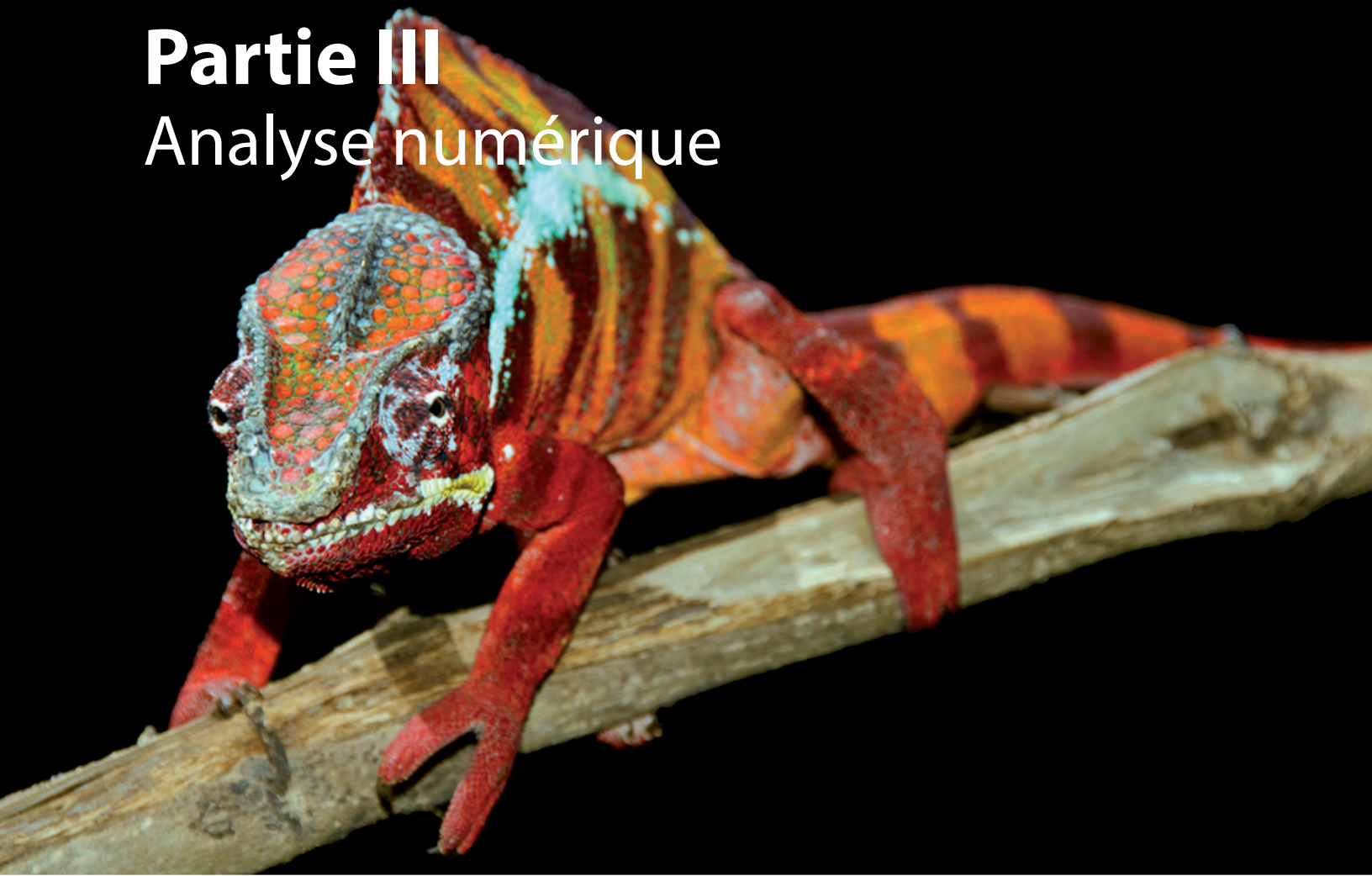


# ACT 2002

## Méthodes numériques en actuariat

### Partie III

#### Analyse numérique



UNIVERSITÉ  
**LAVAL**

Faculté des sciences et de génie  
École d'actuariat



# **ACT 2002**

## Méthodes numériques en actuariat

### **Partie III**

#### Analyse numérique

**Vincent Goulet**

Professeur titulaire | École d'actuariat | Université Laval

© 2013 Vincent Goulet



Cette création est mise à disposition selon le contrat [Paternité-Partage à l'identique 2.5 Canada](#) de Creative Commons. En vertu de ce contrat, vous êtes libre de :

- **partager** — reproduire, distribuer et communiquer l'œuvre ;
- **remixer** — adapter l'œuvre ;
- utiliser cette œuvre à des fins commerciales.

Selon les conditions suivantes :



**Attribution** — Vous devez attribuer l'œuvre de la manière indiquée par l'auteur de l'œuvre ou le titulaire des droits (mais pas d'une manière qui suggérerait qu'ils vous soutiennent ou approuvent votre utilisation de l'œuvre).



**Partage à l'identique** — Si vous modifiez, transformez ou adaptez cette œuvre, vous n'avez le droit de distribuer votre création que sous une licence identique ou similaire à celle-ci.

### Code source

Le code source  $\LaTeX$  de ce document est disponible à l'adresse [https://svn.fsg.ulaval.ca/svn-pub/vgoulet/documents/methodes\\_numeriques/](https://svn.fsg.ulaval.ca/svn-pub/vgoulet/documents/methodes_numeriques/) ou en communiquant directement avec l'auteur.

### Couverture

Le reptile en couverture est un caméléon panthère Ambilobe (*Furcifer pardalis*) originaire de Madagascar.

Crédit photo : Stefan Overmann ; <http://fc-foto.de/25896357>

# Introduction

Les ordinateurs ne savent pas compter ou, en fait, très peu. Ils ne savent traiter que des 0 et des 1 et sont incapables de représenter tous les nombres réels — chose qu’un humain peut faire, du moins conceptuellement. Cela signifie qu’à peu près tout calcul effectué dans un ordinateur comporte une part d’erreur d’arrondi et de troncature. Comme on ne souhaite généralement pas que cette erreur devienne trop grande, il importe de connaître ses sources afin de la diminuer le plus possible. C’est, entre autres choses, l’objet du [chapitre 10](#).

Les procédures numériques pour résoudre des équations à une variable, optimiser une fonction ou calculer une intégrale définie sont aujourd’hui aisément accessibles dans une foule de logiciels à connotation mathématique ou même dans une simple calculatrice. Or, comment ces calculs sont-ils effectués, quels sont les algorithmes à l’œuvre en arrière-scène ? Le [chapitre 11](#) se penche sur les méthodes de base de résolution d’équations à une variable et le [chapitre 12](#) sur celles d’intégration numérique.

On présente également au [chapitre 11](#) les principales fonctions d’optimisation disponibles dans Excel et dans R.

L’étude de ce document implique quelques allers-retours entre le texte et les sections de code informatique présentes dans chaque chapitre. Les sauts vers ces sections sont clairement indiqués dans le texte par des mentions mises en évidence par le symbole ►.

Chaque chapitre comporte des exercices. Les réponses de ceux-ci se trouvent à la fin de chacun des chapitres et les solutions complètes, en annexe.

Je tiens à souligner la précieuse collaboration de MM. Mathieu Boudreault, Sébastien Auclair et Louis-Philippe Pouliot lors de la rédaction des exercices et des solutions. Je remercie également Mmes Marie-Pier Laliberté et Véronique Tardif pour l’infographie des pages couvertures.

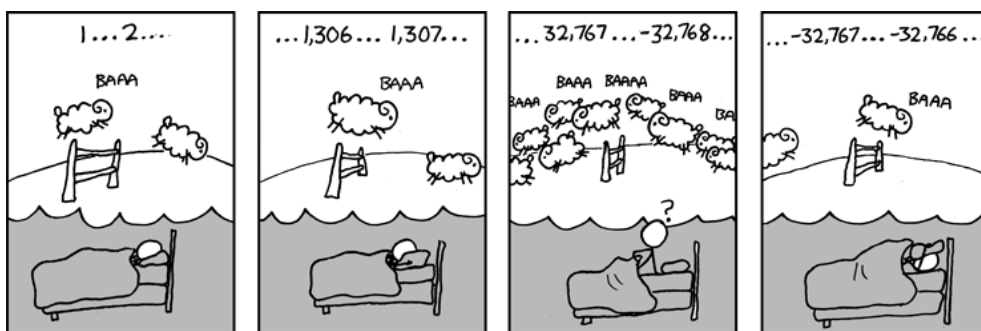


# Table des matières

<b>Introduction</b>	<b>v</b>
<b>10 Arithmétique des ordinateurs</b>	<b>3</b>
10.1 Les ordinateurs ne savent pas compter . . . . .	3
10.2 Conversion de base . . . . .	4
10.3 Unités de mesure . . . . .	11
10.4 Représentation en virgule flottante . . . . .	12
10.5 Éléments d'arithmétique en virgule flottante . . . . .	18
10.6 Codage de caractères . . . . .	23
10.7 Code informatique . . . . .	25
10.8 Exercices . . . . .	26
<b>11 Résolution d'équations à une variable</b>	<b>31</b>
11.1 Mise en contexte . . . . .	31
11.2 Méthode de bisection . . . . .	32
11.3 Méthode du point fixe . . . . .	36
11.4 Méthode de Newton–Raphson . . . . .	40
11.5 Fonctions d'optimisation dans Excel et R . . . . .	49
11.6 Astuce Ripley . . . . .	51
11.7 Outils additionnels . . . . .	52
11.8 Code informatique . . . . .	52
11.9 Exercices . . . . .	61
<b>12 Intégration numérique</b>	<b>65</b>
12.1 Polynômes d'interpolation de Lagrange . . . . .	65
12.2 Principes généraux d'intégration numérique . . . . .	67
12.3 Méthode du point milieu . . . . .	69
12.4 Méthode du trapèze . . . . .	69
12.5 Méthode de Simpson . . . . .	70

12.6	Méthode de Simpson 3/8 . . . . .	70
12.7	Code informatique . . . . .	72
12.8	Exercices . . . . .	73
<b>A</b>	<b>Solutions des exercices</b>	<b>75</b>
	Chapitre 10 . . . . .	75
	Chapitre 11 . . . . .	81
	Chapitre 12 . . . . .	103
	<b>Bibliographie</b>	<b>107</b>





Tiré de XKCD.com



# 10 Arithmétique des ordinateurs

## Objectifs du chapitre

- ▶ Concevoir l'effet de la représentation interne des nombres dans un ordinateur sur la précision des calculs.
- ▶ Convertir un nombre décimal vers une base quelconque, et vice versa.
- ▶ Exprimer un nombre décimal dans la représentation interne d'un ordinateur selon la norme IEEE 754.
- ▶ Planifier ses calculs avec un ordinateur de manière à diminuer l'impact des erreurs d'arrondis et de troncature et optimiser le temps de calcul.

## 10.1 Les ordinateurs ne savent pas compter

Le type de bogue le plus fréquemment rapporté dans les forums de discussion de R a trait au fait que le logiciel ne retourne pas le bon résultat lors d'opérations arithmétiques simples.

- ▶ Exécuter le code informatique de la [section 10.7](#) pour voir quelques exemples de calculs simples que R n'arrive apparemment pas à effectuer correctement. Tous les exemples sont tirés de messages envoyés à la liste de discussion `r-help`<sup>1</sup>.

On devine que les créateurs de R n'ont pas négligé une fonctionnalité aussi fondamentale pour un langage mathématique que le calcul arithmétique. Les supposées erreurs ci-dessus relèvent toutes de la représentation interne des nombres dans un ordinateur et d'erreurs d'arrondi inhérentes aux opérations arithmétiques

---

1. <https://stat.ethz.ch/mailman/listinfo/r-help/>

avec ces nombres. En effet, un des plus célèbres principes de programmation proposés par Kernighan et Plauger (1978) est :

10,0 fois 0,1 ne donne jamais vraiment 1,0.

D'ailleurs, les auteurs des faux rapports de bogues dans `r-help` sont invariablement renvoyés à l'entrée 7.31 de la foire aux questions de R<sup>2</sup>.

Pour quiconque travaille régulièrement avec un ordinateur pour faire du calcul scientifique, connaître globalement le fonctionnement interne d'un ordinateur peut s'avérer précieux pour les raisons suivantes, entre autres :

- ▶ éviter les erreurs d'arrondi et de troncature ;
- ▶ éviter les dépassements et sous-dépassements de capacité (*overflow* et *underflow*) ;
- ▶ optimiser le code informatique.

Les ordinateurs ne savent pas faire d'arithmétique à proprement parler. Ils ne connaissent que deux états : ouvert (1) et fermé (0). Nous en profiterons donc d'abord pour réviser la conversion des nombres décimaux vers, et de, n'importe quelle base. Après avoir, à la section 10.3, introduit les principales unités de mesure en informatique, nous expliquerons la représentation interne des nombres dans un ordinateur en nous concentrant sur la double précision. Nous pourrions ainsi justifier que les ordinateurs ne peuvent représenter qu'un sous-ensemble des nombres réels, d'où les erreurs d'arrondi et de troncature. La section 10.5 explique comment ces erreurs surviennent. On y donne également quelques pistes pour éviter ou pour diminuer l'impact de ces erreurs. Enfin, le chapitre se clôt par un survol d'un sujet quelque peu périphérique : la représentation interne des caractères.

## 10.2 Conversion de base

La *base*, dans un système de numération, est le nombre de symboles (habituellement les chiffres) qui pourront servir à exprimer des nombres. L'humain s'est habitué à compter en base 10, le système *décimal*, où les symboles sont 0, 1, ..., 9. De nos jours, la grande majorité des ordinateurs travaillent toutefois en base 2, le système *binaire*. Les autres systèmes d'usage courant, principalement en informatique, sont l'*octal* (base 8) et l'*hexadécimal* (base 16). Dans ce dernier système, les symboles utilisés pour les nombres 10–15 sont généralement les lettres A–F. (C'est pourquoi l'on retrouve ces six lettres sur le pavé des calculatrices scientifiques).

Cette section étudie la conversion des nombres entre la base 10 et une autre base quelconque.

---

2. <http://cran.r-project.org/doc/FAQ/R-FAQ.html>

### 10.2.1 Notation et définitions

De manière générale, soit  $x$  un nombre (entier pour le moment) dans la base de numération  $b$  composé de  $m$  chiffres ou symboles, c'est-à-dire

$$x = x_{m-1}x_{m-2} \cdots x_1x_0,$$

où  $0 \leq x_i \leq b-1$ . On a donc

$$x = \sum_{i=0}^{m-1} x_i b^i. \quad (10.1)$$

Lorsque le contexte ne permet pas de déterminer avec certitude la base d'un nombre, celle-ci est identifiée en indice du nombre par un nombre décimal. Par exemple,  $10011_2$  est le nombre binaire 10011.

**Exemple 10.1.** Soit le nombre décimal 348. Selon la notation ci-dessus, on a  $x_0 = 8$ ,  $x_1 = 4$ ,  $x_2 = 3$  et  $b = 10$ . En effet,

$$348 = 3 \times 10^2 + 4 \times 10^1 + 8 \times 10^0.$$

Ce nombre a les représentations suivantes dans d'autres bases. En binaire :

$$\begin{aligned} 101011100_2 &= 1 \times 2^8 + 0 \times 2^7 + 1 \times 2^6 \\ &\quad + 0 \times 2^5 + 1 \times 2^4 + 1 \times 2^3 \\ &\quad + 1 \times 2^2 + 0 \times 2^1 + 0 \times 2^0. \end{aligned}$$

En octal :

$$534_8 = 5 \times 8^2 + 3 \times 8^1 + 4 \times 8^0.$$

En hexadécimal :

$$15C_{16} = 1 \times 16^2 + 5 \times 16^1 + 12 \times 16^0.$$

Des représentations ci-dessus, l'hexadécimale est la plus compacte : elle permet de représenter avec un seul symbole un nombre binaire comptant jusqu'à quatre chiffres. C'est, entre autres, pourquoi c'est une représentation populaire en informatique.  $\square$

Dans un ordinateur réel (par opposition à théorique), l'espace disponible pour stocker un nombre est fini, c'est-à-dire que  $m < \infty$ . Le plus grand nombre que l'on

peut représenter avec  $m$  chiffres ou symboles en base  $b$  est

$$\begin{aligned} x_{\max} &= \sum_{i=0}^{m-1} (b-1)b^i \\ &= (b-1) \sum_{i=0}^{m-1} b^i \\ &= (b-1) \left( \frac{b^m - 1}{b - 1} \right) \\ &= b^m - 1. \end{aligned}$$

Par exemple, le plus grand nombre décimal représentable avec  $m = 4$  symboles est

$$x_{\max} = 9999 = 10\,000 - 1 = 10^4 - 1,$$

alors que le plus grand nombre binaire est seulement

$$x_{\max} = 1111 = 10\,000 - 1 = 2^4 - 1 = 15_{10}.$$

Par une extension naturelle de ce qui précède, un nombre composé de  $m \geq 0$  symboles dans sa partie entière et  $n \geq 0$  symboles dans sa partie fractionnaire est représenté en base  $b$  comme

$$\begin{aligned} x &= x_{m-1}x_{m-2} \cdots x_1x_0, x_{-1}x_{-2} \cdots x_{-n} \\ &= \sum_{i=-n}^{m-1} x_i b^i. \end{aligned}$$

Le symbole qui sépare les parties entière et fractionnaire du nombre est la *séparation fractionnaire* (une virgule en français, un point en anglais).

### 10.2.2 Conversion vers une base quelconque

Avant de discuter de la conversion de nombres décimaux vers une base quelconque, il convient de définir les notions de *quotient* et de *reste* d'une division.

Le quotient est la partie entière de la division de deux entiers  $a$  et  $d$  ; sa représentation mathématique habituelle est

$$q = \left\lfloor \frac{a}{d} \right\rfloor, \quad (10.2)$$

où  $\lfloor x \rfloor$  est la fonction qui retourne le plus grand entier inférieur ou égal à  $x$ . Le reste de la division est simplement la valeur

$$r = a - d \left\lfloor \frac{a}{d} \right\rfloor. \quad (10.3)$$

Évidemment, on a  $r \in \{0, 1, \dots, d-1\}$ . Le reste est le résultat de l'opération modulo, notée  $r = a \bmod d$ .

On remarquera que le premier chiffre en partant de la droite d'un entier décimal est le reste de la division de ce nombre par 10, que le second chiffre est le reste de la division par 10 du quotient de la division précédente, et ainsi de suite. La conversion d'un nombre décimal en une base  $b$  implique simplement de diviser par  $b$  plutôt que par 10 et de déterminer le symbole dans la base  $b$  correspondant au reste de chaque division.

L'algorithme suivant reprend ses idées sous forme plus formelle et en ajoutant le traitement de la partie fractionnaire. Nous démontrerons plus loin qu'il n'est pas réellement nécessaire de savoir effectuer la conversion de la partie fractionnaire d'un nombre réel.

**Algorithme 10.1** (Conversion de la base 10 vers la base  $b$ ). *Soit  $x$  un nombre réel en base 10.*

1. Poser  $i \leftarrow 0$  et  $v \leftarrow \lfloor x \rfloor$ .
2. Répéter les étapes suivantes jusqu'à ce que  $v = 0$  :
  - a) Poser  $d_i \leftarrow v \bmod b$  et trouver  $x_i$ , le symbole dans la base  $b$  correspondant à  $d_i$  ;
  - b) Poser  $v \leftarrow \lfloor v/b \rfloor$  ;
  - c) Poser  $i \leftarrow i + 1$ .
3. Poser  $i \leftarrow i + 1$  et  $v \leftarrow x - \lfloor x \rfloor$ .
4. Répéter les étapes suivantes jusqu'à ce que  $v = 0$  ou que  $i = n$  (le nombre voulu ou maximal de chiffres après la séparation fractionnaire) :
  - a) Poser  $d_{-i} \leftarrow \lfloor bv \rfloor$  et trouver  $x_{-i}$ , le symbole dans la base  $b$  correspondant à  $d_{-i}$  ;
  - b) Poser  $v \leftarrow bv - d_{-i}$  ;
  - c) Poser  $i \leftarrow i + 1$ .
5. Retourner

$$x_b = x_{m-1}x_{m-2} \cdots x_1x_0, x_{-1}x_{-2} \cdots x_{-n}.$$

**Exemple 10.2.** Soit le nombre décimal 23,31 que l'on convertit en binaire (base 2) avec un maximum de cinq chiffres après la séparation fractionnaire ( $n = 5$ ). Le [tableau 10.1](#) montre le processus en suivant les étapes de l'[algorithme 10.1](#). On obtient le résultat final en combinant la dernière colonne du [10.1](#) lue de bas en haut avec la dernière colonne du [10.1](#) lue de haut en bas. Ainsi, la représentation binaire du 23,31 limitée à cinq chiffres après la virgule est 10111,01001.

□

TAB. 10.1: Conversion du nombre décimal 23,31 en binaire.

(a) partie entière					(b) partie fractionnaire				
$i$	$v$	$\lfloor v/2 \rfloor$	$v \bmod 2$	$x_i$	$i$	$v$	$2v$	$\lfloor 2v \rfloor$	$x_{-i}$
0	23	11	1	1	1	0,31	0,62	0	0
1	11	5	1	1	2	0,62	1,24	1	1
2	5	2	1	1	3	0,24	0,48	0	0
3	2	1	0	0	4	0,48	0,96	0	0
4	1	0	1	1	5	0,96	1,92	1	1



Tel que mentionné ci-dessus, il est tout à fait possible de se passer de la conversion de la partie fractionnaire. Pour convertir le nombre réel  $x$  en base  $b$  avec un maximum de  $n$  chiffres dans la partie fractionnaire, il suffit de :

1. multiplier  $x$  par  $b^n$  ;
2. convertir l'entier  $\lfloor xb^n \rfloor$  en base  $b$  ;
3. déplacer la virgule de  $n$  positions vers la gauche.

**Exemple 10.3.** Soit de nouveau la conversion de 23,31 en binaire avec une partie fractionnaire d'au plus cinq chiffres. En suivant les étapes de la remarque ci-dessus, on a :

1.  $23,31 \times 2^5 = 23,31 \times 32 = 745,92$  ;
2. la représentation binaire de  $745 = 2^9 + 2^7 + 2^6 + 2^5 + 2^3 + 1$  est 1011101001 ;
3. en déplaçant la virgule de cinq positions vers la gauche, on obtient 10111,01001, comme à l'exemple 10.2.

□

On tire une conclusion intéressante de la procédure ci-dessus. S'il existe un entier  $k \leq n$  tel que la partie fractionnaire de  $x$  est un multiple de  $1/b^k$ , alors la multiplication de  $x$  par  $b^k$  résultera en un entier. Par conséquent, la représentation de  $x$  en base  $b$  aura une partie fractionnaire de longueur finie. Ainsi, pour obtenir une représentation binaire finie, la partie fractionnaire d'un nombre réel doit être un multiple de  $\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \dots$ . Ce n'était pas le cas à l'exemple 10.3 (la partie fractionnaire étant 0,31), d'où la nécessité de limiter le nombre de chiffres après la virgule.



### 10.2.3 Conversion en décimal

La conversion d'un nombre en base  $b$  vers la base 10 repose essentiellement sur la définition (10.1), mais avec chaque symbole  $x_i$  remplacé pour son équivalent décimal. Un algorithme de conversion est le suivant.

**Algorithme 10.2** (Conversion de la base  $b$  vers la base 10). Soit  $x$

$$x_b = x_{m-1}x_{m-2} \cdots x_1x_0, x_{-1}x_{-2} \cdots x_{-n}$$

un nombre réel en base  $b$ .

1. Poser  $x \leftarrow 0$  et  $y \leftarrow 0$ .
2. Pour  $i = m-1, m-2, \dots, 0$ , faire les étapes suivantes :
  - a) Trouver  $d_i$ , le nombre décimal correspondant au symbole  $x_i$  ;
  - b) Poser  $x \leftarrow xb + d_i$ .
3. Pour  $i = -n, -n+1, \dots, -1$ , faire les étapes suivantes :
  - a) Trouver  $d_i$ , le nombre décimal correspondant au symbole  $x_i$  ;
  - b) Poser  $y \leftarrow (y + d_i)/b$ .
4. Retourner  $x + y$ .

On peut, ici aussi, aisément éviter de devoir convertir la partie fractionnaire. Il suffit de déplacer la virgule de  $n$  positions vers la droite dans le nombre  $x_b$  de manière à obtenir un entier, de convertir ce nombre en décimal avec les étapes 1 et 2 de l'algorithme 10.2 et, finalement, de diviser le nombre obtenu par  $b^n$ .

**Exemple 10.4.** Soit le nombre binaire 101011,011 que l'on convertit en décimal. Le tableau 10.2 illustre le processus de conversion selon l'algorithme 10.2 (le chiffre en surbrillance est celui traité lors de chaque étape de l'algorithme). Le résultat final est 43,375.

De manière équivalente, mais plus simple, on peut déplacer la virgule de trois positions vers la droite pour obtenir l'entier binaire 101011011. Ce nombre est  $2^8 + 2^6 + 2^4 + 2^3 + 2 + 1 = 347$  en décimal. Enfin,  $347 \div 2^3 = 43,375$ .  $\square$

**Exemple 10.5.** Soit le nombre hexadécimal AC2,3D8. On peut le convertir en base 10 à partir de la définition :

$$\begin{aligned} \text{AC2,3D8} &= 10 \times 16^2 + 12 \times 16 + 2 + 3 \times 16^{-1} + 13 \times 16^{-2} + 8 \times 16^{-3} \\ &= 2754,240234375. \end{aligned}$$

$\square$

TAB. 10.2: Conversion d'un nombre binaire en décimal

(a) partie entière			(b) partie fractionnaire		
$i$	$d_i$	$x$	$i$	$d_{-i}$	$y$
5	101011,011	1	3	101011,011	0,5
4	101011,011	2	2	101011,011	0,75
3	101011,011	5	1	101011,011	0,375
2	101011,011	10			
1	101011,011	21			
0	101011,011	43			

#### 10.2.4 Conversion avec des bases générales

Il existe de nombreuses bases d'usage courant où chaque symbole est (potentiellement) exprimé dans une base différente. On a qu'à penser à l'heure ou l'ensemble du système de mesure impérial, par exemple. Or, il peut s'avérer utile de savoir convertir de et vers une base quelconque. Le matériel de cette section se retrouve dans très peu d'ouvrages d'analyse numérique ou de statistique numérique. Pourtant, il existe quelques applications intéressantes de la conversion de et vers des bases générales, comme nous le verrons.

On peut généraliser la notion de nombre présentée dans les sections précédentes à une collection de «symboles», chacun dans une base différente. On restreint la discussion aux entiers sans perte de généralité. On a donc

$$x = x_{m-1}x_{m-2} \cdots x_1x_0,$$

où  $x_{m-1}$  est un nombre en base  $b_{m-1}$ ,  $x_{m-2}$  est un nombre en base  $b_{m-2}$ , etc. Nous dirons que le nombre  $x$  est exprimé en base  $[b_{m-1} b_{m-2} \dots b_0]$ . On a alors

$$x = \sum_{i=0}^{m-1} x_i \prod_{j=-1}^{i-1} b_j, \quad (10.4)$$

avec  $b_{-1} = 1$ .

Les algorithmes de conversion 10.1 et 10.2 demeurent essentiellement valides ici. Il suffit de remplacer chaque mention de  $b$  par  $b_i$ .

**Exemple 10.6.** Le jour de l'année et l'heure du jour est un «nombre» exprimé en base  $[365 \ 24 \ 60 \ 60]$ . En utilisant directement l'équation (10.4), le nombre de secondes correspondant à 1 jour, 2 heures, 33 minutes et 20 secondes est

$$\begin{aligned} 1 \text{ j } 2 \text{ h } 33 \text{ min } 20 \text{ sec} &= 1 \times (24)(60)(60) + 2 \times (60)(60) + 33 \times 60 + 20 \\ &= 95,600 \text{ secondes.} \end{aligned}$$

TAB. 10.3: Conversion d'un nombre décimal dans la base générale [365 24 60 60].

