

POLYTECH SORBONNE - MAIN3

PROJET D'INITIATION

Bibliothèque efficace d'arithmétique d'intervalle en C

BEN ALI Hala
LIU Vincent
NGUYEN Quoc Son
SAYAH Samir

Encadrant : Stef Graillat
Coordinatrice : Cecile Braunstein

Table des matières

Remerciement	2
Motivation	3
1 Représentation des nombres	4
1.1 La norme IEEE 754	4
1.1.1 Représentation des nombres spéciaux et modes d'arrondi	5
1.1.2 Erreurs d'arrondi, absorption et cancellation	5
1.1.3 Perte de l'associativité	5
1.2 L'arithmétique d'intervalles	7
1.2.1 La représentation FP-Intervalle	7
2 Opérations	10
2.1 Addition	10
2.2 Soustraction	11
2.3 Multiplication	11
2.4 Inversion	11
2.5 Division	12
2.6 Intersection	12
3 Algorithme de Newton par intervalle	14
3.1 Principe de la méthode de Newton par intervalle	14
3.2 Calcul de $\sqrt{2}$ par intervalle	14
3.3 Résultat obtenu	14
4 Fonctions usuelles	15
4.1 Exponentielle	15
4.2 Logarithme	15
4.3 Racine	16
4.4 Sinus / Cosinus / Tangente	16
4.4.1 Sinus	16
4.4.2 Cosinus	17
4.4.3 Tangente	18
5 Matrices	18
5.1 Produit Matriciel	18
Conclusion	19
6 Bibliographie	20
7 Annexe	21

Remerciement

Nous tenons à remercier notre encadrant Monsieur **Stef Grailat** pour avoir contribué au succès de notre projet et de nous avoir fait partager son expertise à chaque séance. Grâce à lui, nous avons beaucoup appris sur la représentation des nombres en machine et l'arithmétique d'intervalles.

Nous remercions également notre coordinatrice **Cécile Braunstein** pour avoir organisé le module du projet d'initiation et avoir été présente lorsqu'il y avait nécessité.

Nous remercions également nos camarades de classe qui ont effectué le peer-review sérieusement. Cela a été très important pour nous d'avoir des critiques constructives. Nous avons alors pu améliorer notre rapport et éclaircir certaines notions ou explications qui n'étaient pas claires.

Motivation

La représentation des nombres en machine fait l'objet de nombreuses études de nos jours. La vitesse, la précision et la fiabilité des calculs sont primordiales. Elles le sont d'autant plus lorsque le domaine d'application est critique (aéronautique, finance, aérospatial). Quand on parle de précision on pense immédiatement aux nombres réels ou nombres à virgule flottante en informatique. Notre projet va porter sur la représentation de ces nombres en machine.

La norme la plus utilisée de nos jours pour représenter ces nombres est la norme IEEE 754. Bien que cette norme soit très pratique, elle présente des inconvénients et les résultats des calculs effectués ne sont pas toujours fiables. L'alternative que nous avons choisi ici est d'effectuer des calculs non pas sur des nombres, mais sur des intervalles. En mathématiques et en informatique, l'arithmétique d'intervalle est une méthode de calcul consistant à manipuler des intervalles dans le but d'obtenir des résultats rigoureux.

Cette approche permet de borner les erreurs d'arrondi et ainsi de développer des méthodes numériques **qui fournissent des résultats fiables**.

Dans l'arithmétique d'intervalle, un nombre réel x est représenté par une paire de nombres flottants (x_{inf}, x_{sup}) . Dire que x est représenté par cette paire signifie que x appartient à l'intervalle $[x_{inf}, x_{sup}]$. On parle de la notation *inf* – *sup*. Il existe une autre façon de représenter les intervalles. On peut définir un intervalle par son centre et son rayon. On note $\langle c, r \rangle$ avec $r > 0$ l'intervalle $[c - r, c + r]$. On parle alors de notation centre-rayon.

On voit qu'il est nécessaire d'utiliser deux nombres pour coder un intervalle. Le travail de ce projet serait de définir un format "compressé" pour stocker les intervalles. Nous allons proposer un nouveau format pour représenter les intervalles. Ce qui va changer sera la représentation que l'on se fait d'un nombre binaire, un nombre ne codera plus un réel mais un **intervalle**.

Notre but sera donc d'implémenter une bibliothèque en langage C avec un format permettant de manipuler des intervalles et de créer toutes les fonctions nécessaires au calcul d'intervalles, à savoir : les opérations de base et les fonctions usuelles.

