

Licence Creative Commons



Mis à jour le 14 mai 2015 à 18:42

## Au delà des réels: méthodes numériques en informatique





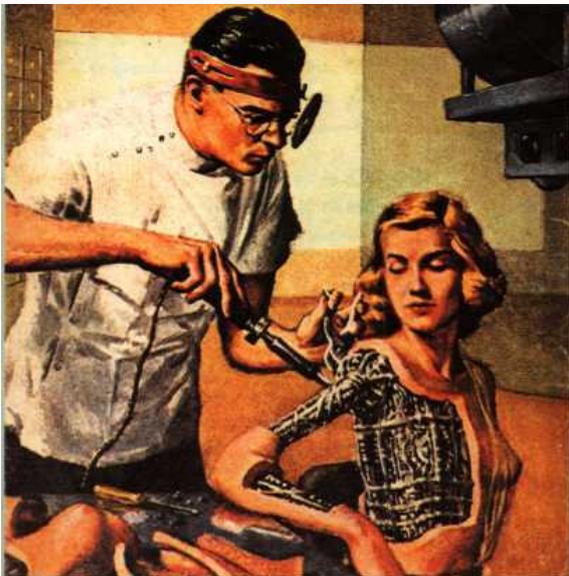
# TABLE DES MATIÈRES

<b>1 IEEE 754-2008</b>	<b>5</b>
1.1 Les nombres à virgule flottante	6
1.1.1 Les problèmes	6
1.1.2 Quelle précision ?	8
1.2 La norme IEEE 754	8
1.2.1 En bref	8
1.2.2 Les normaux, les sous-normaux et les paranormaux	9
1.2.3 Aparté : de l'importance pour un programmeur de bien penser son message d'erreur	12
1.2.4 Successeur d'un VF	13
1.2.5 Reconnaissance des flottants	13
1.2.6 Tableau récapitulatif	13
1.3 Algèbre des nombres VF	14
1.3.1 L'ensemble des VF	14
1.3.2 Comparaison	14
1.3.3 Addition	14
1.3.4 Multiplication	15
1.4 Réels, arrondis et flottants	15
1.4.1 Un problème, dix problèmes, mille problèmes	16
1.4.2 Maîtriser les erreurs	17
1.5 Que la force de l'erreur soit avec vous	19
1.5.1 Atelier Padawan # 1 : majorer l'erreur	19
1.5.2 Atelier Padawan # 2 : quitter la réalité (un peu difficile...)	20
1.5.3 Atelier Padawan # 3 : somme de flottants proches	22
1.5.4 Atelier Padawan # 4 : somme compensée de flottants quelconques	24
1.6 Et la multiplication ? Et la division ? Et le sinus ? Et...	26
1.7 EXERCICES	27
<b>2 Un soupçon de calcul différentiel</b>	<b>37</b>
2.1 Le calcul infinitésimal	38
2.1.1 Infiniment petits	38
2.1.2 Règles de dérivation	39
2.1.3 Dérivée d'un produit	39
2.2 Comment calculer un logarithme sur machine..au XVII <sup>e</sup> siècle ?	40
2.2.1 Calcul d'une approximation d'une racine carrée par la méthode de HÉRON	40
2.2.2 Calcul d'une approximation de log 2 par la méthode de BRIGGS	40
2.2.3 Polynômes interpolateurs	42
2.3 Polynômes et erreurs	44
2.4 Série de TAYLOR	45
2.4.1 Dérivées d'ordre supérieur	45
2.4.2 Relations de domination	46
2.4.3 Polynôme de TAYLOR	48
2.4.4 Formule de TAYLOR-LAGRANGE	49
2.4.5 Formule de Taylor-Mac Laurin	49
2.4.6 Formule de Taylor avec reste intégral	49
2.4.7 Formule de Taylor-Young	49
2.5 Développements limités	49
2.5.1 Voisinage	49
2.5.2 Définition	50
2.6 EXERCICES	51
<b>3 Méthodes itératives</b>	<b>63</b>
3.1 Une boucle sous toutes ses formes	64
3.1.1 Généralité	64
3.1.2 Complexité de la somme	64
3.1.3 Procédure et processus	64
3.1.4 Itération non linéaire : Fibonacci	67
3.2 Exemple : l'algorithme de Babylone	68

3.3	Méthodes itératives pour résoudre une équation	70
3.3.1	Dichotomie (ou méthode de bi-section)	70
3.3.2	Ordre d'une suite et formules de Taylor	71
3.3.3	Méthode de Newton-Raphson-Halley-etc.	71
3.3.4	Étude de la suite associée à l'équation $x^3-2x-5=0$ .	72
3.3.5	Test d'arrêt	72
3.4	Étude expérimentale et théorique de la complexité d'un algorithme	73
3.4.1	Méthode « scientifique »	73
3.4.2	Complexité des boucles	76
3.5	Sommes	78
3.5.1	Notation	78
3.5.2	Somme, boucle et récurrence	79
3.5.3	Manipulation de sommes	79
3.5.4	Sommes multiples	80
3.6	Sommes infinies (séries)	80
3.6.1	Les dangers des ... mal maîtrisés	80
3.6.2	Séries	81
3.6.3	Série harmonique	81
3.6.4	Séries télescopiques	82
3.6.5	Série géométrique	82
3.6.6	Série exponentielle	82
3.7	EXERCICES	83

# 1

# IEEE 754-2008



Des avions qui s'écrasent, des fusées qui s'écrasent, des missiles qui tuent les amis, les indices boursiers totalement faux, les résultats des élections totalement faux, les jeux qui plantent,...

La méconnaissance du traitement des nombres sur machine (et non pas le traitement lui-même !) peut être tragique ou au mieux catastrophique.

Nous allons initier notre exploration de ce monde mystérieux que vous utilisez chaque jour en plongeant dans la machine.

## 1

## Les nombres à virgule flottante

## 1.1 Les problèmes



W.KAHAN (1933 -)

William KAHAN est un mathématicien et informaticien canadien, lauréat du prix TURING en 1989 et principal artisan de la norme IEEE 754 publiée pour la première fois en 1985.

Les problèmes liés à la manipulation des VF (nombres à virgule flottante, ou *floating point (FP)* in english of speciality), sont nombreux.

Citons notamment le célèbre cas du missile américain *Patriot* qui manqua sa cible, un missile *Scud* irakien, durant la première guerre du Golfe en février 1991. Le *Patriot* poursuit sa cible en mesurant le temps mis par les ondes radars pour revenir après rebond sur la cible. Or le système mesurait le temps en dizaines de seconde (drôle d'idée) puis multipliait ce temps par 0,1 pour obtenir le temps en secondes et utilisait un registre à 24 bits en virgule fixe.

Après cent heures d'utilisation (l'armée US ne chôme pas), un écart de 0,34 seconde est apparu (pourquoi? cf en TD...). Un *Scud* volant à 1.676 m/s, le *Patriot* a visé 500 m à côté. Cette erreur causa la mort de 28 soldats tandis que 98 étaient blessés.

Il y a eu beaucoup d'autres cas célèbres, même chez les plus puissants :

Google 5.0 - 4.9 - 0.1

Web Shopping Vidéos Images Actualités Plus Outils de recherche

Environ 68 500 000 résultats (0,11 secondes)

Conseil : Recherchez des résultats uniquement en français. Vous pouvez indiquer votre langue de recherche sur la page Préférences.

5.0 - 4.9 - 0.1 =

-3.6082248e-16

Rad x! ( ) % AC

Bon, ne soyons pas de mauvaise foi :

```
1 *Main> 5 - 4.9 - 0.1
2 -3.608224830031759e-16
```

Même des logiciels spécialisés, très chers (Maple ©), se trompent :

```
1 > evalf(sin(2^100));
2 .4491999480
3 > evalf(sin(2^100), 20);
4 -.58645356896925826300
5 > evalf(sin(2^100), 30);
6 .199885621653625738215132811525
7 > evalf(sin(2^100), 40);
8 -.8721836054182673097807197782134705593243
```

Heureusement, le libre sauve la mise...Giac/XCAS :

```
1 1>> evalf(sin(2^100))
2 -0.872183605418
```

Haskell

```

1 *Main> sin(2^100)
2 -0.8721836054182673

```

Un exemple rigolo : les élections au Schleswig-Holstein en 1992 (merci à Paul ZIMMERMANN).  
Les partis ayant moins de 5% des voix n'ont pas de siège.

Au dépouillement, le parti écologiste obtient exactement 5%.

Après l'annonce officielle des résultats, on découvre que le parti écologiste n'a en fait que 4.97% des voix !

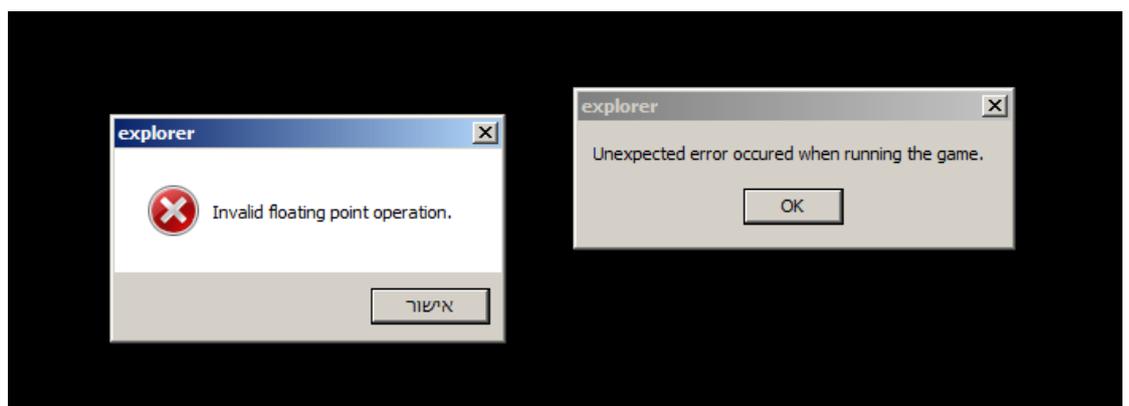
Le programme affichant les résultats arrondi à un seul chiffre après la virgule, et donc a arrondi 4.97 à 5.0.

L'erreur corrigée, le siège est retiré aux verts, et attribué au SPD, qui avec 45 sièges sur 89 obtient la majorité absolue.

Un autre exemple qui fait pitié (merci à William KAHAN) :

	A	B
1	B1 = 4/3	1,33333333333333000000
2	B2 = B1 - 1	0,33333333333333000000
3	B3 = B2 * 3	1,00000000000000000000
4	B4 = B3 - 1	0,00000000000000000000
5	B5 = B4 * 2 <sup>52</sup>	0,00000000000000000000
6	(4/3 - 1) * 3 - 1	0,00000000000000000000
7	((4/3 - 1) * 3 - 1)	-0,000000000000000022204
8	((4/3 - 1) * 3 - 1) * 2 <sup>52</sup>	-1,00000000000000000000

Un dernier qui vous touchera plus :



Le problème des VF est apparu sur machine en 1941 : Konrad ZUSE a proposé son Z3 avec une mantisse de 14 bits, un exposant de 7 bits, un bit de signe le tout en base 2.

En effet, les humains préfèrent compter en base 10...et n'arrêtent pas de faire des erreurs d'arrondis suite à leur manie (pourquoi ? cf en TD...). Certaines machines calculent donc en base 10 (les calculettes, Maple©, les machines dédiées à la Finance,...). La norme IEEE 754-2008 introduit aussi des calculs en base 10.

La machine pour l'instant est basée sur une logique à deux états et préfère donc la base 2 (ou 8 ou 16).



Николай Петрович  
Брусенцов(1925 -)

Certaines machines ont pourtant été conçues selon une logique à trois états (Vrai, Faux, j'en sais rien) comme le *Сетунь* soviétique mis au point en 1958 par une équipe dirigée par Николай Петрович Брусенцов et le célèbre mathématicien Сергей Львович Соболев. L'histoire de cette machine est édifiante : la production des *Сетунь* fut confiée à l'*usine des machines mathématiques* qui en fabriqua très peu car le *Сетунь* était trop bon marché ! De plus il a fonctionné des années durant sans aucune panne ni besoin de pièces de rechange ce qui n'est pas très bien vu dans le milieu informatique :-). Enfin le *Сетунь* réglait le problème crucial du *tiers-exclus* et obéissait aux canons les plus exigeants d'E. DIJKSTRA. Bref, les autorités soviétiques ont décidé de le découper en morceaux et de mettre les restes dans une décharge. De zélés informaticiens ont toutefois réussi à en sauver un exemplaire en le cachant dans un grenier.

Vous utilisez vous-même une logique ternaire en SQL pour gérer NULL. Les machines quantiques peuvent travailler avec des *qtrits*. Vous verrez ça un jour, jeunes padawans.

Autre problème : le profusion d'appellations qui sèment la confusion ! En C (cf dinechin) :

- char est un entier codé sur 8 bits. C'est l'abréviation du nom commun *character* ;
- int est l'abréviation du non commun *integer* qui désigne un nombre entier comme en mathématiques mais cependant  $2147483647 + 1 = -2147483648...$
- short et long sont des adjectifs désignant des entiers de taille différentes ;
- float est un verbe ;
- single ? Simple quoi ?
- double ? Mais le double de quoi ?
- long double : deux adjectifs à suivre ?

Alors il faut clarifier tout ça !

### 1 2 Quelle précision ?

En physique, on ne peut guère aller en deçà de grandeurs de l'ordre de  $10^{-15}$  alors pourquoi s'embêter ?

```

1 *Main> 3 * 0.1
2 0.300000000000000004
3 *Main> sum $ take 10000000 $ repeat 0.1
4 999999.9998389754
```

Il faut parfois aller bien plus loin que  $10^{-15}$  pour avoir une précision de  $10^{-15}$  : c'est bien là le problème. On cherchera avec attention l'exercice [Recherche 1 - 6 page 27](#)

Pourtant, avec 50 bits on peut coder la distance Terre-Lune avec une erreur de l'ordre de la taille d'une bactérie !

Qu'est-ce qui se passe ? On ne peut plus faire de calcul ? On ne pourra pas fabriquer de jeux vidéos pour dégommer plein de méchants ?

Pas de panique : il faut juste avoir en tête que **LES NOMBRES RÉELS N'EXISTENT PAS : TOUT CE QUE VOUS AVEZ VU AU LYCÉE N'EST QU'ILLUSION !**

Oubliez donc les réels : les nombres que vous allez manipuler sont différents mais sont tout aussi maîtrisables. Il suffit de connaître leur arithmétique, l'**arithmétique des nombres à virgule flottante** qui obéissent à la norme IEEE 754.

Des problèmes subsistent néanmoins. Par exemple, on sait très bien ce que donnera un calcul selon un format donné...le problème est qu'on ne sait pas toujours quel format est utilisé...

## 2 La norme IEEE 754

### 2 1 En bref

Nous nous contenterons, sauf mention contraire, de travailler en base 2. Vous en avez parlé en archi : un VF est représenté par un bit de signe, une mantisse (ou significande en français) et un exposant :

$$v = (-1)^s \times m \times 2^E$$

Les formats habituels sont le *binary32* ou *single* avec  $(\#E, \#m) = (8, 24)$  et le *binary64* ou *double* avec  $(\#E, \#m) = (11, 53)$ .

Pour simplifier nos présentations, on illustrera souvent nos propos avec des représentations *toy7* plus petites, comme par exemple  $(\#E, \#m) = (3, 4)$ .

La mantisse doit dans la mesure du possible vérifier  $1 \leq m < 2$  ce qui minimise l'exposant et apporte plus d'informations. Par exemple, en *toy7* :

$$\begin{aligned} v &= 0,010\mathbf{01} \times 2^0 \\ v &= 0,100\mathbf{1} \times 2^{-1} \\ v &= 1,00\mathbf{1} \times 2^{-2} \end{aligned}$$

C'est la *forme normale* d'un VF. De plus, cela se permet de se passer du premier 1 qui est implicite. On notera

$$m = 1, f$$

On peut également gagner de la place en ne stockant pas le signe de l'exposant : il suffit de le translater de  $2^{(\#E)-1} - 1$ ...POURQUOI ?

$$e = E + 2^{(\#E)-1} - 1$$

Représentons  $0,75_{10}$  en *toy7*.

$$0,75_{10} = \frac{3_{10}}{2^2_{10}} = 11_2 \times 2^{-2} = 1,100_2 \times 2^{-1}$$

L'exposant stocké  $e$  vérifie avec le décalage :

$$e - (2^{3-1} - 1) = E = -1 \leftrightarrow e = -1 + 3 = 2 = 10_2$$

Ainsi :

s (1 bit)	e (3 bits)			f (3 bits)		
0	0	1	0	1	0	0

Exemple

On doit également s'occuper des exposants extrêmes qui ne correspondent pas à des formes normales (cf TD).

**Petit algo pour convertir un nombre de valeur absolue inférieure à 1 en base 2**

Soit  $x = a_1 \times 2^{-1} + a_2 \times 2^{-2} + a_3 \times 2^{-3} + \dots + a_n \times 2^{-n}$  l'écriture de  $x$  tel que  $|x| < 1$  en base 2 avec les  $a_i$  égaux à 0 ou 1.

Alors  $2x = a_1 + a_2 \times 2^{-1} + a_3 \times 2^{-2} + \dots + a_n \times 2^{-n+1}$  :  $a_1$  est donc la partie entière de  $2x$ .

Ensuite  $2(2x - a_1) = a_2 + a_3 \times 2^{-1} + \dots + a_n \times 2^{-n+2}$  :  $a_2$  est alors la partie entière de  $2(2x - a_1)$  und so weiter...

À retenir

## 2.2 Les normaux, les sous-normaux et les paranormaux

Nous venons de voir que l'exposant  $E$ , qui est un entier signé, est stocké sur  $\#E$  bits sous la forme  $e = E + 2^{(\#E)-1} - 1$ .

Comme  $e$  est un entier naturel stocké sur  $\#E$  bits, il varie entre  $e_{\min} = 0$  et  $e_{\max} = 2^{\#E} - 1$ .

Par exemple, en *toy7*,  $\#E = 3$  donc  $e_{\min} = 000$  et  $e_{\max} = 111 = 1000 - 1$ .

CEPENDANT les valeurs extrêmes de  $E$  sont réservées pour des cas spéciaux.

Cela donne :

$$E_{\min} = (e_{\min} - 2^{(\#E)-1} + 1) + 1 = 0 - 2^{3-1} + 1 + 1 = -2$$

$$E_{\max} = (e_{\max} - 2^{(\#E)-1} + 1) - 1 = (2^3 - 1) - 2^{3-1} + 1 - 1 = 7 - 4 + 1 - 1 = 3$$

En effet, des problèmes apparaissent avec les formes normales :

- on ne peut pas représenter zéro sous forme normale ;



— des nombres du type  $0.\underbrace{000\dots00}_{\#E}11 \times 2^0 = 1.1 \times 2^{-(\#E+1)}$  ne sont pas représentables sous forme normale car il faudrait un exposant plus petit que  $E_{\min}$ .

La norme IEEE 754 a donc introduit les nombres *sous-normaux* pour remplir le vide qui existerait entre 0 et le plus petit nombre normal (cf exercice Recherche 1 - 10 page 28).

Pour les signaler et les distinguer des normaux, la norme impose que leur exposant décalé  $e$  soit 0 :  $e_{\min}$  est donc une valeur réservée.

Un cas particulier : si  $e = e_{\min}$  et  $f = 0$ , alors la norme ne considère ce nombre ni comme un nombre normal ni comme un nombre sous-normal. C'est la représentation de zéro.

Il y a deux zéros! En *toy7* :

+0 :

s (1 bit)	e (3 bits)	f (3 bits)
0	0 0 0	0 0 0

-0 :

s (1 bit)	e (3 bits)	f (3 bits)
1	0 0 0	0 0 0

Des tests  $x == 0$  pourraient alors échouer ! La norme impose donc que les deux zéros soient considérés comme égaux :

**À retenir**

```

1 *Main> let x = 1 / 1e500
2 *Main> x
3 0.0
4 *Main> x == 0
5 True
6 *Main> let y = -1 / 1e500
    
```

```

1 *Main> y
2 -0.0
3 *Main> y == 0
4 True
5 *Main> x == y
6 True
    
```



Reprenons à présent un exemple classique : vous voulez calculer une norme  $N = \sqrt{x^2 + y^2}$  avec  $x = 1 \times 2^3$  et  $y = 1,1 \times 2^3$  en *toy7*.

Les carrés de ces nombres posent problème : leur exposant est trop grand. On est dans un cas de *dépassement* ou d'*overflow* en anglais de spécialité.

Une première idée serait de remplacer  $x^2$  et  $y^2$  par le plus grand nombre disponible  $1,111 \times 2^4$  soit :

s (1 bit)	e (3 bits)	f (3 bits)
0	1 1 1	1 1 1

Alors  $x^2$  puis  $y^2$  puis  $x^2 + y^2$  seraient remplacés par  $1,111 \times 2^4$ .

On aurait donc  $N = (1,111 \times 2^4)^{1/2}$  arrondi à  $1,010 \times 2^2$  ce qui est dramatiquement faux et pourraient envoyer des *Patriots* sur le toit du bâtiment B.

La norme introduit donc deux infinis codés en *toy7* :

$+\infty$  :

s (1 bit)	e (3 bits)	f (3 bits)
0	1 1 1	0 0 0

$-\infty$  :

s (1 bit)	e (3 bits)	f (3 bits)
1	1 1 1	0 0 0

Soit, plus généralement,  $e = e_{\max}$  et  $f = 0$ .

On retrouve alors les limites habituelles vues au lycée (ah, quand même...) :

```

1 *Main> 1 + 1e500
2 Infinity
3 *Main> let x = 1e500
4 *Main> x
5 Infinity
6 *Main> 1 + x
7 Infinity
    
```

```

1 *Main> x^3
2 Infinity
3 *Main> 1 / x
4 0.0
5 *Main> 4 - x
6 -Infinity
    
```

Nous verrons bien d'autres exemples en exercice.



Vous avez aussi parlé au lycée des **FORMES INDÉTERMINÉES** lors de l'étude des limites. Par exemple, que dire de :

$$\lim_{x \rightarrow +\infty} x^2 - x$$

Vous êtes muets ? Demandons à la machine :

```
1 *Main> x^2 - x
2 NaN
```

Une manière exotique de régler le problème est d'étudier l'algorithme suivant :

<http://www.food.com/recipe/madhur-jaffreys-naan-bread-446809>

Sinon, il suffit de savoir que la Hitroizeu introduit le NaN (comme dans *Not a Number*) dans les cas où le calcul demandé est « indéfini ».

<pre>1 *Main&gt; let x = 1e500 2 *Main&gt; x - x 3 NaN 4 *Main&gt; x / x 5 NaN 6 *Main&gt; sqrt(-1) 7 NaN 8 *Main&gt; 0 / 0 9 NaN</pre>	<pre>1 -- Bizarre ? 2 *Main&gt; x - x 3 NaN 4 *Main&gt; x - x == x - x 5 False 6 *Main&gt; x^2 - x 7 NaN 8 *Main&gt; x * (x - 1) 9 Infinity</pre>
---	---

Un NaN n'est pas comparable. Trier un tableau comportant un NaN ne pose donc pas de problème :

```
1 *Main> let x = 0/0
2 *Main> x
3 NaN
4 *Main> x > 3
5 False
6 *Main> x < 3
7 False
8 *Main> x == 3
9 False
10 *Main> x == x
11 False
```

Il faut penser à un problème célèbre dont fut cette fois victime la marine étasunienne.

#### Aparté

L'affaire est beaucoup moins dramatique mais plus croustillante que celle des *Patriots*. La Marine US s'équipe de PC montés avec...Windows NT 4.0 et en profite pour réduire le 10% le nombre de marins à bord du USS Yorktown. Mais, comme le dira un des experts civils après la panne : « *Using Windows NT, which is known to have some failure modes, on a warship is similar to hoping that luck will be in our favor* »

Que s'est-il passé le 21 septembre 1997 ? Un opérateur a tapé par erreur un zéro sur son clavier et tout le navire a été bloqué en pleine mer :-)

« *Your \$2.95 calculator, for example, gives you a zero when you try to divide a number by zero, and does not stop executing the next set of instructions. It seems that the computers on the Yorktown were not designed to tolerate such a simple failure.* »

On espère que les missiles atomiques sont mieux programmés...

Revenons sur cette anecdote et transplantons là dans un environnement plus pacifique.

Vous cherchez les zéros de, disons, comme ça, au hasard,  $ax^2 + bx + c$ . Comme vous avez de vagues souvenirs de votre jeunesse, vous programmez : retourne-moi  $-b \pm \frac{\sqrt{\Delta}}{2a}$  avec  $\Delta = b^2 - 4ac$ .

On suppose que toutes les valeurs sont proches de zéro et les mesures difficiles à ce niveau de précision.

Un premier essai donne :

```

1 *Main> let a = 1e-6
2 *Main> let b = 1e-4
3 *Main> let c = 1e-10
4 *Main> let d = b^2 - 4*a*c
5 *Main> d
6 9.9999996e-9
7 *Main> -b - sqrt(d)/(2*a)
8 -50.000099
9 *Main> sqrt(d)
10 9.999999799999999e-5

```

Un deuxième essai, très proche, donne :

```

1 Main> let a = 1e-6
2 *Main> let b = 1e-10
3 *Main> let c = 1e-4
4 *Main> let d = b^2 - 4*a*c
5 *Main> d
6 -3.9999999999e-10
7 *Main> -b - sqrt(d)/(2*a)
8 NaN
9 *Main> sqrt(d)
10 NaN

```

Ouf! Le *Yorktown* ne coulera pas!

Ainsi, on n'a plus besoin de surcharger le programme de conditionnelles. Le calcul continue.

Reprenons nos formes indéterminées :

```

1 *Main> let x = 1e500
2 *Main> x^2 - x
3 NaN

```

En effet,  $\infty - \infty$  est une forme indéterminée. Que faisiez-vous au lycée pour lever cette indétermination ?

```

1 *Main> x*(x - 1)
2 Infinity

```

C'est réglé... Autre curiosité :

```

1 *Main> x - x
2 NaN
3 *Main> x - x == x - x
4 False

```

On peut débusquer un NaN en effectuant le test  $x \neq x$ .

Selon la Hitroizeu, un NaN est représenté avec l'exposant maximum et une mantisse non nulle : il y a donc toute une famille de NaN! En voici un exemple en *toy7* :

s (1 bit)	e (3 bits)	f (3 bits)
1	1 1 1	0 1 0

### 2 3 Aparté : de l'importance pour un programmeur de bien penser son message d'erreur

Encore un drame...Le 1<sup>er</sup> juin 2009, le vol Air France 447 finit tragiquement : les 228 personnes à bord meurent englouties au large des côtes du Brésil.

Beaucoup de bruit a été fait autour de cet accident, beaucoup de gens ont été blâmés, en particulier les pilotes.

Voici le problème : pour économiser du carburant, les avions volent maintenant très haut et en limite de décrochage. Il faut un guidage informatique très fin pour le permettre. Ce soir-là, l'avion est pris dans une tempête et le froid extrême obstrue les capteurs qui ne peuvent plus mesurer la vitesse du vent.

Le pilote automatique envoie alors comme message à ses collègues humains des messages trop vagues :

*Les messages de panne successivement affichés sur l'ECAM n'ont pas permis à l'équipage de faire un diagnostic rapide et efficace de la situation dans laquelle l'avion se trouvait, en particulier de l'obstruction des sondes Pitot. Il n'a jamais été en mesure de faire le lien entre les messages qui sont apparus et la procédure à appliquer, alors que la lecture*



de l'ECAM et des messages doit faciliter l'analyse de la situation et permettre d'organiser le traitement des pannes. Plusieurs systèmes avaient pourtant identifié l'origine du problème mais n'ont généré que des messages (extrait du rapport final du Bureau d'Enquêtes et d'Analyses pour la sécurité de l'aviation civile).

Donner les conséquences et pas les causes! Sans indications, les pilotes perdent de précieuses secondes et ne peuvent effectuer la manœuvre habituelle permettant de contrer un décrochage.

**2 4 Successeur d'un VF**

Il n'est pas possible de définir un successeur d'un réel mais un VF, lui, en possède un. Écrivons un flottant avec sa mantisse entière en posant

$$n = \#f$$

**Remarque**

Dans la littérature, on considère souvent  $p = n + 1$ .

Par exemple, pour *toy7*,  $n = 3$ , pour *binary32*,  $n = 23$ , pour *binary64*,  $n = 52$ .

$$v = M(v) \times 2^{E(v)-n}$$

$M(v)$  étant un entier, il admet un successeur. Si l'on reste dans les limites des VF normaux du système :

$$\text{succ}(v) = (M(v) + 1) \times 2^{E(v)-n} = v + 2^{E(v)-n}$$

On note

$$\text{ulp}(v) = 2^{E(v)-n}$$

**À retenir**

**ULP**

Par exemple, en *binary64*, soit  $x$  un réel et  $y = RN(x)$ .

On a  $y = (-1)^s \times 1, b_1 b_2 \dots b_{52} \times 2^E$ .

Alors  $\text{ulp}(y) = 0, 000 \dots 001 \times 2^E = 2^{E-52}$

**2 5 Reconnaissance des flottants**

La machine traite une chaîne de bits 

$s$	$e$	$f$
-----	-----	-----

 ainsi :

$e = 0 \dots 0$	$f = 0$	$\rightarrow v = (-1)^s \times 0.0$
$e = 0 \dots 0$	$f \neq 0$	$\rightarrow v = (-1)^s \times 0.f \times 2^{1-E_{\max}}$
$e = 1 \dots 1$	$f = 0$	$\rightarrow v = (-1)^s \times \infty$
$e = 1 \dots 1$	$f \neq 0$	$\rightarrow \text{NaN}$
$0 \dots 0 < e < 1 \dots 1$	$f \neq 0$	$\rightarrow v = (-1)^s \times 1.f \times 2^{e-E_{\max}}$

**2 6 Tableau récapitulatif**

On notera  $\epsilon_m = 2^{-n}$  l'epsilon de la machine, c'est-à-dire le successeur de 1.

On notera  $\lambda$  le plus petit VF normal positif.

On notera  $\mu$  le plus petit VF sous-normal positif.

On notera  $\Omega$  le plus grand VF normal.

**Recherche**

Quelle relation existe-t-il entre  $\mu$ ,  $\epsilon_m$  et  $\lambda$  ?

Format	$\#E$	$\#f$	$E_{\min}$	$E_{\max}$	$\epsilon_m$	$\lambda$	$\mu$	$\Omega$
<i>toy7</i>	3	3	-2	+3	$2^{-3} = 1/8$	$2^{-2} = 1/4$	$2^{-5} = 1/32$	$1, 111 \times 2^3 = 15$
<i>binary32</i>	8	23	-126	+127	$2^{-23} \approx 1,2 \times 10^{-7}$			
<i>binary64</i>	11	52	-1022					

### 3 Algèbre des nombres VF

#### 3 1 L'ensemble des VF

Nous allons à présent travailler dans notre ensemble de flottants normaux et sous-normaux. Nous noterons par exemple  $\mathbb{V}_{32}$  les normaux et sous-normaux de *binary32* et nous noterons  $\overline{\mathbb{V}_{32}}$  ce même ensemble enrichi des deux infinis.

Dans le cas général, nous noterons tout simplement  $\mathbb{V}_b$  voire  $\mathbb{V}$ .

À retenir

Ainsi  $\overline{\mathbb{V}_b}$  est un ensemble FINI.

On va munir cet ensemble d'opérations usuelles et étudier les structures algébriques produites.

#### 3 2 Comparaison

Recherche

Il est très simple et rapide de comparer deux VF : comment la machine procède-t-elle ? Quel est l'avantage de ce stockage des VF ?

#### 3 3 Addition

On décompose l'action en trois étapes :

1. on commence par ramener les deux nombres au même exposant, en l'occurrence le plus grand des deux ;
2. on ajoute les deux mantisses *complètes* en tenant compte du signe ;
3. on renormalise le nombre obtenu.

Par exemple, en *toy7*, additionnons 1,1 et 0,0111 ou plutôt  $1,1 \times 2^0$  et  $1,11 \times 2^{-2}$  ou plutôt (en *toy7*, le décalage d'exposant est de 3) :

1,1 : 

0	0	1	1	1	0	0
---	---	---	---	---	---	---

0,0011 : 

0	0	0	1	1	1	0
---	---	---	---	---	---	---

On réécrit 0,00101 avec le même exposant que 1,1 mais attention à la précision de la mantisse !

1,100

0,00111

----

1,10111

On ne peut pas garder les deux derniers 1 du plus petit nombre car on n'a pas la place de les stocker !

Que faire ? Il faut arrondir en choisissant selon quatre méthodes usuelles :

Les arrondis

1. l'arrondi au plus proche (RN) qui arrondit au VF...le plus proche. En cas d'égalité, on choisit la valeur paire (donc qui se termine par un 0 en binaire) ;
2. l'arrondi vers 0 (RZ) qui arrondit à la valeur de plus petite valeur absolue : c'est la troncature ;
3. l'arrondi vers  $+\infty$  (RU) qui arrondit à la valeur supérieure la plus petite ;
4. l'arrondi vers  $-\infty$  (RD) qui arrondit à la valeur inférieure la plus grande.

