

Chrysostome Baskiotis

Laurence Lamoulié

Fascicule 1

ANALYSE NUMÉRIQUE

Analyse des erreurs



Année 2012

INTRODUCTION

L'objectif de l'analyse numérique est d'obtenir pour un problème qui est exposé en termes mathématiques, des réponses numériques. Par réponse numérique on entend une solution « approchée » du problème, par opposition à la solution « exacte » obtenue par des méthodes mathématiques, à savoir algébriques ou analytiques. Il ne faut pas toutefois restreindre le domaine d'application de l'analyse numérique uniquement aux problèmes qui ont une solution mathématique précise, car, en effet, un autre grand domaine d'utilisation de l'analyse numérique concerne la résolution des problèmes mathématiques pour lesquels on ne connaît pas de façon analytique la solution « exacte ».

La matière première de l'analyse numérique est le nombre, qu'il s'agisse d'un entier, d'un réel ou, encore, d'un complexe. Nous commencerons donc l'étude de l'analyse numérique par les nombres et nous nous apercevrons, avec étonnement, que les relations des nombres avec l'ordinateur sont souvent difficiles. En effet, le passage d'un nombre par un ordinateur numérique peut modifier de façon significative sa précision et, par voie de conséquence, son exactitude et sa valeur. Un ordinateur « comprend » donc avec des erreurs les nombres réels que nous lui fournissons et aussi quand il effectue des calculs numériques, il n'est pas rare que ces calculs soient entachés d'erreurs. Expliquer, dans le cadre de l'algèbre linéaire, les raisons de ces erreurs, les analyser et les minimiser, autant que faire se peut, est le but de ce cours.

Le cours se divise en trois parties : La première concerne l'étude des erreurs provoquée par la représentation tronquée des valeurs numériques par l'ordinateur. Les deux autres parties sont consacrées à l'algèbre linéaire. Ainsi, à la deuxième partie on étudiera les effets des perturbations des valeurs des vecteurs et matrices à la résolution des systèmes d'équations linéaires par des méthodes directes et itératives. Enfin à la troisième partie on étudiera le calcul numérique des valeurs et vecteurs propres. Quelques applications, comme la méthode des moindres carrés ou l'approximation des images par la décomposition en valeurs singulières, seront aussi présentées.

1

ANALYSE DES ERREURS

1.1	Introduction	3
1.2	L'ordinateur et les nombres	5
1.2.1	Conversion décimale-binaire	6
1.2.2	Élaboration de la forme flottante	7
1.3	Nombres normalisés et spéciaux	9
1.4	Arithmétique arrondie	12
1.4.1	Exercices	14
1.5	Les nombres réels et les nombres-machine	15
1.5.1	Étude de l'erreur de précision relative	15
1.5.2	Exercices	17

Un ordinateur M utilise pour la représentation des valeurs numériques un mot machine qui a un nombre fini de bits. Par conséquent l'ensemble des nombre réels \mathbb{R} sera « vu » par l'ordinateur comme un sous-ensemble fini $M(\mathbb{R}) \subset \mathbb{R}$ (c-à-d. l'ordinateur opère un échantillonnage – au sens traitement du signal du terme – de \mathbb{R}). Cette représentation est une application $f_l : \mathbb{R} \rightarrow M(\mathbb{R})$, telle que pour tout réel x l'application f_l fournit une valeur $f_l(x)$, appelée nombre-machine. Dans la mesure où, en général, nous avons $x \neq f_l(x)$, il s'ensuit qu'une erreur est commise chaque fois que nous demandons à un ordinateur de manipuler une valeur numérique.

L'objectif de ce chapitre est d'étudier, de façon détaillée, les erreurs numériques provoquées par un ordinateur.

1.1 Introduction

Un nombre réel quelconque en représentation décimale est caractérisé par deux propriétés :

- L'exactitude qui est le nombre de digits⁽¹⁾ significatifs
- La précision qui est l'ordre de grandeur du dernier digit significatif.

Ainsi pour le nombre 0.0017892 nous avons une exactitude de 5 et une précision de 10^{-7} . Le passage de ce nombre par un ordinateur numérique peut modifier de façon significative sa précision et, par voie de conséquence son exactitude et sa valeur. Un ordinateur « comprend »

1. Pour un chiffre décimal on utilisera en abrégé le mot digit.

donc avec des erreurs les nombres réels que nous lui fournissons et aussi quand il effectue des calculs numériques.

Utiliser donc un ordinateur pour représenter un nombre réel, a comme conséquence de faire (presque) inévitablement une erreur. Pourquoi ? C'est assez simple de comprendre les raisons. La droite des réels $]-\infty, +\infty[$ a la puissance du continu. Cet ensemble infini de nombres nous essayons de le représenter par un ordinateur qui, pour ce faire, utilise des bits en nombre fini et même limité ! Dès lors on est devant une situation où entre deux nombres réels consécutifs de l'ordinateur il y a une infinité de nombres réels que l'ordinateur ne peut pas représenter et ne représentera jamais. Il tentera seulement de les approcher en faisant ainsi une erreur sur la valeur exacte du nombre que l'on appellera *erreur de représentation*.

Plus concrètement, considérons un nombre binaire :

$$b_q b_{q-1} \dots b_1 . b_{-1} \dots b_{-p}$$

Nous écrivons ce nombre sous forme normalisée, à savoir en partie entière on place un 0 et le reste se trouve en partie non entière (après le point) suivi d'une puissance de 10. Par exemple, pour le nombre précédent, sa forme normalisée est

$$0.b_q b_{q-1} \dots b_1 b_{-1} \dots b_{-p} \times 10^q \quad (2)$$

La partie après le point $b_q b_{q-1} \dots b_1 b_{-1} \dots b_{-p}$ s'appelle *mantisse*⁽³⁾ et pour un nombre réel il a , en général, une infinité d'éléments. Dans notre cas la mantisse est composée de $p+q$ bits, où p peut être infini. Pour stocker un réel dans un ordinateur, il faut pouvoir stocker sa mantisse (et aussi son exposant q , mais pour l'instant on ne se préoccupe de celui-ci). Étant donné qu'un nombre est stocké dans un mot de la mémoire de l'ordinateur et le mot est composé d'un nombre fini de bits, on aboutit à la conclusion que le stockage d'un réel dans un ordinateur dont la mémoire à m bits, avec $m < p+q$, provoque la troncature de sa valeur à $0.b'_1 \dots b'_{-m}$, avec $b'_1 = b_q, b'_2 = b_{q-1}, \dots$

Remarquons que nous procédons de la même façon quand on travaille avec des nombres réels exprimés en base 10, et qu'après une opération on décide de s'arrêter après, par exemple, cinq chiffres décimaux.

L'objectif de la première partie du cours est l'étude des erreurs de représentation et leur influence sur les résultats des calculs lors du déroulement d'un algorithme. Nous commençons par l'étude des nombres réels et de leur représentation par un ordinateur.

2. En général dans la forme normalisée, on place après le point tous les chiffres significatifs, c'est-à-dire le premier chiffre après la virgule doit être différent de zéro. Par exemple la forme normalisée de 0.0001 est 0.1×10^{-3} .

3. Notons que dans des temps plus anciens on appelait cette partie *significande*. Apparemment l'aspect sémantique dans la terminologie informatique n'est pas au goût du jour. D'où l'introduction du terme *mantisse* qui signifie partie décimale d'un logarithme.

EN SUBSTANCE

- **Écriture sous forme normalisée d'un nombre :**
 Soit le nombre 34.10285. En écriture normalisée on a 0.3410285×10^2 , c'est-à-dire :
 - 0 en partie entière ;
 - passage du contenu de la partie entière au début de la partie décimale, et
 - multiplication avec une puissance de la base de numérotation $10.$ ^(a)
- **Dans un ordinateur les nombres sont stockés sous forme normalisée.**
- **Un ordinateur dispose d'un nombre limité de places pour stocker les chiffres de la partie décimale d'une valeur sous forme normalisée. Étant donné qu'un réel a, le plus souvent, un nombre infini de chiffres en partie décimale, le stockage exact d'un tel nombre dans un ordinateur est impossible. Le nombre réellement stocké est une approximation du nombre fourni par l'utilisateur.**

a. La base de numérotation sera notée toujours 10, indépendamment du système de numérotation (décimal, binaire, hexadécimal, ...).

1.2 L'ordinateur et les nombres

D'abord il faut se rappeler qu'un ordinateur traite des nombres en base binaire et non pas décimale. Mais ce fait ne doit pas créer un problème, car un nombre binaire est comme un nombre décimal. Il a un point que l'on qualifiera de binaire (équivalent du point décimal) et il est composé de 0 et de 1. Par exemple le nombre décimal 13.125 est l'écriture condensée du nombre $1 \times 10^1 + 3 \times 10^0 + 1 \times 10^{-1} + 2 \times 10^{-2} + 5 \times 10^{-3}$ en base 10. De même le nombre binaire 11011.001 est l'écriture condensée du nombre $2^4 + 2^3 + 2^1 + 2^0 + 2^{-3}$ en base 2. Nous savons aussi que les nombres réels peuvent être représentés sous une forme qui est appelée notation scientifique. Ainsi le nombre 13.125 peut se mettre sous la forme 1.3125×10^1 . De la même façon nous pouvons avoir une notation scientifique pour les nombres binaires. Par exemple le nombre binaire 11011.001 peut s'écrire à l'aide de cette notation $1.1011001 \times 10^{100}$.

Un ordinateur dispose, pour représenter les nombres, des registres dont le nombre de bits est limité habituellement à 16, 32 ou 64 bits. Donc pour pouvoir représenter les nombres réels par un ordinateur nous devons, si on veut tenir compte du nombre fini de bits d'un registre, introduire deux limitations à la représentation par notation scientifique des nombres :

- (1) Le nombre de chiffres qu'on utilise après le point est fixé d'avance. Si par exemple on fixe à 4 ce chiffre, le nombre 11011.001 sera représenté par 1.1011×10^{100} .
- (2) Le nombre de chiffres qu'on utilise pour la valeur de l'exposant de la base est fixé d'avance. Par exemple si on fixe, sur la représentation précédente, le nombre de chiffres de l'exposant à 2, on aura finalement comme représentation 1.1011001×10^{10} .

Cette représentation – qui peut être tronquée – des nombres s'appelle forme flottante. Les ordinateurs actuels utilisent pour les nombres flottants le standard IEEE-754 dont les caractéristiques

pour les nombres en simple précision sont données par la table 1.1.

Signe s	Exposant e	Mantisse m
$s = 1 \text{ bit}$	$q = 8 \text{ bits}$	$p = 23 \text{ bits}$
□	□□□□□□□□	□□□□□□□□□□□□□□□□□□□□□□□□
\pm	$a_1 a_2 a_3 a_4 a_5 a_6 a_7 a_8$	$b_1 b_2 b_3 b_4 b_5 b_6 b_7 b_8 b_9 b_{10} b_{11} b_{12} b_{13} b_{14} b_{15} b_{16} b_{17} b_{18} b_{19} b_{20} b_{21} b_{22} b_{23}$

TABLE 1.1 – Standard IEEE-754. Simple précision

1.2.1 Conversion décimale – binaire

Nous allons examiner cette structure du nombre flottant en simple précision à l'aide d'un exemple. Prenons le nombre 465.463. Pour le convertir en flottant binaire normalisé on doit

- d'abord convertir séparément la partie entière et la partie décimale;
- ensuite réunir les deux parties et
- à la fin, de normaliser le résultat de la réunion.

