

Méthodes numériques I : IEEE 754

INFO1 - Semaines 18 & 19

Guillaume CONNAN

Dernière mise à jour : 3 mai 2015 à 11:02

IUT de Nantes - Dpt d'informatique

Sommaire

- 1 Preamble
- 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
- 6 MPFR

Sommaire

1 Preamble

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

6 MPFR











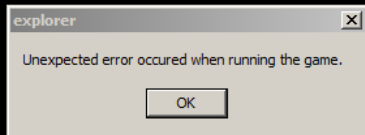
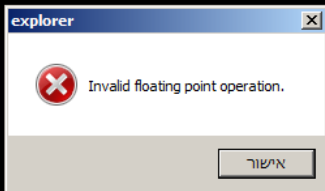
Microsoft Excel - Classeur1

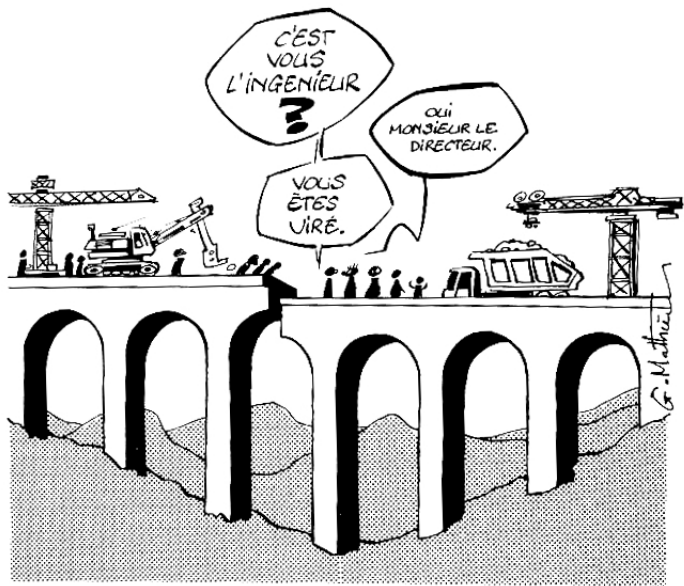
Fichier Edition Affichage Insertion Format Outils Do

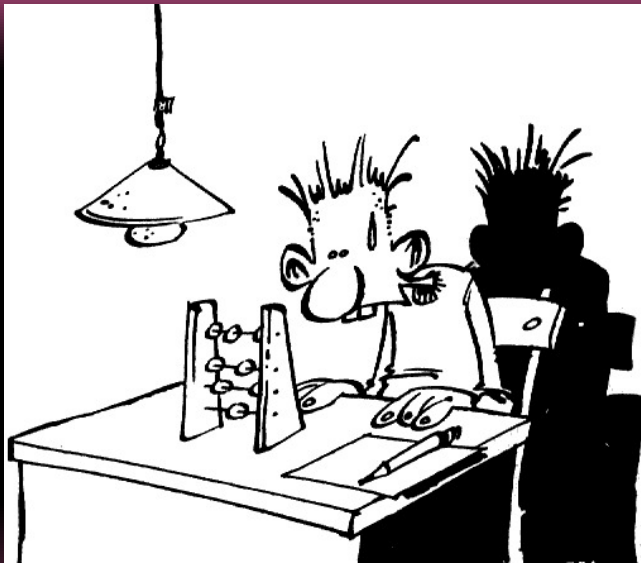
Arial 10 G I S

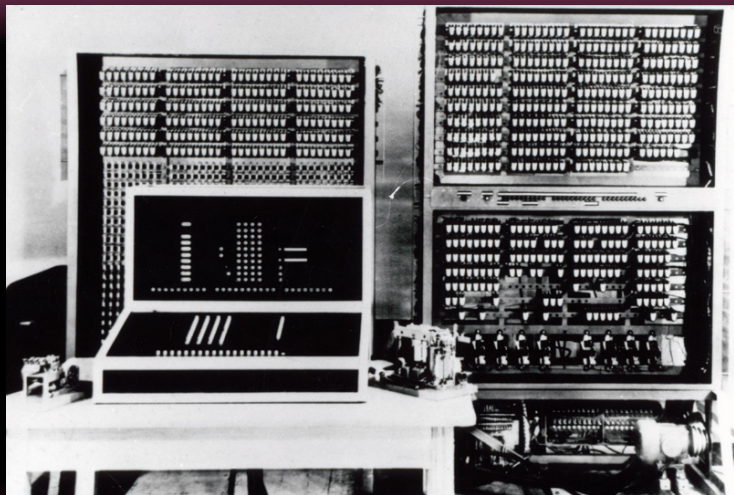
C15 fx

	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 $2147483647 + 1 = -2147483648\dots$
- short long

- char
- int $2147483647 + 1 = -2147483648\dots$
- short long

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

* float

* double

- char
- int $2147483647 + 1 = -2147483648\dots$
- short long
- float
- single
- double
- long double
- \dots

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

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

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

- char
- int $2147483647 + 1 = -2147483648\dots$
- short long
- float
- single
- double
- long double
- 3 3.0 3.0000 3.0e+00 3.0f

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

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

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

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

```
*Main> 3 * 0.1
```

-
-
-


```
*Main> 3 * 0.1  
0.30000000000000004
```

```
▪  
▪
```

```
*Main> 3 * 0.1
```

```
0.30000000000000004
```

```
*Main> sum $ take 10000000 $ repeat 0.1
```

```
.
```

```
*Main> 3 * 0.1
```

```
0.30000000000000004
```

```
*Main> sum $ take 10000000 $ repeat 0.1
```

```
999999.9998389754
```

LES NOMBRES RÉELS
N'EXISTENT PAS.

TOUT CE QUE VOUS AVEZ VU
AU LYCÉE N'EST
QU'ILLUSION !

LES NOMBRES RÉELS
N'EXISTENT PAS.

TOUT CE QUE VOUS AVEZ VU
AU LYCÉE 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
- 6 MPFR



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

* binary 32

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

- *binary32* ($\#E, \#m$) = (8, 24)

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

- *binary32* ($\#E, \#m$) = (8, 24)

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

- *binary32* ($\#E, \#m$) = (8, 24)

- *binary64*

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

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

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

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

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

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

• *float*

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

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

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

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

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

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

• *toy8*

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

- *binary32* ($\#E, \#m$) = (8, 24)
- *binary64* ($\#E, \#m$) = (11, 53).
- *toy7* ($\#E, \#m$) = (3, 4).
- *toy8* ($\#E, \#m$) = (3, 5).

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

- *binary32* ($\#E, \#m$) = (8, 24)
- *binary64* ($\#E, \#m$) = (11, 53).
- *toy7* ($\#E, \#m$) = (3, 4).
- *toy8* ($\#E, \#m$) = (3, 5).

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

- *binary32* ($\#E, \#m$) = (8, 24)
- *binary64* ($\#E, \#m$) = (11, 53).
- *toy7* ($\#E, \#m$) = (3, 4).
- *toy8* ($\#E, \#m$) = (3, 5).

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

- *binary32* ($\#E, \#m$) = (8, 24)
- *binary64* ($\#E, \#m$) = (11, 53).
- *toy7* ($\#E, \#m$) = (3, 4).
- *toy8* ($\#E, \#m$) = (3, 5).

$$v = 0,01001 \times 2^0$$

$$v = 0,1001 \times 2^{-1}$$

$$v = 1,001 \times 2^{-2}$$

$$1 \leq m < 2$$

$$m = 1, f$$

$$v = 0,01001 \times 2^0$$

$$v = 0,1001 \times 2^{-1}$$

$$v = 1,001 \times 2^{-2}$$

$$1 \leq m < 2$$

$$m = 1, f$$

$$v = 0,01001 \times 2^0$$

$$v = 0,1001 \times 2^{-1}$$

$$v = 1,001 \times 2^{-2}$$

$$1 \leq m < 2$$

$$m = 1, f$$

On peut également gagner de la place en ne stockant pas le signe de l'exposant : il suffit de le translater de $2^{(\#E)-1} - 1$...POURQUOI?