Le mode d'arrondi par défaut est le premier.

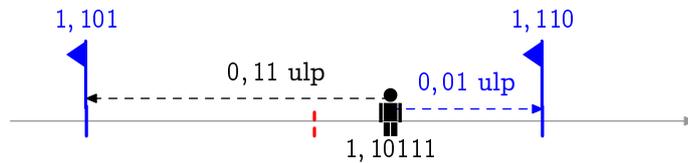
Définition 1 - 1

Posons  $x = 1,10111$  (qui n'est pas un VF!) : il est compris entre les deux VF 1,101 et 1,110.

$$1,10100 \leq \underbrace{1,10111}_x \leq 1,11000$$

$$\underbrace{1,100 + 1,00ulp(x)}_v \leq \underbrace{1,100 + 1,11ulp(x)}_x \leq \underbrace{1,100 + 10,00ulp(x)}_{\text{succ}(v)}$$

$$|x - v| = 0,11ulp(x), \quad |x - \text{succ}(v)| = 0,01ulp(x) \quad \text{donc } RN(x) = \text{succ}(v)$$



$$RN(1, 1 + 0,001) = 1, 110$$

Il y a énormément de choses à dire sur ces problèmes d'arrondis. Nous leur consacrerons une section.

Il y a beaucoup de choses à dire sur les sommes itérées. Nous leur consacrerons une section.

Recherche

Est-ce que l'addition des VF est associative ?  
 Est-on sûr de l'ordre dans lequel un compilateur calcule  $a + b + c + d$  ?

**3 4 Multiplication**

On suit un ordre logique :

1. on « xore » les bits de signe ;
2. on additionne les exposants réels et on décale ou plutôt on additionne les exposants décalés et on retire la valeur d'un décalage ;
3. on multiplie les mantisses ;
4. on normalise.

Effectuons par exemple le produit de 101,1 par -10,01, i.e. :

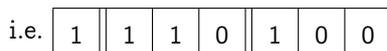


1. 0 xor 1 donne 1 : le bit de signe est 1 ;
2.  $101 + 100 - 11 = 1001 - 11 = 110$  : l'exposant décalé est 110 donc l'exposant est 3 ;
3.  $1,011 \times 1,001 = 1,011 + 1,011 \times 0,001 = 1,011 + 0,001011 = 1,100011$  ;
4. le produit est donc  $1,100011 \times 2^3$ . Le passage à la forme normale va arrondir le produit.  
 On a

$$1, 100 < x < 1, 101$$

avec  $|x - 1, 100| = 0,011ulp(x)$  et  $|x - 1, 101| = 0,101ulp(x)$  donc

$$RN(101, 1 \times (-10,01)) = -1,100 \times 2^3$$



Ici encore, le VF ne correspond pas au réel...

Il est temps de prendre le problème à bras le corps, mais avant :

Recherche

Est-ce que la multiplication des VF est associative ?  
 Est-ce que la multiplication des VF est distributive sur l'addition ?

**4 Réels, arrondis et flottants**

Dans la suite de nos aventures, nous ne considérerons que le mode d'arrondi par défaut sauf mention contraire.

**4 1 Un problème, dix problèmes, mille problèmes**

Nous commençons à bien connaître le problème :

```
1 *Main> 3*0.1 == 0.3
2 False
```

On pourrait alors se dire qu'il ne faut jamais comparer deux VF  $x$  et  $y$  avec  $==$  mais plutôt les considérer égaux s'ils vérifient  $|x - y| < d$  avec  $d$  un seuil d'erreur. Mais résoudre ce problème ainsi va en créer de nombreux autres :

- comment fixer  $d$  ?
- si  $x$  et  $y$  ont des signes différents ?
- $x \equiv y \leftrightarrow |x - y| < d$  définit-elle une relation d'équivalence ?
- ...

**Recherche**

Répondez à la troisième question...

Il faut penser différemment !

Beaucoup de propriétés des réels ne sont plus vraies dans  $\mathbb{V}$ . Par exemple, le programme suivant, que nous allons étudier, n'a aucun intérêt dans  $\mathbb{R}$  mais est primordial dans  $\mathbb{V}$  :

```
1 f1 a =
2   if (a + 1.0) - a /= 1.0 then a
3   else f1 (2.0 * a)
4
5 f2 a b =
6   if (a + b) - a == b then b
7   else f2 a (b + 1.0)
8
9 b = f2 (f1 1) 1
```

**Danger**

Mais comme souvent en mathématique et en informatique, un petit dessin aide à visualiser une situation complexe.

**Danger**

Attention : tout VF est en fait un rationnel mais tout rationnel n'est pas un VF ! Par exemple...

La seule certitude est qu'un VF se représente correctement...lui-même.

Mais je vous rappelle que ce que vous tapez sur votre clavier n'est pas forcément un VF si vous tapez en base 10.

En particulier, si  $x$  est un VF,  $x = RN(x)$ .

La norme assure de plus que pour toute opération arithmétique standard,

$$RN(x) \top RN(y) == RN(x \top y)$$

Il faut aussi penser que rien n'est gratuit :

**À retenir**

S'il y a un *epsilon* dans votre calcul, c'est à vous de vous en occuper et de le contrôler : la norme vous donne les moyens de contrôler mais ce n'est pas la machine qui fait le boulot...c'est VOUS !

N'oubliez pas que vous avez un avantage sur la machine : vous connaissez la valeur réelle du réel que vous manipulez ou faites manipuler.

Voici une petite réflexion encore plus subtile...

**Danger**

$0,001 \times 2^0$  ce n'est pas tout à fait pareil que sa forme normalisée  $1,000 \times 2^{-3}$ .

Les 3 zéros rajoutés à droite n'ont sûrement aucune valeur !

Par exemple si on calcule  $11,111 - 11,110$ , mais sur 6 bits, on aurait peut-être un autre résultat qui dépend de bits ici inconnus...

Il y a donc plusieurs échelle de certitude !

Quand des nombres sont très proches, les soustraire peut engendrer ce genre de phénomène qu'on appelle *annulation catastrophique* : l'*erreur relative* est bien plus grande que l'*erreur absolue* (nous reparlerons de ces notions un peu plus loin.)

On pourra s'intéresser à l' [Recherche 1 - 26 page 30](#) pour méditer sur le danger relatif de certaines annulations (cancellations in anglische) : toute annulation n'est pas catastrophique !

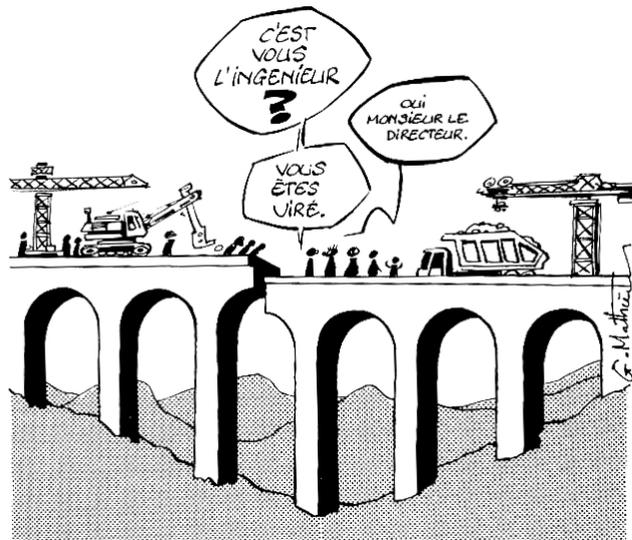
**Danger**

Les machines calculent vite, c'est un bien mais aussi un mal : si on effectue une somme avec une toute petite erreur, une grande somme peut nous faire perdre tous nos bits en une seconde car la somme de beaucoup de très petites erreurs peut créer une grosse erreur...

Mais le plus grand danger vient de :

*95 % of the folks out there are completely clueless about floating-point*

James GOSLING (M. Java) - 28 février 1998

**4 2 Maîtriser les erreurs**

Tout d'abord, une déception : il n'existe pas de traitement automatique des erreurs sur les VF.

Cependant un espoir : il existe des moyens de traiter les problèmes selon le contexte...mais qui nécessitent un bagage minimum en mathématique;-)

D'abord un peu de vocabulaire. Il faudra bien distinguer :

**Précision** (*precision*) : c'est le nombre de bits utilisés pour représenter un nombre. La précision concerne donc le format utilisé pour écrire ou stocker ou arrondir un nombre. Par exemple 3, 3.0, 3.0000, 3.0e0 n'ont pas la même précision, précision qui dépend en plus du langage.

**Exactitude** (*accuracy*) : c'est ce qui relie un nombre au contexte dans lequel il est employé. Connaître la distance qui sépare votre domicile de votre lieu de travail au mètre près peut être considéré comme très proche du résultat exact mais connaître votre taille au mètre près l'est moins...L'exactitude est liée à la mesure de l'erreur.

Nous définirons deux types de mesure d'erreur :

#### Erreurs

Soit  $x$  un nombre et  $\hat{x}$  le nombre qui le représente. On distingue :

#### Définition 1 - 2

l'erreur absolue  $|x - \hat{x}|$

l'erreur relative  $\eta = \frac{x - \hat{x}}{x}$  alors  $\hat{x} = x(1 + \eta)$

commise en prenant  $\hat{x}$  à la place de  $x$ .

On prend souvent comme approximation de l'erreur relative le rapport  $|\log(x) - \log(\hat{x})|$ . Nous justifierons cette approximation lors de l'étude du calcul différentiel.

Pour reprendre un exemple de W. ?,  $3,1777777777777777$  est une approximation plutôt précise (16 décimales) mais inexacte (2 décimales) de  $\pi$ .

Comme nous allons le voir dans de nombreux exemples (e.g. [Recherche 1 - 26 page 30](#)), analyser les erreurs avec un peu de mathématique permet de contrôler son programme.

On peut mesurer les erreurs à deux niveaux. Lorsqu'on étudie un algorithme qui donne un résultat  $r$  en fonction d'une donnée  $d$ , on peut considérer que  $r$  est une fonction de  $d$  :  $r = f(d)$ .

On peut analyser l'erreur commise sur  $r$  ou l'erreur commise sur  $d$ .

#### Définition 1 - 3

L'erreur *aval* (*forward error*) est la différence entre le résultat mesuré  $\hat{r}$  et le résultat théorique  $r$ .

L'erreur *amont*, ou erreur *inverse* (*backward error*) est le plus petit  $|\Delta d|$  tel que  $f(d + \Delta d) = \hat{r}$ .

*Feel nervous, but feel in control. It's not dark magic, it's science.*

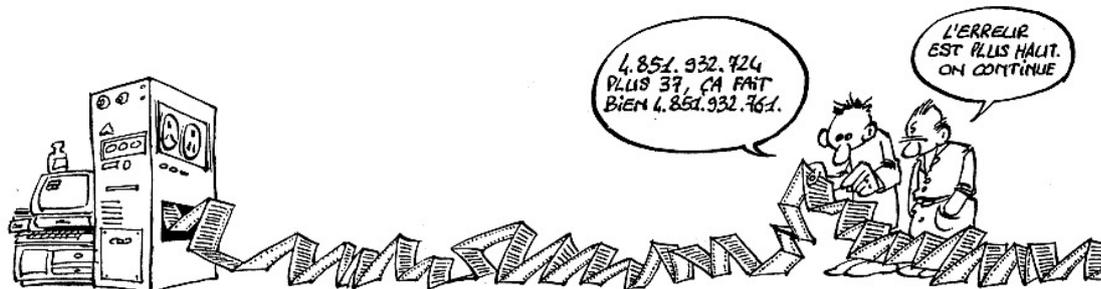
Florent DE DINECHIN

C'est ce peu de mathématique que nous allons découvrir dans ce cours...

Mais un point non négligeable que nous n'aborderons pas (nous l'évoquerons) est comment le processeur, l'OS, le langage (mal)traitent la norme IEEE 754...

#### À retenir

La norme ne vous transforme pas en Harry POTTER : il ne suffit pas de se dire qu'il y a une norme et que tout va bien se passer en appuyant sur son clavier-baguette magique ! La norme vous donne les moyens de maîtriser les forces du mal mais ça peut demander un sacré effort de réflexion (c'est-à-dire de mathématique :- ) .



## 5

## Que la force de l'erreur soit avec vous

## 5 1 Atelier Padawan # 1 : majorer l'erreur

## 5 1 1 Précision machine



Soit  $x$  un nombre. On peut l'encadrer par deux VF successifs :

$$M \times 2^{E-n} \leq x < (M+1) \times 2^{E-n} = M \times 2^{E-n} + 2^{E-n}$$

avec  $M$  un entier,  $n$  la longueur de la pseudo-mantisse.

L'approximation  $\hat{x}$  est donc la plus proche de ces deux bornes de cet intervalle de longueur  $2^{E-n}$ . Avec l'arrondi par défaut :

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E-n}$$

donc l'erreur relative peut être majorée

$$\left| \frac{x - \hat{x}}{x} \right| \leq \left| \frac{x - \hat{x}}{m \times 2^E} \right| \leq \frac{1}{2} \frac{2^{E-n}}{m \times 2^E} = \frac{1}{2} \frac{2^{-n}}{m} = \frac{1}{2} \frac{\varepsilon_M}{m}$$

avec  $m$  la mantisse normalisée.

Si  $\hat{x}$  n'est pas un VF sous-normal, alors  $1 \leq m < 2$  donc

$$\left| \frac{x - \hat{x}}{x} \right| \leq \frac{1}{2} \varepsilon_M$$

L'erreur relative commise est donc majorée par la moitié du successeur de 1, i.e. l'*epsilon* de la machine.

Si  $\hat{x}$  est sous-normal, alors l'erreur absolue est majorée :

$$|x - \hat{x}| \leq \frac{1}{2} 2^{E_{\min} - 1 - n}$$

Ainsi, dans tous les cas :

## Précision

$\hat{x} = x(1 + \eta_1) + \eta_2$  avec :

- si  $\hat{x}$  est normal  $|\eta_1| \leq \frac{1}{2} \varepsilon_M$  est l'erreur relative et  $\eta_2 = 0$
- si  $\hat{x}$  est sous-normal  $\eta_1 = 0$  et  $|\eta_2| \leq \frac{1}{2} 2^{E_{\min} - 1 - n}$

On appelle  $u = \frac{1}{2} \varepsilon_M$  l'unité d'arrondi (*unit roundoff*) ou la précision machine.

Par exemple, la précision machine en *binary64* vaut  $2^{-52}/2 \approx 1.1 \times 10^{-16}$ .

Précision  $\neq$  exactitude !

$R((a \cdot b) \cdot c) = R(a \cdot b) \cdot R(c) = a \cdot b(1 + \eta_{ab}) \cdot c(1 + \eta_c) = a \cdot b \cdot c(1 + \eta_{ab})(1 + \eta_c) \approx a \cdot b \cdot c(1 + \eta_{ab} + \eta_c)$   
La précision est toujours  $u$  mais l'exactitude est d'environ  $2u$ ...

**L'exactitude n'est pas assurée par la précision !**

## Remarque

On préfère étudier l'erreur relative car elle ne dépend pas du facteur d'échelle : si  $x$  et  $\hat{x}$  sont multipliés par un facteur  $k$ , l'erreur absolue aussi mais l'erreur relative reste inchangée.

Oui, bon, OK, mais cette majoration est parfois très pessimiste et ne permet pas d'aller très loin.

**5 1 2 Élimination**

L'élimination (*cancellation* in english of speciality) est le phénomène qui apparaît lorsqu'on soustrait deux nombres très proches.

Avec les notations habituelles,  $\widehat{a} = a(1 + \eta_a)$ ,  $\widehat{b} = b(1 + \eta_b)$ ,  $x = a - b$  et  $\widehat{x} == \widehat{a - b\widehat{a}} - \widehat{b}$ . Alors

$$\left| \frac{x - \widehat{x}}{x} \right| = \left| \frac{-a\eta_a + b\eta_b}{a - b} \right| \leq \max(|\eta_a|, |\eta_b|) \frac{|a| + |b|}{|a - b|}$$

On n'obtient ici qu'une majoration, très pessimiste d'ailleurs.

Cependant on constate que l'erreur est d'autant plus grande que  $|a - b| \ll |a| + |b|$  et de plus cette erreur augmente les erreurs déjà présentes dans  $\widehat{a}$  et  $\widehat{b}$ .

Ainsi, l'erreur causée par l'élimination remettent au premier plan de petites erreurs faites précédemment.

Parfois cette erreur n'est pas forcément dramatique. D'abord parce que les données initiales sont peut-être exactes et donc aucune erreur n'est amplifiée. Ensuite, l'erreur causée peut être « écrasée » par une valeur plus grande. Nous étudierons ce phénomène par exemple dans l' [Recherche 1 - 26 page 30](#).

**5 2 Atelier Padawan # 2 : quitter la réalité (un peu difficile...)**

Dans l'atelier précédent, nous avons travaillé avec des outils « réels » : majorations, minorations avec des opérations définies sur  $\mathbb{R}$  et nous avons oublié qu'en fait nous travaillons dans  $\mathbb{V}$  : il est temps d'enlever nos scaphandres pour être plus à l'aise...

*Exceptionnellement dans cette section nous travaillerons en base  $\beta$  quelconque.*

**5 2 1 Trouver la base (algorithme de Malcolm-Gentleman)**

Observez ce qui se passe et expliquez pourquoi cela se passe sachant que sur JavaScript, il n'y a qu'un seul type numérique : *binary64*...

1	Rhino 1.7 release 3	1	js> Math.pow(2, 53) + 3	1	js> Math.pow(2, 54) + 2
2		2	9007199254740996	2	18014398509481984
3	js> Math.pow(2, 53)	3	js> Math.pow(2, 53) + 4	3	js> Math.pow(2, 54) + 3
4	9007199254740992	4	9007199254740996	4	18014398509481988
5	js> Math.pow(2, 53) + 1	5	js> Math.pow(2, 53) + 5	5	js> Math.pow(2, 55)
6	9007199254740992	6	9007199254740996	6	36028797018963970
7	js> Math.pow(2, 53) + 2	7	js> Math.pow(2, 54)	7	js> Math.pow(2, 55) + 4
8	9007199254740994	8	18014398509481984	8	36028797018963970

Et maintenant, voici un algorithme en impératif qui, si on travaillait dans  $\mathbb{R}$  ne terminerait pas mais que vous allez pouvoir comprendre après notre petite observation sur JavaScript :

```

1  Algorithme de Malcolm-Gentleman
2  a ← 1.0, b ← 1.0
3  TantQue R(R(a + 1.0) - a) == 1.0 Faire
4  |   a ← R(2 × a)
5  FinTantQue
6  TantQue R(R(a + b) - a) ≠ b Faire
7  |   b ← b + 1
8  FinTantQue
9  Retourner b

```

```

1  base = boucle2 (boucle1 1) 1
2  where
3  boucle1 a =
4  |   if (a + 1.0) - a /= 1.0 then a
5  |   else boucle1 (2.0 * a)
6  boucle2 a b =
7  |   if (a + b) - a == b then b
8  |   else boucle2 a (b + 1.0)

```

Regardons la première boucle. On y construit une suite  $(a_i)$  vérifiant pour tout naturel  $i$   $a_{i+1} = 2a_i$  avec  $a_0 = 1.0$ . Une récurrence immédiate démontre que  $a_i = 2^i$  pour tout entier naturel  $i$ .  
...Mouais, c'est exact SI l'on travaille dans  $\mathbb{R}$ ...Mais nous travaillons dans  $\mathbb{V}_\beta$  !

*Dans quel espace tu travailles tu me diras et si ta proposition est vraie je saurai*

Mathemator - 44 ap. GC

En fait notre espace de travail est  $\mathbb{V}_\beta$  avec  $\beta$  notre base (qui est donc un entier supérieur ou égal à 2) dont la mantisse (pas la pseudo-mantisse...) est de longueur  $p$ .

Les lignes 4 à 6 construisent une suite  $a_i = R(2a_{i-1})$  avec  $a_0 = 1.0$  mais est-on sûr que l'on a toujours  $a_i = 2^i$  dans ces conditions ?

On démontre d'abord par récurrence que  $a_i = R(2^i) = 2^i$  tant que  $i$  vérifie  $2^i \leq \beta^p - 1$  :  $2^i = 2 \times 2^{i-1}$  est un entier qui s'écrit avec moins de  $p$  chiffres de la base  $\beta$  et en deçà de  $\beta^p - 1$ , on peut passer d'un flottant à son successeur par des pas inférieurs à 1.

Alors  $R(a_i + 1.0) = a_i + 1.0$  : pas d'arrondi.

Puis  $R(R(a_i + 1.0) - a_i) = R(a_i + 1.0 - a_i) = 1.0$  : la condition de la première boucle est donc vérifiée tant que  $2^i < \beta^p$ . Que se passe-t-il après ?

Dès que  $\beta^p - 1$  est atteint, les choses changent.

Considérons la première itération  $i_s$  telle que  $2^{i_s} \geq \beta^p$ .

On a

$$a_{i_s} = R(2 \times a_{i_s-1}) = R(2 \times 2^{i_s-1}) < R(2 \times \beta^p) \leq R(\beta \times \beta^p) = \beta^{p+1}$$

car d'une part  $a_{i_s-1} = 2^{i_s-1}$ , d'autre part  $a_{i_s-1} \leq \beta^p - 1 < \beta^p$  et enfin  $2 \leq \beta$ .

Ainsi

$$\beta^p \leq a_{i_s} < \beta^{p+1}$$

L'exposant de la forme normale de  $a_{i_s}$  est donc  $p$  et d'après le résultat suivant démontré à la section 1.2.4 page 13, le successeur d'un VF  $v$  vérifie (en notant  $n = p - 1$  la longueur de la pseudo-mantisse) :

$$\text{succ}(v) = v + \beta^{E(v)-n} = v + \beta^{E(v)-p+1}$$

Ici  $\text{succ}(a_{i_s}) = a_{i_s} + \beta^{E(a_{i_s})-p+1} = a_{i_s} + \beta^{p-p+1} = a_{i_s} + \beta$ .

On en déduit que  $a_{i_s} + 1.0$  est entre  $a_{i_s}$  et son successeur  $a_{i_s} + \beta$ .

Donc, selon l'arrondi,  $R(a_{i_s} + 1.0)$  vaut  $a_{i_s}$  ou  $a_{i_s} + \beta$  et finalement  $R(R(a_{i_s} + 1.0) - a_{i_s})$  vaut 0 ou  $\beta$  mais en aucun cas 1.0 : on sort donc de la boucle.

*On en déduit que boucle1 1 vaut  $a_{i_s}$  donc que  $\beta^p \leq \text{boucle1 1} < \beta^{p+1}$*

Notons  $a = \text{boucle1 1}$ . Son successeur est  $a + \beta$ .

Passons à la seconde boucle.

Tant que  $b < \beta$ ...mais je vous laisse conclure la démonstration :-)

#### Moralité

Cette démonstration est plutôt technique et difficile en première lecture. Cependant elle est pleine d'enseignements :

- elle nous montre bien que nous changeons de monde : un test du style `if (a + 1.0) - a != 1.0` paraîtrait bien hors de propos si nous raisonnions avec des réels ;
- nous avons vu que l'on pouvait raisonner rigoureusement sur les VF ;
- nous pouvons malgré tout retenir la trame intuitive de la démonstration : tant que notre nombre entier est représentable avec  $p$  chiffres, on reste exact. Les problèmes arrivent lorsqu'on n'a plus assez de place pour stocker tous les chiffres : il y a de la perte d'information...
- Ce genre de raisonnement se retrouve très souvent pour travailler sur les erreurs commises : il faudra donc en retenir la substantifique moelle...

À retenir

N'oubliez pas que certains langages ne travaillent pas en base 2, par exemple Maple© :

```

1 > a := 1.0:
2 > b := 1.0:
3 > while evalf(evalf(a + 1.0) - a) - 1.0 = 0.0 do
4     a := 2.0 * a;
5     end_while:
6 > while evalf(evalf(a + b) - a) - b <> 0.0 do
7     b := b + 1.0
8     end_while:
9 > b;
10
10                                10.0

```

### 5 2 2 Nouvelles opérations

Pour alléger les notations et bien comprendre que nous changeons de monde mais que la structure des calculs sur les flottants reste algébrique (mais des propriétés intéressantes comme l'associativité disparaissent), nous noterons :

#### Notations

Les opérations arithmétiques de base sont bien définies sur  $\mathbb{V}$  donc nous noterons à présent, si l'arrondi n'a pas besoin d'être précisé :

- $R(x + y)$  sous la forme  $x \oplus y$  ;
- $R(x - y)$  sous la forme  $x \ominus y$  ;
- $R(x \times y)$  sous la forme  $x \otimes y$  ;
- $R(x/y)$  sous la forme  $x \oslash y$ .

Ainsi on notera par exemple  $a \oplus b = a + b + \text{err}(a \oplus b)$ .

### 5 2 3 Trouver la précision

Encore un algorithme de M.A. MALCOLM proposé en 1972.

Cette fois, vous vous occuperez de la démonstration tout(e) seul(e)...

```

1 precision = toPrec 1.0 0
2   where
3     toPrec a i =
4       if (a + 1.0) - a /= 1.0 then i
5       else toPrec (base * a) (i + 1)

```

## 5 3 Atelier Padawan # 3 : somme de flottants proches

### 5 3 1 Quelques lemmes utiles



Un lemme en mathématique est un résultat intermédiaire qui sert à prouver un résultat plus important. Un lemme peut être cependant très compliqué à démontrer : la démonstration du *lemme fondamental* a valu à Bao Châu NGÔ d'obtenir en 2010 la médaille FIELDS...

Fini de jouer...Les démonstrations des algorithmes précédents sont intéressantes mais les résultats obtenus sont anecdotiques : on sait bien dans quelle base on travaille et avec quelle précision.

Passons à présent à des résultats beaucoup plus utiles maintenant que nous sommes au parfum des démonstrations sur les flottants.

Un des grands dangers de la manipulation des flottants est d'effectuer des soustractions car on peut perdre beaucoup d'information par annulation (*cancellation*).

Pour avoir un meilleur contrôle de la chose, nous allons étudier un lemme énoncé par Pat H. ? (1927 - 2008)...mais pour le démontrer, nous aurons besoin de quelques sous-lemmes ;-)

#### Convention

Pour éviter de déduire des lemmes suivants des théorèmes totalement faux, n'oubliez pas que **DANS CE QUI SUIT X ET Y SONT DES NOMBRES À VIRGULE FLOTTANTE !**

**Majoration de l'erreur d'une somme**

Posons  $x \oplus y = x + y + \text{err}(x \oplus y)$ . Alors, s'il n'y a pas de dépassement de capacité,

Lemme 1 - 1

$$|\text{err}(x \oplus y)| \leq \min(|x|, |y|)$$

On a bien sûr un résultat analogue pour la différence

On a  $x \oplus y = x + y + \text{err}(x \oplus y)$  (on est dans  $\mathbb{R}...$ )

De plus  $x = x + y + (-y)$  (on est toujours dans  $\mathbb{R}...$ )

$x$  et  $x \oplus y$  sont donc deux VF distants respectivement de  $|y|$  et de  $|\text{err}(x \oplus y)|$  de  $x + y$  (que de de, de deux et de 2 : pffff...)

Mais la norme nous assure que c'est  $x \oplus y$  le VF le plus proche de  $x + y$  donc l'erreur commise en le choisissant est moindre que celle faite en choisissant  $-y$ .

On en déduit que  $|\text{err}(x \oplus y)| \leq |y|$ .

On montre de même que  $|\text{err}(x \oplus y)| \leq |x|$ .

Remarque

De l'importance de disposer avec la IEEE 754 de la **meilleure approximation** !

Corollaire 1 - 1

**Møller, Knuth, Dekker**

L'erreur commise  $|\text{err}(x \oplus y)|$  peut être exprimée exactement sur  $p$  bits.

Ce n'est qu'un corollaire d'un sous-lemme, mais c'est un résultat très important !

Supposons, sans perdre de généralité, que  $|x| \geq |y|$ .

Comme  $x$  et  $y$  sont des VF, le plus petit bit significatif de  $\text{err}(x \oplus y)$  est au moins de magnitude celle de  $\text{ulp}(y)$ .

De plus  $|\text{err}(x \oplus y)| \leq |y|$ , donc la mantisse entière de  $\text{err}(x \oplus y)$  a une longueur inférieure à  $p$  bits.

Finalement le mantisse entière de  $|\text{err}(x \oplus y)| \leq |y|$  est exactement exprimable sur  $p$  bits.

Lemme 1 - 2

Supposons que  $|x + y| \leq \min(|x|, |y|)$ , alors  $x \oplus y = x + y$ .

On obtient un résultat analogue pour la soustraction.

La démonstration va ressembler à la précédente.

Supposons, sans perdre de généralité, que  $|x| \geq |y|$ .

Le plus petit bit significatif de  $x + y$  est au moins de magnitude celle de  $\text{ulp}(y)$ .

Or  $|x + y| \leq |y|$  donc la mantisse entière de  $x + y$  a une longueur inférieure à  $p$  bits.

Finalement le mantisse entière de  $x + y$  est exactement exprimable sur  $p$  bits.

**N'oubliez pas la convention !**

Nous avons démontré nos lemmes en considérant des VF  $x$  et  $y$ .

Est-ce que le résultat suivant contredit notre dernier lemme ?

Danger

```
1 *Main> 1 - 0.9
2 9.999999999999998e-2
```

**5 3 2 Lemme de Sterbenz**

**Lemme de Sterbenz (1973)**

Soit  $(x, y) \in \mathbb{V}^2$  vérifiant  $\frac{x}{2} \leq y \leq 2x$ . On a

Lemme 1 - 3

$$x \ominus y = x - y$$

La différence de deux VF suffisamment proches est donc exacte.

Une petite démonstration?... Allez, le lemme précédent va sûrement être utile pour avoir un petit schéma de démonstration qui tient sur une ligne :-)

Pour une démonstration plus complète, on pourra par exemple se référer à ? page 123.

On peut se restreindre à des nombres positifs.

On a deux cas :

**1.  $x < y$**  : alors  $x < y \leq 2x$  donc  $0 < y - x \leq x \leq y$ ;

**2.  $x \geq y$**  : alors  $\frac{x}{2} \leq y \leq x$  donc  $-\frac{x}{2} \leq y - x \leq 0$  et par suite  $0 \leq x - y \leq \frac{x}{2} \leq y \leq x$ .

Dans tous les cas  $|x - y| \leq \min(|x|, |y|)$  et on peut utiliser le lemme 1 - 2 page précédente.

#### Remarques

Concrètement, à quoi correspond cette condition  $\frac{x}{2} \leq y \leq 2x$  ?

Demandez-vous ce que l'on peut dire de l'écart maximum entre les exposants de  $x$  et  $y$ .

On admettra que le résultat reste vrai même en cas de sous-capacité.

#### 5 4 Atelier Padawan # 4 : somme compensée de flottants quelconques

Les lemmes précédents ne permettent de s'assurer de l'exactitude de la somme que dans des cas restreints.

Que se passe-t-il quand on ne peut pas s'assurer à chaque instant que les opérandes sont suffisamment proches ?

#### 5 4 1 Fast2Sum

Voici une première version officiellement due à T.J. DEKKER en 1971 mais qui avait déjà été introduite par l'incontournable W. KAHAN en 1965.



Fast2Sum - Dekker & Kahan

On considère deux VF  $x$  et  $y$  tels que  $|x| \geq |y|$  et l'algorithme suivant :

```

1  s ← x ⊕ y
2  yv ← s ⊖ x
3  d ← y ⊖ yv
4  Retourner (s, d)

```

#### Théorème 1 - 1

Alors  $x + y = s + d$  avec  $s = x \oplus y$  et  $d = \text{err}(x \oplus y)$ . De plus  $s$  et  $d$  ne se chevauchent pas.

ce qui donne dans notre langage favori :

```

1  fastTwoSum x y
2  | abs x >= abs y =
3  | let s = x + y in
4  | let yv = s - x in
5  | let d = y - yv in
6  | (s, d)
7  | otherwise = fastTwoSum y x

```

La ligne 1 donne  $s = x + y + \text{err}(x \oplus y)$ .

De plus, comme  $|x| \geq |y|$ , c'est  $x$  « qui gagne » donc  $s$  et  $x$  ont le même signe.

La ligne 2 calcule un  $y$  virtuel ( $y_v$ ) qui vaut  $s \ominus x$  : on aimerait bien que ça fasse  $s - x$  c'est-à-dire  $y + \text{err}(x \oplus y)$ .

Enfin la ligne 3 calcule  $y \ominus y_v$  : on aimerait bien aussi que ça fasse  $y - y_v$  c'est-à-dire  $-\text{err}(x \oplus y)$  comme ça on récupérerait exactement l'erreur commise.

En plus, comme  $d = -\text{err}(x \oplus y)$ , on a  $|d| \leq \frac{1}{2} \text{ulp}(s)$  donc  $s$  et  $d$  ne se chevauchent pas dans leur somme.