Nous avons ainsi pour la partie entière :

- (1) $465 / 2 = 232$ Reste = 1
- (2) $232 / 2 = 116$ Reste = 0
- (3) $116 / 2 = 58$ Reste = 0
- (4) $58 / 2 = 29$ Reste = 0
- (5) $29 / 2 = 14$ Reste = 1
- (6) $14 / 2 = 7$ Reste = 0
- (7) $7 / 2 = 3$ Reste = 1
- (8) $3 / 2 = 1$ Reste = 1

d'où $(465)_{10} = (111010001)_2$.

Pour la partie décimale, nous avons :

- (1) $0.4063 * 2 = 0.926$ Partie entière = 0
- (2) $0.926 * 2 = 1.852$ Partie entière = 1
- (3) $0.852 * 2 = 1.704$ Partie entière = 1
- (4) $0.704 * 2 = 1.408$ Partie entière = 1
- (5) $0.408 * 2 = 0.816$ Partie entière = 0
- (6) $0.816 * 2 = 1.632$ Partie entière = 1
- (7) $0.632 * 2 = 1.264$ Partie entière = 1

(8)	$0.264 * 2 = 0.528$	Partie enti\`{e}re = 0
(9)	$0.528 * 2 = 1.056$	Partie enti\`{e}re = 1
(10)	$0.056 * 2 = 0.112$	Partie enti\`{e}re = 0
(11)	$0.112 * 2 = 0.224$	Partie enti\`{e}re = 0
(12)	$0.224 * 2 = 0.448$	Partie enti\`{e}re = 0
(13)	$0.448 * 2 = 0.896$	Partie enti\`{e}re = 0
(14)	$0.896 * 2 = 1.792$	Partie enti\`{e}re = 1
(15)	$0.792 * 2 = 1.584$	Partie enti\`{e}re = 1
(16)	$0.584 * 2 = 1.168$	Partie enti\`{e}re = 1
(17)	$0.168 * 2 = 0.336$	Partie enti\`{e}re = 0
(18)	$0.336 * 2 = 0.672$	Partie enti\`{e}re = 0
(19)	$0.672 * 2 = 1.344$	Partie enti\`{e}re = 1
(20)	$0.344 * 2 = 0.688$	Partie enti\`{e}re = 0
(21)	$0.688 * 2 = 1.376$	Partie enti\`{e}re = 1
(22)	$0.376 * 2 = 0.752$	Partie enti\`{e}re = 0
(23)	$0.752 * 2 = 1.504$	Partie enti\`{e}re = 1
(24)	$0.504 * 2 = 1.008$	Partie enti\`{e}re = 1
(25)	$0.008 * 2 = 0.016$	Partie enti\`{e}re = 0
(26)	$0.016 * 2 = 0.032$	Partie enti\`{e}re = 0
(27)	$0.032 * 2 = 0.064$	Partie enti\`{e}re = 0
(28)	$0.064 * 2 = 0.128$	Partie enti\`{e}re = 0
(29)	$0.128 * 2 = 0.256$	Partie enti\`{e}re = 0
(30)	$0.256 * 2 = 0.512$	Partie enti\`{e}re = 0
(31)	

d'où $(0.463)_{10} = (0.0111011010000011100101011000000 \dots)_2$.

Ainsi nous avons la conversion $(465.463)_{10} = (111010001.011101101000001110010101100000\dots)_2$ qui sous forme normalisée devient $0.11101000101110110100000\dots \times 10^{1001}$. Il faut maintenant transformer ce nombre binaire, qui a un nombre de chiffres infini, en forme binaire flottante en simple précision, c'est-à-dire calculer les valeurs du signe, de l'exposant et de la mantisse.

1.2.2 Élaboration de la forme flottante

Selon le standard IEEE-754, le signe est codé sur un bit qui a la valeur 0 si le nombre est positif et 1 dans le cas contraire. Ici nous avons un nombre positif, donc le bit du signe est égal à 0.

Pour l'exposant le standard utilise 8 bits. Nous pouvons donc disposer, pour le codage de l'exposant, de 256 valeurs différentes. Il faut aussi tenir compte du fait que l'exposant peut être positif ou négatif. On peut, bien sûr, avoir ici aussi un bit de signe. Mais dans ce cas les opérations d'addition et de soustraction entre différents nombres deviendraient ardues. Comme donc, on

ne veut pas avoir un codage pour les exposants avec des nombres positifs et négatifs, on partage l'intervalle de variation de l'exposant en deux parties égales : la première de 0 à 126 consacrée au codage des exposants négatifs et l'autre partie de 128 et 255 consacrée au codage des exposants positifs, la valeur 127 étant réservée à l'exposant 0.

Dans l'exemple précédent, l'exposant est égal à $1001_2 = 9_{10}$. C'est donc une valeur positive. Par conséquent pour l'exposant 9 on rajoute 127 et on obtient ainsi $9 + 127 = 136 = (10001000)_2$.

La mantisse doit être 11101000101110110100000. Remarquons maintenant que si, dans le cas des nombres décimaux, on décidait de placer le premier chiffre significatif (chiffre non nul) avant le point décimal, ce chiffre serait soit 1, soit 2, ..., soit 9. Mais dans le cas des nombres binaires ce chiffre serait obligatoirement 1. Par conséquent on peut toujours, décider de placer le premier chiffre significatif avant le point binaire et le reste après. Ce chiffre est obligatoirement égal à 1. on peut donc ne pas le stocker mais on doit, dans ce cas, se rappeler de tenir compte lors du décodage. Le chiffre qui n'apparaît pas au codage s'appelle le bit caché et cette technique permet de gagner un bit supplémentaire à la mantisse pour le stockage du nombre. Ainsi la mantisse s'écrit 111010001011101101000001 et l'exposant diminue d'une unité et devient égal à $8 + 127 = 135 = (10000111)_2$ et on obtient donc pour la forme flottante binaire normalisée du nombre 465.463 la représentation suivante :

0	10000111	1110100010111011010000000000
---	----------	------------------------------

Si maintenant on fait la démarche inverse, on peut calculer la valeur décimale du nombre binaire que nous venons de créer en utilisant la conversion en simple précision. Nous avons

$$\begin{aligned}
 1.1110100010111011010000000000 \times 10^{10000111-01111111} &= 111010001.011101101101000001 \\
 &= 2^8 + 2^7 + 2^6 + 2^4 + 2^0 \cdot 2^{-2} + 2^{-3} + 2^{-4} + 2^{-6} + 2^{-7} + 2^{-9} + 2^{-15} \\
 &= 465.462921142578125
 \end{aligned}$$

Remarquons qu'en vue du décodage en décimal, nous avons incorporer en première position, avant le point, le bit caché qui ne figurait pas dans la représentation flottante.

Bien sûr la question qui vient à l'esprit concerne la capacité de ce système de représenter, même approximativement, n'importe quel nombre. En effet nous pouvons envisager que, comme conséquence d'une opération arithmétique, nous avons un nombre très grand (ou très petit) qui ne peut pas être représenté comme un binaire flottant normalisé. Avant d'examiner en détail cette éventualité, il faut établir les valeurs extrêmes que ce système peut représenter ainsi que la résolution du système, c'est-à-dire la valeur de la plus petite différence entre deux nombres que le système est capable de détecter.

EN SUBSTANCE

- **Forme flottante en simple précision selon le standard IEEE-754 :**

$$x = (-1)^S \times M^E \text{ où}$$

- x le nombre,
- S le signe stocké sur un bit,
- M la mantisse stockée sur 23 bits, et
- E l'exposant stocké sur 8 bits.

- **Conversion d'un nombre $V.D$ en base décimale vers un nombre en base β .**-
Notons

– par $(V;v,r)_\beta$ le résultat de la division de V par β , c'est-à-dire $V = v \times \beta + r$.

– par $[D;d,f]_\beta$ le résultat de la multiplication de $0.D$ par β , c'est-à-dire $d.f = 0.D \times \beta$.

Posons $v_0 = V$ et considérons la suite de divisions

$$(v_0;v_1,r_0)_\beta, (v_1;v_2,r_1)_\beta, \dots, (v_{p-1};v_p,r_{p-1})_\beta$$

avec $r_{p-1} < \beta$.

Considérons aussi la suite de multiplications

$$[D;d_1,f-1]_\beta, [f_1;d_2,f_2]_\beta, \dots, [q-1;d_q,f_q]_\beta, \dots$$

Alors la conversion du nombre décimal $V.D$ en base B est donnée par

$$v_p r_p r_{p-1} \dots r_1 r_0 . d_1 d_2 \dots d_q \dots$$

- **Conversion d'un nombre $v_p r_p r_{p-1} \dots r_1 r_0 . d_1 d_2 \dots d_q \dots$ en base β vers un nombre $V.D$ en base décimale.**- **Supposons que la valeur β de la base est exprimée en base décimale, c'est-à-dire $\beta = \beta_{10}$. Alors**

$$V.D = v_p \times \beta^p + r_{p-1} \times \beta^{p-1} + \dots + r_1 \times \beta^1 + r_0 \times \beta^0 . d_1 \times \beta^{-1} + d_2 \times \beta^{-2} + \dots + d_q \times \beta^{-q} + \dots$$

- **Particularité du codage en forme flottante binaire. Le chiffre le plus significatif n'est pas stocké dans un bit (bit caché), mais il est toujours pris en compte lors des calculs et du passage vers la forme flottante décimale.**

1.3 Nombres normalisés et spéciaux

Nous venons de voir que la représentation flottante d'un nombre x est donnée par

$$x = \pm (b_0 . b_1 b_2 \dots b_p) \times 10^{e_1 \dots e_q}$$

avec $b_0 = 1$ et $1 \leq b_1 b_2 \dots b_p < 2$. Le fait que nous avons toujours $b_0 = 1$ permet de ne pas stocker b_0 et de réserver les $p = 23$ bits de la mantisse pour la partie décimale et d'augmenter ainsi la précision du codage de 2^{-22} à 2^{-23} . b_0 est appelé le bit caché de la normalisation. Cette technique, qui a été utilisée pour la première fois sur un ordinateur Vax dans les années 70, a aussi un inconvénient. Celui de ne pas pouvoir stocker de façon simple la valeur 0. En effet si nous suivons la procédure pour le codage flottant binaire, qui vient d'être décrite, on s'aperçoit qu'il faut poser $b_0 = 0$! Comme cette solution est inapplicable on conviendra que le nombre

0	00000000	000000000000000000000000
---	----------	--------------------------

avec $b_0 = 1$ représente effectivement la valeur 0, mais, pour éviter toute ambiguïté, l'exposant 00000000 ne sera utilisé pour aucune autre représentation. Donc le plus petit nombre positif que nous pouvons stocker est

0	00000001	000000000000000000000000
---	----------	--------------------------

ce qui donne $1.0 \times 2^{1-127} = 2^{-126}$ et qui signifie que le domaine de variation de l'exposant E n'est pas entre -127 et 127 mais entre -126 et 127.