$i$	$v$	$b_i$	$\lfloor v/b_i \rfloor$	$v \bmod b_i$	$x_i$
0	91 492	60	1 524	52	52
1	1 524	60	25	24	24
2	25	24	1	1	1
3	1	365	0	1	1

À l'inverse, le nombre de jours, heures, minutes et secondes correspondant à 91 492 secondes est, en utilisant l'algorithme 10.1 avec la base [365 24 60 60] : 1 jour, 1 heure, 24 minutes et 52 secondes (voir le tableau 10.3 pour les calculs).  $\square$

Les bases générales sont particulièrement utiles pour assigner ou extraire les éléments d'une matrice ou d'un tableau dans un ordre non séquentiel. Typiquement, l'index de l'élément à traiter est le résultat d'un calcul. L'exemple suivant illustre cette idée.

**Exemple 10.7.** Soit  $A$  une matrice  $4 \times 5$  que l'on suppose remplie en ordre lexicographique (par ligne). Il est simple de déterminer, ici, que le 14<sup>e</sup> élément est  $a_{34}$ . Or, on observe que la conversion du nombre  $14 - 1 = 13$  dans la base [4 5] donne :

$$13 \div 5 = 2 \text{ reste } 3 \Rightarrow x_0 = 3$$

$$2 \div 4 = 0 \text{ reste } 2 \Rightarrow x_1 = 2,$$

soit le «nombre» (2,3). En additionnant 1 à ce résultat, on obtient précisément la position du 14<sup>e</sup> élément dans la matrice. Les opérations  $-1$  et  $+1$  sont rendues nécessaires par le fait que les lignes et les colonnes sont numérotées à partir de 1, alors que les systèmes de numération débutent à 0.  $\square$

► La fonction `arrayInd` de R permet de faire exactement ce qui est décrit dans l'exemple précédent, soit retourner les coordonnées du  $i^{\text{e}}$  élément d'un tableau. Voir le code informatique de la section 10.7 pour quelques exemples d'utilisation.

## 10.3 Unités de mesure

Les ordinateurs que nous utilisons couramment fonctionnent en base 2. La plus petite quantité d'information qu'un ordinateur peut traiter est un *bit* (compression

de l'anglais *binary digit*). En informatique, un bit est égal à  $\boxed{0}$  ou à  $\boxed{1}$ . Le symbole du bit est b.

Un *octet* (*byte*) est un groupe de huit bits. Son symbole est o ou B. L'octet est l'unité de mesure la plus fréquemment utilisée en informatique, en grande partie parce que le codage d'un caractère dans la plupart des langues occidentales requiert un octet ; voir la [section 10.6](#). Les termes pour de plus grandes quantités de bits ou d'octets utilisent généralement les préfixes usuels du système international d'unités. Par exemple :

- 1 kilooctet (ko) est  $2^{10} = 1,024$  octets ;
- 1 mégaoctet (Mo) est  $2^{20} = 1,048,576$  octets ;
- 1 gigaoctet (Go) est  $2^{30} = 1,073,741,824$  octets.

On voit que cette pratique fort répandue entraîne des distorsions par rapport aux définitions usuelles de kilo, méga, giga, etc., et que cette distorsion s'amplifie au fur et à mesure que l'on grimpe dans l'échelle des tailles d'objets. Pour cette raison, on a défini les préfixes kibi, mébi, gibi, etc., mais ceux-ci demeurent peu utilisés à ce jour.

Dans l'industrie de l'informatique, on joue beaucoup avec les unités pour montrer son produit sous un jour favorable. Deux exemples :

1. Les fournisseurs d'accès Internet expriment la vitesse de téléchargement en Mb/sec, alors que la taille des fichiers est généralement affichée en octets. Il faut donc diviser par 8 la vitesse annoncée pour avoir une idée plus juste du temps requis pour télécharger un fichier.
2. Les fabricants de disques durs divisent la capacité réelle de leurs disques par  $10^9$  pour l'exprimer en gigaoctets. Ainsi, un disque dont la capacité annoncée est de 100 Go ne peut contenir, en fait, que  $100 \times 10^9 \div 2^{30} = 93,132$  Go de données. Cette confusion serait éliminée si les fabricants affichaient plutôt une capacité de 93,132 gibioctets. Moins vendeur...

## 10.4 Représentation en virgule flottante

Cette section décrit la manière standard de représenter les nombres réels dans les ordinateurs d'usage courant. Il est nécessaire de connaître les grandes lignes afin de comprendre pourquoi les nombres réels n'ont pas tous une représentation exacte dans un ordinateur et comment surviennent et se propagent les erreurs d'arrondi. L'étude du sujet est également intéressante en soi pour constater comment les ingénieurs et les informaticiens sont parvenus à stocker un maximum d'information dans un espace limité.

Tel que mentionné précédemment, la capacité de stockage d'un ordinateur, bien que vaste de nos jours, n'en demeure pas moins limitée. Par conséquent, les nombres  $y$  sont représentés par un ensemble de  $m$  bits. Habituellement,  $m$  est un multiple de 2 et sa valeur détermine le type de nombre. Par exemple, des définitions usuelles sont  $m = 2^4 = 16$  pour un entier,  $m = 2^5 = 32$  pour un nombre réel en simple précision (*float* dans plusieurs langages de programmation) et  $m = 2^6 = 64$  pour un nombre réel en double précision (*double*).

Il existe deux grandes façons de représenter les nombres (réels) à l'aide de  $m$  bits.

**Virgule fixe** Dans cette représentation, la position de la virgule dans le nombre est prédéterminée et, par conséquent, n'a pas besoin d'être conservée en mémoire. On réserve alors  $m_e$  positions pour la partie entière et  $m_f$  pour la partie fractionnaire (avec  $m = m_e + m_f$ ). La représentation en virgule fixe est particulièrement utile dans les applications financières.

**Virgule flottante** Dans la représentation des nombres en virgule flottante, celle-ci peut être placée n'importe où dans le nombre. L'étendue de cette représentation est beaucoup plus grande que la virgule fixe, mais ceci au détriment de la précision puisque quelques bits doivent être réservés pour stocker la position de la virgule dans le nombre.

La représentation en virgule flottante est celle utilisée dans les applications scientifiques et celle sur laquelle nous nous concentrons dans la suite.

La représentation en virgule flottante est analogue à la notation scientifique. Le nombre réel  $x$  est entièrement défini par un bit de signe  $S$ , un exposant positif  $E$  et une mantisse  $M$  tel que

$$x = (-1)^S \times B^{E-e} \times M, \quad (10.5)$$

où  $B$  est la base de la représentation et  $e$  est un entier prédéterminé appelé le décalage, ou le biais, de l'exposant. Le recours au décalage facilite les calculs internes et évite de devoir réserver un autre bit pour le signe de l'exposant en le stockant sous forme non signée (c'est-à-dire sous forme d'entier positif).

La norme IEEE 754 définit la manière standard de représenter les nombres en virgule flottante dans les ordinateurs (IEEE, 2003; Wikipedia, 2012, pour une excellente présentation). En premier lieu, la norme stipule que la base dans la représentation est  $B = 2$ . Ceci est implicite et n'est pas stocké où que ce soit dans l'ordinateur. En second lieu, la norme suppose que, pour les nombres dits *normalisés*, la mantisse est toujours de la forme  $M = 1, F$ , où  $F$  est un entier binaire. Le bit à gauche de la virgule est appelé le *bit caché* puisqu'il est lui aussi implicite et non stocké dans l'ordinateur.

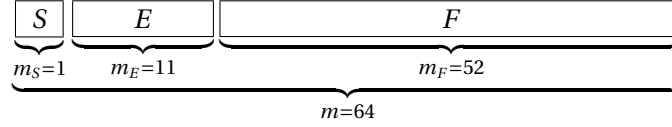


FIG. 10.1: Représentation schématique d'un nombre en double précision dans la norme IEEE 754.

TAB. 10.4: Valeurs spéciales dans la norme IEEE 754

S	E	F	Représentation binaire		Valeur spéciale
0	0	0	0	000000000000 00 ... 00	0
1	0	0	1	000000000000 00 ... 00	-0
0	2047	0	0	111111111111 00 ... 00	$+\infty$
1	2047	0	1	111111111111 00 ... 00	$-\infty$
0 ou 1	2047	$F \neq 0$	S	111111111111 F	NaN <sup>a</sup>

<sup>a</sup> Not a number, par exemple  $\frac{0}{0}$  ou  $\frac{\infty}{\infty}$ .

Nous décrivons plus en détail la norme pour les nombres réels en *double précision* puisque c'est le type avec lequel R travaille toujours<sup>3</sup>. Tout d'abord, un nombre en double précision est stocké dans un mot de  $m = 64$  bits (8 octets) divisé ainsi de gauche à droite :  $m_S = 1$  bit pour le signe,  $m_E = 11$  bits pour l'exposant et  $m_F = 52$  bits pour la partie fractionnaire de la mantisse. Le premier bit du mot utilisé pour le signe est appelé le *bit fort*. Voir la [figure 10.1](#) pour une représentation schématique.

On remarque que le bit caché confère à la mantisse une longueur de mot effective (ou précision) de 53 bits.

Le décalage est  $e = 2^{m_E-1} - 1 = 2^{10} - 1 = 1023$ . Avec une longueur de mot de 11 bits, les valeurs possibles de  $E$  vont de 0 à  $2^{11} - 1 = 2047$ . Cependant, les valeurs 0 et 2047 sont réservées pour des usages spéciaux ; voir le [tableau 10.4](#). Par conséquent, les valeurs possibles de l'exposant (une fois décalé) pour les nombres en double précision vont de  $-1022$  à  $1023$ .

**Exemple 10.8.** Tel qu'expliqué à la [sous-section 10.2.2](#), le nombre  $\frac{1}{2}$  possède une

3. Il est possible de créer des vrais entiers dans R en ajoutant un suffixe L à un nombre entier. Ainsi, 1L est un nombre entier dans le sens où `is.integer(1L)` est TRUE. Cela est relativement peu utilisé en programmation normale

représentation binaire exacte :

$$\begin{aligned}\frac{1}{2} &= 2^{-1} \\ &= (-1)^0 \times 2^{1022-1023} \times 1,0.\end{aligned}$$

Puisque la représentation de l'exposant est  $1022_{10} = 111111110_2$ , la représentation interne de  $\frac{1}{2}$  selon la norme IEEE 754 est

$$\boxed{0} \boxed{0111111110} \boxed{00 \cdots 00}.$$

En revanche, le nombre  $\frac{1}{10}$  n'a pas une représentation binaire finie. En utilisant l'algorithme 10.1, on trouve (le côté droit de l'égalité étant en binaire) :

$$\begin{aligned}\frac{1}{10} &= 0,0001100110011001100110011001 \dots \\ &= 2^{-4} \times 1,1001100110011001 \dots \\ &= (-1)^0 \times 2^{1019-1023} \times 1,1001100110011001 \dots\end{aligned}$$

Or,  $1019_{10} = 111111011_2$ , d'où la représentation interne finie de  $\frac{1}{10}$  est

$$\boxed{0} \boxed{0111111011} \boxed{1001 \cdots 1001}.$$

Si l'on reconvertit ce nombre en décimal, on obtient

$$2^{-4} \times \left( 1 + \sum_{k=0}^{12} (2^{-4k-1} + 2^{-4k-4}) \right),$$

un nombre près de, mais pas exactement égal à,  $\frac{1}{10}$  (le programme de calcul en multiprécision bc donne 0,0999999999999999166). Dès lors, la multiplication de ce nombre par 10 ne donnera pas 1. Voilà qui justifie le principe de programmation énoncé à la page 4.  $\square$

Les principales caractéristiques des nombres en double précision sont les suivantes.

1. Soit  $x_{\max}$  le plus grand nombre représentable. Parce que l'exposant  $E = 2047$  est réservé, on a

$$\begin{aligned}x_{\max} &= \boxed{0} \boxed{1111111110} \boxed{11 \cdots 11} \\ &= (-1)^0 \times 2^{1023} \times 1,11 \cdots 11 \\ &= 2^{1023} \times (2 - 2^{-52}) \\ &= 1,797693 \times 10^{308}.\end{aligned}$$

2. Soit  $x_{\min}$  le plus petit nombre normalisé représentable. Parce que l'exposant  $E = 0$  est réservé, on a

$$\begin{aligned} x_{\min} &= \boxed{0} \boxed{00000000001} \boxed{00 \cdots 00} \\ &= (-1)^0 \times 2^{-1022} \times 1,00 \cdots 00 \\ &= 2^{-1022} \\ &= 2,225\,074 \times 10^{-308}. \end{aligned}$$

Afin de rendre la transition vers 0 moins abrupte, la norme IEEE 754 définit également des nombres *dénormalisés*. Ceux-ci sont identifiés par  $E = 0$  et une partie fractionnaire  $F$  non nulle. Cependant, par définition, les nombres dénormalisés ont  $E - e = -1022$  et leur bit caché est 0. Par conséquent, le plus petit nombre dénormalisé est stocké sous la forme

$$\boxed{0} \boxed{00000000000} \boxed{00 \cdots 01}$$

et sa valeur est

$$\begin{aligned} x_{\min} &= (-1)^0 \times 2^{-1022} \times 0,00 \cdots 01 \\ &= 2^{-1022} \times 2^{-52} \\ &= 2^{-1074} \\ &= 4,940\,656 \times 10^{-324}. \end{aligned}$$

3. On a la même étendue pour les nombres négatifs, soit

$$[-1,797\,693 \times 10^{308}, -2,225\,074 \times 10^{-308}]$$

ou

$$[-1,797\,693 \times 10^{308}, -4,940\,656 \times 10^{-324}].$$

en considérant les nombres dénormalisés.

4. Soit  $\varepsilon$  la plus petite valeur tel que  $1 + \varepsilon \neq 1$  dans la représentation en virgule flottante. Cette valeur est appelée l'*epsilon de la machine* ou la *précision de la machine*. Or, puisque

$$1 = (-1)^0 \times 2^0 \times 1,00 \cdots 00,$$

le nombre suivant est

$$(-1)^0 \times 2^0 \times 1,00 \cdots 01.$$

Par conséquent, on a  $\varepsilon = 2^{-52} = 2,220\,446 \times 10^{-16}$ .



5. Tout nombre  $x$  représente en fait un intervalle de  $\mathbb{R}$ . Par exemple, tout nombre dans l'intervalle  $[1, 1 + \varepsilon)$  est représenté par le nombre 1 dans l'ordinateur.
6. On a un ensemble fini de nombres pour représenter  $\mathbb{R}$ . Or, cet ensemble est plus dense près de 0 que pour les grandes valeurs. En effet, le nombre suivant  $x_{\min}$  est

$$\begin{aligned} x_{\min}^+ &= (1 + \varepsilon) \times 2^{-1022} \\ &= x_{\min} + \varepsilon \times 2^{-1022} \\ &= x_{\min} + 2^{-1074}, \end{aligned}$$

alors que le nombre précédant  $x_{\max}$  est

$$\begin{aligned} x_{\max}^- &= (2 - 2^{-52} - \varepsilon) \times 2^{1023} \\ &= x_{\max} - \varepsilon \times 2^{1023} \\ &= x_{\max} - 2^{971}. \end{aligned}$$

Les deux plus petits nombres représentables sont donc distants de  $2^{-1024} \approx 10^{-324}$ , alors que les deux plus grands le sont de  $2^{971} \approx 10^{292}$  !

Tout calcul impliquant un nombre  $x > x_{\max}$  ( $x < -x_{\max}$ ) entraîne un *dépassement de capacité*. Habituellement, cela entraîne l'arrêt immédiat des calculs et un résultat de  $+\infty$  ( $-\infty$ ). À l'autre bout du spectre, un calcul impliquant un nombre  $x < x_{\min}$  entraîne un *souppassement de capacité* et le résultat est habituellement considéré égal à 0.

R conserve dans une liste nommée `.Machine` les valeurs mentionnées ci-dessus — ainsi que plusieurs autres — pour l'architecture de l'ordinateur courant. Par exemple, on confirme les valeurs de  $\varepsilon_m$ ,  $x_{\min}$ ,  $x_{\max}$ ,  $B$ ,  $m_f$  et  $m_e$ , dans l'ordre, pour l'architecture x86 :

```
> .Machine[c(1, 3:6, 11)]
$double.eps
[1] 2.220446e-16

$double.xmin
[1] 2.225074e-308

$double.xmax
[1] 1.797693e+308

$double.base
```

TAB. 10.5: Représentation en virgule flottante simplifiée

$x$	$\text{fl}(x)$
2,5	$0,25000 \times 10^1$
-42,182	$-0,42182 \times 10^2$
0,214356	$0,21436 \times 10^0$

[1] 2

\$double.digits

[1] 53

\$double.exponent

[1] 11

## 10.5 Éléments d'arithmétique en virgule flottante

La section précédente expliquait pourquoi les nombres stockés dans un ordinateur ne sont pas toujours ceux que l'on croit — ou que l'on voudrait. Puisque l'ordinateur ne peut représenter tous les nombres réels, toute opération arithmétique avec des nombres en virgule flottante implique une erreur d'arrondi ou de troncature (selon l'architecture de l'ordinateur) dont il importe de tenir compte lors de la mise en oeuvre de certains algorithmes. Cette section présente quelques grands principes d'arithmétique en virgule flottante ainsi que de bons usages pour minimiser les erreurs d'arrondi. Consulter les ouvrages [Monahan \(2001\)](#); [Burden et Faires \(1988\)](#); [Knuth \(1997\)](#), entre autres, pour une discussion plus complète de ce vaste sujet.

Aux fins de cette section, on note  $\text{fl}(x)$  la représentation en virgule flottante du nombre  $x$  et l'on définit la représentation simplifiée

$$\text{fl}(x) = (-1)^S \times 10^E \times 0,F, \quad (10.6)$$

où  $F$  compte cinq chiffres significatifs, dont le dernier est arrondi. Le [tableau 10.5](#) fournit quelques exemples.

*Remarque.* Dans la norme IEEE 754, les règles d'arrondi de base sont :

1. arrondir à la valeur la plus près (0,14 devient 0,1 et 0,16 devient 0,2) ;
2. arrondir un nombre exactement à mi-chemin au nombre avec un dernier chiffre significatif pair ou nul (0,05 devient 0,0 alors que 0,15 devient 0,2).

### 10.5.1 Erreurs d'arrondi ou de troncature

Lors de l'addition et de la soustraction de nombres en virgule flottante, on rend les exposants égaux en déplaçant les bits de la mantisse du plus petit nombre vers la droite. Il en résulte une perte de bits significatifs et donc une *erreur d'arrondi*. Dans les cas extrêmes où les deux opérandes sont de grandeur très différente, le plus petit opérande devient nul suite à la perte de tous ses bits significatifs.

La situation pour la multiplication et la division est quelque peu différente, mais la source d'erreur d'arrondi demeure la même.

À toute fin pratique, tout calcul fait naître de l'erreur d'arrondi et celle-ci augmente avec le nombre d'opérations. Les quelques principes de programmation ci-dessous, lorsque suivis par le programmeur prudent et consciencieux, permettent d'atténuer l'impact de l'erreur d'arrondi.

- (P1) L'addition et la soustraction en virgule flottante ne sont pas associatives. Additionner les nombres en ordre croissant de grandeur afin d'accumuler des bits significatifs.
- (P2) Éviter de soustraire un petit nombre d'un grand, ou alors le faire le plus tard possible dans la chaîne des calculs.
- (P3) Pour calculer une somme alternée, additionner tous les termes positifs et tous les termes négatifs, puis faire la soustraction.
- (P4) Multiplier et diviser des nombres d'un même ordre de grandeur. Si  $x$  et  $y$  sont presque égaux, le calcul  $x/y$  sera plus précis en virgule flottante que  $y^{-1}x$ .

Les exemples ci-dessous illustrent ces principes.

**Exemple 10.9.** Soit  $x = 7$ ,  $y = 4$  et  $z = 100\,000$ . Dans la représentation simplifiée (10.6), on a  $\text{fl}(x) = 0,70000 \times 10^1$ ,  $\text{fl}(y) = 0,40000 \times 10^1$  et  $\text{fl}(z) = 0,10000 \times 10^6$ . Or,  $x + y + z = 100\,011$  et

$$\begin{aligned} [\text{fl}(x) + \text{fl}(y)] + \text{fl}(z) &= (0,70000 \times 10^1 + 0,40000 \times 10^1) + 0,10000 \times 10^6 \\ &= 0,11000 \times 10^2 + 0,10000 \times 10^6 \\ &= 0,00001 \times 10^6 + 0,10000 \times 10^6 \\ &= 0,10001 \times 10^6, \end{aligned}$$

soit un résultat raisonnablement précis. Cependant,

$$\begin{aligned} \text{fl}(x) + [\text{fl}(y) + \text{fl}(z)] &= 0,70000 \times 10^1 + (0,00000 \times 10^6 + 0,10000 \times 10^6) \\ &= 0,00000 \times 10^6 + 0,10000 \times 10^6 \\ &= 0,10000 \times 10^6. \end{aligned}$$

En effectuant les opérations dans un mauvais ordre, on obtient  $x + y + z = z$  même si  $x \neq 0$  et  $y \neq 0$ .  $\square$

**Exemple 10.10.** Soit  $x = 7$ ,  $y = 100\,006$  et  $z = 100\,002$ . On a  $z - y - x = -11$ . Or,

$$\begin{aligned}\text{fl}(z) - [\text{fl}(y) + \text{fl}(x)] &= 0,10000 \times 10^6 - (0,10000 \times 10^6 + 0,00000 \times 10^6) \\ &= 0,00000 \times 10^0,\end{aligned}$$

alors que

$$\begin{aligned}[\text{fl}(z) - \text{fl}(y)] - \text{fl}(x) &= (0,10000 \times 10^6 - 0,10000 \times 10^6) - 0,70000 \times 10^1 \\ &= 0,00000 \times 10^1 - 0,70000 \times 10^1 \\ &= -0,70000 \times 10^1.\end{aligned}$$

$\square$

L'exemple suivant est adapté de [Monahan \(2001\)](#).

**Exemple 10.11.** L'évaluation numérique de probabilités loin dans les queues d'une fonction de répartition peut s'avérer imprécise selon la technique employée. Par exemple, la fonction de répartition de la distribution logistique est

$$F(x) = (1 + e^{-x})^{-1}.$$

La vraie valeur de  $\Pr[X > 6]$  est  $1 - F(6) = 1 - (1 + e^{-6})^{-1} = 0,002\,472\,623$ . L'évaluation de cette probabilité dans notre représentation en virgule flottante avec la formule ci-dessus est

$$\begin{aligned}1 - [1 \div (1 + \text{fl}(e^{-6}))] &= 1 - [1 \div (1 + 0,24788 \times 10^{-2})] \\ &= 1 - [0,10000 \times 10^1 \div 0,10025 \times 10^1] \\ &= 1 - 0,99751 \times 10^0 \\ &= 0,10000 \times 10^1 - 0,09975 \times 10^1 \\ &= 0,25000 \times 10^{-2}.\end{aligned}$$

(Nous n'avons pas écrit les entiers en virgule flottante ci-dessus pour alléger la notation.) Toutefois, si l'on utilise plutôt la formule équivalente

$$\begin{aligned}1 - F(6) &= \frac{e^{-6}}{1 + e^{-6}} \\ &= \frac{1}{1 + e^6}\end{aligned}$$

dans laquelle l'opération de soustraction entre deux nombres presque égaux est déjà effectuée, on obtient le résultat bien plus précis

$$\begin{aligned} 1 \div (1 + \text{fl}(e^6)) &= 1 \div (1 + 0,40343 \times 10^3) \\ &= 1 \div (0,00100 \times 10^3 + 0,40343 \times 10^3) \\ &= 1 \div 0,40443 \times 10^3 \\ &= 0,24726 \times 10^{-2}. \end{aligned}$$

□



Toutes les fonctions de calcul de fonction de répartition de R permettent d'évaluer plus précisément les probabilités  $1 - F(x)$  pour une grande valeur de  $x$  en spécifiant l'option `lower.tail = FALSE`.

Les deux principes ci-dessous permettent d'éviter des dépassements ou des soupassements de capacité et d'améliorer la précision des calculs.

- (P5) Chercher des formules mathématiques équivalentes évitant de devoir traiter des très grands ou des très petits nombres.
- (P6) Travailler en échelle logarithmique. Par exemple, le produit de deux grands nombres  $x$  et  $y$  dépassera la capacité plus tard et demeurera précis plus longtemps s'il est calculé sous la forme  $e^{\log x + \log y}$ .

**Exemple 10.12.** Soit  $X_1, \dots, X_n$  un échantillon aléatoire tiré d'une distribution de Pareto translatée, dont la fonction de répartition est

$$F(x; \mu, \alpha) = 1 - \left(\frac{\mu}{x}\right)^\alpha, \quad x > \mu.$$

On peut démontrer que l'estimateur du maximum de vraisemblance de  $\alpha$  est

$$\hat{\alpha} = \frac{n}{\ln\left(\frac{X_1 \cdots X_n}{X_{(1)}^n}\right)},$$

où  $X_{(1)} = \min(X_1, \dots, X_n)$ . Soit  $x$  un objet R contenant un échantillon de 1 000 valeurs d'une distribution Pareto translatée de moyenne 5 000 (le contenu de cet objet n'est pas affiché ici pour des raisons évidentes). Le calcul de l'estimateur  $\hat{\alpha}$  directement par la formule entraîne rapidement un dépassement de capacité, tant lors du produit  $X_1 \cdots X_n$  que lors de l'opération  $X_{(1)}^n$  :

```
> prod(x)
```

```
[1] Inf
> min(x)^length(x)
[1] Inf
```

Le résultat est donc NaN :

```
> length(x)/log(prod(x)/min(x)^length(x))
[1] NaN
```

Selon la grandeur des données dans l'échantillon, la formule équivalente

$$\hat{\alpha} = \frac{n}{\ln \prod_{i=1}^n \frac{X_i}{X_{(1)}}}$$

peut éviter le dépassement de capacité. Elle n'est toutefois d'aucun secours dans le présent exemple :

```
> length(x)/log(prod(x/min(x)))
[1] 0
```

On remarquera de plus que cette approche nécessite un grand nombre de multiplications (voir la [sous-section 10.5.2](#)), en plus d'ouvrir la porte à des erreurs d'arrondi.

Pour obtenir une réponse on utilisera plutôt une autre formule algébriquement équivalente :

$$\hat{\alpha} = \frac{n}{\sum_{i=1}^n \ln X_i - n \ln X_{(1)}}.$$

On obtient alors

```
> length(x)/(sum(log(x)) - length(x) * log(min(x)))
[1] 1.398743
```

Cet exemple illustre combien des calculs *algébriquement* équivalents ne sont pas nécessairement *numériquement* équivalents. □

TAB. 10.6: Coût relatif de quelques opérations en virgule flottante

Opération arithmétique	Coût relatif
Addition et soustraction	1,0
Multiplication	1,3
Division	3,0
Racine carrée	4,0
Logarithme	15,4

### 10.5.2 Coût des opérations

Les ordinateurs modernes disposent d'une unité de calcul en virgule flottante (FPU) intégrée au processeur. Cette unité est évidemment très optimisée et, par conséquent, elle accélère beaucoup le calcul en virgule flottante. Néanmoins, certaines opérations sont plus coûteuses à réaliser que d'autres en termes de temps de calcul. À titre indicatif, on trouve au [tableau 10.6](#) le coût relatif approximatif de quelques opérations en virgule flottante.

À titre d'exemple, considérons le calcul d'une simple moyenne arithmétique

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

Si l'on effectue le calcul tel qu'écrit ci-dessus, cela requiert  $n - 1$  additions et une division. En revanche, le calcul équivalent

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n},$$

requiert  $n - 1$  divisions et  $n - 1$  additions. Pour  $n$  grand utiliser la première approche plutôt que la seconde fera une différence.

## 10.6 Codage de caractères

Pendant que l'on discute de la représentation interne des nombres dans un ordinateur, on peut toucher un mot de la représentation interne, ou le codage, des caractères. Ce sujet demeure une cible mouvante puisque de nouveaux standards apparaissent encore après de nombreuses années de systèmes incompatibles et basés sur la langue anglaise.

Peu importe le système retenu, les caractères doivent être codés en binaire dans l'ordinateur. La partie la plus difficile consiste à créer un système standard

permettant aux ordinateurs et autres appareils numériques de communiquer entre eux. Le premier standard d'usage courant fut le Code américain normalisé pour l'échange d'information, mieux connu sous son acronyme ASCII ([ANSI, 1986](#)). La norme ASCII définit des codes numériques pour 128 caractères, soit 95 affichables (lettres, chiffres, symboles divers et l'espace) et 33 non affichables (essentiellement des caractères de contrôle tels que saut de ligne, retour de chariot ou espacement arrière). Par exemple, les lettres majuscules A–Z occupent les cases 65–90 (1000001–1011010 en binaire) alors que les lettres minuscules a–z occupent les cases 97–122 (1100001–1101010 en binaire). Ainsi, on constate que les versions majuscule et minuscule d'une même lettre ne diffèrent, dans leur représentation binaire, que par leur second bit le plus significatif. Le changement de casse est donc une opération très simple et rapide.