On aimerait bien, on aimerait bien...En fait tout s'arrange car nous allons pouvoir utiliser le lemme de STERBENZ.

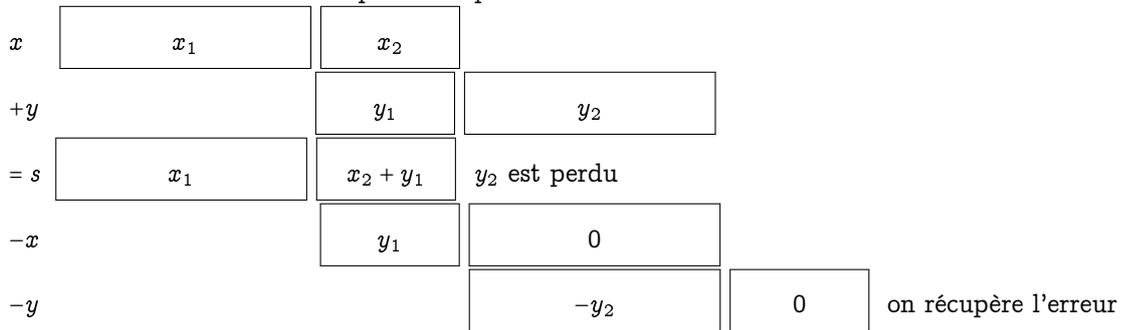
Il suffit de distinguer deux cas :

— si  $x$  et  $y$  sont de même signe OU si  $|y| \leq \frac{|x|}{2}$  : alors  $\frac{x}{2} \leq s \leq 2x$  et on peut appliquer le lemme ;

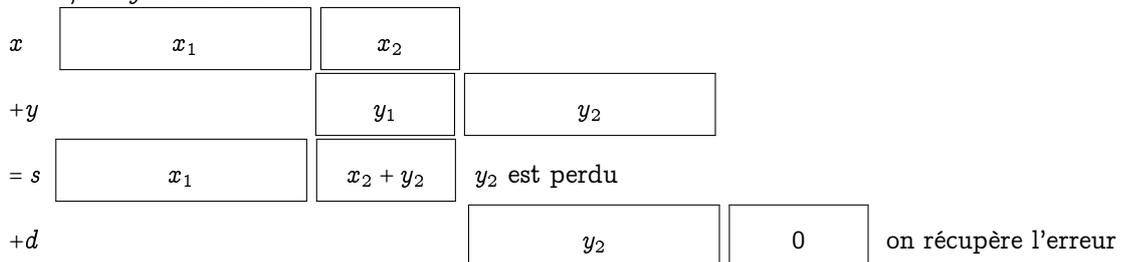
— **sinon** : on a  $x$  et  $y$  de signes opposés ET  $y > \frac{|x|}{2}$  alors si  $y$  est négatif  $x/2 < -y < x$  et sinon  $-x/2 < y < -x$ . Dans les deux cas, d'après le lemme de STERBENZ,  $s$  est calculée exactement et alors  $y_v = y$ .

On illustrera cet algorithme par un exemple simple en TD dans l' [Recherche 1 - 23 page 29](#)

On illustrera encore mieux le problème par cette... illustration :



Ainsi,  $x + y = s + d$  se lit :



**Danger** s et d ne se chevauchent pas :  $s \oplus d = s$  !

**5 4 2 2Sum**

Il s'avère que sur les machines modernes, la comparaison préliminaire du Fast2Sum peut s'avérer plus coûteuse que les trois additions supplémentaires qui vont suivre.

De plus, le Fast2Sum ne fonctionne qu'en base 2 ou 3 alors que l'algorithme suivant fonctionne dans toutes les bases.

Mais il faut encore rentrer dans plus de détails comme l'éventualité d'un parallélisme (lignes 2 et 3), etc...vous verrez ça en Master...

**2Sum - Knuth**

On considère deux VF  $x$  et  $y$  et l'algorithme suivant :

**Théorème 1 - 2**

```

1  s ← x ⊕ y
2  y_v ← s ⊖ x, x_v ← s ⊖ y
3  e_y ← y ⊖ y_v, e_x ← x ⊖ x_v
4  d ← e_x ⊕ e_y
5  Retourner (s, d)
    
```

Alors  $x + y = s + d$  avec  $s = x \oplus y$  et  $d = \text{err}(x \oplus y)$ . De plus  $s$  et  $d$  ne se chevauchent pas.

**5 4 3 Somme compensée d'un nombre quelconque de flottants**

Voici qu'entre à nouveau dans l'arène W. KAHAN avec son algorithme *compensated summation* de 1965 (un peu en avance sur son temps...).

Le « *one liner* » en Haskell :

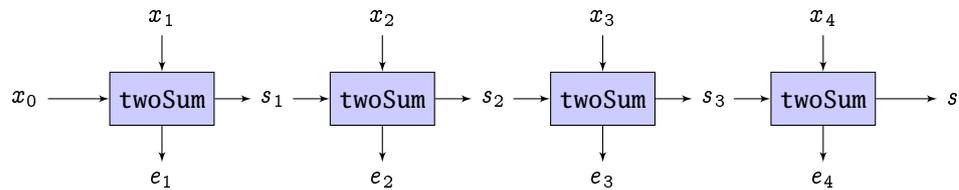
```

1  sommeKahan liste =
2  let (s,e) = foldl (\ (s,c) x -> twoSum s (x + c)) (0,0) liste in s
    
```

Sauriez-vous le traduire en impératif avec une boucle pour ? Il vaudrait mieux pour répondre aux questions de l' [Recherche 1 - 25 page 29](#).

Est-ce optimum ? En 1972, Michèle PICHAT[?] mais aussi Arnold NEUMAIER (mais l'article est en allemand) ont l'idée de sommer à part les erreurs et de soustraire le tout à la fin. Cet algorithme (dit de la *sommation en cascade*) est repris en 2008 par Siegfried M. RUMP, Takeshi OGITA et Shin'Ichi OISHI[?] mais en introduisant le fastSum...faites de même !

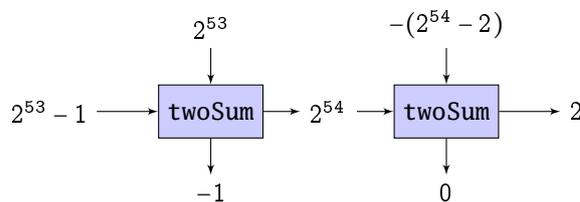
Le schéma suivant vous inspirera sûrement :



Voyons tout de suite un exemple. On veut effectuer la somme en *binary64* de  $x_0 = 2^{53} - 1$ ,  $x_1 = 2^{53}$  et  $x_2 = -(2^{54} - 2)$ .

Que vaut la somme exacte (sans machine car la machine se trompe...)?

Voyons maintenant le schéma du fonctionnement sur machine :



**Recherche**

Expliquez les résultats trouvés. Que va faire l'algorithme de somme compensée ? Et celui de somme en cascade ?

On note  $u = \frac{\epsilon_M}{2}$  et  $\gamma_n = \frac{nu}{1-nu}$ . Alors nos trois amis ? ont démontré le théorème suivant :

**Théorème 1 - 3**

En appliquant l'algorithme de PICHAT-RUMP-OGITA-OISHI à une famille  $(x_i)_{1 \leq i \leq n}$  de flottants et si  $nu < 1$  alors, même en cas de sous-capacité mais sans sur-capacité, le résultat  $s$  retourné par l'algorithme vérifie :

$$\left| s - \sum_{i=1}^n x_i \right| \leq u \left| E \sum_{i=1}^n x_i \right| + \gamma_{n-1}^2 \sum_{i=1}^n |x_i|$$

Ça fait peur, hein ?

**6 Et la multiplication ? Et la division ? Et le sinus ? Et...**

Pour la multiplication, on procède à peu près de la même manière. On peut gagner du temps si le processeur utilise la FMA (Fused Multiple Add) mais ceci est une autre histoire...

Et la division ? On cherche l'inverse  $x$  de  $a$  : on est donc amené à résoudre l'équation  $\frac{1}{x} - a = 0$  d'inconnue  $a$  qui n'est pas linéaire.

Et la racine carrée ? Là aussi, on quitte le linéaire.

Il nous faut donc explorer des outils un peu plus sophistiqués...ce que nous allons faire dans le chapitre suivant.

## EXERCICES

### Recherche 1 - 1

Comment est codé 2 en tant qu'entier 32 bits? En tant que `binary32`? Comment est-il codé en mémoire sous forme hexadécimale?

### Recherche 1 - 2

Donnez au format `binary32` et `binary64` les représentations mémoires hexadécimales des nombres 125,25 puis 0,375 puis -0,375 puis 0,44.

### Recherche 1 - 3

Donnez la représentation mémoire hexadécimale au format `binary32` de  $2^{-130}$ .

### Recherche 1 - 4

Comment traduiriez-vous par une égalité qu'un nombre  $x \in [0, 1[$  s'écrit  $b_1b_2\dots b_n$  dans une base  $\beta$ ?  
Soit l'algorithme :

```

Fonction multSuccessives(x : réel base : entier naturel) : liste d'entiers naturels
nb ← x
dec = []
TantQue nb ≠ 0 Faire
  | bi ← nb × base
  | dec ← dec ++ liste([bi])
  | nb ← {bi}
FinTantQue
Retourner dec
  
```

On notera  $\{x\}$  la partie fractionnaire de  $x$ .

Que pensez-vous de sa précondition? De sa postcondition? de sa terminaison? De sa correction? De sa complexité?

### Recherche 1 - 5

On suppose qu'on dispose d'une fonction `FloatToIntBits` qui renvoie l'écriture machine en base 2 sur 32 bits d'un réel et des fonctions bit à bit habituelles.

Donnez des fonctions renvoyant le signe, l'exposant, la mantisse d'un réel 32 bits.

Pour vous faire plaisir, vous pourrez coder tout ça en Java...Par exemple, voici une fonction qui donne la représentation mémoire hexadécimale d'un flottant 32 bits :

```

1 static String memoire(float reel){
2     int bits;
3     bits = Float.floatToIntBits(reel);
4     return (Integer.toHexString(bits));
5 }
  
```

### Recherche 1 - 6

Que pensez-vous de ça :

```

1 der f h x =
2     (f(x + h) - f(x)) / h
3
4 cosApp = der sin 0.0000001
5
6 sinApp = \x -> - (der cosApp 0.0000001 x)
7
8 cosApp2 = der sinApp 0.00000001
  
```

Qui donne :

```

1 *Main> sinApp (pi/6)
2 0.5051514762044462
3
4
5
6
7 *Main> cosApp2 (pi/2)
8 1110223.0246251565
  
```





iii. Voici les premières bonnes décimales : 12.090 146 129 863 428 0. Estimez les erreurs commises en les exprimant en ulps.

iv. Et  $\sum_{k=1}^{100\,000} \frac{1}{k}$  dans tout ça ?

On pourra utiliser le squelette suivant :

```

1 import qualified Data.List as L -- pour avoir reverse et sort
2
3 -- twoSum : somme compensée de 2 flottants sans branchement
4 twoSum x y =
5   let s      = x + y          in
6   let [xv,yv] = [s - y , s - x ] in
7   let [ex,ey] = [x - xv, y - yv] in
8   (s , ex + ey)
9
10 -- somme basique (+) des éléments d'une liste classée dans l'ordre croissant
11 sommeC liste =
12
13 -- somme basique (+) des éléments d'une liste classée dans l'ordre décroissant
14 sommeD liste =
15
16 -- somme compensée (Kahan) : on rajoute les raisins perdus après chaque cep
17 -- secoué
18 sommeK liste =
19   let (sn,en) = foldl (\ (sk, ek) xk' -> twoSum sk (xk' + ek)) (0,0) liste in sn
20
21 -- somme en cascade (Pichat) : on collecte les raisins perdus à part et on
22 -- les rajoute à la fin
23 sommeP liste =

```

### Recherche 1 - 26

J'ose espérer que vous savez résoudre une équation du type  $ax^2 + bx + c = 0$  avec  $a \neq 0$ ...

Les racines, si elles existent, sont données par une formule bien connue dépendant de  $a$ ,  $b$  et  $c$  :

$$r = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} = \frac{-b \pm \sqrt{\Delta}}{2a} \quad \text{avec} \quad \Delta = b^2 - 4ac$$

Où peuvent se cacher d'éventuelles annulations catastrophiques ? Étudiez ces cas avec attention, voyez si vous pouvez éviter les éliminations catastrophiques en réécrivant les formules un peu dans l'esprit de la « levée d'indétermination ».

1. Que se passe-t-il lorsque  $b^2 \gg |4ac|$  ? En quoi la formule  $\frac{-(b + \text{signe}(b)\sqrt{\Delta})}{2a}$  peut aider ?
2. Que se passe-t-il lorsque  $b^2 \approx 4ac$  ? Peut-on y remédier ? Que peut-on dire de  $\Delta$  par rapport à  $b^2$  ? Y a-t-il élimination catastrophique ?
3. Que se passe-t-il dans le cas de l'équation  $10^{200}x^2 - 3 \times 10^{200}x + 2 \times 10^{200} = 0$  ?
4. Et dans le cas de  $10^{-200}x^2 - 3x + 2 \times 10^{200} = 0$  ?
5. Moralité ?

### Recherche 1 - 27 Bits de garde, d'arrondi, collant

Soustrayez  $1,101 \times 2^{-2}$  puis  $1,111 \times 2^{-2}$  à  $1 \times 2^0$  en *toy7* en gardant tous les bits, même ceux qui dépassent le format de la mantisse puis refaites le calcul en faisant disparaître les bits qui dépassent de la mantisse.

Pourquoi parle-t-on de *bit de garde* ?

En fait, il y a toujours trois bits à droite de la mantisse : G le bit de garde (*guard*), R le bit d'arrondi (*rounding*) et S le bit collant (*sticky*).

Les deux premiers sont comme des paniers qui recueillent les bits translétés (*shifted*). Le dernier est un panier absorbant qui est un ou inclusif des bits translétés au-delà de R.

### Recherche 1 - 28 Calcul de exp 1 et puissances de 10

On veut calculer une approximation de  $e$  ( $= \exp(1)$ ) comme limite au voisinage de l'infini de  $(1 + \frac{1}{n})^n$ .

1. Rappel :  $f$  est dérivable au voisinage de  $x$  si, et seulement si,  $\lim_{t \rightarrow x} \frac{f(t) - f(x)}{t - x} \in \mathbb{R}$ . Dans ce cas on note  $f'(x)$  cette valeur.

2. Utiliser ce rappel pour calculer  $\lim_{t \rightarrow 0} \frac{\ln(1+t)}{t}$ .
3. D duisez-en  $\lim_{n \rightarrow +\infty} n \ln\left(1 + \frac{1}{n}\right)$  puis  $\lim_{n \rightarrow +\infty} \left(1 + \frac{1}{n}\right)^n$ .
4. On note  $\exp n = (1 + 1/n)^n$ . Comparez et commentez :

---

```
1 *Main> [abs (exp 1 - (expn (10.0**k))) | k <- [0..20]]
2 *Main> [abs (exp 1 - (expn (2.0**(3*k)))) | k <- [0..20]]
```

---

### Recherche 1 - 29 0,1 en base 2

Nous d montrerons dans le chapitre suivant que, pour tout r el  $q$  v rifiant  $|q| < 1$  on a

$$\sum_{k=k_0}^{+\infty} q^k = \frac{q^{k_0}}{1-q}$$

1. D montrez que  $\sum_{i=1}^{+\infty} (2^{-4i} + 2^{-4i-1}) = \frac{1}{10}$ .
2. D duisez-en le d veloppement de  $1/10$  en base 2.
3. Montrez que, en *binary32*,  $\frac{\widehat{x}-x}{x} = -\frac{\varepsilon_M}{8}$ .

### Recherche 1 - 30 Max IEEE

On peut calculer le maximum de deux r els ainsi :

Fonction  $\max(x, y : \text{flottants}) : \text{flottant}$

Si  $x > y$  Alors

    |  $\max = x$

Sinon

    |  $\max = y$

FinSi

Est-ce adapt e   la recherche du maximum de deux flottants en arithm tique IEEE ?

### Recherche 1 - 31 Options GCC et double arrondi

Voici un petit code C pour effectuer une multiplication [?] :

---

```
1 #include <stdio.h>
2
3 int
4 main(void)
5 {
6     double a = 1848874847.0;
7     double b = 19954562207.0;
8     double c; c = a * b;
9     printf("c = %20.19e\n", c);
10    return 0;
11 }
```

---

et voici les r sultats obtenus en utilisant deux options de compilation diff rentes :

— Avec `mfpmath=387` :

---

```
1 $ gcc -mfpmath=387 multi.c -o multi_387
2 $ ./multi_387
3 c = 3.6893488147419103232e+19
```

---







**Recherche 1 - 36 4 points**

On considère la fonction suivante :

Haskell

```
p = toP 1.0 0
  where
    toP a cpt =
      if (a + 1.0) - a /= 1.0 then cpt
      else toP (2 * a) (cpt + 1)
```

Imaginons que Haskell travaille uniquement en *toy7* (1.0 est donc un VF codé sur 7 bits). Que vaut alors p ? Justifiez votre réponse.

**Recherche 1 - 37**

Exemple dû à Jean-Michel Muller (in ?)

M. X a récemment été à sa banque (Chaotic Bank Society), pour connaître les nouvelles offres proposées aux meilleurs clients. Son banquier lui propose l'offre suivante : « vous déposez tout d'abord e - 1 euros sur votre compte (où e = 2.7182818... est la base du logarithme népérien). La première année, nous prenons 1 euro sur votre compte de frais de gestion. Par contre, la deuxième année est plus intéressante pour vous, car nous multiplions votre capital restant par 2 et prenons 1 euro de frais de gestion. La troisième année est encore plus intéressante, car nous multiplions votre capital par 3 et prenons 1 euro de frais de gestion. Et ainsi de suite : la n-ième année, nous multiplions votre capital par n et prenons 1 euro de frais de gestion. Intéressant, non ? » Pour prendre sa décision, M. X décide de demander l'aide d'un informaticien et se retourne vers vous.

**Recherche 1 - 38 Bac**

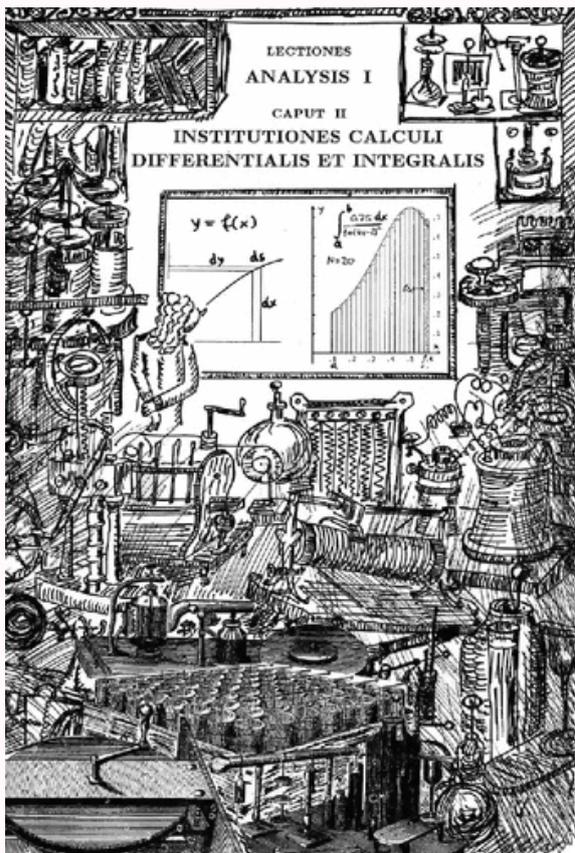
Le sujet du Bac S des Centres Étrangers 2012 fait étudier la suite :

$$I_{n+2} = \frac{1}{2}e - \frac{n+1}{2}I_n, \quad I_1 = \frac{1}{2}(e-1)$$

Calculez  $I_{59}$  à la machine...



# Un soupçon de calcul différentiel



La notion d'*infinitement petit* a créé des guerres entre les savants pendant deux siècles. Elle n'a été formalisée rigoureusement qu'à la fin du XIX<sup>e</sup> siècle. Cependant, nous allons rester fidèle au cheminement de l'esprit humain et nous attacher à cette notion qui est pédagogiquement plus intéressante et se rapproche de l'arithmétique des flottants sur machine. Cette approche historique et intuitive du calcul a été magistralement mise en œuvre dans ?.

# 1 Le calcul infinitésimal

## 1.1 Infiniment petits

La notion de dérivée est née bien avant celle de limite : dès le XVII<sup>e</sup> FERMAT, NEWTON, LEIBNIZ, Jean BERNOULLI, vont développer le calcul infinitésimal, avant que BOLZANO et WEIERSTRASS ne perfectionnent la notion de limite deux siècles plus tard.

Mais les approches des savants du XVII<sup>e</sup> s'adaptent bien à l'utilisation des flottants sur machine. L'informatique nous incite donc à un retour aux sources historiques.

### Analyse non standard et hyperréels

Il est à noter que dans les années 1960, un logicien américain d'origine allemande (encore un !...), Abraham ROBINSON, mit au point l'*analyse non standard* en utilisant les *nombre hyperréels*. Il s'agit de nos réels auxquels on adjoint des nombres *infinitésimaux* qui sont strictement inférieurs en valeur absolue à tout nombre réel non nul, et les nombres *infinitement grands* dont les inverses sont infinitésimaux. Ceci est fait de manière très rigoureuse et permet de justifier une incroyable intuition que FERMAT avait présentée trois siècles plus tôt.

C'est également assez proche de la définition des flottants selon la norme IEEE 754

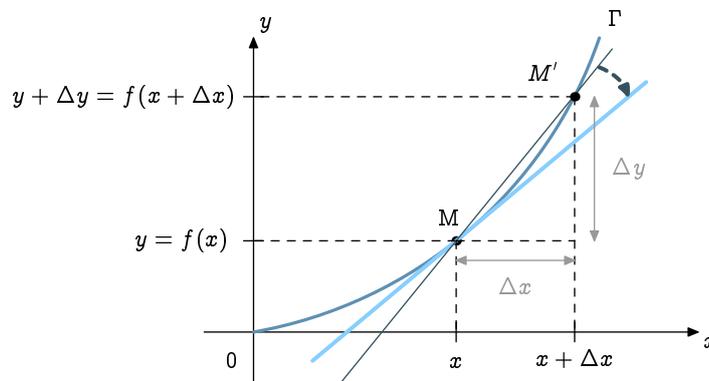


### Remarque

Pour Jean BERNOULLI (1691), les *infiniment petits* sont des quantités qui peuvent être ajoutées à des quantités finies sans changer leurs valeurs. Les courbes sont des polygones à côtés infiniment courts.

Voici qui s'adapte à la notion de sous-capacité et au tracé de courbes sur machine.

Voyons ce que cela donne pour les tangentes à la courbe d'équation  $y = x^2$  :



Si  $x$  augmente de  $\Delta x$ , alors  $y$  devient  $y + \Delta y = (x + \Delta x)^2 = x^2 + 2x\Delta x + (\Delta x)^2$ . Or  $y = x^2$ .

Pour des valeurs de  $\Delta x$  non nulles, cela donne donc :

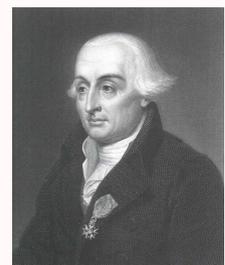
$$\frac{\Delta y}{\Delta x} = 2x + \Delta x$$

C'est ici qu'on a besoin de la notion intuitive d'infiniment petits...qui commence pour nous avec la sous-capacité.

C'est Joseph-Louis LAGRANGE (1736-1813) ou plutôt Giuseppe Lodovico LAGRANGIA, un mathématicien Piémontais d'origine française, qui introduisit les termes et les notations que nous utilisons actuellement :

Nous appellerons la fonction  $f_x$ , *fonction primitive*, par rapport aux fonctions  $f'_x$ ,  $f''_x$ , &c. qui en dérivent, et nous appellerons celles-ci, *fonctions dérivées*, par rapport à celle-là.

Il a aussi réfuté la notion d'infiniment petit et a tenté de fonder l'analyse sur les séries de TAYLOR : nous allons en reparler.



### Remarque

**1 2 Règles de dérivation**

**1 2 1 Dérivée d'une combinaison linéaire**

Soit  $y(x) = a \cdot u(x) + b \cdot v(x)$  avec  $a$  et  $b$  des constantes réelles. En posant comme précédemment  $y + \Delta(y) = y(x + \Delta x)$  puis  $u + \Delta(u) = u(x + \Delta x)$  et  $v + \Delta(v) = v(x + \Delta x)$ , nous obtenons  $\Delta y = a \cdot \Delta u + b \cdot \Delta v$ . Alors :

Dérivation d'une combinaison linéaire

**Théorème 2 - 1**

$$y = au + bv \implies \frac{dy}{dx} = a \cdot \frac{du}{dx} + b \cdot \frac{dv}{dx} \text{ ou bien } y' = a \cdot u' + b \cdot v'$$

**1 3 Dérivée d'un produit**

**Recherche**

Essayez de prouver de même une règle de dérivation d'un produit  $y = u \cdot v$ .

Dérivation d'un produit

**Théorème 2 - 2**

$$y = u \cdot v \implies$$

**1 3 1 Dérivée d'un quotient**

Nous allons d'abord établir un résultat préliminaire. Pour  $x$  « petit » on « sent » que  $\frac{1}{1-x} \approx 1$ . Par exemple  $\frac{1}{1-10^{-5}} \approx 1,0000100001 = 1 + 10^{-5} + 10^{-10}$ . Essayons d'être plus précis sur l'approximation. Posons

$$\frac{1}{1-x} = 1 + \delta$$

et essayons de déterminer ce  $\delta$  sachant qu'on peut le négliger devant  $x$ . Multiplions chaque membre de l'équation par  $1-x$ . On obtient :

$$1 = 1 - x + \delta(1 - x)$$

c'est-à-dire, après simplifications :

$$\delta = x + \delta x$$

Ainsi, comme  $\delta x$  est négligeable devant  $x$ , on peut considérer que  $\delta \approx x$  donc :

$$\frac{1}{1-x} \approx 1 + x$$

Série géométrique

En 1593, François VIÈTE établit avec une méthode similaire que, pour tout réel  $x$  tel que  $|x| < 1$ , on a :

**Remarque**

$$\frac{1}{1-x} = 1 + x + x^2 + x^3 + x^4 + x^5 + \dots$$



**Recherche**

Vous pouvez maintenant essayer d'établir une formule donnant la dérivée d'un quotient  $y = \frac{u}{v}$ .

Dérivation d'un quotient

**Théorème 2 - 3**

$$y = \frac{u}{v} \implies$$

**1 3 2 Fonctions composées**

Soit trois fonctions  $f$ ,  $g$  et  $h$  telles que  $h = f \circ g$ . Ainsi,  $h(x) = f(g(x))$ .

Posons  $g(x + \Delta x) = z + \Delta z$  et  $h(x + \Delta x) = y + \Delta y$ .

Alors  $y + \Delta y = f(g(x + \Delta x)) = f(z + \Delta z)$ . Or  $\frac{\Delta y}{\Delta x} = \frac{\Delta y}{\Delta z} \cdot \frac{\Delta z}{\Delta x}$ . On obtient donc :

Dérivation de fonctions composées

Théorème 2 - 4

$$\frac{dy}{dx} = \frac{dy}{dz} \cdot \frac{dz}{dx} \text{ ou } h'(x) = f'(g(x)) \cdot g'(x)$$

Par exemple, si  $h(x) = (3x - 5)^2$ , alors  $y = h(x) = f(g(x))$  avec  $z = g(x) = 3x - 5$  et  $f(z) = z^2$ . On en déduit que  $h'(x) = \frac{dy}{dz} \cdot \frac{dz}{dx} = 2z \cdot 3 = 6(3x - 5)$ .

**2****Comment calculer un logarithme sur machine..au XVII<sup>e</sup> siècle ?****2 1 Calcul d'une approximation d'une racine carrée par la méthode de Héron**

Le mathématicien HÉRON d'Alexandrie n'avait pas attendu NEWTON et le calcul différentiel pour trouver une méthode permettant de déterminer une approximation de la racine carrée d'un nombre entier positif puisqu'il a vécu seize siècles avant Sir Isaac.

Si  $x_n$  est une approximation strictement positive par défaut de  $\sqrt{a}$ , alors  $a/x_n$  est une approximation par excès de  $\sqrt{a}$  (pourquoi ?) et vice-versa.

La moyenne arithmétique de ces deux approximations est  $\frac{1}{2} \left( x_n + \frac{a}{x_n} \right)$  et constitue une meilleure approximation que les deux précédentes.

On peut montrer c'est une approximation par excès (en développant  $(x_n - \sqrt{a})^2$  par exemple).

```
1 heron a x n
2 | n == 0 = x
3 | otherwise = heron a ((x + a / x) / 2) (n - 1)
```

Le résultat est assez bluffant :

```
1 *Main> heron 2 1 5
2 1.414213562373095
3 *Main> sqrt 2
4 1.4142135623730951
```

Il y a cependant des discussions au sujet de l'exactitude des calculs dans  $\mathbb{V}$  : nous en reparlerons dans un prochain chapitre quand nous saurons diviser deux flottants :-)

**2 2 Calcul d'une approximation de log 2 par la méthode de Briggs**

Michael STIFEL s'est intéressé dès 1544 aux fonctions vérifiant  $\ell(x \cdot y) = \ell(x) + \ell(y)$  afin de simplifier les calculs des astronomes et des navigateurs. Cependant, cela nécessite de disposer de tables *logarithmiques* ( $\lambda\omicron\gamma\omicron\varsigma$  : mot, relation,  $\alpha\rho\iota\theta\mu\omicron\varsigma$  : nombre ; les logarithmes sont donc des relations entre les nombres avant d'être l'anagramme d'algorithmes...).

Nous allons nous occuper comme Henry BRIGGS du logarithme de base 10.

Sachant que  $\log_{10}(10) = 1$  et que  $\log_{10}(a^p) = p \cdot \log_{10}(a)$  pour tout rationnel  $p$ , on en déduit alors que  $\log_{10}(\sqrt{10}) = \frac{1}{2}$ ,  $\log_{10}(\sqrt{\sqrt{10}}) = \frac{1}{4}$ , etc.

Nous pouvons utiliser l'algorithme de HORNER pour calculer les approximations des racines successives de 10 et on obtient le tableau suivant :

Nombres	10,0000	3,1623	1,7783	1,3335	1,0000
Logarithmes	1	0,5	0,25	0,125	0



Le problème, c'est que connaître le logarithme de 3,1623 n'est pas très intéressant. On aimerait plutôt connaître ceux des entiers de 1 à 10.

C'est ici qu'intervient l'anglais Henry BRIGGS qui publia en 1617 ses premières tables de logarithmes décimaux contenant 1000 valeurs avec quatorze décimales.

Voyons comment il a procédé pour le calcul de  $\log(2)$ .

Il a commencé par obtenir des valeurs approchées des nombres suivant :

$$\begin{array}{ll} \sqrt{10} = 10^{\frac{1}{2}} & \sqrt{2} = 2^{\frac{1}{2}} \\ \sqrt{\sqrt{10}} = 10^{\frac{1}{2^2}} & \sqrt{\sqrt{2}} = 2^{\frac{1}{2^2}} \\ \sqrt{\sqrt{\sqrt{10}}} = 10^{\frac{1}{2^3}} & \sqrt{\sqrt{\sqrt{2}}} = 2^{\frac{1}{2^3}} \\ \vdots & \vdots \\ \sqrt{\sqrt{\dots\sqrt{10}}} = 10^{\frac{1}{2^{54}}} & \sqrt{\sqrt{\dots\sqrt{2}}} = 2^{\frac{1}{2^{54}}} \end{array}$$

Il a finalement abouti à :

$$10^{\frac{1}{2^{54}}} \approx \underbrace{1.000\ 000\ 000\ 000\ 000\ 127\ 819\ 149\ 320\ 032\ 35}_{1+a}$$

$$2^{\frac{1}{2^{54}}} \approx \underbrace{1.000\ 000\ 000\ 000\ 000\ 038\ 477\ 397\ 965\ 583\ 10}_{1+b}$$

Or on cherche à déterminer  $x = \log(2)$  mais  $x = \log(2) \Leftrightarrow 10^x = 2$ .

D'après les calculs de BRIGGS :

$$1 + b = 2^{\frac{1}{2^{54}}} = (10^x)^{\frac{1}{2^{54}}} = (1 + a)^x$$

Un résultat classique que nous allons très bientôt étudier montre que  $(1 + a)^x \approx 1 + ax$  pour  $a$  assez proche de zéro et  $x > 0$ . On en déduit que  $1 + b \approx 1 + ax$  et donc que :

$$\log(2) = x \approx \frac{b}{a} = \frac{3\ 847\ 739\ 796\ 558\ 310}{12\ 781\ 914\ 932\ 003\ 235}$$

```
1 *Main> 3847739796558310 / 12781914932003235
2 0.30102999566398114
```

Ainsi,  $\log(2) \approx 0.30102999566398114$ .

