

RÉVISIONS 3 : REPRÉSENTATION DES NOMBRES EN MACHINE

Sommaire

1	Représentation des entiers naturels	2
1.1	Pour l'homme	2
	a - La base 10	2
	b - Les autres bases	2
1.2	...et pour la machine : la base 2 obligatoire!!!	3
	a - Principe	3
	b - Conversion	4
2	Représentation des entiers relatifs	5
2.1	Notation en complément à deux	5
2.2	Le problème du dépassement de capacité - cas particulier de Python	7
3	Représentation des réels	8
3.1	La notation scientifique décimale	8
	a - Nécessité des "flottants" en machine	8
	b - Décomposition en notation scientifique décimale	9
3.2	Décomposition en base binaire	9
	a - Principe	9
	b - Ecriture normalisée	9
3.3	Norme IEEE754 - limitation des nombres en virgule flottante	11
3.4	Exemples de représentation en mémoire et conversions (partie optionnelle!)	12

1 Représentation des entiers naturels

1.1 Pour l'homme

a - La base 10

La base de représentation décimale ou **base 10** des nombres fut probablement adoptée dans le cours de la préhistoire par les peuples *indo-européens*. Elle découle très naturellement du nombre total de doigts présents sur les mains humaines, les hommes s'en servant pour faire du dénombrement élémentaire.

Le principe de décomposition d'un entier en base 10 est d'écrire le nombre en une somme d'autant de puissances de 10 que nécessaire. Par exemple :

$$6753 = 3 \cdot 10^0 + 5 \cdot 10^1 + 7 \cdot 10^2 + 6 \cdot 10^3$$

En généralisant à tout nombre entier naturel n , cela donne :

$$n = \left(\begin{array}{cccc} d_{n_d} & \dots & d_1 & d_0 \\ \uparrow & & \uparrow & \\ \text{poids le} & & \text{poids le} & \\ \text{plus fort} & & \text{faible} & \end{array} \right)_{10} = \sum_{k=0}^{k=n_d} d_k \cdot 10^k \quad \text{avec } d_k \in \{0, 1, 2, \dots, 9\}$$

EXERCICE N°1: Interpréter le fonctionnement du script python suivant :

Listing III.1 – Que fais-je ?

```

1 def Chiffre (n) :
2     """n est un entier naturel donné par son écriture décimale"""
3     c=[]
4     while n!=0 :
5         c.append(n%10)
6         n//=10
7     c.reverse() # inverse l'ordre des éléments d'une liste
8     return c
9 print Chiffre (4587)

```

b - Les autres bases

On peut représenter un nombre entier dans n'importe quelle base selon ce même principe. Très généralement, un nombre en base k peut s'écrire :

$$(d_{n_k} \dots, d_1, d_0)_k = d_0 \cdot k^0 + d_1 \cdot k^1 + \dots + d_{n_k} \cdot k^{n_k}$$

REMARQUE : la conversion en base décimale est immédiate en réalisant à la main la somme décomposée précédente. Voyons cela sur un exemple, le nombre suivant en base 5

$$(204003)_5 = 3 \times 5^0 + 0 \times 5^1 + 0 \times 5^2 + 4 \times 5^3 + 0 \times 5^4 + 2 \times 5^5 = 6753$$

NB : les TD seront l'occasion de revenir sur les algorithmes simples de transformation des nombres entre bases.

1.2 ...et pour la machine : la base 2 obligatoire!!!

a - Principe

Nous avons déjà évoqué dans le cours Révisions 1 la limitation des ordinateurs à la manipulation de **deux états de potentiels haut et bas**. Par exemple, les potentiels des "pattes" d'une mémoire DDR peuvent prendre uniquement des valeurs $\{0\text{ V} ; 1,5\text{ V}\}$. S'agissant d'une base à deux états, nous choisissons d'écrire ces derniers **0** et **1**. Ce sont les chiffres de la base **2** ou **base binaire**.

Les puces mémoires/processeurs sont constituées de pattes, chacune pouvant être au chiffre **0** ou **1**. On appelle cela un **bit** pour l'anglais **B**inary **I**gT. Combien d'états et donc de valeurs peut recevoir une puce constituée de n pattes ?

