

Calcul numérique sur ordinateur

Christophe Jermann*

Journée académique de l'IREM PdL - 13/04/2023



* exposé conçu en lien avec un fascicule éponyme, réalisé avec Frédéric Goualard, disponible sur le site de l'IREM

Au cœur d'outils pédagogiques (et du quotidien)...



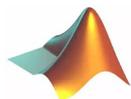
Java, javascript,
C/C++, ...



GeoGebra



Maple™



MATLAB

LibreOffice

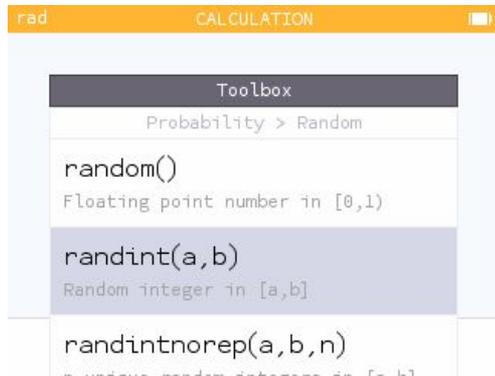


Google
Sheets

Excel

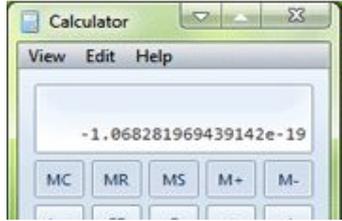
... à la fiabilité questionable...

Exemple : résultat de cent millions de lancers de dé à 6 faces au moyen de la calculatrice Numworks¹ ... quelque chose cloche ?



[1] résultat similaire reproductible avec Scratch, Snap, et d'autres outils et langages de programmation.

... avec des impacts importants (voire catastrophiques) !



Calculatrice Windows
(2001-2018)



Missiles Patriot
(1991)



Intel Pentium IV
(1994)

PayPal		Statement period:
Balance Summary		June 1, 2013 - June 30, 2013
		USD
Beginning Balance		140.25
Ending Balance		-92,233,720,368,547,800.00

Paypal
(2013)



Compteur de vues Youtube
(2014)



Premier vol Ariane 5
(1996)



Bug de l'an 2000
et bientôt, 2038 !

Pourquoi cela ?

Calcul mathématique = ensembles (et nombres) **infinis**

VS

Ordinateur = machine **finie**

Exploiter *au mieux* les possibilités des ordinateurs = **coder** l'information (et les opérations) aussi **efficacement** que possible (espace mémoire, temps de calcul)

⇒ Le codage numérique n'a pas les propriétés mathématiques des ensembles (partiellement) représentés

Le calcul numérique sur ordinateur² obéit à ses propres règles !

[2] ou calculatrice, tablette, smartphone, ...

Le langage des ordinateurs

Langage des ordinateurs = le **binaire**

1 bit (**binary digit**) = 2 valeurs, **0** et **1**.

Chaîne binaire = séquence de bits (**exemple : 00100111**)

Unité typique d'information = octet = 8 bits

⇒ Toute donnée informatique (nombre, texte, image,...) est une chaîne binaire

Exemple : que représente 00110001 00110011 01000010 00110010 ?

Quatre entiers (49 51 66 50) ? Un entier (825 442 866) ? Une date Unix (27/02/1996

17:41:06) ? Un réel ($\sim 2.608 \times 10^{-9}$) ? Quatre caractères ASCII (1 3 A 2) ? Deux

caractères UTF16 (𐀀(自)) ?

Codage “numérique”

Pour **interpréter** une chaîne binaire, il faut un **codage**, souvent associé à une **taille**.

- Codage : interprétation d’une chaîne binaire comme un nombre, un texte, une image, ...
- Taille : nombre de bits (ou d’octets) à considérer pour (dé)coder une chaîne binaire
 - ⚠ Taille \Rightarrow nombre de valeurs possibles (exemples : 8 bits \rightarrow 256 chaînes binaires distinctes ; 16 bits \rightarrow 65536 chaînes ; N bits \rightarrow 2^N chaînes)

Exemples :

- texte : ASCII (8bits), UTF-8/16/32 (8/16/32 bits), ...
- image : RGB (24 bits/pixel), RGBA (32 bits/pixel), ...

Nombres : le problème de base

En France (et dans de nombreux pays), les nombres sont notés au moyen des chiffres arabes (0,1,2,3,4,5,6,7,8,9) selon une **notation positionnelle en base 10**

Exemples :

- $6438 = 6 \times 10^3 + 4 \times 10^2 + 3 \times 10^1 + 8 \times 10^0$
- $-1,4142... = -1 \times 10^0 + 4 \times 10^{-1} + 1 \times 10^{-2} + 4 \times 10^{-3} + 2 \times 10^{-4} + ...$

Pour faire du calcul numérique sur ordinateur, il faut donc **convertir la notation décimale en binaire** (base 2), mais cela induit **des arrondis** :

- dec→bin : certaines notations décimales finies sont codées en notations binaires trop grandes pour l'ordinateur, voire infinies (*on en reparlera*)
- bin→dec : l'affichage décimal d'un nombre binaire stocké dans l'ordinateur peut être inexact (tronqué/arrondi)

⇒ Il ne faut pas croire tout ce que l'on voit à l'écran !