Cela nous permet alors d'avoir également  $\log(4) = 2 \times \log(2)$  et  $\log(8) = 3 \times \log(2)$  mais aussi  $\log(5) = \log(10) - \log(2)$ .

Il nous faudrait donc obtenir de même  $\log(3)$  et  $\log(7)$  pour compléter notre collection car  $\log(6) = \log(3) + \log(2)$  et  $\log(9) = 2 \cdot \log(3)$ .

On peut donc à titre d'exercice calculer ces deux logarithmes pour obtenir :

$$\log(3) \approx 0.47712125471966244 \quad \log(7) \approx 0.845098040142567$$

La route est longue jusqu'aux 1000 logarithmes de BRIGGS[?] :

Num. absolu.	Logarithmi.	Num. absolu.	Logarithmi.	Num. absolu.	Logarithmi.
1	0.00000.00000.0000 30102.99956.6398	34	1.53147.89170.4226 1258.91273.0802	67	1.82607.48027.0083 643.41100.0541
2	0.30102.99956.6398 17609.12590.5568	35	1.54406.80443.5028 1223.44564.1701	68	1.83250.89127.0624 634.01780.3102
3	0.47712.12547.1966 12493.87366.0830	36	1.55630.25007.6729 1189.92232.9970	69	1.83884.90907.3726 624.89492.7700
4	0.60205.99913.2796 9691.00130.0806	37	1.56820.17240.6699 1158.18725.4982	70	1.84509.80400.1426 616.03087.0482
5	0.69897.00043.3602 7918.12460.4762	38	1.57978.35966.1681 1128.10104.0969	71	1.85125.83487.1908 607.41477.1219
6	0.77815.12503.8364 6694.67896.3062	39	1.59106.46070.2650 1099.53843.0146	72	1.85733.24964.3127 599.03636.8919
7	0.84509.80400.1426 5799.19469.7768	40	1.60205.99913.2796 1072.38653.9178	73	1.86332.28601.2046 590.88596.1052
8	0.90308.99869.9194 5115.25224.4738	41	1.61278.38567.1974 1046.54336.7816	74	1.86923.17197.3098 582.95436.6072
9	0.95424.25094.3932 4575.74905.6068	42	1.62324.92903.9790 1021.91651.8169	75	1.87506.12633.9170 575.23288.8909

On peut retrouver ces résultats sur Sage qui permet de travailler très simplement avec MPFR :

---

```

1 # Travail commode avec la bibliothèque MPFR de C via sage
2 # On ouvre dans emacs un fichier d'extension sage qui met emacs en mode sage
3 # La syntaxe de programmation est celle de Python
4 #
5
6 # On travaille sur 128 bits (on donne la précision = taille de la mantisse)
7 RN = RealField(113)
8
9 # Algorithme de Héron pour calculer ?..
10 def racine(a,x,n):
11     if n == 0:
12         return RN(x)
13     else:
14         return racine(RN(a), ??? * RN(0.5), n - 1)
15
16 # quid ?
17 def racine_n(a,n):
18     if n == 0:
19         return RN(a)
20     else:
21         return racine_n(racine(RN(a),1.0,10), n - 1)
22
23 # Calcul de log 2
24 a = ???
25 b = ???
26
27 logDeux = RealField(53)(b / a) # résultat projeté sur 64 bits
28 diffLog = logDeux - (log(2.)/ log(10.)) # On compare avec le log 2 de sage
29
30 # Calcul de log 3 et log 7 ? Généralisation ? log 0 ? log (-10) ?

```

---

Alors :

---

```

1 sage: logDeux
2 0.301029995663981182
3 sage: diffLog
4 5.55111512312578e-17

```

---

## 2 3 Polynômes interpolateurs

### 2 3 1 Qu'est-ce qu'un polynôme ?

Les polynômes en informatique sont un outil très puissant et très utile, notamment en théorie de l'information (détection et correction d'erreur). Les entiers en précision infinie sont considérés comme des polynômes. Les exemples abondent mais il s'agit d'une vision discrète de la notion : un polynôme est une liste finie de nombres et on munit l'ensemble des polynômes d'opérations spécifiques.

Nous allons ici étudier plutôt les fonctions polynomiales, i.e. les fonctions de la forme :

$$p : x \mapsto p(x) = a_n x^n + a_{n-1} x^{n-1} + \dots + a_0$$

où les  $a_i$  sont des constantes arbitraires qui seront pour nous des nombres à virgule flottante. Si  $a_n$  est non nul, on dira que la fonction polynôme est de degré  $n$ .

Cette fonction correspond au polynôme :

$$P = [a_n, a_{n-1}, \dots, a_0]$$

Dans  $\mathbb{F}_2$  nous verrons que la distinction entre polynôme et fonction polynomiale est très importante. Dans  $\mathbb{R}$  ou  $\mathbb{V}$ , c'est moins grave mais bon...

L'intérêt ? Il s'agit de fonctions très rapidement évaluables sur machine car faisant intervenir uniquement des additions et des multiplications.

**2 3 2 Problème d'interpolation**

On connaît quelques valeurs du logarithme donnée par la méthode de BRIGGS. Cependant, cette méthode est très fastidieuse : il faudrait trouver une autre idée pour calculer les valeurs intermédiaires.

Par exemple, on connaît pour l'instant des approximations des logarithmes de 1, 2, 3 et 4. Comment faire pour déterminer des approximations de tous les nombres entre 1 et 4 avec un pas de 0,25 ?

Nous allons chercher un polynôme de degré 3 qui passe par ces quatre points. On cherche donc  $a, b, c$  et  $d$  tels que :

$$\begin{cases} p(x) = a + bx + cx^2 + dx^3 \\ \text{map } p [1, 2, 3, 4] = \text{map } \log [1, 2, 3, 4] \end{cases}$$

C'est un simple système linéaire de 4 équations à quatre inconnues. Cependant, les équations ont une forme très spéciale. Voyons une première éthode proposée par NEWTON en 1676.

**2 3 3 Interpolation de Newton**

Notons  $[y_1, y_2, y_3, y_4] = \text{map } \log [1, 2, 3, 4]$ .

**P R O P. IV.**

*Si recta aliqua in partes quocunq; inæquales AA2, A2A3, A3A4, A4A5, &c. dividatur, & ad puncta divisionum erigantur parallelæ AB, A2B2, A3B3, &c. Invenire Curvam Geometricam generis Parabolici quæ per omnium erectarum terminos B, B2, B3, &c. transibit.*

*Sunto puncta data B, B2, B3, B4, B5, B6, B7, &c. et ad Abscissam quamvis AA7 demitte Ordinatæ perpendiculariter BA, B2A2, B3A3, &c.*

Et fac  $\frac{AB - A2B2}{AA2} = b, \frac{A2B2 - A3B3}{A2A3} = b2,$   
 $\frac{A3B3 - A4B4}{A3A4} = b3, \frac{A4B4 - A5B5}{A4A5} = b4,$   
 $\frac{A5B5 - A6B6}{A5A6} = b5, \frac{A6B6 - A7B7}{A6A7} = b6,$   
 $\frac{A7B7 - A8B8}{A7A8} = b7.$   
 Deinde  $\frac{b - b2}{AA3} = c, \frac{b2 - b3}{A2A4} = c2, \frac{b3 - b4}{A3A5} = c3, \&c.$   
 Tunc  $\frac{c - c2}{AA4} = d, \frac{c2 - c3}{A2A5} = d2, \frac{c3 - c4}{A3A6} = d3, \&c.$   
 Et  $\frac{d - d2}{AA5} = e, \frac{d2 - d3}{A2A6} = e2, \frac{d3 - d4}{A3A7} = e3, \&c.$   
 Sic pergendum est ad ultimam differentiam.

B b

Problème d'interpolation dans *Methodus Differentialis* publié par NEWTON en 1676

Appliquons l'algorithme de Gauß à ce système même si le génie allemand ne naîtra qu'un siècle plus tard :

$x = 1$	$p(x) = a + b + c + d = y_1$
$x = 2$	$p(x) = a + 2b + 4c + 8d = y_2$
$x = 3$	$p(x) = a + 3b + 9c + 27d = y_3$
$x = 4$	$p(x) = a + 4b + 16c + 64d = y_4$

En soustrayant les lignes deux par deux, on fait disparaître  $a$  :

$$\begin{cases} y_2 - y_1 = b + 3c + 7d = \Delta y_1 \\ y_3 - y_2 = b + 5c + 19d = \Delta y_2 \\ y_4 - y_3 = b + 7c + 37d = \Delta y_3 \end{cases}$$

En soustrayant les lignes deux par deux, on fait disparaître  $b$  :

$$\begin{cases} \Delta y_2 - \Delta y_1 = 2c + 12d = \Delta^2 y_1 \\ \Delta y_3 - \Delta y_2 = 2c + 18d = \Delta^2 y_2 \end{cases}$$

En soustrayant les lignes deux par deux, on fait disparaître  $c$  :

$$\Delta^2 y_2 - \Delta^2 y_1 = 6d = \Delta^3 y_1$$

ce qui donne :

$$d = \frac{1}{6} \Delta^3 y_1$$

$$c = \frac{1}{2} \Delta^2 y_1 - \Delta^3 y_1$$

$$b = \Delta y_1 - \frac{3}{2} \Delta^2 y_1 + 3 \Delta^3 y_1 - \frac{7}{6} \Delta^3 y_1 = \Delta y_1 - \frac{3}{2} \Delta^2 y_1 + \frac{11}{6} \Delta^3 y_1$$

$$a = y_1 - \Delta y_1 + \frac{3}{2} \Delta^2 y_1 - \frac{11}{6} \Delta^3 y_1 - \frac{1}{2} \Delta^2 y_1 + \Delta^3 y_1 - \frac{1}{6} \Delta^3 y_1 = y_1 - \Delta y_1 + \Delta^2 y_1 - \Delta^3 y_1$$

Finalement :

$$P(x) = y_1 + (x-1)\Delta y_1 + \frac{1}{2}(x-1)(x-2)\Delta^2 y_1 + \frac{1}{6}(x-1)(x-2)(x-3)\Delta^3 y_1$$

NEWTON avait l'habitude de réunir les résultats dans un schéma de cette forme :

$$\begin{array}{cccc} & & & y_1 \\ & & & \Delta y_1 \\ y_2 & & & \Delta^2 y_1 \\ & \Delta y_2 & & \Delta^3 y_1 \\ y_3 & & \Delta^2 y_2 & \\ & \Delta y_3 & & \\ y_4 & & & \end{array}$$

avec  $\Delta y_i = y_{i+1} - y_i$ ,  $\Delta^2 y_i = \Delta y_{i+1} - \Delta y_i$  et  $\Delta^3 y_i = \Delta^2 y_{i+1} - \Delta^2 y_i$ .

On peut généraliser à des abscisses successives de 1 à  $n$  :

Interpolation de Newton passant par les points  $(1, y_1), (2, y_2), \dots, (n, y_n)$

Théorème 2 - 5

$$p(x) = y_1 + (x-1)\Delta y_1 + \frac{1}{2!}(x-1)(x-2)\Delta^2 y_1 + \dots + \frac{1}{(n-1)!}(x-1)\dots(x-(n-1))\Delta^{n-1} y_1$$

voire à d'autres valeurs d'entiers successifs :

Interpolation de Newton passant par les points  $(\alpha, y_1), (\alpha+1, y_2), \dots, (\alpha+n-1, y_n)$

Théorème 2 - 6

$$p(x) = y_1 + (x-\alpha)\Delta y_1 + \frac{1}{2!}(x-\alpha)(x-\alpha-1)\Delta^2 y_1 + \dots + \frac{1}{(n-1)!}(x-\alpha)\dots(x-\alpha-(n-2))\Delta^{n-1} y_1$$

On peut généraliser à des suites quelconques d'abscisses mais c'est un peu plus compliqué.

### 3 Polynômes et erreurs

Nous allons étudier à l' [Recherche 2 - 8 page 53](#) un moyen d'évaluer une fonction polynomiale en un nombre particulier.

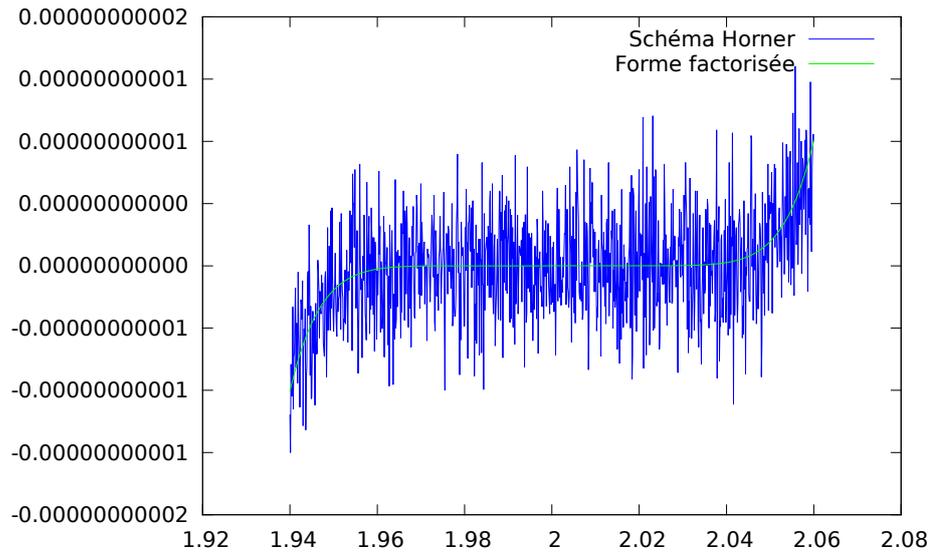
Cependant, nous allons être à nouveau confrontés au problème des erreurs suite au calcul sur les nombres à virgule flottante.

Considérons par exemple la fonction  $x \mapsto (x-2)^9$  au voisinage de 2.

Une petite illustration vaut mieux qu'un long discours...Avec le code suivant :

```
1 comparePoly =
2   plot X11 [
3     Data2D [Title "Schéma Horner", Color Blue, Style Lines] [] [(x, exHorn x) | x <-
4       → [1.94, 1.940125 .. 2.06]],
5     Function2D [Title "Forme factorisée", Color Green] [Range 1.94 2.06, Step 0.000125]
6       → (\x -> exFac x)
7   ]
```

On obtient :



Bon, évidemment, on est allé chercher l'exemple qui fâche. Ailleurs, les choses marchent bien mieux.

Il existe un algorithme de Horner compensé comme il en existe un pour la somme mais c'est un peu plus compliqué. Les curieux pourront se référer à un article de Philippe LANGLOIS et Nicolas LOUVET [?].

À retenir

Comme nous ne sommes qu'en premier cycle, nous mettrons de côté à partir de maintenant ces problèmes d'erreur dans les évaluations.

Il faudra tout de même retenir que ces problèmes existent et sont omniprésents dans la machine : est-ce que mon joueur est sous le feu de l'ennemi ? Je dois savoir si, en 3D, il est à gauche ou à droite d'une zone ce qui revient sur machine à calculer un déterminant dépendant de 3 paramètres, bref évaluer un polynôme de trois variables. Si je me trompe dans mes calculs, quand je suis prêt de la zone frontière, mon joueur peut mourir dans d'atroces souffrances (sur mon écran seulement, bien sûr).



## 4 Série de Taylor

### 4.1 Dérivées d'ordre supérieur

Nous avons eu une approche intuitive de la dérivation jusqu'à maintenant. Elle nous suffira la plupart du temps. Cependant il ne faut pas oublier qu'un manque de rigueur dans l'étude de ces notions peut entraîner de graves erreurs.

Nous rappelons donc quelques définitions qui nous serviront à mieux comprendre certains problèmes.

Définition 2 - 1

#### Dérivabilité

Soit  $f$  une fonction numérique définie sur un intervalle  $I$  de  $\mathbb{R}$ .

Soit  $a \in I$ . On dit que  $f$  est dérivable en  $a$  si, et seulement si, le rapport :

$$\tau_a(h) = \frac{f(a+h) - f(a)}{h}$$

admet une limite finie lorsque  $a$  tend vers 0.

Dans ce cas on note  $f'(a)$  cette limite.

Rappelons que  $f'(a)$  est aussi le coefficient directeur de la tangente à la courbe représentative de  $f$  au point d'abscisse  $a$ .

#### Dérivée $k^{\text{ème}}$

Définition 2 - 2

Si  $f'$  est dérivable sur  $I$  on note  $f''$  ou  $f^{(2)}$  la dérivée de  $f'$  et on l'appelle la dérivée seconde de  $f$ . De même si  $f''$  est dérivable sur  $I$ ,  $f'''$  ou  $f^{(3)}$ , la dérivée de  $f''$ , est la dérivée troisième de  $f$ . Plus généralement on note  $f^{(k)}$  la dérivée  $k^{\text{ème}}$  de  $f$  avec la convention  $f^{(0)} = f$ ,  $f^{(k+1)} = (f^{(k)})'$ .

#### Fonctions de classe $C^n$ , de classe $C^\infty$

Définition 2 - 3

On dit que  $f$  est de classe  $C^n$ ,  $n \in \mathbb{N}$ , sur  $I$  ssi  $f$  est  $n$  fois dérivable et ses  $n$  dérivées sont continues sur  $I$ . On remarquera que si  $f^{(n+1)}$  existe sur  $I$  alors  $f$  est de classe  $C^n$  sur  $I$  qui s'écrit  $f \in C^n(I)$ . On dit que  $f$  est de classe  $C^\infty$  sur  $I$  pour exprimer que  $f$  est indéfiniment dérivable et, a fortiori, que toutes ses dérivées sont continues.

### 4 2 Relations de domination

*Brooks's Law [prov.]*

« Adding manpower to a late software project makes it later » – a result of the fact that the expected advantage from splitting work among  $N$  programmers is  $O(N)$ , but the complexity and communications cost associated with coordinating and then merging their work is  $O(N^2)$

in « The New Hacker's Dictionary »

[http://outpost9.com/reference/jargon/jargon\\_17.html#SEC24](http://outpost9.com/reference/jargon/jargon_17.html#SEC24)

Les notations de LANDAU(1877-1938) ont en fait été créées par Paul BACHMANN(1837-1920) en 1894, mais bon, ce sont tous deux des mathématiciens allemands.

Par exemple, si l'on considère l'expression :

$$f(n) = n + 1 + \frac{1}{n} + \frac{75}{n^2} - \frac{765}{n^3} + \frac{\cos(12)}{n^{37}} - \frac{\sqrt{765481}}{n^{412}}$$

Quand  $n$  est « grand », disons 10 000, alors on obtient :

$$f(10\,000) = 10\,000 + 1 + 0,0001 + 0,00000000075 - 0,000000000000765 + \text{peanuts}$$

Tous les termes après  $n$  comptent pour du beurre quand  $n$  est « grand » : on a expérimenté ceci sur machine avec la notion de *sous-capacité*.

Donnons une définition pour plus de clarté :

#### « Grand O »

Définition 2 - 4

Soit  $f$  et  $g$  deux fonctions de  $\mathbb{N}$  dans  $\mathbb{R}$ . On dit que  $f$  est un « grand O » de  $g$  et on note  $f = O(g)$  ou  $f(n) = O(g(n))$  si, et seulement si, il existe une constante strictement positive  $C$  telle que

$$|f(n)| \leq C|g(n)|$$

pour tout  $n \in \mathbb{N}$ .

Dans l'exemple précédent,  $\frac{1}{n} \leq \frac{1}{1} \times 1$  pour tout entier  $n$  supérieur à 1 donc  $\frac{1}{n} = O(1)$ .

De même,  $\frac{75}{n^2} \leq \frac{75}{1^2} \times 1$  donc  $\frac{75}{n^2} = O(1)$  mais on peut dire mieux :  $\frac{75}{n^2} \leq \frac{75}{1} \times \frac{1}{n}$  et ainsi on prouve que  $\frac{75}{n^2} = O\left(\frac{1}{n}\right)$ .

En fait, un grand O de  $g$  est une fonction qui est au maximum majorée par un multiple de  $g$ .

Point de vue algorithmique, cela nous donne un renseignement sur la complexité au pire.

Pour le tri fusion, on obtient pour  $n > 1$  :

$$K(n) \leq 4nc([\log_2(n)] + 1) \leq 4nc \log_2(n) \left(1 + \frac{1}{\log_2(n)}\right) \leq 8cn \log_2(n)$$

On en déduit que  $K(n) = O(n \log_2(n))$  : la complexité est au pire en  $n \log_2(n)$ , c'est ce qui nous importe. Le rôle des constantes est accessoire car chercher trop de précisions serait illusoire : l'« oubli » de ces constantes correspond en fait aux différences entre langages, processeurs, etc.

On peut cependant faire mieux avec le tri fusion car on a aussi une minoration. C'est le moment d'introduire une nouvelle définition :

Définition 2 - 5

« Grand Oméga »

Soit  $f$  et  $g$  deux fonctions de  $\mathbb{R}$  dans lui-même. On dit que  $f$  est un « grand Oméga » de  $g$  et on note  $f = \Omega(g)$  ou  $f(n) = \Omega(g(n))$  si, et seulement si, il existe une constante strictement positive  $C$  telle que

$$|f(n)| \geq C|g(n)|$$

pour tout  $n \in \mathbb{N}^*$ .

Remarque

Comme  $\Omega$  est une lettre grecque, on peut, par esprit d'unification, parler de « grand omicron » au lieu de « grand O »...

Remarque

$$f = \Omega(g) \iff g = O(f) \dots$$

On montre donc facilement pour le tri fusion que  $K(n) = \Omega(n \log_2(n))$  grâce à l'inégalité  $2cn \lfloor \log_2(n) \rfloor \leq K(n)$ .

Ainsi on a dans ce cas en même temps  $f = O(n \log_2(n))$  et  $f = \Omega(n \log_2(n))$  : c'est encore plus précis et nous incite à introduire une nouvelle définition :

Définition 2 - 6

« Grand Théta »

$$f = \Theta(g) \iff \begin{cases} f = O(g) \\ f = \Omega(g) \end{cases}$$

Le coût de l'algorithme se trouve donc coincé entre deux valeurs de même ordre. On peut ainsi dire que la complexité du tri fusion est en  $n \log_2(n)$  ce qui est *souvent* mieux que le tri par insertion dont la complexité *au pire* est en  $n^2$ . Mais attention ! *Au pire* ne signifie pas *toujours* : si la liste est déjà triée, le tri par insertion est plus économique.

Voici maintenant une petite table pour illustrer les différentes classes de complexité rencontrées habituellement :

coût \ $n$	100	1000	$10^6$	$10^9$
$\log_2(n)$	$\approx 7$	$\approx 10$	$\approx 20$	$\approx 30$
$n \log_2(n)$	$\approx 665$	$\approx 10\ 000$	$\approx 2 \cdot 10^7$	$\approx 3 \cdot 10^{10}$
$n^2$	$10^4$	$10^6$	$10^{12}$	$10^{18}$
$n^3$	$10^6$	$10^9$	$10^{18}$	$10^{27}$
$2^n$	$\approx 10^{30}$	$> 10^{300}$	$> 10^{10^5}$	$> 10^{10^8}$

Gardez en tête que l'âge de l'Univers est environ de  $10^{18}$  secondes...

Il existe également deux autres « comparateurs » que l'on utilise peu en algorithmique mais qui peuvent s'avérer utiles dans d'autres domaines.

Définition 2 - 7

« Petit o »

$f = o(g)$  si, et seulement si, pour toute constante positive  $\varepsilon$ , il existe un entier  $n_0$  tel que, pour tout  $n \geq n_0$ ,

$$|f(n)| < \varepsilon |g(n)|$$

On dit alors souvent que  $f$  est *négligeable* devant  $g$ .

Contrairement au grand  $O$ , la majoration doit se faire quelque soit la constante  $\varepsilon$  et non pas seulement pour une constante arbitraire.

## Définition 2 - 8

## Fonctions équivalentes

$$f(n) \sim g(n) \iff f(n) = g(n) + o(g(n))$$

Voici quelques propriétés fort utiles que nous démontrerons à l'occasion :

Manipulation des  $O$ 

- $O(f) + O(g) = O(f + g)$  ;
- $f = O(f)$  ;
- $k \cdot O(f) = O(f)$  si  $k$  est une constante ;
- $O(O(f)) = O(f)$  ;
- $O(f) \cdot O(g) = O(f \cdot g)$  ;
- $O(f \cdot g) = f \cdot O(g)$ .

## Propriétés 2 - 1

Par exemple,  $2n^3 + 3n^2 + 5n = O(n^3) + O(n^3) + O(n^3) = O(n^3)$ .

## 4 3 Polynôme de Taylor

Brook TAYLOR fut un savant anglais qui reprit les travaux de ses aînés NEWTON, LEIBNIZ, Jacques BERNOULLI et GREGORY.

Musicien, philosophe, peintre, il est surtout célèbre pour les formules qui portent son nom. Il les obtint en étudiant le problème de KEPLER qui concerne le calcul de certaines anomalies dans la révolution des planètes.

Son but est encore une fois d'obtenir une interpolation d'une fonction compliquée à l'aide d'un polynôme simple.

Il obtient le résultat suivant :

## Polynôme de Taylor

Soit  $I$  un intervalle, soit  $f$  une fonction de Ivers  $\mathbb{R}$  de classe  $C^n$  avec  $n \in \mathbb{N}$  et soit  $a$  un élément de  $I$ . Il existe alors un *unique* polynôme  $T_{n,a}$  de degré au plus  $n$  dont les dérivées jusqu'à l'ordre  $n$  coïncident en  $a$  avec celles de  $f$ . De plus :

$$T_{n,a}(x) = f(a) + (x-a)f'(a) + \frac{(x-a)^2}{2!}f''(a) + \frac{(x-a)^3}{3!}f^{(3)}(a) + \dots + \frac{(x-a)^n}{n!}f^{(n)}(a)$$

Nous admettrons l'unicité mais *vous vérifierez* tout de même que pour tout entier  $0 \leq k \leq n$ ,

$$T_{n,a}^{(k)}(a) = f^{(k)}(a)$$

Par exemple, la fonction exponentielle admet comme polynôme de TAYLOR à l'ordre 2 en 0 :

$$\begin{aligned} T_{n,0}(x) &= \exp(0) + (x-0)\exp'(0) + \frac{(x-0)^2}{2!}\exp''(0) \\ &= 1 + x + \frac{x^2}{2} \end{aligned}$$

Le problème principal de sa formule est d'estimer l'erreur commise en considérant  $T_{n,a}$  à la place de  $f$  :

$$f(x) = T_{n,a}(x) + R_{n,a}(x)$$

Voilà pourquoi on retrouve son nom associé à trois formules qui donnent une meilleure idée de ce reste (*quelle erreur commet-on ?*) et qui ont été établies par trois de ses collègues.

Nous ne ferons pas les preuves qui sont assez techniques.

Nous admettrons également que le polynôme de TAYLOR approche *mieux*<sup>a</sup>  $f$  au voisinage de  $a$  que tous les autres polynômes de degré au plus  $n$ .

a. Ce qui signifie que  $|f(x) - T_{n,a}(x)| \leq |f(x) - P(x)|$  pour tout polynôme  $P$  de degré au plus  $n$ .



Brook TAYLOR  
(1685-1731)

## Théorème 2 - 7

**4 4 Formule de Taylor-Lagrange**

Si  $f$  est de classe  $C^n$  sur  $[a, x]$  et si  $f^{(n+1)}$  existe sur  $]a, x[$  alors il existe au moins un réel  $c \in ]a, x[$  tel que :

Théorème 2 - 8

$$f(x) = T_{n,a}(x) + \frac{(x-a)^{n+1}}{(n+1)!} f^{(n+1)}(c)$$

On dit que l'on a écrit la formule de Taylor Lagrange à l'ordre  $n$ .

Il est souvent pratique de réécrire cette formule en posant  $x = a + h$ . Le nombre  $c$  appartient donc à l'intervalle  $]a, a + h[$ . Il existe donc un réel  $\theta \in ]0, 1[$  tel que :

$$f(a+h) = \sum_{k=0}^n \frac{h^k}{k!} f^{(k)}(a) + \frac{h^{n+1}}{(n+1)!} f^{(n+1)}(a+\theta h)$$

**4 5 Formule de Taylor-Mac Laurin**

Dans la formule de Taylor-Lagrange on remplace  $a$  par 0. On obtient :

Théorème 2 - 9

$$f(x) = \sum_{k=0}^n \frac{x^k}{k!} f^{(k)}(0) + \frac{x^{n+1}}{(n+1)!} f^{(n+1)}(\theta x)$$

**4 6 Formule de Taylor avec reste intégral**

Si  $f$  est de classe  $C^{n+1}$  sur  $I$  contenant  $a$  :

Théorème 2 - 10

$$\forall x \in I, f(x) = T_{n,a}(x) + \int_a^x \frac{(x-t)^n}{n!} f^{(n+1)}(t) dt$$

**4 7 Formule de Taylor-Young.**

Si  $f$  est de classe  $C^n$  sur un intervalle  $I$  contenant  $a$  alors :

Théorème 2 - 11

$$\forall x \in I, f(x) = T_{n,a}(x) + o((x-a)^n)$$

**5 Développements limités**

**5 1 Voisinage**

$a \in \mathbb{R}$ , nous appellerons voisinage de  $a$  toute partie de  $\mathbb{R}$  contenant un intervalle ouvert contenant  $a$ . L'intervalle  $]a - \alpha, a + \alpha[$ , avec  $\alpha \in \mathbb{R}^{+*}$ , est un voisinage de  $a$ , nous le noterons  $V_a$  et  $V_a - \{a\}$  sera noté  $V_a^*$ . Un voisinage à droite de  $a$  est un intervalle du type  $]a, a + \alpha[$  et  $]a - \alpha, a[$  est un voisinage à gauche de  $a$ . Dans ce qui suit nous ne distinguerons pas ces types de voisinages, en conséquence  $V_a^*$  pourra tout aussi bien aussi désigner un voisinage à droite ou à gauche, lors des applications, le contexte permettra de les distinguer. Un voisinage de  $+\infty$  est un voisinage du type  $[A, +\infty[$  avec  $A$  réel positif aussi grand que l'on veut que nous pourrions noter  $V_{+\infty}^*$ . Notez que  $]-\infty; -A]$  est un voisinage de  $-\infty$ .

**5 2 Définition**

On dit que  $f$ , définie sur  $V_a$ , admet un développement limité d'ordre  $n \in \mathbb{N}$  au voisinage de  $a$  (on écrit en abrégé : «  $f$  admet un  $DL_n V(a)$  ») s'il existe un polynôme  $P_n \in \mathbb{R}_n[X]$  (ensemble des polynômes de degré inférieur ou égal à  $n$ ) et une fonction  $\varepsilon$  vérifiant :

Définition 2 - 9

$$\begin{cases} \forall x \in V_a, f(x) = P_n(x-a) + (x-a)^n \varepsilon(x) \\ \varepsilon(a) = 0 \\ \lim_{x \rightarrow a} \varepsilon(x) = 0 \end{cases}$$

c'est-à-dire  $f(x) = P_n(x-a) + o((x-a)^n)$ .

Il faut bien remarquer que :

$$q > p \rightarrow (x-a)^q = o_a((x-a)^p)$$

Plus la puissance  $p$  est grande, plus  $(x-a)^p$  est petit.

Nous n'irons pas très loin dans l'étude des développements limités mais elle nous permettra d'avoir des approximations polynomiales de nombreuses fonctions ce qui nous permettra de comprendre de nombreuses méthodes employées en informatique.

## EXERCICES

### Recherche 2 - 1 Formule de dérivation

En utilisant les résultats de la section 2.1.2 page 39, démontrez la formule de dérivation de  $x \mapsto x^n$  et  $x \mapsto 1/x^n$  pour tout entier positif  $n$ .

### Recherche 2 - 2 Méthode de Briggs

Inspirez-vous de la méthode présentée à la section 2.2.2 page 40 pour calculer une approximation de  $\log(3)$  et  $\log(7)$ .

### Recherche 2 - 3 Interpolation de Newton(1) : le gag

Déterminez le polynôme de NEWTON passant par les points (1,1), (2,8), (3,27) et (4,64).

### Recherche 2 - 4 Interpolation de Newton(2)

Déterminez le polynôme de NEWTON passant par les points (1,5), (2,2), (3,5), (4,2) et (5,2).

