

Manipulations des nombres à virgule flottante

Informatique pour tou(te)s - Semaines 3 & 4

Guillaume CONNAN

Lycée Clemenceau - MP/MP*

Dernière mise à jour : 11 janvier 2015 à 12:58

Sommaire

1 Préambule

2 La norme IEEE 754

3 Algèbre des nombres VF

4 Réels, arrondis et flottants

5 Que la force de l'erreur soit avec vous

Sommaire

1 Préambule

2 La norme IEEE 754

3 Algèbre des nombres VF

4 Réels, arrondis et flottants

5 Que la force de l'erreur soit avec vous











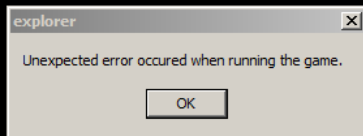
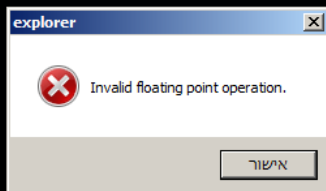
Microsoft Excel - Classeur1

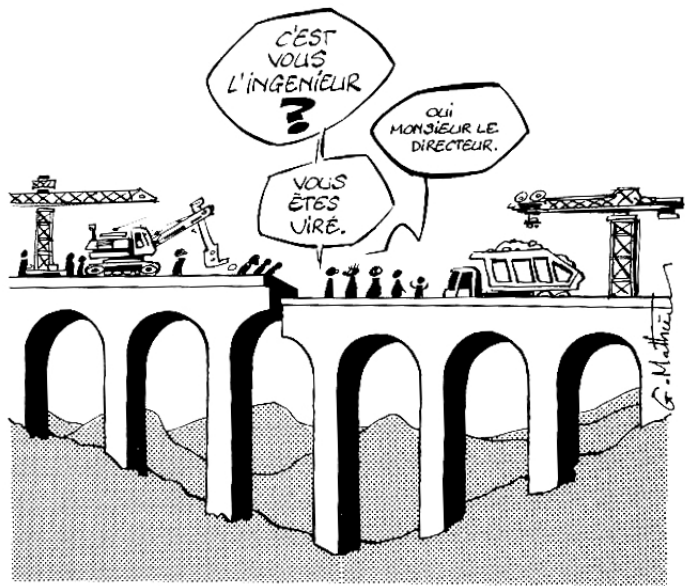
Fichier Edition Affichage Insertion Format Outils Do

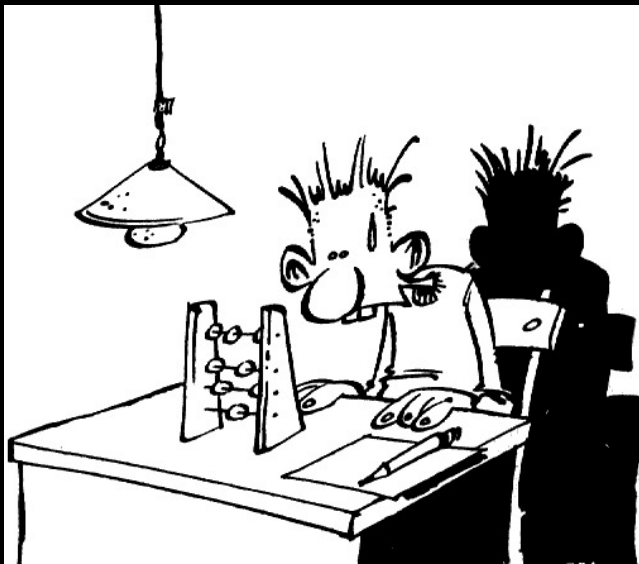
Arial 10 G I S

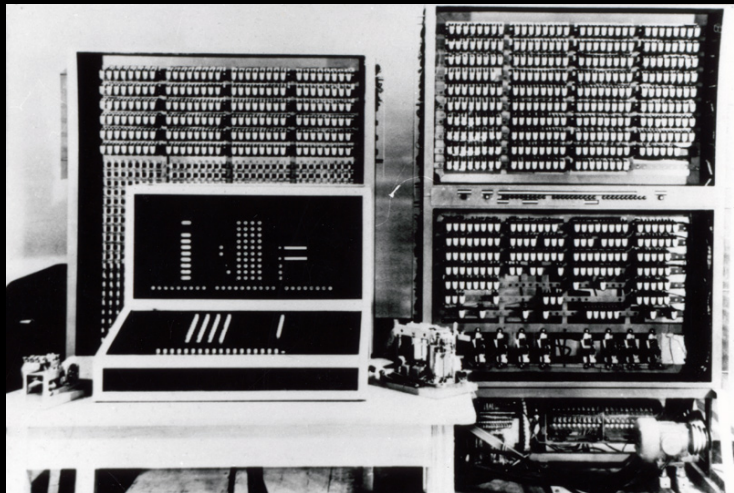
C15 fx

	A	B
1	$B1 = 4/3$	1,33333333333333000000
2	$B2 = B1 - 1$	0,33333333333333300000
3	$B3 = B2 * 3$	1,00000000000000000000
4	$B4 = B3 - 1$	0,00000000000000000000
5	$B5 = B4 * 2^{52}$	0,00000000000000000000
6	$(4/3 - 1) * 3 - 1$	0,00000000000000000000
7	$((4/3 - 1) * 3 - 1)$	-0,000000000000000022204
8	$((4/3 - 1) * 3 - 1) * 2^{52}$	-1,00000000000000000000











- char
- int

- char
- int $2147483647 + 1 = -2147483648...$

- char
- int $2147483647 + 1 = -2147483648...$

- char
- int $2147483647 + 1 = -2147483648...$
- short

- char
- int $2147483647 + 1 = -2147483648...$
- short long

- char
- int $2147483647 + 1 = -2147483648...$
- short long

- char
- int $2147483647 + 1 = -2147483648...$
- short long

• float

• double

- char
- int $2147483647 + 1 = -2147483648...$
- short long
- float
 - single
 - double
 - long double
 - 3

- char
- int $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3

- char
- int $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3

- char
- int $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3

- char
- int $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1

- char
- int $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1

- char
- int $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1

- char
- int $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1

- char
- int $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1

- char
- int $2147483647 + 1 = -2147483648...$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e0 0.3e1
- ...

```
In [1]: 3 * 0.1
```

-
-
-


```
In [1]: 3 * 0.1
```

```
Out[1]: 0.30000000000000004
```

```
▪
```

```
▪
```

```
In [1]: 3 * 0.1
```

```
Out[1]: 0.30000000000000004
```

```
In [2]: sum([0.1 for k in range(10000000)])
```

```
.
```

```
In [1]: 3 * 0.1
```

```
Out[1]: 0.30000000000000004
```

```
In [2]: sum([0.1 for k in range(10000000)])
```

```
Out[2]: 999999.9998389754
```

LES NOMBRES RÉELS
N'EXISTENT PAS.

TOUT CE QUE VOUS AVEZ VU
EN MATHS N'EST
QU'ILLUSION !

LES NOMBRES RÉELS
N'EXISTENT PAS.

TOUT CE QUE VOUS AVEZ VU
EN MATHS N'EST
QU'ILLUSION !

Sommaire

1 Préambule

2 La norme IEEE 754

3 Algèbre des nombres VF

4 Réels, arrondis et flottants

5 Que la force de l'erreur soit avec vous



$$v = (-1)^s \times m \times 2^E$$

- * un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon
- * 11 bits d'exposant E

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en français) m , $1 \leq m$.

• binary32

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.
- *binary32* $(\#E, \#m) = (8, 24)$

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.
- *binary32* $(\#E, \#m) = (8, 24)$

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en français) m , $1 \leq m$.
- *binary32* ($\#E, \#m$) = (8, 24)

• *binary64*

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.
- *binary32* $(\#E, \#m) = (8, 24)$
- *binary64* $(\#E, \#m) = (11, 53)$.

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.
- *binary32* $(\#E, \#m) = (8, 24)$
- *binary64* $(\#E, \#m) = (11, 53)$.

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.
- *binary32* ($\#E, \#m$) = (8, 24)
- *binary64* ($\#E, \#m$) = (11, 53).

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.

- *binary32* ($\#E, \#m$) = (8, 24)
- *binary64* ($\#E, \#m$) = (11, 53).
- *toy7* ($\#E, \#m$) = (3, 4).

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.

- *binary32* ($\#E, \#m$) = (8, 24)
- *binary64* ($\#E, \#m$) = (11, 53).
- *toy7* ($\#E, \#m$) = (3, 4).

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.
- *binary32* ($\#E, \#m$) = (8, 24)
- *binary64* ($\#E, \#m$) = (11, 53).
- *toy7* ($\#E, \#m$) = (3, 4).

$$v = (-1)^s \times m \times 2^E$$

- un bit de signe s qui vaut 0 si le nombre est positif, 1 sinon ;
- 11 bits d'exposant décalé e ;
- 53 bits de mantisse (ou « significande » en anglais) m , $1 \leq m$.

- *binary32* ($\#E, \#m$) = (8, 24)
- *binary64* ($\#E, \#m$) = (11, 53).
- *toy7* ($\#E, \#m$) = (3, 4).

Oups! $1 + 11 + 53 = 65$...En fait, on représente un nombre de \mathbb{V}_{64} sous forme normale :

$$v = (-1)^s \times 1, f \times 2^E$$

On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur $\#E$ bits on peut coder $2^{\#E}$ nombres.

La première moitié va donc de 0 à $2^{\#E-1} - 1$.

Elle correspond aux exposants réels de E_{\min} jusqu'à 0.

La deuxième de $2^{\#E-1}$ à $2^{\#E} - 1$.

Elle correspond aux exposants de 1 jusqu'à E_{\max} .

Il suffit donc de translater les exposants réels de $2^{\#E-1} - 1$.

$$E_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$

On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur $\#E$ bits on peut coder $2^{\#E}$ nombres.

La première moitié va donc de 0 à $2^{\#E-1} - 1$.

Elle correspond aux exposants réels de E_{\min} jusqu'à 0.

La deuxième de $2^{\#E-1}$ à $2^{\#E} - 1$.

Elle correspond aux exposants de 1 jusqu'à E_{\max} .

Il suffit donc de translater les exposants réels de $2^{\#E-1} - 1$.

$$E_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$

On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur $\#E$ bits on peut coder $2^{\#E}$ nombres.

La première moitié va donc de 0 à $2^{\#E-1} - 1$.

Elle correspond aux exposants réels de E_{\min} jusqu'à 0.

La deuxième de $2^{\#E-1}$ à $2^{\#E} - 1$.

Elle correspond aux exposants de 1 jusqu'à E_{\max} .

Il suffit donc de traduire les exposants réels de $2^{\#E-1}$ à $2^{\#E} - 1$ en

$$E_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1}$$

On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur $\#E$ bits on peut coder $2^{\#E}$ nombres.

La première moitié va donc de 0 à $2^{\#E-1} - 1$.

Elle correspond aux exposants réels de E_{\min} jusqu'à 0.

La deuxième de $2^{\#E-1}$ à $2^{\#E} - 1$.

Elle correspond aux exposants de 1 jusqu'à E_{\max} .

Il suffit donc de tradater les exposants réels de $2^{\#E-1} - 1$...

$$E_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$

On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur $\#E$ bits on peut coder $2^{\#E}$ nombres.

La première moitié va donc de 0 à $2^{\#E-1} - 1$.

Elle correspond aux exposants réels de E_{\min} jusqu'à 0.

La deuxième de $2^{\#E-1}$ à $2^{\#E} - 1$.

Elle correspond aux exposants de 1 jusqu'à E_{\max} .

Il suffit donc de traduire les exposants réels de $2^{\#E-1} - 1 \dots$

$$e_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$

On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur $\#E$ bits on peut coder $2^{\#E}$ nombres.

La première moitié va donc de 0 à $2^{\#E-1} - 1$.

Elle correspond aux exposants réels de E_{\min} jusqu'à 0.

La deuxième de $2^{\#E-1}$ à $2^{\#E} - 1$.

Elle correspond aux exposants de 1 jusqu'à E_{\max} .

Il suffit donc de traduire les exposants réels de $2^{\#E-1} - 1 \dots$

$$e_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$

On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur $\#E$ bits on peut coder $2^{\#E}$ nombres.