Les entiers

Entiers mathématiques VS informatiques

En mathématiques :

- Naturels (\mathbb{N}) ou Relatifs (\mathbb{Z}), ensembles infinis dénombrables
- Munis des opérations arithmétiques usuelles (+, -, ×, ÷)
 - Particularités : la soustraction doit rester dans le domaine de définition pour les Naturels ; la division fournit deux résultats, un quotient et un reste

En informatique :

- type signé ou non-signé, de taille fixe (e.g., 32 bits) ou variable (*≠ infinie*)
⇒ ensembles **finis**, codages **multiples**
- Munis des opérations arithmétiques par extension du **calcul booléen**

Un premier calcul sur entiers

On veut convertir en kilomètres la distance du Soleil à l'Étoile polaire (Alpha Ursæ Minoris) qui vaut approximativement 431 années-lumière.

- Mathématiquement :
 $431 \text{ al} \times 300\,000 \text{ km}\cdot\text{s}^{-1} \times 365 \text{ j} \times 24 \text{ h} \times 60 \text{ min} \times 60 \text{ s} = 4\,077\,604\,800\,000\,000 \text{ km}$
- Informatiquement :

Scratch

A Scratch script showing a sequence of multiplication blocks. Each block starts with 'mettre distance à' followed by a multiplication operation. The final block shows the result '4077604800000000'.

```
mettre distance à 431
mettre distance à distance * 300000
mettre distance à distance * 365
mettre distance à distance * 24
mettre distance à distance * 60
mettre distance à distance * 60
distance 4077604800000000
```

Numworks

The Numworks calculator interface showing the calculation: $431 \times 300\,000 \times 365 \times 24 \times 60 \times 60 = 4.0776048E9$.

LibreOffice

al	km/s	j	h	m	s	km
431	300000	365	24	60	60	4077604800000000

Python

```
>>> 431*300000*365*24*60*60
4077604800000000
```

C/C++

```
int a12km(int d) {
    return d * 300000 * 365 * 24 * 60 * 60;
}
code$ g++ distance.cpp -o distance.exe
code$ ./distance.exe
1503883264 !?!
```

Un autre !

On veut calculer le nombre de donnes possibles à la belote, c-à-d. le nombre de permutations d'un paquet de 32 cartes.

- Mathématiquement : la fonction factorielle fait cela très bien ! :-)
 $32! = 263\,130\,836\,933\,693\,530\,167\,218\,012\,160\,000\,000$ ($\sim 2.63 \times 10^{35}$)
- Informatiquement :

Scratch



The Scratch code consists of several blocks: a 'mettre' block for 'n' to 0, another 'mettre' block for 'n!' to 1, a 'répéter' block for 32 iterations, an 'ajouter' block for 1 to 'n', and a 'mettre' block for 'n!' to 'n * n!'. Below the code, the variable 'n' is set to 32 and 'n!' is shown as 2.631308369336935e+35. Red boxes highlight the 'n!' value and the 'e+35' exponent, with '?!?' written next to them.

Numworks



The Numworks calculator shows calculations for 13!, 20!, and 32!. The result for 32! is 2.631308369336935e+35. Red boxes highlight the scientific notation and the exponent, with '?!?' written below.

LibreOffice

N	N!
10	3628800
20	2432902008176640000
32	2631308369336935301672180121600000000

Python

```
>>> from math import factorial
>>> factorial(32)
263130836933693530167218012160000000
```

C/C++

```
int factorial(int N) {
    if (N<2) return 1;
    else return N*factorial(N-1);
}
exposés ./fact.exe
12 : 479001600
13 : 1932053504
32 : -2147483648
```

Mais qu'est-ce qui se passe ?

Quelques éléments d'explication avant d'entrer dans le détail :

- Python, Numworks : représentation des entiers à **taille variable**
 - + représentation exacte de très grands entiers ; opérations exactes
 - espace mémoire et temps de calcul importants
- Scratch, LibreOffice (*et de nombreux autres outils*) : **pas de représentation des entiers**, utilisation de “réels” \pm proprement arrondis à l’affichage
- C/C++ (*et Numpy, et beaucoup de langages de programmation*) : représentation des entiers à **taille fixe**
 - + représentation exacte ; espace mémoire et temps de calcul optimisés
 - bornes sur les entiers représentables assez *petites* ; calcul à maîtriser