### Recherche 2 - 5

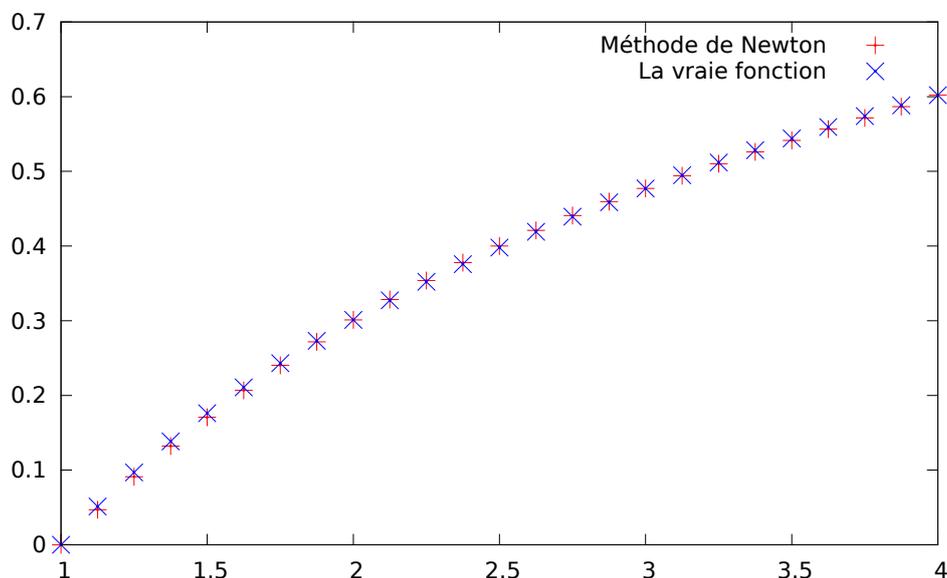
Écrivez la formule du théorème 2 - 6 page 44 sans ... mais avec  $\Sigma$  et  $\Pi$

### Recherche 2 - 6 Calculs des coefficients des polynômes de Newton

Reprenez les calculs effectués dans la section 2.2.3.3 page 43 mais avec  $n$  points d'abscisses 0 à  $n - 1$ .

### Recherche 2 - 7 Polynôme de Newton : le programme

Nous voudrions programmer l'interpolation du logarithme décimal sur [1,4] pour obtenir un tracé tel que celui-ci :



obtenu en utilisant la bibliothèque Easyplot (`import Graphics.EasyPlot`) :

```

1 compareTrace f a b pas =
2   plot X11
3   -- ou bien (PDF ("plot_" ++ nomf ++ "_newton.pdf")) pour un export PDF
4   [
5     Data2D [Title "Méthode de Newton", Color Red] [] [(x, interpol liste x) | x <- [a, (a + pas) .. b]],
6     Data2D [Title "La vraie fonction", Color Blue] [] [(x, f x) | x <- [a, (a + pas) .. b]]
7   ]
8   where liste = map f [a .. b]
```

On pourra créer 5 *one liners* :

```

1  -- Calcule la liste des différences 2 à 2 d'une liste de nbs [1,3,6] -> [2,3]
2  réduitL :: Num a => [a] -> [a]
3  réduitL liste = zipWith ???
4
5  -- Renvoie la liste des listes de différences : le triangle de Newton
6  réduit :: (Eq a, Num a) => [a] -> [[a]]
7  réduit liste = takeWhile ???
8
9  -- Renvoie la tête du triangle de Newton
10 diffNewton :: Num a => [a] -> [a]
11 diffNewton liste = map ???
12
13 -- Renvoie la liste des évaluations des (x-x0)(x-x0-1).../n! du théorème 2.6 accumulées dans une liste
14 baseNewton :: (Enum a, Fractional a) => a -> a -> a -> [a]
15 baseNewton x0 x n = scanl (\ acc k -> ??? ) 1 [1 .. n]
16
17 -- renvoie l'évaluation du poly de Newton en x en partant de x0 avec une liste de valeurs
18 -- l'incrément des abscisses est 1
19 interpol :: (Enum a, Eq a, Fractional a) => [a] -> a -> a -> a
20 interpol liste x0 x = ???

```

- Commencez par créer une fonction `reduitL :: Num a => [a] -> [a]` qui renvoie la liste des écarts entre les éléments d'une liste donnée :

```

1  *Main> réduitL [1,2,5,7,12]
2  [1,3,2,5]

```

- Déterminez ensuite une fonction `reduit :: Num a => [a] -> [[a]]` qui la liste des différentes listes des écarts à partir d'une liste donnée :

```

1  *Main> réduit [5,2,5,2,2]
2  [[-3,3,-3,0],[6,-6,3],[-12,9],[21]]

```

- Déterminez une fonction `diffNewton :: Num a => [a] -> [a]` qui renvoie la liste des  $\Delta^i y_1$  du « triangle » de NEWTON (ceux en gras...) :

```

1  *Main> diffNewton [5,2,5,2,2]
2  [5,-3,6,-12,21]

```

- La formule du théorème 2 - 5 page 44 laisse entrevoir un `zipWith`. Nous disposons déjà de la liste des  $\Delta^i$ . Il nous reste à construire la *base* des polynômes de NEWTON `baseNewton :: (Enum a, Eq a, Fractional a) => a -> a -> a -> [a] ::`

```

1  baseNewton :: (Enum a, Fractional a) => a -> a -> a -> [a]
2  baseNewton x0 x n = scanl (\ acc k -> ??????? ) 1 [1 .. n]

```

Il n'y a plus qu'à zipper et additionner (bref faire un produit scalaire : une liste est un vecteur...) pour créer une fonction `interpol :: (Enum a, Eq a, Fractional a) => [a] -> a -> a -> a` :

```

1  *Main> interpol [1,8,27,64] 1 1.5
2  3.375

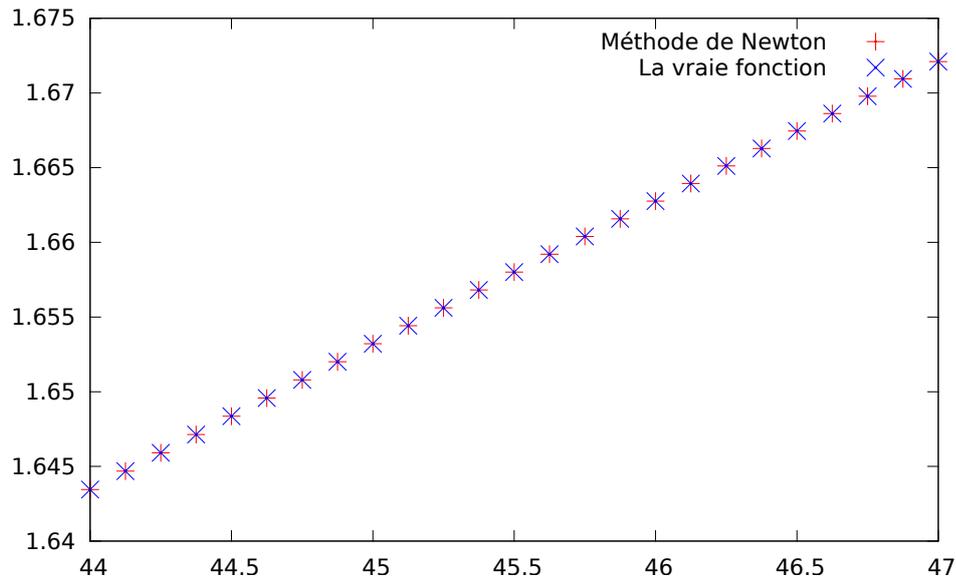
```

- Reprendre toutes les fonctions précédentes pour être compatible avec le théorème 2 - 6 page 44. Par exemple :

```

1  *Main> compareTrace log10 44 47 0.125

```



**Recherche 2 - 8 Algorithme de Horner**

La méthode que nous allons voir porte le nom du britannique William George HORNER (1786 - 1837) mais en fait elle fut publiée presque 10 ans auparavant par un horloger londonien, Theophilus HOLDRED et simultanément par l'italien Paolo RUFFINI (1765 - 1822) mais fut déjà utilisée par NEWTON 150 ans auparavant et par le chinois ZHU SHIJE cinq siècles plus tôt (vers 1300) et avant lui par le Persan SHARAF AL-DIN AL-MUZAFFAR IBN MUHAMMAD IBN AL-MUZAFFAR AL-TUSI vers (1100) et avant lui par le Chinois LIU HUI (vers 200) révisant un des résultats présent dans *Les Neuf Chapitres sur l'art mathématique* publié avant la naissance de JC...



Soit  $P$  un polynôme. Il s'agit de calculer l'évaluation de  $P$  en une valeur particulière. Prenons l'exemple de  $P(x) = 3x^5 - 2x^4 + 7x^3 + 2x^2 + 5x - 3$ . Le calcul classique nécessite 5 additions et 15 multiplications (vérifiez-le!).

On peut faire pas mal d'économies de calcul en suivant le schéma suivant :

$$\begin{aligned}
 P(x) &= \underbrace{a_n x^n + \dots + a_2 x^2 + a_1 x + a_0}_{\text{on met } x \text{ en facteur}} \\
 &= \left( \underbrace{a_n x^{n-1} + \dots + a_2 x + a_1}_{\text{on met } x \text{ en facteur}} \right) x + a_0 \\
 &= \dots \\
 &= (\dots(((a_n x + a_{n-1})x + a_{n-2})x + a_{n-3})x + \dots)x + a_0
 \end{aligned}$$

Ici cela donne  $P(x) = (((((3x) - 2)x + 7)x + 2)x + 5)x - 3$  c'est-à-dire 5 multiplications et 5 additions. En fait il y a au maximum  $2 \times \text{degré}$  de  $P$  opérations (voire moins avec les zéros). Que pensez-vous de ce *one liner* ?

```
1 horner p x = foldl (\ pol coeff -> coeff + x * pol) 0 p
```

Notez de plus qu'une machine disposant de la FMA effectue ce calcul encore plus exactement.

```
1 *Main> horner [3,-2,7,2,5,-3] 1
2 12
```

Écrivez un programme similaire en C. Reproduisez le graphe de la section 2.3 page 44.

**Recherche 2 - 9 Méthode d'Euler**

On ne dispose que de la dérivée d'une fonction : peut-on tracer malgré cela l'allure de sa courbe représentative ?

On créera une fonction Haskell permettant par exemple de tracer la courbe représentative de la fonction dont la dérivée est la fonction  $x \mapsto x$  et qui s'annule pour  $x = 1$ .

Généralisez votre fonction Haskell pour faire de même avec une fonction qui vérifie une équation différentielle

$$y' = f(x, y)$$

On commencera par créer la liste des points `listeEuler f x y pas n` qui renvoie la liste des coordonnées des points de la courbe approchée.

On l'utilisera ensuite dans la fonction :

---

```

1 euler f x0 y0 xf n =
2   let pas = (xf - x0) / n in
3   plot X11 (
4     Data2D [Title "Méthode d'Euler", Color Blue, Style Lines] [] (listeEuler f x0 y0 pas n)
5   )

```

---

On obtient par exemple avec :

---

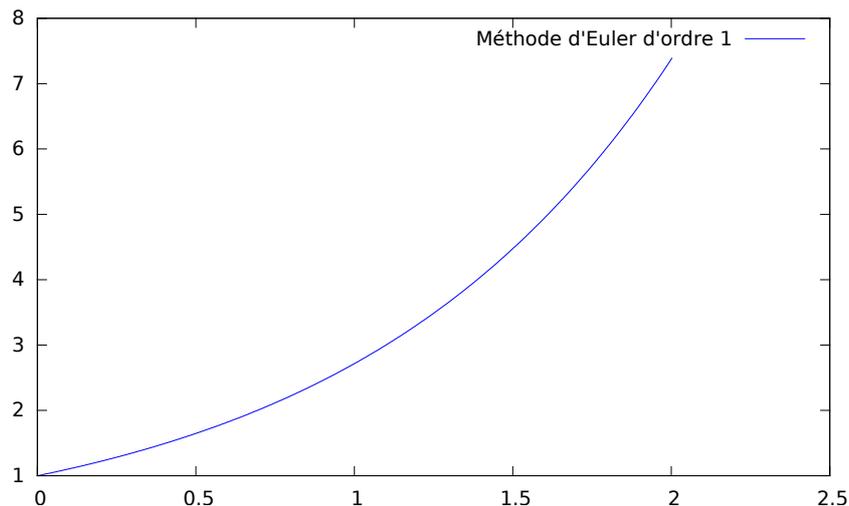
```

1 *Main> euler (\x y -> y) 0 1 2 (2**10)

```

---

la courbe suivante :



Essayez d'étendre au cas des équations différentielles d'ordre 2 :  $y'' = f(x, y, y')$ .

Par exemple :

---

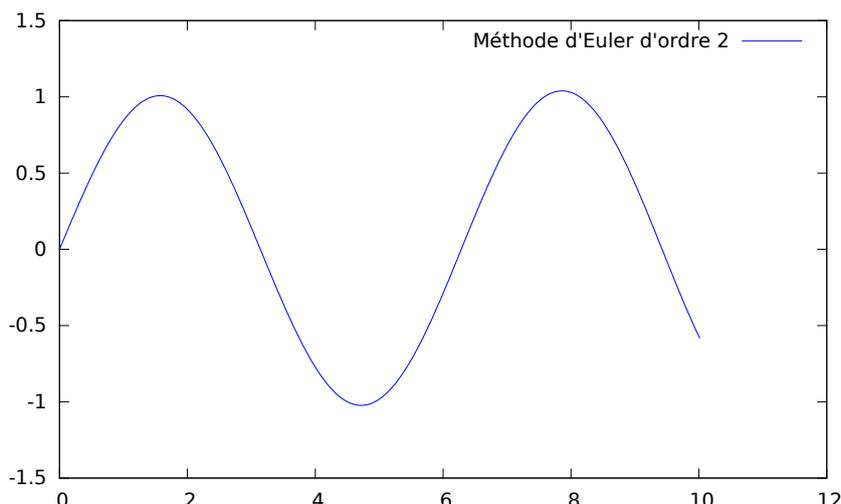
```

1 *AnalyseNumerique> euler2 (\ x y y' -> - y) 0 0 1 10 (2**10)

```

---

donne :



## Recherche 2 - 10 Construction de l'exponentielle

### Une équation différentielle

L'équation différentielle  $f' = kf$  se retrouve dans de nombreux problèmes : désintégration des noyaux des atomes d'un corps radioactif, datation au carbone 14, évolution d'une population où la croissance est proportionnelle au nombre d'habitants, etc. Le problème est de trouver une fonction la satisfaisant.

Par exemple, certains phénomènes en mécanique conduisent à étudier l'équation différentielle  $f'' = -f$ . Nous connaissons au moins deux fonctions la satisfaisant : cosinus et sinus.

### Construction approchée du graphe d'une solution par la méthode d'Euler

Soit  $f$  une fonction dérivable sur  $\mathbb{R}$  vérifiant  $f(0) = 1$  et, pour tout  $x$ ,  $f'(x) = f(x)$ . Utilisez la méthode d'Euler pour déterminer une représentation graphique de  $f$  sur  $[0, 2]$ .

### Analyse : étude des propriétés mathématiques d'une solution

Soit  $f$  une fonction dérivable sur  $\mathbb{R}$  vérifiant  $f(0) = 1$  et, pour tout  $x$ ,  $f'(x) = kf(x)$ , avec  $k \neq 0$ .

1. Montrez que  $f'(0) = k$ .
2. Soit  $y$  un réel fixé et  $g_y$  la fonction définie par  $g_y(x) = f(x+y)f(-x)$ 
  - i. Montrez que  $g_y$  est dérivable sur  $\mathbb{R}$  et calculez  $g'_y(x)$ .
  - ii. Calculez  $g_y(0)$  et déduisez-en que pour tous  $x$  et  $y$  réels,  $f(x+y)f(-x) = f(y)$  (1)
3. Montrez alors successivement que :
  - i. pour tout réel  $x$ ,  $f(x)f(-x) = 1$
  - ii.  $f$  ne s'annule pas sur  $\mathbb{R}$
  - iii. pour tous réels  $x$  et  $y$ ,  $f(x+y) = f(x)f(y)$ .

### Unicité de la fonction solution

Peut-on trouver une autre fonction,  $\varphi$ , distincte de  $f$ , et vérifiant les mêmes propriétés que  $f$ , à savoir :  $\varphi$  est une fonction dérivable sur  $\mathbb{R}$ , vérifiant  $\varphi(0) = 1$  et, pour tout  $x$ ,  $\varphi'(x) = k\varphi(x)$ , avec  $k \neq 0$  ?

Comme  $f$  ne s'annule pas sur  $\mathbb{R}$ , on peut définir la fonction  $\psi = \varphi/f$ .

Vérifiez que  $\psi$  est dérivable sur  $\mathbb{R}$ , calculez sa dérivée. Que peut-on en déduire pour  $\psi$  ? Montrez alors que  $f = \varphi$ .

### Synthèse

Nous avons cherché des solutions au problème :

*$f$  une fonction dérivable sur  $\mathbb{R}$  vérifiant  $f(0) = 1$  et, pour tout  $x$ ,  $f'(x) = kf(x)$ , avec  $k \neq 0$ .*

Nous avons montré que, si une telle fonction existe (ce que nous prouverons dans un prochain chapitre), alors elle est unique et elle vérifie nécessairement la relation  $f(x+y) = f(x)f(y)$  (2).

Il reste à vérifier que, réciproquement, une fonction dérivable, non nulle, vérifiant la relation (2) est nécessairement telle que  $f(0) = 1$  et vérifie pour tout réel  $x$   $f'(x) = kf(x)$ , avec  $k$  un réel non nul.

Cette vérification n'est pas anodine et conclut notre raisonnement d'*analyse-synthèse*.

1. Montrez que  $f$  ne s'annule pas et que  $f$  est à valeurs strictement positives.
2. Montrez que, comme  $f$  n'est pas la fonction nulle, alors  $f(0) = 1$  en utilisant la relation (2).
3. Soit  $a$  un réel fixé. On définit la fonction  $\varphi : x \mapsto f(x+a)$  et la fonction  $\psi : x \mapsto f(x) \times f(a)$ .  
Montrez que  $f'(x+a) = f(a) \times f'(x)$ , puis que, pour tout réel  $a$ ,  $f'(a) = k f(a)$ , où  $k$  est un réel que vous déterminerez.

### Recherche 2 - 11



La loi de Newton sur le refroidissement dit que la vitesse de refroidissement d'un objet est proportionnelle à l'écart de température entre l'objet et le milieu ambiant.

L'inspecteur CLOUSEAU arrive sur les lieux d'un meurtre à 9h00. Il commence par prendre la température de la victime : 30 °C. Une heure plus tard, la température du corps est tombée à 29 °C. Sachant que la température normale du corps d'une personne vivante en bonne santé est de 37 °C, que la victime était syldave, qu'elle aimait les films de gladiateurs et se trouvait dans une pièce maintenue à 0 °C, estimez l'heure du décès de la victime et la couleur de ses yeux.

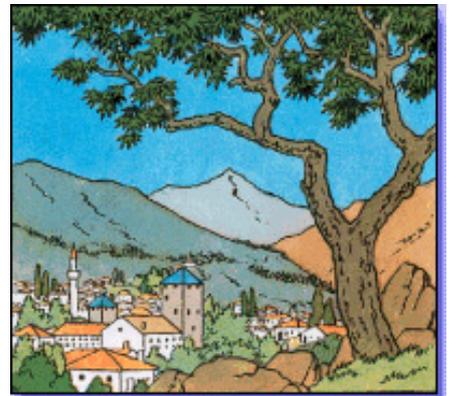
On pourra avoir une idée graphique avec :

```
1 *Main> euler (\ x y -> k * y ) 10 29 2 100 where k = ...
```

### Recherche 2 - 12

Dans la forêt syldave, des débris naturels ( feuilles, branches, animaux morts, cadavres d'espions, etc.) tombent sur le sol et s'y décomposent. La quantité  $Q(t)$  exprimée en  $g \cdot m^{-2}$  de débris jonchant le sol varie avec le temps  $t$ . On suppose que de nouveaux débris tombent au sol à un taux constant de  $200 g \cdot m^{-2}$  par année et que les débris accumulés au sol se décomposent au taux de 50% de la quantité de débris jonchant le sol.

1. On note  $f(t) = Q(t) - 200$ . Déterminez une équation différentielle vérifiée par  $f$ .
2. En déduire l'expression générale de  $Q(t)$ .
3. Exprimer  $Q(t)$  sachant qu'au temps  $t = 0$ , on comptait  $50 g \cdot m^{-2}$



### Recherche 2 - 13



Si vous allez vous promener dans la ville de Saint-Louis dans le Missouri aux États-Unis, vous pourrez y admirer la célèbre *Gateway Arch to the West* conçue en 1947 par l'architecte finlandais Ero SAARINEN et l'ingénieur Hannskarl BANDEL et dont la construction s'acheva en 1965. Cette arche a la forme d'une *chaînette* pondérée d'équation

$$y = 212 - 21 \operatorname{ch}(0,033x)$$

où  $x$  et  $y$  sont mesurés en mètres et rend hommage aux pionniers partis à la conquête de l'Ouest.

Pouvez-vous déterminer la hauteur de l'arche et la distance entre les deux pieds ? On rappelle que  $\operatorname{ch} x = \frac{e^x + e^{-x}}{2}$ .

### Recherche 2 - 14

L'équation de la hauteur  $h$  par rapport au sol d'un fil électrique suspendu entre deux poteaux s'obtient en résolvant l'équation différentielle :

$$h''(x) = k\sqrt{1 + (h'(x))^2}$$

où  $k$  est un paramètre qui dépend de la densité et de la tension du fil et  $x$  est mesuré en mètres horizontalement à partir d'une origine située sur le sol en-dessous du point où la hauteur du fil est la plus faible.

1. Vérifiez que  $h : x \mapsto \frac{1}{k} \operatorname{ch}(kx)$  satisfait cette équation différentielle.
2. Quelle est la hauteur minimale du fil si le paramètre  $k$  vaut 0,05 ?
3. Quelle est la hauteur des poteaux (de même hauteur) s'ils sont distants de 30 m et que le paramètre  $k$  vaut 0,05 ?



**Recherche 2 - 15 DL**

Écrire la formule de Taylor-Mac Laurin et la formule de Taylor-Young à l'ordre  $n$  pour les fonctions suivantes et dans un voisinage de 0 :

- |                           |                           |                               |
|---------------------------|---------------------------|-------------------------------|
| 1. $f(x) = e^x$           | 3. $f(x) = \frac{1}{1-x}$ | 5. $f(x) = \sin(x)$           |
| 2. $f(x) = \frac{1}{1+x}$ | 4. $f(x) = \ln(1+x)$      | 6. $f(x) = \sqrt{1+x}, n = 4$ |

**Recherche 2 - 16 DL**

Donner un  $DL_3(0)$  pour les fonctions suivantes :

- |                        |                      |                               |
|------------------------|----------------------|-------------------------------|
| 1. $f_1(x) = \ln(1+x)$ | 3. $f_2(x) = e^x$    | 5. $f_3(x) = f_1(x)f_2(x)$    |
| 2. $g(x) = f_1(x^2)$   | 4. $h(x) = f_2(-2x)$ | 6. $f_4(x) = f_1(x) + f_2(x)$ |

**Recherche 2 - 17 Fonctions homographiques**

Une fonction homographique est une fonction de la forme  $x \mapsto \frac{ax+b}{cx+d}$  avec  $a, b, c$  et  $d$  des réels. Notons  $g$  la fonction homographique correspondant au cas  $a = 0, b = d = 1$  et  $c = -1$ .

Quel est l'ensemble de définition de  $g$  ?  
 Calculez  $g^{(1)}(x), g^{(2)}(x), g^{(3)}(x)$  et plus généralement  $g^{(n)}(x)$ .  
 Écrivez le plus simplement possible  $T_{5,0}(g)(x)$ .

Sur un même graphique, donnez l'allure des courbes représentatives de  $g$  et de  $T_{5,0}(g)$  avec Haskell.

Soit  $h(x) = \frac{5-3x}{1-x}$ . Montrez qu'il existe deux réels  $\alpha$  et  $\beta$  tels que  $h(x) = \alpha + \frac{\beta}{1-x}$ .  
 Déduisez-en  $T_{5,0}(h)(x)$  pour  $x$  appartenant à un bon intervalle.

**Recherche 2 - 18 Question ouverte...**

Les polynômes de NEWTON et de TAYLOR se ressemblent non ? Dans quelle mesure ?

**Recherche 2 - 19 Courbes de Bézier**

Dans les années 60, les ingénieurs Pierre BÉZIER et Paul DE CASTELJAU travaillant respectivement chez Renault et Citroën, réfléchissent au moyen de définir de manière la plus concise possible la forme d'une carrosserie.

Le principe a été énoncé par BÉZIER mais l'algorithme de construction par son collègue de la marque aux chevrons qui n'a d'ailleurs été dévoilé que bien plus tard, la loi du secret industriel ayant primé sur le développement scientifique...

Pour la petite histoire Pierre BÉZIER (diplômé de l'ENSAM et de SUPÉLEC) fut à l'origine des premières machines à commandes numériques et de la CAO ce qui n'empêcha pas sa direction de le mettre à l'écart : il se consacra alors presque exclusivement aux mathématiques et à la modélisation des surfaces et obtint même un doctorat en 1977.

Paul DE CASTELJAU était lui un mathématicien d'origine, ancien élève de la Rue d'ULM, qui a un temps été employé par l'industrie automobile.

Aujourd'hui, les courbes de BÉZIER sont très utilisées en informatique.

Une Courbe de BÉZIER est une courbe paramétrique aux extrémités imposées avec des points de contrôle qui définissent les tangentes à cette courbe à des instants donnés.



**Algorithme de Casteljau**

Soit  $t$  un paramètre de l'intervalle  $[0, 1]$  et  $P_1, P_2$  et  $P_3$  les trois points de contrôle.

On construit le point  $M_1$  barycentre du système  $\{(P_1, 1-t), (P_2, t)\}$  et  $M_2$  celui du système  $\{(P_2, 1-t), (P_3, t)\}$ .

On construit ensuite le point  $M$ , barycentre du système  $\{(M_1, 1-t), (M_2, t)\}$ .

Exprimez  $M$  comme barycentre des trois points  $P_1, P_2$  et  $P_3$ .

Faites la construction à la main avec  $t = 1/3$  par exemple puis utilisez Haskell.

Nous allons assimiler les points  $M_1, M_2$  et  $M$  à des courbes paramétrées.

Ainsi  $M_1(t) = (1-t)P_1 + tP_2$ ,  $M_2(t) = (1-t)P_2 + tP_3$  puis

$$M(t) = (1-t)M_1(t) + tM_2(t) = (1-t)^2P_1 + 2t(1-t)P_2 + t^2P_3$$

Vérifiez que  $M'(t) = 2(M_2(t) - M_1(t))$ . Comment l'interpréter ?

**Avec 4 points de contrôle**

Faites une étude similaire (« à la main ») avec 4 points de contrôle.

On pourra utiliser une représentation similaire aux arbres de probabilité.

Créez une fonction **bezier3** qui trace la courbe de BEZIER de degré 3 passant par une liste de 4 points de contrôle.

On pourra créer un type synonyme `type Point = (Float,Float)` et utiliser EasyPlot :

---

```

1 type Point = (Float,Float)
2
3 infixr 4 @* -- produit externe
4 (@*) :: Float -> Point -> Point
5 (@*) k p =
6
7 infixr 2 @+ -- somme de points
8 (@+) :: Point -> Point -> Point
9 (@+) p1 p2 =
10
11 infixr 2 @- -- différence de points
12 (@-) :: Point -> Point -> Point
13 (@-) p1 p2 =
14
15 bezier3 :: [Point] -> [Graph2D Float Float]
16 bezier3 [p0,p1,p2,p3] =
17     let m = \ t -> ... in
18         [
19             Data2D [Title "Bezier 3", Color Blue, Style Lines] [] [m t | t <- [0, 0.01 .. 1]],
20             Data2D [Title "Départ", Color Green, Style Lines] [] [p0,p1],
21             Data2D [Title "Arrivée", Color Red, Style Lines] [] [p2,p3]
22         ]
23 bezier3 _ = ...

```

---

Par exemple :

---

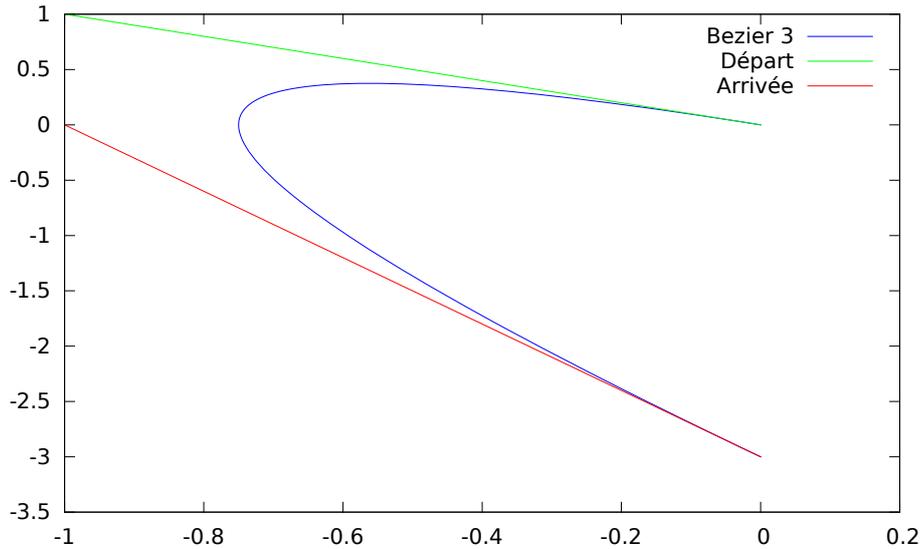
```

1 *Main> plot X11 $ bezier3 [(0,0),(-1,1),(-1,0),(0,-3)]

```

---

donne :



Quel est le rôle des différents points de contrôle ?

**Courbe de Bézier du 3<sup>e</sup> degré avec un nombre quelconque de points de contrôle**

Il est pratique de travailler avec des polynômes de degré trois pour avoir droit à des points d'inflexion.

Augmenter le nombre de points de contrôle implique a priori une augmentation du degré de la fonction polynomiale. Pour remédier à ce problème, on découpe une liste quelconque en liste de listes de 4 points.

Cependant, cela est insuffisant pour obtenir un raccordement de classe  $C^1$  (pourquoi ? pourquoi est-ce important d'avoir un raccordement de classe  $C^1$  ?)

Pour assurer la continuité tout court, il faut que le premier point d'un paquet soit le dernier du paquet précédent.

Le dernier « vecteur vitesse » de la liste  $[P_1, P_2, P_3, P_4]$  est  $\overrightarrow{P_3P_4}$ . Il faut donc que ce soit le premier vecteur vitesse du paquet suivant pour assurer la continuité de la dérivée.

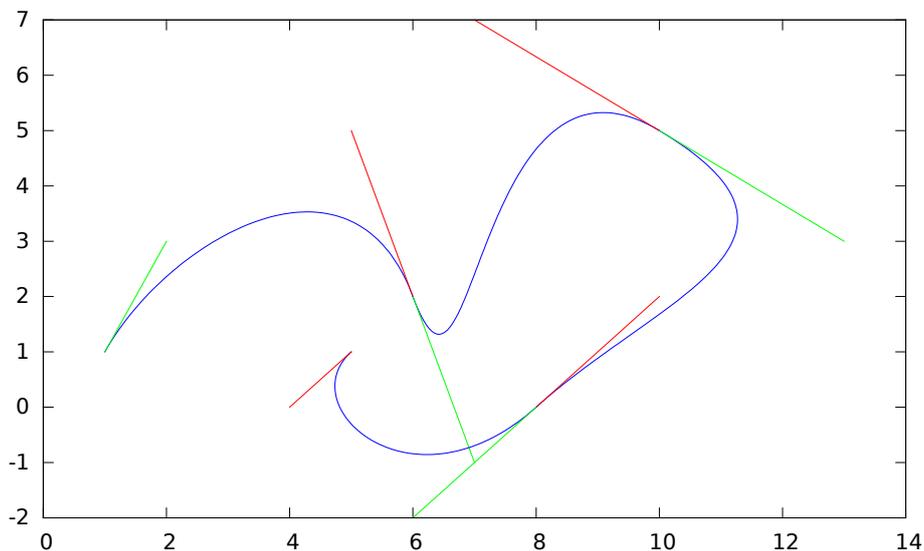
Appelons provisoirement le paquet suivant  $[P'_1, P'_2, P'_3, P'_4]$ . On a d'une part  $P'_1 = P_4$  et d'autre part  $\overrightarrow{P_3P_4} = \overrightarrow{P'_1P'_2}$ , i.e.  $P'_2 = P_4 + \overrightarrow{P_3P_4}$ .

On a donc  $P'_3 = P_5$  et  $P'_4 = P_6$ .

Connaissant `bezier3`, construire une fonction qui trace une courbe de BÉZIER cubique avec un nombre quelconque de points de contrôle (on prendra un nombre pair de points pour se simplifier la vie).

