



République Algérienne Démocratique et Populaire
Ministère de l'Enseignement Supérieure et de la Recherche Scientifique
Université Ibn Khaldoun – Tiaret
Faculté des Mathématiques et de l'Informatique
Département d'Informatique



METHODES NUMERIQUES

COURS POUR :

2ème ANNEE LICENCE INFORMATIQUE

Par : Mr. KARIM MEZZOUG

karim.mezzoug@univ-tiaret.dz

Table des Matières

Chapitre 1 : Généralités sur l'analyse numérique et le calcul scientifique

- 1.1 Motivations.
- 1.2 Arithmétique en virgule flottante et erreurs d'arrondis
 - 1.2.1 Représentation des nombres en machine
 - 1.2.2 Erreurs d'arrondis
- 1.3 Stabilité et analyse d'erreur des méthodes numériques et conditionnement d'un problème

Chapitre 2 : Méthodes directes de résolution des systèmes linéaires

- 2.1 Remarques sur la résolution des systèmes triangulaires
- 2.2 Méthodes d'élimination de Gauss et Méthodes d'élimination de JORDAN
- 2.3 Interprétation matricielle de l'élimination de Gauss : la factorisation LU

Chapitre 3 : Méthodes itératives de résolution des systèmes linéaires

- 3.1 Généralités
- 3.2 Méthodes de Jacobi .
- 3.3 Méthodes de Gauss-Seidel et de sur-relaxation successive
- 3.4 Remarques sur l'implémentation des méthodes itératives
- 3.5 Convergence des méthodes de Jacobi et Gauss-Seidel

Chapitre 4 : Calcul de valeurs et de vecteurs propres

- 4.1 Localisation des valeurs propres
- 4.2 Méthode de la puissance et de Déflation.

Remarque : Le Chapitre 5 : Analyse matricielle a été abordée d'une façon intégrale et progressive à travers les chapitres 2,3 et 4.

Chapitre I : Généralités sur l'Analyse Numérique et le Calcul Scientifique

1. Motivations :

L'analyse numérique est devenue indispensable dans de très nombreux domaines. Les différents disciplines qui ont recours au calcul numérique sont trop variées pour être décrites ici. Elle va de l'astronomie à la sociologie en passant par la physique, la chimie, la biologie, la médecine, le génie civil, l'archéologie, la gestion des stocks, la psychologie expérimentale, les jeux vidéo, etc...

À travers ce chapitre, on étudiera essentiellement en premier lieu la notion d'erreur ensuite l'arithmétique en virgule flottante en passant par la propagation de l'erreur puis nous terminons par la notion de conditionnement d'un problème et la stabilité des algorithmes numériques. Tout ceci sera accompagné par quelques exemples démonstratifs.

2. Notions d'Erreurs :

La notion d'erreur n'est pas considérée au sens de faute dans la méthodologie, ou instruction fautive dans le programme. En fait elle concerne toutes erreurs inévitables. Ainsi les erreurs se distinguent en quatre catégories :

- **Les erreurs humaines** : Toutes erreurs commises à cause de l'inattention ou l'oubli, etc.. des êtres humains pendant des mesures effectuées à la main et à l'œil nu.
- **Les erreurs sur les données** : Généralement dues à l'imprécision des mesures effectuées par des capteurs ou des instruments de mesures électroniques, etc...
- **Les erreurs d'arrondi**. Un ordinateur numérique ne peut représenter les nombres réels qu'avec un nombre fini de chiffres, d'où à chaque opération mathématique élémentaire, il y aura une perte de chiffres significatifs, surtout quand il s'agit de nombre d'opérations très importantes.
- **Les erreurs d'approximation** : ou de discrétisation. Dues essentiellement à la méthodologie et l'algorithme utilisés eux-mêmes par exemple, lorsqu'on calcule une intégrale à l'aide d'une somme finie, une dérivée à l'aide de différences finies ou bien la somme d'une série finie à l'aide d'un nombre fini de ses termes. On analyse numériquement généralement on cherche à évaluer ces erreurs de discrétisation pour chaque algorithme proposé.