En utilisant les notations de la table 1.1 on établit la table 1.2.

Bits de l'exposant $a_1 \cdots a_8$	La valeur numérique correspondante
$(00000000)_2 = 0$	$\pm (0.b_1 b_2 \cdots b_{23})_2 \times 2^{-127}$
$(00000001)_2 = 1$	$\pm (1.b_1 b_2 \cdots b_{23})_2 \times 2^{-126}$
$(00000010)_2 = 2$	$\pm (1.b_1 b_2 \cdots b_{23})_2 \times 2^{-125}$
\vdots	\vdots
$(01111111)_2 = 127$	$\pm (1.b_1 b_2 \cdots b_{23})_2 \times 2^0$
$(10000000)_2 = 128$	$\pm (1.b_1 b_2 \cdots b_{23})_2 \times 2^1$
\vdots	\vdots
$(11111101)_2 = 253$	$\pm (1.b_1 b_2 \cdots b_{23})_2 \times 2^{126}$
$(11111110)_2 = 254$	$\pm (1.b_1 b_2 \cdots b_{23})_2 \times 2^{127}$
$(11111111)_2 = 255$	$\pm \infty$ si $b_1 = b_2 = \cdots = b_{23} = 0$, NaN sinon.

TABLE 1.2 – Nombres du standard IEEE-754, simple précision

Tous les nombres de ce tableau, sauf ceux de la première et de la dernière ligne, sont des nombres flottants binaires normalisés en simple précision. La première ligne indique le codage de la valeur 0. Notons qu'en fonction de la valeur du bit du signe (0 ou 1) nous avons deux zéros, un positif et un négatif. On peut aussi, à l'aide de ce tableau, calculer le plus petit et le plus grand nombre de la représentation en simple précision.

Le plus petit nombre normalisé peut être représenté par

0	00000001	000000000000000000000000
---	----------	--------------------------

ce qui équivaut à $(1.000\dots 0)_2 = 1 \times 10^{00000001} \approx 1.2 \times 10^{-38}$. On le notera par la suite m_p .

Le plus grand nombre normalisé est représenté par

0	11111110	11111111111111111111111111111111
---	----------	----------------------------------

ce qui équivaut à $(1.11\dots 1)_2 \times 10^{11111110} = 2 - 2^{-23} \times 2^{127} \approx 3.4 \times 10^{38}$. On le notera par la suite m_G .

Un examen attentif de cette table montre qu'il existe des nombres soit plus petits, soit plus grands que ceux qui sont normalisés. Ce sont justement les nombres spéciaux que seul l'ordinateur peut utiliser. Par exemple le plus petit nombre non nul qu'on peut stocker est $2^{-149} = [0.00\dots 01] = 0.00\dots 1 \times 2^{00000001}$ dont la représentation, en suivant la technique pour la représentation de zéro, doit être

0	00000000	000000000000000000000000000000000000000001
---	----------	--------------------------------------------

Ces nombres, qu'ils sont plus petits que le plus petit nombre normalisé, c-à-d. plus petits que 1.2×10^{-38} , sont appelés nombres sous-normalisés. Ils ne peuvent pas être normalisés car, dans ce cas, le résultat serait un exposant qui n'est pas dans le domaine de variation entre -126 et 127. Les nombres sous-normalisés sont moins précis que les nombres normalisés, c'est-à-dire que leur utilisation peut provoquer des erreurs de calcul plus importantes. Si un résultat d'une opération arithmétique conduit à un nombre sous-normalisé, on dit qu'on est en présence d'un *underflow*.

Nous pouvons aussi envisager des résultats des opérations arithmétiques qui sont supérieurs au plus grand nombre normalisé, à savoir 3.4×10^{38} . Dans ce cas tous les bits de l'exposant sont égaux à 1. Cette situation est décrite par dernière ligne du tableau 1.2. Deux cas peuvent se présenter : soit la mantisse est nulle et dans ce cas on dit que la valeur représentée par ce flottant est l'infini – positif ou négatif, selon le bit du signe, soit la mantisse n'est pas nulle et dans ce cas on est en présence d'un non-nombre, noté NaN pour *Not a Number*.

EN SUBSTANCE

Il y a cinq catégories de nombres flottants binaires selon les valeurs des bits de l'exposant et de la mantisse.

- **Nombres normalisés** : l'exposant a au moins un bit différent de 0 et aussi au moins un bit égal à 0.
- **Nombres sous-normalisés** : L'exposant est nul et la mantisse a au moins un bit différent de 0.
- **Zéro** : L'exposant et la mantisse sont nuls. Il y a un zéro positif et un zéro négatif en fonction de la valeur du bit du signe.
- **Nan** : L'exposant est composé exclusivement de 1 et la mantisse a au moins un bit différent de 0.
- **Infini** : L'exposant est composé exclusivement de 1 et la mantisse de 0. Il y a un infini positif et un négatif, en fonction de la valeur du bit du signe.

1.4 Arithmétique arrondie

D'après ce qui vient d'être présenté aux deux sections précédentes, il s'ensuit que le résultat x d'une opération arithmétique avec des réels, peut rarement faire partie des nombres représentés par la table 1.2. Cinq cas peuvent se présenter :

- (1) Le résultat est un des nombres normalisés de la table 1.2.
- (2) Le résultat x est entre m_p et m_G mais ne fait pas partie des nombres normalisés du tableau 1.2.
- (3) Le résultat x est un nombre sous-normalisé.
- (4) Le résultat x est une valeur infinie.
- (5) Le résultat x est un NaN.

Les trois derniers cas ne nous intéressent pas. L'ordinateur indiquera le problème et, soit il poursuivra l'exécution du programme, soit il s'arrêtera. Le premier cas ne pose aucun problème. Nous examinerons donc le deuxième cas. Il est évident que pour pouvoir poursuivre l'exécution du programme, l'ordinateur assimilera le résultat x , dont la représentation binaire est $0.b_1b_2 \dots b_{22}b_{23}b_{24} \dots$, à un des nombres normalisés qui sont présents dans la table 1.2. La raison pour laquelle x ne fait pas partie des nombres flottants normalisés est que sa représentation binaire contient des bits b_i non nuls pour $i > p = 23$. L'opération donc de normalisation du nombre x consiste à trouver une représentation binaire avec 23 bits $0.b_1b_2 \dots b_{22}b_{23}$ avec

$$b_i = b_i; 1 \leq i \leq 22$$

Le 23-ème bit et, éventuellement, l'exposant E seront modifiés. Il y a quatre possibilités pour cette modification :

- (1) Soit l'ordinateur ne fait aucune modification : $b_{23} = b_{23}$ et l'exposant reste inchangé.
- (2) Soit il utilisera le nombre flottant x_- de la table 1.2 qui est le plus proche de x et qui est plus petit que x : $b_{23} = \begin{cases} b_{23}, & \text{si } s = 0 \\ b_{23} + 1, & \text{si } s = 1 \end{cases}$, avec modification de l'exposant dans le cas de propagation de la retenue.

- (3) Soit il utilisera le nombre flottant x_+ de la table 1.2 qui est le plus proche de x et qui est plus grand que x : $b_{23} = \begin{cases} b_{23} + 1 & \text{si } s = 0 \\ b_{23}, & \text{si } s = 1 \end{cases}$, avec modification de l'exposant dans le cas de propagation de la retenue.
- (4) Soit il utilisera le nombre flottant x_{\pm} de la table 1.2 qui est le plus proche de x : $b_{23} = b_{23} + b_{24}$ avec modification de l'exposant dans le cas de propagation de la retenue.

Le standard IEEE-754 préconise d'utiliser l'arrondi vers le nombre x_{\pm} le plus proche de x qu'on appellera par la suite représentation par l'arrondi (le plus proche). Nous noterons par $f_l(x)$ le résultat de cet arrondi et on aura, selon ce qui vient d'être dit, $|f_l(x) - x| = \min\{|x_- - x|, |x_+ - x|\}$. Il y a des systèmes d'exploitation qui utilisent aussi la représentation par x_- qu'on appellera représentation par troncature.

Afin d'évaluer l'erreur de calcul provoquée par l'arrondi nous devons connaître la résolution de notre système de représentation des flottants normalisés ou, comme on dit, la précision de la machine.

DÉFINITION 1.4.1 La précision ϵ_{ps} d'un ordinateur est le plus petit nombre positif normalisé tel que

$$1 + \epsilon_{ps} = 1$$

On peut aussi envisager la précision de la machine comme étant la plus grande erreur relative de précision qui nous pouvons avoir lors de la représentation par troncature d'un nombre réel.

Il ne faut pas confondre ce nombre avec le plus petit nombre flottant positif que l'ordinateur peut représenter qui, comme nous avons vu précédemment, est un nombre flottant sous-normalisé.

Pour calculer la précision de la machine on commence par la représentation binaire normalisée, avec le bit caché, de la valeur $1 = (1.00\dots0)_2 \times 10^0$. On a donc

0	01111111	000000000000000000000000
---	----------	--------------------------

ce qui montre que ce nombre fait partie du tableau 1.2. En utilisant cette table on peut représenter le nombre normalisé qui est immédiatement supérieur à 1. Il est donné par

0	01111111	000000000000000000000001
---	----------	--------------------------

qui équivaut en décimal à $1 + 2^{-23}$. La différence entre ce deux nombres est égale à 2^{-23} qui est la précision de la machine. Nous avons donc pour la précision de la machine le fait suivant :

FAIT 1.1 La précision eps d'une machine simple précision qui suit la norme IEEE-754 est

$$\text{eps} = 2^{-23} \approx 1.192 \times 10^{-7}$$

Nous pouvons maintenant avoir une idée concernant l'erreur d'arrondi. Si x est la valeur qu'on veut représenter selon le standard IEEE-754, alors on la remplacera soit par x_- , soit par x_+ . Cette représentation est notée $\text{fl}(x)$. La valeur absolue de la différence entre $\text{fl}(x)$ et x n'est pas supérieure à la moitié de la différence entre x_+ et x_- . Ainsi, si le nombre x s'écrit en binaire

$$x = (b_0.b_1b_2\dots b_{23}b_{24}\dots)_2 \times 2^E$$

avec $b_0 = 1$, alors sa représentation machine sera

$$\text{fl}(x) = (b_0.b_1b_2\dots b_{23})_2 \times 10^E$$

Donc

$$|\text{fl}(x) - x| \leq 2^{-24} \times 10^E$$

ce qui donne

$$|\text{fl}(x) - x| \leq \frac{1}{2} \text{eps} \times 10^E$$

EN SUBSTANCE

- **Il y a quatre types de nombres flottants binaires : nombres normalisés, sous-normalisés, infini et NaN.**
- **À tout nombre réel correspond un nombre machine, noté $\text{fl}(x)$ et qui, pour le standard IEEE-754, est donné par l'approximation $\text{fl}(x) = x_{\pm}$.**
- **La précision eps de la machine est le plus petit nombre normalisé tel que $1 + \text{eps} = 1$. Pour la simple précision on a $\text{eps} = 2^{-23} \approx 1.192 \times 10^{-7}$.**

1.4.1 Exercices

EXERCICE 1.1 Soit un système de représentation flottant avec base β , mantisse à p places, exposant E , avec $E_m \leq E \leq E_M$.

