

CH 4 浮點數 Floating Point Number

- A. 定點數(Fixed Point Number)
 - B. 浮點數(Floating Point Number) – 正規化後，第一個位數為 0.xxxx 表示
 - C. 一般浮點數(Floating Point Number)格式 – sign bit、exponential、matnissa
 - D. IEEE754(國際格式)(有保留特殊值) – 1.xxx 表示；超碼： $2^{n-1} - 1$ ；單：1S 8E 23M；雙：1S 11C 52M；全 1 全 0 不能用
 - E. 誤差 – 固有(Inherent)、截尾(truncation)、捨棄(rounding)
-

1. 定點數(Fixed Point Number)

規定：小數點後的小數位數為固定位數

例：1234.5678901 \Rightarrow 定點取 3 位小數 \Rightarrow 1234.567
4 位小數 \Rightarrow 1234.5678

可參考台大 OCW CH02 計概

<http://ocw.aca.ntu.edu.tw/ntu-ocw/ocw/cou/101S210/2>

2. 浮點數(Floating Point Number)

(1) 數學之定義： $\pm (0.F)_b \times b^{\text{exp}}$ $0 < (0.F)_b < 1$, exp=指數, b為基底(base)

(b 幾位，例如:10 進位，b = 10)

(二) 例: $+(1234.5678)_{10}$

$$= +(0.12345678)_{10} \times 10^4 = +(0.012345678)_{10} \times 10^5 = +(0.0012345678)_{10} \times 10^6$$

(三) 上例可知, 數字上有很多種格式表示同一數值。

但對電腦儲存, 希望只有一種標準格式, 須有一個“正規化”動作。

(四) 以一般型之浮點數格式而言, 其正規化要求:

$(0.F)_b$ 中小數點後, 第一個位數值不為 0 (要 > 0)

例: (1) $-(0.0001011)_2 \xrightarrow{\text{正規化}} -(0.1011)_2 \times 10^{-3}$

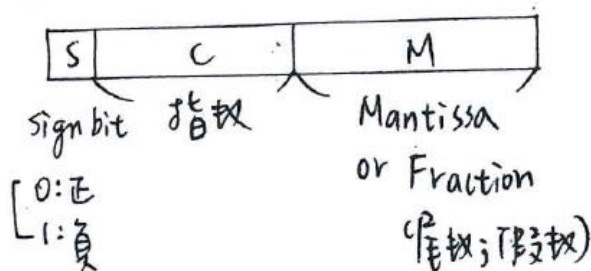
(2) $+(45.73)_8 \Rightarrow +(0.4573)_8 \times 8^2$

(3) $-(0.05A7)_{16} \Rightarrow -(0.5A7)_{16} \times 16^{-1}$

3. 一般浮點數(Floating Point Number)格式

(一) 先正規化, 得到 $\pm(0.F)_b \times b^{\text{exp}}$

存入 Computer 內部之格式如下:



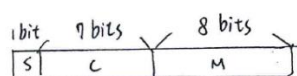
C: 指數

一般採用

[2's complement
or
Excess-Code
表示正負指數]

- 換成 2 進位, 再換成 2 補數(下列例子), 再換成對應之進位系統, Ex: 16 進位等等

(二) 例: 浮點數格式如下

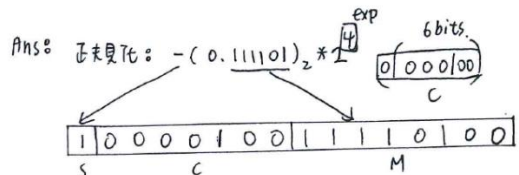


採 2 進制, C 用 2's 補數

表示負指數值。

先化到該進制再正規化。

(1) $-(1111.01)_2$ 之存入內容。



(2) $+(0.00001011)_2$ 之存入內容

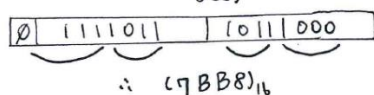
並以 16 進制 (Hex) 表示。

正規化: $+(0.1011)_2 \times 2^{-5}$

exp = -5 \Rightarrow

1	111011
---	--------

5 2's



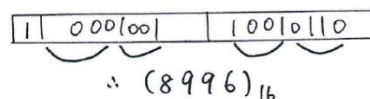
(3) $-(300)_{10}$ 之存入內容 (Hex)

$$-(300)_{10} = -(100101100)_2$$

正規化: $-(0.100101100)_2 \times 2^9$

exp = 9 \Rightarrow

0	001001
---	--------



(4) 內容 (CD)A8₁₆ 表值為?