3. Représentation des nombres en Machine et Arithmétique en virgule flottante :

Les opérations arithmétiques que nous effectuons tous les jours sont généralement exactes, par exemple payer sa facture au supermarché ou au guichet d'une salle de spectacle. Dans de tels cas nous manipulons des nombres entiers et des réels ayant peu de chiffres décimaux en effectuant que de simples additions ou soustractions.

Mais réellement ce n'est pas toujours le cas et se rend compte que nos calculs sont parfois trompeuses : par exemple le (19/09/2019), j'ai fait le plein d'essence routier à 41,97 DA le litre, ce qui m'a coûté 1768 DA, alors que mon réservoir a accepté 42.12 litres de carburant. J'aurais dû, en réalité, déboursier 1767.7764 DA.

Mais, ai-je bien reçu la quantité affichée au compteur de la pompe ? C'est peu probable. La mesure du volume fourni est évidemment entachée d'une certaine erreur (qui varie d'ailleurs en fonction de la température ambiante, de la pression atmosphérique et du débit de fluide) et j'ai peut-être payé mon carburant à 42 DA ou 40 DA le litre !

Cet exemple concret nous permet de mettre en lumière deux sources d'erreurs différentes. L'une est directement liée à la manière d'exprimer les nombres : on utilise deux chiffres décimaux pour estimer le volume de carburant fourni, deux pour le prix au litre et aucun pour le prix total à payer. L'autre source d'erreur est directement liée aux données qui peuvent elles-mêmes n'être qu'approchées comme c'est le cas ici pour le volume de diesel effectivement délivré par la pompe.

3.1. Représentation des nombres en machine :

A- La virgule flottante :

Un nombre X réel peut être écrit sous la forme standardisé comme suit :

$$x \approx \pm m b^P$$

$$\approx \pm .a_1 a_2 \dots a_{ms} b^P$$

Il est constitué d'un signe « + » ou « - », du point décimal « . », d'une mantisse « $a_1 a_2 \dots a_{ms}$ », d'une base « b » et d'un exposant « M », sachant que :

- La base « b » doit être un entier ≥ 2
- Les « a_i » doivent être des entiers tels que $0 \leq a_i < b$
- L'exposant « P » doit respecter la condition : $P_{\min} \leq P \leq P_{\max}$

Un nombre réel non nul est dit normalisé si $a_1 \neq 0$.

Exemple :

$x = -0.0231 \times 10^2$ n'est pas normalisé, mais, $x = -0.2310 \times 10^1$ l'est.

Le nombre de réels représentables est limité et vaut, pour m, b, P_{\max} et P_{\min} fixés :

$$2 \cdot (b - 1) \cdot b^{(ms-1)} \cdot (P_{\max} - P_{\min} + 1) + 1 \quad (\text{ici le point . exprime la multiplication})$$

Le nombre de réels normalisés représentables est :

Exemple :

pour, $b=10$, $ms=2$, $P_{\min}=0$, $P_{\max}=1$:

$$2(10-1)10^{(2-1)}(1-(0)+1)+1 = 180 \times 2 + 1 = 361$$

$$2(b - 1)b^{(ms-1)}(P_{\max} - P_{\min} + 1) + 2$$

car deux représentation de zéro sont prises en compte $+0.0000\dots$ et $-0.0000\dots$

B- Règle d'arrondissement :

Pour arrondir un nombre jusqu'à n chiffres significatifs (c.s), il faut éliminer les chiffres à droite du N éme chiffre significatif conservé si on se trouve après la virgule, sinon on remplace par des zéros :

1. Si le (N + 1) éme chiffre significatif est > 5, on augmente le n éme chiffre de 1.
2. Si le (N + 1) éme chiffre significatif est < 5, les chiffres retenus restent inchangés.
3. Si le (N + 1) éme chiffre significatif est 5, alors deux cas sont possibles :
 - ✓ Tous les chiffres rejetés, situés après le (N + 1) éme c.s, sont des zéros. les chiffres retenus restent inchangés
 - ✓ Parmi les chiffres rejetés, situés après le (N + 1) éme c.s, il existe au moins un qui soit non nul : On ajoute 1 au N éme chiffre.