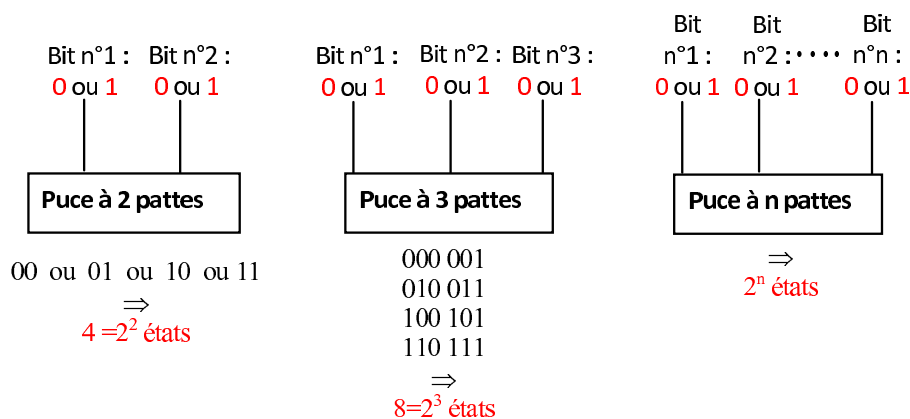


FIGURE III.1 – Etats possibles d'une puce à n pattes

Ainsi donc, il nous faudra coder les nombres en base 2 à l'aide des deux chiffres $\{0, 1\}$. Selon le même principe que précédemment, un entier naturel n s'écrit donc :

$$n = (d_{n_b} \dots d_1 d_0)_2 = \sum_{k=0}^{k=n_b} d_k \cdot 2^k \quad \text{avec } d_k \in \{0, 1\}$$

REMARQUE - (1.2) - 1: IMPORTANT : Les circuits informatiques telles que les puces mémoires ou les processeurs regroupent souvent les nombres binaires par lots de 8 bits, on appelle cela **un octet**. Les nombres informatiques sont ainsi exprimés sur 1, 2, 4, ou encore 8 octets, soit respectivement 8 bits, 16 bits, 32 bits ou 64 bits. Dans la structure des mémoires actuelles des ordinateurs courants, la plus unité adressable que l'on appelle **Bytes** en anglais, est justement codée sur 8 bits soit 1 octet (ce n'était par exemple pas le cas dans les années 70 où le byte pouvait correspondre à 6,7,8 ou 9 bits).

Ces différentes "longueurs" de "mots" offrent les possibilités de codage d'entiers suivantes :

Profondeur en octets	intervalle de codage entiers base 2		intervalle de codage entiers base 10	
1 (8 bits)	00000000	→	11111111	0 → 255
2 (16 bits)	00000000 00000000	→	11111111 11111111	0 → 65535
4 (32 bits)	00000000 00000000 00000000 00000000	→	11111111 11111111 11111111 11111111	0 → 4294967295

A RETENIR : Une profondeur d'écriture binaire de N bits permet de représenter les nombres entiers en base 10 de l'intervalle $[0..2^N - 1]$

b - Conversion

QUESTION : comment représenter en base 2 format 16 bits un nombre entier naturel donné en base 10 ?

RÉPONSE : il existe essentiellement 2 méthodes : la méthode dite «par soustraction» et la méthode dite «par division» (entière). On donne ici le principe de la seconde :

On souhaite décomposer le nombre n_{10} en base 2. Pour cela :

- On réalise la division entière de n_{10} par 2. Si le quotient n'est pas nul, la valeur du reste (0 ou 1) correspond au 1^{er} bit : 0 ou 1.
- On réalise ensuite la division entière du précédent quotient par 2. De même si le quotient n'est pas nul, la valeur du reste (0 ou 1) correspond au 2nd bit.
- Ainsi de suite jusqu'au bit de poids le plus fort.

EXEMPLE SIMPLE : cherchons à décomposer 24_{10} en base 2 par division entières successives :

$$\begin{aligned}
 24 &= (12 \times 2) + 0 \\
 &= ((6 \times 2 + 0) \times 2) + 0 \\
 &= (((3 \times 2 + 0) \times 2 + 0) \times 2) + 0 \\
 &= ((((1 \times 2 + 1) \times 2 + 0) \times 2 + 0) \times 2) + 0 \\
 &= ((((((0 \times 2 + 1) \times 2 + 1) \times 2 + 0) \times 2 + 0) \times 2) + 0) \times 2) + 0 \\
 &= \mathbf{0} \times 2^5 + \mathbf{1} \times 2^4 + \mathbf{1} \times 2^3 + \mathbf{0} \times 2^2 + \mathbf{0} \times 2^1 + \mathbf{0} \times 10^0
 \end{aligned}$$

Ainsi la représentation binaire de $(24)_{10}$ est donc :

$$(1100)_2$$

Par exemple, on souhaite coder en binaire le nombre entier 1515 (Date de la bataille de Marignan¹ !)

