

## TP n° 02 – Représentation des nombres

### I Dépassement de capacité

**Question 1 :** Écrire dans le tableau suivant le float (simple précision) le plus grand exprimable.

| S | Exposant | Mantisse |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
|---|----------|----------|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|--|
|   |          |          |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |

**Question 2 :** Calculer sa valeur dans la base décimale.

**Question 3 :** Calculer la valeur dans la base décimale du float (double précision) le plus grand exprimable.

**Question 4 :** Entrer dans une console python les commandes suivantes et en déduire si le type de float utilisé est de précision simple ou double.

```
>>> 2.0**(1023)
>>> 2.0**(1024)
>>> import sys
>>> sys.float_info.max
```

**Question 5 :** Calculer la valeur minimale ( $> 0$ ) pour le type de variable utilisé par votre système.

**Question 6 :** Vérifier votre calcul grâce à la commande suivante.

```
>>> import sys
>>> sys.float_info.min
```

### II Approximation de calcul

Lors de la première séance, nous avons remarqué que certains calculs étaient approximatifs. L'exemple suivant avait été utilisé.

```
>>> 1-1/3.-1/3.-1/3.
```

**Question 7 :** Calculer le nombre binaire permettant de définir le réel le plus proche de  $1/3$ .

**Question 8 :** Écrire ce nombre sous la forme suivante  $A * 2^{exp}$ , où  $A$  est un entier, et  $exp$ , l'exposant le plus petit qui permet à  $A$  d'être un entier. Vous ferez ce calcul pour un float simple et un float double. (Le calcul sera plus simple si une similitude entre les deux calculs est trouvée)

**Question 9 :** Déterminer dans ces deux cas la valeur du décimal le plus proche de  $1/3$ , dans le cas des deux types de float.

**Question 10 :** Calculer alors la valeur de l'opération  $1 - 1/3. - 1/3. - 1/3.$  en prenant ce nombre approché et comparer cette valeur à celle trouvée en faisant le calcul directement dans la console python.

### III Correction

#### III.1 Dépassement de capacité

Question 1 :

| S | Exposant |   |   |   |   |   |   | Mantisse |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |  |  |  |
|---|----------|---|---|---|---|---|---|----------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|--|--|--|
| 0 | 1        | 1 | 1 | 1 | 1 | 1 | 0 | 1        | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |   |  |  |  |
|   | 7        |   |   | F |   |   |   | 7        |   |   | F |   |   |   | F |   |   |   | F |   |   |   | F |   | F |  |  |  |

Question 2 :

| Signe | Exposant         | Mantisse         |
|-------|------------------|------------------|
| 0     | <u>111...110</u> | <u>111...111</u> |
|       | 8bits            | 23bits           |

- Exposant :  $2^8 - 2 = 254$ , exposant simple  $254 - 127 = 127$ ,
- Le chiffre à calculer est donc 111...111000...000,
- Ce qui donne en décimal  $(2^{24} - 1) * 2^{104} = 3.4028234663852886.10^{38}$

Question 3 :

| Signe | Exposant         | Mantisse         |
|-------|------------------|------------------|
| 0     | <u>111...110</u> | <u>111...111</u> |
|       | 11bits           | 52bits           |

- Exposant :  $2^{11} - 2 = 2046$ , exposant simple  $2046 - 1023 = 1023$ ,
- Le chiffre à calculer est donc 111...111000...000,
- Ce qui donne en décimal  $(2^{53} - 1) * 2^{971} = 1.7976931348623157.10^{308}$

Question 5 :

| Signe | Exposant         | Mantisse         |
|-------|------------------|------------------|
| 0     | <u>000...001</u> | <u>000...000</u> |
|       | 11bits           | 52bits           |

- Exposant : 1, exposant simple  $1 - 1023 = -1022$ ,
- Le chiffre à calculer est donc 1,000...000,
- Ce qui donne en décimal  $1 * 2^{-1022} = 2.2250738585072014.10^{-308}$

$$\begin{aligned}
 0,33.. \times 2 &= 0,66.. = 0 + 0,66.. \\
 0,66.. \times 2 &= 1,33.. = 1 + 0,33.. \\
 0,33.. \times 2 &= 0,66.. = 0 + 0,66.. \\
 0,66.. \times 2 &= 1,33.. = 1 + 0,33.. \\
 0,33.. \times 2 &= 0,66.. = 0 + 0,66.. \\
 &\dots
 \end{aligned}$$

On remarque une récurrence dans l'écriture du  $0,33_{10}$  en binaire :  $0,33_{10} = 0,01010_{2}$

Question 8 :

- 32 bits :  $1, \underbrace{0101..}_{23bits} * 2^{-2} = \underbrace{10101010....}_{24bits} * 2^{-25}$
- 64 bits :  $1, \underbrace{0101..}_{52bits} * 2^{-2} = \underbrace{10101010....}_{53bits} * 2^{-54}$

Le passage de 10101010...0 à 10101010...1 se fait par un décalage des bits vers la droite ce qui revient à diviser par deux.

Le passage de 10101010...1 à 10101010...0 se fait par un décalage des bits vers la gauche ce qui revient à multiplier par deux.

**Question 9 :**

$$\text{— 32 bits : } \underbrace{10101010\dots0}_{24\text{bits}} = \underbrace{1111111\dots}_{24\text{bits}} - \underbrace{10101010\dots1}_{23\text{bits}}$$

$$\text{— 64 bits : } \underbrace{10101010\dots1}_{53\text{bits}} = \underbrace{1111111\dots}_{54\text{bits}} - \underbrace{10101010\dots0}_{54\text{bits}}$$

Le calcul de  $\underbrace{1111111\dots}_{24\text{bits}}$  se fait en ajoutant 1.

$$\text{— 32 bits : } \underbrace{10101010\dots0}_{24\text{bits}} = (2^{24} - 1) - \frac{\underbrace{10101010\dots0}_{24\text{bits}}}{2}$$

$$\text{— 64 bits : } \underbrace{10101010\dots1}_{53\text{bits}} = (2^{54} - 1) - 2 * \underbrace{10101010\dots1}_{53\text{bits}}$$

Regroupement des  $\underbrace{10101010\dots0}_{24\text{bits}}$  et  $\underbrace{10101010\dots1}_{53\text{bits}}$ .

$$\text{— 32 bits : } \underbrace{10101010\dots0}_{24\text{bits}} = \frac{2}{3} \cdot (2^{24} - 1)$$

$$\text{— 64 bits : } \underbrace{10101010\dots1}_{53\text{bits}} = \frac{1}{3} \cdot (2^{54} - 1)$$

Le résultat est alors.

$$\text{— 32 bits : } \underbrace{10101010\dots0}_{24\text{bits}} * 2^{-25} = \frac{2}{3} \cdot (2^{24} - 1) * 2^{-25} =$$

$$\text{— 64 bits : } \underbrace{10101010\dots1}_{53\text{bits}} * 2^{-54} = \frac{1}{3} \cdot (2^{54} - 1) * 2^{-54} =$$

**Question 10 :**

$$\text{— 32 bits : } 1 - 1/3 - 1/3 - 1/3 = 1 - 2 \cdot (2^{24} - 1) * 2^{-25} = 1 - (2^{24} - 1) * 2^{-24} = 2^{-24} = 5.960464477539063 \cdot 10^{-8}$$

$$\text{— 64 bits : } 1 - 1/3 - 1/3 - 1/3 = 1 - 1 \cdot (2^{54} - 1) * 2^{-54} = 1 - (2^{54} - 1) * 2^{-54} = 2^{-54} = 5.55111512312578 \cdot 10^{-17}$$