Entiers non-signés à taille fixe

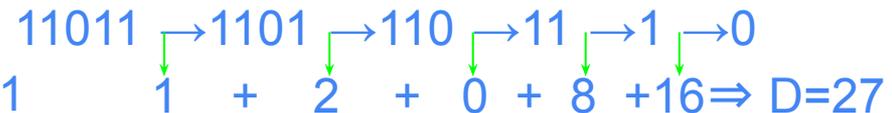
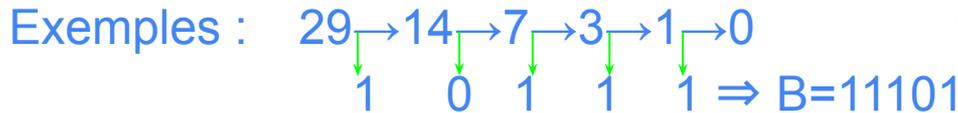
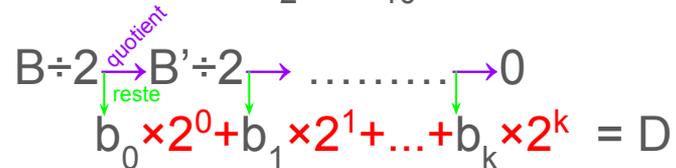
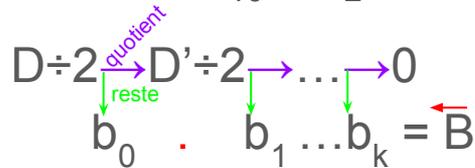
- Tailles usuelles = 8, 16, **32**, 64, 128 bits
 \Rightarrow borne max = 255, 65 535, **4 294 967 295**, 18 446 744 073 709 551 615
 ($<2 \times 10^{19}$), 340 282 366 920 938 463 463 374 607 431 768 211 455 ($<4 \times 10^{38}$)
! Nombres pas si “astronomiques” !!! (cf. exemples précédents)

- Codage “naturel” = positionnel en base 2

Codage $D_{10} \rightarrow B_2$

Décodage $B_2 \rightarrow D_{10}$

Principe



Remarque : $k \geq \text{taille} \Rightarrow \text{non représentable}$

Opérations sur ce codage binaire “naturel”

Addition, soustraction, multiplication et division réalisables directement sur le code binaire.

Exemple : addition bit à bit *avec retenues*

$b + b' + \text{retenue précédente} = \text{somme} + \text{retenue suivante}$

b	b'	p	s	r
0	0	0	0	0
1	0	0	1	0
0	1	0	1	0
1	1	0	0	1
0	0	1	1	0
1	0	1	0	1
0	1	1	0	1
1	1	1	1	1

Illustration :

```
  1 0 1 0
+  1 0 1 1
+ 1 0 1 0 ← retenues
= 1 0 1 0 1
```

(en décimal : $10 + 11 = 21$)

Exception lors d'un calcul

Si le résultat dépasse la borne de taille, on parle de **dépassement de capacité** (*overflow*) : le résultat obtenu est incorrect (*mais l'ordinateur mémorise la situation*)

Exemples : sur 4 bits

1 1 1 1

+ 0 0 0 1

1 1 1 1 (*retenues*)

= 0 0 0 0

0 0 0 0

- 0 0 0 1

1 1 1 1 (*retenues*)

= 1 1 1 1

1 0 1 1

× 0 0 1 1

1 0 1 1

1 0 1 1

= 0 0 0 1

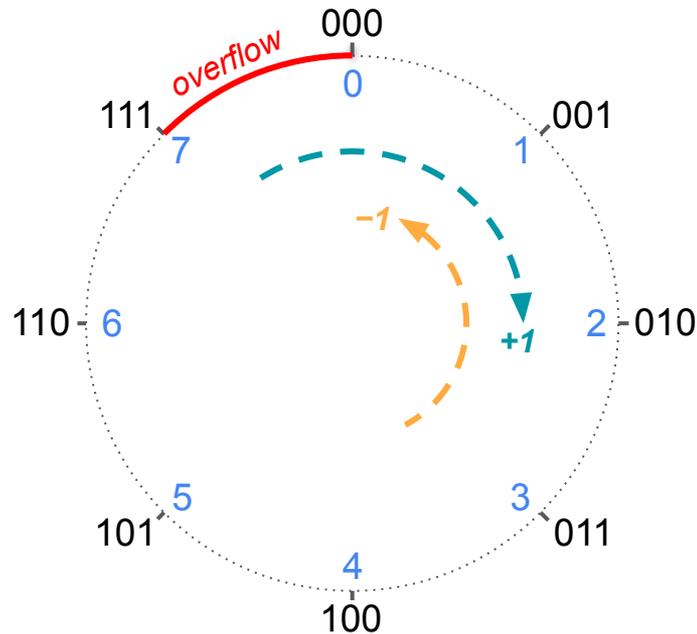
À l'écran : $15 + 1 = 0$

$0 - 1 = 15$

$11 \times 3 = 1$

Une arithmétique modulaire !