La première moitié va donc de 0 à $2^{\#E-1} - 1$.

Elle correspond aux exposants réels de E_{\min} jusqu'à 0.

La deuxième de $2^{\#E-1}$ à $2^{\#E} - 1$.

Elle correspond aux exposants de 1 jusqu'à E_{\max} .

Il suffit donc de traduire les exposants réels de $2^{\#E-1} - 1 \dots$

$$e_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$

On peut également gagner de la place en ne stockant pas le signe de l'exposant.

Sur $\#E$ bits on peut coder $2^{\#E}$ nombres.

La première moitié va donc de 0 à $2^{\#E-1} - 1$.

Elle correspond aux exposants réels de E_{\min} jusqu'à 0.

La deuxième de $2^{\#E-1}$ à $2^{\#E} - 1$.

Elle correspond aux exposants de 1 jusqu'à E_{\max} .

Il suffit donc de traduire les exposants réels de $2^{\#E-1} - 1 \dots$

$$e_{\text{stocké}} = E_{\text{réel}} + 2^{\#E-1} - 1$$

0,75 en toy7

0,75₁₀ en toy7.

$$0,75_{10} = \frac{3_{10}}{2^2_{10}} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \rightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)	f (3 bits)
0	0 1 0	1 0 0

0,75 en toy7

0,75₁₀ en toy7.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)	f (3 bits)
0	0 1 0	1 0 0

0,75 en toy7

0,75₁₀ en toy7.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0

0,75 en toy7

0,75₁₀ en toy7.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0

0,75 en toy7

0,75₁₀ en toy7.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0

0,75 en toy7

0,75₁₀ en toy7.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0

0,75 en toy7

0,75₁₀ en toy7.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0

0,75 en toy7

0,75₁₀ en toy7.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0

0,75 en toy7 : méthode générale

Comment écrire en base 2 la partie fractionnaire d'un nombre en base 10 ?

$$0,75 = b_1 \times 2^{-1} + b_2 \times 2^{-2} + \dots$$

$$2 \times 0,75 = b_1 + b_2 \times 2^{-1} + \dots$$

On voit poindre un bel algo...

0,75 en toy7 : méthode générale

Comment écrire en base 2 la partie fractionnaire d'un nombre en base 10 ?

$$0,75 = b_1 \times 2^{-1} + b_2 \times 2^{-2} + \dots$$

$$2 \times 0,75 = b_1 + b_2 \times 2^{-1} + \dots$$

On voit poindre un bel algo...

0,75 en toy7 : méthode générale

Comment écrire en base 2 la partie fractionnaire d'un nombre en base 10 ?

$$0,75 = b_1 \times 2^{-1} + b_2 \times 2^{-2} + \dots$$

$$2 \times 0,75 = b_1 + b_2 \times 2^{-1} + \dots$$

On voit poindre un bel algo...

0,75 en toy7 : méthode générale

Comment écrire en base 2 la partie fractionnaire d'un nombre en base 10 ?

$$0,75 = b_1 \times 2^{-1} + b_2 \times 2^{-2} + \dots$$

$$2 \times 0,75 = b_1 + b_2 \times 2^{-1} + \dots$$

On voit poindre un bel algo...

Exercice 1

Pourquoi le Patriot a raté son Scud de 600m ?

Représentez $1/10$ avec une mantisse de 24 bits et tronquez le résultat.

Calculez l'erreur en seconde puis après 100 heures d'utilisation. Sachant qu'un Scud vole à 1676 m s^{-1} , quelle est environ l'erreur commise en mètres ?

Et si on avait arrondi au plus proche ?

QUELQUES HUMANOÏDES RARES

**CENTRALIEN****NORMALIEN****ALIEN**

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

toy7, $\#E = 3$, $e_{\min} = 000$, $e_{\max} = 111 = 1000 - 1$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1 - 2^{3-1} + 1) - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = \left(e_{\min} - 2^{(\#E)-1} + 1 \right) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = \left(e_{\max} - 2^{(\#E)-1} + 1 \right) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = \left(e_{\min} - 2^{(\#E)-1} + 1 \right) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = \left(e_{\max} - 2^{(\#E)-1} + 1 \right) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = \left(e_{\min} - 2^{(\#E)-1} + 1 \right) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = \left(e_{\max} - 2^{(\#E)-1} + 1 \right) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = \left(e_{\min} - 2^{(\#E)-1} + 1 \right) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = \left(e_{\max} - 2^{(\#E)-1} + 1 \right) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = \left(e_{\min} - 2^{(\#E)-1} + 1 \right) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = \left(e_{\max} - 2^{(\#E)-1} + 1 \right) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

$$e_{\min} = 0$$

$$e_{\max} = 2^{\#E} - 1$$

$$e = E + 2^{(\#E)-1} - 1$$

$$\text{toy7, } \#E = 3, e_{\min} = 000, e_{\max} = 111 = 1000 - 1$$

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

$$E_{\min} = \left(e_{\min} - 2^{(\#E)-1} + 1 \right) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = \left(e_{\max} - 2^{(\#E)-1} + 1 \right) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

- zéro sous forme normale ?
- $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)} ?$

Sous-normaux

- zéro sous forme normale ?
- $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)} ?$

Sous-normaux...

- zéro sous forme normale ?
- $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)} ?$

Sous-normaux....

- zéro sous forme normale ?
- $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)} ?$

Sous-normaux...

- zéro sous forme normale ?
- $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)} ?$

Sous-normaux....

Exercice 2

En utilisant vos carreaux, représentez les nombres à virgule flottante normaux en toy_7 et rajoutez les sous-normaux.

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	0	0	0	0	0

s (1 bit)	e (3 bits)			f (3 bits)		
1	0	0	0	0	0	0

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	0	0	0	0	0

s (1 bit)	e (3 bits)			f (3 bits)		
1	0	0	0	0	0	0

```
In [1]: x = 1e500
```

```
•  
•  
•  
•  
•  
•  
•  
•  
•  
•
```

```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
•  
•  
•  
•  
•  
•  
•
```

```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
•  
•  
•  
•  
•  
•
```

```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
·  
·  
·  
·  
·
```

```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
Out[3]: inf
```

-
-
-
-

```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
Out[3]: inf
```

```
In [4]: 1/x
```

```
.
```

```
.
```

```
.
```

```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
Out[3]: inf
```

```
In [4]: 1/x
```

```
Out[4]: 0.0
```

```
.
```

```
.
```



```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
Out[3]: inf
```

```
In [4]: 1/x
```

```
Out[4]: 0.0
```

```
In [5]: 4 - x
```

```
.
```

```
In [1]: x = 1e500
```

```
In [2]: 1+x
```

```
Out[2]: inf
```

```
In [3]: x**3
```

```
Out[3]: inf
```

```
In [4]: 1/x
```

```
Out[4]: 0.0
```

```
In [5]: 4 - x
```

```
Out[5]: -inf
```

$$\lim_{x \rightarrow +\infty} x^2 - x$$

```
In [9]: x**2 - x
```

```
.
```

$$\lim_{x \rightarrow +\infty} x^2 - x$$

```
In [9]: x**2 - x
```

```
.
```

$$\lim_{x \rightarrow +\infty} x^2 - x$$

```
In [9]: x**2 - x  
Out[9]: nan
```

$$\lim_{x \rightarrow +\infty} x^2 - x$$

```
In [9]: x**2 - x
```

```
Out[9]: nan
```



```
In [10]: x
```

```
Out[10]: inf
```

```
In [11]: x - x
```

```
Out[11]: nan
```

```
In [12]: x / x
```

```
Out[12]: nan
```

```
In [13]: x - x == x - x
```

```
Out[13]: False
```

```
In [14]: x**2 - x
```

```
Out[14]: nan
```

```
In [15]: x * (x - 1)
```

```
Out[15]: inf
```

```
In [16]: x * (x - 1) == x**2 - x
```

```
Out[16]: False
```


Exercice 3

① *Comment expliquer les résultats suivants :*

```
1 *Main> let f(x) = x^2 / sqrt(x^3 + 1)
2 *Main> f(1e100)
3 1.0e50
4 *Main> f(1e150)
5 0.0
6 *Main> f(1e200)
7 NaN
```

② *Comment éviter le dernier NaN ?*

```
In [16]: f = lambda x: x**2 / sqrt(x**3 + 1)
```

```
In [17]: f(1e100)
```

```
Out[17]: 1e+50
```

```
In [18]: f(1e150)
```

```
-----  
OverflowError
```

```
Traceback (most recent
```

```
→ call last)
```

```
<ipython-input-18-79775d85f31a> in <module>()  
----> 1 f(1e150)
```

```
<ipython-input-16-eed0b81ef932> in <lambda>(x)
```

```
----> 1 f = lambda x: x**2 / sqrt(x**3 + 1)
```

```
OverflowError: (34, 'Numerical result out of range')
```

```
In [19]: f(1e200)
```

```
-----  
OverflowError
```

```
Traceback (most recent
```

```
In [53]: x = 1e-500
```

```
In [54]: x
```

```
Out[54]: 0.0
```

```
In [55]: x / x
```

```
-----  
ZeroDivisionError
```

```
Traceback (most recent
```

```
→ call last)
```

```
<ipython-input-55-fd52a7f8b5f1> in <module>()  
----> 1 x / x
```

```
ZeroDivisionError: float division by zero
```

```
In [56]: sqrt(-1.0)
```

```
-----  
ValueError
```

```
Traceback (most recent
```

```
  call last)
```

```
<ipython-input-56-d1c09f21b443> in <module>()  
----> 1 sqrt(-1.0)
```

```
ValueError: math domain error
```

```
*Main> let x = 1e500
```

```
*Main> x - x
```

```
NaN
```

```
*Main> x / x
```

```
NaN
```

```
*Main> sqrt(-1)
```

```
NaN
```

```
*Main> 0 / 0
```

```
NaN
```

USS Yorktown 1998



Warning
 Windows has categorized you as a government employee. Switching to sleep mode.
OK

Windows 98 Update Wizard
 Windows has detected that your computer is more than 12 months old. Windows 98 requires a newer machine to run properly. Please update your configuration.
OK

Random Error
 Windows 95 has detected a random error. This error occurs every once in a while. Please wait.

Keyboard not plugged
 Windows 95 was unable to detect your keyboard. Press F1 to retry or F2 to abort.

Printer Wizard
 The new and incredible 32bit intelligent wizard has obtained a solution to your printing problem: do not print.

Dialog box
 Sorry, Windows 95 was unable to comply with the "go to hell" command you specified.