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

On peut également gagner de la place en ne stockant pas le signe de l'exposant : il suffit de le translater de $2^{(\#E)-1} - 1$...POURQUOI?

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

On peut également gagner de la place en ne stockant pas le signe de l'exposant : il suffit de le translater de $2^{(\#E)-1} - 1$...POURQUOI?

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

On peut également gagner de la place en ne stockant pas le signe de l'exposant : il suffit de le translater de $2^{(\#E)-1} - 1$...POURQUOI?

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

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)	f (4 bits)
0	010	1000

Soit `0b00101000` stocké en mémoire sous forme hexadécimale : `0x28`

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)	f (4 bits)
0	10	1000

Soit `0b00101000` stocké en mémoire sous forme hexadécimale : `0x28`

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit `0b00101000` stocké en mémoire sous forme hexadécimale : `0x28`

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit `0b00101000` stocké en mémoire sous forme hexadécimale : `0x28`

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit 0000101000 stocké en mémoire sous forme hexadécimale : 0x28

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit 0600101000 stocké en mémoire sous forme hexadécimale : 0x28

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit `0b00101000` stocké en mémoire sous forme hexadécimale : `0x28`

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit `0b00101000` stocké en mémoire sous forme hexadécimale : `0x28`

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit 0b00101000 stocké en mémoire sous forme hexadécimale : 0x28

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit $0b00101000$ stocké en mémoire sous forme hexadécimale : 0×28

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit $0b00101000$ stocké en mémoire sous forme hexadécimale : $0x28$

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit 0600101000 stocké en mémoire sous forme hexadécimale : 0x28

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit $0b00101000$ stocké en mémoire sous forme hexadécimale : 0×28

Exemple en toy8

$0,75_{10}$ en *toy8*.

$$0,75_{10} = \frac{3_{10}}{2_{10}^2} = 11_2 \times 2^{-2} = 1,1000_2 \times 2^{-1}$$

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

s (1 bit)	e (3 bits)			f (4 bits)			
0	0	1	0	1	0	0	0

Soit $0b00101000$ stocké en mémoire sous forme hexadécimale : $0x28$

Exemple en binary32

Et en *binary32* ?

QUELQUES HUMANOÏDES RARES



CENTRALIEN



NORMALIEN



ALIEN

$$e_{\min} = 0$$

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

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

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

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

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

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

$$e_{\min} = 0$$

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

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

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

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

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

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

$$e_{\min} = 0$$

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

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

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

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

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

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

$$e_{\min} = 0$$

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

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

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

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

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

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

$$e_{\min} = 0$$

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

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

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

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

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

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

$$e_{\min} = 0$$

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

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

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

Les valeurs extrêmes de E sont réservées pour des cas spéciaux

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

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

$$e_{\min} = 0$$

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

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

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

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

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

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

$$e_{\min} = 0$$

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

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

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

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 en 1D

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

Sous-normaux...en TD

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

Sous-normaux...en TD

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

Sous-normaux...en TD

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

Sous-normaux....en TD

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


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

```
*Main> x
```

```
·  
·  
·  
·  
·  
·  
·  
·  
·  
·  
·  
·  
·  
·  
·
```

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

```
*Main> x
```

```
0.0
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

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

```
*Main> x
```

```
0.0
```

```
*Main> x == 0
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```



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

```
*Main> x
```

```
0.0
```

```
*Main> x == 0
```

```
True
```

```
*Main> let y = -1 / 1e500
```

```
*Main> y
```

```
-0.0
```

```
.  
.   
.   
.
```

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

```
*Main> x
```

```
0.0
```

```
*Main> x == 0
```

```
True
```

```
*Main> let y = -1 / 1e500
```

```
*Main> y
```

```
-0.0
```

```
*Main> y == 0
```

```
.
```

```
.
```

```
.
```

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

```
*Main> x
```

```
0.0
```

```
*Main> x == 0
```

```
True
```

```
*Main> let y = -1 / 1e500
```

```
*Main> y
```

```
-0.0
```

```
*Main> y == 0
```

```
True
```

```
.
```

```
.
```



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

```
*Main> x
```

```
0.0
```

```
*Main> x == 0
```

```
True
```

```
*Main> let y = -1 / 1e500
```

```
*Main> y
```

```
-0.0
```

```
*Main> y == 0
```

```
True
```

```
*Main> x == y
```

```
.
```

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

```
*Main> x
```

```
0.0
```

```
*Main> x == 0
```

```
True
```

```
*Main> let y = -1 / 1e500
```

```
*Main> y
```

```
-0.0
```

```
*Main> y == 0
```

```
True
```

```
*Main> x == y
```

```
True
```

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1,1 \times 2^3$$

toy7

$$1,111 \times 2^3 = 7$$

x (1 bit) + y (3 bits) = z (3 bits)