1 Représentation des nombres

1.1 La norme IEEE 754

Sur machine, il est impossible de représenter \mathbb{R} tout entier. On peut seulement représenter un sous-ensemble \mathbb{F} de cardinal fini de \mathbb{R} . Les nombres qui appartiennent à \mathbb{F} sont appelés nombres à virgule flottante.

La norme IEEE 754 permet de représenter les nombres à virgule flottante en binaire avec 3 valeurs : le signe, l'exposant et la mantisse.

Soit $x \in \mathbb{F}$, alors :

$$x = (-1)^{\text{Signe}} \times 1, \text{Mantisse} \times 2^{\text{exposant}} \quad \text{où } \text{exposant} = E - \text{décalage}$$

Le bit de poids fort est le bit de signe. Le bit de signe vaut 0 si le nombre est positif, 1 si le nombre est négatif.

L'exposant du nombre peut être positif ou négatif. Pour représenter un nombre signé en binaire en général, on peut utiliser le complément à deux ou bien le décalage. Dans la norme IEEE 754, on n'utilise pas le complément à deux pour représenter l'exposant, mais le décalage.

L'exposant est décalé de telle manière à ce qu'il soit toujours positif. On va stocker l'exposant dans un entier non signé même si il est négatif ! Il est codé sur 8 bits en précision simple, et sur 11 bits en précision double. L'exposant décalé est obtenu en ajoutant le décalage à l'exposant initiale. Le décalage vaut 127 en précision simple, et 1023 en précision double.

La mantisse représente les chiffres après la virgule. Les nombres jusqu'au dernier 1 sont appelés les chiffres significatifs de x . Ce nombre est codé sous 23 bits en simple précision et 64 bits en double précision.

Précision	Signe	E	Mantisse	Décalage
simple(32 bits)	1 bit	8 bits	23 bits	127
double(64 bits)	1 bit	11 bits	52 bits	1023

La représentation d'un nombre est unique sous la norme IEEE 754.

Exemple

Comment représenter le nombre décimal 8,625 en norme IEEE 754 et en précision simple ?

1. 8,625 est positif, donc le bit de signe vaut **0**.
2. 8,625 en binaire vaut : $8,625_{10} = 1000,101_2$ ($8,625 = 2^3 + 2^{-1} + 2^{-3}$).
3. 1000,101 en écriture scientifique vaut : $1000,101_2 = 1,000101.2^3$.
4. Les chiffres après la virgule sont les chiffres significatifs : **000101**.
La mantisse vaut : **000101000000000000000000** (Les 0 sont ajoutés à droite après les chiffres significatifs pour obtenir une mantisse de taille 23).
5. L'exposant vaut **3**.
L'exposant décalé vaut : $E = \text{exposant} + \text{décalage} = 3 + 127 = 130$.
(en précision simple, décalage = 127).
E en binaire vaut : **10000010**.

La représentation de $8,625_{10}$ en norme IEEE 754 est :

Signe	E	Mantisse
0	10000010	000101000000000000000000

1.1.1 Représentation des nombres spéciaux et modes d'arrondi

Il y a deux zéros et deux infinis en norme IEEE 754 en simple précision :

Nombre	Signe	E	Mantisse
∞	0	11111111	000000000000000000000000
$+0$	0	00000000	000000000000000000000000
-0	1	00000000	000000000000000000000000
$-\infty$	1	11111111	000000000000000000000000

Une opération entre deux nombres qui appartient à l'ensemble \mathbb{F} n'appartient pas forcément à \mathbb{F} , il faut alors arrondir le résultat.

La norme IEEE 754 dispose de quatre modes d'arrondi :

- Vers moins l'infini
- Vers plus l'infini
- Vers zéro
- Au plus près

La manipulation de ces modes d'arrondi nous sera très utile par la suite, en particulier pour borner les éventuelles erreurs d'arrondi.

1.1.2 Erreurs d'arrondi, absorption et cancellation

Lorsqu'on fait des calculs sur machine, la précision des nombres est limitée, la machine renvoie un nombre arrondi du résultat exact du calcul, car il ne peut renvoyer qu'un nombre fixe de chiffres selon le nombre de bits sur lequel le nombre est codé. L'erreur d'arrondi produite quand on remplace un réel x par son représentant $fl(x)$ dans \mathbb{F} est bornée.

Le nombre $u = 2^{-53}$ est l'erreur relative maximale que l'ordinateur peut commettre en représentant un nombre réel en arithmétique finie. Ce nombre va beaucoup nous servir pour borner les erreurs d'arrondi lors de nos calculs d'intervalles.

