

# Leçon d'informatique : Représentation des flottants sur des mots de tailles finies

S. Benlhajlahsen - PCSI<sub>1</sub>



## Sommaire

I Nombres réels, nombres décimaux et flottants	1
II Représentation des flottants sur des mots de taille fixe.	2
III Précision des calculs en flottants	3

### Extrait du programme :

Notions	Commentaires
Distinction entre nombres réels, décimaux et flottants.	On montre sur des exemples l'impossibilité de représenter certains nombres réels ou décimaux dans un mot machine.
Représentation des flottants sur des mots de taille fixe. Notion de mantisse, d'exposant.	On signale la représentation de 0 mais on n'évoque pas les nombres dénormalisés, les infinis ni les NaN. Aucune connaissance liée à la norme IEEE-754 n'est au programme.
Précision des calculs en flottants.	On insiste sur les limites de précision dans le calcul avec des flottants, en particulier pour les comparaisons. Le comparatif des différents modes d'arrondi n'est pas au programme.

## I Nombres réels, nombres décimaux et flottants

### I.A définitions

**À retenir :** En mathématiques, on définit :

- l'ensemble  $\mathbb{Q}$  des rationnels.  $x \in \mathbb{Q}$  s'il peut se mettre sous la forme  $x = p/q$  avec  $(p,q) \in \mathbb{Z} \times \mathbb{Z}^*$ .  $\frac{1}{3}$  est un rationnel mais  $\frac{1}{\sqrt{2}}$  ne l'est pas.
- l'ensemble  $\mathbb{D}$  des nombres décimaux. Un nombre décimal est un nombre rationnel qui peut s'écrire sous la forme d'une fraction **dont le dénominateur est une puissance de 10**. Ainsi,  $2,5 = \frac{25}{10}$  est un nombre décimal et rationnel mais  $1/3$  n'est pas un nombre décimal.

**Remarque :** Un nombre décimal  $x$  pourra s'écrire sous la forme :

$$x = m \times 10^e$$

L'écriture scientifique correspond à l'unique écriture de  $x$  telle que  $m \in [1, 10[$ . On appellera alors *mantisse* le réel  $m$  et  $e$  sera l'exposant en base 10.

**Remarque :** En python, le flottant `x = 1.5e-4` ou `x = 1.5*10**(-4)` correspond au nombre décimal  $1,5 \times 10^{-4}$ .

**Nombre à virgule :** C'est un nombre dans lequel la partie entière est séparée de la partie décimale par une virgule. Les nombres à virgule sont les nombres réels écrits en notation décimale.

### I.B Écriture en base 2

Pour faire le parallèle avec les entiers, les nombres à virgules sont aussi en base 2.

**Idée :** Seuls les nombres qui s'écrivent en base 2 avec un nombre fini de chiffres après la virgule seront représentables en machine.

### Exemples :

- 9,0 est un nombre à virgule. On constate qu'il s'écrit :

$$9,0 = 1 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 0 \times 2^{-1} = \underline{1001,0}_2$$

- 1,375 est aussi un nombre à virgule. On constate qu'il s'écrit :

$$1,375 = 1 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3} = 1 \times 2^0 + 0 \times \frac{1}{2^1} + 1 \times \frac{1}{2^2} + 1 \times \frac{1}{2^3} = \underline{1,011}_2$$

- 2,25 est aussi un nombre à virgule. On constate qu'il s'écrit :

$$2,25 = 1 \times 2^1 + 0 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} = 1 \times 2^1 + 0 \times 2^0 + 0 \times \frac{1}{2^1} + 1 \times \frac{1}{2^2} = \underline{10,01}_2$$

**Écriture « scientifique » en base 2 :** Par analogie avec l'écriture scientifique en base 10, on propose que la mantisse  $m$  soit dans  $[1;2[$ . Ainsi, on pourra noter :

- $9,0 = \underline{1,0010}_2 \times 2^3$ ;
- $2,25 = \underline{1,001}_2 \times 2^1$ .