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

$$x = 1,111 \times 2^3$$

$$y = 1,111 \times 2^3$$

$$x^2 + y^2 = 1,111 \times 2^4$$

$$N = (1,111 \times 2^4)^{1/2} \approx 1,010 \times 2^2$$

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1, 1 \times 2^3$$

toy7

$$1, 111 \times 2^4 ?$$

$x = (1 \text{ bit}) \times e (3 \text{ bits}) = f (3 \text{ bits})$

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

$$x = 1, 111 \times 2^3$$

$$y = 1, 111 \times 2^3$$

$$x^2 + y^2 = 1, 111 \times 2^4$$

$$N = (1, 111 \times 2^4)^{1/2} \approx 1, 010 \times 2^2$$

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1,1 \times 2^3$$

toy7

$$1,111 \times 2^4 ?$$

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

$$x = 1,111 \times 2^3$$

$$y = 1,110 \times 2^3$$

$$x^2 + y^2 = 1,111 \times 2^4$$

$$N = (1,111 \times 2^4)^{1/2} \approx 1,010 \times 2^2$$

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1, 1 \times 2^3$$

toy7

$$1, 111 \times 2^4 ?$$

<i>s</i> (1 bit)	<i>e</i> (3 bits)			<i>f</i> (3 bits)		
0	1	1	1	1	1	1

$$x^2 = 1, 111 \times 2^4$$

$$y^2 = 1, 111 \times 2^4$$

$$x^2 + y^2 = 1, 111 \times 2^4$$

$$N = (1, 111 \times 2^4)^{1/2} \approx 1, 010 \times 2^2$$

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1, 1 \times 2^3$$

toy7

$$1, 111 \times 2^4 ?$$

<i>s</i> (1 bit)	<i>e</i> (3 bits)			<i>f</i> (3 bits)		
0	1	1	1	1	1	1

$$x^2 = 1, 111 \times 2^4$$

$$y^2 = 1, 111 \times 2^4$$

$$x^2 + y^2 = 1, 111 \times 2^4$$

$$N = (1, 111 \times 2^4)^{1/2} \approx 1, 010 \times 2^2$$

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1,1 \times 2^3$$

toy7

$$1,111 \times 2^4 ?$$

<i>s</i> (1 bit)	<i>e</i> (3 bits)			<i>f</i> (3 bits)		
0	1	1	1	1	1	1

$$x^2 = 1,111 \times 2^4$$

$$y^2 = 1,111 \times 2^4$$

$$x^2 + y^2 = 1,111 \times 2^4$$

$$N = (1,111 \times 2^4)^{1/2} \approx 1,010 \times 2^2$$

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1,1 \times 2^3$$

toy7

$$1,111 \times 2^4 ?$$

<i>s</i> (1 bit)	<i>e</i> (3 bits)			<i>f</i> (3 bits)		
0	1	1	1	1	1	1

$$x^2 = 1,111 \times 2^4$$

$$y^2 = 1,111 \times 2^4$$

$$x^2 + y^2 = 1,111 \times 2^4$$

$$N = (1,111 \times 2^4)^{1/2} \approx 1,010 \times 2^2$$

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1,1 \times 2^3$$

toy7

$$1,111 \times 2^4 ?$$

<i>s</i> (1 bit)	<i>e</i> (3 bits)			<i>f</i> (3 bits)		
0	1	1	1	1	1	1

$$x^2 = 1,111 \times 2^4$$

$$y^2 = 1,111 \times 2^4$$

$$x^2 + y^2 = 1,111 \times 2^4$$

$$N = (1,111 \times 2^4)^{1/2} \approx 1,010 \times 2^2$$

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1,1 \times 2^3$$

toy7

$$1,111 \times 2^4 ?$$

<i>s</i> (1 bit)	<i>e</i> (3 bits)			<i>f</i> (3 bits)		
0	1	1	1	1	1	1

$$x^2 = 1,111 \times 2^4$$

$$y^2 = 1,111 \times 2^4$$

$$x^2 + y^2 = 1,111 \times 2^4$$

$$N = (1,111 \times 2^4)^{1/2} \approx 1,010 \times 2^2$$

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1, 1 \times 2^3$$

toy7

$$1, 111 \times 2^4 ?$$

<i>s</i> (1 bit)	<i>e</i> (3 bits)			<i>f</i> (3 bits)		
0	1	1	1	1	1	1

$$x^2 = 1, 111 \times 2^4$$

$$y^2 = 1, 111 \times 2^4$$

$$x^2 + y^2 = 1, 111 \times 2^4$$

$$N = (1, 111 \times 2^4)^{1/2} \approx 1, 010 \times 2^2$$

$$N = \sqrt{x^2 + y^2}$$

$$x = 1 \times 2^3 \quad y = 1, 1 \times 2^3$$

toy7

$$1, 111 \times 2^4 ?$$

<i>s</i> (1 bit)	<i>e</i> (3 bits)			<i>f</i> (3 bits)		
0	1	1	1	1	1	1

$$x^2 = 1, 111 \times 2^4$$

$$y^2 = 1, 111 \times 2^4$$

$$x^2 + y^2 = 1, 111 \times 2^4$$

$$N = (1, 111 \times 2^4)^{1/2} \approx 1, 010 \times 2^2$$

$+\infty$:

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

$-\infty$:

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

$e = e_{\max}$ et $f = 0$

$+\infty$:

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

$-\infty$:

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

$e = e_{\max}$ et $f = 0$

$+\infty$:

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

$-\infty$:

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

$e = e_{\max}$ et $f = 0$


```
*Main> 1 + 1e500
```

```
Infinity
```

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

```
·  
·  
·  
·  
·  
·  
·  
·  
·  
·  
·  
·  
·  
·  
·
```



```
*Main> 1 + 1e500
```

```
Infinity
```

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

```
*Main> x
```

```
Infinity
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
*Main> 1 + 1e500
```

```
Infinity
```

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

```
*Main> x
```

```
Infinity
```

```
*Main> 1 + x
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
.
```

```
*Main> 1 + 1e500
```

```
Infinity
```

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

```
*Main> x
```

```
Infinity
```

```
*Main> 1 + x
```

```
Infinity
```

```
·  
·  
·  
·  
·  
·  
·
```

```
*Main> 1 + 1e500
```

```
Infinity
```

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

```
*Main> x
```

```
Infinity
```

```
*Main> 1 + x
```

```
Infinity
```

```
*Main> x^3
```

```
·
```

```
·
```

```
·
```

```
·
```

```
·
```



```
*Main> 1 + 1e500
```

```
Infinity
```

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

```
*Main> x
```

```
Infinity
```

```
*Main> 1 + x
```

```
Infinity
```

```
*Main> x^3
```

```
Infinity
```

```
·  
·  
·  
·
```

```
*Main> 1 + 1e500
```

```
Infinity
```

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

```
*Main> x
```

```
Infinity
```

```
*Main> 1 + x
```

```
Infinity
```

```
*Main> x^3
```

```
Infinity
```

```
*Main> 1 / x
```

```
·
```

```
·
```

```
·
```

```
*Main> 1 + 1e500
```

```
Infinity
```

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

```
*Main> x
```

```
Infinity
```

```
*Main> 1 + x
```

```
Infinity
```

```
*Main> x^3
```

```
Infinity
```

```
*Main> 1 / x
```

```
0.0
```

```
.
```

```
.
```

```
*Main> 1 + 1e500
```

```
Infinity
```

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

```
*Main> x
```

```
Infinity
```

```
*Main> 1 + x
```

```
Infinity
```

```
*Main> x^3
```

```
Infinity
```

```
*Main> 1 / x
```

```
0.0
```

```
*Main> 4 - x
```

```
.
```

```
*Main> 1 + 1e500
```

```
Infinity
```

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

```
*Main> x
```

```
Infinity
```

```
*Main> 1 + x
```

```
Infinity
```

```
*Main> x^3
```

```
Infinity
```

```
*Main> 1 / x
```

```
0.0
```

```
*Main> 4 - x
```

```
-Infinity
```

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

```
*Main> x^2 - x
```

```
.
```

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

```
*Main> x^2 - x
```

```
.
```

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

```
*Main> x^2 - x  
NaN
```




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

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

```
NaN
```

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

```
NaN
```

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

```
NaN
```

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

```
NaN
```

```
-- Bizarre ?
```

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

```
NaN
```

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

```
False
```

```
*Main> x^2 - x
```

```
NaN
```

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

```
Infinity
```



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

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

```
*Main> x^2 - x
```

```
NaN
```

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

```
Infinity
```

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

```
NaN
```

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

```
False
```

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

```
*Main> x^2 - x
```

```
NaN
```

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

```
Infinity
```

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

```
NaN
```

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

```
False
```

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

```
*Main> x^2 - x
```

```
NaN
```

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

```
Infinity
```

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

```
NaN
```

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

```
False
```



```
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
```

Air France 447 - 1^{er} juin 2009



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

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$v = M(v) = 2^{E(v)} \cdot v$$

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

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$v = M(v) = 2^{E(v)} \cdot 2^{F(v)}$$

$$\text{succ}(v) = (M(v) + 1) = 2^{E(v)} \cdot (1 + 2^{F(v)})$$

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$v = M(v) = 2^{E(v)} \cdot v$$

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

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

$$v = M(v) = 2^{\lfloor n/2 \rfloor} v$$

$$\text{succ}(v) = (M(v) + 1) = 2^{\lfloor n/2 \rfloor} (v + 2^{\lfloor n/2 \rfloor})$$

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

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

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

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

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

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

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

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

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

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

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

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

$$n = \#f$$

toy7, $n = 3$

binary32, $n = 23$

binary64, $n = 52$

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

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

Si v est un VF, entre v et son successeur, il y a...