La représentation ASCII ne requiert en soi que sept bits. Néanmoins, on a rapidement codé les caractères sur un octet (huit bits) puisqu'il s'agit du type de donnée natif des ordinateurs depuis les années 1970. L'ajout d'un bit a créé de l'espace pour 128 caractères additionnels (cases 128–255). Une myriade de systèmes de codage différents sont alors apparus pour supporter les caractères des langues autres que l'anglais (les caractères accentués, par exemple) ainsi que d'autres symboles graphiques. La norme ISO 8859-1, ou Latin 1 ([ISO, 1998](#)), a fini par s'imposer comme l'un des standards les plus répandus. Ces listes de codes fixes se révèlent toutefois problématiques lors de l'apparition d'un nouveau symbole. Par exemple, pour faire de la place pour le symbole de l'euro à la fin des années 1990, il a fallu retirer un symbole de ISO 8859-1. Avec quelques autres changements, la nouvelle liste de code est devenue ISO 8859-15. De plus, comment doit-on traiter les langues CJC (chinois, japonais, coréen) qui comptent des milliers de symboles différents ?

Depuis le début des années 1990, la mondialisation de l'informatique a suscité un effort concerté de création d'un système de codage standard couvrant la presque totalité des systèmes d'écriture du monde. Le système devait aussi permettre la composition de textes en plusieurs langues, par exemple en français et dans une écriture de droite à gauche comme l'hébreu ou l'arabe. Le standard ainsi créé est Unicode ([Unicode Consortium, 2007](#)). Le plus populaire système de codage capable de représenter l'ensemble des caractères Unicode est *Unicode Transformation Format* 8 bits, ou UTF-8 ([Unicode Consortium, 2007](#), section 3.9). Dans l'UTF-8, chaque caractère est codé sur une suite d'un à quatre octets et les premiers 128 caractères sont identiques à la norme ASCII. L'UTF-8 est le système de codage par défaut dans MacOS X et les plus récentes distributions GNU/Linux.

Le système R supporte les caractères multi-octets de manière assez exhaustive. Sans doute aidés en cela par le fait qu'il s'agit d'un projet international, les développeurs de R ont mis beaucoup d'effort pour assurer l'internationalisation et la



localisation du logiciel ; voir [Ripley \(2005\)](#).

## 10.7 Code informatique

```
###
### EXEMPLES DE CALCULS EN APPARENCE ERRONÉS
###

## Opérateurs d'addition et de soustraction ne respectant pas les
## règles d'arithmétique de base.
1.2 + 1.4 + 2.8          # 5.4
1.2 + 1.4 + 2.8 == 5.4    # non?!?
2.8 + 1.2 + 1.4          # encore 5.4
2.8 + 1.2 + 1.4 == 5.4    # donc addition non commutative?
2.6 - 1.4 - 1.2          # devrait donner 0
2.6 - 1.5 - 1.1          # correct cette fois-ci

## Division donnant parfois des résultats erronés.
0.2/0.1 == 2              # ok
(1.2 - 1.0)/0.1 == 2      # pourtant équivalent
0.3/0.1 == 3              # fonctionnait pourtant avec 0.2

## Distance plus grande entre 3,2 et 3,15 qu'entre 3,10 et
## 3,15.
3.2 - 3.15 > 0.05         # écart supérieur à 0,05...
3.15 - 3.1 < 0.05         # et cette fois inférieur à 0,05

## Valeurs inexactes dans les suites de nombres générées avec
## 'seq', qui sont générées algébriquement (par additions ou
## soustractions successives).
(a <- seq(0.9, 0.95, by = 0.01))
a == 0.94                  # 0.94 n'est pas dans le vecteur!
(b <- c(0.90, 0.91, 0.92, 0.93, 0.94, 0.95))
b == 0.94                  # mais ici, oui!
a - b                     # remarquer le 5e élément

###
### CONVERSION DANS DES BASES GÉNÉRALES
###

## La fonction 'arrayInd', dont la syntaxe (simplifiée) est
##
##   arrayInd(ind, dim)
##
```

```

## retourne les coordonnées des éléments aux positions 'ind'
## dans un tableau de dimensions 'dim'. Par exemple, le 14e
## élément d'une matrice 4 x 5 remplie par colonne (ordre R)
## est en position (2, 4).
arrayInd(14, c(4, 5))      # comparer avec l'exemple 10.7

## C'est toujours un peu plus compliqué avec les tableaux. Par
## exemple, un tableau 3 x 4 x 5 doit être vu comme cinq
## matrices 3 x 4 placées les unes derrière les autres. Où se
## trouve le 'i'ème élément?
arrayInd(8, 3:5)           # élément (2, 3) de 1ère matrice
arrayInd(13, 3:5)          # élément (1, 1) de 2e matrice
arrayInd(59, 3:5)          # élément (2, 4) de 5e matrice
arrayInd(c(8, 13, 59), 3:5) # en un seul appel

```

## 10.8 Exercices

**10.1** Convertir les nombres décimaux suivants en base 6, puis en binaire.

- a) 119
- b) 343
- c) 96
- d) 43

**10.2** Convertir les nombres hexadécimaux suivants en nombres décimaux.

- a) A1B
- b) 12A
- c) B41
- d) BAF FE

**10.3** a) Utiliser l'algorithme de conversion des nombres en base  $b$  vers la base 10 et les idées de l'exemple 10.7 pour trouver une formule générale donnant la position de l'élément  $a_{ijk}$  d'un tableau de dimensions  $I \times J \times K$  dans l'ordre de la liste des éléments du tableau. Utiliser l'ordre lexicographique, où le tableau est rempli dans l'ordre

$$a_{111}, a_{112}, \dots, a_{11K}, a_{121}, a_{122}, \dots$$

- b) Répéter la partie a) en utilisant l'ordre R, où le tableau est plutôt rempli dans l'ordre  $a_{111}, a_{211}, \dots, a_{I11}, a_{121}, a_{221}, \dots$ . Comparer la réponse avec celle de l'exercice 3.7 b) de la partie I.

**10.4** La plupart des langages de programmation — dont R, voir la note au bas de la page 14 — permettent de définir des nombres entiers «stricts», par opposition à des nombres réels qui s'avèrent simplement être entiers. La représentation de ces nombres entiers est complètement différente de celle de la norme IEEE 754, sauf en ce qui a trait au bit fort, qui sert toujours à stocker le signe du nombre.

Les nombres entiers sont stockés dans une représentation dite «en complément à deux» (Wikipedia, 2013). Dans cette notation :

- i) un entier positif est représenté en binaire de manière usuelle, avec le bit fort égal à 0 ;
- ii) pour un entier négatif, on inverse les bits de la représentation de sa valeur absolue et on ajoute 1 au résultat.

Le bit fort est automatiquement mis à 1 (pour représenter un nombre négatif) par l'opération d'inversion. Les deux opérations ci-dessus sont équivalentes à stocker un nombre négatif  $x$  en tant que  $2^n - |x|$ .

Par exemple, la représentation sur 8 bits du nombre 5 est

$$\boxed{0} \boxed{0000101}.$$

Pour représenter  $-5$ , on inverse ces bits et on ajoute 1 :

$$\boxed{1} \boxed{1111011}.$$

On remarque que  $11111011_2 = 251_{10} = 2^8 - 5$ .

- a) L'un des avantages de la notation en complément à deux est qu'elle évite d'avoir des représentations pour  $+0$  et  $-0$ . Pour des entiers de 8 bits,  $\boxed{0} \boxed{0000000}$  correspond tout naturellement à 0. À quel nombre correspond  $\boxed{1} \boxed{0000000}$  selon cette notation ?
- b) Un autre avantage, plus important que celui mentionné en a), de la notation en complément à deux est qu'elle permet d'effectuer les opérations arithmétiques en binaire naturellement. Pour constater ce fait, effectuer directement sur 8 bits l'opération  $15 - 5$  en utilisant les notations suivantes :
  - i) la notation naturelle où le bit fort représente le signe et les sept autres bits pour représentent la valeur absolue du nombre ;
  - ii) la notation en complément à deux.

Dans les deux cas, ignorer les éventuels dépassements de capacité vers la gauche.

- c) C'est maintenant le temps d'expliquer le trait d'humour dans le strip de XKCD.com au début du document. Pour ce faire, convertir en décimal les représentations en complément à deux sur 8 bits de  $\boxed{00000000}$  à  $\boxed{11111111}$ . Répéter l'exercice sur 16 bits. (En fait, convertir non pas toutes les représentations, mais bien seulement quelques valeurs clés qui permettront d'observer la séquence des nombres décimaux correspondants.)
- d) En plus des types Single et Double pour représenter des nombres en virgule flottante, le VBA dispose des types Byte, Integer et Long pour représenter des nombres entiers. Comme son nom l'indique, le type Byte utilise huit bits d'espace mémoire et ne sert que pour les entiers positifs. Les types Integer et Long requièrent 16 et 32 bits, respectivement, et peuvent contenir des nombres négatifs. La notation en complément à deux est alors utilisée pour ces deux types. Déterminer l'ensemble des entiers représentables pour chacun des types de données Byte, Integer et Long.

**10.5** La norme IEEE 754 pour les nombres en virgule flottante ( $S, E, F$ ) en simple précision est le suivant :

- ▶ longueur totale de  $m = 32$  bits ;
- ▶ 1 bit pour le signe  $S$  (valeur de 0 pour un nombre positif) ;
- ▶ 8 bits pour l'exposant  $E$ , avec un biais de 127 ;
- ▶ 23 bits pour la partie fractionnaire  $F$ .

Un nombre  $x$  est donc représenté comme

$$x = (-1)^S \times 2^{E-127} \times 1.F.$$

Trouver les valeurs  $\epsilon$ ,  $x_{\max}$  et  $x_{\min}$  pour les nombres en simple précision. Comparer les résultats avec les limites du type Single en VBA.

**10.6** Représenter les nombres suivants comme des nombres en virgule flottante en simple précision selon la norme IEEE 754.

- a)  $-1\,234$
- b)  $55$
- c)  $8\,191$
- d)  $-10$
- e)  $\frac{2}{3}$
- f)  $\frac{1}{100}$

**10.7** Les nombres ci-dessous sont représentés en format binaire selon la norme IEEE 754 pour les nombres en simple précision. Convertir ces nombres en décimal.

- a) 

0
---

00111101
----------

100100001000000000000000
--------------------------
- b) 

1
---

00111101
----------

100100001000000000000000
--------------------------
- c) 

0
---

10000100
----------

100100001000000000000000
--------------------------
- d) 

1
---

10000100
----------

100100001000000000000000
--------------------------

**10.8** Trouver, pour les nombres des parties a) et c) de l'exercice 10.7, le nombre suivant et le nombre précédent en représentation binaire.

## Réponses

La notation  $x_b$  signifie que le nombre  $x$  est en base  $b$ . On omet généralement  $b$  pour les nombres en base 10.

**10.1** a)  $315_6, 1110111_2$

b)  $1331_6, 101010111_2$

c)  $240_6, 1100000_2$

d)  $111_6, 101011_2$

**10.2** a) 2 587 b) 298 c) 2 881 d) 765 950

**10.3** a)  $k + K(j - 1 + J(i - 1))$

b)  $i + I(j - 1 + J(k - 1))$

**10.4** a) -128

b) i) -20 ii) 10

d) Type Byte :  $\{0, \dots, 255\}$  ; type Integer :  $\{-32768, \dots, 32767\}$  ;

type Long :  $\{-2\,147\,483\,648, \dots, 2\,147\,483\,647\}$

**10.5**  $\varepsilon = 2^{-23} = 1,192 \times 10^{-7}$ ,  $x_{\max} = (2 - 2^{-23}) \times 2^{127} = 3,403 \times 10^{38}$ ,  $x_{\min} = 2^{-126} = 1,175 \times 10^{-38}$  (nombre normalisé) ou  $x_{\min} = 2^{-149} = 1,401 \times 10^{-45}$  (nombre dénormalisé)

**10.6** a) 

1
---

10001001
----------

001101001000000000000000
--------------------------

b) 

0
---

10000100
----------

101110000000000000000000
--------------------------

c) 

0
---

10001011
----------

111111111110000000000000
--------------------------

d) 

1
---

10000010
----------

010000000000000000000000
--------------------------

e) 

0
---

01111110
----------

0101010101010101010101010
---------------------------

f) 

0
---

01111000
----------

01000111101011100001010
-------------------------

**10.7** a)  $2,120229346 \times 10^{-20}$  b)  $-2,120229346 \times 10^{-20}$  c) 50,0625 d)  $-50,0625$

**10.8** a)  $2,120229508 \times 10^{-20}$  et  $2,120229185 \times 10^{-20}$

b) 50,062503815 et 50,062496185

# 11 Résolution d'équations à une variable

## Objectifs du chapitre

- ▶ Résoudre une équation à une variable à l'aide des méthodes de la bisection, du point fixe et de Newton-Raphson.
- ▶ Faire un choix parmi l'une des méthodes ci-dessous pour résoudre un problème donné.
- ▶ Déterminer si une fonction admet un point fixe dans un intervalle donné.
- ▶ Utiliser les fonctions de résolution d'équations et d'optimisation de Excel et R.

## 11.1 Mise en contexte

Ce chapitre passe en revue les méthodes les plus populaires pour résoudre numériquement une équation à une variable. Ce problème, simple en apparence seulement, survient dans une foule de domaines. Débutons par deux exemples, tirés des mathématiques financières et de l'inférence statistique.

**Exemple 11.1.** Un investisseur est disposé à payer  $k$  \$ aujourd'hui en retour d'une série de  $n$  versements de 1 \$, le premier dans un an. Pour déterminer le taux de rendement qu'il exige, on doit trouver la valeur de  $i$  tel que  $a_{\overline{n}|i} = k$  ou, énoncé autrement, la valeur de  $i$  tel que

$$f(i) = \frac{1 - (1 + i)^{-n}}{i} - k = 0.$$

□

**Exemple 11.2.** On a un échantillon aléatoire  $X_1, \dots, X_n$  provenant d'une distribution binomiale négative avec fonction de masse de probabilité

$$\Pr[X = x] = \binom{x+r-1}{r-1} p^r (1-p)^x, \quad x = 0, 1, \dots,$$

où  $p$  est supposé connu. On souhaite trouver l'estimateur du maximum de vraisemblance du paramètre  $r$ .

La fonction de vraisemblance est

$$\begin{aligned} L(r) &= \prod_{i=1}^n \Pr[X = x_i] \\ &= \prod_{i=1}^n \binom{x_i+r-1}{r-1} p^r (1-p)^{x_i} \end{aligned}$$

et la fonction de log-vraisemblance est

$$\begin{aligned} l(r) &= \sum_{i=1}^n [\ln \Gamma(x_i + r) - \ln \Gamma(r) - \ln \Gamma(x_i - 2) \\ &\quad + r \ln p + x_i \ln(1-p)]. \end{aligned}$$

Par conséquent, l'estimateur du maximum de vraisemblance  $\hat{r}$  est la solution de

$$f(r) = \frac{d}{dr} l(r) = \sum_{i=1}^n (\Psi(x_i + r) - \Psi(r) + \ln p) = 0,$$

où  $\Psi(x) = \Gamma'(x)/\Gamma(x)$  est la *fonction digamma*. □

Le problème sur lequel nous allons nous pencher consiste donc à trouver la racine d'une fonction  $f(x)$ , c'est-à-dire la solution de  $f(x) = 0$ .

## 11.2 Méthode de bisection

Toutes les méthodes numériques de résolution d'une équation qui seront étudiées dans ce chapitre reposent sur une procédure systématique d'essais et d'erreurs. La procédure — ou algorithme — utilisée déterminera l'efficacité de la méthode. À ce chapitre, la méthode de bisection est la plus simple et la plus intuitive méthode de résolution, mais aussi la moins efficace.

Soit  $f$  une fonction continue sur un intervalle  $[a, b]$  choisi de telle sorte que  $f(a)$  et  $f(b)$  sont de signes opposés. Par le théorème de la valeur intermédiaire, il existe un point  $x^*$  tel que  $f(x^*) = 0$ . En d'autres mots, si une fonction continue



est d'un côté de l'axe des abscisses au début d'un intervalle et de l'autre côté à l'autre extrémité de l'intervalle, alors elle a forcément croisé l'axe quelque part dans l'intervalle.

L'idée de la méthode de bisection consiste à trouver la solution  $x^*$  par essais successifs en utilisant le point milieu d'intervalles de plus en plus petits, mais que l'on sait toujours contenir le point cherché. L'algorithme 11.1 systématise cette idée et la figure 11.1 l'illustre pour quelques itérations de l'algorithme.

**Algorithme 11.1** (Méthode de la bisection). *Soit  $f$  une fonction continue sur l'intervalle  $[a, b]$  où  $f(a)$  et  $f(b)$  sont de signes opposés. On cherche  $x^*$  tel que  $f(x^*) = 0$  en un maximum de  $N_{\max}$  itérations.*

1. Poser  $n = 1$ .
2. Répéter les étapes suivantes.
  - 2.1. Poser  $x_n = (a + b)/2$ .
  - 2.2. Si  $f(a)f(x_n) > 0$ , poser  $a = x_n$ , sinon poser  $b = x_n$ .
  - 2.3. Si  $|x_n - x_{n-1}|/|x_n| < \varepsilon$ , poser  $x^* = x_n$  et terminer.
  - 2.4. Si  $n \geq N_{\max}$ , terminer sans convergence.
  - 2.5. Poser  $n = n + 1$ .

L'application de l'algorithme 11.1 génère une suite  $\{x_n\}$ ,  $n = 1, 2, \dots$  qui, si tout se passe bien, converge vers la valeur  $x^*$ .

*Remarques.*

1. À l'étape 2.2 de l'algorithme, le produit  $f(a)f(x_n)$  permet de déterminer si l'on se trouve du même côté de l'axe des abscisses aux points  $a$  et  $x_n$  (produit de deux valeurs négatives ou deux valeurs positives) ou de part et d'autre de l'axe (produit de valeurs de signes opposés).
2. L'étape 2.3 permet de valider si l'algorithme a convergé vers une réponse. La valeur de  $\varepsilon$  sera donc en général «petite».
3. Nous avons utilisé l'erreur relative comme critère d'arrêt à l'étape 2.3. Il s'agit en général du meilleur choix. En effet, un critère tel que  $|f(x_n)| < \varepsilon$  peut être satisfait même si  $x_n \gg x^*$  ou  $x_n \ll x^*$  (voir la figure 11.2 pour deux exemples). Quant à l'erreur absolue  $|x_n - x_{n-1}| < \varepsilon$ , l'inégalité peut être satisfaite même si la suite  $\{x_n\}$  diverge.

**Exemple 11.3.** On cherche la racine de  $f(x) = x^3 + 4x^2 - 10$  dans l'intervalle  $[1, 2]$ . On peut vérifier que la fonction  $f$  n'a qu'une seule racine dans cet intervalle et que, de plus,  $f(1) = -5$  et  $f(2) = 14$ . Ainsi,  $a = 1$  et  $b = 2$  constituent des valeurs de départ adéquates.

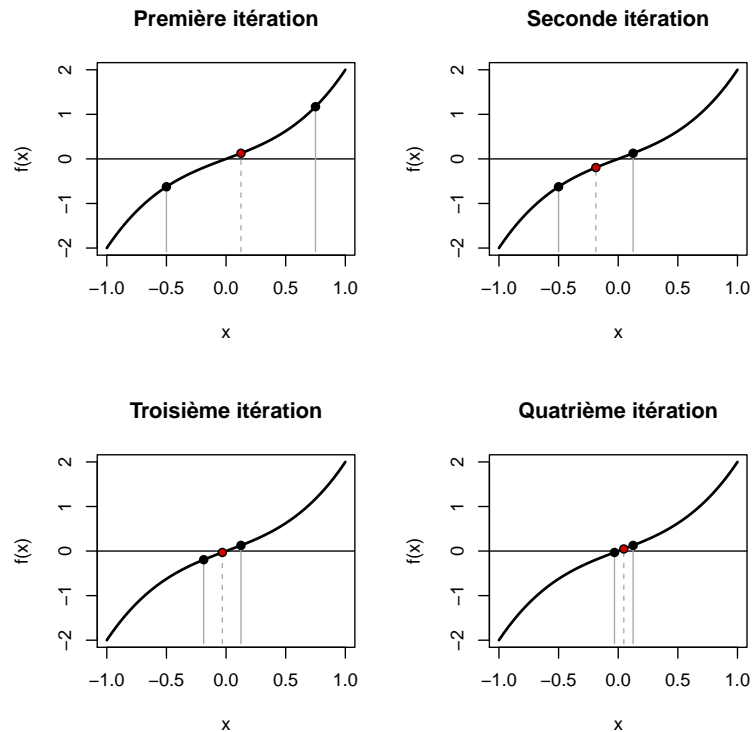


FIG. 11.1: Illustration de la méthode de bisection. À chaque itération, les points noirs reliés par des segments verticaux en ligne pleine identifient les bornes de l'intervalle dans lequel on sait se trouver la solution de l'équation  $f(x) = 0$ . Le point rouge relié par un segment en pointillé est l'essai de l'itération, c'est-à-dire le point milieu de l'intervalle. On remarquera comment ce point devient la borne inférieure ou supérieure de l'intervalle à l'itération suivante.

Le code R de la [section 11.8](#) définit une fonction `bisection` dont les arguments sont :

1. la fonction  $f(x)$ ;
2. la valeur de  $a$ ;
3. la valeur de  $b$ ;
4. la valeur de  $\varepsilon$  ( $10^{-6}$  par défaut) ;
5. le nombre maximal d'itérations  $N_{\max}$  (100 par défaut) ;

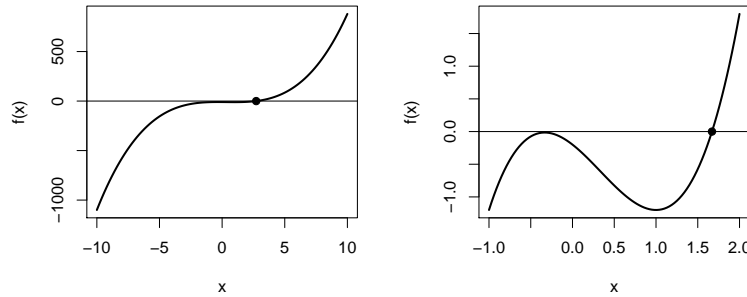


FIG. 11.2: Deux exemples de fonctions pour lesquelles on pourrait avoir  $|f(x_n)| < \varepsilon$  même si  $x_n \gg x^*$  ou  $x_n \ll x^*$ . Le point identifie la véritable racine dans chaque graphique.

6. une valeur booléenne indiquant si les calculs intermédiaires doivent ou non être affichés à l'écran : si l'argument est TRUE, les valeurs de  $a$ ,  $b$ ,  $x_n$  et  $f(x_n)$  sont affichées à chaque itération.

La fonction retourne une liste dont l'élément `$root` contient la réponse et l'élément `$nb.iter` le nombre d'itérations nécessaire pour obtenir cette réponse.

► Exécuter le code informatique de la [section 11.8](#) pour définir la fonction `bisection` puis résoudre cet exemple.

Après 17 itérations de la fonction `bisection`, la réponse obtenue est, à neuf décimales, 1,365 226 746. L'erreur d'approximation est au maximum de

$$|x^* - x_{17}| < |b_{18} - a_{18}| = |1,365234 - 1,365227| = 0,000007$$

et, puisque  $|a_{18}| < |x^*|$ ,

$$\frac{|x^* - x_{17}|}{|x^*|} < \frac{|b_{18} - a_{18}|}{|a_{18}|} = 5,13 \times 10^{-6}.$$

La réponse est donc exacte à au moins cinq décimales près. □

**Exemple 11.4.** La valeur présente d'une série de 10 paiements de fin d'année est 8,2218. Pour déterminer le taux d'intérêt (ou taux de rendement, selon le point de vue), on cherche  $i$  tel que  $f(i) = a_{\overline{10}|i} - 8,2218 = 0$ .

► Exécuter le code informatique de la [section 11.8](#) correspondant à cet exemple pour la solution.

□

### 11.3 Méthode du point fixe

La solution de l'équation

$$g(x) = x$$

est un *point fixe* de la fonction  $g$ . Il s'agit du point où les fonctions  $g(x)$  et  $y = x$  sont égales, c'est-à-dire où la fonction  $g$  croise la droite à 45° dans le plan.

Tout problème de recherche de racine peut être exprimé comme un problème de point fixe. Par exemple, en définissant

$$g(x) = x - f(x),$$

les problèmes de trouver  $x$  tel que  $g(x) = x$  ou  $f(x) = 0$  sont équivalents.

La méthode du point fixe repose sur le théorème du même nom. Celui-ci fournit les conditions garantissant, d'une part, l'existence d'un point fixe pour une fonction  $g$  dans un intervalle donné et, d'autre part, l'unicité de ce point fixe. Un point fixe unique peut néanmoins exister sans que ces conditions ne soient remplies.

**Théorème 11.1** (Théorème du point fixe). *Soit  $g$  une fonction continue sur  $[a, b]$ .*

1. *Si  $g(x) \in [a, b]$  pour tout  $x \in [a, b]$ , alors  $g$  a un point fixe dans  $[a, b]$ .*
2. *Si, de plus,  $g'(x)$  existe sur  $(a, b)$  et qu'il existe une constante  $k$  tel que*

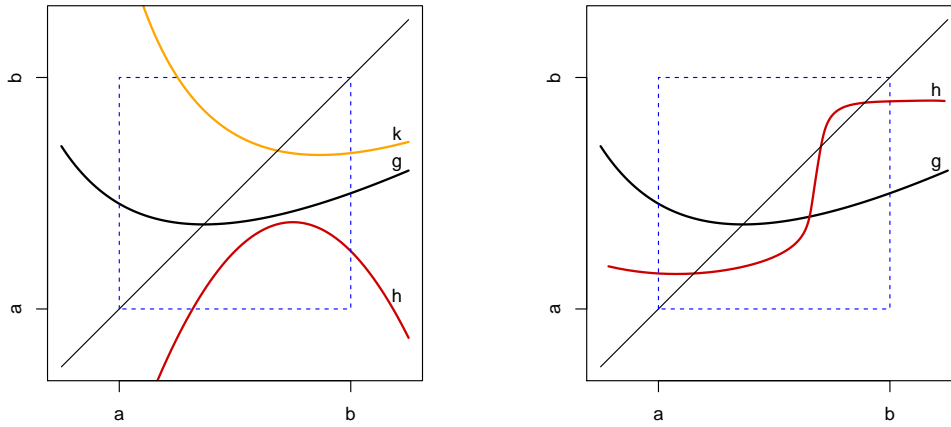
$$|g'(x)| \leq k < 1$$

*pour tout  $x \in (a, b)$ , alors  $g$  a un point fixe unique dans  $[a, b]$ .*

*Démonstration.* Le théorème compte deux volets. Nous en faisons une démonstration simplifiée ne reposant que sur des arguments graphiques.

1. Le premier volet garantit l'existence d'un point fixe sur  $[a, b]$  si l'image de la fonction  $g$  se trouve aussi dans un intervalle  $[a, b]$ . Cette paire d'intervalles engendre le carré en pointillé que l'on retrouve dans les graphiques de la [figure 11.3](#).

En définitive, le premier volet du théorème établit que si la fonction continue  $g$  entre dans l'intervalle  $[a, b]$  par la gauche du carré ( $g(x) \in [a, b]$ ) et en ressort par la droite (idem), alors elle doit forcément croiser la droite  $y = x$ . Observons le graphique (a) de la [figure 11.3](#) :



(a) existence du point fixe

(b) unicité du point fixe

FIG. 11.3: Exemples de graphiques permettant de faire la démonstration du théorème du point fixe

- ▶  $g(x) \in [a, b]$  pour tout  $x \in [a, b]$ , donc un point fixe existe nécessairement ;
  - ▶  $h(x) \notin [a, b]$  pour tout  $x \in [a, b]$ , donc la fonction n'a pas nécessairement de point fixe dans  $[a, b]$  ;
  - ▶  $k(x) \notin [a, b]$  pour tout  $x \in [a, b]$ , mais un point fixe demeure néanmoins possible.
2. Le second volet garantit l'unicité du point fixe si la pente de la fonction sur  $(a, b)$  n'excède jamais 1. Cette condition empêche la fonction de croiser la droite  $y = x$  une première fois et de la croiser de nouveau suite à une forte augmentation. Comparons les deux fonctions du graphique (b) de la [figure 11.3](#) :
- ▶  $g'(x) < 1$  pour tout  $x \in (a, b)$ , donc  $g(x)$  croise  $y = x$  en un seul point et le point fixe est unique ;
  - ▶  $h(x) \in [a, b]$  pour tout  $x \in [a, b]$ , mais  $h'(x) > 1$  pour  $x \in (a, b)$ , ce qui ouvre la porte à des points fixes multiples, tel qu'illustré.