Les erreurs d'arrondi sont en général négligeables, mais peuvent s'accumuler au cours de longs algorithmes. Une grande accumulation d'erreurs d'arrondi peut fausser un résultat et c'est ici tout l'intérêt de coder un nombre avec un intervalle dans lequel on sera sûr de retrouver le résultat de notre calcul.

En plus des erreurs dues à l'arithmétique, les deux sources principales d'erreurs d'arrondi sont : l'absorption et la cancellation.

Lorsqu'on effectue le calcul $x + y$, dès que y est beaucoup plus petit que x alors $x + y$ est arrondi à x , on dit qu'il y a eu absorption : y a été absorbé par x .

La cancellation survient lors de la soustraction de deux nombres relativement proches, lorsqu'il n'y a presque plus de bits significatifs, on parle de cancellation catastrophique.

1.1.3 Perte de l'associativité

L'ensemble \mathbb{F} ne possède pas les mêmes caractéristiques que \mathbb{R} : en passant du premier au second, on perd l'associativité caractéristique des nombres réels.

Il faut en fait faire attention à l'ordre lorsqu'on fait des opérations sur des nombres, surtout les nombres qui ont des ordres de grandeur très différents. Cela peut provoquer des résultats incohérents.

Exemple

Pour $2^{-40} + 2^2$, la machine va renvoyer 4. Ainsi si on fait : $(2^{-40} + 2^2) - 2^2$, la machine va renvoyer 0 au lieu de 2^{-40} en simple précision, c'est le phénomène d'absorption.

De la même façon, le calcul de $(2 \times 10^{20} + 2.5) - 2 \times 10^{20}$ renvoie 0 au lieu de 2.5, le résultat est alors complètement faux, il y a eu cancellation catastrophique.

1.2 L'arithmétique d'intervalles

Le principe de l'arithmétique d'intervalles consiste à remplacer tous les nombres par des intervalles les contenant. Ainsi comme nous l'avons mentionné auparavant, deux notations existent. La notation inf-sup classique qui consiste à représenter un intervalle sous la forme : $I = [c - r, c + r]$ c étant le centre de l'intervalle et r étant son rayon (par exemple $[0,6]$). Cette notation présente l'inconvénient de devoir stocker deux nombres en machine, ce n'est donc pas celle qui va nous intéresser. L'autre notation possible est la notation centre rayon qui peut s'écrire de la manière suivante : $I = \langle c, r \rangle$ (avec le même exemple : $\langle 3, 3 \rangle$).

1.2.1 La représentation FP-Intervalle

La représentation FP-Intervalle que nous proposons consiste à coder le centre et le rayon d'un intervalle $I = [c - r, c + r] = \langle c, r \rangle$ en un seul nombre. Nous allons utiliser le format IEEE 754. À partir du centre c de l'intervalle, nous allons modifier la mantisse de sorte à ajouter une information supplémentaire sur la valeur du rayon. La nouvelle mantisse permet d'avoir un nouveau nombre qui peut être interpréter comme un FP-Intervalle ou un flottant.

Nous allons présenter deux méthodes :

- Coder un FP-Intervalle à partir d'un centre et d'un rayon.
- Extraire le centre et le rayon à partir d'un FP-Intervalle.

Comment coder un intervalle à partir du centre et du rayon ?

1. Coder le centre c sous le format IEEE 754. (C'est-à-dire connaître les nombres s , E et m)
2. Ajouter un **1** après les derniers chiffres significatifs de la mantisse m .
La position du **1** va indiquer la valeur du rayon.
On a : $r = 2^n$ où $n = \text{exposant} - \text{position}$.

Exemple

Comment coder l'intervalle $I = \langle 8,625, 2^{-4} \rangle$?

1. Coder le centre $c = 8,625$. On a vu dans l'exemple précédent que $c = 8,625$ peut se décomposer ainsi :

Nombre	Signe	E	Mantisse m
8,625	0	10000010	000101000000000000000000

2. Coder le rayon $r = 2^{-4}$.

L'exposant du centre est égale à 3. Il faut alors mettre le **1** à la position 7 ($= 3 - (-4)$) dans la mantisse.

Intervalle	Signe	E	Mantisse m'
I	0	10000010	000101 1 000000000000000000

Les 0 après le **1** en rouge n'ont pas de signification particulière.

On peut aussi donc écrire les intervalles sous cette forme :

Intervalle	Signe	E	Mantisse m'
I	0	10000010	000101_____

Ou en écriture scientifique :

Il est important de comprendre que le nouveau nombre caractérisé par le signe S, l'exposant décalé E et la nouvelle mantisse m' est interprété comme un FP-Intervalle, on a $I = \langle 8,625, 2^{-4} \rangle$.