Windows FAT
 Windows 98 has detected too much FAT already in your system. Please do not upgrade to FAT32 and stay with FAT16.
OK

Closing Windows 95
 You are about to exit Windows 95. Do you want to play another game?
OK Cancel

Internal Error
 Your windows has been running for 10h 37m 23s. Microsoft does not allow a windows system to run longer than that. That is why your computer will now crash.
OK

Air France 447 - 1^{er} juin 2009



s (1 bit)	e (3 bits)			f (3 bits)		
1	1	1	1	0	1	0

Arbre de lecture

s	e	f
-----	-----	-----

$e = 0\dots 0$	$f = 0$	\rightarrow	$v = (-1)^s \times 0.0$
$e = 0\dots 0$	$f \neq 0$	\rightarrow	$v = (-1)^s \times 0.f \times 2^{1-E_{\max}}$
$e = 1\dots 1$	$f = 0$	\rightarrow	$v = (-1)^s \times \infty$
$e = 1\dots 1$	$f \neq 0$	\rightarrow	NaN
$0\dots 0 < e < 1\dots 1$	$f \neq 0$	\rightarrow	$v = (-1)^s \times 1.f \times 2^{e-E_{\max}}$

Arbre de lecture

s	e	f
-----	-----	-----

$$e = 0\dots 0 \quad f = 0 \quad \rightarrow \quad v = (-1)^s \times 0.0$$

$$e = 0\dots 0 \quad f \neq 0 \quad \rightarrow \quad v = (-1)^s \times 0.f \times 2^{1-E_{\max}}$$

$$e = 1\dots 1 \quad f = 0 \quad \rightarrow \quad v = (-1)^s \times \infty$$

$$e = 1\dots 1 \quad f \neq 0 \quad \rightarrow \quad \text{NaN}$$

$$0\dots 0 < e < 1\dots 1 \quad f \neq 0 \quad \rightarrow \quad v = (-1)^s \times 1.f \times 2^{e-E_{\max}}$$

Successesseur

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$x = M(x) \times 2^{E(x)}$$

$$\text{succ}(x) = (M(x) + 1) \times 2^{E(x)} \quad (x > 0)$$

Successesseur

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$x = M(x) \times 2^{E(x)}$$

$$\text{succ}(x) = (M(x) + 1) \times 2^{E(x)} \quad (x > 0)$$

SuccessEUR

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$x = M(x) \times 2^{E(x)}$$

$$\text{succ}(x) = (M(x) + 1) \times 2^{E(x)} \quad (x > 0)$$

Successeur

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$x = M(x) \times 2^{E(x)}$$

$$\text{succ}(x) = (M(x) + 1) \times 2^{E(x)} \quad (x > 0)$$

SuccessEUR

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$v = M(v) \times 2^{E(v)-n}$$

$$\text{succ}(v) = (M(v) + 1) \times 2^{E(v)-n} = v + 2^{E(v)-n}$$

Successeur

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$v = M(v) \times 2^{E(v)-n}$$

$$\text{succ}(v) = (M(v) + 1) \times 2^{E(v)-n} = v + 2^{E(v)-n}$$

Successesseur

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$v = M(v) \times 2^{E(v)-n}$$

$$\text{succ}(v) = (M(v) + 1) \times 2^{E(v)-n} = v + 2^{E(v)-n}$$

Successeur

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$v = M(v) \times 2^{E(v)-n}$$

$$\text{succ}(v) = (M(v) + 1) \times 2^{E(v)-n} = v + 2^{E(v)-n}$$

Successeur

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$v = M(v) \times 2^{E(v)-n}$$

$$\text{succ}(v) = (M(v) + 1) \times 2^{E(v)-n} = v + 2^{E(v)-n}$$

Tableau récapitulatif

- On notera ε_m l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera λ le plus petit VF normal positif
- On notera μ le plus petit VF sous-normal positif
- On notera Ω le plus grand VF normal.

Quelle relation existe-t-il entre μ , ε_m et λ ?

Tableau récapitulatif

- On notera ε_m l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera λ le plus petit VF normal positif
- On notera μ le plus petit VF sous-normal positif
- On notera Ω le plus grand VF normal.

Quelle relation existe-t-il entre μ , λ et Ω ?

Tableau récapitulatif

- On notera ε_m l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera λ le plus petit VF normal positif
- On notera μ le plus petit VF sous-normal positif
- On notera Ω le plus grand VF normal.

Quelle relation existe-t-il entre μ , ε_m et λ ?

Tableau récapitulatif

- On notera ε_m l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera λ le plus petit VF normal positif
- On notera μ le plus petit VF sous-normal positif
- On notera Ω le plus grand VF normal.

Quelle relation existe-t-il entre μ , ε_m et λ ?

Tableau récapitulatif

- On notera ε_m l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera λ le plus petit VF normal positif
- On notera μ le plus petit VF sous-normal positif
- On notera Ω le plus grand VF normal.

Quelle relation existe-t-il entre μ , ε_m et λ ?

Tableau récapitulatif

- On notera ε_m l'*epsilon* de la machine, c'est-à-dire le successeur de 1
- On notera λ le plus petit VF normal positif
- On notera μ le plus petit VF sous-normal positif
- On notera Ω le plus grand VF normal.

Quelle relation existe-t-il entre μ , ε_m et λ ?

Nom	#E	#f	E_{\min}	E_{\max}	ε_m	λ	μ	Ω
<i>toy7</i>	3	3	-2	3	$2^{-3} = 1/8$	$2^{-2} = 1/4$	$2^{-5} = 1/32$	$1,111 \times 2^3 = 15$
<i>bin32</i>	8	23	-126	127	2^{-23} $\approx 1,2 \times 10^{-7}$	2^{-126} $\approx 1,2 \times 10^{-38}$	$2^{-126-23}$ $\approx 1,4 \times 10^{-45}$	$(2^{24} - 1) \times 2^{127-2}$ $\approx 3,4 \times 10^{38}$
<i>bin64</i>	11	52	-1022	1023	2^{-52} $\approx 2,2 \times 10^{-16}$	2^{-1022} $\approx 2,2 \times 10^{-308}$	2^{-1074} $\approx 5 \times 10^{-324}$	$(2^{53} - 1)2^{1023-52}$ $\approx 1,8 \times 10^{308}$

```
In [1]: 1 + 1e-16 == 1
```

```
•  
•  
•  
•  
•  
•  
•  
•  
•  
•  
•  
•  
•  
•  
•
```

```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
•  
•  
•  
•  
•  
•  
•  
•  
•  
•  
•  
•  
•  
•  
•
```

```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
•  
•  
•  
•  
•  
•  
•  
•
```

```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
Out[3]: 0.0
```

-
-
-
-
-


```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
Out[3]: 0.0
```

```
In [4]: x - 1e-16
```

```
.
```

```
.
```

```
.
```

```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
Out[3]: 0.0
```

```
In [4]: x - 1e-16
```

```
Out[4]: 0.9999999999999999
```

```
•  
•
```

```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
Out[3]: 0.0
```

```
In [4]: x - 1e-16
```

```
Out[4]: 0.9999999999999999
```

```
In [5]: 1e-16 + 1e-18
```

```
.
```

```
In [1]: 1 + 1e-16 == 1
```

```
Out[1]: True
```

```
In [2]: x = 1 + 1e-16
```

```
In [3]: x - 1
```

```
Out[3]: 0.0
```

```
In [4]: x - 1e-16
```

```
Out[4]: 0.9999999999999999
```

```
In [5]: 1e-16 + 1e-18
```

```
Out[5]: 1.01e-16
```

Sommaire

1 Préambule

2 La norme IEEE 754

3 Algèbre des nombres VF

4 Réels, arrondis et flottants

5 Que la force de l'erreur soit avec vous

$$\mathbb{V}_b$$

$$\overline{\mathbb{V}_b}$$

$$\mathbb{V}_b$$

$$\overline{\mathbb{V}_b}$$

Comparaison

Il est très simple et rapide de comparer deux VF : comment la machine procède-t-elle ? Quel est l'avantage de ce stockage des VF ?

Addition

- ① on commence par ramener les deux nombres au même exposant, en l'occurrence le plus grand des deux ;
- ② on ajoute les deux mantisses *complètes* en tenant compte du signe ;
- ③ on renormalise le nombre obtenu.

Addition

- 1 on commence par ramener les deux nombres au même exposant, en l'occurrence le plus grand des deux ;
- 2 on ajoute les deux mantisses *complètes* en tenant compte du signe ;
- 3 on renormalise le nombre obtenu.

Addition

- 1 on commence par ramener les deux nombres au même exposant, en l'occurrence le plus grand des deux ;
- 2 on ajoute les deux mantisses *complètes* en tenant compte du signe ;
- 3 on renormalise le nombre obtenu.

Addition

En *toy7* : $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

$1,1$:

0	0	1	1	1	0	0
---	---	---	---	---	---	---

$0,0011$:

0	0	0	1	1	1	0
---	---	---	---	---	---	---

$1,100$

$0,00111$

$1,10111$

$1,10111$

Addition

En *toy7* : $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 :

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 :

0	0	0	1	1	1	0
---	---	---	---	---	---	---

1,100

0,00111

1,10111

1,10111

Addition

En *toy7* : $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

$1,1$:

0	0	1	1	1	0	0
---	---	---	---	---	---	---

$0,0011$:

0	0	0	1	1	1	0
---	---	---	---	---	---	---

```

1,100
0,00111
-----
1,10111

```

Addition

En *toy7* : $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 :

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 :

0	0	0	1	1	1	0
---	---	---	---	---	---	---

```

1,100
0,00111
-----
1,10111

```

Addition

En *toy7* : $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 :

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 :

0	0	0	1	1	1	0
---	---	---	---	---	---	---

```

1,100
0,00111
-----
1,10111

```


Addition

En *toy7* : $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 :

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 :

0	0	0	1	1	1	0
---	---	---	---	---	---	---

```

1,100
0,00111
-----
1,10111

```

Addition

En *toy7* : $1,1 + 0,0111 \rightarrow 1,1 \times 2^0 + 1,11 \times 2^{-2}$

1,1 :

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 :

0	0	0	1	1	1	0
---	---	---	---	---	---	---

```

1,100
0,00111
-----
1,10111

```

Les arrondis

Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers $+\infty$ (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers $-\infty$ (RD) qui arrondit à la valeur inférieure la plus grande ;

Le mode d'arrondi par défaut est le premier.

Les arrondis

Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers $+\infty$ (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers $-\infty$ (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.

Les arrondis

Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers $+\infty$ (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers $-\infty$ (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.

Les arrondis

Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers $+\infty$ (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers $-\infty$ (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.

Les arrondis

Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers $+\infty$ (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers $-\infty$ (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.

Les arrondis

Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers $+\infty$ (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers $-\infty$ (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.

Les arrondis

Définition 1

- 1 l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
- 2 l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
- 3 l'arrondi vers $+\infty$ (RU) qui arrondit à la valeur supérieure la plus petite ;
- 4 l'arrondi vers $-\infty$ (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.

$$\begin{array}{r}
 1,100 \\
 0,00111 \\
 \text{---} \\
 1,10111
 \end{array}$$

$x = 1,10111$ VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underline{1,100 + 1,00 \text{ulp}(x)} \leq 1,100 + 1,11 \text{ulp}(x) \leq 1,100 + 10,00 \text{ulp}(x)$$

soit y

$$y = 1,100 + 0,11 \text{ulp}(x) \quad x = \text{succ}(y) = 1,100 + 0,01 \text{ulp}(x) \quad \text{donc } RN(x) = \text{succ}(y)$$

$$1,100 + 0,00111 = 1,10111$$

1,100
 0,00111

 1,10111

$x = 1,10111$ VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ulp}(x)}_{\text{succ}(v)}$$

$|x - v| = 0,11 \text{ulp}(x)$ $|x - \text{succ}(v)| = 0,01 \text{ulp}(x)$ donc $RN(x) = \text{succ}(v)$

$$1,1 + 0,00111 = 1,110$$

1,100
 0,00111

 1,10111

$x = 1,10111$ VF en *toy7*?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ulp}(x)}_{\text{succ}(v)}$$

$|x - v| = 0,11 \text{ulp}(x)$ $|x - \text{succ}(v)| = 0,01 \text{ulp}(x)$ donc $RN(x) = \text{succ}(v)$

$$1,1 + 0,00111 = 1,110$$

1,100
 0,00111
 - - -
 1,10111

$x = 1,10111$ VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ulp}(x)}_{\text{succ}(v)}$$

$|x - v| = 0,11 \text{ulp}(x)$ $|x - \text{succ}(v)| = 0,01 \text{ulp}(x)$ donc $RN(x) = \text{succ}(v)$

$$1,100 + 0,00111 = 1,10111$$

$$\begin{array}{r}
 1,100 \\
 0,00111 \\
 \text{---} \\
 1,10111
 \end{array}$$

$x = 1,10111$ VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ulp}(x)}_{\text{succ}(v)}$$

$$|x - v| = 0,11 \text{ulp}(x) \quad |x - \text{succ}(v)| = 0,01 \text{ulp}(x) \quad \text{donc } RN(x) = \text{succ}(v)$$

$$1,1 \oplus 0,00111 = 1,110$$

1,100
 0,00111

 1,10111

$x = 1,10111$ VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ulp}(x)}_{\text{succ}(v)}$$

$|x - v| = 0,11 \text{ulp}(x)$ $|x - \text{succ}(v)| = 0,01 \text{ulp}(x)$ donc $RN(x) = \text{succ}(v)$

$$1,1 \oplus 0,00111 = 1,110$$

1,100
 0,00111

 1,10111

$x = 1,10111$ VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ulp}(x)}_{\text{succ}(v)}$$

$|x - v| = 0,11 \text{ulp}(x)$ $|x - \text{succ}(v)| = 0,01 \text{ulp}(x)$ donc $RN(x) = \text{succ}(v)$

$$1,1 \oplus 0,00111 = 1,110$$

$$\begin{array}{r}
 1,100 \\
 0,00111 \\
 \text{---} \\
 1,10111
 \end{array}$$

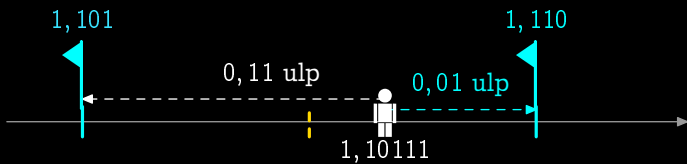
$x = 1,10111$ VF en *toy7* ?

