

2004/2005 (Fiumi)

①

- Note alla rappresentazione in virgola mobile:
(De Fiumi, prefazione Depitole, pg. 50-52)

Un numero in virgola mobile viene rappresentato nel formato

$$x = \pm m \cdot 2^{\pm E}$$

dove m è la mantissa
 E è l'esponente.

- Rappresentazione della mantissa:

- si rappresenta in assoluto e segno
- si rappresenta normalizzata, ossia, per la base 2, nel formato $1, \dots$, in modo da evitare di rappresentare l'1.

- inoltre, se ho $1,010$, è $1 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} + 0 \cdot 2^{-3}$
(De Feb. 2, 12). etc

(2)

• Rappresentazione dell'esponente.

• Per evitare il segno, Non si rappresenta

E ma $B + Bias$, dove $Bias$ è un valore di polarizzazione.

• Per cui si rappresenta come un N bit solo.

• ~~Formula~~ Definiamo $B + Bias = E_{eff}$ codificati su n bit

Il valore di

• $E_{eff} = 0$ indica: \rightarrow Numeri normalizzati.

• $E_{eff} = 2^m - 1$ indica \rightarrow Positive-infinito
 \rightarrow Negative-infinito
 \rightarrow N.o.N.

• Supponiamo di avere 11 bit.

I valori ~~da~~ rappresentabili sono

da $[0, 2047]$

• Gli estremi 0 sono censurati $\rightarrow 0$ \rightarrow per lo 0 e lo zero normalizzato
 $\rightarrow 2047$ \rightarrow per l'infinito e N.o.N.

• Restano 2046 Valori.

(5)

i 2 valori in esponente di potenza all'espone negativo,
 dunque con 2046 valori si rappresenta
 l'intervallo $[-1023, 1023]$ che contiene appunto
 2046 valori.

\Rightarrow allora, dato un esponente $E \in [-1023, 1023]$
 che numero devo sommare ad E per ottenere
 un numero nell'intervallo $[1, 2046]$,
 in quanto ϕ e 2047 sono compresi?

Risposta: 1023. Questo è il BIAS
 per i double.

• ~~Q~~ RAPPRESENTAZIONE DI ϕ e numeri
 Denormalizzati:

- se $B_{eff} = \phi$, $m = \phi$ $\bar{x} \pm \phi$ ~~(-1)~~ $\frac{1}{2} \phi$
- se $B_{eff} = \phi$ e $m \neq 0$ è un caso denormalizzato
 (vedi testo per il
 significato)
- se $B_{eff} = 2^n - 1$ e $m = \phi$ è ~~il~~ Positive \rightarrow Inf
 Negative \rightarrow Inf

• $x \in \mathbb{B} = 2^m - 1$ e $m \neq 0$ allora \bar{x} N. a. N. (4)

• $x \in 0 < \mathbb{B}_{\text{eff}} < 2^m - 1$

allora il numero \bar{x}

$$(-1)^s \cdot (1, m) \cdot 2^{\mathbb{B}_{\text{eff}} - \text{Bias}}$$

ovvero $\text{Bias} = 2^{m-1} - 1$.

• float STANDARD 7828 754

Bit della mantissa: 23

Bit dell'esponente: 8

$$\text{Bias} = 2^{8-1} - 1 = 127$$

• esponente minimo rappresentabile: -126

• esponente max rappresentabile: 127

Double standard IEEE 754

(5)

Bit della mantissa: 52

Bit dell'espone: 11

$$Ries = 2^{11-1} = 1023$$

- espone MAX = $Ries = 1023$
- espone MIN = $-Ries + 1 = -1022$

Da qui otterremo i valori per il double.

RICAPITOLAZIONE DELLO STANDARD IEEE 754 (chiuso)

51

• FLOAT (NUMERI A PRECISIONE SINGOLA)

• rappresentazione su 32 bit

- 1 segno
- 8 esponenti
- 23 mantissa

• intervalli di rappresentazione

$$[-10^{38}, -10^{-38}] \cup [10^{-38}, 10^{38}]$$

• precisione a 7 cifre decimali

• DOUBLE (NUMERI A PRECISIONE DOPIA)

• rappresentazione su 64 bit

- 1 segno
- 11 esponenti
- 52 mantissa

• intervalli di rappresentazione

$$[-10^{308}, -10^{-308}] \cup [10^{-308}, 10^{308}]$$

• precisione a 17 cifre decimali significative.

• Costituito di $N \times M$, positive-infinita, negative-infinita. ⑥

• $m \neq 0$ indicazione

• Se $\text{Eff} = z^m - 1 \Rightarrow$ $m=0$; $S = \emptyset$ indicazione positive-infinita

• $m = \emptyset$; $S = 1$ indicazione negative-infinita

• Se $\text{Eff} = \beta$ e $m \neq \beta \Rightarrow$ non determinabile

• Se $\text{Eff} = \beta$ e $m = \beta \rightarrow S = \emptyset \Rightarrow$ indicazione ?