1. Tout le monde connaît cette célèbre date de bataille, mais peu de gens savent qu'elle opposa François 1^{er} et les vénitiens aux mercenaires Suisses à **Marignan**, aujourd'hui Melegnano en Italie.

1515	2										
1	757	2									
	1	378	2								
		0	189	2							
			1	94	2						
				0	47	2					
					1	23	2				
						1	11	2			
							1	5	2		
								1	2		
									0	2	
										1	2
										1	0

La représentation binaire de $(1515)_{10}$ est donc :

$$(10111101011)_2$$

2 Représentation des entiers relatifs

2.1 Notation en complément à deux

Une première idée d'écriture des nombres entiers relatifs est de représenter un nombre entier naturel sur une certaine profondeur, par exemple 7 bits ou 15 bits, et de réserver un bit pour le signe, pour obtenir dans notre exemple 8 bits ou 16 bits.

Pour un codage sur N bits, le bit réservé pour le signe S est d_{N-1} avec :

$$S = (-1)^{d_{N-1}} \implies \begin{cases} S = +1 \text{ (n positif) si } d_{N-1} = 0 \\ S = -1 \text{ (n négatif) si } d_{N-1} = 1 \end{cases}$$

La représentation binaire de n est donc :

$$n = (-1)^{d_{N-1}} \sum_{k=0}^{N-2} d_{n_b} 2^k$$

Afin d'être distingué d'un entier naturel, un entier relatif est représenté avec le bit de signe est **souligné** :

$$(\underline{1}0010011)_2 = (-19)_{10}$$

Cette convention de représentation présente cependant deux inconvénients majeurs (par exemple en 8 bits) :

- présence de deux zéros : $(\underline{0}0000000)_2 = (0)_{10}$ et $(\underline{1}0000000)_2 = (-0)_{10}$
- la somme en binaire ne fonctionne plus, avec par exemple :

$$\underbrace{(\underline{0}0000011)_2}_{=(+3)_{10}} + \underbrace{(\underline{1}0000010)_2}_{=(-2)_{10}} = \underbrace{(\underline{1}0000101)_2}_{=(-5)_{10} \neq 1}$$

CONCLUSION : l'objectif d'une machine étant essentiellement de faire des calculs non limités au entiers naturels, et ce rapidement, cette convention de représentation ne peut donc convenir.

IDÉE : On garde la même notation pour les entiers positifs, et on adopte la notation dite en **complément à deux** pour les entiers négatifs. Le principe est le suivant, on veut coder un entier relatif négatif n :

- On garde le $N^{\text{ième}}$ bit d_{N-1} pour le signe et on prend les $N - 1$ bits restants qui représentent la valeur absolue du nombre
- On inverse les $N - 1$ bits de la valeur absolue du nombre (opération booléenne *NON*). Cette opération s'appelle **le complément à 1**.
- On ajoute le $N^{\text{ième}}$ bit de signe
- On ajoute 1 à ce résultat. **c'est le complément à 2**

EXEMPLE :

Codons par exemple $n = (-7)_{10}$ avec cette méthode sur 8 bits :

- Valeur absolue $\rightarrow |n| = 7 = (0000111)_2$
- Complément à 1 $\rightarrow \overline{|n|} = (1111000)_2$
- Bit de signe (ici 1 à la place de l'*underscore*) $\rightarrow (11111000)_2$
- Complément à 2 $\rightarrow (11111001)_2$

REMARQUE : les défauts de la précédente méthode de codage sont éliminés :

- codage de zéro (toujours sur 8 bits) $\left\{ \begin{array}{l} (+0)_{10} = (00000000)_2 \xrightarrow{\text{inv 7 bits}} (01111111)_2 \xrightarrow{+1} 10000000_2 = (-0)_{10} \\ (-0)_{10} = (10000000)_2 \xrightarrow{\text{inv 7 bits}} (11111111)_2 \xrightarrow{+1} 00000000_2 = (+0)_{10} \end{array} \right.$
- $(+3)_{10} + (-2)_{10} = (1)_{10}$ soit en binaire avec $(+3)_{10} = (00000011)_2$ et $(-2)_{10} = (11111110)_2$

$$\begin{array}{r} (00000011)_2 \\ + (11111110)_2 \\ = \underbrace{(00000001)_2}_{=(+1)_{10}} \end{array}$$