Arbre de lecture

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

$$e = 0\dots0 \quad f = 0 \quad \rightarrow \quad v = (-1)^s \times 0.0$$

$$e = 0\dots0 \quad f \neq 0 \quad \rightarrow \quad v = (-1)^s \times 0.f \times 2^{1-E_{\max}}$$

$$e = 1\dots1 \quad f = 0 \quad \rightarrow \quad v = (-1)^s \times \infty$$

$$e = 1\dots1 \quad f \neq 0 \quad \rightarrow \quad \text{NaN}$$

$$0\dots0 < e < 1\dots1 \quad f \neq 0 \quad \rightarrow \quad v = (-1)^s \times 1.f \times 2^{e-E_{\max}}$$

Arbre de lecture

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

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 ?

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 ?

Format	#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} 5×10^{-324}	$(2^{53} - 1)2^{1023-5}$ $\approx 1,8 \times 10^{308}$

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
- 6 MPFR

$$\mathbb{V}_b$$

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

$$\mathbb{V}_b$$

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

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
0,0011

1,1011

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

0,0011

1,1011

1,1011

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.

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

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

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

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

oui

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

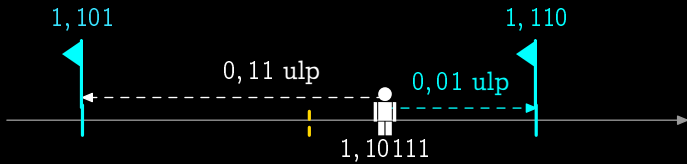
$x = 1,10111$ VF en *toy7*?

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

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

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

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



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

```
Prelude> (1 + 1e-16) + 1e-16
1.0
Prelude> 1 + (1e-16 + 1e-16)
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 ?

```
Prelude> (1 + 1e-16) + 1e-16  
1.0  
Prelude> 1 + (1e-16 + 1e-16)  
1.000000000000000002
```

- 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 :

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

-10,01 :

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

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

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

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

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

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

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

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

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

- 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 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

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

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

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

- 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

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

avec $|x - 1,100| = 0,001011011 \dots$ et $|x - 1,101| = 0,00011011 \dots$ donc

$$|x - 1,101| < |x - 1,100| \Rightarrow x \text{ est plus proche de } 1,101$$

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

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

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

- 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

$$1,100 < x < 1,1005$$

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

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

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

- 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

$$1,100 < x < 1,1005$$

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

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

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

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

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

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

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

```
Prelude> (1.000000001 * 1e-15) * 1e15
1.000000001000000002
Prelude> 1.000000001 * (1e-15 * 1e15)
1.000000001
```

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

```
Prelude> 1e15 * (1e-16 + 1)
1.0e15
Prelude> (1e15 * 1e-16) + (1e15 * 1)
1.00000000000000001e15
```

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

```
Prelude> (1.000000001 * 1e-15) * 1e15
1.000000001000000002
Prelude> 1.000000001 * (1e-15 * 1e15)
1.000000001
```

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

```
Prelude> 1e15 * (1e-16 + 1)
1.0e15
Prelude> (1e15 * 1e-16) + (1e15 * 1)
1.00000000000000001e15
```


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

```
Prelude> (1.000000001 * 1e-15) * 1e15
1.000000001000000002
Prelude> 1.000000001 * (1e-15 * 1e15)
1.000000001
```

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

```
Prelude> 1e15 * (1e-16 + 1)
1.0e15
Prelude> (1e15 * 1e-16) + (1e15 * 1)
1.000000000000000001e15
```

Sommaire

- 1 Preamble
- 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
- 6 MPFR

```
f1 a =  
  if (a + 1.0) - a /= 1.0 then a  
  else f1 (2.0 * a)  
  
f2 a b =  
  if (a + b) - a == b then b  
  else f2 a (b + 1.0)  
  
b = f2 (f1 1) 1
```

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

```
f1 a =  
  if (a + 1.0) - a /= 1.0 then a  
  else f1 (2.0 * a)  
  
f2 a b =  
  if (a + b) - a == b then b  
  else f2 a (b + 1.0)  
  
b = f2 (f1 1) 1
```

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

SAN-ANTONIO

San-Antonio



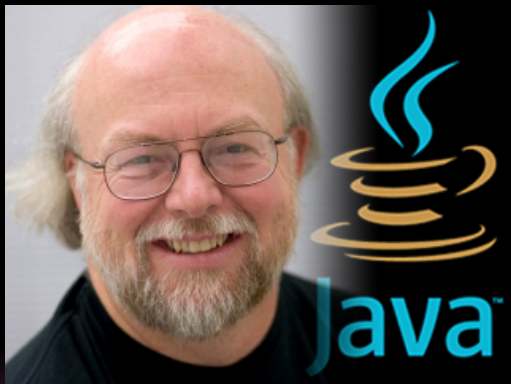
Si MA
TANTE en
AVAIT

85

FLEU
VE
NOIR

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.

3.177777777777 est une approximation plutôt précise (10 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.

3.177777777777 est une approximation plutôt précise (10 décimales) mais imprecise (2 décimales) en fait.

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.

3,1777777777777777 est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de π .

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

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

l'erreur relative

commise en prenant \hat{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 $\eta = \frac{x - \widehat{x}}{x}$

comme on le voit, on passe de l'erreur absolue à l'erreur relative en prenant x à la place de

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 .

$$r = f(d)$$

Définition 2

- L'erreur *aval* (*forward error*) est la différence entre le résultat mesuré \widehat{r} et le résultat théorique r .
- L'erreur *amont*, ou erreur *inverse* (*backward error*) est le plus petit $|\Delta d|$ tel que $f(d + \Delta d) = \widehat{r}$.

$$r = f(d)$$

Définition 2

- L'erreur *aval* (*forward error*) est la différence entre le résultat mesuré \widehat{r} et le résultat théorique r .
- L'erreur *amont*, ou erreur *inverse* (*backward error*) est le plus petit $|\Delta d|$ tel que $f(d + \Delta d) = \widehat{r}$.

$$r = f(d)$$

Définition 2

- L'erreur *aval* (*forward error*) est la différence entre le résultat mesuré \hat{r} et le résultat théorique r .
- L'erreur *amont*, ou erreur *inverse* (*backward error*) est le plus petit $|\Delta d|$ tel que $f(d + \Delta d) = \hat{r}$.

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**
- 6 MPFR



Précision machine

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

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

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

- 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

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

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

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

- 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

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

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

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

- 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

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

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

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

- 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

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

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

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

- 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

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

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

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

- 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

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

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

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

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

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_ε
MACHINE !

Précision machine

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

avec :

- si \widehat{x} est normal $|\eta_1| \leq \frac{1}{2}\epsilon_M$ 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

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

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

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

```
deriv :: (Float -> Float) -> Float -> Float -> Float
deriv  f                x      h      =
  ( f(x + h) - f(x) ) / h
```

```
derivn :: (Float -> Float) -> Float -> Float -> Int -> Float
derivn  f                x      _      0  = f x
derivn  f                x      h      n  =
  derivn (\ t -> deriv f t h) x h (n - 1)
```

Catastrophe

```
deriv :: (Float -> Float) -> Float -> Float -> Float
deriv  f          x          h      =
  ( f(x + h) - f(x) ) / h
```

```
derivn :: (Float -> Float) -> Float -> Float -> Int -> Float
derivn  f          x          _      0 = f x
derivn  f          x          h      n =
  derivn (\ t -> deriv f t h) x h (n - 1)
```

Catastrophe

```
R> [ [ derivn (\ t -> t**4) 1 (10**(-d)) k | k <- [0..5]] | d <-  
  [1..8]]
```