$\rightarrow S = 1 \Rightarrow$ indicazione 0

(Da INTERACT)

• Per il calcolo dell'esperto affettivo, esclusi i valori $\beta \in z^m - 1$, vale la relazione

$$\text{Eff} = z + (z^m - 1)$$

\uparrow ← z
 esperto di $\pm 1. m 2$

(7)

Calcolo delle mantisse:

$$\cdot \text{cio } m = 1, a_1, a_2, a_3, a_4, a_5, \dots$$

$$\begin{aligned} \text{Significato: } m &= 1 \cdot 2^0 + a_1 \cdot 2^{-1} + a_2 \cdot 2^{-2} + \dots \\ &= 1 \cdot 2^0 + \frac{a_1}{2} + \frac{a_2}{4} + \dots \end{aligned}$$

$$\text{da cui } m - 1 = \frac{a_1}{2} + \frac{a_2}{4} + \dots$$

$$(m-1) \cdot 2 = a_1 + \frac{a_2}{2} + \frac{a_3}{4} + \dots$$

$$\text{dunque se } (m-1) \cdot 2 = 0 \rightarrow a_1 = 0$$

$$\text{e } (m-1) \cdot 2 = 1 \rightarrow a_1 = 1$$

$$\text{per cui } a_1 \in \{0, 1\}$$

e così via

(vedi libro sui DSP)

Esempi:

- Sono 8 i bit per l'esponente
- Sono 3 i bit per la mantissa.

la costante da scendere è $2^{-1} = 2^{-4} = 127$

per cui la corrispondenza fra Eff ed E è la seguente

Eff	E
0	nessuna nessuna (eff = 0 indica ±0)
1	-126
2	-125
3	
⋮	
253	123
254	127
255	nessuna (eff = 255 indica \pm infinito)

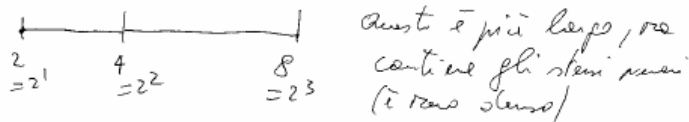
Max positive infinite
 Max negative infinite

• l'intervallo per gli esponenti vale sempre $[-126, +127]$

• si fa asimmetrico per la densità:

i numeri più piccoli sono "meno densi"

• se $m = 3$ bit, ho 8 numeri nell'intervallo $[2, 4]$,
ma anche $[4, 8]$.



- ovviamente, sia m che Bff vengono rappresentati in binario. (9)

Esempio:

- Calcolare la rappresentazione fp di 123, con $m=3$ e $e=5$

- PASSO 1: occorre normalizzare il numero, cioè rappresentarlo in forma ~~normale~~ $1.m \cdot 2^e$
Per normalizzare si moltiplica e divide per 2.

$$123 = 61,5 \cdot 2 = 30,75 \cdot 2^2 = 15,375 \cdot 2^3 = 7,6875 \cdot 2^4$$

$$= 3,84375 \cdot 2^5 = 1,921875 \cdot 2^6$$

Quindi $1,921875 \cdot 2^6$ è la rappresentazione normalizzata di 123

- PASSO 2: calcolo di Bff

~~eff~~ ~~calcolo~~
• l'esponente di $1,921875 \cdot 2^6$ è 6

$$\bullet Bff = 6 + 127 = 133$$

$$\bullet Eff = \underset{12 \text{ bit}}{10000101} \quad (133)_{10}$$

- PASSO 3: calcolo di m

\bar{e} 4,921875 $\Rightarrow Q_0 = 1$ (non si può per lo standard IEEE, non si rappresenta). (10)

$$4,921875 - Q_0 = 0,921875$$

$$0,921875 \times 2 = 1,84375 \Rightarrow Q_{-1} = 1$$

$$\bullet 1,84375 - Q_{-1} = 0,84375$$

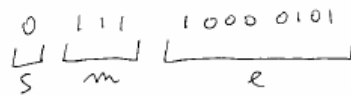
$$0,84375 \times 2 = 1,6875 \Rightarrow Q_{-2} = 1$$

$$0,6875 \times 2 = 1,375 \Rightarrow Q_{-3} = 1$$

\Rightarrow Ci si ferma, perché ~~stessa~~ lo resto non è a 3 cifre.

Devo ~~in~~ rappresentare in fp, con $m=3$, ol
 $e=8$

123 e^{-}



- Osservazione: qual è l'errore che si commette nel rappresentare 123 in fp con i dati che sono?

Il numero che si rappresenta è

$$4.111|_2 \cdot 2^6 = \left(1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8}\right) \cdot 64$$

$$= 64 + 32 + 16 + 8 = 120$$

• Dunque si commette un errore di 3

\Rightarrow per aumentare la precisione, occorre aumentare i bit delle mantisse (per punti, in IEEE 754, sono molti i bit riservati alle mantisse).

• Esercizi:

1) calcolare l'errore commesso nel rappresentare 123 in fp, con $e=8$ ed $m=3, m=4, m=7$

2) calcolare la rappresentazione fp, con $e=8$ ed $m=3$, di:

1.0

0,027

423,5

$3,1 \cdot 10^{-4}$

calcolando anche l'errore commesso.