

Chapitre 5

Etude des nombres réels en machine

La représentation standard (norme *IEEE 754*) des nombres à virgule flottante dits flottants en base 2 sur machine (par exemple Linux), notés X est $(-1)^s \times (1 + m) \times 2^{e-d}$ (flottants normalisés).

Avec :

- s le signe (1 pour négatif, 0 pour positif) sur un bit ;
- m la partie fractionnaire en base 2 sur 23 bits en simple précision (SP, 32 bits) et 52 en double précision (DP, 64 bits), $\frac{1}{2} \leq m < 1$;
- e l'exposant entier non signé sur 8 bits en SP (respectivement 11 en DP) ;
- d le décalage de l'exposant pour coder un nombre entier non signé. Si l'on note q le nombre de bits codant l'exposant, on a $d = 2^{q-1} - 1$ (en SP 127, en DP 1023).

Le premier bit de la mantisse $M = 1 + m$, ($1 \leq M < 2$) qui est implicitement 1 en normalisé n'est pas codé.

Pour densifier les nombres codés autour de 0 on prend des flottants dénormalisés de représentation standard $(-1)^s \times m \times 2^{1-d}$ (l'exposant vrai est en SP -126 et en DP -1022 avec le codage spécial : exposant codé à 0, mantisse codée $\neq 0$).

Le premier bit de la mantisse $0 + m$ qui est implicitement 0 en dénormalisé n'est pas codé.

Les flottants caractéristiques positifs en norme *IEEE 754* sur machine sont les suivants :

- l' $\epsilon_{machine}$, le plus petit flottant positif tel que sur machine $1 \oplus \epsilon_{machine} > 1$;
- X_{min}^n normalisé, $X_{min}^d = 0^+$ (voisin par valeur supérieure de 0) dénormalisé ;
- X_{max}^n normalisé, $X_{max}^d = X_{min}^n$ (voisin par valeur inférieure de X_{min}^n) dénormalisé ;
- le zéro : 0 (codage spécial : exposant codé à 0, mantisse codée à 0) et l'unité : 1 normalisé ;
- l' ∞ (codage spécial : exposant codé avec des 1, mantisse codée à 0) et *NaN* (Not a Number) (codage spécial : exposant codé avec des 1, mantisse codée $\neq 0$).
- l'unité d'arrondi, u égal à $\epsilon_{machine}/2$ (en arrondi au plus près) est tel que pour tout réel x , avec $X_{min}^n \leq |x| \leq X_{max}^n$ on ait $X = x(1 + \delta)$ avec $|\delta| \leq u$.
- Soit X un flottant normalisé, $X^+ = X(1 + \xi)$ avec $|\xi| \leq \epsilon_{machine}$.

Représentation mémoire des flottants :

$$\underline{|s|e_7 \cdots e_0|m_1 \cdots m_{23}|} \text{ sur 32 bits (SP)}$$

$$\underline{|s|e_{10} \cdots e_0|m_1 \cdots m_{52}|} \text{ sur 64 bits (DP)}$$

$$e = \begin{cases} \sum_{i=0}^7 e_i 2^i & e_i \in \{0, 1\} \text{ en SP} \\ \sum_{i=0}^{10} e_i 2^i & e_i \in \{0, 1\} \text{ en DP} \end{cases}$$

$$m = \begin{cases} \sum_{i=1}^{23} m_i 2^{-i} & m_i \in \{0, 1\} \text{ en SP} \\ \sum_{i=1}^{52} m_i 2^{-i} & m_i \in \{0, 1\} \text{ en DP} \end{cases}$$

Étude de flottants caractéristiques

1 $L'\epsilon_{machine}$

$L'\epsilon_{machine}$ est le plus petit flottant positif tel que sur machine $1 \oplus \epsilon_{machine} > 1$.

1.1 Simple précision

1.1.1 Codage binaire

Sur 32 bits, on a :

```
(gdb) x/tw &eps1
0x7fffffffddfd4: 00110100000000000000000000000000
(gdb) █
```

FIGURE 1 – Codage binaire de $l'\epsilon_{machine}$ en SP

Le nombre flottant, en SP sur 32 bits lus de gauche à droite :

- est positif (bit de signe à 0) ;
- est normalisé (l'exposant n'est pas codé par 00000000) ;
- a un exposant codé par 01101000 ;
- a une mantisse codée par 000000000000000000000000.

1.1.2 Valeur exacte

- le bit de signe $s = 0$;
- la valeur de la mantisse codée $m = 0$.

A partir de l'expression binaire de l'exposant codé $(e)_2 = (01101000)_2 = (e_7e_6 \dots e_0)_2$, sa valeur est égale à :

$$\sum_{i=0}^7 e_i 2^i = 2^3 + 2^5 + 2^6 = 8 + 32 + 64 = 104$$

La valeur de l'exposant vrai $e - d = 104 - 127 = -23$.