$$\frac{1}{5} \underbrace{100110}_C \underbrace{10101000}_M$$

$$C = \frac{1}{5} \underbrace{100110}_5 \downarrow \text{2's 还原}$$

$$\therefore -(110011)_2 = -51$$

$$\therefore -(0.10101)_2 * 2^{-51}$$

(5) (AE97)₁₆

$$\frac{1}{5} \underbrace{10101110}_C \underbrace{10010111}_M$$

$$C = \frac{1}{5} \underbrace{10101110}_5 \downarrow$$

$$+ (32+14) = +46$$

$$\therefore -(0.10010111)_2 * 2^{46}$$

- 換成 2 進位，再轉成對應之 **Excess(+2ⁿ⁻¹)**，再換成對應之進位系統，
Ex: 16 進位等等

(1) $+(0.000001011)_2$

正規化: $+(0.1011)_2 * 2^{-5}$

0	011011	10110000
S	C	M
1	7 bits	8 bits

(依是負)

exp = -5 用超碼表示，因為 C 佔 7 bits.

\therefore Excess-64 code

$\therefore C = -5 + 64 = 63 = 4$

$$\begin{array}{r} 011111 \\ - 100 \\ \hline 011011 \\ \hline C \end{array}$$

(2) $+(23 \frac{7}{8})_{10}$ 之 儲存內容 (以 Hex 表示)

$\Rightarrow +(10111.111)_2$

exp = 5 $\therefore 5 + 64 = 69 = 10001011$
C

正規化: $+(0.10111111)_2 * 2^5$

0	10001011	10111111
---	----------	----------

$\therefore (45BF)_{16}$

(1) 不看正負(最左位元)正數最大與負數最小相同，反之

(-) 求 ① 最大值 (最大正數)

以數線來看.

② 最小正數

③ 最大負數

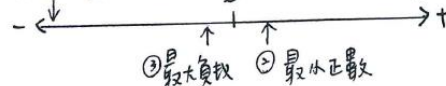
④ 最小值 (最小負數)

④ 最小值

③ 最大負數

② 最小正數

① 最大值



$|①| = |④|$, $|②| = |③|$

\therefore 只要求出 ④ 及 ②，自然得

$③ = -②$ 及 $④ = -①$

- (2) 且 2 補數與 excess 值域相同 $\rightarrow -(2^{n-1}) \sim 2^{n-1} - 1$
- (3) 以最大數為例 $0.11111111 * 2^{127} = (1 - 0.00000001) * 2^{127} = (1 - 2^{-7}) * 2^{127}$

(=) 不論 C 用 2's 補數 or Excess-code 表 EXP, 其 exp 之值域皆相同.

(三) 例 1 8 bits 9 bits 才 2 進制, C 用 2's 補數表示 EXP (or Excess-code)

求 ① 最大正數 ② 最小正數 ③ 最大負數 ④ 最小負數

Ans: 先求 exp 值域 $\because C$ 佔 8 bits $\therefore -2^{8-1} \leq x \leq + (2^{8-1} - 1) \Rightarrow \boxed{-128} \leq x \leq \boxed{+127}$

$$\textcircled{1} \text{ 最大正數 } + (0.1111111)_2 * 2^{+127} = + (1 - 2^{-7}) * 2^{127} = 2^{127} - 2^{120} = 2^{120} * (2^7 - 1) = 2^{120} * 127$$

$$\textcircled{2} \text{ 最小正數 } + (0.1)_2 * 2^{-128}$$

$$\textcircled{3} \text{ 最大負數 } = - (0.1)_2 * 2^{-128} = - \textcircled{2}$$

$$\textcircled{4} \text{ 最小負數 } = - (0.1111111)_2 * 2^{+127} = - \textcircled{1}$$

4. IEEE - 754 (國際格式) (有保留特殊值)

- 該標準之固定值為 $2^{Exp-1} - 1 \rightarrow$ 跟一般 excess 相比還要多減

1

- Mantissa 越多越精確

- 單精度 single precision 32 bits (紅點為小數點)

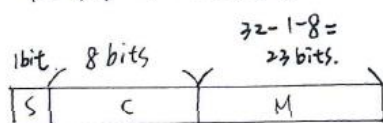
(1) Exp 部分加 $2^{8-1} - 1 = 127$

(2) Exp: 全 1 或全 0 為特殊值, 表示範

$$00000001(1) - 127 \leq \text{exp} \leq 11111110(254) - 127 \rightarrow \boxed{-126 \leq \text{exp} \leq}$$

127

(一) 單精度 (32bits) 格式



其中, C 是用 Excess-127 Code 表示指數.