Attention, ce nombre peut être également interprété comme un simple flottant codé sous la norme IEEE 754 avec $I = 8,6875$.

Autre exemple

Comment coder le rayon $r = 2^{-19}$? L'exposant du centre est égale à 3. Il faut alors mettre le 1 à la position 22 ($= 3 - (-19)$) dans la mantisse.

Signe	E	Mantisse m'
0	10000010	000101000000000000000000 1 0

Pour le centre 8,625 on ne peut seulement placer le **1** entre les position 7 et position 23 dans la mantisse pour coder le rayon car, au delà de la position 7, les bits de poids forts donnent l'information sur le centre.

On ne peut représenter seulement les rayons $r = 2^n$ pour $n \in \{-20, -19, \dots, -4\}$ pour $c = 8,625$ en simple précision.

Comment extraire le rayon et le centre d'un intervalle ?

1. Obtenir l'exposant du nombre.
2. Extraire le dernier **1** dans la mantisse et noter sa position.
3. La puissance de 2 du rayon de l'intervalle est caractérisé par : $r = 2^{\text{exposant} - \text{position}}$.
4. Le reste du nombre binaire sans le **1** est le centre de l'intervalle.

Exemple

Comment extraire le centre et le rayon de l'intervalle $I = 4.5625$?

Intervalle	Signe	E	Mantisse m'
4,5625	0	10000001	00100 1 00000000000000000000

1. E vaut **10000001** en binaire, c'est-à-dire 129.
L'exposant est : $\text{exposant} = E - \text{decalage} = 129 - 127 = 2$.
2. On extrait le dernier **1** dans la mantisse. Il se trouve à la position **6**.
3. Le rayon est alors $r = 2^{2-6} = 2^{-4}$.
4. Le reste du nombre binaire sans le 1 est le centre de l'intervalle :

Nombre	Signe	E	Mantisse m'
4,5	0	10000001	001000000000000000000000

Donc $c = 4.5$.

Le flottant 4,5625 peut donc être interprété comme un FP-Intervalle où $I = < 4.5, 2^{-4} >$.

Restriction de la représentation

Étant donnée qu'on représente tout sur un nombre, il y a des restrictions à propos de la représentation.

1. Le rayon est une puissance de 2.
2. Tous les intervalles ne sont pas représentables.

En particulier, il n'est pas possible de représenter un intervalle dont le rayon est plus grand que le centre.

$$\begin{aligned}\tilde{c}_1 &= ((-1)/(-|c_2| - r_2)) \\ \tilde{c}_2 &= ((-1)/(-|c_2| + r_2))\end{aligned}$$
$$c_3 = (\tilde{c}_1 + 0.5 \times (\tilde{c}_2 - \tilde{c}_1))$$

$r_3 = (c_3 - \tilde{c}_1)$ Afin de bien retrouver l'inverse de notre intervalle C, il faut lui redonner son signe initial, ce qui se traduit par la multiplication suivant sur le nouveau rayon obtenu r_3 : $r_3 = \text{signe}(c_2) \times r_3$

Nombre	Ecriture Scientifique
$A = < 0.25, 2^{-30} >$	+1.00000000000000000000000000000000 _____ * 2^{-2}
$1/A = < 4, 2^{-25} >$	+1.00000000000000000000000000000000 _____ * 2^2

$$r_3 = \frac{(c_1 - c_2 + r_1 + r_2) + u \times 2^{\text{exposant}(c_3)}}{2}$$

Nombre	Ecriture Scientifique
$A = < 42, 2^{-52} >$	+1.010100000000000000000000 _____ * 2^5
$B = < 42, 2^{-40} >$	+1.010100 _____ * 2^5
$A \cap B = < 42, 2^{-52} >$	+1.010100 _____ * 2^5

3 Algorithme de Newton par intervalle

3.1 Principe de la méthode de Newton par intervalle

Nous avons testé notre bibliothèque d'intervalles constituée des opérations élémentaires afin d'approcher $\sqrt{2}$ par la méthode de Newton.

Sur \mathbb{R} , on aurait approcher la valeur $\sqrt{2}$ en résolvant l'équation

$$f(x) = 0$$

avec

$$f(x) = x^2 - 2$$

où f est une fonction de \mathbb{R} dans \mathbb{R} .

La méthode de Newton s'écrit :

$$x_{k+1} = x_k - \frac{x_k^2 - 2}{2 \times x_k}$$

en prenant x_0 assez proche de $\sqrt{2}$.