Remarque Importante :

Supposons que le nombre est représenté sur machine comme suit avec N chiffres significatifs :

$$x \approx \pm m b^p$$

D'où :

$$b^{-1} \leq m < 1.$$

Alors la différence entre la valeur exacte d'un nombre x et la valeur x* de sa représentation sur machine est appelée erreur d'arrondi. Elle est majorée par :

$$\frac{1}{2} b^{-N} b^p,$$

Soit en valeur relative :

$$\left| \frac{x - x^*}{x} \right| \leq \frac{b^{-N} b^p}{2x} \leq \frac{b^{-N} b^p}{2b^{p-1}} = \frac{b^{1-N}}{2}.$$

Exemple : voir section 1.4

3.2. Arithmétique en virgule flottante :

- **Cas 1 :**

b = 10, ms = 3, P_{min} = -15, P_{max} = 16

x = +.125 × 10⁶, y = -.128 × 10⁶, z = +.437 × 10¹², u = +.215 × 10⁻¹⁰

x + y = (+.125 × 10⁶) + (-.128 × 10⁶) = (+.125 - .128) × 10⁶ = -.003 × 10⁶

= -.300 × 10⁴ (normalisation)

x + z = (+.125 × 10⁶) + (+.437 × 10¹²) = (+.000000125 + .437) × 10¹² (alignement)

= +.437000125 × 10¹² = +.437 × 10¹² (standardisation)

Remarque 1 : x + z = z, avec x ≠ 0 !!!

x . y = (+.125 × 10⁶) × (-.128 × 10⁶)

= (+.125) × (-.128) × 10⁶⁺⁶

$$\begin{aligned}
 &= -.016000 \times 10^{12} = -.160 \times 10^{11} \text{ (normalisation)} \\
 x \cdot z &= (+.125 \times 10^6) \times (+.437 \times 10^{12}) \\
 &= (+.125) \times (.437) \times 10^{6+12} \\
 &= +.054625 \times 10^{18} = +.546 \times 10^{17} \text{ (normalisation)}
 \end{aligned}$$

Remarque 2 : l'exposant est supérieur à M_{\max} et le produit $x \cdot z$ est donc supérieur au plus grand nombre représentable ($+.999 \times 10^{16}$) on dit qu'il y a surdépassement. Ceci se traduira par un message d'erreur affiché par le calculateur indiquant qu'il y a **overflow**.

$$\begin{aligned}
 x \div z &= (+.125 \times 10^6) \div (+.437 \times 10^{12}) \\
 &= (+.125) \div (+.437) \times 10^{6-12} \\
 &= +.286041189... \times 10^{-6} = +.286 \times 10^{-6} \text{ (normalisation)} \\
 u \div z &= (+.215 \times 10^{-10}) \div (+.437 \times 10^{12}) \\
 &= (+.215) \div (+.437) \times 10^{-10-12} \\
 &= +.491990846... \times 10^{-22} = +.492 \times 10^{-22} \text{ (normalisation et arrondissement)}
 \end{aligned}$$

remarque 3 : l'exposant est inférieur à M_{\min} et le quotient $u \div z$ est donc plus petit que le plus petit nombre positif représentable ($+100 \times 10^{-15}$) ceci se traduira par un message d'erreur affiché par le calculateur indiquant qu'il y a **underflow**. Dans pareil cas, selon les logiciels et matériels utilisés, le résultat est remplacé par zéro ou par $\varepsilon(0)$. « L' $\varepsilon(0)$ machine » dépend du calculateur et du logiciel employés et correspond au plus petit nombre distinguable de zéro (dans la version 7.0.1 de Matlab, implémentée sur PC portable, l' $\varepsilon(0)$ machine vaut habituellement $4.940656458412465 \times 10^{-324}$).

- **Cas 2 :**

Soient $b = 10$, $m = 3$, $x = y = +.400 \times 10^0$ et $z = +.100 \times 10^3$

On a :

$$y + z = +.100 \times 10^3, \Rightarrow x + (y + z) = +.100 \times 10^3 = z$$

$$x + y = +.800 \times 10^0, \Rightarrow (x + y) + z = .101 \times 10^3 \neq x + (y + z).$$

Remarque 4 : l'ordre dans lequel sont effectuées les opérations peut alterner le résultat final. Par exemple quand on additionne des nombres positifs, il faut commencer par les plus petits le plus possible.