CONCLUSION : ainsi donc pour un nombre binaire codés sur N bits, on peut coder :

$$\begin{array}{ll} 2^{N-1} & \text{ nombres positifs } \rightarrow n \in [0, 2^{N-1} - 1] \\ 2^{N-1} & \text{ nombres négatifs } \rightarrow n \in [-2^{N-1}, -1] \end{array}$$

soit un intervalle totale de représentation : $n \in [-2^{N-1}; 2^{N-1} - 1]$

et un total de $2^{N-1} + \underbrace{1}_{\text{pour 0}} + 2^{N-1} - 1 = 2^N$ valeurs

Donc par exemple pour un codage sur 8 bits, on peut représenter l'intervalle d'entiers relatifs : $[-128; +127]$ soit 256 valeurs comme attendu. Plus généralement :

Profondeur (octets)	intervalle de codage entiers base 2		intervalle de codage entiers base 10	
1 (8 bits)	10000000	→	01111111	-128 → 127
2 (16 bits)	10000000 00000000	→	01111111 11111111	-32768 → 32767
4 (32 bits)	10000000 00000000 00000000 00000000	→	01111111 11111111 11111111 11111111	-2147483648 → 2147483647

2.2 Le problème du dépassement de capacité - cas particulier de Python

QUESTION (LÉGITIME!!!) : que se passe-t-il lorsque l'on demande à une machine (via un langage) de traiter des nombres plus grand que le dimensionnement prévu pour les variables ?

RÉPONSE :

Certains langages nécessitent une déclaration préalable du type de chacune des variables employées en préambule de tout programme. C'est par exemple le cas du Pascal ou du VisualBasic, ou encore Fortran ; on appelle cela le **typage statique**. En particulier, la profondeur des entiers doit être précisée. Toute opération de dépassement d'un entier sera signalée par une erreur de capacité de type *overflow* (Pascal), ou conduira à une erreur de calcul, parfois sans aucun signalement ("écrasement" mémoire sans précaution!!!) (Fortran 77).

EXERCICE N°2: Tester la "profondeur" mémoire de votre calculatrice en insérant au clavier le plus grand nombre binaire entier positif²

Dans le cas du langage Python, aucune limite de profondeur n'est fixée!!! Lorsque l'on tente de manipuler un nombre plus grand que la profondeur maximale prévu (32 ou 64 bits suivant les machines et les systèmes d'exploitation), Python découpe l'entier en lots de 15 bits et place les fragments dans un tableau d'entier chacun de largeur 16 bits. Ce tableau contient donc le codage du nombre en base 2^{15} ; on ajoute également dans ce tableau un entier codé sur 16 bits indiquant le nombre de fragments et dont le signe est celui du nombre codé.

Lorsque Python passe dans ce type "propriétaire" de représentation des nombres, il le signale par l'écriture d'un suffixe "L" pour *Large Integer* à la suite du nombre. Par exemple en se plaçant à la limite de codage des entiers relatifs en 32 bits, soit :

$$(01111111111111111111111111111111)_2 = (2147483647)_{10}$$

et en ajoutant 1 :

```
>>> 2147483647+1
2147483648L
>>>
```

2. sur les modèles TI nspire CX/CX CAS, il suffit par exemple d'inscrire "0b" comme préfixe à toute écriture de nombre binaire. C'est le cas sur bon nombre de calculatrices.

REMARQUE - (2.2) - 2: La notation des entiers négatifs en complément à 2 **nécessite de connaître la profondeur en bits disponible sur la machine**, puisque l'on devra inscrire obligatoirement 1 sur le bit de signe qui se trouve **complètement à gauche**. Pour éviter ce souci, les calculatrices permettent de placer un signe $-$ devant la valeur absolue d'un entier relatif binaire pour coder le nombre entier négatif correspondant.

3 Représentation des réels

3.1 La notation scientifique décimale

a - Nécessité des "flottants" en machine

Dans les problèmes de sciences en général, nous sommes souvent amenés à exploiter **des nombres à virgule**. Par exemple, lors de la mesure d'une épaisseur de lame de verre à l'aide d'un Palmer³, une valeur peut-être $e = 1,12 \text{ mm}$.