Testez avec :

```
1 plot x11 [(1, 1), (2, 3), (5, 5), (6, 2), (7, 7), (10, 5), (10, 2), (8, 0), (4, 0), (5, 1)]
```



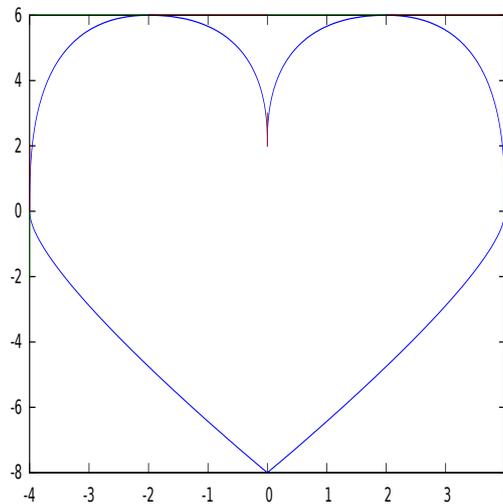
On pourra, par exemple, construire quelques fonctions intermédiaires :

```

1  bezier :: [Point] -> [Point]
2  -- Rajoute les points pour que la courbe soit de classe C1
3  bezier (p1:p2:p3:p4:xs) =
4      ...
5  bezier _          = []
6
7  bezierP :: [Point] -> [Graph2D Float Float]
8  -- Colle les bouts de bezier3
9  bezierP []       = []
10 bezierP listeP = ...
11
12 bezier :: TerminalType -> [Point] -> IO Bool
13 -- crée le graphique avec le type de terminal en option (X11, PDF, ...)
14 bezier option listeP = plot option (bezierP $ bezier listeP)

```

C'était il y a peu la Saint-Valentin alors dessinez un cœur...



### B-splines uniformes

Tout ceci est très beau mais il y a un hic : en changeant un point de contrôle, on modifie grandement la figure.

On considère  $m$  nœuds  $t_0, t_1, \dots, t_m$  de l'intervalle  $[0, 1]$ .

Introduisons une nouvelle fonction :

$$S(t) = \sum_{i=0}^{m-n-1} P_i b_{i,n}(t), \quad t \in [0, 1]$$

les  $P_i$  étant les points de contrôle et les fonctions  $b_{j,n}$  étant définies récursivement par

$$b_{j,0}(t) = \begin{cases} 1 & \text{si } t_j \leq t < t_{j+1} \\ 0 & \text{sinon} \end{cases}$$

et pour  $n \geq 1$

$$b_{j,n}(t) = \frac{t - t_j}{t_{j+n} - t_j} b_{j,n-1}(t) + \frac{t_{j+n+1} - t}{t_{j+n+1} - t_{j+1}} b_{j+1,n-1}(t).$$

On ne considérera par la suite que des nœuds équidistants : ainsi on aura  $t_k = \frac{k}{m}$ .

On parle de B-splines uniformes et on peut simplifier la formule précédente en remarquant également des invariances par translation.

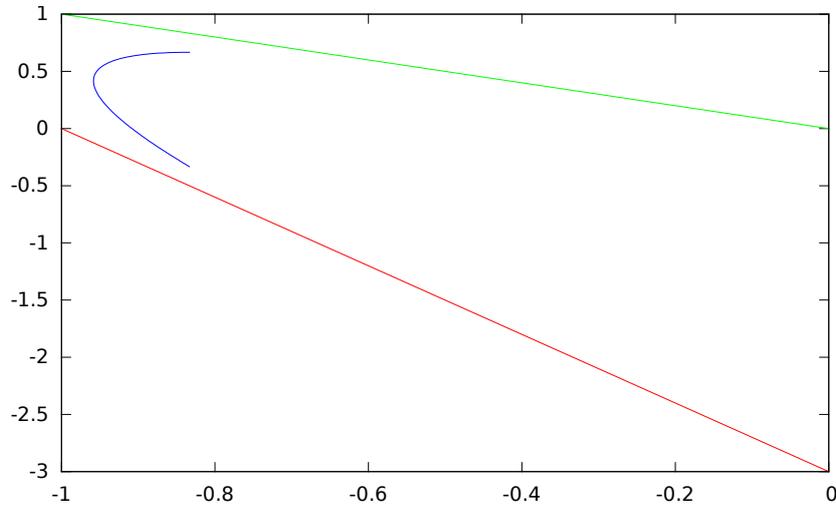
À l'aide des formules précédentes, on peut prouver que dans le cas de 4 points de contrôle on obtient :

$$S(t) = \frac{1}{6} \left( (1-t)^3 P_0 + (3t^3 - 6t^2 + 4) P_1 + (-3t^3 + 3t^2 + 3t + 1) P_2 + t^3 P_3 \right)$$

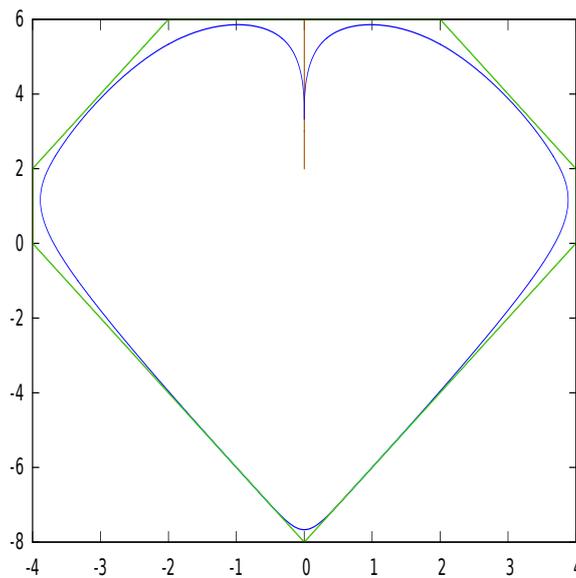
Calculez  $S(0)$ ,  $S(1)$  puis  $S'(0)$  et  $S'(1)$  : que peut-on en conclure ?

Reprenez l'étude faite avec les courbes de BÉZIER

```
1 *Main> plot X11 $ bspline3 [(0,0),(-1,1),(-1,0),(0,-3)]
```



Obtiendrez-vous un plus joli cœur ?



**Recherche 2 - 20 2 points**

Déterminez la dérivée de la fonction  $x \mapsto \sqrt{(\sqrt{(\sqrt{(\sqrt{x})})})}$ . Vous donnerez l'expression finale sous la forme la plus simple possible.

**Recherche 2 - 21 6 points**

Déterminez, en détaillant vos calculs, les polynômes de TAYLOR à l'ordre 4 au voisinage de 0 de :

- 1.  $f_1(x) = \ln(1 - x)$
- 2.  $f_2(x) = \exp(-x)$
- 3.  $f_3(x) = \cos(2x)$



# Méthodes itératives

Comme leur nom l'indique, les méthodes itératives (i.e. les suites définies par une relation  $x_{n+1} = f(x_n)$ ) sont intimement liées à la notion de boucle. Nous verrons ensuite comment elles permettent d'effectuer certains calculs sur machine.

Mais ce n'est qu'un point de départ... Henri POINCARÉ s'est intéressé à la fin du XIX<sup>e</sup> siècle à l'évolution des systèmes dynamiques en partant d'un problème de mécanique céleste : le soleil est beaucoup plus massif que les planètes, c'est pourquoi les planètes décrivent des trajectoires elliptiques autour de lui. Cependant, chaque planète est également influencée par ses congénères mais dans une mesure bien moindre : les trajectoires elliptiques sont certes déformées, mais de manière extrêmement lente. Que va-t-il se passer dans le long terme : ces micro-changements vont-ils avoir une grande influence ou bien rester négligeables ? L'étude des systèmes dynamiques était née : on étudie un système très simple a priori, décrit par un petit nombre de paramètres dont la loi d'évolution est déterminée, et on étudie son évolution à long terme : la simplicité des systèmes cache parfois une incroyable complexité de l'évolution à long terme.

Cette théorie a été popularisée par la médiatisation de certains résultats de la théorie du chaos, de l'étude des fractales. En informatique, elle a influencé par exemple certains modes de compression d'images.

Elle pose aussi le problème de la prédictibilité : le futur peut-il être déduit du présent ? C'est un domaine passionnant, riche mathématiquement et en lien avec la physique, la biologie, les sciences de l'ingénieur, l'informatique... C'est enfin un domaine mathématique très vivant : l'un des gagnants de la médaille Fields 2014 est en effet le franco-brésilien **Artur AVILA** dont le domaine de recherche est justement lié aux systèmes dynamiques.



# 1 Une boucle sous toutes ses formes

## 1 1 Généralité

Regardons une simple boucle :

```

1  Algorithme Une boucle
2  r ← r0
3  Pour k de 1 à n Faire
4  |   r ← f(r)
5  FinPour
6  Retourner r

```

On construit ainsi une *suite* de valeurs prises par la variable  $r$  qui varie avec l'incrément  $k$ . Si on note  $r_k$  la valeur de la variable  $r$  lorsque l'incrément vaut  $k$  :

$$r_{k+1} = f(r_k)$$

L'algorithme boucle peut d'ailleurs être écrit récursivement :

```

1  Fonction boucle(r : type 1 n : entier naturel) : type 2
2  Si n == 0 Alors
3  |   Retourner r
4  Sinon
5  |   Retourner boucle(f(r), n-1)
6  FinSi

```

Ce qui peut s'écrire de manière plus concise à l'aide d'un *pliage* :

```

1  boucle f r0 n = foldl (\ r k -> f(r)) r0 [1 .. n]

```

et pour avoir la liste de toutes les valeurs intermédiaires on remplace `foldl` par `scanl`.

Le programmeur voudrait savoir si sa boucle va effectivement renvoyer ce qui est prévu dans sa spécification et surtout en combien de temps : s'il s'agit de calculer la vitesse d'une voiture en fonction de la distance qui la sépare de la voiture qui est en face, il faut faire vite.

## 1 2 Complexité de la somme

Nous allons étudier le plus classique des problèmes : il s'agit de calculer la somme des entiers de 1 à un entier naturel  $n$  non nul donné.

La *spécification* est immédiate : on crée une fonction ayant pour argument l'entier naturel  $n$  et renvoyant l'entier naturel correspondant à la somme des entiers de 1 à  $n$ .

Il reste à s'occuper de la *réalisation*.

## 1 3 Procédure et processus

Appelons  $S(n)$  la somme des entiers de 1 à  $n$  avec  $n \in \mathbb{N} \setminus \{0\}$ .

On a  $S(1) = 1$  et  $S(n) = n + S(n-1)$ .

```

1  s1 :: Integer -> Integer
2  s1 n
3  | n == 1    = 1
4  | n > 1    = n + s1 (n - 1)
5  | otherwise = error "l'argument doit être un entier naturel non nul"

```

Que se passe-t-il quand on calcule  $s1(5)$  ? On utilise le modèle de *substitution* :

```
s1(5)
5 + s1(4)
5 + (4 + s1(3))
5 + (4 + (3 + s1(2)))
5 + (4 + (3 + (2 + s1(1))))
5 + (4 + (3 + (2 + 1)))
5 + (4 + (3 + 3))
5 + (4 + 6)
5 + 10
15
```

On commence par quatre expansions suivies de cinq réductions.

La *procédure* est récursive car syntaxiquement  $s1$  apparaît dans sa propre définition (`return s1(n - 1) + n`).

Le *processus* est également récursif : les évaluations sont retardées. Intuitivement, cela est le cas car on progresse du cas compliqué vers le cas simple : on « descend » ici de  $n$  vers 1. L'interpréteur doit *empiler* les résultats des expansions tant qu'il ne peut pas les réduire. La taille de la pile nécessaire étant proportionnelle à l'argument  $n$ , on dit qu'il s'agit d'un processus linéairement récursif.

On peut avoir une autre idée : partir du cas simple et « monter » vers le cas compliqué (on parle dans ce cas habituellement d'*induction*).

Ici, la somme vaut 0 puis on ajoute 1, puis on ajoute au résultat 2, puis on ajoute au résultat 3, etc jusqu'à atteindre  $n$ .

Techniquement, on entretient un *accumulateur* de termes qui évolue et un compteur qui va évoluer de 1 jusque  $n$ .

L'accumulateur et le compteur évoluent selon la règle :

```
accumulateur ← accumulateur + compteur
compteur ← compteur + 1
```

la somme cherchée étant la valeur de l'accumulateur lorsque le compteur a dépassé  $n$ .

---

```
1  s2 :: Integer -> Integer
2  s2 n =
3      let s2rec acc compt =
4          if compt > n then
5              acc
6          else
7              s2rec (acc + compt) (compt + 1)
8      in s2rec 0 1
```

---

La procédure `s_iter` est syntaxiquement récursive car elle est définie à partir d'elle-même mais le processus n'est plus récursif ! les évaluations ne sont plus retardées :

```
s2(5)
s_iter(0 , 1)
s_iter(1 , 2)
s_iter(3 , 3)
s_iter(6 , 4)
s_iter(10, 5)
s_iter(15, 6)
15
```

Cette fois-ci, il n'y a plus de phase d'expansion suivie d'une phase de réduction. À chaque étape, il suffit de garder une trace des états de `acc` et `compt`. On parle alors de processus *itératif*.

En général, un processus itératif est un processus pouvant être décrit par un nombre fixe de *variables d'état* et par une règle fixe qui décrit l'évolution de chaque variable à chaque étape du processus ainsi qu'un éventuel test d'arrêt qui indique sous quelles conditions le processus doit s'arrêter.

Ici, le nombre d'étapes nécessaire est proportionnel à la taille de l'argument  $n$ . On parle de processus itératif linéaire.

En fait, une autre manière de distinguer un processus itératif d'un processus récursif est de remarquer que si le processus itératif est interrompu, on peut le relancer en donnant comme arguments les valeurs des variables au moment de l'interruption. Ce n'est pas possible pour un processus récursif car l'interpréteur garde des informations « cachées » dans une pile non accessibles par les variables du programme.

Attention ! Ici, les deux procédures sont syntaxiquement récursives mais les processus sont de nature différentes. Cependant, dans des langages comme Pascal, Ada, Python, un processus itératif décrit par une procédure récursive consommera malgré tout autant de mémoire qu'un processus récursif. C'est pourquoi ces langages utilisent des *boucles* (for, while).

On peut utiliser un pliage :

---

```
1 s3 :: Integer -> Integer
2 s3 n = foldl (+) 0 [1..n]
```

---

...ou directement la fonction `sum` :

---

```
1 sum ([1..100000] :: [Integer])
```

---

Faisons une petite expérience, un « benchmark » grâce au module `Criterion` :

---

```
1 import Criterion.Main
2
3 main :: IO()
4 main = defaultMain [
5     bgroup "somme" [ bench "récursif" $ whnf s1 100000
6                     , bench "itératif" $ whnf s2 100000
7                     , bench "pliage" $ whnf s3 100000
8                     , bench "sum" $ whnf sum ([1..100000] :: [Integer])
9                     ]
10 ]
```

---

et on obtient :

---

```
1 $ ghc -O --make Boucles
2 [1 of 1] Compiling Main          ( Boucles.hs, Boucles.o )
3 Linking Boucles ...
4 $ ./Boucles --output boucles.html
5 benchmarking somme/récursif
6 time                6.129 ms   (5.929 ms .. 6.386 ms)
7                    0.992 R2   (0.986 R2 .. 0.999 R2)
8 mean                5.951 ms   (5.885 ms .. 6.069 ms)
9 std dev             256.4 μs    (164.0 μs .. 372.7 μs)
10 variance introduced by outliers: 22% (moderately inflated)
11
12 benchmarking somme/itératif
13 time                1.439 ms   (1.426 ms .. 1.460 ms)
14                    0.999 R2   (0.998 R2 .. 0.999 R2)
15 mean                1.492 ms   (1.480 ms .. 1.507 ms)
16 std dev             46.41 μs    (42.22 μs .. 52.99 μs)
17 variance introduced by outliers: 18% (moderately inflated)
18
19 benchmarking somme/pliage
20 time                2.194 ms   (2.181 ms .. 2.213 ms)
21                    0.999 R2   (0.998 R2 .. 1.000 R2)
22 mean                2.201 ms   (2.189 ms .. 2.220 ms)
23 std dev             51.77 μs    (34.80 μs .. 69.81 μs)
24 variance introduced by outliers: 10% (moderately inflated)
```

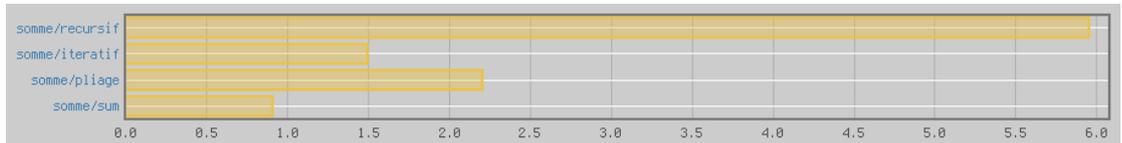
```

25
26 benchmarking somme/sum
27 time          908.1 μs   (902.8 μs .. 916.9 μs)
                1.000 R2   (0.999 R2 .. 1.000 R2)
28
29 mean          904.8 μs   (902.7 μs .. 909.4 μs)
30 std dev       10.03 μs   (4.870 μs .. 16.51 μs)

```

L'option `output` permet d'avoir une sortie html avec des graphiques.

Voici un extrait :



#### 1 4 Itération non linéaire : Fibonacci

Tout le monde connaît l'exemple classique de la suite des nombres de FIBONACCI dans laquelle chaque nombre est la somme des deux précédents :

0, 1, 1, 2, 3, 5, 8, 13, 21, 34, ...

On pense alors naturellement à :

```

1 fib1 :: Integer -> Integer
2 fib1 n
3   | n <= 1  = 1
4   | n > 1   = fib1 (n - 1) + fib1 (n - 2)
5   | otherwise = error "Indice : entier naturel"

```

Mais...

```

1 $ ./FiboBench --output fiboBench.html
2 benchmarking Fibonacci/fib 5
3 time          266.1 ns   (263.9 ns .. 268.7 ns)
4               0.999 R2   (0.999 R2 .. 1.000 R2)
5 mean          264.7 ns   (263.4 ns .. 266.5 ns)
6 std dev       4.963 ns   (3.266 ns .. 6.869 ns)
7 variance introduced by outliers: 23% (moderately inflated)
8
9 benchmarking Fibonacci/fib 10
10 time          3.200 μs   (3.163 μs .. 3.250 μs)
11               0.999 R2   (0.999 R2 .. 0.999 R2)
12 mean          3.192 μs   (3.162 μs .. 3.225 μs)
13 std dev       102.9 ns   (88.03 ns .. 132.9 ns)
14 variance introduced by outliers: 42% (moderately inflated)
15
16 benchmarking Fibonacci/fib 20
17 time          399.6 μs   (393.1 μs .. 405.9 μs)
18               0.998 R2   (0.998 R2 .. 0.999 R2)
19 mean          402.5 μs   (398.9 μs .. 406.1 μs)
20 std dev       12.12 μs   (10.49 μs .. 14.21 μs)
21 variance introduced by outliers: 23% (moderately inflated)

```

Qu'en pensez-vous ?

On peut avoir une autre idée en utilisant deux variables :

```

suivant ← suivant + courant
courant ← suivant

```

Construisez un processus itératif donnant plus efficacement le  $n$ -eme nombre de FIBONACCI (avec et sans boucle).

On pourra de plus gagner du temps en remplaçant le type d'entrée par `Int` : pourquoi ?

```

1 $ ./FiboBench --output fiboBench.html
2 benchmarking Fibonacci/fib rec 35
3 time                211.8 ms   (206.9 ms .. 214.6 ms)
4                    1.000 R²   (0.999 R² .. 1.000 R²)
5 mean                214.4 ms   (212.6 ms .. 217.0 ms)
6 std dev             2.726 ms   (1.363 ms .. 3.603 ms)
7 variance introduced by outliers: 14% (moderately inflated)
8
9 benchmarking Fibonacci/fib it 35
10 time                256.8 ns   (254.6 ns .. 259.8 ns)
11                    0.999 R²   (0.999 R² .. 1.000 R²)
12 mean                256.6 ns   (255.1 ns .. 258.9 ns)
13 std dev             6.053 ns   (4.368 ns .. 7.957 ns)
14 variance introduced by outliers: 33% (moderately inflated)

```

Commentaires ?

On peut encore être plus efficace en utilisant la *mémoïsation* :

```

1 fib3 :: Int -> Integer
2 fib3 =
3     let fib n
4         | n <= 1     = 1
5         | n > 1     = fib3 (n - 1) + fib3 (n - 2)
6         | otherwise = error "Indice : entier naturel"
7     in (map fib [0 ..] !!)

```

Comment ça marche ?

```

1 benchmarking Fibonacci/fib memo 35
2 time                107.0 ns   (106.3 ns .. 107.8 ns)
3                    1.000 R²   (0.999 R² .. 1.000 R²)
4 mean                106.7 ns   (106.3 ns .. 107.4 ns)
5 std dev             1.814 ns   (1.096 ns .. 2.617 ns)
6 variance introduced by outliers: 21% (moderately inflated)

```

## 2 Exemple : l'algorithme de Babylone



Tablette YBC 7289  
XVII<sup>e</sup> avant JC

Il y a presque 4000 ans, les petits Babyloniens calculaient la racine carrée de 2 avec 6 bonnes décimales en utilisant un algorithme vu précédemment (cf section 2.2.1 page 40).

Il s'agit de reprendre ce qui a été vu précédemment avec par exemple :

$$r_0 = 1.0 \text{ et } f : x \mapsto (x + 2/x) * 0.5$$

C'est-à-dire :

```

1 *Main> foldl (\ r k -> (r + 2 / r) * 0.5) 1 [1 .. 6]
2 1.414213562373095

```

ou si vous préférez :

```

1  Algorithme Babylone
2  r ← r0
3  Pour k de 1 à n Faire
4  |   r ← (r+a/r)*0.5
5  FinPour
6  Retourner r

```

**Premier problème** : est-ce que la suite des valeurs calculées par la boucle converge vers  $\sqrt{2}$  ? Nous aurons besoin de deux théorèmes que nous admettrons et dont nous ne donnerons que des versions simplifiées.

**Théorème 3 - 1**

**Théorème de la limite monotone**

Toute suite croissante majorée (ou décroissante minorée) converge

Un théorème fondamental est celui du point fixe. Il faudrait plutôt parler des théorèmes du point fixe car il en existe de très nombreux avatars qui portent les noms de mathématiciens renommés : BANACH, BOREL, BROUWER, KAKUTANI, KLEENE, ... Celui qui nous intéresserait le plus en informatique est celui de KNASTER-TARSKI. On peut en trouver une présentation claire dans ?.

Ils donnent tous des conditions d'existence de points fixes (des solutions de  $f(x) = x$  pour une certaine fonction).

Nous nous contenterons pour l'instant de la version *light* vue au lycée.

**Théorème 3 - 2**

**Théorème du point fixe : version (très) édulcorée du lycée**

Soit  $I$  un intervalle fermé de  $\mathbb{R}$ , soit  $f$  une fonction de  $I$  vers  $I$  et soit  $(r_n)$  une suite d'éléments de  $I$  telle que  $r_{n+1} = f(r_n)$  pour tout entier naturel  $n$ .

SI  $(r_n)$  est convergente ALORS sa limite est UN *point fixe* de  $f$  appartenant à  $I$ .

Cela va nous aider à étudier notre suite définie par  $r_{n+1} = f(r_n) = \frac{r_n + \frac{2}{r_n}}{2}$  et  $r_0 = 1$ .

**Recherche**

Démontrez que pour tout entier naturel non nul  $n$ , on a  $r_n \geq \sqrt{2}$  puis que la suite est décroissante.

Démontrez que  $\sqrt{2}$  est l'unique point fixe positif de  $f$  et conclure.

**Deuxième problème** : quelle est la vitesse de convergence ? Ici, cela se traduit par « combien de décimales obtient-t-on en plus à chaque itération ? ».

**Recherche**

Calculez  $(r_{n+1} - \sqrt{2})$  en fonction de  $(r_n - \sqrt{2})$  et essayez de répondre à la question précédente à l'aide de ce résultat.

On peut donc être sûr a priori que le résultat de notre boucle donnera le résultat spécifié (avec la convergence). Le nombre d'itérations pour une précision donnée peut lui aussi être calculé à l'avance en étudiant la vitesse.

Essayons d'être un peu plus rigoureux et général en introduisant la notion d'ordre d'une suite :

**Définition 3 - 1**

**Ordre d'une suite - Constante asymptotique d'erreur**

Soit  $(r_n)$  une suite convergeant vers  $\ell$ . S'il existe un entier  $k > 0$  tel que :

$$\lim_{n \rightarrow +\infty} \frac{|r_{n+1} - \ell|}{|r_n - \ell|^k} = C$$

avec  $C \neq 0$  et  $C \neq +\infty$  alors on dit que  $(r_n)$  est d'ordre  $k$  et que  $C$  est la constante asymptotique d'erreur.

**Recherche**

Déterminez par exemple l'ordre et la constante asymptotique d'erreur de la suite de Babylone donnant une approximation de  $\sqrt{2}$ .

Si l'on connaît l'ordre d'une suite, on va pouvoir prévoir l'évolution de la précision atteinte au fur et à mesure de l'avance de la boucle.

Ainsi, pour une suite d'ordre  $k$  avec comme constante asymptotique d'erreur  $C$ , on a pour  $n$  suffisamment grand :

$$|r_{n+1} - \ell| \approx C|r_n - \ell|^k$$

## Recherche

Quel est le rapport entre  $\log_{10} x$  et le nombre de décimales d'un nombre plus petit que 1 en base 10 ?

En posant  $d_n = -\log_{10} |r_n - \ell|$  et  $D = -\log_{10} C$  on obtient :

$$d_{n+1} = \dots$$

## Recherche

Que peut-on en conclure ?

Une application : on peut ainsi calculer rapidement sur un petit système embarqué (voire un *smart phone*) la vitesse de la voiture en fonction de la distance la séparant de la voiture qui la précède.



Application *safety sight* sur *smart phone*

## 3

## Méthodes itératives pour résoudre une équation

Résoudre une équation du premier ou du second degré, vous savez faire...Mais dans les autres cas ?

Dans cette section, nous nous intéresserons à la résolution d'une équation du type  $\varphi(x) = \ell$  où  $\varphi$  est une fonction de  $\mathbb{R}$  dans lui-même.

On se ramènera même à une équation de la forme  $f(x) = 0$  en posant  $f(x) = \varphi(x) - \ell$ .

### 3 1 Dichotomie (ou méthode de bi-section)

On veut résoudre une équation du type  $f(x) = 0$  sur un intervalle  $I$  avec  $f$  une fonction continue qui change de signe sur un intervalle  $[a, b]$  inclus dans  $I$ .

On sait qu'il existe au moins une solution  $\alpha$  de notre équation sur  $[a, b]$  : pourquoi ?

Soit  $m$  le centre de l'intervalle. Trois cas se rencontrent :

- si  $f(m) = 0$  alors la vie est belle ;
- si  $f(m)$  a le même signe que  $f(a)$  alors  $\alpha \in [m, b]$  (pourquoi ?) ;
- si  $f(m)$  a le même signe que  $f(b)$  alors  $\alpha \in [a, m]$  (pourquoi ?) ;

## Recherche

On itère ce mécanisme jusqu'à obtenir un intervalle d'amplitude correspondant à la précision demandée sur  $\alpha$  : pourquoi cela suffit-il ? Quel est l'avantage de cette méthode informatiquement parlant ? Quel est son principal désavantage ?

Il est à noter que les divisions par successives peuvent engendrer des problèmes selon la norme IEEE 754 (cf ?).

On parle de *dichotomie* du grec *διχοτομία* qui signifie *division en deux parties égales*.

**3 2** **Ordre d'une suite et formules de Taylor**

Considérons une suite convergente générée par un procédé itératif  $x_{n+1} = \varphi(x_n)$ . Supposons que  $\varphi$  est de classe  $C^1$  au voisinage de sa limite  $\ell$ . Si l'on remarque que  $x_{n+1} = \varphi(x_n) = \varphi(\ell + (x_n - \ell))$  alors :

$$x_{n+1} = \varphi(\ell) + (x_n - \ell)\varphi'(\ell) + O(|x_n - \ell|^2)$$

On en déduit que :

$$x_{n+1} - \ell = (x_n - \ell)\varphi'(\ell) + O(|x_n - \ell|^2)$$

Recherche

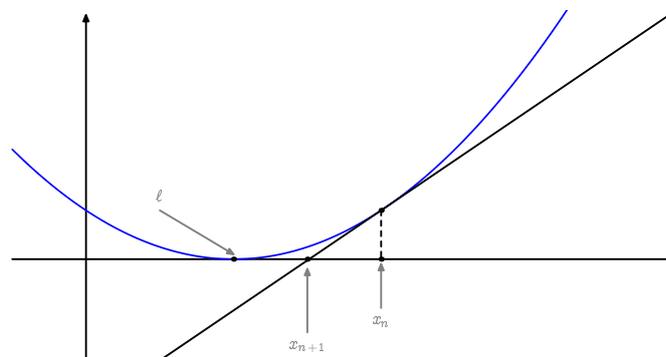
On peut alors en conclure des choses sur l'ordre de la suite selon que  $\varphi'(\ell)$  est nul ou pas : lesquelles ?

**3 3** **Méthode de Newton-Raphson-Halley-etc.**

La méthode de résolution des équations numériques que nous allons voir à présent a été initiée par Isaac NEWTON vers 1669 sur des exemples numériques mais la formulation était fastidieuse. Dix ans plus tard, Joseph RAPHSON met en évidence une formule de récurrence. En 1694, la formule est encore améliorée par Edmund HALLEY, l'homme de la comète. Un siècle plus tard, MOURAILLE et LAGRANGE étudient la convergence des approximations successives en fonction des conditions initiales par une approche géométrique. Cinquante ans plus tard, FOURIER et CAUCHY s'occupe de la rapidité de la convergence.

**3 3 1** **Approche intuitive graphique**

- On part d'un nombre quelconque  $x_0$  ;
- à partir de  $x_0$ , on calcule un nouveau nombre  $x_1$  de la manière suivante (voir figure) : on trace la tangente au graphe de  $f$  au point d'abscisse  $x_0$ , et on détermine le point d'intersection de cette tangente avec l'axe des abscisses. On appelle  $x_1$  l'abscisse de ce point d'intersection ;
- et on recommence : on calcule un nouveau nombre  $x_2$  en appliquant le procédé décrit au point 2 où l'on remplace  $x_0$  par  $x_1$  ;
- etc.



À partir de cette description graphique de la méthode de Newton, trouver la formule, notée (1), donnant  $x_1$  en fonction de  $x_0$ , puis  $x_{n+1}$  en fonction de  $x_n$ .

Quelles hypothèses doit-on faire sur  $f$  et les  $x_n$  pour que la formule ait un sens ?

**3 3 2** **Approche intuitive « taylorienne »**

Nous voulons résoudre l'équation  $f(x) = 0$ . Utilisons à nouveau la ruse :

$$f(x_{n+1}) = f(x_n + (x_{n+1} - x_n)) = f(x_n) + (x_{n+1} - x_n)f'(x_n) + O((x_{n+1} - x_n)^2)$$

Intuitivement, pour que  $f(x_{n+1})$  soit de plus en plus petit, il serait pratique que la partie linéaire de cette expression,  $f(x_n) + (x_{n+1} - x_n)f'(x_n)$ , soit la plus petite possible, c'est-à-dire nulle et on obtient à nouveau :

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

On remarque alors que si la suite converge vers une limite  $\ell$  et si  $f'(\ell)$  est non nul, alors on obtient bien  $f(\ell) = 0$ .

Maintenant, posons  $\varphi(x) = x - \frac{f(x)}{f'(x)}$  et supposons que  $f$  soit de classe  $C^2$  sur un intervalle ouvert  $V$  tel que l'équation  $f(x) = 0$  admette une unique solution d'ordre 1 (c'est-à-dire que  $f(\ell) = 0$  mais  $f'(\ell) \neq 0$ ).

La méthode de NEWTON revient à construire la suite définie par  $x_{n+1} = \varphi(x_n)$  avec  $x_0 \in V$ .

Recherche

Vérifiez que  $\varphi$  est dérivable sur  $V$  et que

$$\varphi'(x) = \frac{f(x)f''(x)}{(f'(x))^2}$$

et concluez en utilisant les résultats de la section 3.3.2 page précédente

À retenir

Il faut bien comprendre que nos approches sont un peu ole ole compte-tenu de nos outils mathématiques limités...

On retiendra que  $x_0$  a besoin d'être suffisamment proche de la solution pour que l'équation  $f(x) = 0$  admette une unique solution et que  $f'(\ell)$  soit non nul. En fait, on pourrait montrer que la méthode marche encore si  $\ell$  est une solution d'ordre 2.

Il faut aussi que la fonction  $f$  soit suffisamment régulière. Bref, il ne faut pas trop bricoler au hasard...

### 3 4 Étude de la suite associée à l'équation $x^3 - 2x - 5 = 0$ .

On veut résoudre l'équation  $x^3 - 2x - 5 = 0$  par la méthode de NEWTON-RAPHSON. On note  $f$  la fonction  $x \mapsto x^3 - 2x - 5$ . C'est en fait sur cet exemple que NEWTON a raisonné.