□

Le théorème suivant, donné sans preuve, établit véritablement la procédure qui sera utilisée par la méthode numérique pour trouver le point fixe d'une fonction.

**Théorème 11.2.** *Si les deux conditions du [théorème 11.1](#) sont satisfaites, alors la série  $\{x_n = g(x_{n-1})\}$  converge vers un point fixe unique dans  $[a, b]$ .*

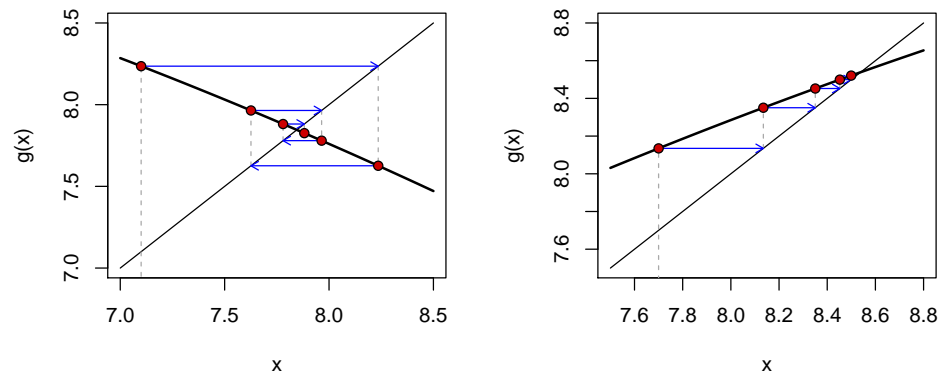


FIG. 11.4: Illustration de la méthode du point fixe pour une fonction décroissante et une fonction croissante. Les points rouges identifient les essais successifs, le premier étant, ici, le point le plus à gauche dans les graphiques. Les flèches illustrent comment est déterminé l'essai suivant. La relation  $x_n = g(x_{n-1})$  implique que l'abscisse d'un essai correspond à la projection de l'ordonnée de l'essai précédent sur la droite  $y = x$ .

Tel que mentionné au début de la [section 11.2](#), toutes les méthodes numériques — et la méthode du point fixe ne fait pas exception — procèdent par essais et erreurs, mais de manière systématique. Le [théorème 11.2](#) nous dit que pour trouver un point fixe, le meilleur essai  $x_n$  à faire est la valeur de la fonction  $g$  à l'essai précédent,  $g(x_{n-1})$ . La [figure 11.4](#) illustre comment la suite des essais successifs peut converger vers le point fixe.

Nous pouvons dès lors composer un algorithme de résolution par la méthode du point fixe.

**Algorithme 11.2** (Méthode du point fixe). *Soit  $g$  une fonction et  $x_0$  une valeur de départ. On cherche  $x^*$  tel que  $g(x^*) = x^*$  en un maximum de  $N_{\max}$  itérations.*

1. Poser  $n = 1$ .
2. Répéter les étapes suivantes.
  - a) Poser  $x_n = g(x_{n-1})$ .
  - b) Si  $|x_n - x_{n-1}|/|x_n| < \varepsilon$ , alors poser  $x^* = x_n$  et terminer.
  - c) Si  $n \geq N_{\max}$ , terminer sans convergence.
  - d) Poser  $n = n + 1$ .

*Remarque.* Comme on peut le déduire des illustrations de la [figure 11.4](#), la rapidité de la convergence est fonction de  $g'(x)$  : plus la pente est *faible*, plus la convergence est *rapide*, et vice versa.

**Exemple 11.5.** On répète l'[exemple 11.4](#) à l'aide de la méthode du point fixe, soit trouver le taux d'intérêt  $i$  tel que  $a_{\overline{10}|i} = 8,2218$ . On a déjà déterminé à l'[exemple 11.4](#) que le taux d'intérêt se trouve dans l'intervalle  $[0,035, 0,040]$ . On peut exprimer le problème sous forme de point fixe ainsi : trouver la valeur de  $i$  tel que

$$g(i) = \frac{1 - (1 + i)^{-10}}{8,2218} = i.$$

Le code R de la [section 11.8](#) définit une fonction `pointfixe` similaire à la fonction `bissection` présentée à la [section 11.2](#). Ses arguments sont d'ailleurs les mêmes, à l'exception que les points  $a$  et  $b$  sont remplacés par une valeur de départ  $x_0$ .

► Exécuter le code informatique de la [section 11.8](#) correspondant à ce bloc de matière pour définir la fonction `pointfixe` puis résoudre cet exemple.

Il faut 34 itérations pour obtenir la convergence avec un critère identique à celui utilisé avec la bissection. Pour ce type de problème, la méthode du point fixe n'est pas la plus efficace. □

**Exemple 11.6.** On répète l'[exemple 11.3](#) à l'aide de la méthode du point fixe, à savoir : trouver la racine de  $f(x) = x^3 + 4x^2 - 10 = 0$  dans l'intervalle  $[1, 2]$ .

Il y a plusieurs façons d'exprimer le problème de recherche de racine en termes de point fixe. On examine cinq différentes fonctions toutes algébriquement équivalentes :

$$\begin{aligned} g_1(x) &= x - x^3 - 4x^2 + 10 \\ g_2(x) &= \left( \frac{10}{x} - 4x \right)^{1/2} \\ g_3(x) &= \frac{1}{2}(10 - x^3)^{1/2} \\ g_4(x) &= \left( \frac{10}{4 + x} \right)^{1/2} \\ g_5(x) &= x - \frac{x^3 + 4x^2 - 10}{3x^2 + 8x}. \end{aligned}$$

Il est laissé en exercice de vérifier que, dans chacun des cas, la racine de  $f(x)$  est le point fixe de  $g_i(x)$ ,  $i = 1, \dots, 5$ .

Les graphiques des cinq fonctions ci-dessus se trouvent à la [figure 11.5](#). L'examen de ces graphiques permet de déterminer, avant de faire tout calcul, pour quelles fonctions la procédure itérative du point fixe convergera ainsi que les taux de convergence relatif, le cas échéant.

► Compléter cet exemple en exécutant le code informatique correspondant à la [section 11.8](#).

□

## 11.4 Méthode de Newton-Raphson

La méthode de Newton-Raphson est l'une des plus populaires et puissantes méthodes numériques de résolution d'équations à une variable. Les algorithmes plus élaborés qui existent aujourd'hui sont d'ailleurs souvent basés sur cet algorithme.

L'étude de la méthode de Newton-Raphson nous ramène au problème de recherche de la racine  $x^*$  d'une fonction  $f(x)$ , c'est-à-dire la solution de  $f(x) = 0$ . Il y a plusieurs manières d'introduire cette méthode de résolution. Nous privilégierons une approche graphique.

Dans la méthode de Newton-Raphson, les essais successifs sont déterminés en utilisant les points où la tangente de la fonction  $f$  à l'essai précédent croise l'axe des abscisses ; voir la [figure 11.6](#) pour une illustration.

Une première justification mathématique de la méthode va comme suit. Soit  $\tilde{x}$  un point «près» de la racine  $x^*$  et tel que  $f(\tilde{x}) \neq 0$ . Comme on peut le voir à la [figure 11.6](#), la tangente de  $f$  en  $\tilde{x}$  croise l'abscisse en un point  $\hat{x}$  «près» de  $x^*$ . Ainsi, le point  $\hat{x}$  est tel que

$$\frac{f(\tilde{x}) - 0}{\tilde{x} - \hat{x}} = f'(\tilde{x}),$$

d'où

$$\begin{aligned} x^* &\simeq \hat{x} \\ &= \tilde{x} - \frac{f(\tilde{x})}{f'(\tilde{x})}. \end{aligned}$$

On peut également justifier la formule précédente à l'aide du développement de Taylor de  $f(x)$  autour de  $\tilde{x}$  :

$$f(x) = f(\tilde{x}) + (x - \tilde{x})f'(\tilde{x}) + \frac{(x - \tilde{x})^2}{2}f''(\tilde{x}) + \dots$$



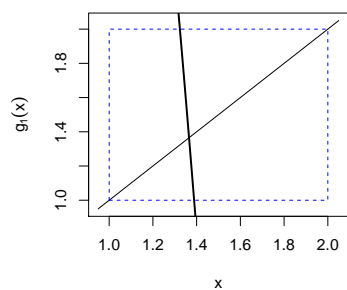
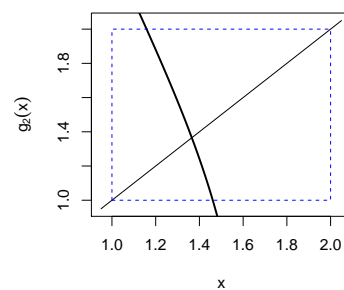
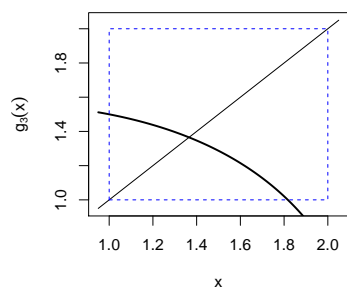
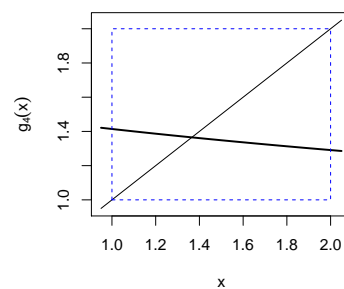
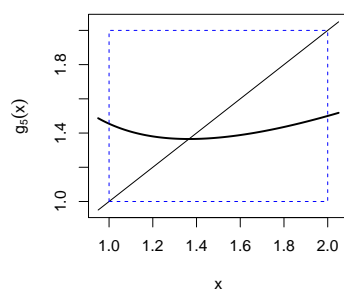
(a) fonction  $g_1(x)$ (b) fonction  $g_2(x)$ (c) fonction  $g_3(x)$ (d) fonction  $g_4(x)$ (e) fonction  $g_5(x)$ 

FIG. 11.5: Graphiques des cinq fonctions de l'exemple 11.6

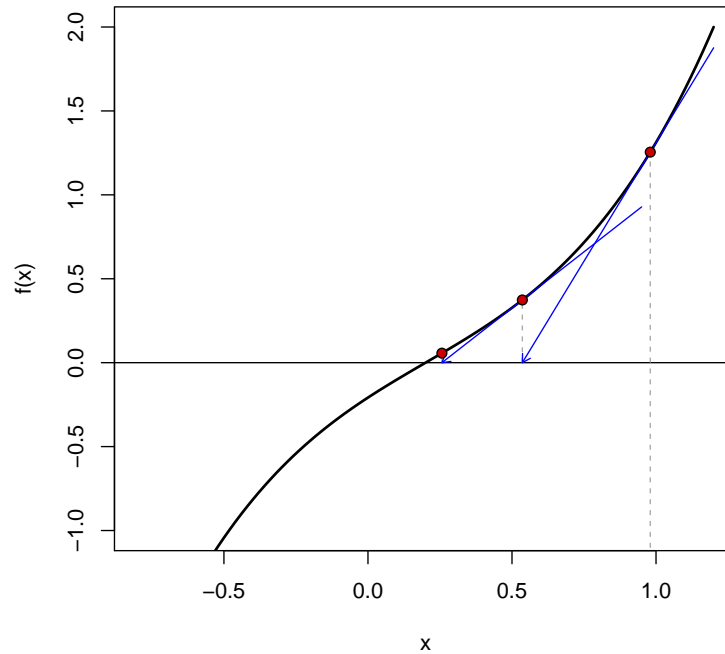


FIG. 11.6: Illustration de la méthode de Newton–Raphson. Les points rouges identifient les essais successifs, le premier étant, ici, le point le plus à droite dans le graphique. Les flèches illustrent comment est déterminé l’essai suivant. On trace la tangente en un point ; le prochain essai est la valeur en abscisse où cette tangente croise l’axe.

Toujours en supposant que  $\tilde{x}$  est près de  $x^*$ , on a

$$f(x^*) = 0 = f(\tilde{x}) + (x^* - \tilde{x})f'(\tilde{x}) + \frac{(x^* - \tilde{x})^2}{2}f''(\tilde{x}) + \dots,$$

où l’on peut considérer les termes de puissance 2 et plus dans le développement comme négligeables. On obtient alors l’approximation

$$x^* \simeq \tilde{x} - \frac{f(\tilde{x})}{f'(\tilde{x})}.$$

Tel qu’illustré à la [figure 11.6](#), on peut répéter la procédure ci-dessus si le point  $\hat{x}$  est trop éloigné de la racine  $x^*$ . On utilise alors  $\hat{x}$  comme nouveau point

où l'on calcule la tangente, et ainsi de suite jusqu'à l'obtention d'une «bonne» approximation.

Formellement, soit  $f$  une fonction différentiable sur  $[a, b]$  et  $x^* \in [a, b]$ . Alors  $x^*$  peut être obtenu comme le point de convergence de la série  $\{x_n\}$  définie par

$$x_n = x_{n-1} - \frac{f(x_{n-1})}{f'(x_{n-1})},$$

avec  $x_0$  une valeur de départ quelconque. L'algorithme 11.3 systématise cette procédure.

**Algorithme 11.3** (Méthode de Newton–Raphson). *Soit  $f$  une fonction continue et différentiable sur l'intervalle  $[a, b]$ . On cherche  $x^*$  tel que  $f(x^*) = 0$  en un maximum de  $N_{\max}$  itérations avec une valeur de départ  $x_0$ .*

1. Poser  $n = 1$ .
2. Répéter les étapes suivantes.
  - a) Poser  $x_n = x_{n-1} - f(x_{n-1}) / f'(x_{n-1})$ .
  - b) Si  $|x_n - x_{n-1}| / |x_n| < \varepsilon$ , alors poser  $x^* = x_n$  et terminer.
  - c) Si  $n \geq N_{\max}$ , terminer sans convergence.
  - d) Poser  $n = n + 1$ .

*Remarques.*

1. L'examen attentif de la méthode de Newton–Raphson révèle qu'il s'agit en fait d'une procédure de point fixe avec

$$g(x) = x - \frac{f(x)}{f'(x)}.$$

La méthode de Newton–Raphson nous fournit donc la «bonne» fonction  $g$  à utiliser dans la méthode du point fixe pour trouver la racine de la fonction  $f$ .

2. On peut démontrer que si  $f$  est doublement différentiable sur  $[a, b]$  et que  $x^* \in [a, b]$  est tel que  $f(x^*) = 0$  et  $f'(x^*) \neq 0$ , alors il existe un  $\delta > 0$  tel que la série  $\{x_n\}$  converge vers  $x^*$  pour tout  $x_0 \in [x^* - \delta, x^* + \delta]$ . Autrement dit :
  - les hypothèses du [théorème 11.1](#) sont satisfaites ;
  - la méthode de Newton–Raphson fonctionne toujours avec un «bon» choix de valeur de départ  $x_0$ .

**Exemple 11.7.** On reprend l'exemple 11.5 de calcul d'un taux de rendement, cette fois à l'aide de la méthode de Newton–Raphson. Vous savez peut-être de votre cours de mathématiques financières que la méthode de Newton–Raphson est la

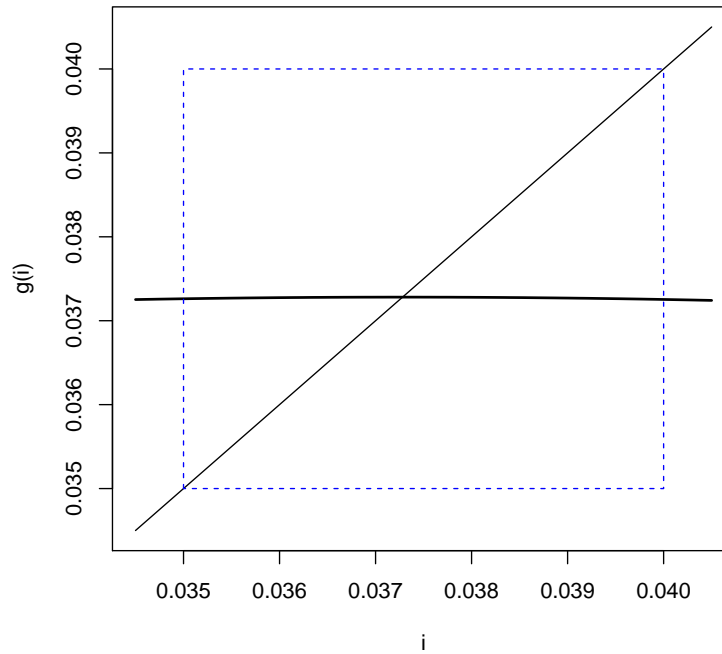


FIG. 11.7: Fonction  $g(i) = i - f(i)/f'(i)$  où  $f(i) = (1 - (1 + i)^{-10})/i - 8,2218$  pour  $0,035 \leq i \leq 0,040$  (ligne épaisse) et la droite  $y = i$  (ligne mince)

plus efficace pour résoudre ce genre de problème. La [figure 11.7](#) montre pourquoi elle est plus efficace que la méthode du point fixe : la fonction

$$g(i) = i - \frac{f(i)}{f'(i)}$$

avec  $f(i) = a_{\overline{n}|i} - k$  est beaucoup plus plate que la fonction

$$g(i) = \frac{1 - v^n}{k}$$

utilisée à l'[exemple 11.5](#).

On a défini en R une fonction `nr` similaire à la fonction `point fixe` mentionnée à la [section 11.3](#), sauf qu'on y a ajouté un argument pour la dérivée de la fonction  $f(x)$ .

► Exécuter le code informatique de la [section 11.8](#) correspondant à ce bloc de matière pour définir la fonction `nr` puis résoudre cet exemple. Prendre note comment l'on peut également utiliser la fonction `pointfixe` pour obtenir le taux de rendement.

La méthode de Newton–Raphson ne requiert donc que 2 itérations pour obtenir un résultat équivalent à celui obtenu par la méthode du point fixe avec la fonction `g` de l'[exemple 11.5](#) en 34 itérations.  $\square$

**Exemple 11.8.** Revisitons l'[exemple 11.3](#) où l'on cherche la racine de la fonction  $f(x) = x^3 + 4x^2 - 10$  dans l'intervalle  $[1, 2]$ . On a  $f'(x) = 3x^2 + 8x$ , d'où  $f'(x) \neq 0$  pour tout  $x \in [1, 2]$ . La méthode de Newton–Raphson dicte en définitive de rechercher le point fixe de la fonction

$$\begin{aligned} g(x) &= x - \frac{f(x)}{f'(x)} \\ &= x - \frac{x^3 + 4x^2 - 10}{3x^2 + 8x}. \end{aligned}$$

Celle-ci est présentée à la [figure 11.8](#). On constate que la procédure itérative de Newton–Raphson convergera rapidement vers une racine unique dans l'intervalle  $[1, 2]$ . En effet :

```
> f <- function(x) x^3 + 4*x^2 - 10
> fp <- function(x) 3*x^2 + 8 * x
> nr(f, fp, start = 1.5)
$root
[1] 1.36523

$nb.iter
[1] 3
```

Le lecteur attentif aura remarqué que la fonction `g` utilisée ici était la fonction `g5` dans l'[exemple 11.6](#), celle pour laquelle la convergence de la méthode du point fixe était la plus rapide. Mettre cette constatation en relation avec la première des remarques suivant l'[algorithme 11.3](#), à la page 43.  $\square$

**Exemple 11.9.** Cet exemple illustre l'importance d'utiliser, pour certaines fonctions, une valeur de départ près de la racine. La fonction  $f(x) = (4x - 7)/(x - 2)$  a une racine en  $x = 1,75$  et une asymptote en  $x = 2$ , ce qui cause quelques soucis. (Voir la [figure 11.9](#) pour un graphique de cette fonction.)

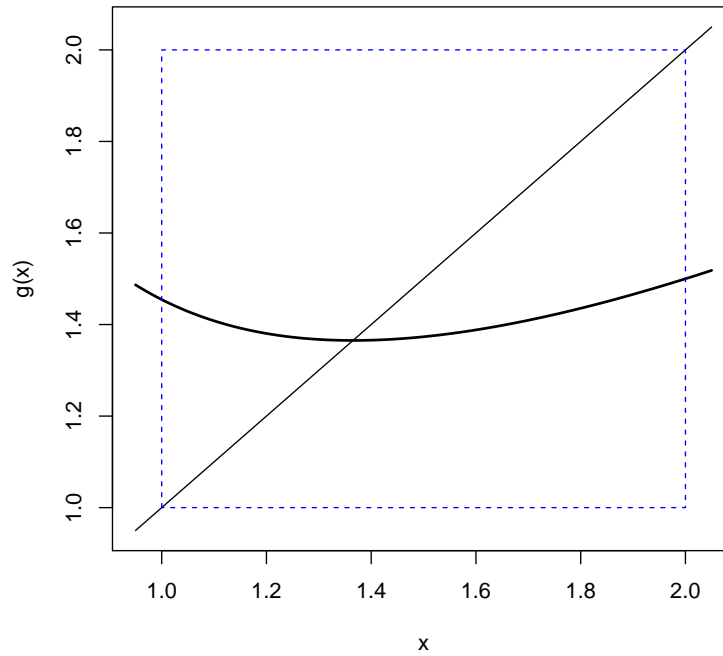


FIG. 11.8: Fonction  $g(x) = x - f(x)/f'(x)$ , où  $f(x) = x^3 + 4x^2 - 10$

On peut tenter de trouver numériquement la racine de la fonction  $f(x)$  avec l'algorithme de Newton-Raphson en utilisant

$$f(x) = \frac{4x-7}{x-2}$$

et

$$f'(x) = -\frac{1}{(x-2)^2},$$

ou encore directement par l'algorithme du point fixe avec

$$\begin{aligned} g(x) &= x - \frac{f(x)}{f'(x)} \\ &= 4x^2 - 14x + 14. \end{aligned}$$

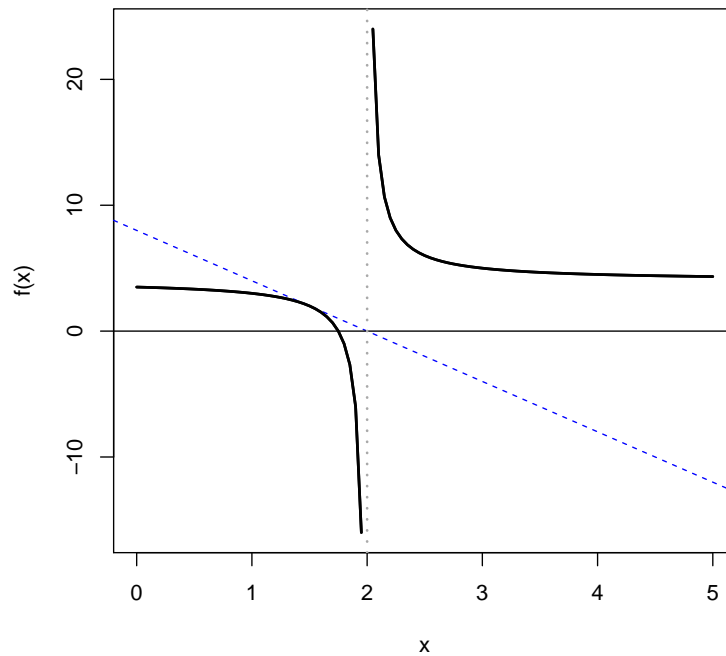


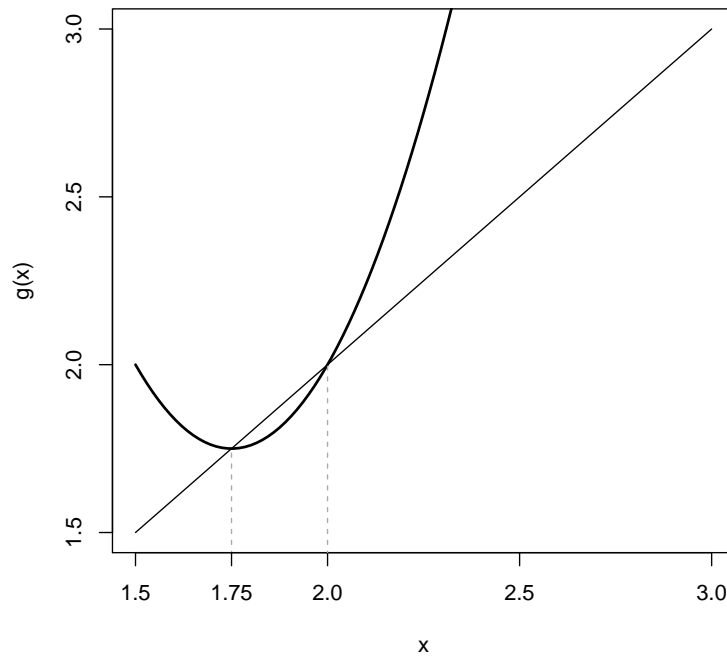
FIG. 11.9: Fonction  $f(x) = (4x-7)/(x-2)$  (ligne épaisse), asymptote (ligne pointillée) et tangente en  $x = 1,5$  (ligne traitillée)

Le graphique de la fonction  $g(x)$  ci-dessus se trouve à la [figure 11.10](#). On notera toutefois que le point fixe en  $x = 2$  n'est pas une solution acceptable dans le problème d'origine.

On constate également dans cette dernière figure que les hypothèses pour l'existence et l'unicité d'un point fixe — et, donc, pour la convergence de la méthode de Newton–Raphson — ne sont pas rencontrées. La valeur de départ utilisée aura donc un impact sur la réponse obtenue.

► Compléter cet exemple en exécutant le code informatique correspondant à la [section 11.8](#).

□

FIG. 11.10: Fonction  $g(x) = 4x^2 - 14x + 14$ 

Le principal inconvénient de la méthode de Newton–Raphson demeure le fait de devoir connaître la dérivée de la fonction  $f$ . Dans les cas où celle-ci s’avère difficile ou inefficace à calculer, on peut utiliser la *méthode de la sécante*. Avec cette méthode, plutôt que d’utiliser la tangente en un essai pour déterminer la valeur de l’essai suivant, on a recours à la sécante entre les deux essais précédents ; voir la [figure 11.11](#). Les valeurs des essais successifs sont alors données par :

$$x_n = x_{n-1} - \frac{f(x_{n-1})(x_{n-1} - x_{n-2})}{f(x_{n-1}) - f(x_{n-2})},$$

avec  $x_0, x_1$  des valeurs de départ.



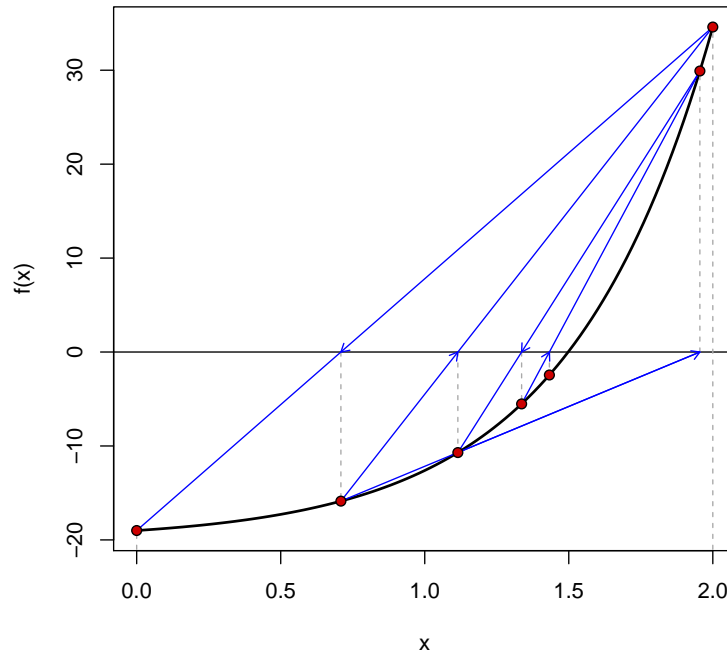


FIG. 11.11: Illustration de la méthode de la sécante. Cette méthode requiert deux essais initiaux. Dans le graphique, il s'agit, dans l'ordre, du point à l'extrême gauche et celui à l'extrême droite. Les flèches illustrent comment est déterminé l'essai suivant. On trace la sécante entre deux points et le prochain essai est la valeur en abscisse où cette sécante croise l'axe.