Note: 比正常起碼少 1.

IEEE 754 之正負表示不同於一般型格式

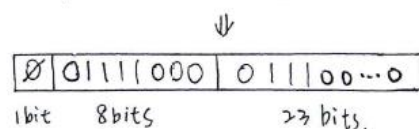
格式: $\pm (1.F)_2 \times 2^{\text{exp}}$

↳ 沒有規定小數點後第一個 bit 非 0

例: 以 IEEE 754 浮點數格式表示:

(1) $+(0.00000010111)_2$ 答案用 16 進制表現.

Ans: 正負表示: $+(1.0111)_2 \times 2^{-7}$ $\text{exp} = -7 \therefore C = -7 + 127 = 120 = (1111000)_2$



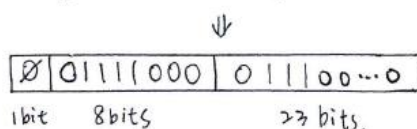
以 16 進制表示:

$\Rightarrow 0.01111000.0111000 \dots 0$ (4 個 0 代表 16 bits)
 $\therefore (3 \text{ C } 3 \text{ 8 } 0000)_{16}$

例: 以 IEEE 754 浮點數格式表示:

(1) $+(0.00000010111)_2$ 答案用 16 進制表現.

Ans: 正負表示: $+(1.0111)_2 \times 2^{-7}$ $\text{exp} = -7 \therefore C = -7 + 127 = 120 = (1111000)_2$



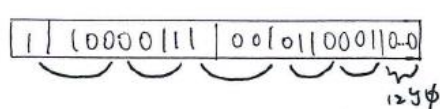
以 16 進制表示:

$\Rightarrow 0.01111000.0111000 \dots 0$ (4 個 0 代表 16 bits)
 $\therefore (3 \text{ C } 3 \text{ 8 } 0000)_{16}$

(2) $-(300 \frac{3}{8})_{10}$ 之內容?

Ans: $\Rightarrow -(100101100.011)_2$

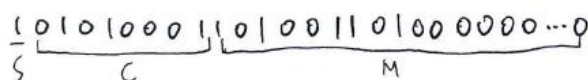
正負表示: $-(1.00101100011)_2 \times 2^8$ $\text{exp} = 8 \therefore C = 8 + 127 = 128 + 7 = (10000111)_2$



$\Rightarrow (C3963000)_{\text{Hex}}$

(3) $(A8D34000)_{\text{Hex}}$ 之值?

$\text{exp} = C - 127 = (01010001)_2 - 127$
 $= 5 \times 2^4 + 1 - 127 = -46$



$\therefore -(1.101001101)_2 \times 2^{-46}$

(4) (EDCA9000) Hex 之值?

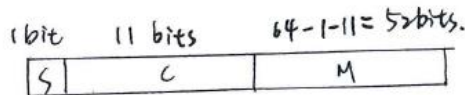
1101101100101010010...0
S C M

$$\begin{aligned} \text{exp} &= (11011011)_2 - 127 \\ &= (1011011)_2 - 128 + 1 \\ &= (1011011)_2 + 1 = 92 \\ \therefore &= (1.10010101001)_2 \times 2^{92} \end{aligned}$$

- 倍精度 double precision 64 bits

- (1) Exp 部分加 $2^{11-1} - 1 = 1023$
- (2) Exp: 全 1 或全 0 為特殊值，表示範圍: 0000000001(1) - $1023 \leq \text{exp} \leq 11111111110(2046) - 1023 \rightarrow -1022 \leq \text{exp} \leq +1023$

倍精度 (64bits) 格式



其中 C 是用 Excess-1023 code 表示指數。

- 特殊值規定

exp	mantissa	result
全 1	為 0	正無窮大
全 1	非 0	NaN(未定義或不可表示)
全 0	為 0	0
全 0	非 0	無法正規化

- 誤差種類

- (1) 固有(Inherent)誤差(Error)，
ex. 圓周率
- (2) 截尾誤差(Truncation Error)
- (3) 捨棄誤差(Rounding Error)，
ex. 四捨五入行程之誤差，
 $5/2=2$