Calculer le nombre de valeurs normalisées qui peuvent être représentées par ce système.

Application : $\beta = 10, p = 3, E_m = -15, E_M = 16$.

EXERCICE 1.2 Soient trois réels $x = 0.125 \times 10^6$, $y = 0.437 \times 10^{12}$, $z = 0.215 \times 10^{-10}$. En utilisant le système de représentation de l'exercice précédent, calculer

- (1) la somme $x + y$ et commenter le résultat.
- (2) le produit $x \times y$ et commenter le résultat.
- (3) le quotient z/y et commenter le résultat.

EXERCICE 1.3 Soient trois réels $x = 0.4 \times 10^0$, $y = 0.4 \times 10^0$, $z = 0.1 \times 10^3$. En utilisant le système de représentation de l'exercice 1.1, calculer les deux sommes

$(x + y) + z$ et $x + (y + z)$

Conclusion.

1.5 Les nombres réels et les nombres-machine

Pour l'étude des différences entre un réel x et le nombre-machine $fl(x)$ qui est sa représentation dans l'ordinateur, nous utilisons trois types d'erreurs :

- Erreur de représentation, notée Δx ou $\varepsilon(x)$:

$$\Delta x = fl(x) - x \quad (1.1)$$

ou encore l'erreur absolue de représentation :

$$|\Delta x| = |fl(x) - x|$$

- Erreur relative de représentation, notée $\iota(x)$:

$$\iota(x) = \frac{\Delta x}{fl(x)} = \frac{fl(x) - x}{fl(x)} \quad (1.2)$$

- Erreur relative de précision, notée $\eta(x)$:

$$\eta(x) = \frac{\Delta x}{x} = \frac{fl(x) - x}{x} = \frac{fl(x)}{x} - 1 \quad (1.3)$$

ou encore l'erreur relative absolue de précision

$$|\eta(x)| = \frac{|\Delta x|}{|x|} = \frac{|fl(x) - x|}{|x|}$$

L'erreur relative de précision permet d'exprimer le nombre-machine comme suit :

$$fl(x) = x(1 + \eta(x)) \quad (1.4)$$

1.5.1 Étude de l'erreur de précision relative

Nous avons vu qu'un ordinateur utilise pour la représentation des valeurs réelles un mot machine dont les bits sont partagés en

- un bit de signe,
- p bits de mantisse, et
- q bits d'exposant.

Dans la suite on considère que la base de numérotation de l'ordinateur est β .

EXEMPLE 1.5.1 Considérons un ordinateur avec base de numérotation $\beta = 10$ et nombre de chiffres de la mantisse $p = 4$. Nous avons $123.4567 = 0.1234567 \times 10^3 = 0.1234 + 0.567 \times 10^{-4} \times 10^3$. Si l'ordinateur approche les valeurs réelles par troncature, alors $m(x) = 0.1234 \times 10^3$. Si l'approximation se fait par arrondi, alors $m(x) = 0.1235 \times 10^3$.

Cet exemple nous aide à comprendre que si nous avons un réel x , alors il sera représenté sous forme normalisée comme suit :

$$x = s \times w \times \beta^n \quad (1.5)$$

où s est le signe. On peut, encore écrire :

$$x = s \times (r + t \times \beta^{-p}) \times \beta^n, \text{ avec } \beta^{-1} \leq |r| < 1 \text{ et } 0 \leq |t| < 1 \quad (1.6)$$

Si le nombre-machine représentant x est obtenu par troncature, nous avons

$$fl(x) = s \times r \times \beta^n \quad (1.7)$$

tandis que pour approximation par arrondi, nous avons

$$fl(x) = \begin{cases} s \times r \times \beta^n, & \text{si } |x - r \times \beta^n| < |x - r \times \beta^n + \beta^{-p}| \\ s \times r \times \beta^n + \beta^{-p}, & \text{sinon} \end{cases}$$

Nous avons, pour l'erreur relative de précision, le résultat suivant :

THÉORÈME 1.5.1 (de la précision relative).- Soit un calculateur avec base de numérotation β et nombre de chiffres de la mantisse p . Alors l'erreur relative de précision $\eta(x)$ est bornée comme suit :

$$|\eta(x)| \leq \begin{cases} \beta^{1-p}, & \text{si approximation par troncature} \\ 0.5 \times \beta^{1-p}, & \text{si approximation par arrondi} \end{cases}; x \in \mathbb{R} \quad (1.8)$$

DÉMONSTRATION. D'après (1.6) et (1.7) on a

$$\Delta x = fl(x) - x = t \times \beta^{-p} \times \beta^n \quad (1.9)$$

d'où

$$|\eta(x)| = \frac{|\Delta x|}{|x|} = \frac{|t \times \beta^{-p} \times \beta^n|}{|r + t \times \beta^{-p} \times \beta^n|} = \frac{|t \times \beta^{-p}|}{|r + t \times \beta^{-p}|} \leq \frac{|t|}{|r|} \times \beta^{-p} \leq \beta^{1-p} \quad (1.10)$$

car $\beta^{-1} \leq |r| < 1$, $0 \leq |t| < 1$ et $\frac{|t|}{|r|} \leq \frac{1}{\beta^{-1}} = \beta$ dans le cas de troncature.

Dans le cas de l'arrondi, nous avons $0 \leq |t| < 1/2$ et donc $\frac{|t|}{|r|} \leq \frac{1}{2} \beta^{1-p}$. ■

Dans les calculs on utilise la valeur absolue de l'erreur de précision relative $|\eta(x)|$.

Parfois cette quantité est considérée comme étant la précision de la machine.

EN SUBSTANCE

- **Il y a trois types d'erreur de représentation :**
 - **Erreur de représentation :** $\Delta x = fl(x) - x$
 - **Erreur relative de représentation :** $\iota(x) = \frac{\Delta x}{fl(x)} = \frac{fl(x) - x}{fl(x)}$
 - **Erreur relative de précision :** $\eta(x) = \frac{\Delta x}{x} = \frac{fl(x) - x}{x} = \frac{fl(x)}{x} - 1$, avec $fl(x) = x(1 + \eta(x))$.
- **En règle générale, pour l'analyse des erreurs on utilise la valeur absolue de l'erreur de précision relative, qui pour le standard IEEE-754 est**

$$\eta(x) \leq \begin{cases} \beta^{1-p}, & \text{si approximation par troncature} \\ 0.5 \times \beta^{1-p}, & \text{si approximation par arrondi} \end{cases}$$

1.5.2 Exercices

EXERCICE 1.4 Soit le nombre 387.62000_{10} .

- (1) Calculer sa valeur en hexadécimal.
- (2) Convertir la valeur hexadécimale en décimale.
- (3) Calculer l'erreur absolue de représentation.
- (4) Calculer l'erreur relative absolue de précision.

EXERCICE 1.5 Considérons une machine décimale avec mantisse à 4 chiffres. Calculer l'erreur de représentation et l'erreur relative de représentation pour les nombres suivants :

9.023506, 158.26 et 0.001588946.

EXERCICE 1.6 Soient les nombres : 15.2750, 358.937 et 5233.618.

- (1) Convertir ces nombres en hexadécimal.
- (2) Convertir les nombres hexadécimaux ainsi obtenus en nombres décimaux avec exactitude de 8.
- (3) Calculer l'erreur relative de cette dernière conversion sous la forme $x.xx \times 10^{-7}$.

N.B. Pour la conversion en hexadécimal on utilisera le même nombre de chiffres hexadécimaux après la virgule qu'en valeur décimale.

EXERCICE 1.7 Soit un nombre réel $x \in \mathbb{R}$. Nous pouvons l'écrire sous la forme :

$$x = s \cdot m \cdot \beta^e, \text{ avec } \frac{1}{\beta} \leq m < 1$$

où s est le signe, m est la mantisse considérée sans limitation de bits, β est la base et e est l'exposant qui est une valeur entière. Dans un ordinateur, x sera représenté comme un flottant normalisé sous la forme

$$fl(x) = s \cdot M \cdot \beta^E, \text{ avec } fl(x) \in M(\mathbb{R})$$

où M est la mantisse limitée à p digits, E est l'exposant limité à q digits. Ainsi pour le codage en nombre flottants, on utilise $N = p + q + 1$ digits.

Dans la suite, on prendra $\beta = 2$.

On se propose de calculer l'erreur de la représentation pour différentes opérations arithmétiques.

- (1) Erreur de l'affectation : Montrer que

$$x - fl(x) = s \cdot 2^{E-p} \cdot \alpha$$

avec

- (a) $\alpha \in [-0.5, 0.5[$ si représentation par arrondi.
 - (b) $\alpha \in [0, 1[$ si représentation par troncature.
- (2) Calculer l'erreur d'une opération d'addition.
 - (3) Même chose pour l'opération de soustraction.
 - (4) Même chose pour l'opération de multiplication.
 - (5) Même chose pour l'opération de division.

2

Le standard IEEE 754

2.1	Formats de IEEE-754	20
2.1.1	Format simple précision	21
2.2	Exercices	22
2.3	Exercices de laboratoire	24

Les développements du chapitre précédent, nous permettent de conclure qu'un ordinateur tente de représenter l'infinité de nombres réels R par des moyens finis, à savoir les nombres-machine M . Cette représentation peut se faire de plusieurs manières et au début de l'informatique c'était effectivement le cas : avec presque chaque ordinateur, on s'ingénuaient à faire une présentation particulière. Bien sûr la comparaison des résultats numériques entre différents ordinateurs devenait impossible et la portabilité des programmes faisait appel plus à la chance qu'à la raison. Cette situation a pris fin en 1985 avec l'apparition du standard IEEE-754 qui équipe aujourd'hui la quasi totalité des ordinateurs, au point de parler d'une norme IEEE-754 et non pas d'un standard. Pour établir ce standard il a fallu choisir, entre autre, le type du codage, les opérations à utiliser et la méthode d'arrondi.

Ainsi le standard IEEE-754 a les caractéristiques suivantes :

- codage binaire en virgule flottante avec bit caché ;
- permet le calcul des quatre opérations arithmétiques et de l'extraction de la racine carrée ;
- préconise l'utilisation de l'arrondi au plus près, mais il accepte les quatre modes d'arrondi, présentés en 1.4, et
- considère chaque résultat comme exact, c'est-à-dire n'effectue pas la propagation des erreurs du calcul actuel sur les calculs qui se suivent⁽¹⁾.

Nous présentons dans la suite de manière détaillée le standard IEEE-754.

1. ce qui ne signifie pas que, lors d'une suite des calculs, il n'y a pas de propagation d'erreurs dues aux arrondis, mais simplement que, si on a fait un premier calcul, le calcul suivant considérera le résultat du premier comme exact et il ne tiendra pas compte, pour le nouveau calcul, de l'erreur d'arrondi du premier. Ainsi l'erreur sera incorporée dans le résultat et le prochain calcul considérera le calcul précédent comme exact.