La valeur exacte de $l'\epsilon_{machine}$ en SP est égale à :

$$(-1)^0 \times (1 + 0) \times 2^{-23} = 2^{-23}$$

1.1.3 Valeur décimale approchée

La valeur décimale approchée de $l'\epsilon_{machine}$ en SP est égale à :

$$2^{-23} \approx 1.1920929 \times 10^{-7}$$

1.2 Double précision

1.2.1 Codage binaire

Sur 64 bits, on a :

Le nombre flottant, en DP sur 64 bits lus de gauche à droite :

- est normalisé (l'exposant n'est pas codé par 00000000);
- a un exposant codé par 00000001;
- a une mantisse codée par 000000000000000000000000.

2.1.2 Valeur exacte

- le bit de signe $s = 0$;
- la valeur de la mantisse codée $m = 0$.

A partir de l'expression binaire de l'exposant codé $(e)_2 = (00000001)_2 = (e_7 e_6 \dots e_0)_2$, sa valeur est égale à :

$$\sum_{i=0}^7 e_i 2^i = 2^0 = 1$$

La valeur de l'exposant vrai $e - d = 1 - 127 = -126$.

La valeur exacte de X_{min}^n en SP est égale à :

$$(-1)^0 \times (1 + 0) \times 2^{-126} = 2^{-126}$$

2.1.3 Valeur décimale approchée

La valeur décimale approchée de X_{min}^n en SP est égale à :

$$2^{-126} \approx 1.1754943 \times 10^{-38}$$

2.2 X_{min}^n normalisé en double précision

2.2.1 Codage binaire

Sur 64 bits, on a :

Le nombre flottant, en DP sur 64 bits lus de gauche à droite :

```
(gdb) x/tg &xMin2
0x7fffffffde20: 000000000001000000000000000000000000000000000000000000000000000000
(gdb)
```

FIGURE 4 – Codage binaire de X_{min}^n en DP

- est positif (bit de signe à 0);
- est normalisé (l'exposant n'est pas codé par 00000000);
- a un exposant codé par 0000000001;
- a une mantisse codée par 00.

2.2.2 Valeur exacte

- le bit de signe $s = 0$;
- la valeur de la mantisse codée $m = 0$.

A partir de l'expression binaire de l'exposant codé $(e)_2 = (0000000001)_2 = (e_{10} e_9 \dots e_0)_2$, sa valeur est égale à :

$$\sum_{i=0}^{10} e_i 2^i = 2^0 = 1$$

sa valeur est égale à :

$$\sum_{i=1}^{52} m_i 2^{-i} = 2^{-52} \frac{1 - 2^{52}}{1 - 2} = 1 - 2^{-52}$$

La valeur exacte de X_{max}^n en DP est égale à :

$$(-1)^0 \times (2 - 2^{-52}) \times 2^{1023} = 2^{1024} - 2^{971}$$

3.2.3 Valeur décimale approchée

La valeur décimale approchée de X_{max}^n en DP est égale à :

$$2^{1024} - 2^{971} \approx 1.7976931348623157 \times 10^{308}$$

3.3 X_{max}^d dénormalisé en simple précision

3.3.1 Codage binaire

Sur 32 bits, on a :

Le nombre flottant, en SP sur 32 bits lus de gauche à droite :

```
(gdb) x/tw &xMax3
0x7fffffffde1c: 00000000011111111111111111111111
(gdb)
```

FIGURE 9 – Codage binaire de X_{max}^d en SP

- est positif (bit de signe à 0) ;
- est dénormalisé (l'exposant est codé par 00000000) ;
- a une mantisse codée par 111111111111111111111111.

3.3.2 Valeur exacte

- le bit de signe $s = 0$;
- la valeur de l'exposant vrai $= -126$;

A partir de l'expression binaire de la mantisse codée $(m)_2 = (0, 111111111111111111111111)_2 = (0, m_1 m_2 \dots m_{23})_2$, sa valeur est égale à :

$$\sum_{i=1}^{23} m_i 2^{-i} = 2^{-23} \frac{1 - 2^{23}}{1 - 2} = 1 - 2^{-23}$$

La valeur exacte de X_{max}^d en SP est égale à :

$$(-1)^0 \times (1 - 2^{-23}) \times 2^{-126} = 2^{-126} - 2^{-149}$$

3.3.3 Valeur décimale approchée

La valeur décimale approchée de X_{max}^d en SP est égale à :

$$2^{-126} - 2^{-149} \approx 1.1754942 \times 10^{-38}$$