**Nombres décimaux en base 2 :** En base 10, les nombres décimaux s'écrivent comme une somme finie de nombre du type :

$$\dots 100, 10, 1, \frac{1}{10}, \frac{1}{100} \dots$$

En base 2, les nombres décimaux s'écrivent comme une somme finie de nombre du type :

$$\dots 4, 2, 1, \frac{1}{2}, \frac{1}{4} \dots$$

**Exemple :**  $\frac{1}{10} = 0,1$  est un nombre décimal en base 10 mais **n'est pas un nombre décimal en base 2**. Comme  $\frac{1}{2^3} = 0,125$  et  $\frac{1}{2^4} = 0,0625$ , on peut commencer par écrire  $0,1 = \frac{1}{2^4} + a$  avec  $a = 0,0375$ . Comme  $\frac{1}{2^5} = 0,03125$ , on peut ensuite écrire  $0,1 = \frac{1}{2^4} + \frac{1}{2^5} + b$  avec  $b = 0,00625$ . Cette décomposition amène à la série numérique :

$$0,1 = \left( \frac{1}{2^4} + \frac{1}{2^5} \right) \sum_{k=0}^{+\infty} \frac{1}{2^{4k}} = \underline{0,0001100110011001\dots}_2$$

**Méthode de décomposition en base 2** La méthode est la même que dans le chapitre précédent. On commence par encadrer le nombre  $x$  entre deux puissances de 2 :  $2^{p-1} \leq x < 2^p$  où  $p$  est un entier relatif. Tant que le reste n'est pas nul, on ajoute à l'écriture binaire les quotients de la division par  $2^{p-1}$  puis  $2^{p-2}$ , ...

## II Représentation des flottants sur des mots de taille fixe.

En 64 bits, l'encodage des flottants est donnée par la norme IEEE-754. Un flottant pourra s'écrire de manière générale :

$$a = \pm m \cdot 2^e$$

On voit donc qu'il faut stocker :

- le signe  $\pm$ ;
- l'information sur la mantisse  $m \in [1;2[$ ;
- l'information sur l'exposant  $e$ .

**Problèmes :** On voit tout de suite que cet norme amène à définir deux zéros !

L'exposant peut être positif pour les grand nombres et négatif pour les petits nombres. Ce problème a déjà été rencontrés pour les entiers signés. Il a été résolu par le complément à deux. Il sera ici résolu dans la norme IEEE-754 par un exposant biaisé.



### III.B Comparaison des flottants

Du fait qu'un nombre réel est codé par un `float` qui en est une valeur approchée, la comparaison des nombres peut être surprenante.

Par exemple `0.1+0.1+0.1-0.3 == 0` prend la valeur `False` car le calcul est fait à l'aide des valeurs arrondies de `0.1` et de `0.3`, le résultat de l'opération `0.1+0.1+0.1-0.3` est petit mais il est de l'ordre de l'erreur faite lors de l'arrondi de `0.3`.

|| **À retenir :** On ne fera donc pas de comparaison entre des flottants.  
C'est une difficulté qui peut apparaître dans un algorithme de recherche de racine d'une équations sur les flottants.

### III.C Comment opérer malgré tout une comparaison entre flottants ?

Si on choisit un flottant  $a$  et qu'on note  $b$  son suivant, la différence relative entre  $a$  et  $b$  est  $\frac{b-a}{\max(|a|,|b|)}$ , elle est de l'ordre de  $2^{-52}$  soit environ  $2 \times 10^{-16}$ . Cette différence relative ne dépend pas de la taille de  $a$ .

Plutôt que comparer directement deux nombres non nuls  $a$  et  $b$ , on pourra tester si leur différence relative est inférieure à un seuil choisi, par exemple  $10^{-9}$ .

On observe alors que `(0.1+0.1+0.1-0.3)/0.3 < 10**(-9)` est vrai alors que `(0.1+0.1+0.1-0.3)/0.3 < 10**(-16)` est faux.

Pour comparer un nombre  $a$  à 0, on cherchera par exemple si la valeur absolue de  $a$  est inférieure à un seuil choisi (qui peut être théoriquement choisi aussi petit que  $10^{-323}$ ).

|| **Remarque :** Si on veut tester si deux nombres  $a$  et  $b$  sont proches, le module `math` contient la fonction `isclose` qui teste si la différence relative entre  $a$  et  $b$  est inférieure à  $10^{-9}$  et renvoie `True` si c'est le cas.

### Conclusion

On retiendra la distinction entre nombres réels, décimaux et flottants. On remarquera que la représentation des flottants se fait, par analogie avec l'écriture scientifique à l'aide d'une mantisse et d'un exposant. Enfin, on n'oubliera pas que le calcul sur les flottants est soumis à des erreurs qui empêche la comparaison entre eux.