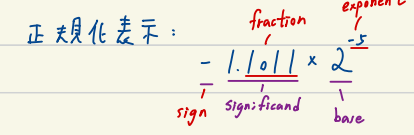


11. IEEE 754 表示法

A designer need to find a compromise between fraction & exponent

設以科學記號表示 - 進位小數

111. fraction 影响了可表示之數的 precision



121. exponent = 可表示之數的 range

∴ 字組大小是固定的, sign + fraction + exponent = fixed length

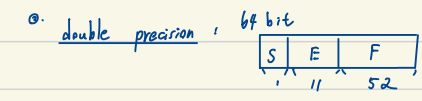
11. Good design demands good compromise

IEEE 754 有兩種精度表示法:

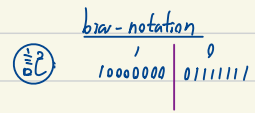
- ①. single precision
- ②. double precision



其中: S: sign, E: exponent, F: fraction



其中: E 欄位數值是利用 bias-notation 表示



Why? 因為在做浮點數加法時,會先比較 exponent 是否相同, 若不同,需要 shift 至相同才能做加法
∴ 需要大量 exponent 的比較, 又 bias-notation 的大小排序 和 unsign number 相同, 可使用無符號比較器來做比較 加速比較速度.

題目: 121. 將 10 進位小數轉至 IEEE 754

121. 給定 IEEE 754, 求 10 進位數

step: I. 先將 10 進位轉至 2 進位
II. 轉至正規化表示, 並以 IEEE 754 表示

eg $-13.8125_{10} \rightarrow -1101.1101_{12} \rightarrow -1.1011101 \times 2^3$

S: 1 E: 3 = 1000010 F: 10111000000000000000000

= 1 1000010 10111000000000000000000

$1.0000 \cdot 00 \times 2^1$

IEEE 754 之編碼情形 & 數值範圍

編碼情形共有 5 種:

- ①. ± 0 : S 000...0, E 0...0, F 0...0
- ②. $\pm \infty$: S 11...1, E 0...0, F 0...0
- ③. NaN: S 11...1, E 0...0, F non-zero
- ④. 正規化數: $\pm 1.x \times 2^E$
 $1 \leq E \leq 254$, F 為任何數
- ⑤. 非正規化數: $\pm 0.x \times 2^{1-B}$

數值範圍:

gradual underflow

(i). 非正規化數最小:
 $0.00...0100...01 \times 2^{-126} = 2^{-149}$

(ii). 非正規化數最大:
 $0.01...011...1 \times 2^{-126} = (1 - 2^{-23}) \times 2^{-126}$

(iii). 正規化數最小:
 $1.0...010...0 \times 2^{-126} = 2^{-126} = 2 \cdot 2^{-127} = 2 \cdot 10^{-38}$

(iv). 正規化數最大:
 $1.1...11 \times 2^{127} = (2 - 2^{-23}) \times 2^{127} \approx 2 \cdot 10^{38}$

其中: $1.1...1 \times 2^{127} = (2 - 2^{-23}) \times 2^{127} \approx 2 \cdot 10^{38}$

例: P217 3 題 已知 single precision

111. $-13.5_{10} \rightarrow -1101.1_{12} \rightarrow -1.101 \times 2^3$

$1 - 10000000 \rightarrow 10000000$

$3 - 10000010 \rightarrow 10000010$

121. 最大正規化數:

$0.111111101111...1 \rightarrow 1.111...1 \times 2^{127}$

$= (2 - 2^{-23}) \times 2^{127}$

131. 最小非正規化數:

$0.000...0100...01 \rightarrow 0.00...01 \times 2^{-126}$

$= 2^{-23} \times 2^{-126} = 2^{-149}$

141. E 全 1

- F 為非 0 = NaN
- F 為全 0 = $\pm \infty$

P261 57 題: double = E 欄位 11 個 bit $0 \sim 2^{11}-1 \Rightarrow 0 \sim 2047$

0 = preserve for ± 0 , denormalized
2047 = preserve for $\pm \infty$, NaN
 $1 \sim 2046$ = normalized number.

∴ 一般數值運算根本用不到非正規化數那麼小的數
∴ MIPS 中不支援, 會為 underflow
∴ 整數不含有 underflow 問題
課本無列到 denormalized number
∴ 小於 2^{-126} 的數皆視為 underflow

補

IBM format:

IBM format 為 excess-64

以 16 為底且 exponent 大小為 7 bit

$\therefore 1000000 \rightarrow 0$

\therefore 為 base 16, 無法保證 significand 小數點左邊必為 1

\therefore 設為 0, 正規化數為: $\pm 0.x \times \dots \times 16^E$

eg $-938.8125 \rightarrow -1110101010.1101 \times 2^0$

base 16 $\rightarrow -3AA.D \times 16^0$

Normalize: $-0.3AAD \times 16^3$

$\text{又 } 1000000 = 0 \therefore 3 = 1000011$

\Rightarrow IBM format: $\begin{array}{c} | 1000011 001101010110100 \dots 0 \\ \hline 1 \quad 7 \quad 24 \end{array}$