



REPRÉSENTATION DES NOMBRES RÉELS

[Frédéric PEURIERE - Marion SZPIEG]

Connâître la représentation approximative des nombres réels (à virgule flottante) et en mesurer les conséquences (capacités et arrondis)

1. Petits tests en Python...

Que fait la commande suivante en Python, si a et b sont deux nombres `>>> a - b == 0` ?

Taper puis noter ce que renvoie Python :

```
>>> 4 - 4 == 0
```

```
>>> 4 - 4.0 == 0
```

```
>>> 0.2 + 0.8 - 1 == 0
```

```
>>> 0.2 + 0.6 - 0.8 == 0
```

```
>>> 0.1 + 0.5 - 0.6 == 0
```

```
>>> 0.1 + 0.2 - 0.3 == 0
```

```
>>> 0.7 + 0.2 - 0.9 == 0
```

Êtes vous surpris ? Pourquoi ?

Chercher d'autres cas surprenants et écrivez en quelques uns ici :

Afin de remédier à cela, nous allons effectuer un test du type $|a-b| < \epsilon$ où ϵ est une valeur proche de zéro à définir. Taper la commande suivante :

```
>>> abs(0.1 + 0.2 - 0.3) < 1E-5 # test adapté aux flottants
```

Quel résultat est affiché ? Quelle information cela donne t il ?

Créer un programme en Python afin de coder le test ci-dessous. Votre programme contiendra une fonction **égal**(a, b, y, eps) qui prendra en argument 4 nombres : a et b dont on fera la somme qu'on comparera avec y le tout avec une précision) eps près (10^{-5} dans l'exemple précédent).

Tester ce programme avec $a=0.1$, $b=0.2$ et $y=0.3$ et différentes valeurs de eps , et trouver à partir de quelle précision (sous la forme d'une puissance de 10) le test renvoie False.

2. Codage des nombres réels en virgule fixe

2.1. Le principe

En base 10, l'expression 652,375 est une manière abrégée d'écrire :

On fait exactement pareil pour la base 2. Ainsi, l'expression 110,101 signifie :

2.2. Passer de la base 2 à la base 10

Pour convertir un nombre réel de la base 2 vers la base 10, on le décompose à l'aide de puissances de 2 avec des exposants positifs et négatifs. Par exemple :

110,101 =

Remarques sur l'exposant :

- Les valeurs extrêmes 0 et 255 pour n sont des valeurs réservées (voir plus loin)
- Dans la norme IEEE 754, il n'a pas été choisi la méthode du complément à 2 pour représenter l'exposant. En mémoire, le nombre n est un nombre entier positif. Or, nous avons besoin aussi d'avoir des exposants négatifs pour les nombres inférieurs à 1. Le nombre n appartenant à l'intervalle $[1,254]$, on choisit un décalage de 127 afin que l'exposant réel $n-d$ appartienne à l'intervalle $[-126,127]$.

Méthode pour encoder un nombre en virgule flottante en simple précision

Convertir 5,25 en binaire selon la norme IEEE 754, en simple précision

étape 1 : écrire le nombre **non signé** en binaire (virgule fixe)

étape 2 : écrire le nombre en « notation scientifique binaire »

étape 3 : calculer s , f et n :

-
- la fraction f de la mantisse est ce qui se trouve derrière la virgule , donc $f=.....$
- l'exposant $n-127=...$ donc $n=$

étape 4 : reprendre les éléments précédents : 5,25 sera représenté en simple précision par :

Méthode pour convertir en base 10 un nombre en virgule flottante en simple précision

Inversement, retrouver le nombre en base 10 qui s'écrit comme ci-dessous en simple précision :

1 0 1 1 1 1 0 1 0 1 0 1 1 0

-
- $n=$
- la mantisse est : $m=$
- le nombre cherché est donc :

Remarque : voici un tableau avec l'encodage des valeurs spéciales :

Bit de signe s	Exposant n	Fraction f	Valeur spéciale
0	0	0	
1	0	0	
0	1111111	0	
1	1111111	0	
0	1111111	pas 0	

Petit récap'

Encodage	Signe s	Exposant n	Fraction f	Écriture selon la norme IEEE 754
32 bits (simple précision)				
64 bits (double précision)				

La 2^e ligne est donnée à titre d'information.

3.3. Conclusion

- Les nombres décimaux et réels n'ont pas tous un codage exact en machine.
- Les calculs et représentations sont nécessairement arrondis et la propagation des erreurs d'arrondi est une problématique délicate.
- Travailler quand on peut avec des entiers plutôt qu'avec des décimaux.
- Il faudra être prudent sur les tests (par exemple , ne pas tester si un nombre flottant $A=0$, mais plutôt tester si $A < 10^{-10}$)
- Les résultats dépendent de propagation d'arrondis faits par la machine.
- On évitera d'additionner deux quantités dont l'écart relatif est très important
- On évitera de soustraire deux quantités très proches

3. A vos machines

Connectez vous alors à <https://baseconvert.com/ieee-754-floating-point>

1. Que pouvez vous dire de la représentation de 0 ?.....
2. a) Copier la représentation décimale de 0.1 en 32 bits (on peut faire de même en 64 bits) et coller la dans une calculatrice (par exemple <https://web2.0calc.fr/>)

b) Additionner de même la représentation de 0.2. Enfin comparez cette valeur avec celle de 0.3 (en faisant pas exemple une soustraction)

c) Expliquez ce que vous avez trouvé dans la partie 1 (avec le programme)
3. Testez et observez aussi les exemples proposés en bas de page

decimal: -0.1

decimal: 1e+100

decimal: NaN