- **Cas 3 :**

La moyenne μ et la variance σ de n nombres s'obtiennent à l'aide des formules suivantes :

$$\begin{aligned}
 \mu &= \frac{1}{n} \sum_{i=1}^n x_i \\
 \sigma^2 &= \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \mu^2
 \end{aligned}$$

Si l'on choisit $b = 10$, $m = 3$ et $n = 2$, les résultats obtenus pour la variance des deux nombres $x_1 = +.310 \times 10^2$ et $x_2 = +.320 \times 10^2$, sont $\sigma^2 = +.250 \times 10^0$ ou $\sigma^2 = -.200 \times 10^1$ selon la formule utilisée.

Remarque 5 : deux formules théoriquement équivalentes peuvent donner des résultats numériques très différents ! L'erreur est due à la soustraction de deux nombres à peu près égaux. On parle de **cancellations** ou d'annulation numérique. Elle se traduit par la perte de chiffres significatifs et peut conduire, comme dans cet exemple, à des résultats aberrants : $\sigma^2 < 0$!!

- **Cas 4 :**

$$A = +.4025 \times 10^1, B = +.4019 \times 10^1 \Rightarrow A - B = C = +.6000 \times 10^{-2}$$

les zéros sont-ils significatifs ? Le six est-il significatif ?

Si A et B sont des valeurs approchées et qu'on a en réalité :

$$A = +.4024786 \times 10^1 \text{ et } B = +.4019432 \times 10^1 \text{ on obtient alors } C = +.5354 \times 10^{-2}$$

Remarque 6 : la soustraction de deux nombres proches peut se traduire par une perte de chiffres significatifs, c'est la cancellation ou l'annulation.

- **Cas 5 :**

Les racines du trinôme du second degré $ax^2 + bx + c$ sont données par les formules équivalentes :

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{2c}{-b \mp \sqrt{b^2 - 4ac}}$$

$$\text{Si } b^2 \gg ac \Rightarrow x = \frac{-b \pm |b|}{2a} \text{ et } x = \frac{2c}{-b \mp |b|}$$

Remarque 7 : on obtient alors une racine «nulle» et l'autre «infinie» ce qui est évidemment aberrant et résulte à nouveau d'une cancellation.

- **Cas 6 :**

$$\sum_{n=1}^{\infty} \frac{1}{n} \text{ est une série divergente mais } \lim_{n \rightarrow \infty} \frac{1}{n} \rightarrow 0.$$

En base 10, le calcul avec un seul chiffre, ($m_s = 1$) donne une somme qui vaut 2. Avec deux chiffres ($m_s = 2$) on obtient 3.9

Remarque 8 : une série divergente peut sembler converger lorsqu'on calcule numériquement la somme de ses termes. Dans l'exemple précédent, en travaillant avec un seul chiffre la série semble converger vers 2. Toutefois, si on augmente le nombre de chiffres des mantisses utilisées, on obtient des sommes différentes. Tout se passe donc comme si la série admettait plusieurs limites distinctes... Contrairement aux apparences il n'y a donc pas convergence

3.3. Erreur Absolue et Erreur Relative :

Soit x la valeur exacte et soit x^* la valeur approchée de x alors l'erreur absolue $\Delta(x)$ se calcule comme suit : $\Delta(x) = |x - x^*|$

Et l'erreur relative $\delta(x)$ se calcule par : $\delta(x) = \frac{\Delta(x)}{|x|}$

Remarque : Généralement l'erreur absolue s'avère plus efficace que l'erreur absolue. Par exemple l'erreur absolue entre $x=1$ et $x^*=1.1$ et la même qu'entre $y=1000$ et $y^*=1000.1$ alors que l'erreur relative 0.1/1 et différent que celle 0.1/1000.