On ne peut pas simplement remplacer les réels x_k par des intervalles X_k car il y a une **dépendance des données**.

En effet, l'algorithme va calculer pour $x, y, z \in X_k$:

$$x - \frac{y^2 - 2}{2 \times z}$$

au lieu de :

$$x - \frac{x^2 - 2}{2 \times x}$$

pour $x \in X_k$.

Donc on ne peut pas simplement remplacer x_k par des intervalles X_k , car l'algorithme va effectuer des itérations avec trois nombres x, y et z appartenant à X_k différents.

3.2 Calcul de $\sqrt{2}$ par intervalle

Nous allons procéder autrement pour la méthode de Newton par intervalle. Nous utiliserons le résultat suivant pour $f(x) = x^2 - 2$.

Soit $f : \mathbb{R} \rightarrow \mathbb{R}$ une fonction dérivable, $X = [x_1, x_2] \subset \mathbb{R}$ et $\tilde{x} \in X$. Supposons de plus que $0 \notin f'(X)$. On définit alors :

$$N(\tilde{x}, X) := \tilde{x} - \frac{f(\tilde{x})}{f'(X)}$$

Si $N(\tilde{x}, X) \subset X$, alors X contient une racine de f .

Si $N(\tilde{x}, X) \cap X = \emptyset$, alors $f(x) \neq 0$ pour tout $x \in X$.

3.3 Résultat obtenu

On souhaite approcher la valeur $\sqrt{2}$, qui est codé sous le format IEEE 754 de cette manière :

Nombre	Ecriture Scientifique
$\sqrt{2}$	$+1.0110101000001001111001100110011111110011101111001101 * 2^0$

On part de $X_0 =]1.4, 2^{-4}]$ qui contient bien $\sqrt{2}$.

Les itérations obtenus sont :

+1.011_____ * 2 ⁰
+1.0110100011_____ * 2 ⁰
+1.01101010000010100010111111101_____ * 2 ⁰
+1.0110101000001001111001100110011111110110011011_____ * 2 ⁰
+1.011010100000100111100110011001111111001110111100110_____ * 2 ⁰

On voit que le résultat converge vers X_4 .

Les chiffres significatifs de X_4 sont les mêmes que $\sqrt{2}$ jusqu'au dernier 1 en position 52. Donc l'intervalle itéré par l'algorithme de Newton X_4 contient $\sqrt{2} : \sqrt{2} \in X_4$.

Le rayon de cette intervalle X_4 est 2^{-52} , c'est-à-dire que nous avons obtenu une approximation de $\sqrt{2}$ avec une précision de 2^{-52} .

On observe que lorsque l'on passe de X_n à X_{n+1} , on double le nombre de chiffres significatifs pour le centre. Cela correspond à un résultat vu en cours de Systèmes non linéaires et Optimisation :

La méthode de Newton converge quadratiquement.

4 Fonctions usuelles

Dans cette partie, nous allons expliquer la manière de procéder afin d'implémenter des fonctions usuelles en utilisant le format FP-Intervalle. En général, lorsque la fonction considérée est strictement monotone, il suffit de passer en notation inf-sup, et de calculer les nouvelles bornes de l'intervalle.

Une subtilité réside pour les fonctions périodiques de type sinus, cosinus où il faudra gérer les cas des intervalles où la fonction est monotone, et là où la fonction ne l'est pas.

4.1 Exponentielle

On cherche à trouver $\exp(I)$ avec $I = \langle C_I, r_I \rangle$

$I = \langle C_I, r_I \rangle = \{z_I : C_I - r_I \leq z_I \leq C_I + r_I\}$

$I = \langle C_I, r_I \rangle = \{z_I : a \leq z_I \leq b\}$ avec $a = C_I - r_I$ et $b = C_I + r_I$

Alors l'intervalle image de I par la fonction exponentielle qui est **strictement croissante** est :

$K = \exp(I) = \langle C_k, r_k \rangle = \{z : \exp(a) \leq z \leq \exp(b)\}$

donc :

$$\boxed{C_k = \frac{\exp(a) + \exp(b)}{2}}$$

$$\boxed{r_k = \frac{\exp(a) - \exp(b)}{2} + u \times 2^{\text{exposant}(C_k)}}$$

Exemple

Nombre	Ecriture Scientifique
$A = \langle 0, 2^{-10} \rangle$	+1._____ * 2 ⁻¹⁰
$e^A = \langle 1, 2^{-21} \rangle$	+1.00000000000000000000_____ * 2 ⁰

