

II. Erreurs d'annulation

Lorsque l'on effectue certaines opérations, par exemple la soustraction de deux nombres petits et proches, on voit apparaître ce qui semble une aberration. Ainsi :

```
-->1.23456789E-11-1.23456788E-11
ans =
1.000000000577265217E-19
```

On devrait dans l'idéal obtenir $1,00000000000000000000.10^{-19}$

On peut expliquer ce phénomène par le fait que l'on ne calcule pas en base 10. En effet, le stockage en binaire constitue une approximation des nombres entrés sous forme décimale :

$$\tilde{a}=a(1+\Delta a); \tilde{b}=b(1+\Delta b); \tilde{x}=\tilde{a}-\tilde{b};$$

Par conséquent : $\left| \frac{\tilde{x}-x}{x} \right| = \left| \frac{a\Delta a - b\Delta b}{a-b} \right|$ et ce nombre peut être très grand si a et b sont très proches.

Si l'on veut estimer l'évolution de l'erreur d'annulation lors de l'application d'une fonction, on peut faire le calcul suivant :

$$y=f(x); \left| \frac{\Delta y}{y} \right| = \left| \frac{\tilde{y}-y}{y} \right| = \left| \frac{f(\tilde{x})-f(x)}{f(x)} \right| = \left| \frac{f(x(1+\Delta x))-f(x)}{f(x)} \right| \approx \left| \frac{f(x)+xf'(x)-f(x)}{f(x)} \right| = \left| \frac{xf'(x)}{f(x)} \right|$$

On obtient ainsi un critère de stabilité : il faut que : $\left| \frac{xf'(x)}{f(x)} \right| \ll 1$ pour que l'erreur n'influence pas trop le résultat.

Somme d'une série

On peut encore constater des erreurs importantes lors de calculs impliquant des nombres d'ordres de grandeur très différent. Par exemple, la somme de la série suivante :

$$S_n = 1 + \sum_{k=1}^n \frac{1}{k(k+1)} = 2 - \frac{1}{n+1} = S_{th}$$

On utilise le programme suivant :

```
n=1000;
Sn=0;
Tn=1;
for k=1:n
    Sn=Sn+1/(k*(k+1));
    Tn=Tn+1/(k*(k+1));
    k=k+1;
end
Sth=2-1/(n+1)
Sn=Sn+1
Tn
erreur1=abs((Sn-Sth)/Sth)
erreur2=abs((Tn-Sth)/Sth)
```

Qui donne le résultat :

```
Sn =  
  1.999001  
Tn =  
  1.999001  
erreur1 =  
  3.332E-16  
erreur2 =  
  1.555E-15
```

Dans ce cas, pour $n=1000$, la somme effectuée avec le membre de gauche fait apparaître une erreur de $1,555 \cdot 10^{-15}$ alors que l'ajout de 1 à la somme des termes seuls ne produit une erreur que de $3,332 \cdot 10^{-16}$, soit 10 fois plus petite.

On explique cela par le fait que le second calcul implique des termes plus petits et plus proches les uns des autres avant d'ajouter le 1 final. Le premier en revanche ajoute des termes très petits à un total déjà de l'ordre de grandeur de 1.

III. Propagation d'erreurs

Enfin, la propagation d'erreurs (ou *smearing*) est la perte de chiffres significatifs due à des erreurs d'arrondis qui se cumulent.

On peut s'en apercevoir dans le programme suivant, qui calcule $\exp(x)$ de deux manières différentes, d'une part avec la fonction exponentielle du logiciel, et d'autre part en utilisant la formule de Taylor-intégrale :

$$\exp(x) = \sum_{k=0}^n \frac{x^k}{k!} + \int_0^x \exp(t) \frac{t^n}{n!} dt$$

```
n=1 ; x=-20 ; ex=1 ; tol=1e-3 ; un=x ;  
while abs(un)/exp(x) >= tol & n<100000  
  ex=ex+un;  
  n=n+1;  
  un=un*x/n;  
end  
ex  
expx=exp(x)  
erreur=abs((ex-exp(x))/exp(x))  
n
```

L'exécution donne :

```
ex =
  5.623089322732710410E-09
expx =
  2.061153622438557870E-09
erreur =
  1.728127230070320763E+00
n =
  7.500000000000000000E+01
```

On peut remarquer qu'au bout de 75 itérations, l'ajout d'autres termes u_n ne contribue presque plus au calcul de $\exp(x)$ (les nouveaux termes étant inférieurs à un millième de la somme). On a donc fait un calcul de $\exp(x)$ censé être précis au millième et pourtant, l'erreur relative calculée à ce point est énorme : de l'ordre de grandeur de la valeur à calculer, puisque le rapport de l'écart sur la valeur théorique vaut environ 1,7 !

Plusieurs solutions sont envisageables pour effectuer ce même calcul sans une erreur trop grande. On peut ainsi :

- regrouper les termes positifs et négatifs (les $(-20)^i$ avec i pair ou impair)
- calculer $\exp(20)$ puis utiliser le fait que $\exp(-20) = \frac{1}{\exp(20)}$

De telles optimisations minimiseront le nombre d'opérations effectuées sur des nombres très différents et donc le nombre d'erreurs d'arrondis.

Conclusion

Nous avons vu que le simple fait de déplacer la somme d'un terme avant ou après un calcul pouvait faire varier l'erreur totale d'un facteur 10. Dans d'autres cas, l'erreur relative atteint une telle taille que le résultat du calcul est simplement faux dès son premier chiffre.

On peut donc affirmer que si l'on souhaite obtenir les calculs les plus précis possibles avec l'outil informatique, il convient de regarder d'assez près les calculs que l'on effectue et de comprendre comment l'ordinateur va les effectuer.

On pourra ainsi repérer les expressions dans lesquelles des erreurs d'arrondis ou de trop grandes différences d'ordre de grandeur pourront apparaître et ainsi optimiser l'ordre des calculs afin d'augmenter la précision de ceux-ci.