```
[[1.0,4.641001,14.540005,27.60004,24.001312,-3.6239624e-2],  
 [1.0,4.0603995,12.24041,24.557114,0.0,1192.0929],  
 [1.0,4.0061474,12.159347,-238.41856,476837.1,-8.3446483e8],  
 [1.0,4.001856,0.0,238418.58,-4.7683717e9,9.536743e13],  
 [1.0,4.005432,0.0,0.0,0.0,0.0],  
 [1.0,3.8146973,0.0,0.0,0.0,0.0],  
 [1.0,4.7683716,0.0,0.0,0.0,0.0],  
 [1.0,0.0,0.0,0.0,0.0,0.0]]
```

Catastrophe

```
R> [ [ derivn (\ t -> t**4) 1 (10**(-d)) k | k <- [0..5]] | d <-  
  [1..8]]
```

```
[[1.0,4.641001,14.540005,27.60004,24.001312,-3.6239624e-2],  
 [1.0,4.0603995,12.24041,24.557114,0.0,1192.0929],  
 [1.0,4.0061474,12.159347,-238.41856,476837.1,-8.3446483e8],  
 [1.0,4.001856,0.0,238418.58,-4.7683717e9,9.536743e13],  
 [1.0,4.005432,0.0,0.0,0.0,0.0],  
 [1.0,3.8146973,0.0,0.0,0.0,0.0],  
 [1.0,4.7683716,0.0,0.0,0.0,0.0],  
 [1.0,0.0,0.0,0.0,0.0,0.0]]
```

Catastrophe

```
deriv :: (Double -> Double) -> Double -> Double -> Double
deriv  f          x      h      =
  ( f(x + h) - f(x) ) / h
```

```
derivn :: (Double -> Double) -> Double -> Double -> Int -> Double
derivn  f          x      _      0 = f x
derivn  f          x      h      n =
  derivn (\ t -> deriv f t h) x h (n - 1)
```

Catastrophe

```
deriv :: (Double -> Double) -> Double -> Double -> Double
deriv  f                x      h      =
  ( f(x + h) - f(x) ) / h
```

```
derivn :: (Double -> Double) -> Double -> Double -> Int -> Double
derivn  f                x      -      0      = f x
derivn  f                x      h      n      =
  derivn (\ t -> deriv f t h) x h (n - 1)
```

Catastrophe

```
> [ [ derivn (\ t -> t**4) 1 (10**(-d)) k | k <- [0..5]] | d <-  
  [1..8]  
[  
[1.0,4.6410000000000004  
  ,14.5400000000000042,27.599999999999874,24.00000000000002,1.06581  
[1.0,4.0604010000000002  
  ,12.241399999999292,24.360000000278603,23.99999994295854,8.8817  
[1.0,4.006004000999486 ,12.024014000244776,24.03599941303014  
  ,24.001245435556484,-2.4424906541753444],  
[1.0,4.000600039999469 ,12.002400162636206,24.00324383700081  
  ,28.86579864025407,-66613.38147750939],  
[1.0,4.0000600004308495,12.00023858061172 ,24.20286193682841  
  ,-22204.46049250313,2.220446049250313e9],  
[1.0,4.000005999760248  
  ,11.999956583963467,0.0,2.220446049250313e8,-6.66133814775094e1  
[1.0,4.000000601855902  
  ,12.012613126444194,-222044.6049250313,2.220446049250313e12,0.0  
[1.0,4.000000042303498  
  ,11.102230246251565,0.0,2.2204460492503128e16,-4.44089209850062
```

Catastrophe

```
> [ [ derivn (\ t -> t**4) 1 (2**(-d)) k | k <- [0..5]] | d <-  
  [1..15]]  
[  
 [1.0,8.125,27.5,42.0,24.0,0.0],  
 [1.0,5.765625,18.875,33.0,24.0,0.0],  
 [1.0,4.814453125,15.21875,28.5,24.0,0.0],  
 [1.0,4.390869140625,13.5546875,26.25,24.0,0.0],  
 [1.0,4.191436767578125,12.763671875,25.125,24.0,0.0],  
 [1.0,4.094730377197266,12.37841796875,24.5625,24.0,0.0],  
 [1.0,4.047119617462158,12.1883544921875,24.28125,24.0,0.0],  
 [1.0,4.023498594760895,12.093963623046875,24.140625,24.0,0.0],  
 [1.0,4.011734016239643,12.046928405761719,24.0703125,24.0,0.0],  
 [1.0,4.005863190628588,12.02345085144043,24.03515625,24.0,0.0],  
 [1.0,4.002930641290732,12.011722087860107,24.017578125,24.0,0.0],  
 [1.0,4.001465082183131,12.005860209465027,24.0087890625,24.0,0.0],  
 [1.0,4.000732481481464,12.002929896116257,24.00439453125,24.0,0.0],  
 [1.0,4.000366225838661,12.001464903354645,24.001953125,32.0,-262144.0],  
 [1.0,4.00018310919404,12.000732421875,24.0,256.0,-2.5165824e7]]  
]
```


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!

Recherche 1

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 $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$?*

Recherche 1

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 $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$?*

Recherche 1

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 $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$?*

Recherche 1

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 $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$?*

Recherche 1

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 $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$?*

Recherche 1

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 ?*

Recherche 2 (Max IEEE)

Fonction $\text{max}(x, y : \text{flottants}) : \text{flottant}$

Si $x \geq y$ Alors

 | $\text{max} = x$

Sinon

 | $\text{max} = y$

FinSi

```
λ> max 1 (1 + 1e-16)
1.0
λ> max 0.1 (1e500 / 1e500)
0.1
λ> 0.1 <= 1e500/1e500
False
λ> 0.1 > 1e500/1e500
False
```



JavaScript

JavaScript

1.7 3 2012 05 18

```
js> Math.pow(2, 53)
9007199254740992
js> Math.pow(2, 53) + 1
9007199254740992
js> Math.pow(2, 53) + 2
9007199254740994
js> Math.pow(2, 53) + 3
9007199254740996
js> Math.pow(2, 53) + 4
9007199254740996
js> Math.pow(2, 53) + 5
9007199254740996
```

JavaScript

```
js> Math.pow(2, 54)
18014398509481984
js> Math.pow(2, 54) + 2
18014398509481984
js> Math.pow(2, 54) + 3
18014398509481988

js> Math.pow(2, 55)
36028797018963970
Math.pow(2, 55) + 4
36028797018963970
```

JavaScript

JavaScript

1.7 3 2012 05 18

```
js> Math.pow(2, 53)
9007199254740992
js> Math.pow(2, 53) + 1
9007199254740992
js> Math.pow(2, 53) + 2
9007199254740994
js> Math.pow(2, 53) + 3
9007199254740996
js> Math.pow(2, 53) + 4
9007199254740996
js> Math.pow(2, 53) + 5
9007199254740996
```

JavaScript

```
js> Math.pow(2, 54)
18014398509481984
js> Math.pow(2, 54) + 2
18014398509481984
js> Math.pow(2, 54) + 3
18014398509481988

js> Math.pow(2, 55)
36028797018963970
Math.pow(2, 55) + 4
36028797018963970
```

```
1  Algorithme de Malcolm-Gentleman
2   $a \leftarrow 1.0$ 
3   $b \leftarrow 1.0$ 
4  TantQue  $R(R(a + 1.0) - a) == 1.0$  Faire
5  |    $a \leftarrow R(2 \times a)$ 
6  FinTantQue
7  TantQue  $R(R(a + b) - a) \neq b$  Faire
8  |    $b \leftarrow b + 1$ 
9  FinTantQue
10 Retourner  $b$ 
```