4.2 Logarithme

On cherche à trouver $\log(I)$ avec $I = \langle C_I, r_I \rangle$

$I = \langle C_I, r_I \rangle = \{z_I : C_I - r_I \leq z_I \leq C_I + r_I\}$

$$I = \langle C_I, r_I \rangle = \{z_I : a \leq z_I \leq b\} \text{ avec } a = C_I - r_I \text{ et } b = C_I + r_I$$

Alors l'intervalle image de I par la fonction logarithme qui est **strictement croissante** est :

$$K = \log(I) = \langle C_k, r_k \rangle = \{z : \log(a) \leq z \leq \log(b)\}$$

donc :

$$\begin{cases} C_k = \frac{\log(a)+\log(b)}{2} \\ r_k = \frac{\log(a)-\log(b)}{2} \end{cases}$$

$$\boxed{C_k = \log(\sqrt{a \times b})}$$

$$r_k = \log(\sqrt{\frac{a}{b}}) + u \times 2^{\text{exposant}(C_k)}$$

Example

[illegible]

4.3 Racine

On cherche à trouver \sqrt{I} avec $I = \langle C_I, r_I \rangle$

$$I = \langle C_I, r_I \rangle = \{z_I : C_I - r_I \leq z_I \leq C_I + r_I\}$$
$$I = \langle C_I, r_I \rangle = \{z_I : a \leq z_I \leq b\} \text{ avec } a = C_I - r_I \text{ et } b = C_I + r_I$$

Alors l'intervalle image de I par la fonction racine qui est **strictement croissante** est :

$$K = \sqrt{I} = \langle C_k, r_k \rangle = \{z : \sqrt{a} \leq z \leq \sqrt{b}\}$$

donc :

$$\boxed{C_k = \frac{\sqrt{a} + \sqrt{b}}{2}}$$

$$\boxed{r_k = \frac{\sqrt{a} - \sqrt{b}}{2} + u \times 2^{\text{exposant}(C_k)}}$$

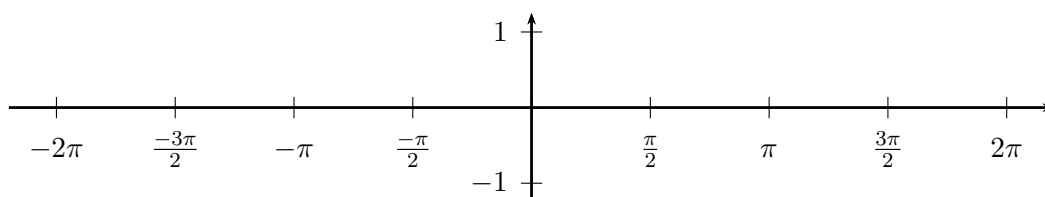
Example

Nombre	Ecriture Scientifique
$A = < 2, 2^{-10} >$	+1.0000000000 _____ * 2 ¹
$\sqrt{A} = < \sqrt{2}, 2^{-52} >$	++ + 1.011010100000100111100110011001111111001110111100110 _____ * 2 ⁰

4.4 Sinus / Cosinus / Tangente

4.4.1 Sinus

on cherche à trouver $\sin(I)$ avec $I = \langle C_I, r_I \rangle = [a, b]$ avec $a = C_I - r_I$ et $b = C_I + r_I$



la fonction $f : \mathbb{R} \rightarrow [-1, 1]$

$x \rightarrow \sin(x)$ est periodique de periode 2π , donc pour simplifier le problème on reduit notre intervalle de telle sorte qu'on se ramène sur $[0, 2\pi]$.

soit $K = \sin(I) = \langle C_k, r_k \rangle = [\min, \max]$ tel que \min est la borne inférieure de K et \max est la

borne supérieure de K .

on a quatre cas :

a) $\mathbf{I} \subset [0, \frac{\pi}{2}]$ ou $\mathbf{I} \subset [\frac{\pi}{2}, \frac{3\pi}{2}]$ ou $\mathbf{I} \subset [\frac{3\pi}{2}, 2\pi]$:

la fonction est monotone donc :

$$\max = \max(\sin(a), \sin(b))$$

$$\min = \min(\sin(a), \sin(b))$$

b) $\frac{\pi}{2}$ et $\frac{3\pi}{2} \in \mathbf{I}$:

la fonction atteint ses bornes supérieure est inférieure en $\frac{\pi}{2}$ et $\frac{3\pi}{2}$ donc :

$$\max = 1$$

$$\min = -1$$

c) **seulement** $\frac{\pi}{2} \in \mathbf{I}$:

la fonction atteint sa borne supérieure en $\frac{\pi}{2}$ donc :

$$\max = 1$$

$$\min = \min(\sin(a), \sin(b))$$

d) **seulement** $\frac{3\pi}{2} \in \mathbf{I}$:

la fonction atteint sa borne inférieure en $\frac{3\pi}{2}$ donc :

$$\max = \max(\sin(a), \sin(b))$$

$$\min = -1$$

Conclusion :