2.1 Formats de IEEE-754

Selon la relation (1.6), la représentation d'un nombre flottant est donnée par

$$x = s \times M \times \beta^E \quad (2.1)$$

où s est le signe, M est la mantisse normalisée stockée sur p bits, β est la base et E est l'exposant qui est un nombre entier non négatif stocké sur q bits. La prise en compte de cette représentation par un ordinateur numérique, impose que :

- la mantisse normalisée ait une précision de $p-1$ bits⁽²⁾, et
- la valeur de E se trouve dans un intervalle $[E_{\min}, E_{\max}]$ déterminé par le nombre de bits q de l'exposant.

Concrètement un réel x peut se mettre sous la forme

$$x = (-1)^s \times (r + t \times \beta^{-(p-1)}) \times \beta^E, \text{ où } 1 \leq r < \beta \text{ et } 0 \leq t < 1$$

En numérotation binaire (c'est-à-dire $\beta = 2$) le réel x s'écrit

$$(-1)^s \times b_0.b_1b_2 \dots b_{p-1} \times \beta^E$$

avec

- $s = 0$ ou 1 ;
- $b_0 = 1$;
- $b_i = 0$ ou $1, i = 1, \dots, p-1$;
- $\beta = 2$, la base et
- E l'exposant, avec $E \in [E_{\min}, E_{\max}]$.

Comme le standard IEEE-754 utilise le codage avec bit caché, le réel x s'écrit selon ce standard

$$(-1)^s \times .b_1b_2 \dots b_{p-1} \times \beta^E$$

Pour calculer l'intervalle de variation $[E_{\min}, E_{\max}]$ de la valeur de l'exposant on utilise le biais qui est défini par la relation :

$$\text{biais} = B = \underbrace{0 \ 1 \dots 1}_{(q-1)\text{fois}} = 2^{(q-1)} - 1 \quad (2.2)$$

Donc le calcul de l'exposant en binaire se fait en rajoutant le biais à la valeur de la conversion en binaire de E , c-à-d. contenu de l'exposant = biais + $(E)_2$, où $(E)_2$ signifie la valeur E en base de numérotation 2.

La table 2.1 de la page suivante présente les différents formats du standard IEEE-754

L'application des formats est différente selon le langage de programmation. Nous donnons dans la table 2.2 les déclarations à faire selon le format choisi :

2. En effet selon le standard IEEE-754 p est le nombre de bits pour la mantisse et le signe. Étant donné que le signe occupe un bit dans tous les formats, il reste pour la mantisse $p-1$ bits.

Paramètre	Simple précision	Simple précision étendue	Double précision	Double précision étendue	Quadruple (Extension du standard)	Étendue (Introduite par Intel)
Nombre de bits de la mantisse plus le signe (p)	24	≥ 32	53	≥ 64	113	64
Nombre de chiffres décimaux exacts ($p / \log_2 10$)	7.22	≥ 9.63	15.95	≥ 19.26	34.01	19.26
Type du codage	Bit caché	Non spécifié	Bit caché	Non spécifié	Bit caché	Bit explicite
Nombre de bits pour la mantisse (p - 1)	23	≥ 31	52	≥ 63	112	64
E_{\max} E_{\min}	+127 -126	≥ +1023 ≤ -1022	+1023 -1022	≥ +16383 ≤ -16382	+16383 -16382	+16383 -16382
Biais de l'exposant	+127	Non spécifié	+1023	Non spécifié	+16383	+16383
Nombre de bits pour l'exposant	8	≥ 11	11	≥ 15	15	15
Nombre de bits pour le signe	1	1	1	1	1	1
Nombre de bits pour le format	32	≥ 43	64	≥ 79	128	80
Valeur max décimale ($2E_{\max} + 1$)	3.4028E+38	≥ 1.7976E+308	1.7976E+308	≥ 1.1897E+4932	1.1897E+4932	1.1897E+4932
Valeur min décimale ($2E_{\min}$)	1.1754E-38	≤ 2.2250E-308	2.2250E-308	≤ 3.3621E-4932	3.3621E-4932	3.3621E-4932
Valeur min décimale dénormalisée ($2E_{\min} - p + 1$)	1.4012E-45	≤ 1.0361E-317	4.9406E-324	≤ 3.6451E-4951	6.4751E-4966	1.8225E-4951

TABLE 2.1 – Les formats du standard IEEE-754

Langage	Simple précision	Double précision	Double étendue	Quadruple précision	Étendue
SCILAB		Par défaut	En interne 80 bits		
FORTRAN	REAL*4	REAL*8		REAL*16	REAL*10
C, C++	float	double		long double	long double

TABLE 2.2 – Déclarations des formats selon le langage de programmation

2.1.1 Format simple précision

Les caractéristiques de ce format sont les suivantes :

Les exposants 00000000 et 11111111 sont à usage réservé et ne sont pas utilisés pour la conversion. Par conséquent nous avons $E_{\min} = -\text{biais} + 1$ et $E_{\max} = \text{biais}$.

Un mot de 32 bits pour lequel tous ses bits sont à 0, représente la valeur décimale 0,0. Il est évident que, en fonction de la valeur du signe s, on a un zéro positif et un zéro négatif.

Pour l'infini nous avons les représentations suivantes :

- 0 11111111 00000000 0000000000000000 = $+\infty$ = 7F800000₁₆

- 1 11111111 00000000 0000000000000000 = $-\infty$ = FF800000₁₆.

Tous les nombres avec représentation

`s 11111111 ffffffffffffffffffffffffffff`

s'appellent NaN – Not a Number – s'il y a au moins un `f` qui est différent de zéro. Ce sont toujours des résultats indéfinis produits par une erreur de calcul, e.g. du type $0/0$ ou $\log(0)$. En particulier nous avons les deux NaN suivants :

– Not-a-Number-Signaling (NaNS) dont la représentation est :

`0 11111111 01111111 1111111111111111 = NaNS = 7FBFFFFFF` ₁₆.

– Not-a-Number-Quiet (NaNQ) dont la représentation est :

`0 11111111 10000000 0000000000000000 = NaNQ = 7FC00000` ₁₆.

NaNS engendre un arrêt anormal (trap) de l'exécution du programme tandis que NaNQ permet la poursuite de l'exécution du programme. Tout calcul numérique avec NaN comme entrée produit NaN comme sortie, c'est-à-dire nous avons, par exemple, $0 \times \text{NaN} = \text{NaN}$.

Nous avons donc pour un nombre décimal V , la représentation flottante suivante en simple précision :

`s xxxxxxxx ffffffffffffffffffffffff = s E F`

où s représente le bit du signe, E représente l'exposant biaisé et F la fraction décimale. Cette valeur est déterminée comme suit :

- Si $E = 255$ et $F = 0$, alors $V = \text{NaN}$.
- Si $E = 255$ et $F = 0$ et $s=0$, alors $V = \infty$.
- Si $E = 255$ et $F = 0$ et $s=1$, alors $V = -\infty$.
- Si $E = 0$ et $F = 0$ et $s=0$, alors $V = 0$.
- Si $E = 0$ et $F = 0$ et $s=1$, alors $V = -0$.
- Si $0 < E < 255$, alors $V = (-1)^s \times 2^{E-127} \times (1.F)$
- Si $E=0$ et $F = 0$, alors $V = (-1)^s \times 2^{-126} \times (1.F)$

Le dernier cas représente ce que nous appelons des nombres sous-normalisés auxquels nous n'avons pas accès mais qui sont utilisés de façon interne pour les calculs.

La table 2.3 présente un sommaire du codage en simple précision.

2.2 Exercices

EXERCICE 2.1 Calculer pour le standard IEEE-754 en simple précision

- (1) le biais;
- (2) les valeurs E_{\min} et E_{\max} ;
- (3) la valeur de la plus grande mantisse;
- (4) le plus grand nombre que nous pouvons représenter ;
- (5) le plus petit nombre positif que nous pouvons représenter ;
- (6) la précision des valeurs numériques représentées par ce standard.

EXERCICE 2.2 Calculer la représentation au standard IEEE-754 simple précision des flottants 5.75 et 0.1. Évaluer l'erreur de représentation.

EXERCICE 2.3 Répéter pour la double précision les calculs de l'ex. 2.1.

EXERCICE 2.4 Répéter pour la double précision les calculs de l'ex. 2.2.

Nom	Signe (s) 1 [31]	Exposant (E) 8 [30-23]	Mantisse (m) 23 [22-0]	Intervalle en hexa	Intervalle	Intervalle décimal
-Nan Quiet	1	11...11	11...11 ⋮ 10...01	FFFFFFF ⋮ FFC00001		
Indéterminé	1	11...11	10...00	FFC00000		
-Nan Signalé	1	11...11	01...11 ⋮ 00...01	FFBFFFF ⋮ FF800001		
-Infinity (Overflow négatif)	1	11...11	00...00	FF800000	$< -(2 - 2^{-23}) \times 2^{127}$	$\leq -3.4028235677973365E + 38$
Normalisé négatif $-1.m \times 2^{(E-127)}$	1	11...10 ⋮ 00...01	11...11 ⋮ 00...00	FF7FFFF ⋮ 80800000	$-(2 - 2^{-23}) \times 2^{127}$ -2^{-126}	$-3.4028234663852886E+38$ $-1.1754943508222875E-38$
Dénormalisé négatif $-0.m \times 2^{-126}$	1	00...00	11...11 ⋮ 00...01	807FFFF ⋮ 80000001	$-(1 - 2^{-23}) \times 2^{-126}$ $-(1 + 2^{-52}) \times 2^{-150}$ $= -2^{-149}/2$	$-1.1754942106924411E-38$ $-7.0064923216240862E-46$ $= -1.4012984643248170E-45$
Underflow négatif	1	00...00	00...00	80000000	2^{-150} ⋮ < -0	< -0
-0	1	00...00	00...00	80000000	-0	-0
0	0	00...00	00...00	00000000	0	0
Underflow positif	0	00...00	00...00	00000000	> 0 ⋮ 2^{-150}	> 0 ⋮ $7.0064923216240861E-46$
Dénormalisé positif $0.m \times 2^{-126}$	0	00...01 ⋮ 00...01	00...01 ⋮ 11...11	00000001 ⋮ 007FFFF	$(1 + 2^{-52}) \times 2^{-150}$ $= 2^{-149}/2$ $(1 - 2^{-23}) \times 2^{-126}$	$7.0064923216240862E-46$ $= 1.4012984643248170E-45$ $1.1754942106924411E-38$
Normalisé positif $1.m \times 2^{E-127}$	0	00...01 ⋮ 11...10	00...00 ⋮ 11...11	00800000 ⋮ 7F7FFFF	2^{-126} ⋮ $> (2 - 2^{-23}) \times 2^{127}$	$1.1754943508222875E-38$ ⋮ $3.4028234663852886E+38$
+Infini (Overflow positif)	0	11...11	00...00	7F800000	$> (2 - 2^{-23}) \times 2^{127}$	$\geq 3.4028235677973365E + 38$
+NaN signalé	0	11...11	00...01 ⋮ 01...11	7F800001 ⋮ 7FBFFFF		
+NaN quiet	0	11...11	10...00 ⋮ 11...11	7FC00000 ⋮ 7FFFFFF		

TABLE 2.3 – Simple précision

EXERCICE 2.5 Considérons un système de représentation binaire pour des nombres réels non négatifs avec trois bits pour l'exposant et trois bits pour la mantisse.