```
1  Algorithme de Malcolm-Gentleman
2   $a \leftarrow 1.0$ 
3   $b \leftarrow 1.0$ 
4  TantQue  $R(R(a + 1.0) - a) == 1.0$  Faire
5  |    $a \leftarrow R(2 \times a)$ 
6  FinTantQue
7  TantQue  $R(R(a + b) - a) \neq b$  Faire
8  |    $b \leftarrow b + 1$ 
9  FinTantQue
10 Retourner  $b$ 
```



```

base = boucle2 (boucle1 1) 1
  where
    boucle1 a =
      if (a + 1.0) - a /= 1.0 then a
      else boucle1 (2.0 * a)
    boucle2 a b =
      if (a + b) - a == b then b
      else boucle2 a (b + 1.0)

```

$a_{i+1} = 2a_i$ $a_0 = 1.0$ $a_i = 2^i$ - Monoton

Dans quel espace je travaille, je me dirais et si la proposition est vraie je suis

Mathematician - Haskell/GH

En fait $a_{i+1} = \wedge (2 \times a_i) \dots$

```

base = boucle2 (boucle1 1) 1
  where
    boucle1 a =
      if (a + 1.0) - a /= 1.0 then a
      else boucle1 (2.0 * a)
    boucle2 a b =
      if (a + b) - a == b then b
      else boucle2 a (b + 1.0)

```

$a_{i+1} = 2a_i$ $a_0 = 1.0$ $a_i = 2^i \dots$ Mouais

Dans quel espace tu travailles tu me diras et si ta proposition est vraie je saurai.

Mathematica 4.1.0p.10

En fait $a_{i+1} = 2(2 \times a_i) \dots$

```
base = boucle2 (boucle1 1) 1
where
  boucle1 a =
    if (a + 1.0) - a /= 1.0 then a
    else boucle1 (2.0 * a)
  boucle2 a b =
    if (a + b) - a == b then b
    else boucle2 a (b + 1.0)
```

$a_{i+1} = 2a_i$ $a_0 = 1.0$ $a_i = 2^i$ Mouais

Dans quel espace tu travailles tu me diras et si ta proposition est vraie je saurai

Mathemator - 44 ap. GC

En fait $a_{i+1} = 2 \times a_i$...

```

base = boucle2 (boucle1 1) 1
  where
    boucle1 a =
      if (a + 1.0) - a /= 1.0 then a
      else boucle1 (2.0 * a)
    boucle2 a b =
      if (a + b) - a == b then b
      else boucle2 a (b + 1.0)

```

$a_{i+1} = 2a_i$ $a_0 = 1.0$ $a_i = 2^i$ Mouais

Dans quel espace tu travailles tu me diras et si ta proposition est vraie je saurai

Mathemator - 44 ap. GC

Exercice 10 (2010)

```

base = boucle2 (boucle1 1) 1
  where
    boucle1 a =
      if (a + 1.0) - a /= 1.0 then a
      else boucle1 (2.0 * a)
    boucle2 a b =
      if (a + b) - a == b then b
      else boucle2 a (b + 1.0)

```

$a_{i+1} = 2a_i$ $a_0 = 1.0$ $a_i = 2^i$ Mouais

Dans quel espace tu travailles tu me diras et si ta proposition est vraie je saurai

Mathemator - 44 ap. GC

En fait $a_{i+1} = R(2 \times a_i)$...

```
base = boucle2 (boucle1 1) 1
where
  boucle1 a =
    if (a + 1.0) - a /= 1.0 then a
    else boucle1 (2.0 * a)
  boucle2 a b =
    if (a + b) - a == b then b
    else boucle2 a (b + 1.0)
```

$a_{i+1} = 2a_i$ $a_0 = 1.0$ $a_i = 2^i$ Mouais

Dans quel espace tu travailles tu me diras et si ta proposition est vraie je saurai

Mathemator - 44 ap. GC

En fait $a_{i+1} = R(2 \times a_i)$...

```
base = boucle2 (boucle1 1) 1
where
  boucle1 a =
    if (a + 1.0) - a /= 1.0 then a
    else boucle1 (2.0 * a)
  boucle2 a b =
    if (a + b) - a == b then b
    else boucle2 a (b + 1.0)
```

$a_{i+1} = 2a_i$ $a_0 = 1.0$ $a_i = 2^i$ Mouais

Dans quel espace tu travailles tu me diras et si ta proposition est vraie je saurai

Mathemator - 44 ap. GC

En fait $a_{i+1} = R(2 \times a_i)$...

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$

• on sort donc de la boucle.

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$

● on sort donc de la boucle.

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.

● On suit donc de la droite :

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.

• Première itération i telle que $2^i > \beta^p$.

• On a donc $a_i < \beta^p$.

• On a donc $a_i + 1.0 < \beta^p$.

• On a donc $a_i + 1.0 - a_i = 1.0$.

• On a donc $R(a_i + 1.0 - a_i) = 1.0$.

• On a donc $R(R(a_i + 1.0) - a_i) = 1.0$.

• On sort donc de la boucle.

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- 4 Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- 5 $a_{i_s} = R(2 \times a_{i_s-1})$

● on sort donc de la boucle.

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- 4 Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- 5 $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) = 2 \times \beta^p$.

6 on sort donc de la boucle.

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- 4 Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- 5 $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p)$

6 on sort donc de la boucle.

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- 4 Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- 5 $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$

6 on sort donc de la boucle.

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- 4 Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- 5 $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$

6 On sort donc de la boucle.

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- 4 Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- 5 $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$

➊ On sort donc de la boucle :

- ① Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- ② Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- ③ $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- ④ Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- ⑤ $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$
- ⑥ $\beta^p \leq a_{i_s} < \beta^{p+1}$
- ⑦ $\text{succ}(a_{i_s}) = a_{i_s} + \beta^{E(a_{i_s}-p+1)}$

⑧ On s'arrête de la boucle.

- ① Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- ② Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- ③ $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- ④ Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- ⑤ $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$
- ⑥ $\beta^p \leq a_{i_s} < \beta^{p+1}$
- ⑦ $\text{succ}(a_{i_s}) = a_{i_s} + \beta^{E(a_{i_s}-p+1)} = a_{i_s} + \beta^{p-p+1} = a_{i_s} + \beta$

⑧ on sort donc de la boucle.

- ① Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- ② Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- ③ $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- ④ Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- ⑤ $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$
- ⑥ $\beta^p \leq a_{i_s} < \beta^{p+1}$
- ⑦ $\text{succ}(a_{i_s}) = a_{i_s} + \beta^{E(a_{i_s}-p+1)} = a_{i_s} + \beta^{p-p+1} = a_{i_s} + \beta$

⑧ On sort donc de la boucle.

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- 4 Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- 5 $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$
- 6 $\beta^p \leq a_{i_s} < \beta^{p+1}$
- 7 $\text{succ}(a_{i_s}) = a_{i_s} + \beta^{E(a_{i_s}-p+1)} = a_{i_s} + \beta^{p-p+1} = a_{i_s} + \beta$

8 $\beta^p \leq a_{i_s} < \beta^{p+1}$

9 la suite est de la forme :

- ① Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- ② Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- ③ $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- ④ Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- ⑤ $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$
- ⑥ $\beta^p \leq a_{i_s} < \beta^{p+1}$
- ⑦ $\text{succ}(a_{i_s}) = a_{i_s} + \beta^{E(a_{i_s}-p+1)} = a_{i_s} + \beta^{p-p+1} = a_{i_s} + \beta$
- ⑧ $R(a_{i_s} + 1.0)$ vaut a_{i_s} ou $a_{i_s} + \beta$
- ⑨ $R(R(a_{i_s} + 1.0) - a_{i_s})$ vaut 0 ou β mais en aucun cas 1.0.
- ⑩ on sort donc de la boucle.

- ① Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- ② Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- ③ $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- ④ Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- ⑤ $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$
- ⑥ $\beta^p \leq a_{i_s} < \beta^{p+1}$
- ⑦ $\text{succ}(a_{i_s}) = a_{i_s} + \beta^{E(a_{i_s}-p+1)} = a_{i_s} + \beta^{p-p+1} = a_{i_s} + \beta$
- ⑧ $R(a_{i_s} + 1.0)$ vaut a_{i_s} ou $a_{i_s} + \beta$
- ⑨ $R(R(a_{i_s} + 1.0) - a_{i_s})$ vaut 0 ou β mais en aucun cas 1.0.
- ⑩ on sort donc de la boucle.

- 1 Par récurrence : $a_i = R(2^i) = 2^i$ tant que i vérifie $2^i \leq \beta^p - 1$.
- 2 Dans ce cas $R(a_i + 1.0) = a_i + 1.0$: pas d'arrondi.
- 3 $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$: la condition de la première boucle est donc vérifiée tant que $2^i < \beta^p$.
- 4 Première itération i_s telle que $2^{i_s} \geq \beta^p$.
- 5 $a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$
- 6 $\beta^p \leq a_{i_s} < \beta^{p+1}$
- 7 $\text{succ}(a_{i_s}) = a_{i_s} + \beta^{E(a_{i_s}-p+1)} = a_{i_s} + \beta^{p-p+1} = a_{i_s} + \beta$
- 8 $R(a_{i_s} + 1.0)$ vaut a_{i_s} ou $a_{i_s} + \beta$
- 9 $R(R(a_{i_s} + 1.0) - a_{i_s})$ vaut 0 ou β mais en aucun cas 1.0.
- 10 on sort donc de la boucle.