## 11.5 Fonctions d'optimisation dans Excel et R

Les méthodes de bisection, du point fixe, de Newton–Raphson et consorts permettent de résoudre des équations à une variable de la forme  $f(x) = 0$  ou  $g(x) = x$ . Il existe également des versions de ces méthodes pour les systèmes à plusieurs variables comme

$$f_1(x_1, x_2, x_3) = 0$$

$$f_2(x_1, x_2, x_3) = 0$$

$$f_3(x_1, x_2, x_3) = 0.$$

De tels systèmes d'équations surviennent plus souvent qu'autrement lors de l'optimisation d'une fonction. Par exemple, en recherchant le maximum ou le minimum d'une fonction  $f(x, y)$ , on souhaitera résoudre le système d'équations

$$\begin{aligned}\frac{\partial}{\partial x} f(x, y) &= 0 \\ \frac{\partial}{\partial y} f(x, y) &= 0.\end{aligned}$$

La grande majorité des suites logicielles comportent des outils d'optimisation de fonctions. Ce document passe en revue les fonctions disponibles dans Excel et R.

### 11.5.1 Solveur de Excel

Le principal outil d'optimisation utilisé dans Excel est le Solveur. La rubrique d'aide de cet outil est complète et on trouvera plusieurs exemples dans le classeur SOLVAMP.XLS livré avec Excel.

Le reste de cette section est dévolu à des fonctions d'optimisation de R.

### 11.5.2 Fonction uniroot

La fonction `uniroot` recherche la racine d'une fonction dans un intervalle. C'est donc la fonction de base pour trouver la solution (unique) de l'équation  $f(x) = 0$  dans un intervalle déterminé.

### 11.5.3 Fonction optimize

La fonction `optimize` recherche le minimum local (par défaut) ou le maximum local d'une fonction dans un intervalle donné.

### 11.5.4 Fonction nlm

La fonction `nlm` minimise une fonction non linéaire sur un nombre arbitraire de paramètres.

### 11.5.5 Fonction nlminb

La fonction `nlminb` est similaire à `nlm`, sauf qu'elle permet de spécifier des bornes inférieure ou supérieure pour les paramètres. Attention, toutefois : les arguments de la fonction ne sont ni les mêmes, ni dans le même ordre que ceux de `nlm`.

### 11.5.6 Fonction `optim`

La fonction `optim` est l'outil d'optimisation tout usage de R. À ce titre, la fonction est souvent utilisée par d'autres fonctions. Elle permet de choisir parmi plusieurs algorithmes d'optimisation différents et, selon l'algorithme choisi, de fixer des seuils minimum et/ou maximum aux paramètres à optimiser.

### 11.5.7 `polyroot`

En terminant, un mot sur `polyroot()`, qui n'est pas à proprement parler une fonction d'optimisation, mais qui pourrait être utilisée dans ce contexte. La fonction `polyroot` calcule toutes les racines (complexes) du polynôme  $\sum_{i=0}^n a_i x^i$ . Le premier argument est le vecteur des coefficients  $a_0, a_1, \dots, a_n$ , dans cet ordre.

► Des exemples d'utilisation des fonctions R ci-dessus sont fournis dans le code informatique de la [section 11.8](#).

## 11.6 Astuce Ripley

J'ai appris le truc suivant dans un message de Brian Ripley — un important développeur de R — posté dans les forums de discussion de R. Il m'a été très utile à de nombreuses reprises, alors je le dissémine.

Une application statistique fréquente de l'optimisation est la maximisation numérique d'une fonction de vraisemblance ou, plus communément, la minimisation de la log-vraisemblance négative

$$-l(\theta) = - \sum_{i=1}^n \ln f(x_i; \theta).$$

Les fonctions d'optimisation sont d'ailleurs illustrées dans ce contexte dans le code informatique de la [section 11.8](#).

En actuariat, nous utilisons principalement des lois de probabilité dont les paramètres sont strictement positifs. Or, en pratique, il n'est pas rare que les fonctions d'optimisation s'égarent dans les valeurs négatives des paramètres. La fonction de densité n'étant pas définie, la log-vraisemblance vaut alors NaN et cela peut faire complètement dérailler la procédure d'optimisation ou, à tout le moins, susciter des doutes sur la validité de la réponse.

Afin de pallier à ce problème, l'Astuce Ripley™ propose d'estimer non pas les paramètres de la loi eux-mêmes, mais plutôt leurs logarithmes. Si l'on définit

$\tilde{\theta} = \ln \theta$ , alors on peut écrire la fonction de log-vraisemblance ci-dessus sous la forme

$$-l(\tilde{\theta}) = -\sum_{i=1}^n \ln f(x_i; e^{\tilde{\theta}}).$$

Dès lors,  $\tilde{\theta}$  (qui peut représenter un ou plusieurs paramètres) demeure valide sur tout l'axe des réels, ce qui permet d'éviter bien des soucis de nature numérique lors de la minimisation de  $-l(\tilde{\theta})$ .

Évidemment, le résultat de l'optimisation est l'estimateur du maximum de vraisemblance de  $\tilde{\theta}$ . Il faudra donc veiller à faire la transformation inverse pour retrouver l'estimateur de  $\theta$ .

L'utilisation de l'astuce est illustrée à la [section 11.8](#).

## 11.7 Outils additionnels

Les packages disponible sur CRAN fournissent plusieurs autres outils d'optimisation pour R. Pour un bon résumé des options disponibles, consulter la *CRAN Task View* consacrée à l'optimisation :

<http://cran.r-project.org/web/views/Optimization.html>

## 11.8 Code informatique

```
###
### MÉTHODE DE BISSECTION
###

## Fonction pour trouver la solution de 'FUN'(x) = 0 par la
## méthode de bisection, où 'FUN' est une fonction continue
## sur l'intervalle ['lower', 'upper'].
bisection <- function(FUN, lower, upper, TOL = 1E-6,
                      MAX.ITER = 100, echo = FALSE)
{
  ## Cas triviaux où une borne est la solution
  if (identical(FUN(lower), 0))
    return(lower)
  if (identical(FUN(upper), 0))
    return(upper)

  ## Bornes de départ inadéquates
  if (FUN(lower) * FUN(upper) > 0)
    stop('FUN(lower) and FUN(upper) must be of opposite signs')
```

```
x <- lower
i <- 1

repeat
{
  xt <- x
  x <- (lower + upper)/2
  fx <- FUN(x)

  if (echo)
    print(c(lower, upper, x, fx))

  if (abs(x - xt)/abs(x) < TOL)
    break

  if (MAX.ITER < (i <- i + 1))
    stop('Maximum number of iterations reached without convergence')

  if (FUN(lower) * fx > 0)
    lower <- x
  else
    upper <- x
}
list(root = x, nb.iter = i)
}
```

### ## EXEMPLE 11.3

##

## *La fonction dont on cherche la racine*

```
f <- function(x) x^3 + 4*x^2 - 10
```

## *Un graphique de la fonction permet de vérifier qu'elle ne*

## *possède qu'une seule racine sur un intervalle donné ainsi*

## *que des valeurs de départ pour l'algorithme.*

```
curve(f(x), xlim = c(0, 3), lwd = 2) # graphique sur [0, 3]
```

```
abline(h = 0) # axe des abscisses
```

## *Résolution. Nous utilisons a = 1 et b = 2 comme valeurs de*

## *départ. Nous pourrions choisir des valeurs de départ plus*

## *près de la racine, ce qui accélérerait la résolution.*

```
bisection(f, lower = 1, upper = 2, TOL = 1E-5, echo = TRUE)
```

### ## EXEMPLE 11.4

##

```

## La fonction dont on cherche la racine
f <- function(x) (1 - (1 + x)^(-10))/x - 8.2218

## Graphique de la fonction sur un intervalle "raisonnable"
## dans lequel on peut pressentir que se trouvera la réponse.
## Ici, on y va pour un taux d'intérêt se trouvant entre 1 %
## et 10 %.
curve(f(x), xlim = c(0.01, 0.10), lwd = 2)
abline(h = 0)

## Le graphique permet de déterminer à l'oeil que la racine se
## trouve entre 3 % et 4 %. Résolution avec la fonction
## 'bissection'.
bissection(f, lower = 0.03, upper = 0.04)

####
#### MÉTHODE DU POINT FIXE
####

## Fonction pour trouver la solution de 'FUN'(x) = x par la
## méthode du point fixe à partir d'un essai initial 'start'.
pointfixe <- function(FUN, start, TOL = 1E-6, MAX.ITER = 100,
                      echo = FALSE)
{
  x <- start

  if (echo)
    expr <- expression(print(xt <- x))
  else
    expr <- expression(xt <- x)

  i <- 0

  repeat
  {
    eval(expr)

    x <- FUN(xt)

    if (abs(x - xt)/abs(x) < TOL)
      break

    if (MAX.ITER < (i <- i + 1))
      stop('Maximum number of iterations reached without convergence')
  }
}

```

```

    list(fixed.point = x, nb.iter = i)
}

## EXEMPLE 11.5
##
## La fonction dont on cherche le point fixe.
g <- function(x) (1 - (1 + x)^(-10))/8.2218

## Un graphique de la fonction permet de vérifier que les
## conditions pour qu'il existe un point fixe unique dans
## l'intervalle [0,035, 0,040] sont satisfaites. Cependant, la
## forme de la courbe (pente près de 1) nous indique que la
## convergence sera relativement lente.
f <- function(x) x          # pour tracer la droite y = x
lim <- c(0.035, 0.040)      # intervalle [a, b]
curve(g, xlim = lim, ylim = lim, lwd = 2,
      xlab = "i", ylab = "g(i)")
curve(f, add = TRUE)
polygon(rep(lim, each = 2), c(lim, rev(lim)),
      lty = "dashed", border = "blue")

## Résolution avec la fonction 'pointfixe'.
pointfixe(g, start = 0.0375, echo = TRUE)

## EXEMPLE 11.6
##
## On cherche la racine de  $f(x) = x^3 + 4 * x^2 - 10$  en
## exprimant le problème sous forme de point fixe. Les cinq
## fonctions ci-dessous sont toutes algébriquement
## équivalentes, c'est-à-dire que  $g(x) = x$  lorsque  $f(x) = 0$ .
g1 <- function(x) x - x^3 - 4 * x^2 + 10
g2 <- function(x) sqrt(10/x - 4*x)
g3 <- function(x) sqrt(10 - x^3)/2
g4 <- function(x) sqrt(10/(4 + x))
g5 <- function(x) x - (x^3 + 4*x^2 - 10)/(3*x^2 + 8*x)

## Si les fonctions sont algébriquement équivalentes, elles ne
## le sont pas numériquement devant la méthode du point fixe.
## Ainsi, avec la fonction 'g1', la procédure diverge.
## (Remarque: la fonction 'pointfixe' ne prévoit pas de message
## d'erreur pour ce cas. Qu'ajouteriez-vous à la fonction?)
pointfixe(g1, 1.5)
pointfixe(g1, 1.5, echo = TRUE) # plus évident ainsi

## Avec la fonction 'g2', la procédure s'arrête lorsqu'il faut

```

```

## calculer la racine carrée d'un nombre négatif.
pointfixe(g2, 1.5)

## Avec les trois autres fonctions, la méthode du point fixe
## est extrêmement rapide et précise. Une rapide analyse des
## graphiques fournis dans le chapitre nous aurait permis de
## déterminer avec quelle fonction la convergence serait la
## plus rapide. En effet, c'est la fonction 'g5' qui a la
## pente la plus faible près de son point fixe.
pointfixe(g3, 1.5)
pointfixe(g4, 1.5)
pointfixe(g5, 1.5)

####
#### MÉTHODE DE NEWTON-RAPHSON
####

## Fonction pour trouver la solution de 'FUN'(x) = x par la
## méthode de Newton-Raphson à partir de sa dérivée 'FUNp' et
## d'un essai initial 'start'.
##
## On ajoute une amélioration par rapport aux fonctions
## 'bissection' et 'pointfixe', soit la possibilité de passer
## des arguments additionnels aux fonctions 'FUN' et 'FUNp'
## via l'argument '...'
nr <- function(FUN, FUNp, start, TOL = 1E-6,
               MAX.ITER = 100, echo = FALSE, ...)
{
  x <- start

  if (echo)
    expr <- expression(print(xt <- x))
  else
    expr <- expression(xt <- x)

  i <- 0

  repeat
  {
    eval(expr)

    x <- xt - FUN(xt, ...)/FUNp(xt, ...)

    if (abs(x - xt)/abs(x) < TOL)
      break
  }
}

```



```

    if (MAX.ITER < (i <- i + 1))
      stop('Maximum number of iterations reached without convergence')
  }
  list(root = x, nb.iter = i)
}

```

#### ## EXEMPLE 11.7

##

*## On répète l'exemple 11.5 avec la méthode de Newton-Raphson.  
 ## La convergence devrait être plus rapide qu'avec la méthode  
 ## du point fixe parce que la fonction  $g(i) = i - f(i)/f'(i)$   
 ## est plus plate que celle utilisée à l'exemple 11.5.  
 ## Premièrement, la fonction dont on cherche la racine.*

```
f <- function(i) (1 - (1 + i)^(-10))/i - 8.2218
```

*## Sa dérivée.*

```
fp <- function(i) (10 * i * (1 + i)^(-11) + (1 + i)^(-10) - 1)/i^2
```

*## Résolution avec la fonction 'nr'.*

```
nr(f, fp, start = 0.0375, echo = TRUE)
```

*## À noter que si la fonction g est définie adéquatement, on  
 ## peut de manière tout aussi équivalente utiliser la fonction  
 ## 'pointfixe' pour effectuer les itérations de la méthode de  
 ## Newton-Raphson.*

```
g <- function(i) i - f(i)/fp(i)
pointfixe(g, 0.0375, echo = TRUE)
```

#### ## EXEMPLE 11.9

##

*## Les fonctions f(x), f'(x) et g(x) définies dans le texte de  
 ## l'exemple.*

```
f <- function(x) ifelse(x == 2, NA, (4 * x - 7)/(x - 2))
fp <- function(x) ifelse(x == 2, NA, -1/(x - 2)^2)
g <- function(x) 4 * x^2 - 14 * x + 14
```

*## On vérifie que, étant donné la forme des fonctions 'f' et  
 ## 'g', la valeur de départ utilisée dans les méthodes de  
 ## Newton-Raphson ou du point fixe a un impact sur la réponse  
 ## obtenue. Avec une valeur de départ 'start' = 1,625 < 1,75,  
 ## la convergence se fait vers la bonne racine.*

```
nr(f, fp, 1.625, echo = TRUE)
pointfixe(g, 1.625)
```

```

## On peut vérifier sur le graphique de la fonction 'g'
## qu'avec une valeur de départ entre 1,75 et 2, la procédure
## itérative convergera aussi vers la bonne réponse.
nr(f, fp, 1.875, echo = TRUE)
pointfixe(g, 1.875, echo=TRUE)

## La tangente en  $x = 1,5$  tracée sur le graphique de la
## fonction 'f' montre que 'start' = 1,5 constitue un mauvais
## choix car  $g(1,5) = 2$ . La procédure itérative converge donc
## vers une valeur non admissible.
nr(f, fp, 1.5)          # division par 0
pointfixe(g, 1.5)       # point fixe non admissible

## Pour toute valeur de départ supérieure à 2, la procédure
## itérative diverge. On peut vérifier ce fait dans les
## graphiques de 'f' et de 'g'.
nr(f, fp, 3, echo = TRUE)
pointfixe(g, 3, echo = TRUE)

###
### FONCTIONS D'OPTIMISATION DE R
###

## FONCTION 'uniroot'
##
## La fonction 'uniroot' recherche la racine d'une fonction
## 'f' dans un intervalle spécifié soit comme une paire de
## valeurs dans un argument 'interval', soit via des arguments
## 'lower' et 'upper'.
##
## On calcule la solution de l'équation  $x - 2^{(-x)} = 0$  dans
## l'intervalle  $[0, 1]$ .
f <- function(x) x - 2^(-x)      # fonction
uniroot(f, c(0, 1))             # appel simple
uniroot(f, lower = 0, upper = 1) # équivalent

## On peut aussi utiliser 'uniroot' avec une fonction anonyme.
uniroot(function(x) x - 2^(-x), lower = 0, upper = 1)

## FONCTION 'optimize'
##
## On cherche le maximum local de la densité d'une loi bêta
## dans l'intervalle (0, 1), son domaine de définition. (Ce
## problème est facile à résoudre explicitement.)
##

```

```
## Les arguments de 'optimize' sont essentiellement les mêmes
## que ceux de 'uniroot'. Ici, on utilise aussi l'argument
## '...' pour passer les paramètres de la loi bêta à 'dbeta'.
##
## Par défaut, la fonction recherche un minimum. Il faut donc
## lui indiquer de rechercher plutôt un maximum.
optimize(dbeta, interval = c(0, 1), maximum = TRUE,
         shape1 = 3, shape2 = 2)

## On pourrait aussi avoir recours à une fonction auxiliaire.
## Moins élégant et moins flexible.
f <- function(x) dbeta(x, 3, 2)
optimize(f, lower = 0, upper = 1, maximum = TRUE)

## FONCTION 'nlm'
##
## Pour la suite, nous allons donner des exemples
## d'utilisation des fonctions d'optimisation dans un contexte
## d'estimation des paramètres d'une loi gamma par la méthode
## du maximum de vraisemblance.
##
## On commence par se donner un échantillon aléatoire de la
## loi. Évidemment, pour ce faire nous devons connaître les
## paramètres de la loi. C'est un exemple fictif.
set.seed(1) # toujours le même échantillon
x <- rgamma(10, 5, 2)

## Les estimateurs du maximum de vraisemblance des paramètres
## 'shape' et 'rate' de la loi gamma sont les valeurs qui
## maximisent la fonction de vraisemblance
##
## prod(dgamma(x, shape, rate))
##
## ou, de manière équivalente, qui minimisent la fonction de
## log-vraisemblance négative
##
## -sum(log(dgamma(x, shape, rate))).
##
## On remarquera au passage que les fonctions de calcul de
## densités de lois de probabilité dans R ont un argument
## 'log' qui, lorsque TRUE, retourne la valeur du logarithme
## (naturel) de la densité de manière plus précise qu'en
## prenant le logarithme après coup. Ainsi, pour faire le
## calcul ci-dessus, on optera plutôt, pour l'expression
##
```

```

##      -sum(dgamma(x, shape, rate, log = TRUE))
##
## La fonction 'nlm' suppose que la fonction à optimiser
## passée en premier argument a elle-même comme premier
## argument le vecteur 'p' des paramètres à optimiser. Le
## second argument de 'nlm' est un vecteur de valeurs de
## départ, une pour chaque paramètre.
##
## Ainsi, pour trouver les estimateurs du maximum de
## vraisemblance avec la fonction 'nlm' pour l'échantillon
## ci-dessus, on doit d'abord définir une fonction auxiliaire
## conforme aux attentes de 'nlm' pour calculer la fonction de
## log-vraisemblance (à un signe près).
f <- function(p, x) -sum(dgamma(x, p[1], p[2], log = TRUE))

## L'appel de 'nlm' est ensuite tout simple. Remarquer comment
## on passe notre échantillon aléatoire (contenu dans l'objet
## 'x') comme second argument à 'f' via l'argument '...' de
## 'nlm'. Le fait que l'argument de 'f' et l'objet contenant
## les valeurs portent le même nom est sans importance. R sait
## faire la différence entre l'un et l'autre.
nlm(f, c(1, 1), x = x)

## === ASTUCE RIPLEY ===
## L'optimisation ci-dessus a généré des avertissements? C'est
## parce que la fonction d'optimisation s'est égarée dans les
## valeurs négatives, alors que les paramètres d'une gamma
## sont strictement positifs. Cela arrive souvent en pratique
## et cela peut faire complètement dérailler la procédure
## d'optimisation (c'est-à-dire: pas de convergence).
##
## L'Astuce Ripley consiste à pallier à ce problème en
## estimant plutôt les logarithmes des paramètres. Pour ce
## faire, il s'agit de réécrire la log-vraisemblance comme une
## fonction du logarithme des paramètres, mais de la calculer
## avec les véritables paramètres.
f2 <- function(logp, x)
{
  p <- exp(logp)          # retour aux paramètres originaux
  -sum(dgamma(x, p[1], p[2], log = TRUE))
}
nlm(f2, c(0, 0), x = x)

## Les valeurs obtenues ci-dessus sont toutefois les
## estimateurs des logarithmes des paramètres de la loi gamma.

```

```

## On retrouve les estimateurs des paramètres en prenant
## l'exponentielle des réponses.
exp(nlm(f2, c(0, 0), x = x)$estimate)
## =====

## FONCTION 'nlminb'
##
## L'utilisation de la fonction 'nlminb' peut s'avérer
## intéressante dans notre contexte puisque l'on sait que les
## paramètres d'une loi gamma sont strictement positifs.
nlminb(c(1, 1), f, x = x, lower = 0, upper = Inf)

### FONCTION 'optim'
##
## La fonction 'optim' est très puissante, mais requiert aussi
## une bonne dose de prudence. Ses principaux arguments sont:
##
## par: un vecteur contenant les valeurs initiales des
##      paramètres;
## fn: la fonction à minimiser. Le premier argument de fn
##     doit être le vecteur des paramètres.
##
## Comme pour les autres fonctions étudiées ci-dessus, on peut
## passer des arguments à 'fn' (les données, par exemple) par
## le biais de l'argument '...' de 'optim'.
optim(c(1, 1), f, x = x)

## FONCTION 'polyroot'
##
## Racines du polynôme  $x^3 + 4x^2 - 10$ . Les réponses sont
## données sous forme de nombre complexe. Utiliser les
## fonctions 'Re' et 'Im' pour extraire les parties réelles et
## imaginaires des nombres, respectivement.
polyroot(c(-10, 0, 4, 1))      # racines
Re(polyroot(c(-10, 0, 4, 1))) # parties réelles
Im(polyroot(c(-10, 0, 4, 1))) # parties imaginaires

```

## 11.9 Exercices

**11.1** En vous basant sur les fonctions bisection, pointfixe et nr présentées dans le code de la [section 11.8](#), écrire une fonction R pour effectuer les calculs de l'algorithme de la sécante. Outre les arguments communs à toutes les fonctions que sont le niveau de tolérance  $\varepsilon$ , le nombre maximal d'itérations  $N_{\max}$  et une valeur booléenne spécifiant si les valeurs successives des itérations

doivent être affichées à l'écran, la fonction doit compter les arguments  $f(x)$ ,  $x_0$  et  $x_1$ .

**11.2** Trouver la solution des équations suivantes par les méthodes de bisection, de Newton–Raphson et de la sécante.

- a)  $x^3 - 2x^2 - 5 = 0$  pour  $1 \leq x \leq 4$
- b)  $x^3 + 3x^2 - 1 = 0$  pour  $-4 \leq x \leq 0$
- c)  $x - 2^{-x} = 0$  pour  $0 \leq x \leq 1$
- d)  $e^x + 2^{-x} + 2 \cos x - 6 = 0$  pour  $1 \leq x \leq 2$
- e)  $e^x - x^2 + 3x - 2 = 0$  pour  $0 \leq x \leq 1$

**11.3** Déterminer la valeur numérique de  $\sqrt{2}$  à l'aide de la méthode de bisection dans l'intervalle  $[0, 2]$  avec 10 itérations. Comparer avec la vraie valeur.

**11.4** Soit  $\{x_n\}$  une suite définie par

$$x_n = \sum_{k=1}^n \frac{1}{k}.$$

Démontrer que  $\lim_{n \rightarrow \infty} (x_n - x_{n-1}) = 0$ , mais que la suite diverge néanmoins. Ceci illustre que l'erreur absolue peut être un mauvais critère d'arrêt dans les méthodes numériques.

**11.5** Considérer la fonction  $g(x) = 2^{-x}$  sur l'intervalle  $[0, 1]$ .

- a) Vérifier si les hypothèses du [théorème 11.1](#) quant à l'existence et l'unicité d'un point fixe dans  $[0, 1]$  sont satisfaites.
- b) Déterminer graphiquement l'existence et l'unicité d'un point fixe de  $g(x)$  dans  $[0, 1]$  puis, le cas échéant, calculer ce point fixe.

**11.6** Vérifier que les cinq fonctions  $g_1, \dots, g_5$  de l'[exemple 11.6](#) ont toutes un point fixe en  $x^*$  lorsque  $f(x^*) = 0$ , où  $f(x) = x^3 + 4x^2 - 10$ .

**11.7** Soit  $g(x) = 4x^2 - 14x + 14$ . Pour quels intervalles de valeurs de départ la procédure de point fixe converge-t-elle et diverge-t-elle ?

**11.8** Les trois fonctions ci-dessous sont toutes des candidates pour faire l'approximation, par la méthode du point fixe, de  $\sqrt[3]{21}$  :

$$g_1(x) = \frac{20x + 21x^{-2}}{21}$$

$$g_2(x) = x - \frac{x^3 - 21}{3x^2}$$

$$g_3(x) = x - \frac{x^4 - 21x}{x^2 - 21}.$$

Classer ces fonctions en ordre décroissant de vitesse de convergence de l'algorithme du point fixe. *Astuce* : comparer les valeurs des dérivées autour du point fixe.

**11.9** Démontrer, à l'aide du théorème du point fixe, que la fonction  $g(x) = 2^{-x}$  possède un point fixe unique dans l'intervalle  $[\frac{1}{3}, 1]$ . Calculer par la suite ce point fixe à l'aide de la fonction point fixe.

**11.10** Vérifier graphiquement que le [théorème 11.1](#) demeure valide si la condition  $|g'(x)| \leq k < 1$  est remplacée par  $g'(x) \leq k < 1$ .

**11.11** Utiliser la méthode de Newton–Raphson pour trouver le point sur la courbe  $y = x^2$  le plus près du point  $(1, 0)$ . *Astuce* : minimiser la distance entre le point  $(1, 0)$  et le point  $(x, x^2)$ .

**11.12** Le taux de rendement interne d'une série de flux financiers  $\{CF_t\}$  est le taux  $i$  tel que

$$\sum_{t=0}^n \frac{CF_t}{(1+i)^t} = 0.$$

Écrire une fonction R permettant de calculer le taux de rendement interne d'une série de flux financiers quelconque à l'aide de la méthode de Newton–Raphson. Les arguments de la fonction sont un vecteur CF et un scalaire erreur.max. Au moins un élément de CF doit être négatif pour représenter une sortie de fonds. Dans tous les cas, utiliser  $i = 0,05$  comme valeur de départ.

**11.13** Refaire l'exercice 11.12 en VBA en composant une fonction accessible dans une feuille Excel et comparer le résultat avec la fonction Excel TRI ( ).

**11.14** a) Écrire une fonction R pour estimer par la méthode du maximum de vraisemblance le paramètre  $\theta$  d'une loi gamma de paramètre de forme  $\alpha = 3$  et de paramètre d'échelle  $\theta = 1/\lambda$  — de moyenne  $3\theta$ , donc — à partir d'un échantillon aléatoire  $x_1, \dots, x_n$ . *Astuce* : maximiser la fonction de log-vraisemblance plutôt que la fonction de vraisemblance.

b) Simuler 20 observations d'une loi gamma avec paramètre de forme  $\alpha = 3$  et moyenne 3000, puis estimer le paramètre d'échelle  $\theta$  par le maximum de vraisemblance.

**11.15** Refaire l'exercice 11.14 à l'aide d'une routine VBA. Comparer votre réponse avec celle obtenue à l'aide du Solveur de Excel.

**11.16** À l'aide de la méthode du point fixe, trouver le taux d'intérêt  $i$  tel que

$$a_{10|i}^{(12)} = \frac{1 - (1+i)^{-10}}{i^{(12)}} = 8,$$

où

$$\left(1 + \frac{i^{(12)}}{12}\right)^{12} = 1 + i.$$

Comparer les résultats obtenus avec la méthode de Newton–Raphson.