Le comportement des entiers non-signés peut-être visualisé de façon modulaire :

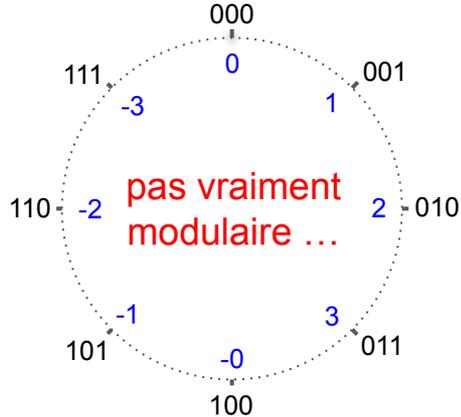


Et les signés ?

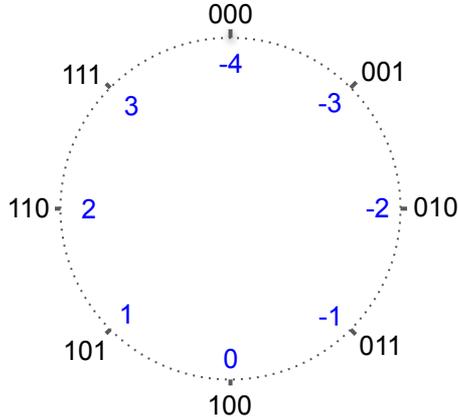
Plusieurs codages possibles :

- bit de signe (*0=positif, 1=négatif*) + codage “naturel”
- codage “naturel” + biais pair ou impair
- complément à 1, ou à 2

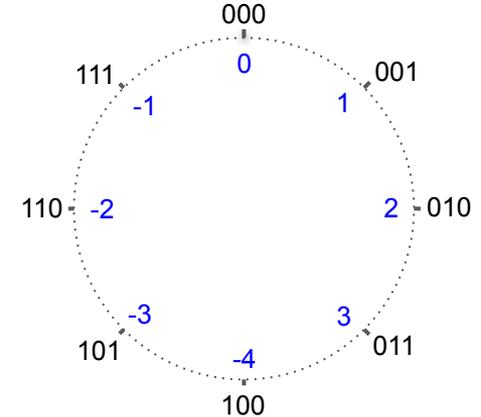
avec bit de signe



avec biais pair



en complément à 2



⚠ Dans tous les cas, | borne inférieure | \approx borne supérieure $\approx 2^k$ pour une taille k

Complément à 2 : propriétés et opérations

- Codage efficace :
 - L'opposé d'un nombre binaire B est $\beta+1$ (β inverse binaire ($0 \leftrightarrow 1$) de B)
 - ⚠ L'opposé du plus petit nombre (codé $100\dots 0_2$) n'est pas représentable
 - Le bit de poids fort donne le signe
 - Zéro représenté de façon unique
 - Nombre maximum de valeurs représentées
- Opérations efficaces :
 - addition et soustraction “naturelles”
 - multiplication quasi-“naturelle” (*vérification de signe*)
 - une division binaire efficace est possible

```
>>> from numpy import int8
>>> int8(127)+int8(1)
-128
>>> int8(-128)-int8(1)
127
>>> abs(int8(-128))
-128
```

Retour sur les exemples

- Entiers de taille fixe (*C/C++*) : l'overflow explique le résultat surprenant.
C'est plutôt la norme (pour des raisons d'efficacité) ⇒ Attention !
- Entiers de taille variable (*Numworks, Python*) : les problèmes sont identiques, mais repoussés à des résultats faramineux (*e.g., 2^{8192} sur Numworks*)
Réservé à des usages particuliers car inefficace
- “Réels” (*Scratch, LibreOffice*) : le nombre de chiffres représentable conditionne l'exactitude du résultat.
Mathématiquement $\mathbb{N} \subset \mathbb{Z} \subset \mathbb{R}$, mais informatiquement non

Les réels

Réels mathématiques VS informatiques

En mathématiques :

- Décimaux (\mathbb{D}), Rationnels (\mathbb{Q}) ensembles infinis dénombrables
- Réels (\mathbb{R}) ensemble infini indénombrable (*irrationnels* : π , $\sqrt{2}$, ...)
- Munis des opérations arithmétiques usuelles (+, -, ×, ÷) et de fonctions élémentaires (racines, logarithmes, trigonométrie,...)

En informatique :

- Virgule fixe ou flottante, de taille fixe (e.g., 32 bits)
⇒ ensembles **finis**, codages **multiples**
- Opérations arithmétiques (*usuellement*) par extension du calcul booléen (*hardware*), fonctions élémentaires par approximation (*software*)

Tracé de polygones

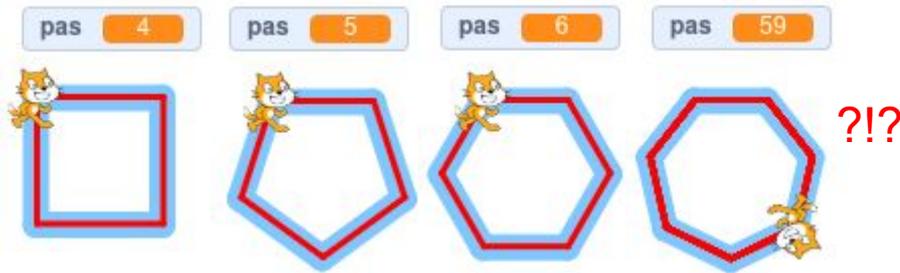
On veut tracer un polygone à N côtés selon deux méthodes :