```
base = boucle2 (boucle1 1) 1
  where
    boucle1 a =
      if (a + 1.0) - a /= 1.0 then a
      else boucle1 (2.0 * a)
    boucle2 a b =
      if (a + b) - a == b then b
      else boucle2 a (b + 1.0)
```

Notons $a = \text{boucle1 } 1$

Son successeur est $a + \beta$.

Tant que $b < \beta$.

```
base = boucle2 (boucle1 1) 1
  where
    boucle1 a =
      if (a + 1.0) - a /= 1.0 then a
      else boucle1 (2.0 * a)
    boucle2 a b =
      if (a + b) - a == b then b
      else boucle2 a (b + 1.0)
```

Notons $a = \text{boucle1 } 1$

Son successeur est $a + \beta$.

Tant que $b < \beta$..

```
base = boucle2 (boucle1 1) 1
  where
    boucle1 a =
      if (a + 1.0) - a /= 1.0 then a
      else boucle1 (2.0 * a)
    boucle2 a b =
      if (a + b) - a == b then b
      else boucle2 a (b + 1.0)
```

Notons $a = \text{boucle1 } 1$

Son successeur est $a + \beta$.

Tant que $b < \beta$..

```
base = boucle2 (boucle1 1) 1
  where
    boucle1 a =
      if (a + 1.0) - a /= 1.0 then a
      else boucle1 (2.0 * a)
    boucle2 a b =
      if (a + b) - a == b then b
      else boucle2 a (b + 1.0)
```

Notons $a = \text{boucle1 } 1$

Son successeur est $a + \beta$.

Tant que $b < \beta$..

- 1 elle nous montre bien que nous changeons de monde : un test du style **if** $(a + 1.0) - a \neq 1.0$ paraîtrait bien hors de propos si nous raisonnions avec des réels ;
- 2 nous avons vu que l'on pouvait raisonner rigoureusement sur les VF ;
- 3 nous pouvons malgré tout retenir la trame intuitive de la démonstration : tant que notre nombre entier est représentable avec p chiffres, on reste exact. Les problèmes arrivent lorsqu'on n'a plus assez de place pour stocker tous les chiffres : il y a de la perte d'information...
- 4 Ce genre de raisonnement se retrouve très souvent pour travailler sur les erreurs commises : il faudra donc en retenir la substantifique moelle...

- 1 elle nous montre bien que nous changeons de monde : un test du style **if** $(a + 1.0) - a \neq 1.0$ paraîtrait bien hors de propos si nous raisonnions avec des réels ;
- 2 nous avons vu que l'on pouvait raisonner rigoureusement sur les VF ;
- 3 nous pouvons malgré tout retenir la trame intuitive de la démonstration : tant que notre nombre entier est représentable avec p chiffres, on reste exact. Les problèmes arrivent lorsqu'on n'a plus assez de place pour stocker tous les chiffres : il y a de la perte d'information...
- 4 Ce genre de raisonnement se retrouve très souvent pour travailler sur les erreurs commises : il faudra donc en retenir la substantifique moelle...

- 1 elle nous montre bien que nous changeons de monde : un test du style **if** $(a + 1.0) - a \neq 1.0$ paraîtrait bien hors de propos si nous raisonnions avec des réels ;
- 2 nous avons vu que l'on pouvait raisonner rigoureusement sur les VF ;
- 3 nous pouvons malgré tout retenir la trame intuitive de la démonstration : tant que notre nombre entier est représentable avec p chiffres, on reste exact. Les problèmes arrivent lorsqu'on n'a plus assez de place pour stocker tous les chiffres : il y a de la perte d'information...
- 4 Ce genre de raisonnement se retrouve très souvent pour travailler sur les erreurs commises : il faudra donc en retenir la substantifique moelle...

- 1 elle nous montre bien que nous changeons de monde : un test du style **if** $(a + 1.0) - a \neq 1.0$ paraîtrait bien hors de propos si nous raisonnions avec des réels ;
- 2 nous avons vu que l'on pouvait raisonner rigoureusement sur les VF ;
- 3 nous pouvons malgré tout retenir la trame intuitive de la démonstration : tant que notre nombre entier est représentable avec p chiffres, on reste exact. Les problèmes arrivent lorsqu'on n'a plus assez de place pour stocker tous les chiffres : il y a de la perte d'information...
- 4 Ce genre de raisonnement se retrouve très souvent pour travailler sur les erreurs commises : il faudra donc en retenir la substantifique moelle...

- 1 elle nous montre bien que nous changeons de monde : un test du style **if** $(a + 1.0) - a \neq 1.0$ paraîtrait bien hors de propos si nous raisonnions avec des réels ;
- 2 nous avons vu que l'on pouvait raisonner rigoureusement sur les VF ;
- 3 nous pouvons malgré tout retenir la trame intuitive de la démonstration : tant que notre nombre entier est représentable avec p chiffres, on reste exact. Les problèmes arrivent lorsqu'on n'a plus assez de place pour stocker tous les chiffres : il y a de la perte d'information...
- 4 Ce genre de raisonnement se retrouve très souvent pour travailler sur les erreurs commises : il faudra donc en retenir la substantifique moelle...

Nouvelles opérations

- $R(x + y)$ sous la forme $x \oplus y$;
- $R(x - y)$ sous la forme $x \ominus y$;
- $R(x \times y)$ sous la forme $x \otimes y$;
- $R(x/y)$ sous la forme $x \oslash y$.

Nouvelles opérations

- $R(x + y)$ sous la forme $x \oplus y$;
- $R(x - y)$ sous la forme $x \ominus y$;
- $R(x \times y)$ sous la forme $x \otimes y$;
- $R(x/y)$ sous la forme $x \oslash y$.

Nouvelles opérations

- $R(x + y)$ sous la forme $x \oplus y$;
- $R(x - y)$ sous la forme $x \ominus y$;
- $R(x \times y)$ sous la forme $x \otimes y$;
- $R(x/y)$ sous la forme $x \oslash y$.

Nouvelles opérations

- $R(x + y)$ sous la forme $x \oplus y$;
- $R(x - y)$ sous la forme $x \ominus y$;
- $R(x \times y)$ sous la forme $x \otimes y$;
- $R(x/y)$ sous la forme $x \oslash y$.

Trouver la précision

M.A. MALCOLM 1972 :