**11.17** Trouver la solution des équations suivantes à l'aide des fonctions R appropriées.

- a)  $x^3 - 2x^2 - 5 = 0$  pour  $1 \leq x \leq 4$
- b)  $x^3 + 3x^2 - 1 = 0$  pour  $-4 \leq x \leq 0$
- c)  $x - 2^{-x} = 0$  pour  $0 \leq x \leq 1$
- d)  $e^x + 2^{-x} + 2 \cos x - 6 = 0$  pour  $1 \leq x \leq 2$
- e)  $e^x - x^2 + 3x - 2 = 0$  pour  $0 \leq x \leq 1$

**11.18** Les fonctions de densité de probabilité et de répartition de la distribution de Pareto sont données à l'exercice 6.3 de la partie I. Calculer les estimateurs du maximum de vraisemblance des paramètres de la Pareto à partir d'un échantillon aléatoire obtenu par simulation avec la commande

```
> x <- lambda * (runif(100)^(-1/alpha) - 1)
```

pour des valeurs de alpha et lambda choisies.

## Réponses

**11.7** Converge pour  $x_0 \in [1,5,2]$  ; diverge pour  $x_0 \leq 1,5$  et  $x_0 > 2$

**11.8**  $g_2, g_1$  ;  $g_3$  ne converge pas

**11.11** 0,589755

**11.16** 0,0470806



# 12 Intégration numérique

## Objectifs du chapitre

- ▶ Avoir une connaissance générale de ce qu'est un polynôme d'interpolation de Lagrange et comment il peut servir dans le contexte de l'intégration numérique.
- ▶ Calculer la valeur approximative d'une intégrale définie à l'aide des méthodes du point milieu, du trapèze, de Simpson et de Simpson 3/8.
- ▶ Développer les formules d'approximation composées de chacune des méthodes d'intégration numériques présentées dans le chapitre.

Il n'est pas rare de devoir calculer l'intégrale d'une fonction qui n'admet pas de primitive sous forme explicite ou dont la primitive est très difficile à calculer. Dans de tels cas, l'intégration numérique permet d'obtenir une approximation — parfois excellente — de l'intégrale cherchée.

Les méthodes d'approximation qui seront étudiées dans ce chapitre reposent toutes sur le remplacement de la fonction à intégrer par une «bonne» approximation sur un intervalle donné. Cette approximation devra évidemment être simple à intégrer, autrement aucun gain n'est réalisé. Or, quelles fonctions sont simples à intégrer sinon les polynômes ?

## 12.1 Polynômes d'interpolation de Lagrange

Un résultat connu en analyse mathématique sous le nom de Théorème d'approximation de Weierstrass établit que l'on peut faire l'approximation de toute fonction continue par un polynôme de degré suffisant. Formellement, soit

$$P_n(x) = a_0 + a_1x + \dots + a_{n-1}x^{n-1} + a_nx^n$$

un polynôme de degré  $n$ . Si  $f$  est une fonction continue sur un intervalle  $[a, b]$ , alors pour tout  $\varepsilon > 0$  il existe un polynôme  $P(x)$  tel que

$$|f(x) - P(x)| < \varepsilon$$

pour tout  $x \in [a, b]$ .

La notion d'approximation d'une fonction par un polynôme est beaucoup utilisée en analyse numérique puisque les polynômes sont simples à dériver et à intégrer, et que leurs dérivées et intégrales sont elles-mêmes des polynômes.

Un type de polynôme d'approximation que vous connaissez déjà est le polynôme (ou développement) de Taylor autour d'un point  $x_0$ . Cependant, les polynômes de Taylor sont conçus pour être précis autour de  $x_0$  et non sur tout un intervalle. Dans le contexte de l'intégration numérique, où nous voudrions remplacer la fonction à intégrer par un polynôme, nous aurons besoin d'une bonne approximation sur tout le domaine d'intégration.

C'est là qu'interviennent les polynômes d'interpolation de Lagrange. Soit  $x_0, x_1, \dots, x_n$  un ensemble de  $n + 1$  points distincts et  $f$  une fonction qui passe par ces points (ou *nœuds*). Le polynôme d'interpolation de Lagrange de degré  $n$  de la fonction  $f$  est

$$P_n(x) = f(x_0)L_0(x) + f(x_1)L_1(x) + \dots + f(x_n)L_n(x) = \sum_{k=0}^n f(x_k)L_k(x),$$

où

$$\begin{aligned} L_k(x) &= \frac{(x - x_0)(x - x_1) \cdots (x - x_{k-1})(x - x_{k+1}) \cdots (x - x_n)}{(x_k - x_0)(x_k - x_1) \cdots (x_k - x_{k-1})(x_k - x_{k+1}) \cdots (x_k - x_n)} \\ &= \prod_{\substack{i=0 \\ i \neq k}}^n \frac{(x - x_i)}{(x_k - x_i)}. \end{aligned}$$

Remarquer que le terme  $(x - x_k)$  n'apparaît pas au numérateur de la fonction  $L_k(x)$ .

On peut démontrer (voir, par exemple, [Burden et Faires, 2011](#)) que le polynôme  $P(x)$  peut être aussi près que l'on veut de la fonction  $f(x)$ .

**Exemple 12.1.** Le polynôme d'interpolation de Lagrange de premier degré de la fonction  $f$  est

$$\begin{aligned} P_1(x) &= f(x_0)L_0(x) + f(x_1)L_1(x) \\ &= f(x_0) \frac{(x - x_1)}{(x_0 - x_1)} + f(x_1) \frac{(x - x_0)}{(x_1 - x_0)}. \end{aligned}$$

Le polynôme de second degré est, quant à lui :

$$\begin{aligned} P_2(x) &= f(x_0)L_0(x) + f(x_1)L_1(x) + f(x_2)L_2(x) \\ &= f(x_0)\frac{(x-x_1)(x-x_2)}{(x_0-x_1)(x_0-x_2)} + f(x_1)\frac{(x-x_0)(x-x_2)}{(x_1-x_0)(x_1-x_2)} \\ &\quad + f(x_2)\frac{(x-x_0)(x-x_1)}{(x_2-x_0)(x_2-x_1)}. \end{aligned}$$

Supposons que  $f(x) = 1/x$  et que les polynômes d'interpolation passent par les nœuds  $x_0 = 2$ ,  $x_1 = 2,75$  et  $x_2 = 4$ . On a alors

$$\begin{aligned} P_1(x) &= \frac{1}{2} \frac{(x-2,75)}{(2-2,75)} + \frac{1}{2,75} \frac{(x-2)}{(2,75-2)} \\ &= -\frac{2}{3}(x-2,75) + \frac{16}{33}(x-2) \\ &= -\frac{6}{33}x + \frac{171}{198} \end{aligned}$$

et

$$\begin{aligned} P_2(x) &= \frac{1}{2} \frac{(x-2,75)(x-4)}{(2-2,75)(2-4)} + \frac{1}{2,75} \frac{(x-2)(x-4)}{(2,75-2)(2,75-4)} \\ &\quad + \frac{1}{4} \frac{(x-2)(x-2,75)}{(4-2)(4-2,75)} \\ &= \frac{1}{3}(x-2,75)(x-4) - \frac{64}{165}(x-2)(x-4) + \frac{1}{10}(x-2)(x-2,75) \\ &= \frac{1}{22}x^2 - \frac{35}{88}x + \frac{49}{44}. \end{aligned}$$

Des approximations de  $f(3) = 1/3$  sont donc  $P_1(3) \approx 0,31818$  et  $P_2(3) \approx 0,32955$ .  $\square$

► La fonction `poly.calc` du package **polynom** permet de calculer le polynôme d'interpolation de Lagrange de degré  $n$  dans R. Exécuter le code informatique de la [section 12.7](#) correspondant à ce bloc de matière pour vérifier les calculs de l'[exemple 12.1](#).

## 12.2 Principes généraux d'intégration numérique

Supposons que l'on cherche à calculer  $\int_a^b f(x) dx$ . Toutes les méthodes d'intégration numérique reposent, en premier lieu, sur le découpage du domaine  $(a, b)$

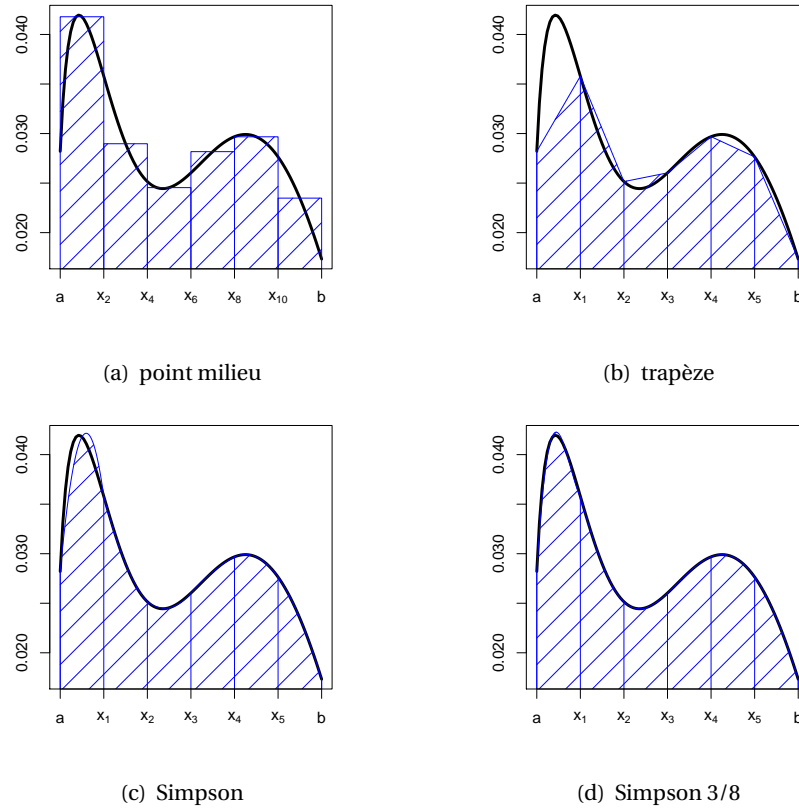


FIG. 12.1: Comparaison de quatre méthodes d'intégration numérique

en  $n$  intervalles. Cela permet d'évaluer l'intégrale comme une somme d'intégrales sur chacun de ces intervalles :

$$\int_a^b f(x) dx = \sum_{j=0}^{n-1} \int_{x_j}^{x_{j+1}} f(x) dx.$$

Les diverses méthodes d'intégration numérique diffèrent par la suite essentiellement par l'approximation de la fonction  $f(x)$  sur l'intervalle  $(x_j, x_{j+1})$  utilisée afin de rendre l'intégrale  $\int_{x_j}^{x_{j+1}} f(x) dx$  simple à calculer.

La [figure 12.1](#) illustre les procédures d'approximation de quatre méthodes courantes d'intégration numérique.

Les formules d'intégration numérique présentées dans la suite découlent de la procédure suivante : le domaine d'intégration  $(a, b)$  est découpé en  $n$  intervalles de longueur égale. Chacun de ces intervalles est à son tour divisé en  $m$  sous-intervalles

de longueur égale, pour un total de  $nm$  points. On a alors

$$\begin{aligned} \int_a^b f(x) dx &= \int_{x_0}^{x_m} f(x) dx + \int_{x_m}^{x_{2m}} f(x) dx + \cdots + \int_{x_{(n-1)m}}^{x_{nm}} f(x) dx \\ &= \sum_{j=0}^{n-1} \int_{x_{jm}}^{x_{(j+1)m}} f(x) dx, \end{aligned}$$

où  $x_0 = a$  et  $x_{nm} = b$ . L'approximation numérique se trouve dans l'évaluation de l'intégrale du côté droit de la dernière équation. Pour toutes les méthodes étudiées dans ce chapitre sauf la méthode du point milieu, la fonction  $f$  sera remplacée par un polynôme d'interpolation de Lagrange.

Afin de ne pas alourdir inutilement la notation, les formules d'approximation de l'intégrale

$$\int_{x_{jm}}^{x_{(j+1)m}} f(x) dx$$

seront présentées pour le cas  $j = 0$  seulement dans les sections suivantes.

## 12.3 Méthode du point milieu

La méthode du point milieu est la plus simple et la plus intuitive méthode d'intégration numérique. Les intervalles sont divisés en  $m = 2$  parties et la valeur de la fonction  $f$  sur l'intervalle  $(x_0, x_2)$  est estimée par  $f(x_1)$  (12.2(a)). Ainsi, on a l'approximation

$$\int_{x_0}^{x_2} f(x) dx \approx 2h f(x_1), \quad (12.1)$$

où  $h = x_2 - x_1 = x_1 - x_0$ .

La formule composée pour l'approximation par la méthode du point milieu de  $\int_a^b f(x) dx$  avec  $n$  intervalles est

$$\int_a^b f(x) dx \approx 2h \sum_{j=0}^{n-1} f(x_{2j+1}) \quad (12.2)$$

où  $h = (b - a)/(2n)$  et  $x_j = a + jh$ .

## 12.4 Méthode du trapèze

Les trois prochaines méthodes d'intégration numérique sont basées sur l'approximation de la fonction  $f$  sur un intervalle par un polynôme d'interpolation de Lagrange de degré  $m$ .

La méthode du trapèze utilise un polynôme du premier degré ( $m = 1$ ) pour faire l'approximation de la valeur de  $f(x)$  sur l'intervalle  $(x_0, x_1)$ , ce qui est équivalent à une simple interpolation linéaire (12.2(b)). On a donc

$$\int_{x_0}^{x_1} f(x) dx \approx \frac{h}{2} [f(x_0) + f(x_1)], \quad (12.3)$$

où  $h = x_1 - x_0$ . Contrairement à ce que peut laisser croire la figure 12.1, cette méthode est généralement plus précise que la méthode du point milieu.

La formule composée pour l'approximation de  $\int_a^b f(x) dx$  par la méthode du trapèze avec  $n$  intervalles est

$$\int_a^b f(x) dx \approx \frac{h}{2} \left[ f(a) + f(b) + 2 \sum_{j=1}^{n-1} f(x_j) \right], \quad (12.4)$$

où  $h = (b - a)/n$  et  $x_j = a + jh$ .

## 12.5 Méthode de Simpson

La méthode d'approximation de Simpson est la plus usuelle des méthodes d'intégration numérique, encore que pas nécessairement la plus précise. La fonction  $f(x)$  est remplacée, sur l'intervalle  $(x_0, x_2)$ , par un polynôme d'interpolation de Lagrange du second degré (12.2(c)). On a donc  $m = 2$  et on peut démontrer que

$$\int_{x_0}^{x_2} f(x) dx \approx \frac{h}{3} [f(x_0) + 4f(x_1) + f(x_2)], \quad (12.5)$$

où  $h = x_2 - x_1 = x_1 - x_0$ . Cette méthode d'approximation numérique s'avère exacte pour les fonctions polynomiales de degré trois ou moins.

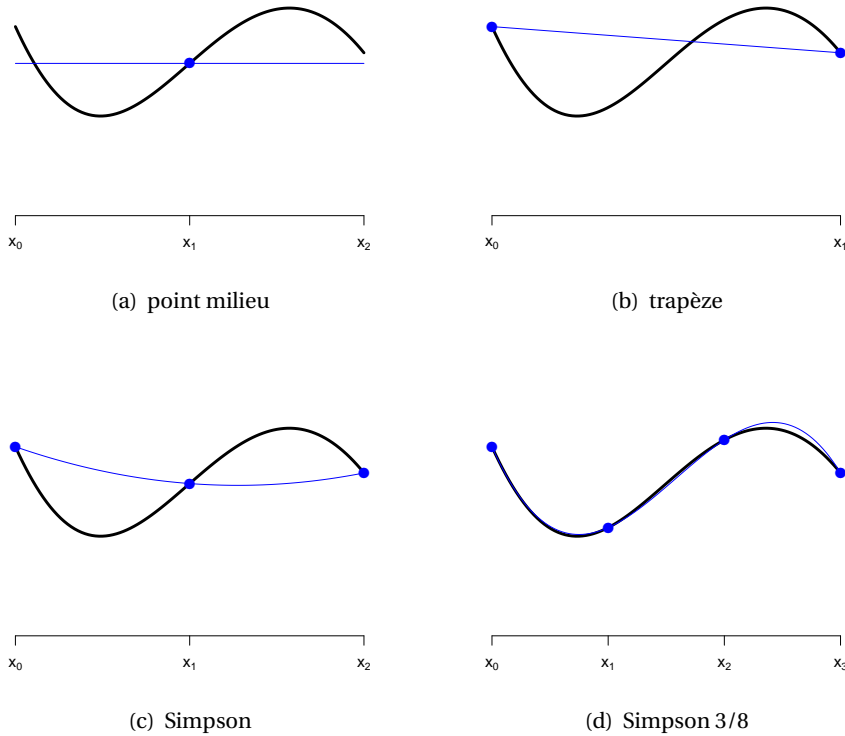
La formule composée pour l'approximation par la méthode de Simpson de  $\int_a^b f(x) dx$  avec  $n$  intervalles est

$$\int_a^b f(x) dx \approx \frac{h}{3} \left[ f(a) + 2 \sum_{j=1}^{n-1} f(x_{2j}) + 4 \sum_{j=1}^n f(x_{2j-1}) + f(b) \right], \quad (12.6)$$

où  $h = (b - a)/(2n)$  et  $x_j = a + jh$ .

## 12.6 Méthode de Simpson 3/8

La méthode de Simpson 3/8 constitue une extension de la méthode de Simpson où la fonction  $f$  est remplacée par un polynôme d'interpolation de degré  $m = 3$

FIG. 12.2: Approximation de  $f(x)$  sur un intervalle

(12.2(d)). On peut alors démontrer que

$$\int_{x_0}^{x_3} f(x) dx \approx \frac{3h}{8} [f(x_0) + 3f(x_1) + 3f(x_2) + f(x_3)], \quad (12.7)$$

où  $h = x_3 - x_2 = x_2 - x_1 = x_1 - x_0$ . La dérivation de la formule composée d'approximation de  $\int_a^b f(x) dx$  avec  $n$  intervalles est laissée en exercice.

Consulter [Burden et Faires \(2011\)](#) pour de plus amples détails sur les polynômes d'interpolation de Lagrange et les méthodes d'intégration numérique présentées ci-dessus.

► Il n'y a pas de fonction dans la distribution de base de R pour le calcul d'intégrales avec précisément l'une ou l'autre des méthodes mentionnées dans ce chapitre. Cependant, la fonction `integrate` réalise à peu près la même chose avec un algorithme différent. Consulter le code informatique de la [section 12.7](#) pour quelques menus exemples.

## 12.7 Code informatique

```
###
### POLYNÔMES D'INTERPOLATION DE LAGRANGE
###

## La fonction 'poly.calc' du package 'polynom' permet de
## calculer le polynômes d'interpolation de Lagrange passant
## par un ensemble de points (x, y), où x et y sont des
## vecteurs de même longueur. Le package n'est pas livré avec
## R, il faut donc l'installer depuis CRAN, puis le charger en
## mémoire. [Décommenter la ligne ci-dessous pour installer
## (une seule fois!) le package.]
#install.packages("polynom", repos = "http://cran.ca.r-project.org")
library(polynom)

## EXEMPLE 12.1
##
## On va vérifier les calculs de cet exemple avec 'poly.calc'.
x <- c(2, 2.75, 4)          # noeuds
(P1 <- poly.calc(x[1:2], 1/x[1:2])) # polynôme de degré 1
(P2 <- poly.calc(x, 1/x))      # polynôme de degré 2

## Pour calculer la valeur des polynômes en un point, il faut
## utiliser la fonction 'predict'.
predict(P1, 3)              # tel qu'obtenu dans l'exemple
predict(P2, 3)              # idem

###
### MÉTHODES D'INTÉGRATION NUMÉRIQUE
###

## Pas de démonstration en R des méthodes étudiées dans le
## chapitre. Il est laissé en exercice d'écrire des fonctions
## 'pointmilieu', 'trapeze', 'simpson' et 'simpson38'.
```



```
## Cependant, il y a dans R une fonction 'integrate' qui
## permet d'intégrer numériquement une fonction 'f' entre des
## bornes 'lower' et 'upper'.
integrate(sin, 0, 2)           # intégrale de sin(x)
f <- function(x) x^2 * exp(-x) # une autre fonction
integrate(f, 0, 1)            # intégrale sur [0, 1]
```

## 12.8 Exercices

**12.1** À partir des formules d'approximation de  $\int_{x_0}^{x_m} f(x) dx$ ,  $m = 1, 2, 3$ , développer les formules composées pour les méthodes d'approximation suivantes.

- a) Trapèze ( $m = 1$ )
- b) Simpson ( $m = 2$ )
- c) Simpson 3/8 ( $m = 3$ )

**12.2** Évaluer numériquement les intégrales suivantes avec les méthodes du point milieu, du trapèze, de Simpson et de Simpson 3/8. Dans chaque cas, n'utiliser qu'un seul sous-intervalle, c'est-à-dire  $n = 1$ .

- a)  $\int_{0,5}^1 x^4 dx$
- b)  $\int_0^{0,5} \frac{2}{x-4} dx$
- c)  $\int_1^{1,5} x^2 \ln x dx$
- d)  $\int_0^{\pi/4} e^{3x} \sin 2x dx$

**12.3** L'approximation de  $\int_0^2 f(x) dx$  avec  $n = 1$  donne 4 avec la méthode du trapèze et 2 avec la méthode de Simpson. Déterminer la valeur de  $f(1)$ .

## Réponses

- 12.2** a) 0,158203 ; 0,265625 ; 0,194010 ; 0,193866  
 b) -0,266667 ; -0,267857 ; -0,267064 ; -0,267063  
 c) 0,174331 ; 0,228074 ; 0,192245 ; 0,192253  
 d) 1,803915 ; 4,143260 ; 2,583696 ; 2,585789

**12.3**  $\frac{1}{2}$



# A Solutions des exercices

## Chapitre 10

La notation  $x_b$  signifie que le nombre  $x$  est en base  $b$ . On omet généralement  $b$  pour les nombres en base 10.

**10.1** L'algorithme de conversion des nombres décimaux en une base  $b$  se résume essentiellement à ceci pour la partie entière :

1. les chiffres du nombre en base  $b$  sont obtenus de droite à gauche en prenant le reste de divisions par  $b$  ;
2. on divise par  $b$  d'abord le nombre décimal d'origine, puis la partie entière de la division précédente, jusqu'à ce que celle-ci soit égale à 0.

On a donc les résultats suivants.

a) Conversion en base 6 :

$$119 \div 6 = 19 \text{ reste } 5$$

$$19 \div 6 = 3 \text{ reste } 1$$

$$3 \div 6 = 0 \text{ reste } 3,$$

d'où  $119 \equiv 315_6$ .

Conversion en binaire :

$$119 \div 2 = 59 \text{ reste } 1$$

$$59 \div 2 = 29 \text{ reste } 1$$

$$29 \div 2 = 14 \text{ reste } 1$$

$$14 \div 2 = 7 \text{ reste } 0$$

$$7 \div 2 = 3 \text{ reste } 1$$

$$3 \div 2 = 1 \text{ reste } 1$$

$$1 \div 2 = 0 \text{ reste } 1,$$

d'où  $119 \equiv 1110111_2$ .

b) Conversion en base 6 :

$$343 \div 6 = 57 \text{ reste } 1$$

$$57 \div 6 = 9 \text{ reste } 3$$

$$9 \div 6 = 1 \text{ reste } 3$$

$$1 \div 6 = 0 \text{ reste } 1,$$

$$\text{d'où } 343 \equiv 1331_6.$$

Conversion en binaire :

$$343 \div 2 = 171 \text{ reste } 1$$

$$171 \div 2 = 85 \text{ reste } 1$$

$$85 \div 2 = 42 \text{ reste } 1$$

$$42 \div 2 = 21 \text{ reste } 0$$

$$21 \div 2 = 10 \text{ reste } 1$$

$$10 \div 2 = 5 \text{ reste } 0$$

$$5 \div 2 = 2 \text{ reste } 1$$

$$2 \div 2 = 1 \text{ reste } 0$$

$$1 \div 2 = 0 \text{ reste } 1,$$

$$\text{d'où } 119 \equiv 10101011_2.$$

Conversion en binaire :

c) Conversion en base 6 :

$$96 \div 6 = 16 \text{ reste } 0$$

$$16 \div 6 = 2 \text{ reste } 4$$

$$2 \div 6 = 0 \text{ reste } 2,$$

$$\text{d'où } 96 \equiv 240_6.$$

$$96 \div 2 = 48 \text{ reste } 0$$

$$48 \div 2 = 24 \text{ reste } 0$$

$$24 \div 2 = 12 \text{ reste } 0$$

$$12 \div 2 = 6 \text{ reste } 0$$

$$6 \div 2 = 3 \text{ reste } 0$$

$$3 \div 2 = 1 \text{ reste } 1$$

$$1 \div 2 = 0 \text{ reste } 1,$$

$$\text{d'où } 96 \equiv 1100000_2.$$

Conversion en binaire :

d) Conversion en base 6 :

$$43 \div 6 = 7 \text{ reste } 1$$

$$7 \div 6 = 1 \text{ reste } 1$$

$$1 \div 6 = 0 \text{ reste } 1,$$

$$\text{d'où } 43 \equiv 111_6.$$

$$43 \div 2 = 21 \text{ reste } 1$$

$$21 \div 2 = 10 \text{ reste } 1$$

$$10 \div 2 = 5 \text{ reste } 0$$

$$5 \div 2 = 2 \text{ reste } 1$$

$$2 \div 2 = 1 \text{ reste } 0$$

$$1 \div 2 = 0 \text{ reste } 1,$$

$$\text{d'où } 43 \equiv 101011_2.$$

**10.2** On fait les deux premières conversions à l'aide de la définition d'un nombre hexadécimal, puis les deux dernières à l'aide de l'algorithme de conversion des nombres en base  $b$  vers la base 10.

- a)  $A1B_{16} = 10 \times 16^2 + 1 \times 16 + 11 = 2587$   
 b)  $12A_{16} = 1 \times 16^2 + 2 \times 16 + 10 = 298$   
 c)  $B41_{16} = (11 \times 16 + 4) \times 16 + 1 = 2881$   
 d)  $BAFFE_{16} = (((11 \times 16 + 10) \times 16) + 15) \times 16 + 15 \times 16 + 14 = 765950$

**10.3** La généralisation de l'algorithme de conversion des nombres en base  $b$  vers la base 10 à la conversion d'un nombre

$$x = x_{m-1}x_{m-2} \cdots x_1x_0$$

en base  $[b_{m-1} \dots b_0]$  vers la base 10 est la suivante (nombre entiers seulement) :

1. Poser  $x = 0$ .
2. Pour  $i = m-1, m-2, \dots, 0$ , faire les étapes suivantes.
  - i) Trouver  $d_i$ , le nombre décimal correspondant au symbole  $x_i$ .
  - ii) Poser  $x = xb_{i-1} + d_i$ , avec  $b_{-1} = 1$ .

Cet algorithme permet de trouver les formules demandées.

- a) On trouve la position de l'élément  $a_{ijk}$  dans l'ordre de la liste des éléments du tableau en convertissant le nombre  $[i-1 \ j-1 \ k-1]$  de la base  $[I \ J \ K]$  à la base 10, puis à additionnant 1. À l'aide de l'algorithme ci-dessus, on obtient

$$(((i-1) \times J + j-1) \times K + k-1) + 1 = k + K(j-1 + J(i-1))$$

- b) Dans l'ordre R, on convertit le nombre  $[k-1 \ j-1 \ i-1]$  exprimé dans la base  $[K \ J \ I]$  en base 10. On obtient alors

$$(((k-1) \times J + j-1) \times I + i-1) + 1 = i + I(j-1 + J(k-1))$$

soit la même réponse qu'à l'exercice 3.7 b) de la partie I.

**10.4** a) Par la valeur du bit fort, on sait que  $\boxed{1} \boxed{0000000}$  est un nombre négatif. Pour retrouver sa valeur absolue, on soustrait d'abord 1 (ce qui donne 0111111) et on inverse les bits : 10000000. Or  $10000000_2 = 2^7 = 128$ . Le nombre décimal correspondant à  $\boxed{1} \boxed{0000000}$  en notation en complément à deux est donc  $-128$ .