- “répéter X fois” : aller au premier sommet ; Répéter N fois : tracer un côté, tourner de $360/N$ degrés ;
- “répéter jusqu’à” : aller au premier sommet ; tracer le premier côté ; Répéter jusqu’à être revenu au premier sommet : tracer un côté, tourner de $360/N$ degrés ;

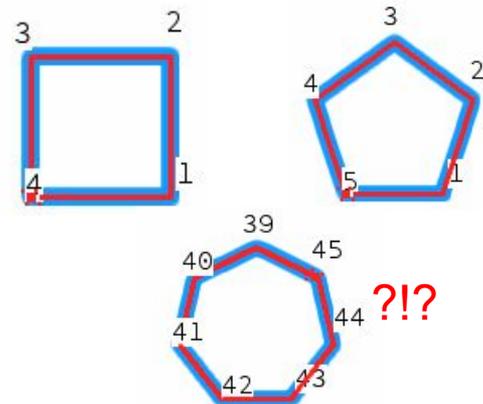
Mathématiquement : c’est équivalent

Informatiquement : pas toujours...

Scratch



Numworks/Turtle



Triangle manquant d'aire

On veut calculer l'aire A d'un triangle isocèle de côtés a, b, c qui se "dégonfle", devenant de plus en plus plat ($a=b, c \rightarrow 0$) en utilisant la formule de Héron :

$$A = \sqrt{s \times (s-a) \times (s-b) \times (s-c)} \text{ avec } s = (a+b+c)/2$$

Mathématiquement : étant donné les propriétés du triangle, on peut réécrire la formule générale* en $A = c/4 \times \sqrt{4a^2 - c^2}$, formulation rigoureusement équivalente

Informatiquement : on calcule avec les deux formules pour $a=b=10^6$ et $c=10^{-k}$

LibreOffice

k	0	1	2	3	4	5	6	7	8	9	10
Héron	500000	50000	5000	500	50,000031	5,0000381	0,5000038	0,0499422	0,0050059	0	0
Simplifiée	500000	50000	5000	500	50	5	0,5	0,05	0,005	0,0005	5E-05

Python

```

exposé$ python triangle.py
k      Héron      Simplifiée
0      499999.9999999375  499999.9999999375
1      50000.00004656606  49999.99999999994
9      0.0004656612873077392  0.0005000000000000001
10     0.0         5.0000000000000016e-05
    
```

* Kahan, W. (2014). *Miscalculating area and angles of a needle-like triangle*.

Sprint en série

Si Achille couvre à chaque instant la moitié de la distance le séparant de la ligne d'arrivée, combien de temps son sprint durera-t-il ?

Mathématiquement : Cette série bien connue converge vers 1 ... à l'infini !

$u_n = u_{n-1}/2, u_0 = 1/2$: suite géométrique de raison $1/2 \Rightarrow$ terme général : $u_n = 1/2^{n+1}$

somme partielle : $S_n = \sum_{k=0..n} 1/2^{k+1} = 1/2 (1 - 1/2^{n+1}) / (1 - 1/2) = 1 - 1/2^n < 1$ pour tout n ; $\lim_{n \rightarrow \infty} S_n = 1$

Informatiquement :

Scratch

The Scratch script shows an initialization phase where variables 'x', 'pas', and 'nb' are set to 0, 0.5, and 0 respectively. An iteration loop follows, where 'nb' is incremented by 1, 'x' is advanced by 'pas' (200 units), a 0.1 second wait is implemented, and both 'x' and 'pas' are updated: 'x' becomes 'x + pas' and 'pas' is divided by 2. The simulation shows the cat character moving towards a green flag at x=1, with the distance 'x' and step size 'pas' decreasing as 'nb' increases.

Python

```
>>> x, pas, nb = 0, 0.5, 0
>>> while x < 1:
...     x, pas, nb = x+pas, pas/2, nb+1
...
>>> print(x,pas,nb)
1.0 2.7755575615628914e-17 54
```

LibreOffice

nb	0	1	10	25	50	51
pas	0,5	0,25	0,00048828125	1,49011611938477E-08	4,4408920985006E-16	2,2204460492503E-16
x	0,00	0,50	0,9990234375	0,99999970197678	0,99999999999999	1

