



## 2. Comparaison avec des flottants

```
c = 0
x = 1.0
y = x + 1.0
while y - x == 1.0:
    x = x * 2.0
    y = x + 1.0
    c = c + 1
```

1/ Si  $x, y \in \mathbb{R}$ , alors à chaque itération  $y - x = 1$  donc la boucle est infinie

2/ Test sur machine : le programme s'arrête

3/ Ajout d'un compteur  $c$

Arrêt pour  $c = 53 = 52 + 1$

4/ Arrêt  $\Leftrightarrow \exists$  une itération  $n$  tq :  $y - x \neq 1$  vrai  
Montre nous qu'il y a **absorption** i.e  $y = x + 1 = x$

À l'itération  $i$  :  
 $x_i = 2^i$   
 $y_i = 2^i + 1$   
donc  $\frac{1}{x_i} = 2^{-i}$

Si  $\frac{1}{x_i} < 2^{-52}$  alors  $x_{i+1} = x$

$\Leftrightarrow 2^{-i} < 2^{-52}$

$\Leftrightarrow i \geq 53 \Rightarrow$  **absorption à la 53<sup>e</sup> itérations!**

Ce que montre bien le compteur.

## 13. Résolution d'une équation du second degré

```
def zero_poly_deg_2(a,b,c):
```

```
    delta = b**2 - 4*a*c
```

```
    if delta > 0:
```

```
        x1 = 1/(2*a)*(-b - sqrt(delta))
```

```
        x2 = 1/(2*a)*(-b + sqrt(delta))
```

```
    elif delta < 0:
```

```
        x1 = 1/(2*a)*(-b - sqrt(-delta)*1j)
```

```
        x2 = 1/(2*a)*(-b + sqrt(-delta)*1j)
```

```
    else:
```

```
        x1 = x2 = -b/(2*a)
```

```
    return x1, x2
```

zéros poly-deg-2 (1.2, 1.5e8, 4e-4)

remise (-125000000, 0.0).

Or il est évident que 0 n'est pas solution.

Mathématiquement, les racines s'écrivent :

$$x_{\pm} = \frac{1}{2a} (-b \pm \sqrt{\Delta})$$

avec  $\Delta = b^2 - 4ac$

Calculons  $\Delta$  : 
$$\Delta = (1,5 \times 10^8)^2 - 4 \times 1,2 \times 10^{-4}$$

$$= \underbrace{2,25 \times 10^{16}}_{\alpha} - \underbrace{9,6 \times 10^{-4}}_{\beta}$$

or  $\frac{\beta}{\alpha} \sim 10^{-20} < 2^{-52} \Rightarrow \text{fl}(\alpha - \beta) = \alpha = \beta^2!$  **Absorption!**

Dans, ici,  $\Delta = \beta^2 \Rightarrow r_{\pm} = \frac{1}{2,4} (-b \pm \beta)$

$$\Leftrightarrow \begin{cases} r_+ = 0 \\ r_- = \frac{2\beta}{2,4} = \frac{1,5}{1,2} \times 10^8 = 1,25 \times 10^8. \end{cases}$$

L'erreur de calcul des racines résulte de l'absorption qui se produit lors du calcul de  $\Delta$ .