Exemple1 : Soit $x=\pi$ et (où $\pi = 3.141592653589793 \dots$) et $\pi^* = 0.31415927 \cdot 10^{+1}$. Supposons que dans un calcul apparaisse la quantité :

$$\begin{aligned} \Delta\pi &= |0.3141592653589793 - 0.31415927| \cdot 10^{+1} = -0.46410207 \cdot 10^{-7} \\ \left| \frac{0.3141592653589793 \cdot 10^{+1} - 0.31415927 \cdot 10^{+1}}{0.3141592653589793 \cdot 10^{+1}} \right| &= \left| \frac{-0.46410207 \cdot 10^{-8}}{0.3141592653589793} \right| \\ &= 0.14772828 \cdot 10^{-7} \leq \frac{10^{-7}}{2} \end{aligned}$$

Exemple 2 : Pour la valeur exacte $x = 2/3$, la valeur approchée $x^*_1 = 0.666667$ est 1000 fois plus précise que la valeur approchée $x^*_2 = 0.667$. En effet, nous avons :

$$\begin{aligned} \Delta_1(x) &= |x - x^*_1| = |2/3 - 0.666667| = |0.66666667 - 0.66666700| = 0.00000033 = 0.33 \cdot 10^{-6} \\ \Delta_2(x) &= |x - x^*_2| = |2/3 - 0.667| = |0.66666667 - 0.66700000| = 0.00033333 = 0.33333 \cdot 10^{-3} \end{aligned}$$

Les erreurs relatives correspondantes sont :

$$\begin{aligned} \delta_1(x) &= \frac{\Delta_1(x)}{|x|} = \frac{0.33 \cdot 10^{-6}}{|0.66666667|} = 0.5 \cdot 10^{-6} = 0.5 \cdot 10^{-4} \% \\ \delta_2(x) &= \frac{\Delta_2(x)}{|x|} = \frac{0.33333 \cdot 10^{-3}}{|0.66666667|} = 0.5 \cdot 10^{-3} = 0.5 \cdot 10^{-1} \% \end{aligned}$$

3.3.1. Majoration d'Erreur :

Soit ε l'epsilon absolue machine et ε_r l'epsilon relative machine alors l'erreur absolue est majorée par :

$$\Delta(x) \leq \varepsilon \quad \text{et} \quad x^* - \varepsilon \leq x \leq x^* + \varepsilon \quad \text{et} \quad x = x^* \pm \varepsilon$$

Et on l'erreur relative est majorée par : $\delta(x) \leq \varepsilon_r$ et $x^* = x(1 \pm \varepsilon_r)$

3.3.2. Propagation de l'Erreur :

Soit f une fonction de n variable réels $(x_1, x_2, x_3, \dots, x_n)$.

Evaluons maintenant l'erreur de $f(x_1, x_2, x_3, \dots, x_n)$:

Etant les valeurs exactes des x_i sont données par la formule approchée :

$$x^* = x + dx_i \quad \text{sachant que } dx_i = x^* - x$$

$$\text{donc : } f(x^*) \approx f(x) + \sum_{i=1}^n dx_i \frac{\partial f(x)}{\partial x_i}$$

$$\text{et : } \Delta f(x) = |df(x)| = |f(x^*) - f(x)| = \left| \sum_{i=1}^n \frac{\partial f(x)}{\partial x_i} dx_i \right| \leq \sum_{i=1}^n \left| \frac{\partial f(x)}{\partial x_i} \right| |dx_i|$$

$$\text{sachant que } |dx_i| = \Delta(x_i) \quad \text{et} \quad df(x) = f(x) - f(x^*)$$

Ainsi pour l'erreur relative on aura :

$$\left| \frac{\Delta f(x)}{f(x)} \right| = \left| \sum_{i=1}^n \frac{x_i}{f(x)} \frac{\partial f(x)}{\partial x_i} \frac{dx_i}{x_i} \right| \leq \sum_{i=1}^n \left| \frac{x_i}{f(x)} \frac{\partial f(x)}{\partial x_i} \right| \left| \frac{dx_i}{x_i} \right|$$

En pratique les valeurs des dérivées de la fonction sont évaluées selon la valeur x^* .

- **Exemple d'application sur l'addition :**

Soit la fonction : $f(x,y) = x \pm y$.

$$\text{Alors : } d(x \pm y) \leq d(x) + d(y) \quad \text{et} \quad \Delta(x \pm y) \leq \Delta(x) + \Delta(y)$$

$$\left| \frac{d(x \pm y)}{(x \pm y)} \right| \leq \left| \frac{x}{x \pm y} \right| \left| \frac{dx}{x} \right| + \left| \frac{y}{x \pm y} \right| \left| \frac{dy}{y} \right|$$

Exemple de calculs : Soient $x^* = 255$ et $y^* = 250$ avec $\delta x = \delta y = 0,1\% = 10^{-3}$.