Début d'explication

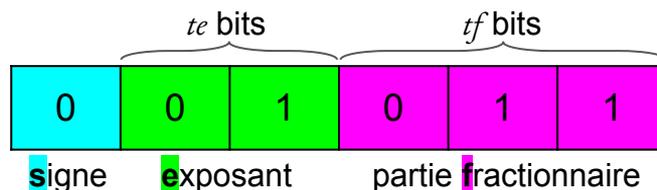
- Pour le calcul réel, les nombres à **virgule flottante** sont utilisés par la majorité des outils informatiques et langages de programmation.
- Ces nombres utilisent un codage standardisé en taille fixe K (e.g., $K=64$ bits) ;
⇒ ils ne peuvent représenter plus de 2^K **valeurs distinctes** (e.g., $<2 \times 10^{19}$) ;
⇒ toutes ces valeurs ont un **nombre fini de chiffres significatifs** (e.g., <17).
- Ce codage permet la représentation de nombres très petit (e.g., 10^{-300}) et très grands (e.g., 10^{300}) ;
⇒ les calculs doivent être réalisés sur des valeurs **homogènes** ;
⇒ les résultats doivent être **arrondis** pour être représentables.

Nombres à virgule flottante IEEE 754 (1/3)

Un standard pour le codage et les opérations : la **norme IEEE 754**

- Codage : deux paramètres, te =taille exposant, tf =taille partie fractionnaire

e = entier signé
avec biais impair
($biais = 2^{te-1} - 1$)



f = codage "naturel" en
puissances négatives
($bit\ unité\ implicite = 1$)

Exemple de décodage : $(-1)^0 \times 2^{(0 \times 2^1 + 1 \times 2^0 - 1)} \times (1 + 0 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3}) =$

1,375

Exemple de codage :

$$-3,5 \rightarrow \begin{cases} s = 1 \\ \text{signifiant} = 11,1 \rightarrow 1,11 \\ e = 10 \end{cases} \Rightarrow \text{flottant} = 1\ 10\ 110$$

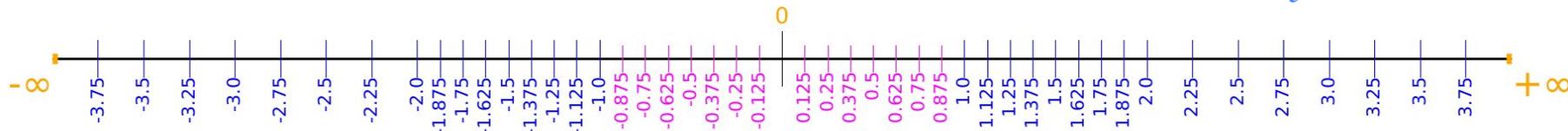
décalage=1 \Rightarrow exposant=1+biais = 2

Nombres à virgule flottante IEEE 754 (2/3)

Certains exposants représentent des valeurs particulières :

- $e = 0$ (valeur d'exposant signé = $-b$)
 - $f = 0 \Rightarrow$ représente la valeur 0 (± 0 selon signe)
 - $f \neq 0 \Rightarrow$ nombre "dénormalisé" : le bit implicite devant la partie fractionnaire est 0 et non plus 1 ; codage/décodage de f inchangé, mais e considéré comme valant 1.
- $e = 11\dots 1$ (valeur d'exposant signé = $b+1$)
 - $f = 0 \Rightarrow$ infini (valeur = $\pm\infty$ selon signe)
 - $f \neq 0 \Rightarrow$ NaN (Not a number)

Exemple : Valeur décimale de tous les flottants sur 6 bits avec $t_e=2$, $t_f=3$



soit $2 \times ((2^2 - 2) \times 2^3 + 2^3 + 1) - 1 = 49$ valeurs distinctes (sans Nan, $-0 = +0$) pour 64 codes

Nombres à virgule flottante IEEE 754 (3/3)

- Tailles standards : 32 bits (*simple* $te=8$, $te=23$), 64 bits (*double* $te=11$, $te=52$), 128 bits (*quadruple* $te=15$, $te=112$) + tailles intermédiaires pour les calculs
⚠ *Des tailles plus petites (8,16 bits), non standardisées, apparaissent aujourd'hui, notamment pour les applications en apprentissage profond*

⇒ Valeurs décimales minimales (non nulle) et maximales encodables :

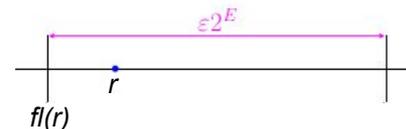
- simple : $\pm 1.4013 \times 10^{-45}$ – $\pm 3.40282 \times 10^{38}$
- double : $\pm 4.94066 \times 10^{-324}$ – $\pm 1.79769 \times 10^{308}$
- quadruple : $\pm 6.47518 \times 10^{-4966}$ – $\pm 1.18973 \times 10^{4932}$