- (1) Présenter les différents exposants.
- (2) évaluer les valeurs binaires et leur équivalent en décimal de toutes les mantisses indépendamment de l'exposant et du bit caché.

- (3) Calculer le plus petit et le plus grand nombre que nous pouvons stocker.
- (4) évaluer la précision de la machine.
- (5) Calculer l'intervalle de la droite des réels qui peut être représenté par ce système et faire un diagramme de la distribution des nombres normalisés flottants de cet intervalle. Commentaires.
- (6) Répétez la même chose pour les nombres sous-normalisés. Commentaires.
- (7) Quel est le nombre de valeurs flottantes différentes que nous pouvons représenter par ce système? Généraliser.
- (8) Comparer, en utilisant cet exemple, le système de représentation en virgule flottante avec celui en virgule fixe et établir les avantages et les inconvénients de chaque système.

2.3 Exercices de laboratoire

Pour résoudre les exercices suivants, il faut utiliser la version 4.4 de Scilab qu'il faut télécharger du site de ScicosLab : <http://www.scicoslab.org>.

EXERCICE 2.6 Écrire un programme en Scilab qui permet de calculer la précision eps de votre ordinateur.

EXERCICE 2.7 Écrire un programme en Scilab qui permet d'évaluer

- (1) le nombre de places p de la mantisse et
- (2) la base de numérotation

de votre ordinateur.

EXERCICE 2.8 (D'APRÈS WALTER GANDER) Voici deux programmes qui calculent le plus petit nombre positif normalisé

```
clear();
format("e",22);

x = 1;
t = x;
while %eps * t/2 > 0
    t = %eps * t /2;
    if (t >0) then
        x = t;
    end;
end;
x

¶
```

```
clear();
format("e",22);
t = 1;
while %eps * t > 0
    ta = t;
    t = t/2;
end;
t = ta
```

Le premier fournit le résultat $x = 8.094771541462983-320$ qui est faux, tandis que le second donne la bonne réponse $t = 2.225073858507201-308$ qui est la réponse correcte.

Expliquer les raisons de ces résultats.

3

PROPAGATION DES ERREURS

3.1	Erreurs lors des opérations arithmétiques	27
3.1.1	Exercice	30
3.2	Formalisation de la notion de l'algorithme	30
3.3	Propagation des Erreurs	31
3.3.1	Exercice	34
3.4	Erreurs directe et inverse	34
3.5	Références	36

3.1 Erreurs lors des opérations arithmétiques

Nous allons examiner la qualité du résultat d'une opération arithmétique c'est-à-dire soit d'une addition, soit d'une soustraction, soit d'une multiplication, soit d'une division ou encore d'une extraction de la racine carrée.

FAIT 3.1 Le résultat d'une opération arithmétique \otimes entre deux valeurs réelles a et b , est entaché d'une erreur qui est due d'une part aux erreurs de représentation des nombres a et b , d'autre part, à l'erreur de représentation du résultat, c'est-à-dire :

$$\text{résultat de l'opération } a \otimes b = f \mid (f \mid (a) \otimes f \mid (b)) \quad (3.1)$$

On peut décomposer cette erreur en deux parties :

- Une erreur d'entrée, due à l'erreur de représentation du résultat de l'opération, notée $\eta^I(a \otimes b)$.
- Une erreur due au calcul et qui sera notée $\eta^C(a \otimes b)$.

L'erreur du calcul est inhérente à la machine et à l'opération et nous ne pouvons pas modifier sa valeur. Par contre l'erreur d'entrée dépend de l'ordre des calculs et, par conséquent, pour chaque série de calculs il y a un ordonnancement de ces calculs qui conduit à une erreur d'entrée minimale.

Les théorèmes suivants précisent les erreurs d'entrée pour les opérations arithmétiques usuelles.

THÉORÈME 3.1.1 (de l'erreur pour les sommes) .- Soit la somme

$$S = x_1 + x_2 + \dots \quad (3.2)$$

L'erreur absolue de l'entrée pour cette somme est

$$\Delta^l S = \Delta^l x_1 + \Delta^l x_2 + \dots \quad (3.3)$$

et l'erreur de précision de l'entrée est

$$\eta^l(S) = \frac{x_1}{S} \eta^l(x_1) + \frac{x_2}{S} \eta^l(x_2) + \dots \quad (3.4)$$

DÉMONSTRATION. D'après (3.2) nous avons

$$m(S) = S + \Delta^l S + \Delta^c S = m(x_1) + m(x_2) + \dots + \Delta^c S = x_1 + \Delta^l x_1 + x_2 + \Delta^l x_2 + \dots + \Delta^c S$$

d'où on obtient (3.3). D'après (1.3) nous avons

$$\Delta x = x \cdot \eta(x)$$

et (3.3) donne :

$$S \cdot \eta^l(S) = x_1 \eta^l(x_1) + x_2 \eta^l(x_2) + \dots$$

d'où on obtient (3.4). ■

En Analyse Numérique on s'intéresse à l'erreur maximale, c-à-d. à la situation qui arrive quand toutes les erreurs se rajoutent. On essaie donc de borner l'erreur d'entrée afin d'obtenir l'erreur maximale. Nous avons dans le cas de la somme :

$$|\eta^l(S)| \leq \left| \frac{x_1}{S} \right| |\eta^l(x_1)| + \left| \frac{x_2}{S} \right| |\eta^l(x_2)| + \dots \leq \frac{1}{|S|} (|x_1| + |x_2| + \dots) \cdot \text{eps}$$

d'où finalement

$$|\eta^l(S)| \leq \frac{|x_i|}{|x_i|} \times \text{eps} \quad (3.5)$$

Cette borne montre que si $|x_i|$ est petit par rapport à $|x_i| \cdot \text{eps}$, l'erreur de précision de l'entrée peut devenir importante.

THÉORÈME 3.1.2 (de l'erreur pour le produit) .- Soit le produit

$$P = x_1 \cdot x_2 \cdot \dots \quad \text{avec } x_i \neq 0 \quad (3.6)$$

L'erreur absolue de l'entrée est

$$\Delta^l P = P \cdot \left(\frac{\Delta^l x_1}{x_1} + \frac{\Delta^l x_2}{x_2} + \dots \right) \quad (3.7)$$

et l'erreur de précision de l'entrée est

$$\eta^l(P) = \eta^l(x_1) + \eta^l(x_2) + \dots \quad (3.8)$$

DÉMONSTRATION. D'après (3.2) nous avons

$$\begin{aligned}
 m(P) &= P + \Delta^I P + \Delta^C P = m(x_1) \cdot m(x_2) \cdot \dots + \Delta^C P \\
 &= (x_1 + \Delta^I x_1) \cdot (x_2 + \Delta^I x_2) \cdot \dots + \Delta^C P \\
 &= x_1 \cdot x_2 \cdot x_3 \cdot \dots + \Delta^I x_1 (x_2 \cdot x_3 \cdot \dots) + \Delta^I x_2 (x_1 \cdot x_3 \cdot \dots) + \dots \\
 &\quad + \Delta^I x_1 \Delta^I x_2 (x_3 \cdot \dots) + \Delta^I x_2 \Delta^I x_3 (x_1 \cdot \dots) + \dots \\
 &\quad + \Delta^I x_1 \cdot \Delta^I x_2 \cdot \Delta^I x_3 \cdot \dots + \Delta^C P
 \end{aligned}$$

Si on néglige les termes qui ont comme facteurs les produits $\Delta^I x_1 \Delta^I x_2, \dots; \Delta^I x_1 \Delta^I x_2 \Delta^I x_3, \dots; \Delta^I x_1 \cdot \Delta^I x_2 \cdot \Delta^I x_3 \cdot \dots$, on obtient

$$\begin{aligned}
 \Delta^I P &= \Delta^I x_1 (x_2 \cdot x_3 \cdot \dots) + \Delta^I x_2 (x_1 \cdot x_3 \cdot \dots) + \dots \\
 &= \frac{\Delta^I x_1}{x_1} (x_1 \cdot x_2 \cdot x_3 \cdot \dots) + \frac{\Delta^I x_2}{x_2} (x_1 \cdot x_2 \cdot x_3 \cdot \dots) + \dots
 \end{aligned}$$

d'où on a (3.7). De cette dernière on a

$$\frac{\Delta^I P}{P} = \frac{\Delta^I x_1}{x_1} + \frac{\Delta^I x_2}{x_2} + \dots$$

qui donne (3.8). ■

De ce théorème nous avons

$$\eta^I(P) \leq N \cdot \epsilon_{ps}$$

où N est le nombre de facteurs dans le produit P .

THÉORÈME 3.1.3 (de l'erreur pour la division).- Soit

$$Q = \frac{a}{b}; b \neq 0$$

L'erreur absolue de l'entrée est

$$\Delta^I Q = \frac{b \Delta^I a - a \Delta^I b}{b^2} \quad (3.9)$$

et l'erreur de précision à l'entrée est

$$\eta^I(Q) = \eta^I(a) - \eta^I(b) \quad (3.10)$$

DÉMONSTRATION. Nous avons

$$m(Q) = Q + \Delta^I Q + \Delta^C Q = \frac{a + \Delta^I a}{b + \Delta^I b} + \Delta^C Q \Rightarrow$$

$$Q + \Delta^I Q \cdot (b + \Delta^I b) = a + \Delta^I a \Rightarrow$$

$$b \Delta^I Q = \Delta^I a - Q \Delta^I b = \Delta^I a - \frac{a}{b} \Delta^I b \Rightarrow$$

$$\Delta^I Q = \frac{b \Delta^I a - a \Delta^I b}{b^2} \Rightarrow \frac{\Delta^I Q}{Q} = \frac{b \Delta^I a - a \Delta^I b}{b^2 \frac{a}{b}}$$

d'où on obtient (3.9) et (3.10). ■

Nous avons de ce théorème que la borne maximale pour l'erreur relative de la division est

$$|\eta^l(Q)| \leq 2 \cdot \text{eps}$$

La table ci-après fournit les erreurs relatives d'entrée pour les opérations élémentaires entre deux nombres.

Opération	Erreur d'entrée
addition $a + b$	$\eta^l(a \pm b) = \frac{a}{a \pm b} \eta^l(a) \pm \frac{b}{a \pm b} \eta^l(b)$
multiplication $a \times b$	$\eta^l(a \times b) = \eta^l(a) + \eta^l(b)$
division a/b	$\eta^l(a/b) = \eta^l(a) - \eta^l(b)$
racine carrée \sqrt{a}	$\eta^l(a) = 0.5 \cdot \eta^l(a)$

TABLE 3.1 – Erreurs des opérations arithmétiques

3.1.1 Exercice

EXERCICE 3.1 Soient trois nombres réels positifs a, b, c avec $a > b > c$. On calcule leur somme de deux façons :

(1) $S = (a + b) + c$

(2) $S = a + (b + c)$

Quel résultat est le meilleur ?