QUESTION : Si nous souhaitons exploiter cette mesure avec un ordinateur, doit-il absolument comprendre ce qu'est un nombre à virgule?

RÉPONSE : **Non!!!** En choisissant comme unité d'écriture le μm , le résultat devient $h = 1120 \mu\text{m}$, **c'est à dire un entier** que l'ordinateur comprend!!!

Cependant, nous savons que l'intervalle des nombres entiers représentables en machine est limité par la profondeur, c'est à dire le nombre de bits choisis. Ainsi, lorsque nous effectuons un calcul en choisissant de "translater" la virgule complètement à droite pour tous les opérandes afin d'obtenir des entiers (en changeant l'unité), nous obtenons des nombres plus grands susceptibles d'entraîner un "overflow".

Imaginons par exemple que nous codions les entiers sur 32 bits dans une machine et qu'à l'aide de celle-ci nous souhaitions calculer le volume de la lame de verre; un objectif plutôt raisonnable pour les ordinateurs modernes et pourtant!!! Cela donne :

$$\left[\begin{array}{ll} \text{épaisseur} & e = 1120 \mu\text{m} \\ \text{largeur} & l = 2,5 \text{ cm} = 25000 \text{ } \mu\text{m} \\ \text{longueur} & L = 6 \text{ cm} = 60000 \text{ } \mu\text{m} \end{array} \right] \Rightarrow V(\mu\text{m}^3) = 1120 * 25000 * 60000 = 1680000000000 \mu\text{m}^3 > \underbrace{2147483647}_{\text{limite 32 bits}}$$

Ainsi, ce simple calcul provoque un dépassement!!!

SOLUTION : on doit pouvoir étendre la fourchette des nombres représentables. Ceci est possible en définissant un type de nombres dits **à virgule flottante** ou type *float* en Python. Le principe est de pouvoir déplacer à loisir la position de la virgule en indiquant un exposant de puissance de 10; c'est la classique **notation scientifique décimale** dans laquelle on distingue d'une part les chiffres significatifs, et d'autre part les ordres de grandeurs à travers la puissance de 10 :

Par exemple : $c = 299792458 \text{ m.s}^{-1} = 3,0.10^8 \text{ m.s}^{-1}$ et $\mu_0 = 4\pi.10^{-7} \text{ H.m}^{-1}$

Le calcul précédent devient possible si nous souhaitons toujours garder le résultat en μm^3 :

$$\left[\begin{array}{ll} \text{épaisseur} & e = 1120 \mu\text{m} = 112.10^1 \mu\text{m} \\ \text{largeur} & l = 2,5 \text{ cm} = 25000 \text{ } \mu\text{m} = 25.10^3 \mu\text{m} \\ \text{longueur} & L = 6 \text{ cm} = 60000 \text{ } \mu\text{m} = 6.10^4 \mu\text{m} \end{array} \right] \Rightarrow V(\mu\text{m}^3) = 112 * 25 * 6 = \underbrace{16800}_{<2147483647} .10^8 \mu\text{m}^3$$

On évite ainsi les dépassements.

3. cf TP

b - Décomposition en notation scientifique décimale

En physique par exemple, il est fréquent d'employer la notation scientifique décimale. Par exemple, elle est bien commode pour écrire certaines constantes :

$$\begin{cases} c = 2,99792458.10^8 \text{ m.s}^{-1} \\ k = 1,3806488.10^{-23} \text{ J.K}^{-1} \\ \mathcal{N}_a = 6,022142.10^{23} \text{ mol}^{-1} \end{cases}$$

Un réel x , dont le nombre de décimales est fini ou non, peut s'écrire sous la forme suivante :

$$x = (-1)^S \times 10^E \times (d_0 + d_1.10^{-1} + d_2.10^{-2} + \dots + d_i.10^{-i} + \dots) \quad \text{avec} \quad \begin{cases} S = 0, 1 \text{ qui fixe le signe de } x \\ E \text{ entier relatif (exposant)} \\ d_i \in \{0, 1, 2, \dots, 9\} \text{ et on fixe } d_0 \neq 0 \end{cases}$$