Question : $\delta(x-y) = ?$

Nous avons d'abord : $\Delta x = x^* \cdot \delta x = 0,255$, $\Delta y = y^* \cdot \delta y = 0,250$.

et puis $x^* - y^* = 5$ avec une erreur relative :

$$\delta(x-y) = \Delta(x-y) / |x^* - y^*| \leq \Delta x + \Delta y / |x^* - y^*| = 0,101 \cdot 10^0 = 10,1\%.$$

On remarque que x^* et y^* sont 101 fois plus précis pour x et y (respectivement) que $(x^* - y^*)$ l'est pour $(x - y)$.

- **Exemple d'application 2 :** reprendre l'exemple précédent en appliquant le calcul suivant :

$$\delta(x,y) = ?$$

4. Stabilité et analyse d'erreur des méthodes numériques et conditionnement d'un problème

4.1. Conditionnement

Le conditionnement décrit la sensibilité de la valeur d'une fonction à une petite variation de son argument, c'est-à-dire :

$$\frac{f(x) - f(x^*)}{f(x)} \quad \text{en fonction de} \quad \frac{x - x^*}{x}$$

Lorsque $(x - x^*)$ est petit. Pour une fonction suffisamment régulière, on a évidemment :

$$\left| \frac{f(x) - f(x^*)}{f(x)} \cdot \frac{x - x^*}{x} \right| \simeq \left| \frac{xf'(x)}{f(x)} \right|.$$

D'où on tire :

Définition : On appelle conditionnement d'une fonction numérique f de classe C^1 en un point x , le nombre :

$$\text{cond}(f)_x := \left| \frac{xf'(x)}{f(x)} \right|.$$

Exemple :

$$f(x) = \sqrt{x}$$

$$\left| \frac{xf'(x)}{f(x)} \right| = \frac{1}{2}.$$

Ceci correspond à un bon conditionnement, puisque l'erreur relative sur f sera au plus moitié d'une erreur relative sur x .

Exemple : $f(x) = a - x$

$$\left| \frac{xf'(x)}{f(x)} \right| = \left| \frac{x}{a - x} \right|.$$

Ici, le conditionnement est très mauvais si x est voisin de a .

4.2. Stabilité

Définition : La stabilité décrit la sensibilité d'un algorithme numérique pour le calcul d'une fonction $f(x)$.

Exemple :

$$f(x) = \sqrt{x+1} - \sqrt{x}.$$

Le conditionnement de cette fonction est égal à :

$$\left| \frac{xf'(x)}{f(x)} \right| = \left| \frac{x(\sqrt{x} - \sqrt{x+1})}{2\sqrt{x(x+1)}} - \frac{1}{\sqrt{x+1} - \sqrt{x}} \right| = \frac{1}{2} \sqrt{\frac{x}{x+1}}$$

Cette dernière expression étant proche de 1/2 pour x grand. Donc, si x est grand, le conditionnement de f est bon. Cependant, dans un calcul à 6 chiffres significatifs, on a :

$$f(12345) = \sqrt{12346} - \sqrt{12345} = 111,113 - 111,108 = 0,500000 \cdot 10^{-2}.$$

Or un calcul précis donne : $f(12345) = 0,4500032 \dots \cdot 10^{-2}$. On a donc une erreur relative de 10% ce qui est important et peu en accord avec le bon conditionnement de f. Ceci est dû à l'algorithme utilisé dans ce calcul que l'on peut expliciter comme suit :

$$\begin{aligned} x_0 & : = 12345 \\ x_1 & : = x_0 + 1 \\ x_2 & : = \sqrt{x_1} \\ x_3 & : = \sqrt{x_0} \\ x_4 & : = x_2 - x_3 \end{aligned}$$

Il y a quatre fonctions à intervenir et, à priori, même si le conditionnement de f est bon, il se peut que le conditionnement d'une ou plusieurs fonctions utilisées dans l'algorithme soit supérieur à celui de f. C'est ce qui se produit ici pour la fonction $x_3 \rightarrow x_4 = x_2 - x_3$ (x_2 étant supposé fixe) dont le conditionnement est grand lorsque x_3 est voisin de x_2 comme on l'a vu précédemment.