- b) Tout d'abord,  $15_{10} = 1111_2$  et  $5_{10} = 101_2$ .  
 i) Selon la notation «naturelle», on aurait pour 15 et  $-5$  les représentations sur 7 bits suivantes, dans l'ordre :

$$\boxed{0} \boxed{0001111} \quad \text{et} \quad \boxed{1} \boxed{0000101}.$$

Or, en additionnant en binaire, on obtient

TAB. A.1: Séquences des nombres en notation en complément à deux sur 8 bits

Bits	Valeur brute	Valeur en complément à deux
0000 0000	0	0
0000 0001	1	1
0000 0010	2	2
0111 1110	126	126
0111 1111	127	127
1000 0000	128	-128
1000 0001	129	-127
1000 0010	130	-126
1111 1110	254	-2
1111 1111	255	-1

$$\begin{array}{r}
 0\ 0001111 \\
 +\ 1\ 0000101 \\
 \hline
 1\ 0010100
 \end{array}$$

et  $10010100_2 = -20_{10} \neq 15 + (-5)$ .

- ii) En notation en complément à deux, on a plutôt les représentations suivantes pour 15 et -5 :

$$\boxed{0} \boxed{0001111} \quad \text{et} \quad \boxed{1} \boxed{1111011}.$$

En additionnant en binaire et en ignorant le neuvième bit qui apparaîtrait à gauche suite à la retenue, on obtient :

$$\begin{array}{r}
 0\ 0001111 \\
 +\ 1\ 1111011 \\
 \hline
 0\ 0001010
 \end{array}$$

Toujours en complément à deux,  $0001010_2 = 10_{10} = 15 + (-5)$ . On voit donc que l'addition en complément à deux donne directement le bon résultat.

- c) Voir les tableaux A.1 et A.2. La troisième colonne des tableaux est obtenue par l'opération  $-(2^n - x)$ , où  $n = 8$  ou  $n = 16$  et  $x$  est la valeur dans la seconde colonne.
- d) L'ensemble des nombres admissibles pour le type Byte est  $\{0, \dots, 2^8 - 1\} = \{0, \dots, 255\}$ . Pour le type Integer, on a de la partie c) que les nombres

TAB. A.2: Séquences des nombres en notation en complément à deux sur 16 bits

Bits	Valeur brute	Valeur en complément à deux
0000 0000 0000 0000	0	0
0000 0000 0000 0001	1	1
0000 0000 0000 0010	2	2
0111 1111 1111 1110	32 766	32 766
0111 1111 1111 1111	32 767	32 767
1000 0000 0000 0000	32 768	-32 768
1000 0000 0000 0001	32 769	-32 767
1000 0000 0000 0010	32 770	-32 766
1111 1111 1111 1110	65 534	-2
1111 1111 1111 1111	65 535	-1

admissibles sont les entiers  $\{-2^{15}, \dots, 2^{15} - 1\} = \{-32\,768, \dots, 32\,767\}$ . Par analogie, on trouve que l'ensemble des entiers admissibles pour le type Long est  $\{-2^{31}, \dots, 2^{31} - 1\} = \{-2\,147\,483\,648, \dots, 2\,147\,483\,647\}$ .

**10.5** Voir la [section 10.4](#). Les calculs sont exactement les mêmes que pour les nombres en double précision.

**10.6** Dans les égalités ci-dessous, le côté droit est en binaire.

a) Premièrement,  $1234 \equiv 10011010010_2$ . On a donc

$$\begin{aligned} -1234 &= (-1)^1 \times 2^{10} \times 1,001101001 \\ &= (-1)^1 \times 2^{137-127} \times 1,1010010. \end{aligned}$$

Or, puisque  $137 \equiv 10001001_2$ , on a la représentation en simple précision

$$\boxed{1} \boxed{10001001} \boxed{001101001000000000000000}$$

b) On a  $55 \equiv 110111_2$ , d'où

$$\begin{aligned} 55 &= (-1)^0 \times 2^5 \times 1,10111 \\ &= (-1)^0 \times 2^{132-127} \times 1,10111. \end{aligned}$$

Or, puisque  $132 \equiv 10000100_2$ , on a la représentation en simple précision

$$\boxed{0} \boxed{10000100} \boxed{101110000000000000000000}$$

c) On a  $8191 \equiv 111111111111_2$  et  $149 \equiv 10001011$ , d'où

$$\begin{aligned} 8191 &= (-1)^0 \times 2^{12} \times 1,111111111111 \\ &= (-1)^0 \times 2^{149-127} \times 1,111111111111 \\ &= \boxed{0} \boxed{10001011} \boxed{111111111111000000000000}. \end{aligned}$$

d) On a  $10 \equiv 1010_2$  et  $130 \equiv 10000010$ , d'où

$$\begin{aligned} -10 &= (-1)^1 \times 2^3 \times 1,010 \\ &= (-1)^1 \times 2^{130-127} \times 1,010 \\ &= \boxed{1} \boxed{10000010} \boxed{010000000000000000000000}. \end{aligned}$$

e) La représentation de  $\frac{2}{3}$  en binaire est  $0,101010\dots$  (La façon la plus simple d'obtenir ce résultat consiste à convertir  $\frac{2}{3} \times 2^n$ , où  $n$  est le nombre de bits souhaité après la virgule). Puisque  $126 \equiv 1111110_2$ , on a

$$\begin{aligned} \frac{2}{3} &= (-1)^0 \times 2^{-1} \times 1,010101010101010101010 \\ &= (-1)^0 \times 2^{126-127} \times 1,010101010101010101010 \\ &= \boxed{0} \boxed{01111110} \boxed{010101010101010101010}. \end{aligned}$$

f) La représentation binaire de  $\frac{1}{100}$  est infinie :  $0,000000101000111101\dots$ .  
Puisque  $120 \equiv 01111000_2$ , on a

$$\begin{aligned} \frac{1}{100} &= (-1)^0 \times 2^{-7} \times 1,01000111101011100001010 \\ &= (-1)^0 \times 2^{120-127} \times 1,01000111101011100001010 \\ &= \boxed{0} \boxed{01111000} \boxed{01000111101011100001010} \end{aligned}$$

**10.7** a) Puisque  $111101_2 \equiv 61$ , on a le nombre

$$\begin{aligned} (-1)^0 \times 2^{61-127} \times 1,100100001 &= (-1)^0 \times 2^{-66} \times 1,100100001 \\ &= 2^{-66}(1 + 2^{-1} + 2^{-4} + 2^{-9}) \\ &\equiv 2,120229346 \times 10^{-20}. \end{aligned}$$

b) Signe inversé par rapport à la partie a).

c) Puisque  $10000100_2 = 2^7 + 2^2 \equiv 132$ , on a le nombre

$$\begin{aligned} (-1)^0 \times 2^{132-127} \times 1,100100001 &= (-1)^0 \times 2^5 \times 1,100100001 \\ &= 2^5(1 + 2^{-1} + 2^{-4} + 2^{-9}) \\ &\equiv 50,0625. \end{aligned}$$



d) Signe inversé par rapport à la partie c).

**10.8** a) Le nombre suivant est

0	00111101	100100001000000000000001
---	----------	--------------------------

,

soit

$$2^{-66}(1 + 2^{-1} + 2^{-4} + 2^{-9} + 2^{-23}) \equiv 2,120\,229\,508 \times 10^{-20}.$$

Le nombre précédent est

0	00111101	100100000111111111111111
---	----------	--------------------------

,

soit

$$2^{-66}(1 + 2^{-1} + 2^{-4} + 2^{-9} - 2^{-23}) \equiv 2,120\,229\,185 \times 10^{-20}.$$

b) Le nombre suivant est

0	10000100	100100001000000000000001
---	----------	--------------------------

,

soit

$$2^5(1 + 2^{-1} + 2^{-4} + 2^{-9} + 2^{-23}) \equiv 50,062\,503\,815.$$

Le nombre précédent est

0	10000100	100100000111111111111111
---	----------	--------------------------

,

soit

$$2^5(1 + 2^{-1} + 2^{-4} + 2^{-9} - 2^{-23}) \equiv 50,062\,496\,185.$$

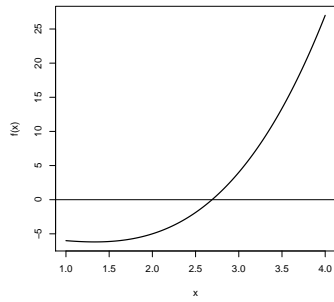
On remarque que les nombres sont beaucoup plus éloignés les uns des autres ici qu'en a).

## Chapitre 11

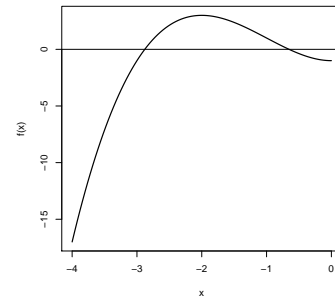
**11.2** Les solutions suivantes ont été obtenues à l'aide de nos fonctions de résolution d'équations à une variable. Les valeurs intermédiaires sont affichées pour montrer la convergence. La [figure A.1](#) contient les graphiques des cinq fonctions pour les intervalles mentionnés dans l'énoncé.

a) La fonction n'a qu'une seule racine dans  $[1, 4]$ .

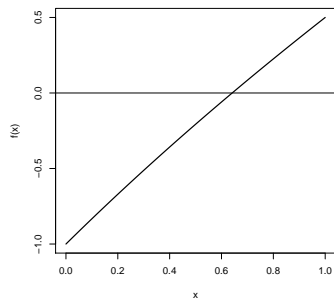
```
> f1 <- function(x) x^3 - 2 * x^2 - 5
> f1p <- function(x) 3 * x^2 - 4 * x
> bisection(f1, lower = 1, upper = 4, echo = TRUE)
```



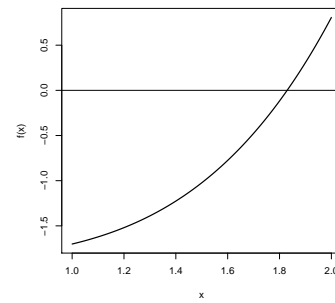
(a)  $f(x) = x^3 - 2x^2 - 5$



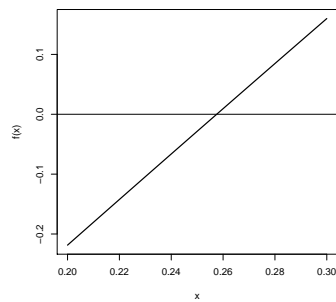
(b)  $f(x) = x^3 + 3x^2 - 1$



(c)  $f(x) = x - 2^{-x}$



(d)  $f(x) = e^x + 2^{-x} + 2 \cos x - 6$



(e)  $f(x) = e^x - x^2 + 3x - 2$

FIG. A.1: Fonctions de l'exercice 11.2

```
[1] 1.000 4.000 2.500 -1.875
[1] 2.5000 4.0000 3.2500 8.2031
[1] 2.5000 3.2500 2.8750 2.2324
[1] 2.500000 2.875000 2.687500 -0.034424
[1] 2.6875 2.8750 2.7812 1.0432
[1] 2.68750 2.78125 2.73438 0.49078
[1] 2.68750 2.73438 2.71094 0.22481
[1] 2.687500 2.710938 2.699219 0.094355
[1] 2.687500 2.699219 2.693359 0.029757
[1] 2.6875000 2.6933594 2.6904297 -0.0023855
[1] 2.690430 2.693359 2.691895 0.013673
[1] 2.6904297 2.6918945 2.6911621 0.0056403
[1] 2.6904297 2.6911621 2.6907959 0.0016266
[1] 2.69042969 2.69079590 2.69061279 -0.00037968
[1] 2.6906128 2.6907959 2.6907043 0.0006234
[1] 2.69061279 2.69070435 2.69065857 0.00012185
[1] 2.69061279 2.69065857 2.69063568 -0.00012892
[1] 2.6906e+00 2.6907e+00 2.6906e+00 -3.5365e-06
[1] 2.6906e+00 2.6907e+00 2.6907e+00 5.9155e-05
[1] 2.6906e+00 2.6907e+00 2.6906e+00 2.7809e-05
[1] 2.6906e+00 2.6906e+00 2.6906e+00 1.2136e-05
$root
[1] 2.6906

$nb.iter
[1] 21
> nr(f1, f1p, start = 2.5, echo = TRUE)
[1] 2.5
[1] 2.7143
[1] 2.691
[1] 2.6906
$root
[1] 2.6906

$nb.iter
[1] 3
> secante(f1, start0 = 1, start1 = 4, echo = TRUE)
[1] 1.5455
[1] 1.9969
```

```
[1] 4.1051
[1] 2.2947
[1] 2.4787
[1] 2.7514
[1] 2.6831
[1] 2.6904
[1] 2.6906
$root
[1] 2.6906
```

```
$nb.iter
[1] 9
```

- b) La fonction possède deux racines dans  $[-4, 0]$ . La convergence se fera vers l'une ou l'autre selon la position de la ou des valeurs de départ par rapport à l'extremum de la fonction dans l'intervalle. Ici, il s'agit d'un maximum  $x = -2$ . Ainsi, avec des valeurs de départ inférieures au maximum, on trouve la première racine :

```
> f2 <- function(x) x^3 + 3 * x^2 - 1
> f2p <- function(x) 3 * x^2 + 6 * x
> bisection(f2, lower = -3, upper = -2.8, echo = TRUE)
[1] -3.000 -2.800 -2.900 -0.159
[1] -2.90000 -2.80000 -2.85000 0.21838
[1] -2.900000 -2.850000 -2.875000 0.033203
[1] -2.900000 -2.875000 -2.887500 -0.062014
[1] -2.887500 -2.875000 -2.881250 -0.014185
[1] -2.8812500 -2.8750000 -2.8781250 0.0095642
[1] -2.8812500 -2.8781250 -2.8796875 -0.0022966
[1] -2.8796875 -2.8781250 -2.8789062 0.0036373
[1] -2.87968750 -2.87890625 -2.87929688 0.00067121
[1] -2.87968750 -2.87929688 -2.87949219 -0.00081245
[1] -2.8795e+00 -2.8793e+00 -2.8794e+00 -7.0567e-05
[1] -2.87939453 -2.87929688 -2.87934570 0.00030034
[1] -2.87939453 -2.87934570 -2.87937012 0.00011489
[1] -2.8794e+00 -2.8794e+00 -2.8794e+00 2.2161e-05
[1] -2.8794e+00 -2.8794e+00 -2.8794e+00 -2.4203e-05
[1] -2.8794e+00 -2.8794e+00 -2.8794e+00 -1.0210e-06
[1] -2.87938538 -2.87938232 -2.87938385 0.00001057
$root
[1] -2.8794
```

```
$nb.iter
[1] 17
> nr(f2, f2p, start = -3, echo = TRUE)
[1] -3
[1] -2.8889
[1] -2.8795
[1] -2.8794
$root
[1] -2.8794

$nb.iter
[1] 3
> secante(f2, start0 = -3, start1 = -2, echo = TRUE)
[1] -2.75
[1] -3.0667
[1] -2.862
[1] -2.8772
[1] -2.8794
[1] -2.8794
$root
[1] -2.8794

$nb.iter
[1] 6
Pour trouver la seconde racine, on utilise des valeurs de départ supérieures
au maximum :
> bisection(f2, lower = -1, upper = 0.5, echo = TRUE)
[1] -1.00000 0.50000 -0.25000 -0.82812
[1] -1.000000 -0.250000 -0.625000 -0.072266
[1] -1.00000 -0.62500 -0.81250 0.44409
[1] -0.81250 -0.62500 -0.71875 0.17850
[1] -0.718750 -0.625000 -0.671875 0.050953
[1] -0.671875 -0.625000 -0.648438 -0.011236
[1] -0.671875 -0.648438 -0.660156 0.019719
[1] -0.6601562 -0.6484375 -0.6542969 0.0042058
[1] -0.6542969 -0.6484375 -0.6513672 -0.0035239
[1] -0.65429688 -0.65136719 -0.65283203 0.00033872
```

```
[1] -0.6528320 -0.6513672 -0.6520996 -0.0015932
[1] -0.65283203 -0.65209961 -0.65246582 -0.00062736
[1] -0.65283203 -0.65246582 -0.65264893 -0.00014435
[1] -6.5283e-01 -6.5265e-01 -6.5274e-01 9.7175e-05
[1] -6.5274e-01 -6.5265e-01 -6.5269e-01 -2.3592e-05
[1] -6.5274e-01 -6.5269e-01 -6.5272e-01 3.6791e-05
[1] -6.5272e-01 -6.5269e-01 -6.5271e-01 6.5995e-06
[1] -6.5271e-01 -6.5269e-01 -6.5270e-01 -8.4961e-06
[1] -6.5271e-01 -6.5270e-01 -6.5270e-01 -9.4828e-07
[1] -6.5271e-01 -6.5270e-01 -6.5270e-01 2.8256e-06
[1] -6.5270e-01 -6.5270e-01 -6.5270e-01 9.3867e-07
[1] -6.527e-01 -6.527e-01 -6.527e-01 -4.804e-09
```

```
$root
```

```
[1] -0.6527
```

```
$nb.iter
```

```
[1] 22
```

```
> nr(f2, f2p, start = -1, echo = TRUE)
```

```
[1] -1
```

```
[1] -0.66667
```

```
[1] -0.65278
```

```
[1] -0.6527
```

```
$root
```

```
[1] -0.6527
```

```
$nb.iter
```

```
[1] 3
```

```
> secante(f2, start0 = -2, start1 = -1, echo = TRUE)
```

```
[1] -0.5
```

```
[1] -0.63636
```

```
[1] -0.65394
```

```
[1] -0.6527
```

```
[1] -0.6527
```

```
$root
```

```
[1] -0.6527
```

```
$nb.iter
```

```
[1] 5
```

On remarquera que les deux valeurs de départ de la méthode de la sécante

n'ont pas à se trouver de part et d'autre de la racine.

- c) La fonction n'a qu'une seule racine dans  $[0, 1]$  et elle est légèrement supérieure à 0,6.

```
> f3 <- function(x) x - 2^(-x)
> f3p <- function(x) 1 + log(2) * 2^(-x)
> bisection(f3, lower = 0.6, upper = 0.65, echo = TRUE)

[1] 0.60000 0.65000 0.62500 -0.02342
[1] 0.6250000 0.6500000 0.6375000 -0.0053259
[1] 0.6375000 0.6500000 0.6437500 0.0037029
[1] 0.63750000 0.64375000 0.64062500 -0.00081001
[1] 0.6406250 0.6437500 0.6421875 0.0014468
[1] 0.6406250 0.6421875 0.6414062 0.0003185
[1] 0.64062500 0.64140625 0.64101562 -0.00024573
[1] 0.64101562 0.64140625 0.64121094 0.00003639
[1] 0.64101562 0.64121094 0.64111328 -0.00010467
[1] 0.64111328 0.64121094 0.64116211 -0.00003414
[1] 6.4116e-01 6.4121e-01 6.4119e-01 1.1251e-06
[1] 6.4116e-01 6.4119e-01 6.4117e-01 -1.6507e-05
[1] 6.4117e-01 6.4119e-01 6.4118e-01 -7.6910e-06
[1] 6.4118e-01 6.4119e-01 6.4118e-01 -3.2830e-06
[1] 6.4118e-01 6.4119e-01 6.4118e-01 -1.0789e-06
[1] 6.4118e-01 6.4119e-01 6.4119e-01 2.3101e-08
[1] 6.4118e-01 6.4119e-01 6.4119e-01 -5.2791e-07
$root
[1] 0.64119

$nb.iter
[1] 17

> nr(f3, f3p, start = 0.6, echo = TRUE)

[1] 0.6
[1] 0.641
[1] 0.64119
$root
[1] 0.64119

$nb.iter
[1] 2

> secante(f3, start0 = 0.6, start1 = 0.65, echo = TRUE)
```

```
[1] 0.64122
[1] 0.64119
$root
[1] 0.64119
```

```
$nb.iter
[1] 2
```

- d) La fonction n'a qu'une seule racine dans [1,2] et elle est légèrement supérieure à 1,8.

```
> f4 <- function(x) exp(x) + 2^(-x) + 2 * cos(x) - 6
> f4p <- function(x) exp(x) - 2^(-x) * log(2) - 2 * sin(x)
> bisection(f4, lower = 1.8, upper = 1.85, echo = TRUE)
[1] 1.800000 1.850000 1.825000 -0.017913
[1] 1.825000 1.850000 1.837500 0.033517
[1] 1.825000 1.837500 1.831250 0.007667
[1] 1.825000 1.831250 1.8281250 -0.0051567
[1] 1.8281250 1.831250 1.8296875 0.0012468
[1] 1.8281250 1.8296875 1.8289063 -0.0019571
[1] 1.82890625 1.82968750 1.82929688 -0.00035568
[1] 1.8292969 1.8296875 1.8294922 0.0004454
[1] 1.8293e+00 1.8295e+00 1.8294e+00 4.4827e-05
[1] 1.82929688 1.82939453 1.82934570 -0.00015544
[1] 1.8293e+00 1.8294e+00 1.8294e+00 -5.5307e-05
[1] 1.8294e+00 1.8294e+00 1.8294e+00 -5.2405e-06
[1] 1.8294e+00 1.8294e+00 1.8294e+00 1.9793e-05
[1] 1.8294e+00 1.8294e+00 1.8294e+00 7.2762e-06
[1] 1.8294e+00 1.8294e+00 1.8294e+00 1.0178e-06
$root
[1] 1.8294
```

```
$nb.iter
[1] 15
> nr(f4, f4p, start = 1.8, echo = TRUE)
[1] 1.8
[1] 1.8301
[1] 1.8294
$root
[1] 1.8294
```



```
$nb.iter
[1] 2
> secante(f4, start0 = 1.8, start1 = 1.85, echo = TRUE)
[1] 1.8289
[1] 1.8294
[1] 1.8294
$root
[1] 1.8294
```

```
$nb.iter
[1] 3
```

- e) Encore ici, la fonction n'a qu'une seule racine dans l'intervalle mentionné et elle se situe autour de 0,25.

```
> f5 <- function(x) exp(x) - x^2 + 3 * x - 2
> f5p <- function(x) exp(x) - 2 * x + 3
> bisection(f5, lower = 0.24, upper = 0.26, echo = TRUE)
[1] 0.240000 0.260000 0.250000 -0.028475
[1] 0.2500000 0.2600000 0.2550000 -0.0095634
[1] 0.25500000 0.26000000 0.25750000 -0.00011444
[1] 0.2575000 0.2600000 0.2587500 0.0046084
[1] 0.2575000 0.2587500 0.2581250 0.0022471
[1] 0.2575000 0.2581250 0.2578125 0.0010664
[1] 0.25750000 0.25781250 0.25765625 0.00047597
[1] 0.25750000 0.25765625 0.25757812 0.00018077
[1] 2.5750e-01 2.5758e-01 2.5754e-01 3.3166e-05
[1] 2.5750e-01 2.5754e-01 2.5752e-01 -4.0637e-05
[1] 2.5752e-01 2.5754e-01 2.5753e-01 -3.7355e-06
[1] 2.5753e-01 2.5754e-01 2.5753e-01 1.4715e-05
[1] 2.5753e-01 2.5753e-01 2.5753e-01 5.4898e-06
[1] 2.5753e-01 2.5753e-01 2.5753e-01 8.7717e-07
[1] 2.5753e-01 2.5753e-01 2.5753e-01 -1.4291e-06
[1] 2.5753e-01 2.5753e-01 2.5753e-01 -2.7598e-07
[1] 2.5753e-01 2.5753e-01 2.5753e-01 3.0059e-07
$root
[1] 0.25753
```

```
$nb.iter
[1] 17
> nr(f5, f5p, start = 0.26, echo = TRUE)
```

TAB. A.3: Valeurs successives de la méthode de bisection pour l'exercice 11.3

$n$	$a_n$	$b_n$	$x_n$	$f(x_n)$
1	0	2	1	-1
2	1,00	2,00	1,50	0,25
3	1,0000	1,5000	1,2500	-0,4375
4	1,250000	1,500000	1,375000	-0,109375
5	1,37500000	1,50000000	1,43750000	0,06640625
6	1,37500000	1,43750000	1,40625000	-0,02246094
7	1,40625000	1,43750000	1,42187500	0,02172852
8	1,4062500000	1,4218750000	1,4140625000	-0,0004272461
9	1,41406250	1,42187500	1,41796875	0,01063538
10	1,41406250	1,41796875	1,41601562	0,00510025

```

[1] 0.26
[1] 0.25753
[1] 0.25753
$root
[1] 0.25753

$nb.iter
[1] 2

> secante(f5, start0 = 0.24, start1 = 0.26, echo = TRUE)

[1] 0.25753
[1] 0.25753
$root
[1] 0.25753

$nb.iter
[1] 2

```

**11.3** On trouve la racine de  $f(x) = x^2 - 2$  dans l'intervalle  $[0, 2]$  par la méthode de bisection avec un maximum de 10 itérations. Le [tableau A.3](#) contient les valeurs successives de  $a$ ,  $b$ ,  $x = (a + b)/2$  et  $f(x)$ . On obtient donc une réponse de 1,41601562, alors que la vraie réponse est le nombre irrationnel  $\sqrt{2} = 1,414214$ .

**11.4** On a que

$$\begin{aligned} x_n - x_{n-1} &= \sum_{k=1}^n \frac{1}{k} - \sum_{k=1}^{n-1} \frac{1}{k} \\ &= \frac{1}{n}, \end{aligned}$$

d'où  $\lim_{n \rightarrow \infty} (x_n - x_{n-1}) = \lim_{n \rightarrow \infty} \frac{1}{n} = 0$ . Or,  $\sum_{k=1}^n k^{-1}$  est la série harmonique qui est connue pour diverger. On peut, par exemple, justifier ceci par le fait que l'intégrale

$$\int_1^{\infty} \frac{1}{x} dx$$

diverge elle-même. Cet exercice illustre donc que le critère  $|x_n - x_{n-1}| < \varepsilon$  peut être satisfait même pour une série divergente.

**11.5** a) On considère la fonction  $g(x) = 2^{-x}$  dans l'intervalle  $[0, 1]$ . En premier lieu,  $1/2 \leq 2^{-x} \leq 1$  pour  $0 \leq x \leq 1$ , donc  $g(x) \in [0, 1]$ . De plus,  $g'(x) = -2^{-x}(\ln 2)$ , d'où

$$|g'(x)| = \frac{\ln 2}{2^x} \leq \ln 2 < 1, \quad 0 \leq x \leq 1.$$

Par le [théorème 11.1](#), la fonction  $g$  possède un point fixe unique dans l'intervalle  $[0, 1]$ .

b) Le graphique de la fonction  $g$  se trouve à la [figure A.2](#). On constate que la fonction a effectivement un point fixe unique dans l'intervalle  $[0, 1]$  et que celui-ci se trouve près de  $x = 0,6$ . On peut donc utiliser cette valeur comme point de départ de l'algorithme du point fixe :

```
> pointfixe(function(x) 2^(-x), start = 0.6)
```

```
$fixed.point
```

```
[1] 0.64119
```

```
$nb.iter
```

```
[1] 15
```

On remarque que la réponse est indépendante de la valeur de départ :

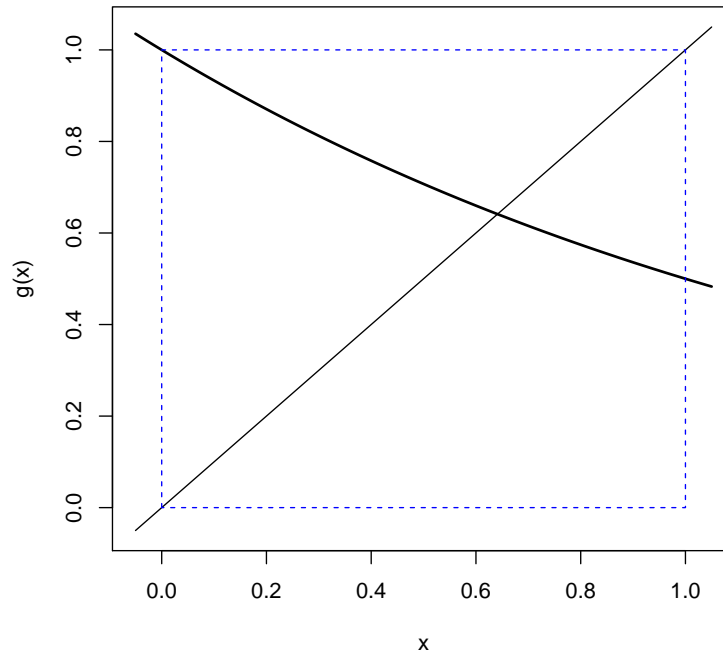
```
> pointfixe(function(x) 2^(-x), start = 0)
```

```
$fixed.point
```

```
[1] 0.64119
```

```
$nb.iter
```

```
[1] 18
```

FIG. A.2: Fonction  $g(x) = 2^{-x}$  dans  $[0, 1]$ 

```
> pointfixe(function(x) 2^(-x), start = 1)
```

```
$fixed.point  
[1] 0.64119
```

```
$nb.iter  
[1] 17
```

```
> pointfixe(function(x) 2^(-x), start = 2)
```

```
$fixed.point  
[1] 0.64119
```

```
$nb.iter  
[1] 18
```

**11.6** Les cinq fonctions sont les suivantes :

$$g_1(x) = x - x^3 - 4x^2 + 10$$

$$g_2(x) = \left( \frac{10}{x} - 4x \right)^{1/2}$$

$$g_3(x) = \frac{1}{2}(10 - x^3)^{1/2}$$

$$g_4(x) = \left( \frac{10}{4+x} \right)^{1/2}$$

$$g_5(x) = x - \frac{x^3 + 4x^2 - 10}{3x^2 + 8x}.$$

Soit  $f(x) = x^3 + 4x^2 - 10$ . On peut réécrire l'équation  $f(x) = 0$  de différentes façons :

- a)  $x = x - f(x)$  ;
- b)  $x^3 = 10 - 4x^2 \Rightarrow x^2 = 10/x - 4x \Rightarrow x = (10/x - 4x)^{1/2}$  ;
- c)  $4x^2 = 10 - x^3 \Rightarrow x^2 = (10 - x^3)/4 \Rightarrow x = \frac{1}{2}(10 - x^3)^{1/2}$  ;
- d)  $4x^2 + x^3 = 10 \Rightarrow x^2 = 10/(4+x) \Rightarrow x = (10/(4+x))^{1/2}$  ;
- e)  $x = x - f(x)/f'(x)$ .