3.2 Formalisation de la notion de l'algorithme

Un algorithme est une transformation $\varphi : D \rightarrow R^m$ avec $D \subseteq R^n$, qui fait passer d'un ensemble des valeurs initiales (données) $x^T = [x_1, \dots, x_n]$ à un ensemble des valeurs finales (résultats) $y^T = [y_1, \dots, y_m]$. La transformation φ est décomposée en une suite des transformations élémentaires $\varphi = \varphi^{(r)} \circ \dots \circ \varphi^{(1)}$, avec $\varphi^{(k)} : D_{k-1} \rightarrow R^{n_k}$ et $D_{k-1} \subseteq R^{n_{k-1}}$, telle que $\varphi^{(k)}(x^{(k-1)}) = x^{(k)}$; $k = 1, \dots, r$.

Si on remplace x par $\text{fl}(x)$, qui est le vecteur composé par les nombres-machine du vecteur x , alors l'application de l'algorithme fournit le résultat $\text{fl}(y) = \text{fl}(\varphi(\text{fl}(x)))$. Nous avons donc

pour le résultat de l'algorithme une erreur absolue $\Delta y = \text{fl}(y) - y$ qui s'exprime à l'aide de la relation :

$$\Delta y_i = m(y_i) - y_i = \varphi_i(\text{fl}(x)) - \varphi_i(x) = \sum_{j=1}^n (m(x_j) - x_j) \frac{\partial \varphi_i(x)}{\partial x_j} = \sum_{j=1}^n \frac{\partial \varphi_i(x)}{\partial x_j} \Delta x_j ; i = 1, \dots, m \quad (3.1)$$

ou, sous forme matricielle :

$$\Delta y = \begin{pmatrix} \Delta y_1 \\ \vdots \\ \Delta y_m \end{pmatrix} = \begin{pmatrix} \frac{\partial \varphi_1}{\partial x_1} & \dots & \frac{\partial \varphi_1}{\partial x_n} \\ \vdots & \dots & \vdots \\ \frac{\partial \varphi_m}{\partial x_1} & \dots & \frac{\partial \varphi_m}{\partial x_n} \end{pmatrix} \begin{pmatrix} \Delta x_1 \\ \vdots \\ \Delta x_n \end{pmatrix} = J[\varphi(x)] \Delta x \quad (3.2)$$

avec $J[\varphi(x)]$ le jacobien de φ . La quantité $\frac{\partial \varphi_i(x)}{\partial x_j}$ représente la sensibilité avec laquelle y_i réagit aux perturbations Δx_j de x_j .

L'erreur relative est donnée par :

$$\eta(y_i) = \sum_{j=1}^n \frac{x_j}{\varphi_i(x)} \cdot \frac{\partial \varphi_i}{\partial x_j} \eta(x_j) ; y_i \neq 0, x_j \neq 0, i = 1, \dots, m \quad (3.3)$$

La valeur $\frac{x_j}{\varphi_i(x)} \cdot \frac{\partial \varphi_i}{\partial x_j}$ indique la manière avec laquelle une erreur relative en x_j affecte l'erreur relative en y_i . Ces valeurs ont l'avantage de ne pas dépendre de la grandeur des valeurs x_j et y_i . On appelle ces valeurs les **facteurs de conditionnement** de l'algorithme. Si un algorithme a, pour quelques facteurs de conditionnement, des valeurs grandes, alors on dit que le problème est mal conditionné. Pour de tels algorithmes des petites erreurs relatives concernant les données en entrée peuvent provoquer des grandes erreurs relatives concernant les résultats.

3.3 Propagation des Erreurs

Nous allons maintenant étudier la propagation des erreurs de calcul lors de l'exécution d'un algorithme.

FAIT 3.2 Le déroulement d'un algorithme se fait selon le schéma suivant :

$$x = x^{(0)} \rightarrow \varphi^{(1)} \quad x^{(0)} = x^{(1)} \rightarrow \dots \rightarrow \varphi^{(r)} \quad x^{(r-1)} = x^{(r)} = y \quad (3.1)$$

Si on s'arrête à la k -ième étape ($1 \leq k \leq r$) nous avons effectué la partie

$$x = x^{(0)} \rightarrow \varphi^{(1)} \quad x^{(0)} = x^{(1)} \rightarrow \dots \rightarrow \varphi^{(k)} \quad x^{(k-1)} = x^{(k)} \quad (3.2)$$

et il nous reste à effectuer la partie :

$$x^{(k)} \rightarrow \varphi^{(k+1)} \quad x^{(k)} = x^{(k+1)} \rightarrow \dots \rightarrow \varphi^{(r)} \quad x^{(r-1)} = x^{(r)} = y \quad (3.3)$$

Notons par $\psi^{(k)}$ l'application qui reste à faire après k étapes, i.e.

$$\psi^{(k)} = \varphi^{(r)} \circ \varphi^{(r-1)} \circ \dots \circ \varphi^{(k+1)} : D_k \rightarrow \mathbb{R}^m ; k = 1, 2, \dots, r-1 \quad (3.4)$$

Notons par $J\varphi^{(k)} \cdot J\psi^{(l)}$ le jacobien de la composition $\varphi^{(k)} \circ \psi^{(l)}$. Nous savons que nous avons :

$$J(f \circ g) = Jf(g(x)) \cdot Jg(x) \quad (3.5)$$

ainsi que

$$\begin{aligned} J\varphi(x) &= J\varphi^{(r)} x^{(r-1)} \cdot J\varphi^{(r-1)} x^{(r-2)} \cdot \dots \cdot J\varphi^{(1)} x^{(r-1)} \\ J\psi^{(k)} x^{(k)} &= J\varphi^{(r)} x^{(r-1)} \cdot J\varphi^{(r-1)} x^{(r-2)} \cdot \dots \cdot J\varphi^{(k+1)} x^{(k)} \end{aligned} \quad k = 0, 1, \dots, r \quad (3.6)$$

Pour étudier la propagation des erreurs lors du déroulement d'un algorithme, examinons ce qui se passe quand on réalise une application élémentaire $\varphi^{(k)}$. Nous avons, à la place de la valeur $x^{(k+1)} = \varphi^{(k+1)} x^{(k)}$, la valeur

$$fl x^{(k+1)} = fl \varphi^{(k+1)} fl x^{(k)} \quad (3.7)$$

L'erreur absolue $\Delta x^{(k+1)} = fl x^{(k+1)} - x^{(k+1)}$ devient donc :

$$\Delta x^{(k+1)} = fl x^{(k+1)} - x^{(k+1)} = fl \varphi^{(k+1)} fl x^{(k)} - \varphi^{(k+1)} x^{(k)} \quad (3.8)$$

qui peut aussi s'écrire sous la forme suivante :

$$\Delta x^{(k+1)} = fl \varphi^{(k+1)} fl x^{(k)} - \varphi^{(k+1)} fl x^{(k)} + \varphi^{(k+1)} fl x^{(k)} - \varphi^{(k+1)} x^{(k)} \quad (3.9)$$

De (3.2) nous avons :

$$\varphi^{(k+1)} fl x^{(k)} - \varphi^{(k+1)} x^{(k)} = J\varphi^{(k+1)} x^{(k)} \Delta x^{(k)} \quad (3.10)$$

La machine ne peut pas forcément coder la valeur $\varphi^{(k+1)} m x^{(k)}$ et nous aurons donc un nombre-machine $m \varphi^{(k+1)} m x^{(k)}$ qui représentera cette valeur. Comme

$$\varphi^{(k+1)} x^{(k)} = \varphi_1^{(k+1)} x^{(k)}, \varphi_2^{(k)} x^{(k)}, \dots, \varphi_{n_{k+1}}^{(k)} x^{(k)} \quad (3.11)$$

nous avons pour la représentation machine les relations :

$$m \varphi_i^{(k+1)} fl x^{(k)} = \varphi_i^{(k+1)} fl x^{(k)} (1 + \eta_i fl x^{(k)}) ; i = 1, \dots, n_k \quad (3.12)$$

où $\eta_i x^{(k)}$ est l'erreur relative d'arrondi, engendrée durant le calcul de la i -ième composante de l'application élémentaire $\varphi^{(k+1)}$ et telle que $\eta_i x^{(k)} \leq \text{eps}$. Sous forme matricielle (3.12) s'écrit :

$$fl \varphi^{(k+1)} fl x^{(k)} = (I + H_{k+1}) \cdot \varphi^{(k+1)} fl x^{(k)} \quad (3.13)$$

avec

$$H_{k+1} = \text{diag } \eta_1 \text{ fl } x^{(k)}, \eta_2 \text{ fl } x^{(k)}, \dots, \eta_{n_k} \text{ fl } x^{(k)} \quad (3.14)$$

Par conséquent le premier terme de la partie droite de (3.9) s'écrit :

$$\text{fl } \varphi^{(k+1)} \text{ fl } x^{(k)} - \varphi^{(k+1)} \text{ fl } x^{(k)} = H_{k+1} \cdot \varphi^{(k+1)} \text{ fl } x^{(k)} \quad (3.15)$$

D'autre part

$$H_{k+1} \cdot \varphi^{(k+1)} \text{ fl } x^{(k)} - H_{k+1} \cdot \varphi^{(k+1)} x^{(k)} = H_{k+1} \cdot x^{(k+1)} \quad (3.16)$$

car les erreurs dues aux différences entre $\text{fl } x^{(k)}$ et $x^{(k)}$ sont multipliées par les erreurs relatives qui sont les éléments de la matrice H_{k+1} et deviennent ainsi négligeables. Par conséquent nous avons :

$$\text{fl } \varphi^{(k+1)} \text{ fl } x^{(k)} - \varphi^{(k+1)} \text{ fl } x^{(k)} = H_{k+1} \cdot x^{(k+1)} \quad (3.17)$$

relation qui permet de calculer H_{k+1} . La quantité $H_{k+1} \cdot x^{(k+1)}$ est l'erreur absolue entre le résultat de la $(k+1)$ -ième application élémentaire et sa représentation machine. Les éléments diagonaux de H_{k+1} représentent les erreurs relatives correspondantes. En utilisant (3.9), (3.10) et (3.17) nous avons :

$$\Delta x^{(k+1)} = J\varphi^{(k+1)} x^{(k)} \cdot \Delta x^{(k)} + H_{k+1} \cdot x^{(k+1)} ; \quad \Delta x^{(0)} = \Delta x \quad (3.18)$$

Pour l'erreur absolue du résultat de l'algorithme on utilise (3.6) et on obtient :

$$\Delta y = J\varphi(x) \cdot \Delta x + J\psi^{(1)}(x^{(1)}) \cdot H_1 \cdot x^{(1)} + \dots + J\psi^{(r-1)}(x^{(r-1)}) \cdot H_{r-1} \cdot x^{(r-1)} + H_r y \quad (3.19)$$

De (3.19) on constate que ce sont les jacobiens $J\psi^{(k)}$ de l'application restante $\psi^{(k)}$ qui sont importants pour l'évaluation des effets des erreurs d'arrondi intermédiaires H_k sur le résultat final, car, pour deux algorithmes différents, $J\varphi(x)$ reste inchangé. Ainsi l'effet total de l'arrondi pour un algorithme donné A est fourni par

$$E_r(A, x) = J\psi^{(1)}(x^{(1)}) \cdot H_1 \cdot x^{(1)} + \dots + J\psi^{(r-1)}(x^{(r-1)}) \cdot H_{r-1} \cdot x^{(r-1)} + H_r y \quad (3.20)$$

FAIT 3.3 Un algorithme A est plus crédible qu'un autre algorithme A si $E_r(A, x) \leq E_r(A, x)$ pour le même ensemble de données x .