```
precision = toPrec 1.0 0
  where
    toPrec a i =
      if (a + 1.0) - a /= 1.0 then i
      else toPrec (base * a) (i + 1)
```

????

Trouver la précision

M.A. MALCOLM 1972 :

```
precision = toPrec 1.0 0
  where
    toPrec a i =
      if (a + 1.0) - a /= 1.0 then i
      else toPrec (base * a) (i + 1)
```

????



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 3 (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

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \ominus y)| \leq |x|$

FAITES UN DESSIN!

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \oplus y)| \leq |x|$

FAITES UN DESSIN!

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \oplus y)| \leq |x|$

FAITES UN DESSIN!

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \ominus y)| \leq |x|$

FAITES UN DESSIN !

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \oplus y)| \leq |x|$

FAITES UN DESSIN !

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \oplus y)| \leq |x|$

FAITES UN DESSIN !

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \oplus y)| \leq |x|$

FAITES UN DESSIN !

De l'importance de disposer avec la IEEE 754 de la meilleure approximation !

- $x \oplus y = x + y + \text{err}(x \oplus y)$
- $x \ominus y = x + y + (-y)$
- $x \oplus y$ le VF le plus proche de $x + y$
- $|\text{err}(x \oplus y)| \leq |y|$
- de même $|\text{err}(x \oplus y)| \leq |x|$

FAITES UN DESSIN !

À noter

De l'importance de disposer avec la IEEE 754 de la **meilleure approximation !**

Corollaire 4 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

Corollaire 4 (Møller, Knuth, Dekker)

L'erreur commise $|\text{err}(x \oplus y)|$ peut être exprimée exactement sur p bits.

Corollaire 4 (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 4 (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 4 (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 $\text{ulp}(y)$
- $|\text{err}(x \oplus y)| \leq |y|$

Corollaire 4 (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 4 (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 4 (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 5

*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 5

*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 5

*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 5

*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 $\forall x$ et y .
Est-ce que le résultat suivant contredit notre dernier lemme ?

```
*Main> 1 - 0.9  
9.999999999999998e-2
```

Lemme 6 (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 $|x - y| \leq \min(|x|, |y|)$ et on peut utiliser le lemme 5 page 346.

Lemme 6 (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 $0 < y - x < x = y$

• Si $x > y$ alors $\frac{x}{2} \leq y < x$ donc $0 < x - y < x - \frac{x}{2} = \frac{x}{2} < y$

Dans tous les cas $|x - y| \leq \min(|x|, |y|)$ et on peut utiliser le lemme 5 page 346.

Lemme 6 (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 5 page 346.

Lemme 6 (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 5 page 346.

- 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 .



Théorème 7 (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 non** : 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 7 (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 $|y| < \frac{|x|}{2}$) alors $|d| < |y|$ 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 7 (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 $|y| \leq \frac{|x|}{2}$: 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 7 (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

```
fastTwoSum x y
| abs x >= abs y =
  let s = x + y in
  let yv = s - x in
  let d = y - yv in
  (s, d)
| otherwise = fastTwoSum y x
```

Théorème 8 (2Sum - KNUTH)

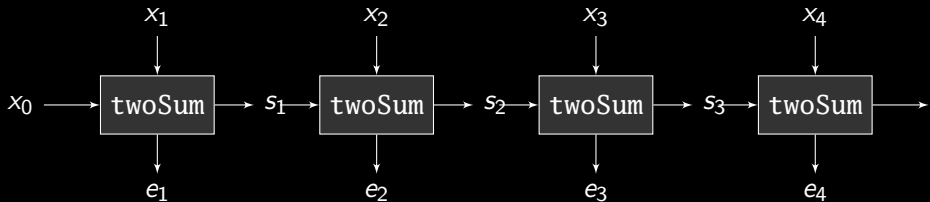
On considère deux VF x et y et l'algorithme suivant :

```
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  $(s, d)$ 
```

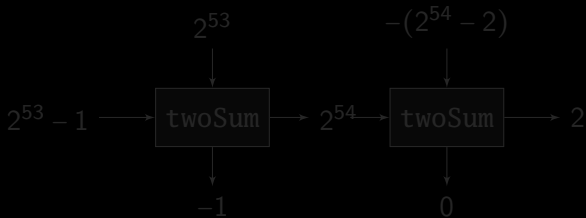
Alors $x + y = s + d$ avec $s = x \oplus y$ et $d = \text{err}(x \oplus y)$. De plus s et d ne se chevauchent pas.

Somme compensée d'un nombre quelconque de flottants

```
sommeKahan liste =  
  let (s,e) = foldl (\ (s,c) x -> twoSum s (x + c)) (0,0) liste in s
```

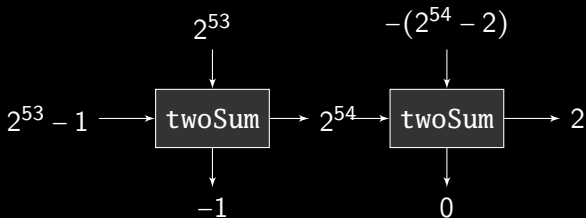


$$x_0 = 2^{53} - 1, x_1 = 2^{53} \text{ et } x_2 = -(2^{54} - 2)$$



Expliquez les résultats trouvés. Que va faire l'algorithme de somme compensée ? Et celui de somme en cascade ?

$$x_0 = 2^{53} - 1, x_1 = 2^{53} \text{ et } x_2 = -(2^{54} - 2)$$



Recherche 1

Expliquez les résultats trouvés. Que va faire l'algorithme de somme compensée ? Et celui de somme en cascade ?

Sommaire

- 1 Preamble
- 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
- 6 **MPFR**

```

#include <stdio.h>
#include <mpfr.h>

int main (void){
    mpfr_t un , dix , trois , zero_1, zero_3, trois_fois_zero_1, diff ; /*
        ↳ Déclarations simultanées */
    int test;

    mpfr_set_default_prec (53) ; /* Fixer la precision par default */
    mpfr_set_emin ( -1073) ; /* Fixer l' exposant emin ( en realite , emin -
        ↳ precision +2) */
    mpfr_set_emax (1024) ; /* Fixer l' exposant emax ( en realite , emax +1) */

    mpfr_inits (zero_1,zero_3, trois_fois_zero_1, (mpfr_ptr) 0); /* Précision par
        ↳ default 53 bits */
    mpfr_inits2 (128, un, trois, dix , diff, (mpfr_ptr) 0) ; /* Précision de 128
        ↳ bits exactement */

    mpfr_set_d (un , 1.0 , MPFR_RNDN ); /* MPFR_RND Division en arrondi au plus
        ↳ près */
    mpfr_set_d (trois , 3.0 , MPFR_RNDN );
    mpfr_set_d (dix , 10. , MPFR_RNDN );

    mpfr_div (zero_1 , un , dix , MPFR_RNDN );
    mpfr_div (zero_3 , trois , dix , MPFR_RNDN );
    mpfr_mul (trois_fois_zero_1, trois, zero_1, MPFR_RNDN );

```

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.