Calcul IEEE 754

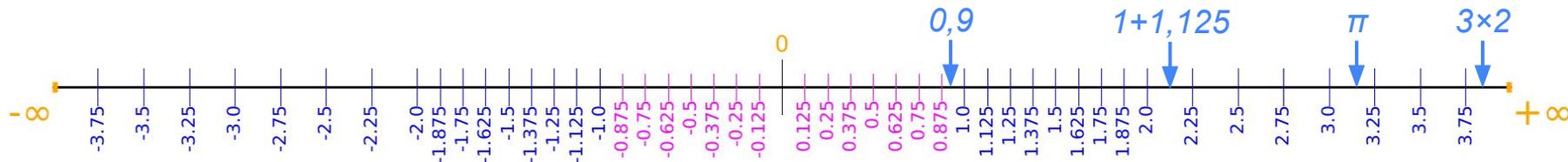
- Le standard impose les opérations “*arithmétiques*” sur les flottants (+, −, ×, ÷, mod, $\sqrt{\quad}$, ...) plus des opérations de comparaison (<, =) (*-0=0, NaN⇒faux*) et de manipulation (suivant/précédent, signe, ...)
- Il recommande des fonctions élémentaires (exponentiation, logarithmes, trigonométriques, ...)
- Toutes les opérations doivent traiter tous les flottants (NaN, dénormalisés, infinis, ...) (*opération “indéfinie” ⇒ NaN ou infini*)
- Le calcul flottant diffère cependant du calcul réel pour trois raisons principales, toutes dues à la finitude de représentation :
 - les erreurs d'**arrondi** ;
 - le phénomène d'**absorption** ;
 - le phénomène de ***cancellation***.

Arrondi

- Lors du codage ou d'un calcul, une valeur r peut être non représentable
⇒ il faut l'arrondir au flottant $fl(r)$ le plus proche,
i.e., vérifiant $|fl(r)-r| \leq \varepsilon 2^{E-1}$ ($\varepsilon = \text{précision machine} = 2^{-tf}$)



Exemples (sur 6 bits): $fl(\pi) = fl(11,00100100001\dots) = 0\ 11\ 101$ (valeur = $3,25_{10}$)
(underflow) $fl(0,9) = fl(0,11100110011\dots) = 0\ 00\ 111$ (valeur = $0,875_{10}$)
 $fl(1+1,125) = fl(10,001) = 0\ 10\ 000$ (valeur = $2,0_{10}$)
(overflow) $fl(3 \times 2) = fl(110,0) = 0\ 11\ 000$ (valeur = $+\infty$)



⚠ L'arrondi doit être pris après chaque opération flottante

Exemple : le calcul réel $\sin(a+b \times c) \rightarrow fl(\sin(fl(fl(a)+fl(fl(b) \times fl(c))))))$

Absorption et *cancellation*

- La méthode d'addition (et de soustraction) flottante consiste à aligner les exposants des opérandes
⇒ Perte de précision de l'opérande la plus petite, pouvant aller jusqu'à son annulation : **absorption** !

Exemple (sur 6 bits) : $\text{fl}(3+0.125) = \text{fl}(1.1 \times 2^1 + 0.001 \times 2^0) = \text{fl}(1.1 \times 2^1 + 0.000 \times 2^1) = 1.1 \times 2^1$ (valeur = 3.0_{10})

- La soustraction de flottants très proches peut annuler la plupart des bits de leurs parties fractionnaires
⇒ Perte de signification du résultat, composé de bits possiblement arrondis lors de précédents calculs : **cancellation** !

Exemple (sur 6 bits) : $\text{fl}(3+0.375-3) = \text{fl}(\text{fl}(1.1 \times 2^1 + 0.011 \times 2^0) - 1.1 \times 2^1) = \text{fl}(1.110 \times 2^1 - 1.1 \times 2^1) = 0.1 \times 2^0$
(valeur = 0.5_{10} , soit 33% d'écart au résultat réel !)

Retour sur les exemples

- Les arrondis expliquent les problèmes observés lors du tracé de polygone ($fl(\sum_{k=1}^7 fl(360 \div 7)) \neq 360$) et du tirage de dé (*Numworks/randint, p3 : sur 32 bits, $fl((2^{32}-k)/2^{32})=1$ pour $k=1..128 \Rightarrow 1$ chance sur 33×10^6 de tirer 7*).
- L'absorption explique le calcul d'aire qui s'annule subitement ($fl(10^6 + 10^6 + 10^{-10}) = 10^6$) et le sprint d'Achille trop tôt interrompu ($fl(1 - 2^{-53} + 2^{-54}) = 1$).
- La cancellation (*et les arrondis*) explique le calcul d'aire complètement erroné ($fl(10^6 + 10^{-10} - 10^6) \neq 1e-10$).

Tout est perdu ?

Le calcul flottant est entaché d'erreur vis-à-vis du calcul réel, doit-on jeter les supercalculateurs, tel Jean Zay à 36,85Pflop/s, à la poubelle ?

Non, car ce calcul est **predictible** et explicable !

Exemple : IEEE 754 impose des opérations arithmétiques arrondies au flottant près.

On peut ainsi prédire le degré d'approximation de l'évaluation d'un calcul C : on parle de **calcul d'erreur**, absolue ($|\text{fl}(C)-C|$) ou relative ($|\text{fl}(C)-C|/|C|$).