D'autre part on peut admettre que $|H_r \cdot y| \leq |y| \cdot \text{eps}$ et, si l'erreur d'entrée des données est $\Delta^1(x)$, nous avons $\Delta^1(x) \leq |x| \cdot \text{eps}$. Par conséquent pour chaque algorithme nous avons une erreur d'entrée

$$\Delta^1 y = [|J\varphi(x)| |x| + |y|] \cdot \text{eps} \quad (3.21)$$

Appelons la quantité $\Delta^1 y$ erreur inhérente du résultat y . Cette erreur est inévitable quand nous utilisons une machine. Si les erreurs intermédiaires d'arrondi $H_k \cdot x^{(k)}$ sont de même ordre de grandeur que l'erreur inhérente, nous disons que l'algorithme φ est numériquement stable.

FAIT 3.4 L'objectif principal de l'Analyse Numérique est d'élaborer des algorithmes qui sont numériquement stables.

3.3.1 Exercice

EXERCICE 3.2 Considérons le calcul

$$y = a^2 - b^2$$

selon deux algorithmes :

- Algorithme 1 : $y \leftarrow a^2 - b^2$
- Algorithme 2 : $y \leftarrow (a + b) \cdot (a - b)$

En utilisant les formules de la propagation des erreurs

- (1) Calculer l'erreur Δy pour chacun de deux algorithmes.
- (2) Trouver s'il y a un algorithme qui est plus crédible que l'autre.
- (3) Étudier la stabilité numérique des algorithmes.

3.4 Erreurs directe et inverse

Nous allons maintenant, pour finir, explorer les relations qui existent entre les nombres-condition et les erreurs lors de l'exécution d'un algorithme.

Supposons qu'en exécutant l'algorithme $y = f(x)$ nous obtenons le résultat $f_l(y)$ qui est une approximation de la valeur de y . La question qui se pose concerne la qualité de $f_l(x)$. D'après Wilkinson, nous pouvons poser cette question de deux façons :

- Mesurer la différence entre y et $f_l(y)$. Il s'agit en réalité de l'erreur que nous commettons quand nous calculons la solution numérique (c'est-à-dire à l'aide d'un ordinateur) du problème. Nous partons de x , nous appliquons f et nous obtenons comme résultat $f_l(y)$ au lieu de y . Il est donc légitime d'appeler cette erreur, *erreur directe*. L'erreur relative de précision $\eta(y)$, donnée par (1.4.3), peut être considérée comme une erreur de la précision relative directe.
- Nous pouvons aussi s'interroger sur les données en entrée $x + \Delta x$ qui ont permis à l'algorithme de nous fournir comme résultat $f_l(y)$, c'est-à-dire on cherche à évaluer $x + \Delta x$ tel que $f_l(y) = f(x + \Delta x)$. C'est une sorte de problème inverse qui n'a pas forcément une solution unique. Nous sommes ici en présence d'une erreur dont la cause se trouve au départ mais dont nous nous rendons compte après le calcul et par un retour en arrière. En effet nous pensons partir de x , nous appliquons f , nous obtenons comme résultat $f_l(y)$ au lieu de y et nous nous demandons, par retour en arrière, quelle était la valeur de départ $x + \Delta x$ que nous avons utilisé, à la place de x , et qui nous a permis d'avoir ce résultat. Cette erreur $|\Delta x|$ (ou, encore, $\frac{|\Delta x|}{|x|}$) sur les données en entrée est appelée *erreur inverse*.

Cette deuxième façon de voir les choses est plus avantageuse que la première car elle permet de considérer que les erreurs d'arrondi ne sont pas des conséquences quasi fortuites du mécanisme de stockage en mémoire des informations numériques par les ordinateurs, mais plutôt

des conséquences des perturbations subies par les données. Ainsi le problème de borner supérieurement l'erreur relative de précision devient un problème de la théorie de perturbations qui dispose d'un arsenal mathématique des techniques et des méthodes très développé.

Faisons l'hypothèse que la fonction $y = f(x)$ est deux fois continûment différentiable. Alors nous avons

$$f(y) - y = f(x + \Delta x) - f(x) = f'(x)\Delta x + \frac{f''(x + \xi\Delta x)}{2!}(\Delta x)^2, \quad \xi \in]0, 1[$$

d'où nous obtenons

$$\frac{f(y) - y}{y} = \frac{xf'(x)}{f(x)} \frac{\Delta x}{x} + \frac{f''(x + \xi\Delta x)}{y \cdot 2!} (\Delta x)^2$$

Comme nous avons déjà vu, le facteur $\frac{xf'(x)}{f(x)}$ du premier terme du second membre est le nombre-condition de f . Nous pouvons donc de cette relation et des définitions des erreurs directe et inverse, obtenir l'inégalité suivante :

FAIT 3.5

$$\text{erreur directe} \leq \text{nombre-condition} \times \text{erreur inverse}$$

Nous avons ainsi comme conclusion qu'une erreur inverse faible ne garantit pas que l'erreur directe soit, à son tour, faible. Il faudrait en plus, que le nombre-condition soit aussi faible, c'est-à-dire que le problème soit bien conditionné.

On dit qu'un algorithme est fiable si

$$\text{erreur inverse} \sim \text{eps}$$

Cette discussion sur les erreurs directe et inverse nous aide à comprendre la notion de l'algorithme numériquement stable que nous avons présenté à la fin de section précédente.

En effet, considérons un réel $x \in \mathbb{R}$ et supposons que nous faisons un calcul $f(x)$ qui a comme résultat $y \in \mathbb{R}$. Si pour le calcul on utilise un ordinateur, le résultat sera $f_l(y) = f(f_l(x))$ que l'on peut mettre sous la forme

$$y(1 + \eta(y)) = f(x(1 + \eta(x)))$$

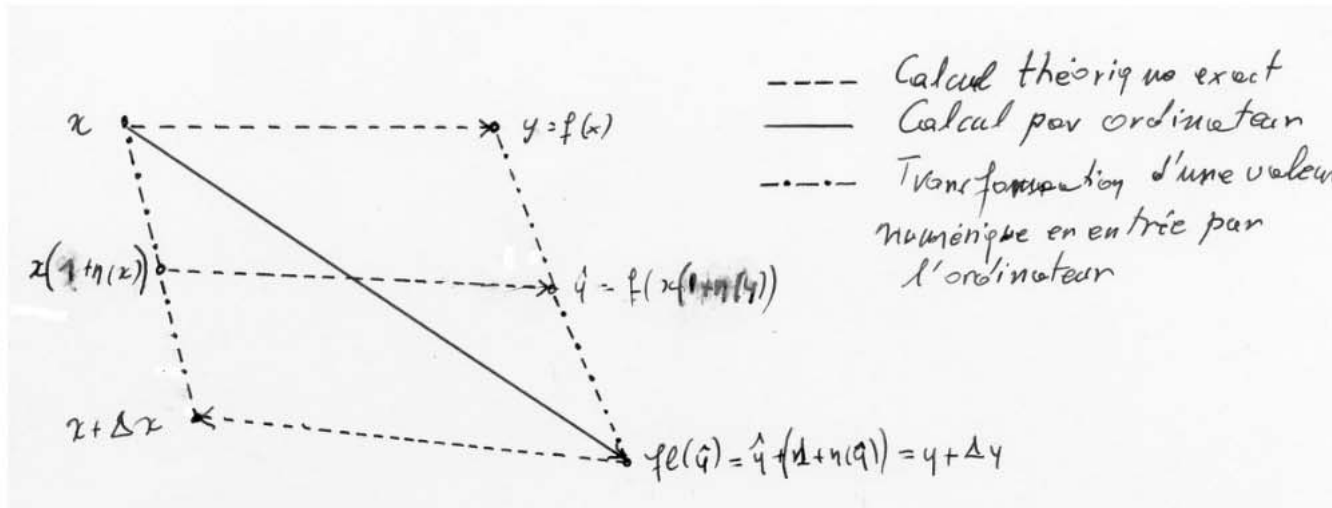
où $\eta(\cdot)$ est l'erreur relative de précision (cf. relation (1.3)).

Cette relation exprime une situation idéale car il n'y a aucune raison, en principe, que le résultat de f sur la donnée particulière $x(1 + \eta(x))$ soit exactement l'erreur d'arrondi $y \cdot \eta(y)$ commise par l'ordinateur sur la valeur numérique y . Donc on a en réalité la relation

$$y + \Delta y = (x + \Delta x)$$

qui exprime un mélange de l'erreur directe et inverse comme on peut le voir sur la figure 3.1.

Sur cette figure Δy représente l'erreur directe et Δx est l'erreur inverse. Ainsi un algorithme numériquement stable signifie que Δx et Δy sont de même ordre de grandeur.

FIGURE 3.1 – Mélange erreur directe et inverse pour $y = f(x)$.

3.5 Références

Les livres et articles utilisés pendant la rédaction de ces notes sur l'Analyse des Erreurs sont les suivants :

- G. ENGELN-MÜLLGES, F. UHLIG : Numerical algorithms with Fortran, Springer, 1996
 - J. STOER, R. BULIRSCH : Introduction to numerical analysis, Second ed., Springer-Verlag, 1992
 - W. PRESS et al. : Numerical recipes in Fortran, Second edition, Cambridge Un.P., 1992
 - J. H. WILKINSON : Rounding errors in algebraic processes, Dover, 1994
 - J. H. WILKINSON : The algebraic eigenvalue problem, Clarendon Press, 1965
 - F. B. HILDEBRAND : Introduction to numerical analysis, Dover, 1987.
 - N. J. HIGHAM : Accuracy and stability of numerical algorithms, Siam, 1996
 - D. E. KNUTH : The art of computer programming, vol. 2, Seminumerical algorithms, second ed., Addison-Wesley, 1981
 - D. GOLDBERG : What every computer scientist should know about floating-point arithmetic, ACM Computing Surveys, v.23, no 1, pp.5-48, 1991
 - W. KAHAN : Lecture Notes on the Status of IEEE Standard 754 for Binary Floating-Point Arithmetic, 1966
- et aussi un site Internet <http://grouper.ieee.org/groups/754/>.

Table des matières

INTRODUCTION	1
1 ANALYSE DES ERREURS	3
1.1 Introduction	3
1.2 L'ordinateur et les nombres	5
1.2.1 Conversion décimale – binaire	6
1.2.2 Élaboration de la forme flottante	7
1.3 Nombres normalisés et spéciaux	9
1.4 Arithmétique arrondie	12
1.4.1 Exercices	14
1.5 Les nombres réels et les nombres-machine	15
1.5.1 Étude de l'erreur de précision relative	15
1.5.2 Exercices	17
2 LE STANDARD IEEE-754	19
2.1 Formats de IEEE-754	20
2.1.1 Format simple précision	21
2.2 Exercices	22
2.3 Exercices de laboratoire	24
3 PROPAGATION DES ERREURS	27
3.1 Erreurs lors des opérations arithmétiques	27
3.1.1 Exercice	30
3.2 Formalisation de la notion de l'algorithme	30
3.3 Propagation des Erreurs	31
3.3.1 Exercice	34
3.4 Erreurs directe et inverse	34
3.5 Références	36