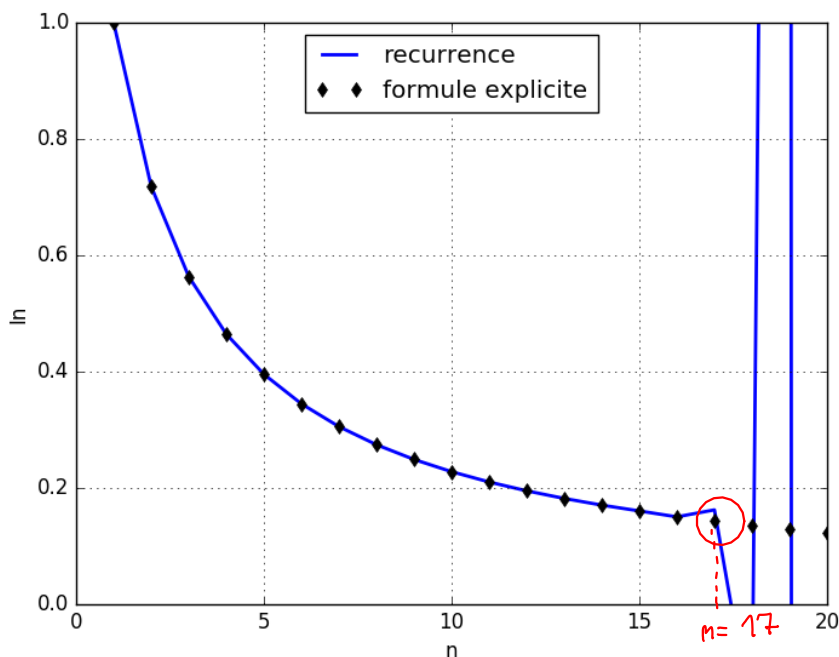
#### I4 - Divergence numérique

Soit la suite  $I_n$  définie :

- explicitement par : 
$$I_n = \int_0^1 x^n e^x dx$$

- par récurrence par : 
$$\begin{cases} I_0 = e - 1 \\ I_n = e - n I_{n-1}. \end{cases}$$

Expérimentalement, on observe les résultats suivants :



On observe une divergence entre le calcul explicite et le calcul itératif fondé sur la relation de récurrence définissant la suite.

Cette divergence est notable à partir de  $n = 17$  et s'amplifie au rangs suivants.

Pour comprendre, calculons les termes de la suite  $I_n$ .

$I_0 = e - 1$  mais  $e$  étant irrationnel, dans  $\mathbb{F}$ ,  $fl(e) \neq e$ .

On peut supposer que l'erreur sur  $e$  est de l'ordre de  $\epsilon = 2^{-32}$

Donc  $fl(e) = e + \epsilon$

D'où  $fl(I_0) = I_0 + \epsilon$  *erreur de représentation de  $I_0$ .*

Calculons  $I_1$  en tenant compte que l'erreur de représentation de  $I_0$ .

$$I_1 = e - I_0 \Rightarrow fl(I_1) = fl(e) - fl(I_0) = e + \epsilon - (I_0 + \epsilon) = e - I_0$$

$= e - (e - 1) = 1$   
*Jusqu'ici tout va bien!*

De même pour  $I_2; I_3 \dots I_n$ .

$$I_2 = e - 2I_1 \Rightarrow fl(I_2) = fl(e) - 2fl(I_1) = e + \epsilon - 2$$

$$I_3 = e - 3I_2 \Rightarrow fl(I_3) = fl(e) - 3fl(I_2) = e + \epsilon - 3(e + \epsilon - 2) = e + \epsilon - 3e - 3\epsilon + 6 = -2e + 6 - 2\epsilon$$

*les erreurs ne se compensent plus!*

$$I_4 = e - 4I_3 \Rightarrow fl(I_4) = fl(e) - 4fl(I_3) = e + \epsilon - 4(-2e + 6 - 2\epsilon) = 9e - 24 + 9\epsilon$$

Plus généralement :

$$\begin{aligned} I_n &= e - nI_{n-1} \\ &= e - n(e - (n-1)I_{n-2}) \\ &= e - n(e - (n-1)(e - (n-2)I_{n-3})) \end{aligned}$$

$$I_n = e - n(e - (n-1)(e - (n-2)(\dots - 2(e - 1 \times (e - I_0)))) \dots)$$

En développant l'expression, il apparaît le terme  $n! I_0$  qui contient l'ordre de grandeur de l'erreur commise à savoir :

$$fl(n! I_0) = n! fl(I_0) = n! I_0 + n! \epsilon \leftarrow \text{ordre de grandeur de l'erreur commise au rang } n : \delta_n$$

A.N. :  $m = 17, \delta_{17} \approx 0,08$   
 $m = 18, \delta_{18} = 1,42$   
 $m = 19, \delta_{19} = 27$  } cohérent avec le graphique!