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00 \text{ulp}(x)}_v \leq \underbrace{1,100 + 1,11 \text{ulp}(x)}_x \leq \underbrace{1,100 + 10,00 \text{ulp}(x)}_{\text{succ}(v)}$$

$$|x - v| = 0,11 \text{ulp}(x) \quad |x - \text{succ}(v)| = 0,01 \text{ulp}(x) \quad \text{donc } RN(x) = \text{succ}(v)$$

$$1,1 \oplus 0,00111 = 1,110$$



- ① Est-ce que l'addition des VF est associative ?

```
In [14]: (1 + 1e-16) + 1e-16  
Out[14]: 1.0
```

```
In [15]: 1 + (1e-16 + 1e-16)  
Out[15]: 1.000000000000000002
```

- ② Est-on sûr de l'ordre dans le quel un compilateur calcule $a + b + c + d$?

- ① Est-ce que l'addition des VF est associative ?

```
In [14]: (1 + 1e-16) + 1e-16
```

```
Out[14]: 1.0
```

```
In [15]: 1 + (1e-16 + 1e-16)
```

```
Out[15]: 1.00000000000000002
```

- ② Est-on sûr de l'ordre dans le quel un compilateur calcule $a + b + c + d$?

- 1 Est-ce que l'addition des VF est associative ?

```
In [14]: (1 + 1e-16) + 1e-16
```

```
Out[14]: 1.0
```

```
In [15]: 1 + (1e-16 + 1e-16)
```

```
Out[15]: 1.00000000000000002
```

- 2 Est-on sûr de l'ordre dans le quel un compilateur calcule $a + b + c + d$?

Multiplication

- 1 on « xore » les bits de signe ;
- 2 on additionne les exposants réels et on décale ou plutôt on additionne les exposants décalés et on retire la valeur d'un décalage ;
- 3 on multiplie les mantisses ;
- 4 on normalise.

Multiplication

- 1 on « xore » les bits de signe ;
- 2 on additionne les exposants réels et on décale ou plutôt on additionne les exposants décalés et on retire la valeur d'un décalage ;
- 3 on multiplie les mantisses ;
- 4 on normalise.

Multiplication

- 1 on « xore » les bits de signe ;
- 2 on additionne les exposants réels et on décale ou plutôt on additionne les exposants décalés et on retire la valeur d'un décalage ;
- 3 on multiplie les mantisses ;
- 4 on normalise.

Multiplication

- 1 on « xore » les bits de signe ;
- 2 on additionne les exposants réels et on décale ou plutôt on additionne les exposants décalés et on retire la valeur d'un décalage ;
- 3 on multiplie les mantisses ;
- 4 on normalise.

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array}$$

$$-10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array}$$

• 0 xor 1 donne 1 : le bit de signe est 1

• 1 xor 1 donne 0 : le bit de signe est 0

• 0 xor 0 donne 0

$$\begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 1 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array} \begin{array}{|c|} \hline 0 \\ \hline \end{array}$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ② $101 + 100 - 11 = 1001 - 11 = 110$: l'exposant décalé est 110 donc l'exposant est 3 ;
- ③ $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$;
- ④ le produit est donc $1,100011 \times 2^3$. Le passage à la forme normale va arrondir le produit. On a

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① $0 \text{ xor } 1$ donne 1 : le bit de signe est 1 ;
- ② $101 + 100 - 11 = 1001 - 11 = 110$: l'exposant décalé est 110 donc l'exposant est 3 ;
- ③ $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$;
- ④ le produit est donc $1,100011 \times 2^3$. Le passage à la forme normale va arrondir le produit. On a

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① $0 \text{ xor } 1$ donne 1 : le bit de signe est 1 ;
- ② $101 + 100 - 11 = 1001 - 11 = 110$: l'exposant décalé est 110 donc l'exposant est 3 ;
- ③ $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$;
- ④ le produit est donc $1,100011 \times 2^3$. Le passage à la forme normale va arrondir le produit. On a

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ② $101 + 100 - 11 = 1001 - 11 = 110$: l'exposant décalé est 110 donc l'exposant est 3 ;
- ③ $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$;
- ④ le produit est donc $1,100011 \times 2^3$. Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec $|x - 1,100| = 0,011 \text{ (up)} < x$ et $|x - 1,101| = 0,101 \text{ (up)} > x$ donc

$$RN(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ② $101 + 100 - 11 = 1001 - 11 = 110$: l'exposant décalé est 110 donc l'exposant est 3 ;
- ③ $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$;
- ④ le produit est donc $1,100011 \times 2^3$. Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec $|x - 1,100| = 0,011u/p(x)$ et $|x - 1,101| = 0,101u/p(x)$ donc

$$RM(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$[1,101] \times [0,101] = 0,110$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ② $101 + 100 - 11 = 1001 - 11 = 110$: l'exposant décalé est 110 donc l'exposant est 3 ;
- ③ $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$;
- ④ le produit est donc $1,100011 \times 2^3$. Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec $|x - 1,100| = 0,011ulp(x)$ et $|x - 1,101| = 0,101ulp(x)$ donc

$$RM(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$101,1 \times (-10,01) = -1,100 \times 2^3$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ② $101 + 100 - 11 = 1001 - 11 = 110$: l'exposant décalé est 110 donc l'exposant est 3 ;
- ③ $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$;
- ④ le produit est donc $1,100011 \times 2^3$. Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec $|x - 1,100| = 0,011ulp(x)$ et $|x - 1,101| = 0,101ulp(x)$ donc

$$RN(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ② $101 + 100 - 11 = 1001 - 11 = 110$: l'exposant décalé est 110 donc l'exposant est 3 ;
- ③ $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$;
- ④ le produit est donc $1,100011 \times 2^3$. Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec $|x - 1,100| = 0,011 \text{ulp}(x)$ et $|x - 1,101| = 0,101 \text{ulp}(x)$ donc

$$RN(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ② $101 + 100 - 11 = 1001 - 11 = 110$: l'exposant décalé est 110 donc l'exposant est 3 ;
- ③ $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$;
- ④ le produit est donc $1,100011 \times 2^3$. Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec $|x - 1,100| = 0,011 \text{ulp}(x)$ et $|x - 1,101| = 0,101 \text{ulp}(x)$ donc

$$RN(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$

$$101,1 \times (-10,01)$$

$$101,1 : \begin{array}{|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 1 & 0 & 1 & 1 \\ \hline \end{array} \quad -10,01 : \begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

- ① 0 xor 1 donne 1 : le bit de signe est 1 ;
- ② $101 + 100 - 11 = 1001 - 11 = 110$: l'exposant décalé est 110 donc l'exposant est 3 ;
- ③ $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$;
- ④ le produit est donc $1,100011 \times 2^3$. Le passage à la forme normale va arrondir le produit. On a

$$1,100 < x < 1,101$$

avec $|x - 1,100| = 0,011ulp(x)$ et $|x - 1,101| = 0,101ulp(x)$ donc

$$RN(101,1 \times (-10,01)) = -1,100 \times 2^3$$

$$\begin{array}{|c|c|c|c|c|c|c|} \hline 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ \hline \end{array}$$

① Est-ce que la multiplication des VF est associative ?

```
In [1]: (1.000000001 * 1e-15) * 1e15
```

```
Out[1]: 1.000000001000000002
```

```
In [2]: 1.000000001 * (1e-15 * 1e15)
```

```
Out[2]: 1.000000001
```

② Est-ce que la multiplication des VF est distributive sur l'addition ?

```
In [3]: 1e15 * (1e-16 + 1)
```

```
Out[3]: 10000000000000000.0
```

```
In [4]: (1e15 * 1e-16) + (1e15 * 1)
```

```
Out[4]: 10000000000000000.1
```

- ① Est-ce que la multiplication des VF est associative ?

```
In [1]: (1.00000001 * 1e-15) * 1e15
```

```
Out[1]: 1.0000000100000002
```

```
In [2]: 1.00000001 * (1e-15 * 1e15)
```

```
Out[2]: 1.00000001
```

- ② Est-ce que la multiplication des VF est distributive sur l'addition ?

```
In [3]: 1e15 * (1e-16 + 1)
```

```
Out[3]: 10000000000000000.0
```

```
In [4]: (1e15 * 1e-16) + (1e15 * 1)
```

```
Out[4]: 10000000000000000.1
```

- 1 Est-ce que la multiplication des VF est associative ?

```
In [1]: (1.000000001 * 1e-15) * 1e15
```

```
Out[1]: 1.00000000100000002
```

```
In [2]: 1.000000001 * (1e-15 * 1e15)
```

```
Out[2]: 1.000000001
```

- 2 Est-ce que la multiplication des VF est distributive sur l'addition ?

```
In [3]: 1e15 * (1e-16 + 1)
```

```
Out[3]: 10000000000000000.0
```

```
In [4]: (1e15 * 1e-16) + (1e15 * 1)
```

```
Out[4]: 10000000000000000.1
```


- 1 Est-ce que la multiplication des VF est associative ?

```
In [1]: (1.000000001 * 1e-15) * 1e15
```

```
Out[1]: 1.00000000100000002
```

```
In [2]: 1.000000001 * (1e-15 * 1e15)
```

```
Out[2]: 1.000000001
```

- 2 Est-ce que la multiplication des VF est distributive sur l'addition ?

```
In [3]: 1e15 * (1e-16 + 1)
```

```
Out[3]: 10000000000000000.0
```

```
In [4]: (1e15 * 1e-16) + (1e15 * 1)
```

```
Out[4]: 10000000000000000.1
```

Sommaire

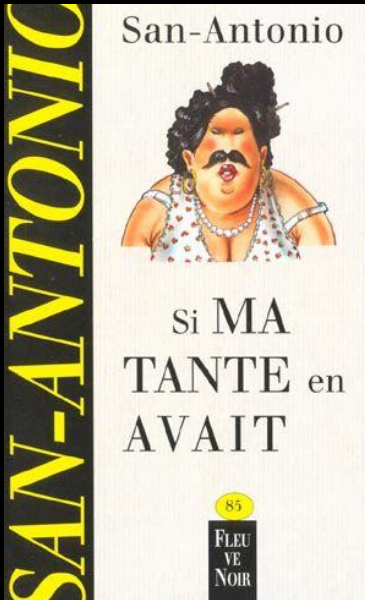
- 1 Préambule
- 2 La norme IEEE 754
- 3 Algèbre des nombres VF
- 4 Réels, arrondis et flottants
- 5 Que la force de l'erreur soit avec vous

```
def base1(a):  
    if (a + 1.0) - a != 1.0:  
        return a  
    else:  
        return base1(2.0 * a)  
  
def base2(a, b):  
    if (a + b) - a == b:  
        return b  
    else:  
        return base2(a, b + 1.0)
```

????!!!!?????!!???

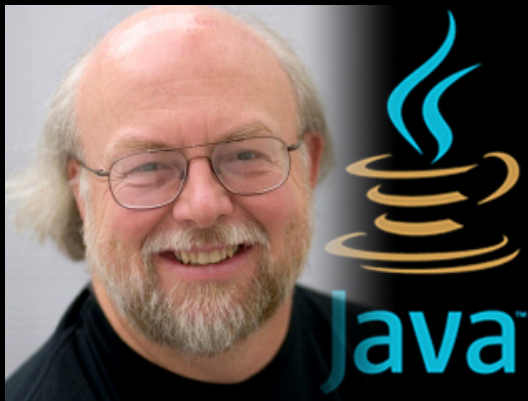
```
def base1(a):  
    if (a + 1.0) - a != 1.0:  
        return a  
    else:  
        return base1(2.0 * a)  
  
def base2(a, b):  
    if (a + b) - a == b:  
        return b  
    else:  
        return base2(a, b + 1.0)
```

????!!!!?????!!???