Les fonctions  $g_1$  à  $g_5$  correspondent, dans l'ordre, au côté droit de chacune des équations ci-dessus. On remarquera que la fonction  $g_5$  est celle préconisée par la méthode de Newton–Raphson — et celle qui converge le plus rapidement.

**11.7** La fonction  $g(x) = 4x^2 - 14x + 14$  est une parabole ayant deux points fixes. On observe graphiquement que l'un se trouve en  $x = 2$  et que  $g(1,5) = g(2)$ . Or, pour toute valeur de départ  $x_0 < 1,5$  de l'algorithme du point fixe, on constate que les valeurs de  $x_1, x_2, \dots$  se retrouveront de plus en plus loin sur la branche droite de la parabole. Il en va de même pour tout  $x_0 > 2$ . La procédure diverge donc pour de telles valeurs de départ. En revanche, il est facile de vérifier graphiquement que la procédure converge pour toute valeur de départ dans l'intervalle  $[1,5, 2]$ . Si  $x_0 = 1,5$ , on obtient le point fixe  $x = 2$  après une seule itération.

**11.8** En premier lieu, on remarque que les fonctions  $g_1, g_2$  et  $g_3$  sont développées à partir de l'équation  $x^3 = 21$  pour calculer  $\sqrt[3]{21}$  par la méthode du point fixe. La [figure A.3](#) présente ces trois fonctions. On a

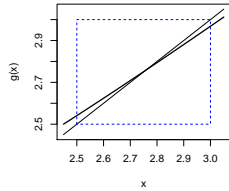
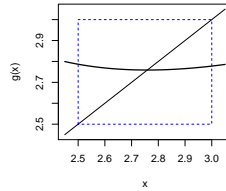
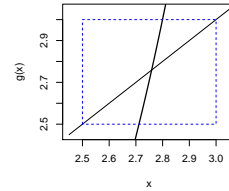
(a) fonction  $g_1(x)$ (b) fonction  $g_2(x)$ (c) fonction  $g_3(x)$ 

FIG. A.3: Fonctions de l'exercice 11.8

$$g'_1(x) = \frac{20}{21} - \frac{2}{x^3}$$

$$g'_2(x) = \frac{2}{3} - \frac{14}{x^3}$$

$$g'_3(x) = 1 - \frac{2x^5 - 84x^3 + 21x^2 + 441}{(x^2 - 21)^2}.$$

Considérons l'intervalle  $[2,5, 3]$  autour du point fixe. On a que  $|g'_2(x)| < |g'_1(x)|$  pour tout  $x$  dans cet intervalle, donc la procédure du point fixe converge plus rapidement pour  $g_2$  (la fonction de la méthode de Newton–Raphson). Cela se vérifie aisément sur les graphiques de la [figure A.3](#). Toujours à l'aide des graphiques, on voit que la procédure diverge avec la fonction  $g_3$ .

- 11.9** On considère la fonction  $g(x) = 2^{-x}$  dans l'intervalle  $[\frac{1}{3}, 1]$ . En premier lieu, on a  $g(x) \in [\frac{1}{2}, 1/\sqrt[3]{2}] \subset [\frac{1}{3}, 1]$  pour  $x \in [\frac{1}{3}, 1]$ . De plus, on a  $g'(x) = -2^{-x-1}(\ln 2)$ , d'où

$$|g'(x)| = \frac{\ln 2}{2^{x+1}} \leq \ln 2 < 1, \quad \frac{1}{3} \leq x \leq 1.$$

Par le théorème du point fixe, la fonction  $g$  possède donc un point fixe unique dans l'intervalle  $[\frac{1}{3}, 1]$ . La valeur de ce point fixe est

```
> pointfixe(function(x) 2^(-x), start = 2/3)
$fixed.point
[1] 0.64119

$nb.iter
[1] 14
```

- 11.10** La condition  $|g'(x)| \leq k < 1$  du [théorème 11.1](#) assure l'unicité d'un point fixe. Cette condition est nécessaire pour que la fonction  $g(x)$  ne puisse repasser au-dessus de la droite  $y = x$  après l'avoir déjà croisée une première fois pour

créer un point fixe. Or, il est seulement nécessaire de limiter la croissance de la fonction — soit les valeurs positives de sa pente. Peu importe la vitesse de décroissance de la fonction dans l'intervalle  $[a, b]$  (mais toujours sujet à ce que  $g(x) \in [a, b]$ ), la fonction ne peut croiser la droite  $y = x$  qu'une seule fois. La condition  $g'(x) \leq k < 1$  est donc suffisante pour assurer l'unicité du point fixe dans  $[a, b]$ .

**11.11** La distance entre deux points  $(x_1, y_1)$  et  $(x_2, y_2)$  de  $\mathbb{R}^2$  est

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}.$$

Par conséquent, la distance entre le point  $(x, x^2)$  et le point  $(1, 0)$  est une fonction de  $x$

$$d(x) = \sqrt{(x-1)^2 + (x^2-0)^2} = \sqrt{x^4 + x^2 - 2x + 1}.$$

On cherche la valeur de minimisant  $d(x)$  ou, de manière équivalente,  $d^2(x)$ . On doit donc résoudre

$$\frac{d}{dx} d^2(x) = f(x) = 4x^3 + 2x - 2 = 0$$

à l'aide de la méthode de Newton–Raphson. Cela requiert également

$$f'(x) = 12x^2 + 2.$$

Le point sur la courbe  $y = x^2$  le plus du point  $(1, 0)$  est donc la limite de la suite

$$x_n = x_{n-1} - \frac{4x_{n-1}^3 + 2x_{n-1} - 2}{12x_{n-1}^2 + 2}, \quad n = 1, 2, \dots,$$

avec  $x_0$  une valeur de départ «près» de la solution. On obtient

```
> nr(function(x) 4 * x^3 + 2 * x - 2,
+   function(x) 12 * x^2 + 2, start = 0.6)
$root
[1] 0.58975

$nb.iter
[1] 2
```

On voit à la [figure A.4](#) que la solution, disons le point  $(x^*, (x^*)^2)$ , est la projection du point  $(1, 0)$  sur la courbe  $y = x^2$ ; autrement dit, la droite passant par  $(x^*, (x^*)^2)$  et  $(1, 0)$  est perpendiculaire à la tangente de  $y = x^2$  en  $x = x^*$ .

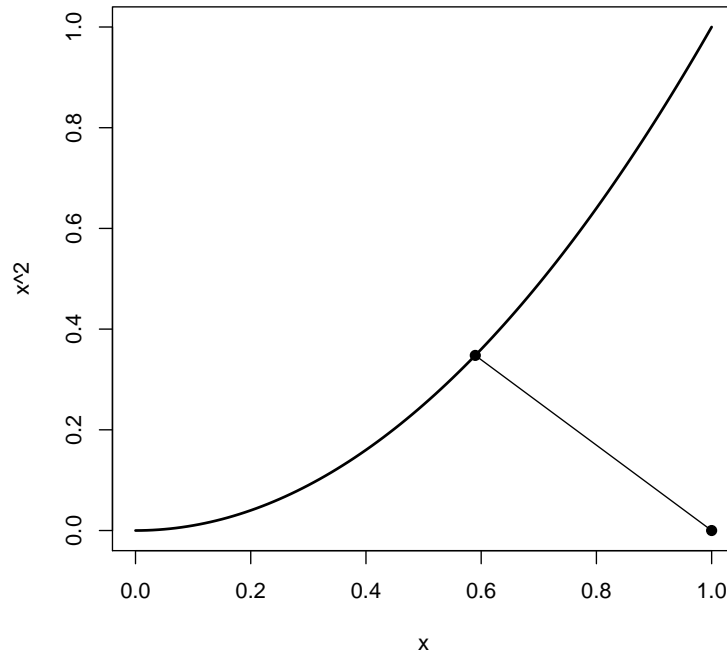


FIG. A.4: Point de  $y = x^2$  le plus près du point (1,0)

- 11.12** La fonction devra utiliser la méthode de Newton–Raphson pour trouver la racine de  $f(i) = 0$ , où

$$f(i) = \sum_{t=0}^n \frac{CF_t}{(1+i)^t}$$

et

$$f'(i) = - \sum_{t=0}^n \frac{tCF_t}{(1+i)^{t+1}}.$$

La fonction `tri` de la [figure A.5](#) calcule le taux de rendement interne d'un vecteur de flux financiers en utilisant la fonction de Newton–Raphson `nr`.

- 11.13** La [figure A.6](#) contient le code d'une fonction VBA `tri2` effectuant le calcul du taux de rendement interne d'un ensemble de flux financiers. Ceux-ci



```

tri <- function(CF, erreur.max = 1E-6)
{
  if (!any(CF < 0))
    stop("au moins un flux financier doit être
          négatif")

  t <- 0:(length(CF) - 1)

  f <- function(i) sum(CF/(1 + i)^t)
  fp <- function(i) sum((-CF * t)/((1 + i)^(t + 1)))

  nr(f, fp, start = 0.05, TOL = erreur.max)$root
}

```

FIG. A.5: Fonction R de calcul du taux de rendement interne d'une série de flux financiers

sont passés en argument à la fonction par le biais d'une plage horizontale ou verticale de montants périodiques consécutifs. La procédure de Newton-Raphson est codée à même cette fonction. On peut vérifier facilement que le résultat de cette fonction est identique à celui de la fonction Excel TRI().

- 11.14** a) On veut estimer par la méthode du maximum de vraisemblance le paramètre  $\theta$  d'une loi Gamma(3, 1/ $\theta$ ), dont la fonction de densité de probabilité est

$$f(x; \theta) = \frac{1}{2\theta^3} x^2 e^{-x/\theta}, \quad x > 0.$$

La fonction de log-vraisemblance d'un échantillon aléatoire  $x_1, \dots, x_n$  tiré de cette loi est

$$\begin{aligned}
 l(\theta) &= \sum_{i=1}^n \ln f(x_i; \theta) \\
 &= \sum_{i=1}^n \left( 2 \ln x_i - \frac{x_i}{\theta} - \ln 2 - 3 \ln \theta \right) \\
 &= 2 \sum_{i=1}^n \ln x_i - \frac{1}{\theta} \sum_{i=1}^n x_i - n \ln 2 - 3n \ln \theta.
 \end{aligned}$$

```

Function tri2(flux As Range, Optional ErreurMax
              As Double = 0.000000001)
    Dim nrow As Integer, ncol As Integer, n As Integer
    Dim CF() As Double, f As Double, fp As Double,
        i As Double, it As Double

    nrow = flux.Rows.Count
    ncol = flux.Columns.Count
    n = IIf(nrow > 1, nrow, ncol) - 1
    ReDim CF(0 To n)

    For t = 0 To n
        CF(t) = flux.Cells(t + 1, 1)
    Next t

    If CF(0) > 0 Then
        tri2 = "Erreur"
        Exit Function
    End If

    i = 0.05
    Do
        f = 0
        fp = 0
        it = i
        For t = 0 To n
            f = f + CF(t) / (1 + i) ^ t
            fp = fp - t * CF(t) / (1 + i) ^ (t + 1)
        Next t
        i = i - f / fp
    Loop Until Abs(i - it) / i < ErreurMax
    tri2 = i
End Function

```

FIG. A.6: Fonction VBA de calcul du taux de rendement interne d'une série de flux financiers

```

emv.gamma <- function(x, erreur.max = 1E-6)
{
  n <- length(x)
  xsum <- sum(x)
  f <- function(theta) xsum/theta^2 -
    (3 * n)/theta
  fp <- function(theta) -2 * xsum / theta^3 +
    (3 * n)/theta^2
  nr(f, fp, xsum/n/3, TOL = erreur.max)$root
}

```

FIG. A.7: Fonction R d'estimation du paramètre d'échelle d'une loi gamma par le maximum de vraisemblance

On cherche le maximum de la fonction  $l(\theta)$ , soit la valeur de  $\theta$  tel que

$$l'(\theta) = \frac{1}{\theta^2} \sum_{i=1}^n -\frac{3n}{\theta} = 0.$$

Résoudre cette équation à l'aide de la méthode de Newton–Raphson requiert également

$$l''(\theta) = -\frac{2}{\theta^3} \sum_{i=1}^n x_i + \frac{3n}{\theta^2}.$$

La [figure A.7](#) présente une fonction R pour effectuer le calcul à l'aide de notre fonction de Newton–Raphson `nr`. On remarquera que la somme des valeurs de l'échantillon — qui ne change pas durant la procédure itérative — est calculée une fois pour toute dès le début de la fonction. De plus, on utilise comme valeur de départ l'estimateur des moments de  $\theta$ ,  $\bar{x}/3$ .

b) On a, par exemple,

```

> x <- rgamma(20, shape = 3, scale = 1000)
> emv.gamma(x)
[1] 809.96

```

**11.15** La [figure A.8](#) présente le code VBA d'une fonction `emvGamma` très similaire à la fonction R de l'exercice précédent. Pour l'utiliser dans Excel, il suffit de lui passer en argument une plage (verticale) contenant un échantillon aléatoire d'une loi  $\text{Gamma}(3, 1/\theta)$ .

```

Function emvGamma(Data As Range, Optional ErreurMax
    As Double = 0.000000001)
    Dim n As Integer
    Dim sData As Double, f As Double,
        fp As Double, x As Double, xt As Double

    n = Data.Rows.Count
    sData = WorksheetFunction.Sum(Data)

    x = sData / n / 3
    Do
        f = sData / x ^ 2 - (3 * n) / x
        fp = -2 * sData / x ^ 3 + 3 * n / x ^ 2
        xt = x
        x = x - f / fp
    Loop Until Abs(x - xt) / x < ErreurMax
    emvGamma = x
End Function

```

FIG. A.8: Fonction VBA d'estimation du paramètre d'échelle d'une loi gamma par le maximum de vraisemblance

#### 11.16 On cherche à résoudre l'équation

$$a_{10|i}^{(12)} = \frac{1 - (1+i)^{-10}}{i^{(12)}} = 8$$

ou, de manière équivalente,

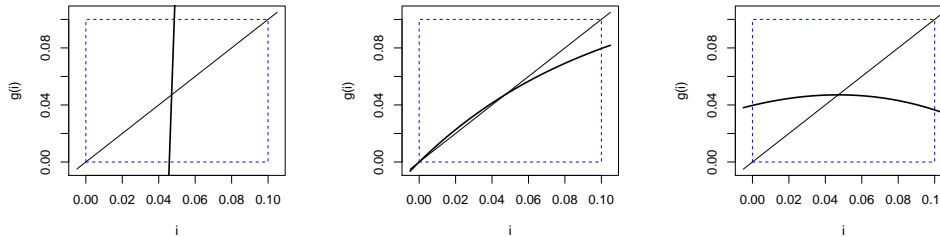
$$f(i) = \frac{1 - (1+i)^{-10}}{12(1+i)^{1/12} - 1} - 8 = 0.$$

Pour résoudre le problème à l'aide de la méthode du point fixe, on peut considérer la fonction

$$g_1(i) = i - f(i),$$

mais, comme le graphique de la A.9(a) le montre, la procédure itérative divergera. En posant

$$\frac{1 - (1+i)^{-10}}{8} = 12[(1+i)^{1/12} - 1],$$

(a) fonction  $g_1(x)$ (b) fonction  $g_2(x)$ (c) fonction  $g_3(x)$ FIG. A.9: Trois fonctions pour résoudre  $a_{10}^{(12)}i = 8$  par la méthode du point fixe

puis en isolant  $i$  du côté droit de l'équation, on obtient une alternative

$$g_2(i) = \left( \frac{97 - (1+i)^{-10}}{96} \right)^{12} - 1.$$

Cette fonction satisfait les hypothèses pour l'existence et l'unicité d'un point fixe, mais sa pente est toutefois près de 1, donc la convergence sera lente ; voir la A.9(b). Enfin, on peut considérer la fonction

$$g_3(i) = i - \frac{f(i)}{f'(i)} = \frac{120(1+i)^{-11}[(1+i)^{1/12} - 1] - [1 - (1+i)^{-10}](1+i)^{-11/12}}{144[(1+i)^{1/12} - 1]^2}.$$

On remarquera que cette dernière est la fonction utilisée dans la méthode de Newton–Raphson et sa pente est presque nulle, d'où une convergence très rapide ; voir la A.9(c). On obtient les résultats suivants :

```
> f <- function(i)
+   (1 - (1 + i)^(-10))/(12 * ((1 + i)^(1/12) - 1)) - 8
> fp <- function(i)
+   (120 * (1 + i)^(-11) * ((1 + i)^(1/12) - 1) -
+   (1 - (1 + i)^(-10)) * (1 + i)^(-11/12)) /
+   (144 * ((1 + i)^(1/12) - 1)^2)
> g2 <- function(i) ((97 - (1 + i)^(-10))/96)^12 - 1
> g3 <- function(i) i - f(i)/fp(i)
> pointfixe(g2, start = 0.05)
```

```

$fixed.point
[1] 0.047081

$nb.iter
[1] 40

> pointfixe(g3, start = 0.05)
$fixed.point
[1] 0.047081

$nb.iter
[1] 2

```

- 11.17** a) `> f <- function(x) x^3 - 2 * x^2 - 5`  
`> uniroot(f, lower = 1, upper = 4)`
- b) Comme un simple graphique le démontre, il y a deux racines dans l'intervalle.
- ```

> f <- function(x) x^3 + 3 * x^2 - 1
> curve(f, xlim = c(-4, 0))
> uniroot(f, lower = -4, upper = -1)
> uniroot(f, lower = -1, upper = 0)

```
- c) `> f <- function(x) x - 2^(-x)`  
`> uniroot(f, lower = 0, upper = 1)`
- d) `> f <- function(x) exp(x) + 2^(-x) + 2 * cos(x) - 6`  
`> uniroot(f, lower = 1, upper = 2)`
- e) `> f <- function(x) exp(x) - x^2 + 3 * x - 2`  
`> uniroot(f, lower = 0, upper = 1)`

**11.18** On a

```

> dpareto <- function(x, alpha, lambda)
+ {
+   (alpha * lambda^alpha)/(x + lambda)^(alpha+1)
+ }
> f <- function(par, x) -sum(log(dpareto(x, par[1], par[2])))
> optim(c(1, 1000), f, x = x)

```

ou, en utilisant l'Astuce Ripley ([section 11.6](#)) consistant à estimer le logarithme des paramètres pour éviter les soucis de convergence :

```

> dpareto <- function(x, logAlpha, logLambda)
+ {
+   alpha <- exp(logAlpha)
+   lambda <- exp(logLambda)
+   (alpha * lambda^alpha)/(x + lambda)^(alpha+1)
+ }
> optim(c(log(2), log(1000)), f, x = x)
> exp(optim(c(log(2), log(1000)), f, x = x)$par)

```

## Chapitre 12

**12.1** On fait le cas de la méthode Simpsons seulement, la procédure à suivre dans les autres cas étant tout à fait similaire.

On souhaite développer la formule (12.6) à partir de l'approximation sur un sous intervalle donnée par l'équation (12.5) pour l'intervalle  $[x_0, x_2]$ . De manière plus générale, nous avons

$$\int_{x_{2j}}^{x_{2(j+1)}} f(x) dx \approx \frac{h}{3} [f(x_{2j}) + 4f(x_{2j+1}) + f(x_{2(j+1)})], \quad j = 0, \dots, n-1,$$

d'où

$$\begin{aligned}
 \int_a^b f(x) dx &\approx \sum_{j=0}^{n-1} \int_{x_{2j}}^{x_{2(j+1)}} f(x) dx \\
 &= \frac{h}{3} \sum_{j=0}^{n-1} [f(x_{2j}) + 4f(x_{2j+1}) + f(x_{2(j+1)})] \\
 &= \frac{h}{3} \left[ f(x_0) + \sum_{j=1}^{n-1} f(x_{2j}) + 4 \sum_{j=0}^{n-1} f(x_{2j+1}) + \sum_{j=0}^{n-2} f(x_{2(j+1)}) + f(x_{2n}) \right] \\
 &= \frac{h}{3} \left[ f(a) + 2 \sum_{j=1}^{n-1} f(x_{2j}) + 4 \sum_{j=0}^{n-1} f(x_{2j+1}) + f(b) \right]
 \end{aligned}$$

puisque, par définition,  $f(x_0) = f(a)$  et  $f(x_{2n}) = f(b)$ .

**12.2** Puisque  $n = 1$ , on peut utiliser directement les formules (12.1), (12.3), (12.5) et (12.7).

a) Point milieu :

$$\int_{0,5}^1 x^4 dx \approx 2(0,25)(0,75)^4 = 0,158203$$

Trapèze :

$$\int_{0,5}^1 x^4 dx \approx \frac{0,5}{2} [(0,5)^4 + 1] = 0,265625$$

Simpson :

$$\int_{0,5}^1 x^4 dx \approx \frac{0,25}{3} [(0,5)^4 + 4(0,75)^4 + 1] = 0,194010$$

Simpson 3/8 :

$$\int_{0,5}^1 x^4 dx \approx \frac{3(0,5/3)}{8} [(0,5)^4 + 3(2/3)^4 + 3(5/6)^4 + 1] = 0,193866$$

b) Point milieu :

$$\int_0^{0,5} \frac{2}{x-4} dx \approx 2(0,25)(-2/3,75) = 0,266667$$

Trapèze :

$$\int_0^{0,5} \frac{2}{x-4} dx \approx \frac{0,5}{2} [(-2/4) + (-2/3,5)] = -0,267857$$

Simpson :

$$\int_0^{0,5} \frac{2}{x-4} dx \approx \frac{0,25}{3} [(-2/4) + 4(-2/3,75) + (-2/3,5)] = -0,267064$$

Simpson 3/8 :

$$\begin{aligned} \int_0^{0,5} \frac{2}{x-4} dx &\approx \frac{3(0,5/3)}{8} [(-2/4) + 3(-12/23) + 3(12/22) + (-2/3,5)] \\ &= -0,267063 \end{aligned}$$

c) Point milieu :

$$\int_1^{1,5} x^2 \ln x dx \approx 2(0,25)(1,75)^2 \ln 1,75 = 0,174331$$

Trapèze :

$$\int_1^{1,5} x^2 \ln x dx \approx \frac{0,5}{2} [0 + (1,5)^2 \ln 1,5] = 0,228074$$



Simpson :

$$\int_1^{1,5} x^2 \ln x \, dx \approx \frac{0,25}{3} [0 + 4(1,25)^2 \ln 1,25 + (1,5)^2 \ln 1,5] = 0,192245$$

Simpson 3/8 :

$$\begin{aligned} \int_1^{1,5} x^2 \ln x \, dx \approx \frac{3(0,5/3)}{8} [0 + (7/6)^2 \ln(7/6) + (4/3)^2 \ln(4/3) \\ + (1,5)^2 \ln 1,5] = 0,192253 \end{aligned}$$

d) Point milieu :

$$\int_0^{\pi/4} e^{3x} \sin 2x \, dx \approx 2(\pi/8) e^{3\pi/8} \sin(\pi/4) = 1,803915$$

Trapèze :

$$\int_0^{\pi/4} e^{3x} \sin 2x \, dx \approx \frac{\pi/4}{2} [0 + e^{3\pi/4} \sin(\pi/2)] = 4,143260$$

Simpson :

$$\int_0^{\pi/4} e^{3x} \sin 2x \, dx \approx \frac{\pi/8}{3} [0 + 4e^{3\pi/8} \sin(\pi/4) + e^{3\pi/4} \sin(\pi/2)] = 2,583696$$

Simpson 3/8 :

$$\begin{aligned} \int_0^{\pi/4} e^{3x} \sin 2x \, dx \approx \frac{3(\pi/12)}{8} [0 + 3e^{3\pi/12} \sin(\pi/6) + 3e^{3\pi/6} \sin(\pi/3) \\ + e^{3\pi/4} \sin(\pi/2)] = 2,585789 \end{aligned}$$

**12.3** Avec la méthode du trapèze et  $n = 1$ , on a  $h = 2 - 0 = 2$  et l'approximation de l'intégrale est

$$f(0) + f(2) = 4.$$

Avec la méthode de Simpson,  $h = (2 - 0)/2 = 1$  et l'approximation est

$$\frac{1}{3} [f(0) + 4f(1) + f(2)] = 2.$$

En remplaçant la première équation dans la seconde et en résolvant, on trouve facilement  $f(1) = 1/2$ .



# Bibliographie

- ANSI. 1986, *ANSI X3.4-1986: Coded Character Set. 7-Bit American Standard Code for Information Interchange*, American National Standards Institute, New York.
- Burden, R. L. et J. D. Faires. 1988, *Numerical Analysis*, 4<sup>e</sup> éd., PWS-Kent, Boston, ISBN 0-5349158-5-X.
- Burden, R. L. et J. D. Faires. 2011, *Numerical Analysis*, 9<sup>e</sup> éd., Brooks/Cole, ISBN 978-0-5387335-1-9.
- Unicode Consortium, T. 2007, *The Unicode Standard, Version 5.0.0*, Addison-Wesley, Boston, ISBN 0-32148091-0.
- IEEE. 2003, *754-1985 IEEE Standard for Binary Floating-Point Arithmetic*, IEEE, Piscataway, NJ.
- ISO. 1998, *ISO/IEC 8859-1:1998. Information Technology – 8-bit Single-Byte Coded Graphic Character Sets – Part 1: Latin Alphabet No. 1*, International Organization for Standardization, Geneva.
- Kernighan, B. W. et P. J. Plauger. 1978, *The Elements of Programming Style*, 2<sup>e</sup> éd., McGraw-Hill, ISBN 0-07034207-5.
- Knuth, D. E. 1997, *The Art of Computer Programming*, vol. 2, Seminumerical Algorithms, Addison-Wesley, Reading, MA.
- Monahan, J. F. 2001, *Numerical Methods of Statistics*, Cambridge University Press, Cambridge, UK, ISBN 0-52179168-5.
- Ripley, B. D. 2005, «Internationalization features of R 2.1.0», *R News*, vol. 5, n<sup>o</sup> 1, p. 2–7. URL <http://cran.r-project.org/doc/Rnews/>.
- Wikipedia. 2012, «IEEE 754-2008 — Wikipedia, The Free Encyclopedia», URL [http://en.wikipedia.org/wiki/IEEE\\_floating\\_point\\_standard](http://en.wikipedia.org/wiki/IEEE_floating_point_standard).

---

Wikipedia. 2013, «Complément à deux — Wikipédia, l'encyclopédie libre», URL  
[http://fr.wikipedia.org/wiki/Complément\\_à\\_deux](http://fr.wikipedia.org/wiki/Complément_à_deux).



Ce document a été produit avec le système de mise en page  $\text{\LaTeX}$ . Le texte principal est en Adobe Utopia 11 points, le code informatique en Bera Mono (un clone de Bitstream Vera) et les titres en Adobe Myriad Pro. Les graphiques ont été réalisés avec R.