Exemple : erreur relative pour la somme S de 3 flottants a, b et c

Trois formulations : $a+(b+c)$; $b+(a+c)$; $c+(a+b)$

$$\text{fl}(a+\text{fl}(b+c)) = (a+(b+c) \times (1+\delta_1)) \times (1+\delta_2) \text{ avec } \delta_i \leq \varepsilon/2 (= \textit{precision machine} \div 2)$$

$$= S \times (1 + \delta_2 + (b+c)/S \delta_1 + (b+c)/S \delta_1 \delta_2)$$

\Rightarrow erreur relative $\approx \delta_2 + (b+c)/S \delta_1$ (*on néglige les erreurs d'ordre 2*)

Cela permet de choisir la formulation selon les magnitudes de $(a+b)$, $(b+c)$ et $(a+c)$.

Et au-delà des flottants ?

Si les performances ne sont pas un souci, il existe de nombreuses alternatives au calcul flottant :

- calcul rationnel (*e.g.*, *Python/fractions*)
- calcul décimal (généralement en précision *arbitraire*³) (*e.g.*, *Python/decimal*)
- calcul à virgule fixe ou flottante en précision *arbitraire*³ (*e.g.*, *Python/gmpy2*)

Il existe aussi des couches logicielles au-dessus du calcul flottant visant à le rendre plus rigoureux

- calcul stochastique (arrondis aléatoires)
- calcul par intervalles (encadrement garanti)

Domaine de recherche encore très actif !

[3] arbitraire≠infinie ... et souvent il faut fixer la précision cible avant le calcul

Conclusion

Mon exposé en 180 (*dixièmes de*) secondes

- L'incarnation physique des nombres en machine les rend obligatoirement finis.
- Le calcul informatique obéit à des règles précises, prédictibles, explicables, mais différentes du calcul mathématique.
- Cela concerne les entiers, les réels, et plus généralement toutes les familles de nombres dans tous les outils et sur tous les supports numériques.

Vous ne pourrez plus dire que vous n'étiez pas prévenus ! 😊

Des questions ?

Des réponses ! →

(et des éléments historiques, des exemples illustratifs, des détails sur les procédures de calcul de Python, Numworks, LibreOffice, ...)



Disponible sur le site de l'IREM PdL : <https://irem.univ-nantes.fr/>

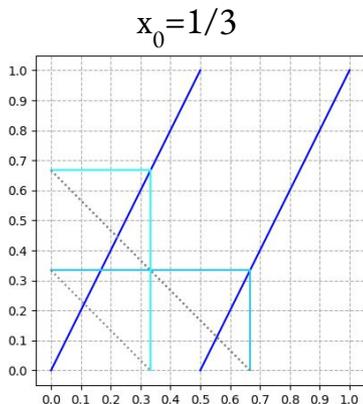
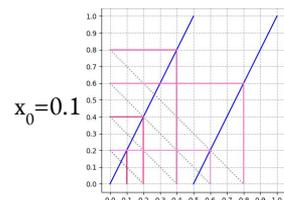
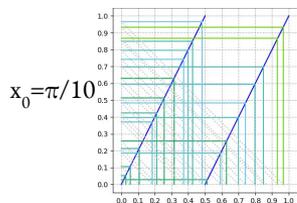
Bonus de fin de générique (de fin)

La fonction dyadique

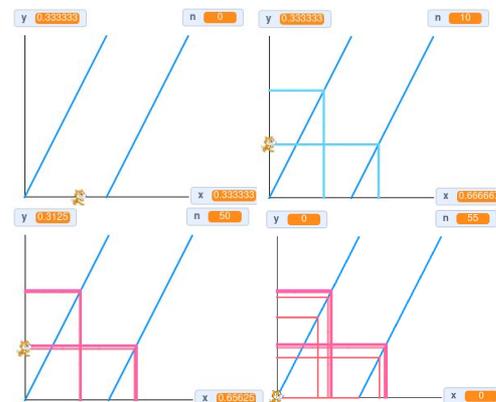
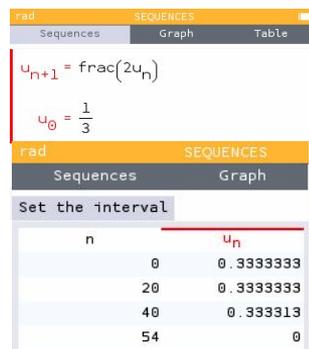
La fonction $f(x) = 2x \text{ mod } 1$ permet de définir un système dynamique $x_{k+1} = f(x_k)$, de condition initiale $x_0 \in [0,1[$, aussi appelé *décalage de Bernoulli*.

Mathématiquement :

- $x_0 \in \mathbb{R} \Rightarrow$ orbite apériodique
- $x_0 \in \mathbb{Q} \Rightarrow$ orbite périodique



Informatiquement :



```

code$ python DyadicMap.py
Using floats:
0 : 0.3333333333333333
20 : 0.3333333333139308
40 : 0.33331298828125
54 : 0.0
Using rationals:
0 : 1/3
20 : 1/3
40 : 1/3
54 : 1/3
    
```