Recherche

1. Montrez rapidement que l'équation  $f(x) = 0$  admet une unique solution  $\alpha$  sur  $\mathbb{R}$ . Montrez que  $2 < \alpha < 3$ .
2. Déterminez la fonction  $\varphi$  telle que  $x_{n+1} = \varphi(x_n)$ , la suite  $(x_n)$  étant celle décrite au paragraphe précédent en prenant  $x_0 = 3$ .
3. Étudiez le sens de variation de la fonction  $\varphi$  puis celui de  $\varphi'$  et déduisez-en que  $[\alpha, 3]$  est stable par  $\varphi$  et que  $\varphi$  est strictement croissante sur  $[\alpha, 3]$ .
4. Que pouvez-vous en déduire sur la convergence de la suite  $(x_n)$  ?

### 3 5 Test d'arrêt

Afin de construire un algorithme donnant une approximation d'une solution d'une équation numérique par la méthode de NEWTON-RAPHSON, il faudrait déterminer un test d'arrêt c'est-à-dire savoir à partir de quel rang  $n$   $|x_n - \alpha|$  restera inférieur à une valeur donnée.

Il suffit de remarquer que  $f(x_n) = f(x_n) - f(\alpha)$ . Nous supposons  $f$  de classe  $C^2$  sur un « bon » voisinage  $I$  de  $\alpha$  (ce critère nous échappe encore à notre niveau). Alors  $f$  est en particulier dérivable en  $\alpha$  donc

$$f(x_n) = f(x_n) - f(\alpha) \sim f'(\alpha) \times (x_n - \alpha)$$

C'est-à-dire, puisque  $f'(\alpha)$  est supposé non nul :

$$x_n - \alpha \sim \frac{f(x_n)}{f'(\alpha)}$$

Or  $f$  étant de classe  $C^2$ , on a  $f'$  continue et non nulle en  $\alpha$  donc  $f'(\alpha) \sim f'(x_n)$ .

Finalement

$$x_n - \alpha \sim \frac{f(x_n)}{f'(x_n)} = x_n - x_{n+1}$$

Nous choisirons donc comme test d'arrêt  $\left| \frac{f(x_n)}{f'(x_n)} \right| < p$  avec  $p$  la précision choisie. Il ne reste plus qu'à écrire l'algorithme.

## 4

## Étude expérimentale et théorique de la complexité d'un algorithme

### 4.1 Méthode « scientifique »

Le problème : on dispose d'une liste de  $N$  nombres. Déterminez les triplets dont la somme est nulle.

Utilisons la force brute

```

1 import Criterion.Main
2
3 troisSommes :: [Int] -> [(Int,Int,Int)]
4 troisSommes xs =
5     let inds = [0 .. (length xs -1)]
6         in [(xs!!i , xs!!j , xs!!k) | i <- inds, j <- drop i inds, k <- drop j inds, xs!!i +
           ↪ xs!!j + xs!!k == 0]
7
8 main :: IO()
9 main = defaultMain [
10     bgroup "TroisSommes" [ bench "3som 100" $ whnf troisSommes [-50..49]
11                           , bench "3som 200" $ whnf troisSommes [-100..99]
12                           , bench "3som 400" $ whnf troisSommes [-200..199]
13                           , bench "3som 800" $ whnf troisSommes [-400..399]
14                           , bench "3som 1600" $ whnf troisSommes [-800..799]
15                           , bench "3som 3200" $ whnf troisSommes [-1600..1599]
16                       ]
17 ]

```

et on observe :

```

1 $ ./TroisSom --output troisSom.html
2 benchmarking TroisSommes/3som 100
3 time           640.9 µs   (636.3 µs .. 645.0 µs)
4               1.000 R²   (0.999 R² .. 1.000 R²)
5 mean          640.4 µs   (637.6 µs .. 644.4 µs)
6 std dev       11.88 µs   (9.047 µs .. 15.88 µs)
7
8 benchmarking TroisSommes/3som 200
9 time           4.322 ms   (4.305 ms .. 4.344 ms)
10              1.000 R²   (0.999 R² .. 1.000 R²)
11 mean          4.358 ms   (4.342 ms .. 4.384 ms)
12 std dev       65.63 µs   (42.33 µs .. 119.0 µs)
13
14 benchmarking TroisSommes/3som 400
15 time           31.72 ms   (31.43 ms .. 32.08 ms)
16              1.000 R²   (0.999 R² .. 1.000 R²)
17 mean          31.58 ms   (31.50 ms .. 31.75 ms)
18 std dev      239.8 µs    (83.75 µs .. 368.0 µs)
19
20 benchmarking TroisSommes/3som 800
21 time           261.7 ms   (259.7 ms .. 263.6 ms)
22              1.000 R²   (1.000 R² .. 1.000 R²)
23 mean          262.8 ms   (262.0 ms .. 263.2 ms)
24 std dev      758.8 µs    (95.48 µs .. 967.9 µs)

```

```

25 variance introduced by outliers: 16% (moderately inflated)
26
27 benchmarking TroisSommes/3som 1600
28 time                2.318 s      (2.289 s .. 2.339 s)
29                    1.000 R2    (NaN R2 .. 1.000 R2)
30 mean                2.320 s      (2.314 s .. 2.323 s)
31 std dev             4.803 ms     (0.0 s .. 4.834 ms)
32 variance introduced by outliers: 19% (moderately inflated)
33
34 benchmarking TroisSommes/3som 3200
35 time                16.47 s     (15.00 s .. 17.51 s)
36                    0.999 R2    (0.998 R2 .. 1.000 R2)
37 mean                17.33 s     (16.82 s .. 17.64 s)
38 std dev             471.1 ms    (0.0 s .. 535.4 ms)
39 variance introduced by outliers: 19% (moderately inflated)

```

Voyons sur Python :

```

1 def trois_sommes(xs):
2     N = len(xs)
3     cpt = 0
4     for i in range(N):
5         for j in range(i + 1, N):
6             for k in range(j + 1, N):
7                 if xs[i] + xs[j] + xs[k] == 0:
8                     cpt += 1
9     return cpt

```

Comparons les temps pour différentes tailles :

```

1 In [23]: %timeit trois_sommes([randint(-10000,10000) for k in range(100)])
2 10 loops, best of 3: 27.5 ms per loop
3
4 In [24]: %timeit trois_sommes([randint(-10000,10000) for k in range(200)])
5 1 loops, best of 3: 216 ms per loop
6
7 In [25]: %timeit trois_sommes([randint(-10000,10000) for k in range(400)])
8 1 loops, best of 3: 1.82 s per loop

```

Qu'en pensez-vous ?

Il ne semble donc pas aberrant de considérer que le temps de calcul est de  $aN^3$  mais que vaut  $a$  ?

```

1 In [62]: temps(range(400))
2 Out[62]: 4.005484320001415

```

$4,00 = a \times 400^3$  donc  $a \approx 6,25 \times 10^{-8}$

Donc pour  $N = 1000$ , on devrait avoir un temps de  $6,25 \times 10^{-8} \times 10^9 = 62,5$

```

1 In [66]: temps(range(1000))
2 Out[66]: 68.54615448799996

```

Voici la même chose en C, sans vraiment chercher à optimiser le code.

```

1 #include <stdio.h>
2 #include <time.h>
3
4 typedef int Liste[13000];

```

```

5
6 int trois_sommes(int N)
7 {
8     int cpt = 0;
9     Liste liste;
10
11     for ( int k = 0; k < N; k++)
12     {
13         liste[k] = k - (N / 2);
14     }
15
16     for (int i = 0; i < N; i++)
17     {
18         for (int j = i + 1; j < N; j++)
19         {
20             for (int k = j + 1; k < N; k++)
21             {
22                 if (liste[i] + liste[j] + liste[k] == 0) {cpt++;}
23             }
24         }
25     }
26     return cpt;
27 }
28
29 void chrono(int N)
30 {
31     clock_t temps_initial, temps_final;
32     float temps_cpu;
33
34     temps_initial = clock();
35     trois_sommes(N);
36     temps_final = clock();
37     temps_cpu = ((double) temps_final - temps_initial) / CLOCKS_PER_SEC;
38     printf("Temps en sec pour %d : %f\n",N, temps_cpu);
39
40 }
41
42 int main(void)
43 {
44     chrono(100);
45     chrono(200);
46     chrono(400);
47     chrono(800);
48     chrono(1600);
49     chrono(3200);
50     chrono(6400);
51     chrono(12800);
52
53     return 1;
54 }

```

Et on obtient :

```

1 $ gcc -std=c99 -Wall -Wextra -Werror -pedantic -O4 -o somm3 Trois_Sommes.c
2 $ ./somm3
3 Temps en sec pour 100 : 0.000000
4 Temps en sec pour 200 : 0.000000
5 Temps en sec pour 400 : 0.020000
6 Temps en sec pour 800 : 0.090000
7 Temps en sec pour 1600 : 0.720000
8 Temps en sec pour 3200 : 5.760000
9 Temps en sec pour 6400 : 45.619999
10 Temps en sec pour 12800 : 360.839996

```

Que pensez-vous de tout ceci ? Points communs ? Différences ? Pourquoi ne pas se soucier des constantes multiplicatives ?

## 4 2 Complexité des boucles

Une bonne lecture de chevet est le troisième volume de « Zi Arte » (?).

Le temps d'exécution est la somme des produits (coût × fréquence) pour chaque opération.

- le coût dépend de la machine, du compilateur, du langage ;
- la fréquence dépend de l'algorithme, de la donnée en entrée.

Dans notre algorithme des 3 sommes en python, on peut donc gagner du temps :



D. E. Knuth  
Né en 1938

```

1 def trois_sommes(xs):
2     N = len(xs)
3     cpt = 0
4     for i in range(N):
5         xi = xs[i]
6         for j in range(i + 1, N):
7             sij = xi + xs[j]
8             for k in range(j + 1, N):
9                 if sij + xs[k] == 0:
10                    cpt += 1
11    return cpt

```

C'est mieux :

```

1 In [28]: xs = list(range(-50,51))
2 In [29]: %timeit trois_sommes(xs)
3 100 loops, best of 3: 12.9 ms per loop
4
5 In [30]: xs = list(range(-100,101))
6 In [31]: %timeit trois_sommes(xs)
7 10 loops, best of 3: 94.4 ms per loop
8
9 In [32]: xs = list(range(-200,201))
10 In [33]: %timeit trois_sommes(xs)
11 1 loops, best of 3: 851 ms per loop
12
13 In [34]: xs = list(range(-400,401))
14 In [35]: %timeit trois_sommes(xs)
15 1 loops, best of 3: 7.04 s per loop

```

En C, ça ne changera rien, car le compilateur est intelligent et a remarqué tout seul qu'il pouvait garder en mémoire certains résultats.

On peut visualiser en échelle log-log :

```

1 tailles = [2**k for k in range(7,11)]*
2 listes = [list(range(- N//2, N//2 + 1)) for N in tailles]
3 t = [temps(xs) for xs in listes]
4 for taille in tailles:
5     plt.loglog(tailles, t, basex = 2, basey = 2)
6 plt.title('Temps de calcul selon la taille de la liste en échelle log-log')
7 plt.xlabel('Taille des listes')
8 plt.ylabel('Temps en s')

```

C'est assez rectiligne.

Regardons de nouveau le code des trois sommes et comptons le nombre d'opérations :

OPÉRATION	FRÉQUENCE
Déclaration de la fonction et du paramètre (l. 1)	2
Déclaration de N, cpt et i (l. 2, 3 et 4)	3
Affectation de N, cpt et i (l. 2, 3 et 4)	3
Déclaration de xi (l. 5)	N
Affectation de xi (l. 5)	N
Accès à xs[i] (l. 5)	N
Déclaration de j (l.6)	N
Calcul de l'incrément de i (l. 6)	N
Affectation de j (l.6)	N
Déclaration de sij (l. 7)	$S_1$
Affectation de sij (l. 7)	$S_1$
Accès à xs[j] (l.7)	$S_1$
Somme (l.7)	$S_1$
Déclaration de k (l.8)	$S_1$
Incrément de j (l. 8)	$S_1$
Affectation de k (l.8)	$S_1$
Accès à x[k] (l.9)	$S_2$
Calcul de la somme (l.9)	$S_2$
Comparaison à 0 (l.9)	$S_2$
Incrément de cpt (l.9)	entre 0 et $S_2$
Affectation de la valeur de retour (l.11)	1

Que valent  $S_1$  et  $S_2$  ?

$$S_1 = \sum_{i=0}^{N-1} N - (i + 1) = \sum_{i'=0}^{N-1} i' = \frac{N(N-1)}{2}$$

$$\begin{aligned} S_2 &= \sum_{i=0}^{N-1} \sum_{j=i+1}^{N-1} N - (j + 1) \\ &= \sum_{i=0}^{N-1} \sum_{j'=0}^{N-(i+2)} j' \\ &= \sum_{i=0}^{N-2} \frac{(N - (i + 2))(N - (i + 1))}{2} \\ &= \sum_{i'=1}^{N-2} \frac{i'(i' + 1)}{2} \\ &= \frac{1}{2} \left( \sum_{i=1}^{N-2} i^2 + \sum_{i=1}^{N-2} i \right) \\ &= \frac{1}{2} \left( \frac{(N-2)(2N-3)(N-1)}{6} + \frac{(N-2)(N-1)}{2} \right) \\ &= \frac{N(N-1)(N-2)}{6} \end{aligned}$$

Notons  $a$  le temps constant d'affectation,  $d$  le temps constant de déclaration,  $x$  le temps constant d'accès à une cellule,  $s$  le temps constant d'une somme,  $c$  le temps constant d'une comparaison. Le temps d'exécution vérifie :

$$\tau(N) \leq (2d + 3d + 3a + a) + (d + a + x + d + s + a)N + (d + a + x + s + d + s + a)S_1 + (x + s + c + s)S_2$$

Or  $S_1 \sim N^2$  et  $S_2 \sim N^3$  quand  $N$  est « grand ».

Finalement...

$$\tau(N) = O(N^3)$$

...et ce, que l'on utilise C, Python, BrainFuck (<https://fr.wikipedia.org/wiki/Brainfuck>), etc. C'est aussi ce que l'on a observé expérimentalement

## 5 Sommes



« *What is one and one ?* »  
 « *I don't know* » said Alice. « *I lost count* ».  
 « *She can't do Addition.* »

Alice in wonderland - Lewis CAROLL

### 5 1 Notation

Si on écrit  $1 + 2 + 3 + 4 + \dots + (n - 1) + n$ , on comprend que les  $\dots$  signifient qu'ils doivent être remplacés selon le motif induit par les termes les entourant.

On peut être amené à étudier des sommes plus générales :  $a_0 + a_1 + \dots + a_n$ , chaque terme  $a_k$  de la somme étant défini d'une manière quelconque.

Si l'on écrit la somme avec des  $\dots$ , il faut que le motif de construction soit présenté de manière suffisamment claire pour être induit :

$$1 + 4 + 9 + 22 + 53 + 128 + 309 + \dots$$

n'est pas évident à comprendre mais :

$$1 + 4 + (2 \times 4 + 1) + (2 \times 9 + 4) + \dots + (2a_{n-1} + a_{n-2})$$

est plus explicite.

En 1820, le mathématicien Joseph FOURIER<sup>a</sup> introduit le fameux sigma majuscule. Citons-le :

Le signe  $\sum_{i=1}^{+\infty}$  indique que l'on doit donner au nombre entier  $i$  toutes ses valeurs  $1, 2, 3, \dots$ , et prendre la somme des termes.

Notez bien, en tant qu'informaticien(ne), que le  $i$  joue ici le rôle d'une variable locale et aurait très bien pu porter un autre nom sans changer le résultat.

Depuis, on a généralisé l'emploi des bornes afin de définir les indices à partir de certaines propriétés. On peut par exemple écrire les sommes suivantes avec ou sans bornes explicitement spécifiées :



Joseph FOURIER  
(1768-1830)

$$\sum_{\substack{1 \leq k < 100 \\ k \text{ impair}}} k^2 = \sum_{k=0}^{49} (2k+1)^2$$

Cette formulation est particulièrement adaptée aux changements d'indices :

$$\sum_{1 \leq k \leq n} a_k = \sum_{1 \leq k+1 \leq n} a_{k+1} = \sum_{k=1}^n a_k = \sum_{i=0}^{n-1} a_{i+1}$$

La dernière forme demande un peu de réflexion alors que la deuxième beaucoup moins.

L'informaticien canadien Kenneth IVERSON, prix TURING en 1979 et inventeur des langages APL et J, introduit une notation fort pratique mais peu utilisée en France, nommée **crochets d'Iverson**.

Soit  $P$  une certaine propriété :

$$[P] = \begin{cases} 1 & \text{si } P \text{ est vraie} \\ 0 & \text{sinon} \end{cases}$$

<sup>a</sup>. Nous aurions pu ne pas étudier les sommes, ne pas compresser les images et le son car FOURIER échappa de peu à la guillotine durant la Terreur : <http://www.youtube.com/watch?v=MwGD13BnH-o>

Par exemple :

$$\sum_{\substack{1 \leq k < 100 \\ k \text{ impair}}} k^2 = \sum_k k^2 [1 \leq k < 100] \times [k \text{ impair}]$$

## 5 2 Somme, boucle et récurrence

Une somme suggère naturellement l'emploi d'une boucle « pour ». En effet  $\sum_{k=1}^n a_k$  se lit « somme des  $a_k$  pour  $k$  variant de 1 à  $n$  » ce qui, dans un langage quelconque, par exemple... OCAML donne, pour la somme des  $n$  premiers carrés d'entiers :

```
# let somme(n)=
  let s = ref 0 in
  for k=0 to n do
    s := !s + k*k
  done;
  !s;;
  val somme : int -> int = <fun>
# somme(3);;
- : int = 14
```

On peut également penser à une somme en terme de suite. Si on note  $S_n = \sum_{k=0}^n k^2$ , alors :

$$S_0 = 0; \quad \forall n > 0, S_n = S_{n-1} + n^2$$

Dans notre langage préféré :

```
# let rec s = fonction
  | 0 -> 0
  | n -> s(n-1) + n*n;;
  val s : int -> int = <fun>
# s(3);;
- : int = 14
```

Cela nous rappelle d'ailleurs qu'une suite est une fonction de  $\mathbb{N}$  dans tout autre ensemble (ici  $\mathbb{N}$  lui-même).

## 5 3 Manipulation de sommes

La principale difficulté du calcul avec des sommes est d'obtenir par diverses manipulation une expression plus simple. Tout dérive le plus souvent des trois lois suivantes :

**Distributivité**  $\sum_{k \in K} \lambda a_k = \lambda \sum_{k \in K} a_k$  ( $\lambda$  ne dépendant pas de  $k$ );

**Associativité**  $\sum_{k \in K} (a_k + b_k) = \sum_{k \in K} a_k + \sum_{k \in K} b_k$ ;

**Commutativité**  $\sum_{k \in K} a_k = \sum_{k \in K} a_{\pi(k)}$  ( $\pi \in S_{\text{Card } K}$ ).

Par exemple, nous allons prouver un résultat fort utile concernant les suites géométriques.

Nous voudrions calculer  $S_n = \sum_{k=0}^n x^k$  pour tout réel  $x \neq 1$ .

Or

$$S_n + x^{n+1} = x^0 + \sum_{1 \leq k \leq n+1} x^k = x^0 + \sum_{0 \leq k-1 \leq n} x^k = x^0 + \sum_{0 \leq k-1 \leq n} x^{(k-1)+1}$$

En posant  $i = k - 1$  on obtient :

$$S_n + x^{n+1} = x^0 + \sum_{0 \leq i \leq n} x^{i+1} = x^0 + x \sum_{0 \leq i \leq n} x^i = 1 + x S_n$$

d'où :

$$S_n(1 - x) = 1 - x^{n+1}$$

Finalement, comme  $x \neq 1$ , on obtient :

$$\sum_{k=0}^n x^k = \frac{1 - x^{n+1}}{1 - x}$$

**5 4 Sommes multiples**

Nous voulons exprimer  $S = a_1b_1 + a_1b_2 + a_1b_3 + a_2b_1 + a_2b_2 + a_2b_3 + a_3b_1 + a_3b_2 + a_3b_3$  avec le symbole  $\sum$  :

$$\begin{aligned}
 S &= \sum_{i,j} a_i b_j [1 \leq i, j \leq 3] \\
 &= \sum_{i,j} a_i b_j [1 \leq i \leq 3][1 \leq j \leq 3] \\
 &= \sum_i \left( \sum_j a_i b_j [1 \leq i \leq 3][1 \leq j \leq 3] \right) \\
 &= \sum_i a_i [1 \leq i \leq 3] \left( \sum_j b_j [1 \leq j \leq 3] \right) \\
 &= \left( \sum_i a_i [1 \leq i \leq 3] \right) \left( \sum_j b_j [1 \leq j \leq 3] \right) \\
 &= \left( \sum_{i=1}^3 a_i \right) \left( \sum_{j=1}^3 b_j \right)
 \end{aligned}$$

Nous voulons cette fois exprimer  $S' = a_1b_1 + a_1b_2 + a_1b_3 + a_2b_2 + a_2b_3 + a_3b_3$  avec le symbole  $\sum$  :

$$\begin{aligned}
 S' &= \sum_{i,j} a_i b_j [1 \leq i \leq j \leq 3] \\
 &= \sum_{i,j} a_i b_j [1 \leq i \leq 3][i \leq j \leq 3] \\
 &= \sum_i \left( \sum_j a_i b_j [1 \leq i \leq 3][i \leq j \leq 3] \right) \\
 &= \sum_i a_i [1 \leq i \leq 3] \left( \sum_j b_j [i \leq j \leq 3] \right) \\
 &= \sum_{i=1}^3 a_i \left( \sum_{j=i}^3 b_j \right)
 \end{aligned}$$

Nous verrons d'autres exemples en exercice.

**6 Sommes infinies (séries)****6 1 Les dangers des ... mal maîtrisés**

À quoi peut bien être égal  $S = \frac{1}{2^0} + \frac{1}{2^1} + \frac{1}{2^2} + \dots$  ?

Multiplions les deux membres par 2, nous obtenons :

$$2S = 2 + \frac{1}{2^0} + \frac{1}{2^1} + \frac{1}{2^2} + \dots = 2 + S$$

donc  $S = 2$ . Waouh, magique !

On recommence avec :

$$T = 2^0 + 2^1 + 2^2 + \dots$$

alors :

$$2T = 2^1 + 2^2 + 2^3 + \dots = T - 2^0 = T - 1 \text{ donc } T = -1.$$

Waouh !...Euh...Quoique, même réveillé brutalement en plein somme il peut sembler bizarre que la somme de nombres strictement positifs donne un nombre strictement négatif...

Il va falloir prendre des précautions.

**6 2 Séries**

Définition 3 - 2

Soit  $u$  une suite réelle définie sur  $I = \{n \in \mathbb{N} \mid n \geq n_0\}$ . On appelle série de terme général  $u_n$  la suite  $S$  définie par  $S_n = \sum_{k=n_0}^n u_k$ . Si la suite  $S$  converge, on dit que la série de terme général  $u_n$  converge et on note  $\sum_{k=n_0}^{+\infty} u_k = \lim_{n \rightarrow +\infty} S_n$ . Une série qui ne converge pas est dite divergente, dans ce cas l'écriture  $\sum_{k=n_0}^{+\infty} u_k$  n'a pas de sens. La série de terme général  $u_n$  est notée  $(\sum u_n)$ .

Si on nous demande d'étudier la série  $(\sum \frac{1}{n^2})$ , on nous demande d'étudier la suite  $S$  définie par  $S_n = \sum_{k=1}^n \frac{1}{k^2}$  avec évidemment  $n \in \mathbb{N}^*$ .

Remarque

Si la suite  $u$  est définie à partir de  $n_0$ , rien ne nous empêche de changer les notations pour travailler avec la même suite définie à partir du rang 0 ou du rang 1 ou ... En effet, considérons la suite  $v$  définie par  $v_n = u_{n+n_0}$ , il est clair que la suite  $v$  est définie sur  $\mathbb{N}$  et, qu'étudier la suite  $v$ , c'est étudier la suite  $u$ . La conséquence de ceci est que, dans ce qui suit, les suites qui interviendront seront, suivant les besoins, définie à partir de  $n = 0$  ou  $n = 1 \dots$

Remarque

$S_n$  est appelé une somme partielle.

Théorème 3 - 3

Si la série  $(\sum u_n)$  converge alors on a  $\lim_{n \rightarrow +\infty} u_n = 0$

La démonstration est très simple.

Considérons  $S_n = \sum_{j=0}^n u_j$ . Par hypothèse  $\lim_{n \rightarrow +\infty} S_n = \lim_{k \rightarrow +\infty} S_k = \ell$ . Si nous remplaçons  $k$  par  $n - 1$ , nous obtenons

$$\lim_{n \rightarrow +\infty} S_n = \lim_{n-1 \rightarrow +\infty} S_{n-1} = \ell$$

mais comme  $S_n - S_{n-1} = u_n$  on en déduit

$$\lim_{n \rightarrow +\infty} u_n = \lim_{n \rightarrow +\infty} (S_n - S_{n-1}) = \ell - \ell = 0$$

Remarque

La réciproque est totalement fautive, ce n'est pas parce que le terme général d'une série a pour limite zéro que cette série converge.

**6 3 Série harmonique**

Nous allons démontrer que la série  $(\sum \frac{1}{n})$ , qui porte le nom de série harmonique, diverge. Nous avons donc  $u_n = \frac{1}{n}$  avec  $n \in \mathbb{N}^*$  évidemment et

$$S_n = \sum_{j=1}^n u_j = 1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n}$$

La suite  $S$  est une suite de termes strictement positifs et comme  $S_{n+1} - S_n = \frac{1}{n+1} > 0$  nous sommes assurés que la suite  $S$  est strictement croissante : ainsi, soit elle diverge vers  $+\infty$ , soit elle converge vers une limite réelle  $\ell$ .

Calculons  $S_{2n} - S_n$ .

$$\begin{aligned} S_n &= \sum_{j=1}^n u_j = 1 + \frac{1}{2} + \dots + \frac{1}{n} \\ S_{2n} &= \sum_{j=1}^{2n} u_j = 1 + \frac{1}{2} + \dots + \frac{1}{n} + \frac{1}{n+1} + \frac{1}{n+2} + \dots + \frac{1}{n+n} \\ S_{2n} - S_n &= \frac{1}{n+1} + \frac{1}{n+2} + \dots + \frac{1}{n+n} \end{aligned}$$

La dernière somme comporte exactement  $n$  termes tous  $\geq \frac{1}{2n}$ . On obtient alors

$$S_{2n} - S_n \geq \frac{1}{2n} + \frac{1}{2n} + \dots + \frac{1}{2n} = \frac{n}{2n} = \frac{1}{2}$$

Supposons que  $(S_n)$  converge vers  $\ell$ , alors  $\lim_{n \rightarrow +\infty} S_n = \lim_{n \rightarrow +\infty} S_{2n} = \ell$ .

En reportant dans l'inégalité obtenue précédemment :  $\ell - \ell \geq \frac{1}{2}$  ce qui est absurde donc  $(S_n)$  diverge vers  $+\infty$ , doucement, mais sûrement.

#### 6 4 Séries télescopiques

Intéressons-nous à la convergence de la série de terme général  $u_n = \frac{1}{n^2 - 1}$  ; il est évident que  $u_n$  n'est définie que pour  $n \geq 2$ . Nous allons calculer  $\lim_{n \rightarrow +\infty} S_n$  avec  $S_n = \sum_{k=2}^n u_k$  en utilisant la remarque indispensable suivante :

$$\frac{1}{x^2 - 1} = \frac{1}{2} \left( \frac{1}{x-1} - \frac{1}{x+1} \right)$$

Cela donne :

$$\begin{aligned} S_n &= \sum_{k=2}^n u_k = \sum_{k=2}^n \frac{1}{k^2 - 1} = \frac{1}{2} \sum_{k=2}^n \left( \frac{1}{k-1} - \frac{1}{k+1} \right) \\ S_n &= \frac{1}{2} \left( \sum_{k=2}^n \frac{1}{k-1} - \sum_{k=2}^n \frac{1}{k+1} \right) \end{aligned}$$

Dans la première sommation remplaçons  $k$  par  $j+1$  et dans la deuxième  $k$  par  $j-1$  :

$$S_n = \frac{1}{2} \left( \sum_{j=1}^{n-1} \frac{1}{j} - \sum_{j=3}^{n+1} \frac{1}{j} \right) = \frac{1}{2} \left( 1 + \frac{1}{2} - \frac{1}{n} - \frac{1}{n+1} \right)$$

On obtient alors

$$\lim_{n \rightarrow +\infty} S_n = \frac{3}{4}$$

Il est alors permis d'écrire que  $\sum_{n=2}^{+\infty} u_n = \sum_{n=2}^{+\infty} \frac{1}{n^2 - 1}$  a un sens ou existe et vaut  $\frac{3}{4}$ . On dit aussi que la somme de la série est  $\frac{3}{4}$ .

#### Remarque

Il ne faut pas croire que l'on est capable de trouver la somme de toute série convergente, c'est plutôt rare, par contre, on est très souvent capable de dire si une série converge ou non mais nous ne présenterons pas les outils évolués qui le permettent.

#### 6 5 Série géométrique

On appelle série géométrique toute série dont le terme général est de la forme  $u_n = q^n$ .

D'après le théorème précédent, il est manifeste (**pourquoi ? !**) qu'une série géométrique converge si, et seulement si,

$$|q| < 1$$

Calculons dans ce cas la somme de la série :

$$\begin{aligned} S_n &= \sum_{j=0}^n q^j = 1 + q + q^2 + q^3 + \dots + q^n = \frac{1 - q^{n+1}}{1 - q} \text{ et} \\ \lim_{n \rightarrow +\infty} S_n &= \frac{1}{1 - q} = \sum_{j=0}^{+\infty} q^j \end{aligned}$$

#### 6 6 Série exponentielle

C'est la série dont le terme général est de la forme  $u_n = \frac{x^n}{n!}$  où  $x$  est un réel quelconque. Nous admettons (la démonstration vous sera donnée peut-être l'année prochaine) qu'une telle série converge et que

$$\sum_{k=0}^{+\infty} \frac{x^k}{k!} = e^x, \text{ pour tout } x \text{ réel.}$$

## EXERCICES

**Recherche 3 - 1**

On considère que op est une opération à temps constant.

Quelle est la complexité des boucles suivantes :

```
Pour i de 1 à n Faire
| op
FinPour
```

```
Pour i de 1 à n Faire
| Pour j de 1 à 90n Faire
| | op
| FinPour
| Pour k de 40n à 1 Faire
| | op
| FinPour
FinPour
```

```
Pour i de 1 à n*n Faire
| Pour j de 1 à n*3 Faire
| | Pour k de 1 à n Faire
| | | op
| | FinPour
| FinPour
FinPour
```

```
Pour i de 1 à n Faire
| Pour j de 1 à i Faire
| | op
| FinPour
FinPour
```

```
Pour i de 1 à n Faire
| Pour j de 1 à n Faire
| | op
| FinPour
| Pour k de 1 à n Faire
| | op
| FinPour
FinPour
```

```
Pour i de 1 à n Faire
| Pour j de 1 à i Faire
| | Pour k de 1 à j Faire
| | | op
| | FinPour
| FinPour
FinPour
```

```
Pour i de 1 à n Faire
| Pour j de 1 à n Faire
| | op
| FinPour
FinPour
```

```
s ← 0
i ← n
TantQue i > 0 Faire
| Pour j de 0 à i - 1 Faire
| | s ← s + 1
| FinPour
| i ← i // 2
FinTantQue
Retourner s
```

```
Pour i de 1 à n Faire
| Pour j de 1 à n Faire
| | Pour k de 1 à n Faire
| | | op
| | FinPour
| FinPour
FinPour
```

```

s ← 0
i ← 1
TantQue i < n Faire
  Pour j de 0 à i - 1 Faire
    | s ← s + 1
  FinPour
  i ← i * 2
FinTantQue
Retourner s

```

```

s ← 0
i ← 1
TantQue i < n Faire
  Pour j de 0 à n - 1 Faire
    | s ← s + 1
  FinPour
  i ← i * 2
FinTantQue
Retourner s

```

### Recherche 3 - 2 CCP MP 2015

Non, non, ce n'est pas un extrait d'un livre de seconde, c'est un sujet 0 proposé pour l'épreuve de Math du concours CCP...

1. Écrire une fonction factorielle qui prend en argument un entier naturel  $n$  et renvoie  $n!$  (on n'acceptera pas bien sûr de réponse utilisant la propre fonction factorielle du module `math` de Python ou Scilab).
2. Écrire une fonction `seuil` qui prend en argument un entier  $M$  et renvoie le plus petit entier naturel  $n$  tel que  $n! > M$ .
3. Écrire une fonction booléenne nommée `est_divisible`, qui prend en argument un entier naturel  $n$  et renvoie `True` si  $n!$  est divisible par  $n + 1$  et `False` sinon.
4. On considère la fonction suivante nommée `mystere` :

---

```

1 def mystere(n):
2     s = 0
3     for k in range(1,n+1):
4         s = s + factorielle(k)
5     return s

```

---

- i. Quelle valeur renvoie `mystere(4)` ?
- ii. Déterminer le nombre de multiplications qu'effectue `mystere(n)`.
- iii. Proposer une amélioration du script de la fonction `mystere` afin d'obtenir une complexité linéaire.