95 % of the folks out there are completely clueless about floating-point

James Gosling (M. Java) - 28 février 1998





Précision (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

3.17777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de $\frac{1}{3}$.

Précision (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

3.17777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de $\frac{1}{3}$.

Précision (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

Exactitude (*accuracy*) : c'est ce qui relie un nombre au contexte dans lequel il est employé.

3,17777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de π .

Précision (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

Exactitude (*accuracy*) : c'est ce qui relie un nombre au contexte dans lequel il est employé.

3,17777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de π .

Précision (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

Exactitude (*accuracy*) : c'est ce qui relie un nombre au contexte dans lequel il est employé.

3,17777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de π .

Précision (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

Exactitude (*accuracy*) : c'est ce qui relie un nombre au contexte dans lequel il est employé.

3,17777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de π .

Soit x un nombre et \widehat{x} le nombre qui le représente. On distingue :

l'erreur absolue $|x - \widehat{x}|$

l'erreur relative
commise en prenant \widehat{x} à la place de x .

Soit x un nombre et \widehat{x} le nombre qui le représente. On distingue :

l'erreur absolue $|x - \widehat{x}|$

l'erreur relative $\eta = \frac{x - \widehat{x}}{x}$

commise en prenant \widehat{x} à la place de x .

Soit x un nombre et \widehat{x} le nombre qui le représente. On distingue :

l'erreur absolue $|x - \widehat{x}|$

l'erreur relative $\eta = \frac{x - \widehat{x}}{x}$ alors $\widehat{x} = x(1 + \eta)$

commise en prenant \widehat{x} à la place de x .

Soit x un nombre et \widehat{x} le nombre qui le représente. On distingue :

l'erreur absolue $|x - \widehat{x}|$

l'erreur relative $\eta = \frac{x - \widehat{x}}{x}$ alors $\widehat{x} = x(1 + \eta)$

commise en prenant \widehat{x} à la place de x .

Soit x un nombre et \widehat{x} le nombre qui le représente. On distingue :

l'erreur absolue $|x - \widehat{x}|$

l'erreur relative $\eta = \frac{x - \widehat{x}}{x}$ alors $\widehat{x} = x(1 + \eta)$

commise en prenant \widehat{x} à la place de x .

Soit x un nombre et \widehat{x} le nombre qui le représente. On distingue :

l'erreur absolue $|x - \widehat{x}|$

l'erreur relative $\eta = \frac{x - \widehat{x}}{x}$ alors $\widehat{x} = x(1 + \eta)$

commise en prenant \widehat{x} à la place de x .

Feel nervous, but feel in control. It's not dark magic, it's science.

Florent de Dinechin



Sommaire

1 Préambule

2 La norme IEEE 754

3 Algèbre des nombres VF

4 Réels, arrondis et flottants

5 Que la force de l'erreur soit avec vous

Précision machine



Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| = \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{u}{m}$$

Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\epsilon}{m}$$

Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{u}{m}$$

Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\epsilon_M}{m}$$

Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

• Si $x \neq 0$ alors l'erreur relative est majorée $\left| \frac{x - \widehat{x}}{x} \right| \leq \frac{1}{2} \frac{\varepsilon_M}{m}$

• Si $x = 0$ alors l'erreur absolue est majorée $|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$

Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

- ① Si \widehat{x} VF alors l'**erreur relative** est majorée $\left| \frac{x - \widehat{x}}{x} \right| \leq \frac{1}{2} \varepsilon_M$
- ② Si \widehat{x} est sous-normal, alors l'**erreur absolue** est majorée

Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

- ① Si \widehat{x} VF alors l'**erreur relative** est majorée $\left| \frac{x - \widehat{x}}{x} \right| \leq \frac{1}{2} \varepsilon_M$
- ② Si \widehat{x} est sous-normal, alors l'**erreur absolue** est majorée

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E_{\min} - 1 - n}$$

Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

- ① Si \widehat{x} VF alors l'**erreur relative** est majorée $\left| \frac{x - \widehat{x}}{x} \right| \leq \frac{1}{2} \varepsilon_M$
- ② Si \widehat{x} est sous-normal, alors l'**erreur absolue** est majorée
 $|x - \widehat{x}| \leq \frac{1}{2} 2^{E_{\min} - 1 - n}$

Précision machine

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

$$|x - \widehat{x}| \leq \frac{1}{2} 2^{E-n}$$

$$\left| \frac{x - \widehat{x}}{x} \right| \leq \left| \frac{x - \widehat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

- ① Si \widehat{x} VF alors l'**erreur relative** est majorée $\left| \frac{x - \widehat{x}}{x} \right| \leq \frac{1}{2} \varepsilon_M$
- ② Si \widehat{x} est sous-normal, alors l'**erreur absolue** est majorée $|x - \widehat{x}| \leq \frac{1}{2} 2^{E_{\min} - 1 - n}$

NE FAITES PAS DES TESTS D'ÉGALITÉ
MAIS DES TESTS D'APPARTENANCE
À DES INTERVALLES DE LARGEUR $L'\varepsilon$
MACHINE!

Précision machine

$$\widehat{x} = x(1 + \eta_1) + \eta_2$$

avec :

- si \widehat{x} est normal $|\eta_1| \leq \frac{1}{2}\varepsilon_M$ et $\eta_2 = 0$
- si \widehat{x} est sous-normal $\eta_1 = 0$ et $|\eta_2| \leq \frac{1}{2}2^{E_{\min}-1-n}$

Mouais...

Précision machine

$$\widehat{x} = x(1 + \eta_1) + \eta_2$$

avec :

- si \widehat{x} est normal $|\eta_1| \leq \frac{1}{2}\varepsilon_M$ et $\eta_2 = 0$
- si \widehat{x} est sous-normal $\eta_1 = 0$ et $|\eta_2| \leq \frac{1}{2}2^{E_{\min}-1-n}$

Mouais...

Précision machine

$$\widehat{x} = x(1 + \eta_1) + \eta_2$$

avec :

- si \widehat{x} est normal $|\eta_1| \leq \frac{1}{2}\varepsilon_M$ et $\eta_2 = 0$
- si \widehat{x} est sous-normal $\eta_1 = 0$ et $|\eta_2| \leq \frac{1}{2}2^{E_{\min}-1-n}$

Mouais...

Précision machine

$$\widehat{x} = x(1 + \eta_1) + \eta_2$$

avec :

- si \widehat{x} est normal $|\eta_1| \leq \frac{1}{2}\varepsilon_M$ et $\eta_2 = 0$
- si \widehat{x} est sous-normal $\eta_1 = 0$ et $|\eta_2| \leq \frac{1}{2}2^{E_{\min}-1-n}$

Mouais...

Élimination

$$\widehat{a} = a(1 + \eta_a), \widehat{b} = b(1 + \eta_b), x = a - b \text{ et } \widehat{x} = \widehat{a} - \widehat{b} = a - b$$

$$\left| \frac{x - \widehat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$

Élimination

$$\widehat{a} = a(1 + \eta_a), \widehat{b} = b(1 + \eta_b), x = a - b \text{ et } \widehat{x} = \widehat{a - b} = \widehat{a} - \widehat{b}.$$

$$\left| \frac{\widehat{x} - x}{x} \right| = \left| \frac{a\eta_a - b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$

Élimination

$$\widehat{a} = a(1 + \eta_a), \widehat{b} = b(1 + \eta_b), x = a - b \text{ et } \widehat{x} = \widehat{a - b} = \widehat{a} - \widehat{b}.$$

$$\left| \frac{x - \widehat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$

Élimination

$$\widehat{a} = a(1 + \eta_a), \widehat{b} = b(1 + \eta_b), x = a - b \text{ et } \widehat{x} = \widehat{a - b} = \widehat{a} - \widehat{b}.$$

$$\left| \frac{x - \widehat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$

Élimination

$$\widehat{a} = a(1 + \eta_a), \widehat{b} = b(1 + \eta_b), x = a - b \text{ et } \widehat{x} = \widehat{a - b} = \widehat{a} - \widehat{b}.$$

$$\left| \frac{x - \widehat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$

Élimination

$$\widehat{a} = a(1 + \eta_a), \widehat{b} = b(1 + \eta_b), x = a - b \text{ et } \widehat{x} = \widehat{a - b} = \widehat{a} - \widehat{b}.$$

$$\left| \frac{x - \widehat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$

Élimination

$$\widehat{a} = a(1 + \eta_a), \widehat{b} = b(1 + \eta_b), x = a - b \text{ et } \widehat{x} = \widehat{a - b} = \widehat{a} - \widehat{b}.$$

$$\left| \frac{x - \widehat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$

Catastrophe

```
def deriv(f,h,x):  
    return (f(x+h) - f(x)) / h
```

```
def derivn(f,x,h,n):  
    if n == 0:  
        return f(x)  
    else:  
        return derivn(lambda x: deriv(f,h,x),x,h,n-1)
```

Catastrophe

```
def deriv(f,h,x):  
    return (f(x+h) - f(x)) / h
```

```
def derivn(f,x,h,n):  
    if n == 0:  
        return f(x)  
    else:  
        return derivn(lambda x: deriv(f,h,x),x,h,n-1)
```

Catastrophe

```
In [1]: [derivn(lambda x: x**4,1,1e-7,k) for k in range(5)]
```

```
Out[1]:
```

```
[1,  
 4.0000000601855902,  
12.012613126444194,  
-222044.6049250313,  
2220446049250.313]
```

```
In [2]: [derivn(lambda x: x**4,1,1e-8,k) for k in range(5)]
```

```
Out[2]: [1, 4.0000000042303498, 11.102230246251565, 0.0,  
- 2.220446049250313e+16]
```


Catastrophe

```
In [1]: [derivn(lambda x: x**4,1,1e-7,k) for k in range(5)]
```

```
Out[1]:
```

```
[1,  
 4.0000000601855902,  
12.012613126444194,  
-222044.6049250313,  
2220446049250.313]
```

```
In [2]: [derivn(lambda x: x**4,1,1e-8,k) for k in range(5)]
```

```
Out[2]: [1, 4.000000042303498, 11.102230246251565, 0.0,  
→ 2.220446049250313e+16]
```

Catastrophe

```
In [3]: [derivn(lambda x: x**4,1,1e-9,k) for k in range(5)]
```

```
Out[3]: [1, 4.0000000330961484, 0.0, 0.0, 0.0]
```

```
In [4]: [derivn(lambda x: x**4,1,1e-3,k) for k in range(6)]
```

```
Out[4]:
```

```
[1,  
 4.006004000999486,  
 12.024014000244776,  
 24.03599941303014,  
 24.001245435556484,  
 -2.4424906541753444]
```

Catastrophe

```
In [3]: [derivn(lambda x: x**4,1,1e-9,k) for k in range(5)]
```

```
Out[3]: [1, 4.0000000330961484, 0.0, 0.0, 0.0]
```

```
In [4]: [derivn(lambda x: x**4,1,1e-3,k) for k in range(6)]
```

```
Out[4]:
```

```
[1,  
 4.006004000999486,  
12.024014000244776,  
24.03599941303014,  
24.001245435556484,  
-2.4424906541753444]
```

Catastrophe

```
In [5]: [derivn(lambda x: x**4,1,2**-10,k) for k in range(6)]
Out[5]: [1, 4.005863190628588, 12.02345085144043, 24.03515625,
         → 24.0, 0.0]
```

```
In [6]: [derivn(lambda x: x**4,1,1.0/1026.,k) for k in range(6)]
Out[6]:
[1,
 4.005851753982072,
 12.023405112200917,
 24.035087928668187,
 23.999573957547014,
 0.755876563500351]
```

Catastrophe

```
In [5]: [derivn(lambda x: x**4,1,2**-10,k) for k in range(6)]
Out[5]: [1, 4.005863190628588, 12.02345085144043, 24.03515625,
         → 24.0, 0.0]
```

```
In [6]: [derivn(lambda x: x**4,1,1.0/1026.,k) for k in range(6)]
Out[6]:
[1,
 4.005851753982072,
 12.023405112200917,
 24.035087928668187,
 23.999573957547014,
 0.755876563500351]
```

Catastrophe

```
In [7]: [derivn(lambda x: x**4,1,2**-20,k) for k in range(6)]  
Out[7]: [1, 4.0000005722045898, 12.0, 0.0, 268435456.0,  
         ↪ -844424930131968.0]
```

Catastrophe

```
In [7]: [derivn(exp,0,2**k-5,k) for k in range(10)]
```

```
Out[7]:
```

```
[1.0,  
 1.0157890399712883,  
 1.0318273737254913,  
 1.0481189373895177,  
 1.0646677284967154,  
 1.081477828323841,  
 1.0985534191131592,  
 1.1158447265625,  
 1.1376953125,  
 0.953125]
```

Catastrophe

```
In [8]: [derivn(exp,0,2**k,k) for k in range(10)]
```

```
Out[8]: [1.0, 1.00000004768371582, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0,
         ↪ 0.0, 0.0]
```


TROP DE PRÉCISION TUE LA PRÉCISION !