Par exemple on peut écrire la constante de Boltzmann k selon cette décomposition :

$$k = (-1)^0 \times 10^{-38} (1 + 3.10^{-1} + 8.10^{-2} + 0.10^{-3} + 6.10^{-4} + 4.10^{-5} + 8.10^{-6} + 8.10^{-7}) \text{ J.K}^{-1}$$

3.2 Décomposition en base binaire**a - Principe**

On peut exploiter l'idée générale de la décomposition précédente pour écrire un réel x mais cette fois en base binaire avec :

$$x = (-1)^S \times 2^E \times (b_0 + b_1.10^{-1} + b_2.10^{-2} + \dots + b_i.10^{-i} + \dots) \quad \text{avec} \quad \begin{cases} S = 0, 1 \text{ qui fixe le signe de } x \\ E \text{ entier relatif (exposant)} \\ b_i \in \{0, 1, 2, \dots, 9\} \text{ et on fixe } b_0 = 1 \end{cases} \quad \text{(III.1)}$$

CONCLUSION : Pour connaître le réel x il faut connaître S , E et les b_i .

b - Ecriture normalisée

La profondeur de codage des nombres en machine étant limitée (en python, c'est la taille mémoire qui est le facteur limitant), on doit fatalement limiter E et la donnée des b_i .

On normalise l'écriture des réels flottants de la manière suivante :

► DÉFINITION DE LA mantisse M

On choisit de ne garder que les $m + 1$ premiers termes b_i ce qui permet de définir **la mantisse** avec $b_0 = 1$:

$$M = 1 + b_1 \times 2^{-1} + b_1 \times 2^{-1} + b_2 \times 2^{-2} + \dots + b_m \times 2^{-m}$$

avec $M \stackrel{b_0=1}{=} \sum_{k=0}^m \frac{b_k}{2^k} < \sum_{k=0}^m 2^{-k} < \sum_{k=0}^{\infty} 2^{-k} = \frac{1}{1 - 2^{-1}} = 2$ ainsi $1 \leq M < 2$

\Rightarrow la mantisse est connue lorsque l'on connaît les m premiers $b_i \rightarrow$ le codage de M occupe m bits dans la mémoire

► DÉFINITION DE L'EXPOSANT E

On réserve e bits pour coder l'exposant.

On écrira en fait en mémoire l'exposant décalé E' tel que :

$$E' = E + 2^{e-1} - 1 \quad \text{en imposant } E' \in [0, 2^e - 1] \text{ donc } E \in \mathbb{N}$$

Ainsi, on peut déterminer l'intervalle des valeurs possibles de l'exposant E :

$$\begin{aligned} 0 &\leq E' \leq 2^e - 1 \\ \Leftrightarrow 0 + 1 - 2^{e-1} &\leq E' + 1 - 2^{e-1} \leq 2^e - 1 + 1 - 2^{e-1} \\ \Leftrightarrow 1 - 2^{e-1} &\leq E \leq 2^{e-1} \end{aligned}$$

$$\longrightarrow E_{norm} \in [1 - 2^{e-1}, 2^{e-1}]$$

REMARQUE - (3.2) - 3:

ATTENTION : cet intervalle est différent de celui employé pour coder un entier relatif sur e bits. La convention adoptée ici ne vaut que pour l'exposant.

Certaines valeurs d'exposant sont réservées pour le codage de cas particulier :

- $E' = 2^e - 1$ soit $E = 2^{e-1}$ est employé pour coder l'infini ou un calcul conduisant à coder un "NaN" i.e. Not A Number, comme par exemple la division par 0.
- $E' = 0$ soit $E = 1 - 2^{e-1}$ est employé pour coder les nombres dénormalisés pour lesquels $M < 1$.

Finalement l'intervalle normalisé de valeurs possibles pour l'exposant est :

$$E_{norm.poss} \in [2 - 2^{e-1}, 2^{e-1} - 1]$$

A RETENIR :