### Recherche 3 - 3 Questions sur le cours

Répondez aux questions de la page page 69

### Recherche 3 - 4 Dichotomie

Déterminez la constante asymptotique de la méthode dichotomique.

Déterminez un algorithme (impératif et récursif) donnant l'approximation d'une solution avec comme arguments la fonction  $f$ , les bornes de l'intervalle de départ  $a$  et  $b$  et la précision  $p$ .

Vous éviterez les *tant que* et les conditionnelles basées sur une comparaison de flottants par sécurité...

---

```

1 *Main> dico2 (\ x -> x**2 - 2) 1 2 1e-15
2 1.414213562373095
3 *Main> sqrt 2

```

---

### Recherche 3 - 5 Babylone

Existe-t-il un lien entre l'algorithme de Babylone et la méthode de NEWTON-RAPHSON ?

### Recherche 3 - 6 Newton 1

Répondez aux questions posées dans les encadré *recherche* de la section consacrée à la méthode de NEWTON-RAPHSON

**Recherche 3 - 7 Newton 2**

Programmez cette méthode...On pourra créer une fonction :

```
1 newtRaph :: (Fractional a, Ord a) => (a -> a) -> (a -> a) -> a -> a -> a
2 newtRaph                f          f'          x    p = ...
```

Et vérifiez par exemple :

```
1 *Main> newtRaph (\ t -> t**2 - 2) (\ t -> 2*t) 1 1e-15
2 1.4142135623730951
```

Commentez ces résultats :

```
1 *Main> newtRaph (\ t -> t**3 - t) (\ t -> 3*t**2 - 1) 0.25 1e-15
2 0.0
3 *Main> newtRaph (\ t -> t**3 - t) (\ t -> 3*t**2 - 1) 0.5 1e-15
4 -1.0
5 *Main> newtRaph (\ t -> t**3 - t) (\ t -> 3*t**2 - 1) (-1/(sqrt 5)) 1e-15
6 C-c C-cInterrupted.
```

**Recherche 3 - 8 Newton 3 : calcul de l'inverse sur machine**

Utilisez la méthode de NEWTON-RAPHSON pour déterminer une approximation de l'inverse d'un flottant sur machine.

```
1 *Main> inverse 1.5 1 1e-15
2 0.6666666666666666
```

Pour étudier les problèmes dus aux arrondis, on pourra lire les pages 157 à 167 de ?.

**Recherche 3 - 9 Newton 4 : influence de la première approximation**

Notons  $\varepsilon_n$  l'erreur relative après  $n$  itérations par la méthode de NEWTON-RAPHSON pour calculer  $\sqrt{a}$ .

Alors  $\varepsilon_n = \frac{x_n - \sqrt{a}}{\sqrt{a}}$  et donc  $x_n = \frac{1 + \varepsilon_n}{\sqrt{a}}$ . Or  $x_{n+1} = \frac{1}{2} \left( x_n + \frac{a}{x_n} \right)$ .

Montrez qu'alors

$$\varepsilon_{n+1} = \frac{\varepsilon_n^2}{2(1 + \varepsilon_n)}$$

Si on est trop loin de la solution au départ,  $\varepsilon_n$  n'est pas négligeable devant 1 : comparez  $\varepsilon_{n+1}$  et  $\varepsilon_n$ .

Que se passe-t-il sinon ?

**Recherche 3 - 10 Méthode de la sécante**

Le calcul de la dérivée peut être embêtant sur machine. Que peut-il se passer si on remplace  $f'(x_n)$  par  $\frac{f(x_n) - f(x_{n-1})}{x_n - x_{n-1}}$  ?

**Recherche 3 - 11 Exponentiations**

On veut calculer  $x^n$  avec  $n$  un entier naturel et  $x$  un élément d'un ensemble muni d'une loi multiplicative associative.

On pose  $p_1 = x$  et  $p_j = x^j$ .

On impose de déterminer un mécanisme qui calcule  $p_i$  si on connaît déjà  $p_1, p_2, \dots, p_{i-1}$  sous la forme :

$$p_i = p_j \cdot p_k \quad \text{avec } 0 \leq j, k \leq i - 1$$

1. Donnez un algorithme naïf. Quel est son coût ?
2. On écrit l'exposant en binaire. On remplace chaque 1 par CX et chaque 0 par C. On enlève le premier couple CX (le plus à gauche).

On traduit C par « mettre au carré » et X par « multiplier par  $x$  ».

Traduisez cet algorithme de manière plus constructive et étudiez sa complexité. Vous donnerez une version impérative et une version récursive.

### Recherche 3 - 12 Recherches séquentielles et par saut

Un fichier séquentiel informatisé est schématiquement constitué ainsi, chaque fiche contient

- L'adresse  $S$  de la fiche suivante si elle existe sinon **nil** pour indiquer la fin du fichier.
- L'adresse  $P$  de la fiche précédente si elle existe sinon **nil** pour indiquer qu'il n'y a rien en deçà.
- Une donnée  $D_i$  alphanumérique
- Une première adresse **début** qui est l'adresse de la première fiche.

Nous supposons que le fichier est trié suivant l'ordre croissant et pour simplifier l'étude nous poserons  $D_i = i$  (la donnée est égale au numéro de la fiche). Le fichier contient  $N$  fiches. Nous appellerons temps d'accès à la fiche numéro  $i$  le nombre de fiches lues en partant de **début** pour arriver à la fiche  $i$  y compris la fiche  $i$ . Ainsi le temps d'accès à la fiche 1 est 1.

#### Recherche séquentielle.

Pour trouver la fiche  $i$  ( $1 \leq i \leq N$ ) on part de **début**, on lit les fiches dans l'ordre des numéros jusqu'à obtenir la fiche  $i$ . Si on cherche successivement les fiches 1,2,3,...,  $N$  en repartant de **début** à chaque fois, quel est le temps moyen d'accès à une fiche ?

#### Recherche par saut.

- Nous supposons ici que  $N$  est le carré d'un entier,  $N = a^2$  ( $a \in \mathbb{N}^*$ ). Soit  $k$  un entier divisant  $N$  et supérieur à 1,  $N = kj$ .
- Nous créons un autre fichier dit de nœuds  $N_1, N_2, \dots, N_j$ .
- Le départ, noté **dép**, contient l'adresse du premier nœud.
- Le premier nœud contient l'adresse de la fiche numéro  $k$ , la donnée de la fiche numéro  $k$  l'adresse du nœud suivant et l'adresse de la fiche  $k - 1$ .
- Le deuxième nœud contient l'adresse de la fiche numéro  $2k$ , la donnée de la fiche numéro  $2k$  l'adresse du nœud suivant et de la fiche numéro  $2k - 1$ .
- .....
- Le dernier nœud (le  $j^{\text{ème}}$ ) contient l'adresse de la fiche numéro  $N$ , et celle de  $N - 1$ , la donnée de la fiche numéro  $N$  et **nil**.

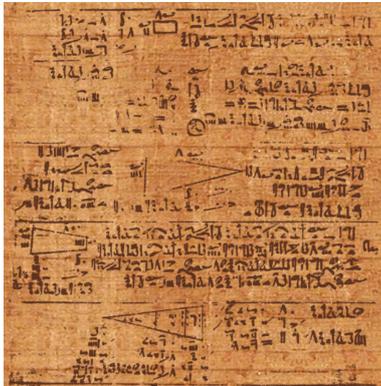
Pour chercher une fiche  $i$  on part de **dép**, on passe au premier nœud et on regarde la donnée  $\alpha$  portée par le nœud. Si  $\alpha < i$  on passe au nœud suivant, sinon on passe à la fiche indiquée par le nœud et, si nécessaire, on revient en arrière dans le fichier jusqu'à atteindre la fiche  $i$ . Un temps d'accès à une fiche sera le nombre de nœuds parcourus + le nombre de fiches parcourues y compris celle cherchée.

Si on cherche successivement les fiches 1,2,3,...,  $N$  en repartant de **dép** à chaque fois, on se propose de déterminer le temps moyen d'accès à une fiche et trouver la valeur de  $k$  qui minimise ce temps moyen.

1. Relier la recherche par saut avec la ?? page ??
2. Déterminer le temps d'accès à la fiche 1.
3. Déterminer le temps d'accès à la fiche 2 (si  $k > 2$ ).
4. Déterminer le temps d'accès à la fiche  $k$ .
5. Déterminer le temps d'accès à la fiche  $k - 1$ .
6. Déterminer le temps d'accès à la fiche  $k + 1$ .
7. Déterminer le temps d'accès à la fiche  $ik$ .
8. Déterminer le temps d'accès à la fiche  $N$ .
9. Déterminer le temps d'accès à la fiche  $j$ .
10. Déterminer le temps total d'accès aux fiches du  $i^{\text{ème}}$  paquet.
11. Déterminer le temps total d'accès à toutes les fiches.
12. Déterminer le temps moyen d'accès à une fiche en fonction de  $k$  et de  $N$ .
13. Déterminer la valeur de  $k$  qui minimise ce temps moyen.
14. Comparer le temps moyen d'accès à une fiche par une recherche séquentielle et une recherche par saut en prenant  $N = 10\,000$  et une unité de temps égale à  $10^{-2}s$ .

**Recherche 3 - 13 Multiplication du paysan russe**

Voici ce qu'apprennaient les petits soviétiques pour multiplier deux entiers. Une variante de cet algorithme a été retrouvée sur le papyrus de Rhind datant de 1650 avant JC, le scribe Ahmes affirmant que cet algorithme était à l'époque vieux de 350 ans. Il a survécu en Europe occidentale jusqu'aux travaux de Fibonacci.



```

1  Fonction MULRUSSE(x:entier ,y: entier, acc:entier): entier
2  Si x==0 Alors
3  | Retourner acc
4  Sinon
5  | Si x est pair Alors
6  | | Retourner MULRUSSE(x/2,y*2,acc)
7  | Sinon
8  | | Retourner MULRUSSE((x-1)/2,y*2,acc+y)
9  | FinSi
10 FinSi

```

Que vaut **acc** au départ ?

Écrivez une version récursive de cet algorithme en évitant l'alternative des lignes 5 à 9.

Écrivez une version impérative de cet algorithme.

Prouvez la correction de cet algorithme.

Étudiez sa complexité.

Vous connaissez  $x \gg 1$  décale l'entier  $x$  d'un bit vers la droite et  $x \ll 1$  décale  $x$  d'un bit vers la gauche en complétant par un zéro à droite,  $x \& y$  renvoie l'entier obtenu en faisant la conjonction logique bit à bit des représentations binaires de  $x$  et  $y$  et  $\sim x$  renvoie le complément à 1 de  $x$ .

Ré-écrivez la multiplication russe en utilisant que ces opérations bit à bit (pas de division ni de multiplication).

**Recherche 3 - 14 Test aléatoire : produit de polynômes correct**

On considère deux polynômes écrits différemment, par exemple l'un sous forme factorisée, l'autre sous forme développée : comment faire pour vérifier que les deux polynômes sont bien égaux ?

Quelle est la complexité en nombre de calculs arithmétiques basiques de votre méthode ?

Imaginons maintenant que les polynômes soient de degré  $d$ . On choisit un nombre entier  $n$  dans l'intervalle  $[0, 100d]$  selon une distribution uniforme et on calcule l'image de ce nombre par chacun des polynômes : comment exploiter le résultat ? Évaluer sa « sûreté » ? Complexité ?

Que pensez-vous de l'issue de l'algorithme par rapport aux issues habituelles ? Pouvez-vous prouver que l'algorithme est correct ?

Que se passe-t-il si on répète cette expérience ? Soyez précis(e) dans la description de votre tirage.

**Recherche 3 - 15 Produit de matrices correct**

On a deux matrices  $A$  et  $B$  dont on a calculé le produit. On a trouvé  $C$ . On voudrait vérifier que le calcul a été bien mené.

Quelle est la complexité habituelle du calcul d'un produit de matrices ?

On préfère choisir un vecteur  $v$  dans  $\{0, 1\}^n$  et effectuer le produit  $D \times v$  en notant  $D = A \times B - C$ .

$D$  étant non nul, il a au moins un coefficient non nul. On peut se ramener au cas  $d_{11} \neq 0$ .

Exprimer alors  $v_1$  en fonction des autres  $v_k$  et des  $d_{1k}$ .

En déduire que si  $D$  est non nulle, la probabilité que  $D \times v = 0_n$  est inférieure à  $1/2$ .

Conclure en terme de vérification et de complexité.

### Recherche 3 - 16 Algorithme de Strassen

En 1969 (« *Gaussian elimination is not optimal* » in *Numerische mathematik*), Konrad STRASSEN propose un algorithme pour calculer le produit de deux matrices carrées de taille  $n$ .

Estimez tout d'abord la complexité en nombre d'additions et en nombre de multiplications de la multiplication usuelle que vous avez vue en Sup.

Considérons à présent le produit de deux matrices de taille  $2 \times 2$  :

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix} \times \begin{pmatrix} e & f \\ g & h \end{pmatrix} = \begin{pmatrix} w & x \\ y & z \end{pmatrix}$$

On note :

$$p_1 = a(f - h), p_2 = (a + b)h, p_3 = (c + d)e, p_4 = d(g - e), p_5 = (a + d)(e + h), p_6 = (b - d)(g + h), p_7 = (c - a)(e + f)$$

alors :

$$\begin{cases} w = p_4 + p_5 + p_6 - p_2 \\ x = p_1 + p_2 \\ y = p_3 + p_4 \\ z = p_1 + p_5 - p_3 + p_7 \end{cases}$$

Comparez le nombre d'opérations effectuées.

Considérez maintenant des matrices carrées de taille  $2^t$  (si ce n'est pas le cas, que peut-on faire?).

On peut découper nos matrices en quatre blocs de taille  $2^{t-1}$  et appliquer récursivement notre algorithme : il faut quand même vérifier quelque chose pour passer des formules sur les coefficients aux formules sur les matrices : quoi donc ?

Comptez le nombre de multiplications et d'additions nécessaires.

C'est une application de *divide et impera* : pourquoi alors ne pas diviser nos matrices en 9 blocs trois fois plus petits ?

### Recherche 3 - 17 Schéma de Horner-Ruffini-Holdred-Newton-Al-Tusi-Liu-Hui...



La méthode que nous allons voir porte le nom du britannique William George HORNER (1786 - 1837) mais en fait elle fut publiée presque 10 ans auparavant par un horloger londonien, Theophilus HOLDRED et simultanément par l'italien Paolo RUFFINI (1765 - 1822) mais fut déjà utilisée par NEWTON 150 ans auparavant et par le chinois ZHU SHIJE cinq siècles plus tôt (vers 1300) et avant lui par le Persan SHARAF AL-DIN AL-MUZAFFAR IBN MUHAMMAD IBN AL-MUZAFFAR AL-TUSI vers (1100) et avant lui par le Chinois LIU HUI (vers 200) révisant un des résultats présent dans *Les Neuf Chapitres sur l'art mathématique* publié avant la naissance de JC...

Il faut cependant noter que RUFFINI l'avait employée en fait comme un moyen de calculer rapidement le quotient et le reste d'un polynôme par  $(X - \alpha)$ .

Dans toute la suite, un polynôme de degré  $n$  sera représenté par le vecteur de ses coefficients. Par exemple,  $[1 \ 2 \ 3]$  correspond au polynôme  $1 + 2x + 3x^2$ .

Prenons l'exemple de  $P(x) = 3x^5 - 2x^4 + 7x^3 + 2x^2 + 5x - 3$ . Le calcul classique nécessite 5 additions et 15 multiplications.

On peut faire pas mal d'économies de calcul en suivant le schéma suivant :

$$\begin{aligned} P(x) &= a_n x^n + \dots + a_2 x^2 + a_1 x + a_0 \\ &= \underbrace{(a_n x^{n-1} + \dots + a_2 x + a_1)}_{\text{on met } x \text{ en facteur}} x + a_0 \\ &= \dots \\ &= (\dots((a_n x + a_{n-1})x + a_{n-2})x + a_{n-3})x + \dots)x + a_0 \end{aligned}$$

Ici cela donne  $P(x) = (((((3x) - 2)x + 7)x + 2)x + 5)x - 3$  c'est-à-dire 5 multiplications et 5 additions.

Comparez les complexités au pire du calcul de  $P(t)$  pour  $t \in \mathbb{K}$ . Vous prendrez comme « unité de complexité » les opérations arithmétiques de base : + - \*.

Déterminez une fonction horner( $P, t$ ) qui évalue le polynôme  $P$  en  $t$  selon le schéma de HORNER.

**Recherche 3 - 18**

On travaille dans  $\mathbb{R}$ . Calculer les sommes :

1.  $\sum_{i=3}^6 ij, \sum_{i=0}^4 i^2, \sum_{k=0}^4 k^2, \sum_{j=0}^4 i^2, \sum_{i=1}^{50} a_j$
2.  $\sum_{i=1}^{n \geq 1} a, \sum_{i=h}^{n \geq h} a, \sum_{i=1}^{n \geq 1} 3$
3.  $\sum_{1 \leq i < j \leq 5} ij, \sum_{1 \leq i < j \leq 5} (i + 2j)$

**Recherche 3 - 19**

On travaille dans  $\mathbb{R}$ , calculer les sommes  $\sum_{i=1}^n i, \sum_{i=50}^{100} j, \sum_{i=h}^k i, \sum_{i=h}^k (i + j), \sum_{i=h}^k ij, \sum_{i=h}^k (\alpha i + j)$

**Recherche 3 - 20**

On considère un tableau de nombres ayant  $n$  lignes et  $p$  colonnes. Les lignes sont numérotées du haut en bas et les colonnes de la gauche vers la droite. Le nombre se trouvant sur la  $i^{\text{ème}}$  ligne et  $j^{\text{ème}}$  colonne est noté  $a_{i,j}$  ou  $a_{ij}$ . Par exemple, si  $n = 4$  et  $p = 3$ , le tableau se présente par

$a_{1,1}$	$a_{1,2}$	$a_{1,3}$
$a_{2,1}$	$a_{2,2}$	$a_{2,3}$
$a_{3,1}$	$a_{3,2}$	$a_{3,3}$
$a_{4,1}$	$a_{4,2}$	$a_{4,3}$

On note  $L_i$  la somme des nombres de la ligne numéro  $i$ ,  $C_j$  la somme des nombres de la colonne  $j$  et  $T$  la somme de tous les nombres du tableau.

1. Exprimer  $L_i$  et  $C_j$  à l'aide d'un symbole  $\sum$ .
2. Exprimer  $T$  de deux façons.

**Recherche 3 - 21**

On reprend les notations de l'exercice précédent avec  $n = p = 5$ . Que calculent les sommes suivantes? On pourra donner la solution en grisant les cases du tableau qui correspondent aux nombres intervenant dans les sommes.

- |                                       |                                       |  |
|---------------------------------------|---------------------------------------|--|
| 1. $\sum_{i=1}^5 \sum_{j=1}^5 a_{ij}$ | 3. $\sum_{i=1}^5 \sum_{j=1}^i a_{ij}$ | 5. $\sum_{i=1}^5 \sum_{j \geq i} a_{ij}$ |
| 2. $\sum_{j=1}^5 \sum_{i=1}^5 a_{ij}$ | 4. $\sum_{j=1}^5 \sum_{i=1}^j a_{ij}$ | 6. $\sum_{i+j=6} a_{ij}$                 |
|                                       |                                       | 7. $\sum_{1 \leq i < j \leq 5} a_{ij}$   |

**Recherche 3 - 22**

Un fichier séquentiel informatisé est schématiquement constitué ainsi, chaque fiche contient

- L'adresse  $S$  de la fiche suivante si elle existe sinon **nil** pour indiquer la fin du fichier.
- L'adresse  $P$  de la fiche précédente si elle existe sinon **nil** pour indiquer qu'il n'y a rien en deçà.
- Une donnée  $D_i$  alphanumérique
- Une première adresse **début** qui est l'adresse de la première fiche.

Nous supposons que le fichier est trié suivant l'ordre croissant et pour simplifier l'étude nous poserons  $D_i = i$  (la donnée est égale au numéro de la fiche). Le fichier contient  $N$  fiches. Nous appellerons temps d'accès à la fiche numéro  $i$  le nombre de fiches lues en partant de **début** pour arriver à la fiche  $i$  y compris la fiche  $i$ . Ainsi le temps d'accès à la fiche 1 est 1.

**Recherche séquentielle.**

Pour trouver la fiche  $i$  ( $1 \leq i \leq N$ ) on part de **début**, on lit les fiches dans l'ordre des numéros jusqu'à obtenir la fiche  $i$ .

Si on cherche successivement les fiches  $1, 2, 3, \dots, N$  en repartant de **début** à chaque fois, quel est le temps moyen d'accès à une fiche ?

**Recherche par saut.**

- Nous supposons ici que  $N$  est le carré d'un entier,  $N = a^2$  ( $a \in \mathbb{N}^*$ ). Soit  $k$  un entier divisant  $N$  et supérieur à 1,  $N = kj$ .

- Nous créons un autre fichier dit de nœuds  $N_1, N_2, \dots, N_j$ .
- Le départ, noté *dép*, contient l'adresse du premier nœud.
- Le premier nœud contient l'adresse de la fiche numéro  $k$ , la donnée de la fiche numéro  $k$  l'adresse du nœud suivant et l'adresse de la fiche  $k - 1$ .
- Le deuxième nœud contient l'adresse de la fiche numéro  $2k$ , la donnée de la fiche numéro  $2k$  l'adresse du nœud suivant et de la fiche numéro  $2k - 1$ .
- .....
- Le dernier nœud (le  $j^{ème}$ ) contient l'adresse de la fiche numéro  $N$ , et celle de  $N - 1$ , la donnée de la fiche numéro  $N$  et **nil**.

Pour chercher une fiche  $i$  on part de *dép*, on passe au premier nœud et on regarde la donnée  $\alpha$  portée par le nœud. Si  $\alpha < i$  on passe au nœud suivant, sinon on passe à la fiche indiquée par le nœud et, si nécessaire, on revient en arrière dans le fichier jusqu'à atteindre la fiche  $i$ . Un temps d'accès à une fiche sera le nombre de nœuds parcourus + le nombre de fiches parcourues y compris celle cherchée.

Si on cherche successivement les fiches  $1, 2, 3, \dots, N$  en repartant de *dép* à chaque fois, on se propose de déterminer le temps moyen d'accès à une fiche et trouver la valeur de  $k$  qui minimise ce temps moyen.

1. Déterminer le temps d'accès à la fiche 1.
2. Déterminer le temps d'accès à la fiche 2 (si  $k > 2$ ).
3. Déterminer le temps d'accès à la fiche  $k$ .
4. Déterminer le temps d'accès à la fiche  $k - 1$ .
5. Déterminer le temps d'accès à la fiche  $k + 1$ .
6. Déterminer le temps d'accès à la fiche  $ik$ .
7. Déterminer le temps d'accès à la fiche  $N$ .
8. Déterminer le temps d'accès à la fiche  $j$ .
9. Déterminer le temps total d'accès aux fiches du  $i^{ème}$  paquet.
10. Déterminer le temps total d'accès à toutes les fiches.
11. Déterminer le temps moyen d'accès à une fiche en fonction de  $k$  et de  $N$ .
12. Déterminer la valeur de  $k$  qui minimise ce temps moyen.
13. Comparer le temps moyen d'accès à une fiche par une recherche séquentielle et une recherche par saut en prenant  $N = 1000$  et une unité de temps égale à  $10^{-2}$  s.

**Recherche 3 - 23 Achille et la tortue**

Le paradoxe suivant a été imaginé par ZÉNON D'ÉLÉE (490-430 Avant JC). Achille fait une course avec la tortue. Il part 100 mètres derrière la tortue, mais il va dix fois plus vite qu'elle. Quand Achille arrive au point de départ de la tortue, la tortue a parcouru 10 mètres. Pendant qu'Achille parcourt ces 10 mètres, la tortue a avancé d'un mètre. Pendant qu'Achille parcourt ce mètre, la tortue a avancé de 10cm... Puisqu'on peut réitérer ce raisonnement à l'infini, Zénon conclut qu'Achille ne peut pas dépasser la tortue...

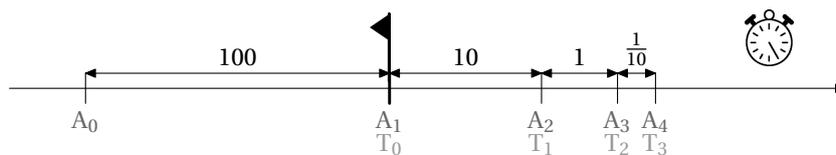
Pour étudier ce problème, nous aurons besoin d'un résultat intermédiaire.

Il est assez simple de démontrer par récurrence, par exemple, que pour tout entier  $n > 1$  et tout réel  $x > -1$ , on a :

$$(1 + x)^n > 1 + nx \quad (\text{Inégalité de BERNOULLI})$$

Qu'en déduisez-vous au sujet de la limite des suites de terme général  $q^n$  ?

Voici la situation :



Intéressons-nous d'abord à la distance Achille-Tortue. Notons  $d_n$  la distance :

$$d_n = T_n - A_n = T_n - T_{n-1} = \frac{1}{10}(T_{n-1} - A_{n-1}) = \frac{1}{10}d_{n-1}$$

La suite  $(d_n)$  est donc géométrique de raison  $1/10$  et de premier terme 100. On en déduit que  $d_n = 100 \times \left(\frac{1}{10}\right)^{n-1}$ .

Or  $|1/10| < 1$ , donc  $\lim_{n \rightarrow +\infty} (1/10)^{n-1} = 0 = \lim_{n \rightarrow +\infty} d_n$ .

Ainsi Achille va rattrapper la tortue, mais au bout d'une infinité de trajets : pour le Grec Zénon, la notion d'infini étant « au-delà du réel » ; pour lui Achille ne rattrapera jamais la tortue, ce qui est absurde.

Intéressons-nous plutôt à la durée du trajet d'Achille pour atteindre la tortue, en supposant sa vitesse constante et égale à  $v$ .

Notons  $t_1 = \frac{d_1}{v} = \frac{100}{v}$ ,  $t_2 = \frac{d_2}{v} = \frac{1}{10} \frac{d_1}{v} = \frac{1}{10} t_1$ , etc. La suite  $(t_n)$  est donc géométrique de raison  $\frac{1}{10}$  de premier terme

$t_1 = 100/v$ , d'où  $t_n = \frac{100}{v} \left(\frac{1}{10}\right)^{n-1}$ . La durée du trajet est alors :

$$\begin{aligned} \tau_n &= t_1 + t_2 + \dots + t_n \\ &= \frac{100}{v} \left( 1 + \frac{1}{10} + \left(\frac{1}{10}\right)^2 + \dots + \left(\frac{1}{10}\right)^{n-1} \right) \\ &= \frac{1000}{9v} \left( 1 - \left(\frac{1}{10}\right)^n \right) \end{aligned}$$

Nous venons de voir qu'Achille atteindra la tortue quand  $n$  tend vers  $+\infty$ . Or nous savons calculer  $\lim_{n \rightarrow +\infty} \tau_n = \frac{1000}{9v}$  qui est un nombre fini. Achille va donc effectuer cette infinité de trajets en un temps fini égal à  $1000/9v$ .

Zénon pensait qu'une somme infinie de termes strictement positifs était nécessairement infinie, d'où le paradoxe à ses yeux. Il a fallu des siècles à l'esprit humain pour dépasser cette *limite*.

Vous pouvez donc prouver le petit théorème suivant :

**Série géométrique**

La suite de terme général

$$S_n = \sum_{k=0}^n q^k$$

converge si, et seulement si,  $|q| < 1$ . Dans ce cas :

$$S = \lim_{n \rightarrow +\infty} S_n = \frac{1}{1 - q}$$

**Théorème 3 - 4**

**Recherche 3 - 24 Tapis de Sierpinski**

Monsieur SIERPINSKI avait ramené d'un voyage en Orient un tapis carré de 1 mètre de côté dont il était très content. Jusqu'au jour où les mites s'introduisirent chez lui.

En 24 heures, elles dévorèrent dans le tapis un carré de côté trois fois plus petit, situé exactement au centre du tapis. En constatant les dégats, Monsieur Sierpinski entra dans une colère noire ! Puis il se consola en se disant qu'il lui restait huit petits carrés de tapis, chacun de la taille du carré disparu. Malheureusement, dans les 12 heures qui suivirent, les mites avaient attaqué les huit petits carrés restants : dans chacun, elles avaient mangé un carré central encore trois fois plus petit. Et dans les 6 heures suivantes elles grignotèrent encore le carré central de chacun des tout petits carrés restants. Et l'histoire se répéta, encore et encore ; à chaque étape, qui se déroulait dans un intervalle de temps deux fois plus petit que l'étape précédente, les mites faisaient des trous de taille trois fois plus petite...

1. Faire des dessins pour bien comprendre la géométrie du tapis troué. Calculer le nombre total de trous dans le tapis de Monsieur Sierpinski après  $n$  étapes. Calculer la surface  $S_n$  de tapis qui n'a pas encore été mangée après  $n$  étapes. Trouver la limite de la suite  $(S_n)_{n \geq 0}$ . Que reste-t-il du tapis à la fin de l'histoire ?
2. Calculer la durée totale du festin « mitique »...

**Recherche 3 - 25**

$|x| < 1$  et on note  $S_n(x) = \sum_{k=0}^n x^k$ ,  $\sigma_n(x) = \sum_{k=0}^n kx^k = \sum_{k=1}^n kx^k$ .

1. Calculer  $S'_n(x)$ .
2. Calculer  $\sum_{k=0}^{+\infty} x^k$ .
3. Calculer la limite de  $S'_n(x)$  lorsque  $n$  tend vers  $+\infty$ .
4. Exprimer  $\sigma_n(x)$  à l'aide de  $S_n(x)$ .
5. Calculer  $\sum_{k=0}^{+\infty} kx^k$ .

**Recherche 3 - 26**

On note  $\Pr([X = k]) = \frac{\lambda^k}{k!} e^{-\lambda}$ . Calculer

1.  $E(X) = \sum_{k=0}^{+\infty} k \Pr([X = k])$ .
2.  $E(X^2) = \sum_{k=0}^{+\infty} k^2 \Pr([X = k])$ .
3.  $E(X^2) - E(X)^2$ .

**Recherche 3 - 27**

On note  $\Pr([X = k]) = p(1-p)^{k-1}$  avec  $p \in ]0, 1[$ . Calculer  $\sum_{k=1}^{+\infty} \Pr([X = k])$  et  $\sum_{k=1}^{+\infty} k \Pr([X = k])$ .

**Recherche 3 - 28**

On considère le tableau  $A = (a_{ij}) = \begin{bmatrix} 1 & 2 & -2 & 4 & 0 \\ 0 & 1 & 2 & -2 & 4 \\ 4 & 0 & 1 & 2 & -2 \\ -2 & 4 & 0 & 1 & 2 \\ 2 & -2 & 4 & 0 & 1 \end{bmatrix}$ . Calculez :

- |   |   |  |
|---|---|--|
| 1. $S_1 = \sum_{i=1}^5 \sum_{j=1}^5 a_{ij}$ | 3. $S_3 = \sum_{i=1}^5 \sum_{j=1}^i a_{ij}$ | 5. $S_5 = \sum_{i=1}^5 \sum_{j \geq i} a_{ij}$ |
| 2. $S_2 = \sum_{j=1}^5 \sum_{i=1}^5 a_{ij}$ | 4. $S_4 = \sum_{j=1}^5 \sum_{i=1}^j a_{ij}$ | 6. $S_6 = \sum_{i+j=6} a_{ij}$                 |
|   |   | 7. $S_7 = \sum_{1 \leq i < j \leq 5} a_{ij}$   |

**Recherche 3 - 29**

Donnez une expression la plus simple possible de  $\sum_{i=p}^{i=s-1} \left( k \times i + \frac{h(h-1)}{2} \right)$  sans symbole  $\Sigma$ , ni  $\dots$ .

**Recherche 3 - 30**

On travaille dans  $\mathbb{R}$ . Calculer les sommes le plus simplement possible sans utiliser ni  $\sum$  ni  $\dots$  :

$S_1 = \sum_{i=3}^6 ij$	$S_{11} = \sum_{i=1}^n i$	$S_{101} = \sum_{i=h}^k (i+j)$
$S_{10} = \sum_{1 \leq i < j \leq 5} (i+2j)$	$S_{100} = \sum_{i=h}^k i$	$S_{110} = \sum_{i=h}^k ij$

**Recherche 3 - 31 9 points**

On considère la fonction  $u$  définie par :

Haskell

```
u n = foldl (\x k -> (x + 4 / x) * 0.5) 2.5 [1 .. n]
```

1. Traduisez cette fonction par une boucle *pour* en pseudo-langage.
2. Traduisez cette fonction sous forme d'une suite  $(u_n)$  définie par récurrence.
3. Démontrez que, pour tout entier naturel,  $u_n \geq 2$ .
4. Démontrez que la suite  $(u_n)$  est décroissante.
5. Déduisez-en le rôle de la fonction Haskell  $u$ .