DIVISEZ PAR DES PUISSANCES DE
2...PAS DES PUISSANCES DE 10 !

TROP DE PRÉCISION TUE LA
PRÉCISION !
DIVISEZ PAR DES PUISSANCES DE
2...PAS DES PUISSANCES DE 10 !

TROP DE PRÉCISION TUE LA
PRÉCISION !
DIVISEZ PAR DES PUISSANCES DE
2...PAS DES PUISSANCES DE 10 !

Exercice 4

Vous savez résoudre une équation du type $ax^2 + bx + c = 0$ avec $a \neq 0$...

Les racines, si elles existent, sont données par une formule bien connue dépendant de a , b et c :

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».

- ① *Que se passe-t-il lorsque $b^2 \gg |4ac|$? En quoi la formule $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$ peut aider ?*
- ② *Que se passe-t-il lorsque $b^2 \approx 4ac$? Peut-on y remédier ? Que peut-on dire de Δ par rapport à b^2 ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$?*
- ④ *Et dans le cas de 10 ?*

Exercice 4

Vous savez résoudre une équation du type $ax^2 + bx + c = 0$ avec $a \neq 0$...

Les racines, si elles existent, sont données par une formule bien connue dépendant de a , b et c :

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».

- ① *Que se passe-t-il lorsque $b^2 \gg |4ac|$? En quoi la formule $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$ peut aider ?*
- ② *Que se passe-t-il lorsque $b^2 \approx 4ac$? Peut-on y remédier ? Que peut-on dire de Δ par rapport à b^2 ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$?*
- ④ *Et dans le cas de $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$?*

Exercice 4

Vous savez résoudre une équation du type $ax^2 + bx + c = 0$ avec $a \neq 0$...

Les racines, si elles existent, sont données par une formule bien connue dépendant de a , b et c :

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».

- ① *Que se passe-t-il lorsque $b^2 \gg |4ac|$? En quoi la formule $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$ peut aider ?*
- ② *Que se passe-t-il lorsque $b^2 \approx 4ac$? Peut-on y remédier ? Que peut-on dire de Δ par rapport à b^2 ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$?*
- ④ *Et dans le cas de $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$?*
- ⑤ *Moralité ?*

Exercice 4

Vous savez résoudre une équation du type $ax^2 + bx + c = 0$ avec $a \neq 0$...

Les racines, si elles existent, sont données par une formule bien connue dépendant de a , b et c :

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».

- ① *Que se passe-t-il lorsque $b^2 \gg |4ac|$? En quoi la formule $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$ peut aider ?*
- ② *Que se passe-t-il lorsque $b^2 \approx 4ac$? Peut-on y remédier ? Que peut-on dire de Δ par rapport à b^2 ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$?*
- ④ *Et dans le cas de $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$?*
- ⑤ *Moralité ?*

Exercice 4

Vous savez résoudre une équation du type $ax^2 + bx + c = 0$ avec $a \neq 0$...

Les racines, si elles existent, sont données par une formule bien connue dépendant de a , b et c :

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».

- ① *Que se passe-t-il lorsque $b^2 \gg |4ac|$? En quoi la formule $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$ peut aider ?*
- ② *Que se passe-t-il lorsque $b^2 \approx 4ac$? Peut-on y remédier ? Que peut-on dire de Δ par rapport à b^2 ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$?*
- ④ *Et dans le cas de $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$?*
- ⑤ *Moralité ?*

Exercice 4

Vous savez résoudre une équation du type $ax^2 + bx + c = 0$ avec $a \neq 0$...

Les racines, si elles existent, sont données par une formule bien connue dépendant de a , b et c :

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».

- ① *Que se passe-t-il lorsque $b^2 \gg |4ac|$? En quoi la formule $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$ peut aider ?*
- ② *Que se passe-t-il lorsque $b^2 \approx 4ac$? Peut-on y remédier ? Que peut-on dire de Δ par rapport à b^2 ? Y a-t-il élimination catastrophique ?*
- ③ *Que se passe-t-il dans le cas de l'équation $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$?*
- ④ *Et dans le cas de $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$?*
- ⑤ *Moralité ?*

Exercice 5 (Max IEEE)

Fonction $\text{max}(x, y : \text{flottants}) : \text{flottant}$

Si $x \geq y$ **Alors**

| $\text{max} = x$

Sinon

| $\text{max} = y$

FinSi

```
In [1]: max(1, 1+1e-16)
```

```
Out[1]: 1
```

```
In [2]: max(0.1, 1e500/1e500)
```

```
Out[2]: 0.1
```

```
In [3]: 0.1 <= 1e500/1e500
```

```
Out[3]: False
```

```
In [4]: 0.1 > 1e500/1e500
```

```
Out[4]: False
```

Élimination

```
def expn(n):  
    return (1. + 1./n)**n
```

Élimination

```
In [1]: [(exp(1) -
         ↪ expn(10.0**k)) for k
         ↪ in range(20)]
```

```
Out[1]:
```

```
[0.7182818284590451,
 0.12453936835904278,
 0.01346799903751661,
 0.0013578962234515046,
 0.000135901634119584,
 1.359126674760347e-05,
 1.359363291708604e-06,
 1.3432696333026684e-07,
 3.011168736577474e-08,
```

```
-2.2355251516614771e-07,
-2.2477574246337895e-07,
-2.248980650598753e-07,
-0.00024166757819266138,
 0.002171794372144209,
 0.0021717943720220845,
-0.31675337809021675,
 1.718281828459045,
 1.718281828459045,
 1.718281828459045,
 1.718281828459045]
```

Élimination

```
In [2]: [(exp(1) -  
         ↪ expn(2.0**(3*k)))  
         ↪ for k in range(20)]
```

Out[2]:

```
[0.7182818284590451,  
 0.1524973145086972,  
 0.020936875893946105,  
 0.0026498282900537795,  
 0.0003317472693793455,  
 4.147652875108321e-05,  
 5.1846929989274315e-06,  
 6.480886076687398e-07,  
 8.101110671177025e-08,
```

```
1.0126388616527038e-08,  
1.2657985770658797e-09,  
1.582245445774788e-10,  
1.977795704988239e-11,  
2.4722446312352986e-12,  
3.090860900556436e-13,  
3.863576125695545e-14,  
4.884981308350689e-15,  
4.440892098500626e-16,  
1.718281828459045,  
1.718281828459045]
```

```
In [89]: [2.0**53 + k for k in range(10)]
```

```
Out[89]:
```

```
[9007199254740992.0,  
 9007199254740992.0,  
 9007199254740994.0,  
 9007199254740996.0,  
 9007199254740996.0,  
 9007199254740996.0,  
 9007199254740998.0,  
 9007199254741000.0,  
 9007199254741000.0,  
 9007199254741000.0]
```

```
In [91]: [2.0**55 + k for k in range(10)]
```

```
Out[91]:
```

```
[3.602879701896397e+16,  
 3.602879701896397e+16,  
 3.602879701896397e+16,  
 3.602879701896397e+16,  
 3.602879701896397e+16,  
 3.6028797018963976e+16,  
 3.6028797018963976e+16,  
 3.6028797018963976e+16,  
 3.6028797018963976e+16,  
 3.6028797018963976e+16]
```


Attention !

Pour éviter de déduire des lemmes suivants des théorèmes totalement faux, n'oubliez pas que **DANS CE QUI SUIT X ET Y SONT DES NOMBRES À VIRGULE FLOTTANTE !**

Lemme 2 (Majoration de l'erreur d'une somme)

Posons $x \oplus y = x + y + \text{err}(x \oplus y)$. Alors, s'il n'y a pas de dépassement de capacité,

$$|\text{err}(x \oplus y)| \leq \min(|x|, |y|)$$

On a bien sûr un résultat analogue pour la différence

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- x est un VF situé à la distance $|y|$ de $x + y$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \ominus y)| \leq |x|$

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- x est un VF situé à la distance $|y|$ de $x + y$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \ominus y)| \leq |x|$

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- x est un VF situé à la distance $|y|$ de $x + y$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \ominus y)| \leq |x|$

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- x est un VF situé à la distance $|y|$ de $x + y$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \ominus y)| \leq |x|$

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- x est un VF situé à la distance $|y|$ de $x + y$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \ominus y)| \leq |x|$

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- x est un VF situé à la distance $|y|$ de $x + y$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \oplus y)| \leq |x|$

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- x est un VF situé à la distance $|y|$ de $x + y$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \oplus y)| \leq |x|$

De l'importance de disposer avec la IEEE 754 de la meilleure approximation

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- x est un VF situé à la distance $|y|$ de $x + y$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \oplus y)| \leq |x|$

À noter

De l'importance de disposer avec la IEEE 754 de la **meilleure approximation !**

Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

- $|x| \geq |y|$
- x et y sont des VF

Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

- $|x| \geq |y|$
- x et y sont des VF : le plus petit bit significatif de $\text{err}(x \oplus y)$ est au moins de magnitude celle de $\text{ulp}(y)$



Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

- $|x| \geq |y|$
- x et y sont des VF : le plus petit bit significatif de $\text{err}(x \oplus y)$ est au moins de magnitude celle de $\text{ulp}(y)$



Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

- $|x| \geq |y|$
- x et y sont des VF : le plus petit bit significatif de $\text{err}(x \oplus y)$ est au moins de magnitude celle de $\text{ulp}(y)$



Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

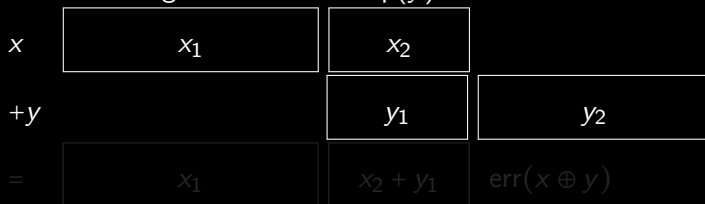
- $|x| \geq |y|$
- x et y sont des VF : le plus petit bit significatif de $\text{err}(x \oplus y)$ est au moins de magnitude celle de $\text{ulp}(y)$



Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

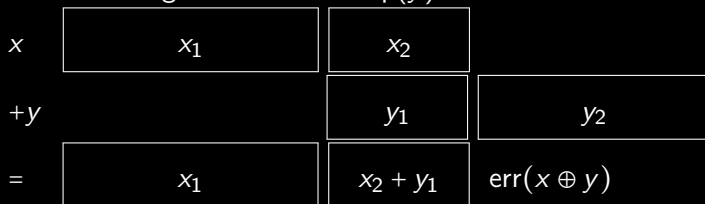
- $|x| \geq |y|$
- x et y sont des VF : le plus petit bit significatif de $\text{err}(x \oplus y)$ est au moins de magnitude celle de $\text{ulp}(y)$



Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

- $|x| \geq |y|$
- x et y sont des VF : le plus petit bit significatif de $\text{err}(x \oplus y)$ est au moins de magnitude celle de $\text{ulp}(y)$

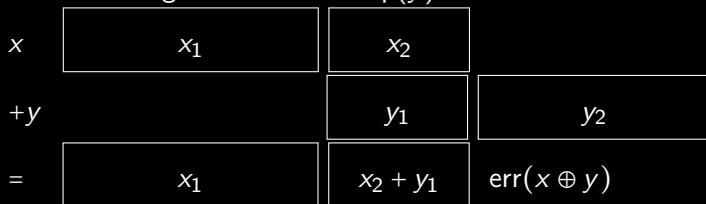


© 2011, 2012, 2013, 2014

Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

- $|x| \geq |y|$
- x et y sont des VF : le plus petit bit significatif de $\text{err}(x \oplus y)$ est au moins de magnitude celle de $\text{ulp}(y)$