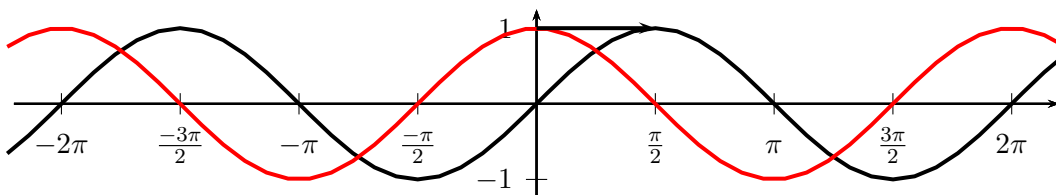
$$\boxed{C_k = \frac{\max + \min}{2}}$$

$$\boxed{r_k = \frac{\max - \min}{2} + u \times 2^{\text{exposant}(C_k)}}$$

4.4.2 Cosinus

il suffit de remarquer que :

$$\boxed{\cos(x) = \sin(x + \pi/2)}$$



Conclusion

En conclusion de ce projet, nous avons pu proposer en langage C Une bibliothèques de fonctions complète et efficace. Elle nous permet de manipuler des intervalles codés avec notre norme, le format FP-INTERV qui est un format compressé basé sur la norme IEEE 754. Les avantages que ce format nous apporte sont nombreux, il nous permet d'avoir des résultats précis et fiables.

Les fonctions que nous avons implémenté couvrent des opérations très classiques telles que l'addition et la multiplication d'intervalles, mais aussi des fonctions plus sophistiquées telles que l'exponentielle et le logarithme d'intervalles, et enfin des matrices d'intervalles.

Un résultat important de notre projet était la possibilité d'obtenir une approximation de $\sqrt{2}$ avec une précision de 2^{-52} à partir de l'algorithme de Newton adapté pour les intervalles.

C'est donc avec une bonne dynamique de groupe que nous avons pu mettre en place notre bibliothèque d'arithmétique d'intervalle, chaque membre a su faire avancer le projet de manière conséquente pour en faire une bibliothèque aboutie et opérationnelle.

Une perspective d'amélioration de notre projet aurait été de manipuler les bibliothèques CBLAS afin de rendre notre produit matriciel plus performant car même s'il est opérationnel, celui-ci reste relativement lent ($O(n^3)$).

Notre projet est disponible sous le lien :

<https://gitlab.com/vinceliu/bibliIntervalleC>

D'un point de vue personnel, ce projet nous a permis de comprendre parfaitement comment l'ordinateur stocke les nombres à virgule flottante. Nous avons pu utiliser de nouvelles bibliothèques en langage C telle que la bibliothèque *fenv* ou *IEEE754*. Cela fut un très bon entraînement afin de mieux maîtriser ce langage de programmation.

Nous avons également appris à mieux travailler en équipe afin de surmonter les difficultés, nous avons pu échanger nos idées, les confronter et tirer le meilleur de ce que chacun avait à proposer, cela a contribué à la réussite de ce projet.

6 Bibliographie

- [1] David Defour. FP-ANR : A representation format to handle floating-point cancellation at run-time, June 2017. working paper or preprint.
- [2] John L. Gustafson. The end of error. Chapman & Hall/CRC Computational Science Series. CRC Press, Boca Raton, FL, 2015.
- [3] Nicholas J. Higham. Accuracy and stability of numerical algorithms. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, second edition, 2002.
- [4] Ramon E. Moore, R. Baker Kearfott, and Michael J. Cloud. Introduction to interval analysis. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2009.

7 Annexe

Addition entre FP-Intervalle

```

1 // Functions to perform addition over FP-INTER format
2 FP_INTERV FpInterAdd(FP_INTERV A, FP_INTERV B)
3 {
4
5     int exprA, exprB, exprC;
6     double rayonC, centreA, centreB, centreC;
7
8     if (A == 0)
9     {
10         return B;
11     }
12     if (B == 0)
13     {
14         return A;
15     }
16
17     // On recupere le rayon et le centre des 2 intervalles
18     centreA = FPInterv2CenterRadius(A, &exprA);
19     centreB = FPInterv2CenterRadius(B, &exprB);
20
21     // nouveau centre : Arrondi de C1 + C2
22     centreC = centreA + centreB;
23
24     int e;
25     // On recupere la puissance de 2 du centre
26     frexp(centreC, &e);
27     e = e - 1;
28     // On change le mode d'arrondi
29     fesetround(FE_UPWARD);
30     // On calcule le nouveau rayon
31     rayonC = pow(2, exprA) + pow(2, exprB) + pow(2, -53 + e);
32     // On change le mode d'arrondi
33     fesetround(FE_TONEAREST);
34     // frexp(rayonC, &e);
35     // On calcule la puissance de 2 du rayon
36     exprC = (int) ceil(log(rayonC) / log(2));
37
38     return CenterRadius2FPInter(centreC, exprC);
39 }
40