En conclusion, le choix d'un bon algorithme numérique est essentiel. Par exemple, ci dessus, un meilleur algorithme est obtenu en utilisant :

$$f(x) = \sqrt{x+1} - \sqrt{x} = \frac{1}{\sqrt{x+1} + \sqrt{x}}.$$

Dans ce cas, toujours avec 6 chiffres significatifs :

$$f(12345) = \frac{1}{\sqrt{12346} + \sqrt{12345}} = \frac{1}{222,221} = 0,450002 \cdot 10^{-2}$$

ce qui donne une erreur relative de 0.0003 %.

Exemple d'instabilité : Calcul de e^{-12} à l'aide de sa série de Taylor

$$e^x = \sum_{n=0}^N \frac{x^n}{n!} (= S_N) \text{ pour } N \text{ grand.}$$

Exercice : Calculer à l'aide de Scilab les valeurs successives de cette somme en fonction de N pour $x = -12$. À dix chiffres significatifs. Puis à refaire les calculs avec la formule suivante :

$$e^{-x} = \frac{1}{e^x}$$

Laquelle des deux formule est plus stables ?

1. Conclusion :

Dans la vie réelle plusieurs catastrophes ont été généralement le résultat inévitable d'une mauvaise gestion de l'arithmétique des ordinateurs ou calculateurs (erreurs d'arrondis, d'annulation ou concellation). Voici quelques exemples qui ont marqués la vie humaine :

1.1. Missile Patriot

En février 1991, pendant la Guerre du Golfe, une batterie américaine de missiles Patriot, à Dharan (Arabie Saoudite), a échoué dans l'interception d'un missile Scud irakien. Le Scud a frappé un baraquement de l'armée américaine et a tué 28 soldats. La commission d'enquête a conclu à un calcul incorrect du temps de parcours, dû à un problème d'arrondi. Les nombres étaient représentés en virgule fixe sur 24 bits, donc 24 chiffres binaires. Le temps était compté par l'horloge interne du système en $1/10$ de seconde. Malheureusement, $1/10$ n'a pas d'écriture finie dans le système binaire : $1/10 = 0,1$ (dans le système décimal) = $0,0001100110011001100110011...$ (dans le système binaire). L'ordinateur de bord arrondissait $1/10$ à 24 chiffres, d'où une petite erreur dans le décompte du temps pour chaque $1/10$ de seconde. Au moment de l'attaque, la batterie de missile Patriot était allumée depuis environ 100 heures, ce qui avait entraîné une accumulation des erreurs d'arrondi de 0,34s. Pendant ce temps, un missile Scud parcourt environ 500 m, ce qui explique que le Patriot soit passé à côté de sa cible. Ce qu'il aurait fallu faire c'était redémarrer régulièrement le système de guidage du missile.

1.2. Explosion d'Ariane 5

Le 4 juin 1996, une fusée Ariane 5, a son premier lancement, a explosé 40 secondes après l'allumage. La fusée et son chargement avaient coûté 500 millions de dollars. La commission d'enquête a rendu son rapport au bout de deux semaines. Il s'agissait d'une erreur de programmation dans le système inertiel de référence. À un moment donné, un nombre codé en virgule flottante sur 64 bits (qui représentait la vitesse horizontale de la fusée par rapport à la plate-forme de tir) était converti en un entier sur 16 bits. Malheureusement, le nombre en question était plus grand que 32768 (overflow), le plus grand entier que l'on peut coder sur 16 bits, et la conversion a été incorrect.

1.3. Bourse de Vancouver

En 1982, la bourse de Vancouver a créé un nouvel indice avec une valeur nominale de 1000. Après chaque transaction boursière, cet indice était recalculé et tronqué après le troisième chiffre décimal et, au bout de 22 mois, la valeur obtenue était 524.881, alors que la valeur correcte était 1098.811. Cette différence s'explique par le fait que toutes les erreurs d'arrondi étaient dans le même sens : l'opération de troncature diminuait à chaque fois la valeur de l'indice.