CONCLUSION : un nombre binaire codé sur $m + e + 1$ bits s'écrit :

$$\underbrace{S}_{\text{signe}} \underbrace{c_{e-1} \dots c_1 c_0}_{\text{exposant } E'} \underbrace{b_1 b_2 \dots b_m}_{\text{mantisse}}$$

et sa valeur décimale est donnée par la relation III.1 une fois normalisée, soit :

$$x = (-1)^s \times 2^E \times M$$

REMARQUE - (3.2) - 4: Comme pour les entiers, l'intervalle des réels représentables en mémoire est limité. Ajoutons que s'il est possible de représenter tous les entiers dans l'intervalle donné, ce n'est pas le cas avec les réels puisqu'il y en a une infinité dans tout intervalle. Le nombre de bits consacrés à la mantisse $b_1 \dots b_m$ étant limité à m cela signifie qu'il existe un intervalle minimal "incompressible" entre deux réels successifs représentables. Ainsi, **la plupart des réels décimaux ne sont pas représentables exactement en mémoire avec ce format et seront inscrits de manière approchée avec une précision dépendant de la profondeur de bits.**

3.3 Norme IEEE754 - limitation des nombres en virgule flottante

Depuis 1985, la norme L'Institute of Electrical and Electronics Engineers (ou IEEE) a défini le principe de représentation en machine des nombres à virgule flottante ; il s'agit de la norme IEEE754. Elle fixe le nombre de bits consacrés à M, E et S en fonction de la profondeur de représentation choisie :

Profondeur	Signe S	Exposant décalé E'	Mantisse M
32 bits	1 bits	8 bits	23 bits
64 bits	1 bits	11 bits	52 bits

EXERCICE N°3: Quelques précisions sur cette norme

- 1 Quel est l'intervalle des exposants représentables en norme IEEE754 pour une profondeur de 64 bits ?
- 2 Déterminer l'expression du plus grand réel positif représentable x_M , et donner sa valeur numérique sous la forme $r \times 10^n$ avec n entier et $|r| < 10$, toujours en 64 bits.
- 3 Déterminer l'expression du plus petit réel positif normalisé représentable x_m , et donner sa valeur numérique sous la forme $x \times 10^n$ encore une fois en 64 bits.

RÉPONSE :

- 1 $E \in [-2^{10-1} + 1, 2^{10}] = [-1023 + 1, 1024]$
- 2 $x_M = (-1)^0 \times 2^{1023} \times \left[1 + \sum_{i=1}^{52} \frac{1}{2^i} \right] = 1,79769 \cdot 10^{308}$
- 3 $x_m = (-1)^0 \times 2^{-1022} \times \underbrace{\left[1 + \sum_{i=1}^{52} \frac{0}{2^i} \right]}_{=1} = 2,22507 \cdot 10^{-308}$

Dans la mesure où le nombre de bits consacré à M est limité, il existe une précision absolue de représentation des nombres réels en virgule flottante. Cette précision absolue est définie par la différence entre deux réels flottants consécutifs représentables :

$$\delta x = |x' - x|$$

EXERCICE N°4: Précision de représentation

- 1 Déterminer la plus petite variation de mantisse possible.
- 2 En déduire la précision absolue δx pour les plus petits nombres représentables en 64 bits et pour les plus grands.

REMARQUE - (3.3) - 5: Il est possible de sortir de la norme IEEE754 pour étendre l'intervalle des flottants représentables ou bien augmenter la précision sur un flottant. Il suffit pour cela de modifier le nombre de bits consacrés à e et m avec la contrainte $e + m = cste = 7, 15, 31, \text{ ou } 63 \text{ bits}$ (le bit manquant étant celui réservé au signe du flottant). On a un compromis à trouver entre portée (intervalle des réels représentés) et précision. Pour augmenter la précision, il faut augmenter m , mais alors e diminue et on peut représenter moins d'entiers. Pour augmenter la portée, il faut augmenter e , mais alors m diminue et la précision également.

$0,6 \times 2 = 1,2$	le coefficient b_1 de 2^{-1} est	1
$0,2 \times 2 = 0,4$	le coefficient b_2 de 2^{-2} est	0
$0,4 \times 2 = 0,8$	le coefficient b_3 de 2^{-3} est	0
$0,8 \times 2 = 1,6$	le coefficient b_4 de 2^{-4} est	1
$0,6 \times 2 = 1,2$	le coefficient b_5 de 2^{-5} est	1
...		

\implies on constate que ce codage est périodique, et qu'il ne sera limité que par la profondeur de représentation choisie, or ici nous codons la somme de la mantisse sur 52 bits. Ainsi 0,4 va s'écrire en norme IEEE754 :

$$\underbrace{0}_{=S} \underbrace{0111111101}_{=E'} \underbrace{1001100110011001100110011001100110011001100110011001100110011001}_{\text{coeff. } b_1, b_2, \dots, b_{52}}$$