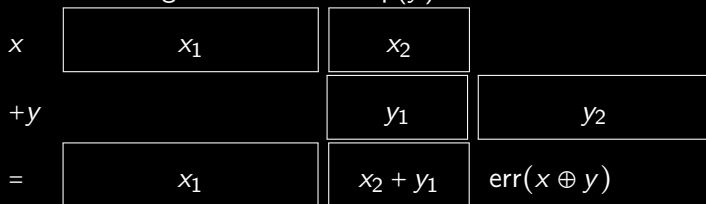


- $|\text{err}(x \oplus y)| \leq |y|$ donc la mantisse entière de $\text{err}(x \oplus y)$ a une longueur inférieure à p bits

Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

- $|x| \geq |y|$
- x et y sont des VF : le plus petit bit significatif de $\text{err}(x \oplus y)$ est au moins de magnitude celle de $\text{ulp}(y)$

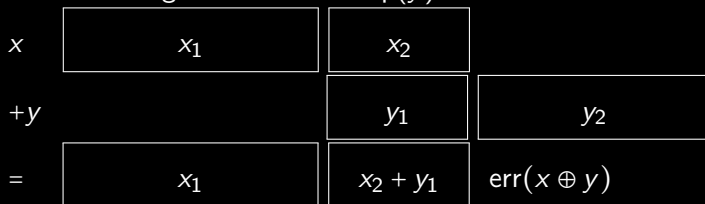


- $|\text{err}(x \oplus y)| \leq |y|$ donc la mantisse entière de $\text{err}(x \oplus y)$ a une longueur inférieure à p bits

Corollaire 3 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

- $|x| \geq |y|$
- x et y sont des VF : le plus petit bit significatif de $\text{err}(x \oplus y)$ est au moins de magnitude celle de $\text{ulp}(y)$



- $|\text{err}(x \oplus y)| \leq |y|$ donc la mantisse entière de $\text{err}(x \oplus y)$ a une longueur inférieure à p bits

Lemme 4

*Supposons que $|x + y| \leq \min(|x|, |y|)$, alors $x \oplus y = x + y$.
On obtient un résultat analogue pour la soustraction.*

- 1 Supposons, sans perdre de généralité, que $|x| \geq |y|$
- 2 Le plus petit bit significatif de $x + y$ est au moins de magnitude celle de $\text{ulp}(y)$
- 3 $|x + y| \leq |y|$ donc la mantisse entière de $x + y$ a une longueur inférieure à p bits

Lemme 4

*Supposons que $|x + y| \leq \min(|x|, |y|)$, alors $x \oplus y = x + y$.
On obtient un résultat analogue pour la soustraction.*

- 1 Supposons, sans perdre de généralité, que $|x| \geq |y|$
- 2 Le plus petit bit significatif de $x + y$ est au moins de magnitude celle de $\text{ulp}(y)$
- 3 $|x + y| \leq |y|$ donc la mantisse entière de $x + y$ a une longueur inférieure à p bits

Lemme 4

Supposons que $|x + y| \leq \min(|x|, |y|)$, alors $x \oplus y = x + y$.

On obtient un résultat analogue pour la soustraction.

- 1 Supposons, sans perdre de généralité, que $|x| \geq |y|$
- 2 Le plus petit bit significatif de $x + y$ est au moins de magnitude celle de $\text{ulp}(y)$
- 3 $|x + y| \leq |y|$ donc la mantisse entière de $x + y$ a une longueur inférieure à p bits

Lemme 4

Supposons que $|x + y| \leq \min(|x|, |y|)$, alors $x \oplus y = x + y$.

On obtient un résultat analogue pour la soustraction.

- 1 Supposons, sans perdre de généralité, que $|x| \geq |y|$
- 2 Le plus petit bit significatif de $x + y$ est au moins de magnitude celle de $\text{ulp}(y)$
- 3 $|x + y| \leq |y|$ donc la mantisse entière de $x + y$ a une longueur inférieure à p bits

DANGER

Nous avons démontré nos lemmes en considérant des VF x et y .
Est-ce que le résultat suivant contredit notre dernier lemme ?

In [86]: 1 - 0.9

Out[86]: 0.09999999999999998

Lemme 5 (Lemme de Sterbenz (1973))

Soit $(x, y) \in \mathbb{V}^2$ vérifiant $\frac{x}{2} \leq y \leq 2x$ alors

$$x \ominus y = x - y$$

La différence de deux VF suffisamment proches est donc exacte.

Dans tous les cas, si vous n'avez rien, on peut utiliser le lemme 4
page 205.

Lemme 5 (Lemme de Sterbenz (1973))

Soit $(x, y) \in \mathbb{V}^2$ vérifiant $\frac{x}{2} \leq y \leq 2x$ alors

$$x \ominus y = x - y$$

La différence de deux VF suffisamment proches est donc exacte.

• Si $x < y$ alors $x < y \leq 2x$ donc $0 < y - x < x$

• Si $x > y$ alors $x/2 < y < x$

Dans tous les cas, le résultat obtenu est exact et on peut utiliser le lemme 4 (page 205).

Lemme 5 (Lemme de Sterbenz (1973))

Soit $(x, y) \in \mathbb{V}^2$ vérifiant $\frac{x}{2} \leq y \leq 2x$ alors

$$x \ominus y = x - y$$

La différence de deux VF suffisamment proches est donc exacte.

- 1 $x < y$: alors $x < y \leq 2x$ donc $0 < y - x \leq x \leq y$;
- 2 $x \geq y$: alors $\frac{x}{2} \leq y \leq x$ donc $-\frac{x}{2} \leq y - x \leq 0$ et par suite $0 \leq x - y \leq \frac{x}{2} \leq y \leq x$.

Dans tous les cas $|x - y| \leq \min(|x|, |y|)$ et on peut utiliser le lemme 4 page 295.

Lemme 5 (Lemme de Sterbenz (1973))

Soit $(x, y) \in \mathbb{V}^2$ vérifiant $\frac{x}{2} \leq y \leq 2x$ alors

$$x \ominus y = x - y$$

La différence de deux VF suffisamment proches est donc exacte.

- 1 $x < y$: alors $x < y \leq 2x$ donc $0 < y - x \leq x \leq y$;
- 2 $x \geq y$: alors $\frac{x}{2} \leq y \leq x$ donc $-\frac{x}{2} \leq y - x \leq 0$ et par suite $0 \leq x - y \leq \frac{x}{2} \leq y \leq x$.

Dans tous les cas $|x - y| \leq \min(|x|, |y|)$ et on peut utiliser le lemme 4 page 295.

Lemme 5 (Lemme de Sterbenz (1973))

Soit $(x, y) \in \mathbb{V}^2$ vérifiant $\frac{x}{2} \leq y \leq 2x$ alors

$$x \ominus y = x - y$$

La différence de deux VF suffisamment proches est donc exacte.

- 1 $x < y$: alors $x < y \leq 2x$ donc $0 < y - x \leq x \leq y$;
- 2 $x \geq y$: alors $\frac{x}{2} \leq y \leq x$ donc $-\frac{x}{2} \leq y - x \leq 0$ et par suite $0 \leq x - y \leq \frac{x}{2} \leq y \leq x$.

Dans tous les cas $|x - y| \leq \min(|x|, |y|)$ et on peut utiliser le lemme 4 page 295.

Lemme 5 (Lemme de Sterbenz (1973))

Soit $(x, y) \in \mathbb{V}^2$ vérifiant $\frac{x}{2} \leq y \leq 2x$ alors

$$x \ominus y = x - y$$

La différence de deux VF suffisamment proches est donc exacte.

- 1 $x < y$: alors $x < y \leq 2x$ donc $0 < y - x \leq x \leq y$;
- 2 $x \geq y$: alors $\frac{x}{2} \leq y \leq x$ donc $-\frac{x}{2} \leq y - x \leq 0$ et par suite $0 \leq x - y \leq \frac{x}{2} \leq y \leq x$.

Dans tous les cas $|x - y| \leq \min(|x|, |y|)$ et on peut utiliser le lemme 4 page 295.

- Concrètement, à quoi correspond cette condition $\frac{x}{2} \leq y \leq 2x$?
- Demandez-vous ce que l'on peut dire de l'écart maximum entre les exposants de x et y .

- Concrètement, à quoi correspond cette condition $\frac{x}{2} \leq y \leq 2x$?
- Demandez-vous ce que l'on peut dire de l'écart maximum entre les exposants de x et y .

```
In [9]: 2.**53
```

```
Out[9]: 9007199254740992.0
```

```
In [10]: 2.**53 + 1
```

```
Out[10]: 9007199254740992.0
```

```
In [11]: 2.**53 + 2
```

```
Out[11]: 9007199254740994.0
```

```
In [12]: 2.**53 + 3
```

```
Out[12]: 9007199254740996.0
```

```
In [13]: 2.**53 + 4
```

```
Out[13]: 9007199254740996.0
```

```
In [14]: 2.**53 + 5
```

```
Out[14]: 9007199254740996.0
```

$$2^{53} = (-1)^0 \times 1,000\dots000 \times 2^{53+1023}$$

0 10000110100 000

$$1 = (-1)^0 \times 1,000\dots000 \times 2^{0+1023}$$

0 01111111111 000

$$1 - \frac{1 - 0,000\dots000 \times 1}{2^{23}}$$

$$2^{53} = (-1)^0 \times 1,000\dots000 \times 2^{53+1023}$$

0 1000110100 00

$$1 = (-1)^0 \times 1,000\dots000 \times 2^{0+1023}$$

0 011111111111 00

$$1 = \underbrace{0,000\dots000}_{{52 \text{ bits}}} 1 \times 2^{53}$$

Unit in the Last Place

Unit in the Last Place

Mauvaise nouvelle : une addition entre deux nombres flottants peut donc, dans certains cas, créer une erreur.

Bonne nouvelle : on peut toujours éviter ce problème.

Unit in the Last Place

Mauvaise nouvelle : une addition entre deux nombres flottants peut donc, dans certains cas, créer une erreur.

Bonne nouvelle : on peut récupérer cette erreur.

Unit in the Last Place

Mauvaise nouvelle : une addition entre deux nombres flottants peut donc, dans certains cas, créer une erreur.

Bonne nouvelle : on peut récupérer cette erreur.



Théorème 6 (Fast2Sum - Dekker & Kahan)

On considère deux VF x et y tels que $|x| \geq |y|$ et l'algorithme suivant :

```
1   $s \leftarrow x \oplus y$   
2   $y_v \leftarrow s \ominus x$   
3   $d \leftarrow y \ominus y_v$   
4  Retourner  $(s, d)$ 
```

Alors $x + y = s + d$ avec $s = x \oplus y$ et $d = \text{err}(x \oplus y)$. De plus s et d ne se chevauchent pas.

Théorème 6 (Fast2Sum - Dekker & Kahan)

On considère deux VF x et y tels que $|x| \geq |y|$ et l'algorithme suivant :

```

1   $s \leftarrow x \oplus y$ 
2   $y_v \leftarrow s \ominus x$ 
3   $d \leftarrow y \ominus y_v$ 
4  Retourner  $(s, d)$ 

```

Alors $x + y = s + d$ avec $s = x \oplus y$ et $d = \text{err}(x \oplus y)$. De plus s et d ne se chevauchent pas.

• si x et y sont de même signe OU si $|y| < \frac{1}{2}$ alors $s = x + 2y$ et on peut appliquer le lemme

Théorème 6 (Fast2Sum - Dekker & Kahan)

On considère deux VF x et y tels que $|x| \geq |y|$ et l'algorithme suivant :

```

1  s ← x ⊕ y
2  yv ← s ⊖ x
3  d ← y ⊖ yv
4  Retourner (s, d)

```

Alors $x + y = s + d$ avec $s = x \oplus y$ et $d = \text{err}(x \oplus y)$. De plus s et d ne se chevauchent pas.

- si x et y sont de même signe OU si $|y| \leq \frac{|x|}{2}$: alors $\frac{x}{2} \leq s \leq 2x$ et on peut appliquer le lemme ;
- sinon : on a x et y de signes opposés ET $y > \frac{|x|}{2}$ alors si y est négatif $x/2 < -y < x$ et sinon $-x/2 < y < -x$. Dans les deux cas, d'après le lemme de Sterbenz, s est calculée exactement et alors $y_v = y$.

