

$$4 + 2 + 1 + 1/2 + 1/8 = 7.625$$

Revenons sur 0,1 + 0,2

0,1 et 0,2 ont des notations décimales *finies* (ce sont des *décimaux*)

Leur notation *dyadique* n'est pas finie !

$$0,1 = (0,00011001100110011001100110011001100110011 \dots)_2$$

En machine elle est tronquée (mais sera très proche de 0,1)

Ce n'est *généralement* pas gênant : on n'a généralement pas besoin d'une telle précision.

Cette approche est intéressante et naïvement, on pourrait penser que la machine stocke ainsi ses nombres.

Problème :

comment manipuler des nombres très grands et des nombres très petits en même temps ?

La taille de l'univers d'un côté, la masse d'un atome de l'autre : il faudrait des milliers de chiffres.

La notation scientifique

$$A = 300000000 \times 0.00000015$$

La notation décimale *n'est pas adaptée*.

On préfère la *notation scientifique* :

$$A = (3 \times 10^8) \times (1.5 \times 10^{-7})$$

Souvenons nous

- on *multiplie 3 et 1,5*
- on *ajoute les exposants 8 et -7*

$$A = (3 \times 1.5) \times 10^{8-7}$$

$$A = 4.5 \times 10^1$$

$$A = 45$$

La machine procède de la même manière en base 2.

Nombre dyadique

Un **nombre dyadique** est s'écrit :

$$\pm (1, b_1 \dots b_k)_2 \times 2^e$$

où b_1, \dots, b_k sont des bits et e est un entier relatif.

La suite de bits $b_1 \dots b_k$ est la *mantisse* du nombre,

La puissance de 2 est *l'exposant* du nombre.

Exemple

$$6,25 = (110,01)_2 = (1,1001)_2 \times 2^2$$

- La mantisse est la suite 1 0 0 1
- L'exposant est 2

Nombres à virgule flottante

Dans cette norme, les nombres dyadiques sont codés sur 64 bits en réservant :

- 1 bit pour le signe ;
- 11 bits pour l'exposant ;
- 52 bits pour la mantisse.

s	e	m
-----	-----	-----

Ce sont *les nombres à virgule flottante*.

Amplitude

Sans entrer dans les détails, en codant sur 64 bits on peut représenter des nombres entre :

- $2^{-1022} \approx 2,23 \times 10^{-308}$ pour le plus petit et
- $2^{1024} - 2^{971} \approx 1,80 \times 10^{308}$ pour le plus grand

Des améliorations sont faites pour les nombres très proches de 0.

Quand un flottant dépasse le plus grand nombre possible il est considéré comme *infini*

```
>>> 2.0 * 10**308 # dépasse le plus grand
inf
```

Quelques surprises avec inf

`inf` se comporte “grosso modo” comme l’infini des mathématiques...

mais l’implémentation révèle quelques surprises :

```
>>> a = float('inf') # pour définir inf
>>> a
inf
>>> -a
-inf # - infini
>>> a + a
inf
>>> a - a # opération interdite
nan # not a number
>>> a + a == a
True
>>> b = 2.0 * 10 ** 309 # b = inf
>>> c = 2 * 10 ** 1000 # un integer
>>> c > b # inf est plus grand que tous les nombres
False
```

Attention donc, les comparaisons entre grands entiers et grands flottants ne sont pas correctes mathématiquement parlant. Il faut absolument les éviter.

Deux problèmes dans les calculs avec les flottants

Absorption

```
>>> (1. + 2.**53) - 2.**53 # = 1
0.0 # 1 a été absorbé par l'enorme nombre 2**53
>>> 2.**53 - 2.**53 + 1 # on change l'ordre...
1 # et ça fonctionne
```

Annulation

Soustraire deux nombres proches fait perdre de la précision

```
>>> a = 2.**53 + 1
>>> b = 2.**53
>>> a - b
0.0
```

Il peut y avoir des conséquences

Les calculs avec des flottants engendrent toujours des erreurs qu’il est possible d’éviter en limitant leur quantité et les répétitions.

- Le 25 février 1991, à Dhara en Arabie Saoudite, un missile Patriot américain a raté l’interception d’un missile Scud irakien, ce dernier provoquant la mort de 28 personnes. L’enquête a mis en évidence le défaut suivant :
- L’horloge interne du missile mesure le temps en 1/10s. Ce nombre *n’est pas dyadique* et est converti avec une erreur d’environ 0,000000095s

- Le missile a été mis en route 100h avant son lancement, ce qui entraîne un décalage de

$$0,000000095 \times 100 \times 3600 \times 10 \approx 0,34s.$$

- C'est assez pour qu'il rate sa cible.

Source