Finalement, les oscillations apparentes de la suite résulte de l'erreur de représentation de  $e$  en machine. (toutes les suites ne sont pas aussi sensible à cette erreur).

15 - Erreur sur les multiplications (facultatif)

Nous raisonnons sur l'exemple de la multiplication (cf question 3)

$$1/ \eta_q : \left| \frac{fl(A \times B) - A \times B}{A \times B} \right| \leq \epsilon$$

Mantisse.

Soit :  $A = (1+m) \times 2^{e_1}$  au pire  $fl(A) = (1+m \pm \epsilon) \times 2^{e_1}$ , avec  $0 \leq m < 1$   
 $B = (1+m') \times 2^{e_2}$   $\Rightarrow$   $fl(B) = (1+m' \pm \epsilon) \times 2^{e_2}$ ,  $0 \leq m' < 1$

Alors :  $fl(A) \times fl(B) = (1+m \pm \epsilon)(1+m' \pm \epsilon) \times 2^{e_1+e_2}$   
 $= (1+m)(1+m') + [\pm(1+m) \pm (1+m')] \epsilon + \epsilon^2 \times 2^{e_1+e_2}$

$fl(A \times B) = ((1+m)(1+m') + [\pm(1+m) \pm (1+m')] \epsilon) \times 2^{e_1+e_2}$

Soit :  $\left| \frac{fl(A \times B) - A \times B}{A \times B} \right| = \frac{([\pm(1+m) \pm (1+m')] \epsilon + \epsilon^2) \times 2^{e_1+e_2}}{(1+m)(1+m') \times 2^{e_1+e_2}}$

$= \left| \frac{[\pm(1+m) \pm (1+m')] \epsilon}{(1+m)(1+m')} \right| \approx \epsilon$

$\sim 1$  car  $0 \leq m < 1$  et  $0 \leq m' < 1$

Finalement,  $\left| \frac{fl(A \times B) - A \times B}{A \times B} \right| \leq \epsilon$

"Semi-rigoureux"  
 - fait  
 l'erreur ne dépasse pas  $\max\left(\frac{2+m+m'}{(1+m)(1+m')}\right) \epsilon$   
 peut dépasser 1 pour  $m = m' = 0.5$  par exemple

2/ Sur 64 bits, 52 sont réservés au codage de la mantisse donc :

$\epsilon = 2^{-52} \approx 2,22 \times 10^{-16}$

3/ Erreur maximale sur  $10^6$  multiplications successives.

On a montré que :

$fl(A) \times fl(B) = A \times B + [\pm(1+m) \pm (1+m')] \epsilon \times 2^{e_1+e_2} + \epsilon^2 \times 2^{e_1+e_2}$

De la même façon, avec  $C = (1+m'') \times 2^{e_3}$  et  $fl(C) = (1+m'' \pm \epsilon) \times 2^{e_3}$

$fl(A) \times fl(B) \times fl(C) = A \times B \times C + [\pm(1+m) \pm (1+m') \pm (1+m'')] \epsilon + (2\epsilon^2 + \beta \epsilon^3) \times 2^{e_1+e_2+e_3}$

$\Rightarrow fl(A \times B \times C) \approx A \times B \times C + [\pm(1+m) \pm (1+m') \pm (1+m'')] \epsilon \times 2^{e_1+e_2+e_3}$

au pire  $\approx 6 = 3 \times 2$   $O(\epsilon)$  : négligeable.  
 au pire  $\approx 6 = 3 \times 2$ .

Soit  $\eta$  le produit de  $10^6$  facteurs, les erreurs d'ordre  $\sim \epsilon^m$  sont évanouies, subsistent l'erreur relative  $\sim \epsilon$ .

$fl(\eta) < \eta + 10^6 \times 2 \times \epsilon \times 2^m \Rightarrow$  erreur relative  $\sim 2 \times 10^6 \times \epsilon \approx 4 \times 10^{-10}$

Les erreurs s'accumulent.