Théorème 6 (Fast2Sum - Dekker & Kahan)

On considère deux VF x et y tels que $|x| \geq |y|$ et l'algorithme suivant :

```

1  s ← x ⊕ y
2  yv ← s ⊖ x
3  d ← y ⊖ yv
4  Retourner (s, d)

```

Alors $x + y = s + d$ avec $s = x \oplus y$ et $d = \text{err}(x \oplus y)$. De plus s et d ne se chevauchent pas.

- si x et y sont de même signe OU si $|y| \leq \frac{|x|}{2}$: alors $\frac{x}{2} \leq s \leq 2x$ et on peut appliquer le lemme ;
- sinon : on a x et y de signes opposés ET $y > \frac{|x|}{2}$ alors si y est négatif $x/2 < -y < x$ et sinon $-x/2 < y < -x$. Dans les deux cas, d'après le lemme de Sterbenz, s est calculée exactement et alors $y_v = y$.

Somme compensée

x	x_1	x_2	
$+y$		y_1	y_2
$= s$	x_1	$x_2 + y_1$	y_2 perdu
$-x = y_v$		y_1	0
$-y = -d$			$-y_2$

0	on récupère l'erreur
-----	----------------------

```
def fast2sum(x,y):
    if abs(x) >= abs(y):
        s = x + y
        yv = s - x
        d = y - yv
        return (s,d)
    else:
        return fast2sum(y,x)
```

```
In [100]: fast2sum(2.0**54,1.0)
```

```
Out[100]: (1.8014398509481984e+16, 1.0)
```

```
In [101]: fast2sum(2.0**54,2.0)
```

```
Out[101]: (1.8014398509481984e+16, 2.0)
```

```
In [102]: fast2sum(2.0**54,3.0)
```

```
Out[102]: (1.8014398509481988e+16, -1.0)
```

```
In [103]: fast2sum(2.0**54,4.0)
```

```
Out[103]: (1.8014398509481988e+16, 0.0)
```

```
In [104]: fast3sum(2.0**54,1.0,1.0)
```

```
Out[104]: (1.8014398509481984e+16, 1.0)
```

```
In [100]: fast2sum(2.0**54,1.0)
```

```
Out[100]: (1.8014398509481984e+16, 1.0)
```

```
In [101]: fast2sum(2.0**54,2.0)
```

```
Out[101]: (1.8014398509481984e+16, 2.0)
```

```
In [102]: fast2sum(2.0**54,3.0)
```

```
Out[102]: (1.8014398509481988e+16, -1.0)
```

```
In [103]: fast2sum(2.0**54,4.0)
```

```
Out[103]: (1.8014398509481988e+16, 0.0)
```

```
In [104]: fast2sum(2.0**54,5.0)
```

```
Out[104]: (1.8014398509481988e+16, 1.0)
```

```
In [100]: fast2sum(2.0**54,1.0)
```

```
Out[100]: (1.8014398509481984e+16, 1.0)
```

```
In [101]: fast2sum(2.0**54,2.0)
```

```
Out[101]: (1.8014398509481984e+16, 2.0)
```

```
In [102]: fast2sum(2.0**54,3.0)
```

```
Out[102]: (1.8014398509481988e+16, -1.0)
```

```
In [103]: fast2sum(2.0**54,4.0)
```

```
Out[103]: (1.8014398509481988e+16, 0.0)
```

```
In [104]: fast2sum(2.0**54,5.0)
```

```
Out[104]: (1.8014398509481988e+16, 1.0)
```



```
In [100]: fast2sum(2.0**54,1.0)
```

```
Out[100]: (1.8014398509481984e+16, 1.0)
```

```
In [101]: fast2sum(2.0**54,2.0)
```

```
Out[101]: (1.8014398509481984e+16, 2.0)
```

```
In [102]: fast2sum(2.0**54,3.0)
```

```
Out[102]: (1.8014398509481988e+16, -1.0)
```

```
In [103]: fast2sum(2.0**54,4.0)
```

```
Out[103]: (1.8014398509481988e+16, 0.0)
```

```
In [104]: fast2sum(2.0**54,5.0)
```

```
Out[104]: (1.8014398509481988e+16, 1.0)
```

```
In [100]: fast2sum(2.0**54,1.0)
```

```
Out[100]: (1.8014398509481984e+16, 1.0)
```

```
In [101]: fast2sum(2.0**54,2.0)
```

```
Out[101]: (1.8014398509481984e+16, 2.0)
```

```
In [102]: fast2sum(2.0**54,3.0)
```

```
Out[102]: (1.8014398509481988e+16, -1.0)
```

```
In [103]: fast2sum(2.0**54,4.0)
```

```
Out[103]: (1.8014398509481988e+16, 0.0)
```

```
In [104]: fast2sum(2.0**54,5.0)
```

```
Out[104]: (1.8014398509481988e+16, 1.0)
```

Théorème 7 (2Sum - Knuth)

On considère deux VF x et y et l'algorithme suivant :

1 $s \leftarrow x \oplus y$

2 $y_v \leftarrow s \ominus x$

3 $x_v \leftarrow s \ominus y_v$

4 $y_a \leftarrow y \ominus y_v$

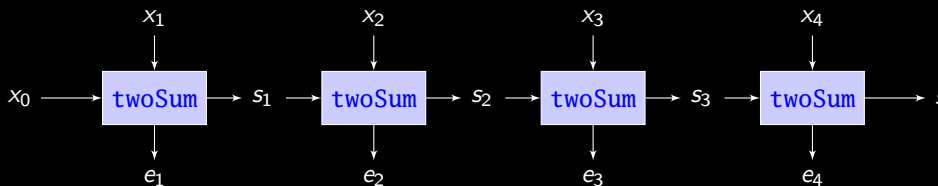
5 $x_a \leftarrow x \ominus x_v$

6 $d \leftarrow x_a \oplus y_a$

7 **Retourner** (s, d)

Alors $x + y = s + d$ avec $s = x \oplus y$ et $d = \text{err}(x \oplus y)$. De plus s et d ne se chevauchent pas.

Somme compensée d'un nombre quelconque de flottants



Somme compensée d'un nombre quelconque de flottants

```
def sommeKahan(liste):  
    (s,c) = (0.,0.)  
    for x in liste:  
        (s,c) = fast2sum(s, x + c)  
    return s
```

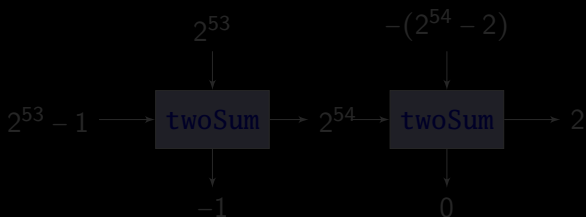
```
def sommePichat(liste):  
    s,e = 0.,0.  
    for x in liste:  
        (s,c) = fast2sum(s, x)  
        e += c  
    return s + e
```

Somme compensée d'un nombre quelconque de flottants

```
def sommeKahan(liste):
    (s,c) = (0.,0.)
    for x in liste:
        (s,c) = fast2sum(s, x + c)
    return s
```

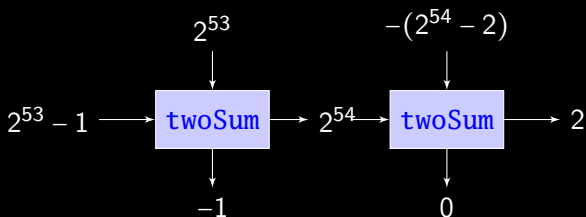
```
def sommePichat(liste):
    s,e = 0.,0.
    for x in liste:
        (s,c) = fast2sum(s, x)
        e += c
    return s + e
```

$$x_0 = 2^{53} - 1, \quad x_1 = 2^{53} \quad \text{et} \quad x_2 = -(2^{54} - 2)$$



Expliquez les résultats trouvés. Que va faire l'algorithme de somme compensée ? Et celui de somme en cascade ?

$$x_0 = 2^{53} - 1, \quad x_1 = 2^{53} \quad \text{et} \quad x_2 = -(2^{54} - 2)$$



Exercice 6

Expliquez les résultats trouvés. Que va faire l'algorithme de somme compensée ? Et celui de somme en cascade ?


```
In [105]: sommePichat([1. / n for n in range(1,100001)])
```

```
Out[105]: 12.090146129863427
```

```
In [106]: sommeKahan([1. / n for n in range(1,100001)])
```

```
Out[106]: 12.090146129863427
```

```
In [107]: sum([1. / n for n in range(1,100001)])
```

```
Out[107]: 12.090146129863335
```

```
In [108]: sum([1. / n for n in range(100000,0,-1)])
```

```
Out[108]: 12.090146129863408
```

```
sage: sum(1. / n, n , 1 , 100000).n(64)
```

```
12.0901461298634279
```

```
In [105]: sommePichat([1. / n for n in range(1,100001)])
```

```
Out[105]: 12.090146129863427
```

```
In [106]: sommeKahan([1. / n for n in range(1,100001)])
```

```
Out[106]: 12.090146129863427
```

```
In [107]: sum([1. / n for n in range(1,100001)])
```

```
Out[107]: 12.090146129863335
```

```
In [108]: sum([1. / n for n in range(100000,0,-1)])
```

```
Out[108]: 12.090146129863408
```

```
sage: sum(1. / n, n , 1 , 100000).n(64)
```

```
12.0901461298634279
```

Banque chaotique

Exemple dû à Jean-Michel Muller.

M. X a récemment été à sa banque (Chaotic Bank Society), pour connaître les nouvelles offres proposées aux meilleurs clients. Son banquier lui propose l'offre suivante : « vous déposez tout d'abord $e - 1$ euros sur votre compte (où $e = 2.7182818\dots$ est la base du logarithme népérien). La première année, nous prenons 1 euro sur votre compte de frais de gestion. Par contre, la deuxième année est plus intéressante pour vous, car nous multiplions votre capital restant par 2 et prenons 1 euro de frais de gestion. La troisième année est encore plus intéressante, car nous multiplions votre capital par 3 et prenons 1 euro de frais de gestion. Et ainsi de suite : la n -ième année, nous multiplions votre capital par n et prenons 1 euro de frais de gestion. Intéressant, non ? » Pour prendre sa décision, M. X décide de demander l'aide d'un informaticien et se retourne vers vous.

Banque chaotique

Exemple dû à Jean-Michel Muller.

M. X a récemment été à sa banque (Chaotic Bank Society), pour connaître les nouvelles offres proposées aux meilleurs clients. Son banquier lui propose l'offre suivante : « vous déposez tout d'abord $e - 1$ euros sur votre compte (où $e = 2.7182818\dots$ est la base du logarithme népérien). La première année, nous prenons 1 euro sur votre compte de frais de gestion. Par contre, la deuxième année est plus intéressante pour vous, car nous multiplions votre capital restant par 2 et prenons 1 euro de frais de gestion. La troisième année est encore plus intéressante, car nous multiplions votre capital par 3 et prenons 1 euro de frais de gestion. Et ainsi de suite : la n -ième année, nous multiplions votre capital par n et prenons 1 euro de frais de gestion. Intéressant, non ? » Pour prendre sa décision, M. X décide de demander l'aide d'un informaticien et se retourne vers vous.

```
def chaotic(n):
    cpt = exp(1) - 1
    for i in range(1, n + 1):
        cpt = i*cpt - 1
    return cpt
```

$$\begin{aligned}
 c_n &= n! \times \left(c_0 - 1 - \frac{1}{2!} - \frac{1}{3!} - \dots - \frac{1}{n!} \right) \\
 &= n! \times (c_0 - (e - 1) + \frac{1}{(n+1)!} + \frac{1}{(n+2)!} + \dots)
 \end{aligned}$$

```
def chaotic(n):
    cpt = exp(1) - 1
    for i in range(1, n + 1):
        cpt = i*cpt - 1
    return cpt
```

$$\begin{aligned}
 c_n &= n! \times \left(c_0 - 1 - \frac{1}{2!} - \frac{1}{3!} - \dots - \frac{1}{n!} \right) \\
 &= n! \times \left(c_0 - (e - 1) + \frac{1}{(n+1)!} + \frac{1}{(n+2)!} + \dots \right)
 \end{aligned}$$