```

Voici un exemple d'implémentation en langage C de l'opération élémentaire qui, à partir de deux intervalles de type "FP_Interv" retourne l'intervalle additionnée. Le code source se trouve dans le fichier "operation.c".

A partir de deux intervalles A et B , on extrait les centres c_A et c_B et les exposants des rayons $exprA$, et $exprB$ avec une fonction de conversion "FPInterv2CenterRadius" que nous avons écrite. On calcule l'arrondi $c_C = c_A + c_B$, puis on prend la puissance de 2 de c_C avec la fonction *frexp* qui se trouve dans la bibliothèque *math*.

L'erreur d'arrondi lors du calcul de c_C est bornée par $u \cdot 2^{exp(c_C)}$, donc par le terme $pow(2, -53 + e)$. Lors du calcul du rayon r_C , on change le mode d'arrondi, pour borner l'erreur d'arrondi de l'opération $r_A + r_B$. On va arrondir vers $+\infty$. On utilise pour cela, la fonction *fsetround* qui se trouve dans la bibliothèque *fenv*.

Enfin après avoir obtenu la puissance de 2 du nouveau rayon, on renvoie la fonction "CenterRadius2FPInter". Cette fonction, à partir du centre et du rayon retourne l'intervalle correspondant.

Exponentielle de FP-intervalle

```

1 // Function to perform Exp function over FP-INTER format
2 FP_INTERV FpInterExp(FP_INTERV I)
3 {
4     double center , centerExp , RadiusExp ;
5     int exprExp , expr , e ;
6     double a , b ;
7
8     center=FPInterv2CenterRadius(I,&expr) ;
9     // On passe en notation inf-sup
10    fesetround(FE_DOWNWARD) ;
11    a = center-pow(2,expr) ;
12    fesetround(FE_UPWARD) ;
13    b = center+pow(2,expr) ;
14
15    // Centre en inf-sup
16    fesetround(FE_TONEAREST) ;
17    centerExp = (exp(b)+exp(a))/2 ;
18    // On récupère la puissance de 2 du centre
19    frexp(centerExp,&e) ;
20    e-- ;
21    // Calcul du nouveau rayon
22    fesetround(FE_UPWARD) ;
23    RadiusExp = (exp(b)-exp(a))/2 + pow(2,-53+e) + pow(2,-1074) ;
24    // Puissance de 2 du nouveau rayon
25    exprExp = (int) ceil(log(RadiusExp)/log(2)) ;
26
27    return CenterRadius2FPInter(centerExp , exprExp) ;
28 }

```

$$w(e_{5,2}) = \frac{|N_5 \cap N_2|}{|N_5 \cup N_2|} = \frac{|\{1\}|}{|\{1,3,4\}|} = \frac{1}{3}$$

Voici un second exemple d'implémentation en langage C d'une de nos fonction : la fonction exponentielle. La fonction exponentielle appliquée à un intervalle I de type "FP_Interv" prend en entrée l'intervalle I et retourne l'intervalle image par la fonction exponentielle. Le code source de la fonction se trouve dans le fichier "*function.c*".

A partir de l'intervalle I , on extrait comme précédemment le centre ainsi que la puissance de 2 du rayon de l'intervalle.

Il faut ensuite calculer la borne inférieur a et la borne supérieur b de notre intervalle I . Pour cela, il faut gérer les modes d'arrondis avec la fonction *fsetround*.

A partir des deux nombres réels a et b , les bornes de l'intervalle $\exp(I)$ sont simplement e^a et e^b . On en déduit le nouveau centre de $\exp(I)$ par $\frac{e^a+e^b}{2}$ sans oublier le mode d'arrondi correspondant. Ensuite, la méthode pour calculer le nouveau rayon se fait de manière analogue que la première fonction.

Enfin, pour calculer le nouveau intervalle, il faut faire simplement appel à la fonction de conversion qui prend en paramètre le nouveau centre et la puissance de 2 du rayon.