 10^5) + 0.848 \times 10^5 \\ &= 1.254 \times 10^5. \\ &= 1.254 \times 10^5 - \boxed{0.125 \times 10^6} \end{aligned}$$

We observe a slight difference in the result.

$\hat{t} \times 5 : (2)$

(a).  $\cos(1.57079) \approx 0.000000632679 = \cos_{ex}$ .

$\cos(1.57079) \approx 0.00001632679 = \cos_{app}$ .

$$\text{error} \% = \left| \frac{\cos_{app} - \cos_{ex}}{\cos_{ex}} \right| \times 100 = 100\%$$

$$\begin{aligned}
 (2) \quad e_{abs} &= \cos(x+h) - \cos x = -2 \sin\left(\frac{x+h}{2}\right) \sin\left(\frac{h}{2}\right) \\
 &= -2 \sin\left(\frac{2x+h}{2}\right) \sin\left(\frac{h}{2}\right) = -2 \sin\left(x + \frac{h}{2}\right) \sin\left(\frac{h}{2}\right) \\
 x &\approx \frac{\pi}{2} \rightarrow \sin\left(x + \frac{h}{2}\right) \approx 1, h \ll 1 \rightarrow \sin\left(\frac{h}{2}\right) \approx \frac{h}{2}.
 \end{aligned}$$

$$e_{abs} \approx -2 \cdot 1 \cdot \frac{h}{2} = -h.$$

$$(3) \quad e_{rel} = \frac{e_{abs}}{\cos_{ex}} \approx \frac{-h}{\cos x} \approx \frac{-h \sin x}{\cos x} = -h \tan x \gg 1 \approx \infty.$$

We conclude that the relative error is large and at the same time the absolute error is small so we conclude that rel error is large when the ex value is small, has no units.

relative error: how occurs to the measurement is compared to the actual value (or exact value).

Or the measur of precision.

Absolute is the distance between ex and app.

Conclusion: Relative error only makes sense when the exact value not zero or close to zero.