

## ARITMETICA di MACCHINA

$$x = 123456.789$$

$\begin{array}{c} \nearrow 10^1 \\ \uparrow 10^0 \\ \searrow 10^{-1} \end{array}$ 
 $\nearrow 10^{-2}$

notazione posizionale

forma esponenziale

$$x = 0.\underbrace{123456789}_{\substack{\text{mantissa} \\ \in \mathbb{Z}}} \cdot \underbrace{10^6}_{\substack{\text{base di rapp.} \\ \text{esponente della base}}}$$

$$x = 1.23456789 \cdot 10^5 \quad \text{2}^{\text{forma esponenziale}}$$

$$\| x = (-1)^s 0.a_1 a_2 \dots a_t \cdot \beta^e \| \quad \begin{array}{l} \text{forma} \\ \text{esponenziale} \end{array}$$

con  $s \in \{0, 1\}$      $a_i \in \{0, 1, \dots, \beta-1\}$

$\beta =$  base di rappresentazione

$t = n^{\circ}$  cifre mantissa

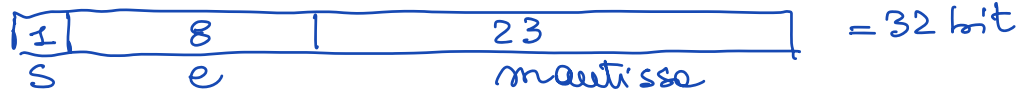
$e \in \mathbb{Z}$  esponente della base

## Sistemi FLOATING POINT

$$F(\beta, t, L, U) = \{ x = (-1)^s 0.a_1 a_2 \dots a_t \cdot \beta^e \mid \begin{array}{l} s \in \{0, 1\}, a_i \in \{0, 1, \dots, \beta-1\} \text{ per } i=1, \dots, t, \underline{a_1 \neq 0} \\ L < 0, U > 0, L \leq e \leq U \cup \{0\} \end{array} \}$$

## ☘ sistema F semplice precisione

$\beta = 2$  registro di 4 Byte



float in C, C++

Se  $a_1 \neq 0 \Rightarrow a_1 = 1$

23 bit per m, ma  $t = 24$

L, U

## ☘ sistema doppia precisione

$\beta = 2$  8 Byte = 64 bit double in C, C++

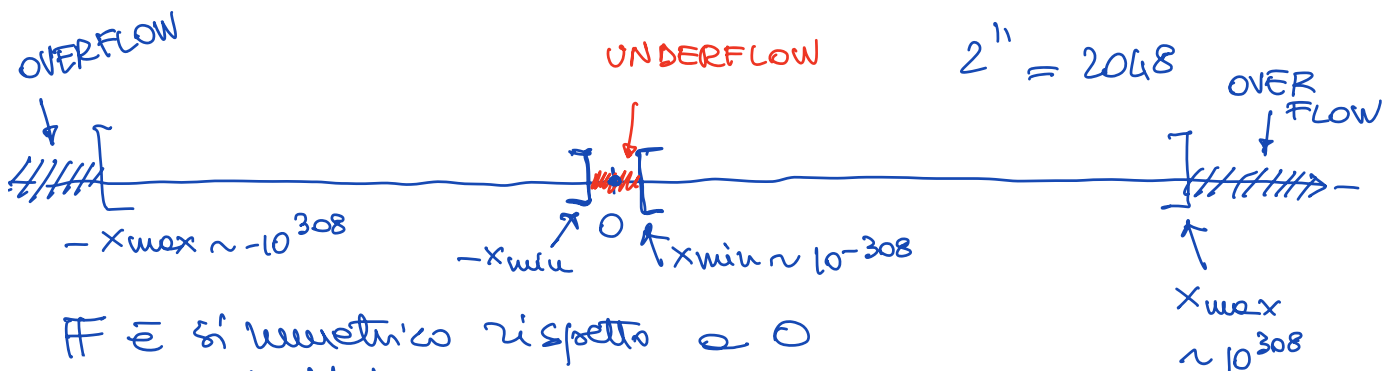


.52 bit per m  $\Rightarrow t = 53$

$L = \underline{\underline{-1021}}$

$U = 1024$

$1021 + 1 + 1024 \neq 2046$



F è simmetrico rispetto a 0

è limitato

è finito (cardinalità finita)

$$\text{card}(F) = 2 \beta^{t-1} (\beta - 1) (t - L + U)$$

i floating point non sono equamente distribuiti

Esempio

$$\mathbb{F}(\beta=2, t=3, L=-2, U=3)$$

$$x = (-1)^s \cdot 0.a_1 a_2 a_3 \cdot \beta^e$$

$-2 \leq e \leq 3$

-leggo in base 10

$$x_{\min} = \underbrace{0.100}_{\text{base 2}} \cdot 2^{-2} = \frac{1}{2} \cdot \frac{1}{4} = \frac{1}{8} = \frac{4}{32}$$

$$\underline{0.101} \cdot 2^{-2} = \left(\frac{1}{2} + \frac{1}{8}\right) \cdot \frac{1}{4} = \frac{5}{8} \cdot \frac{1}{4} = \frac{5}{32}$$

$$0.110 \cdot 2^{-2} = \left(\frac{1}{2} + \frac{1}{4}\right) \cdot \frac{1}{4} = \frac{3}{4} \cdot \frac{1}{4} = \frac{3}{16} = \frac{6}{32}$$

$$0.111 \cdot 2^{-2} = \left(\frac{1}{2} + \frac{1}{4} + \frac{1}{8}\right) \cdot \frac{1}{4} = \frac{7}{8} \cdot \frac{1}{4} = \frac{7}{32}$$

---


$$\underline{0.100} \cdot 2^{-1} = \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4} = \frac{8}{32} = \frac{4}{16}$$

$$\underline{0.101} \cdot 2^{-1} = \frac{5}{8} \cdot \frac{1}{2} = \frac{5}{16}$$

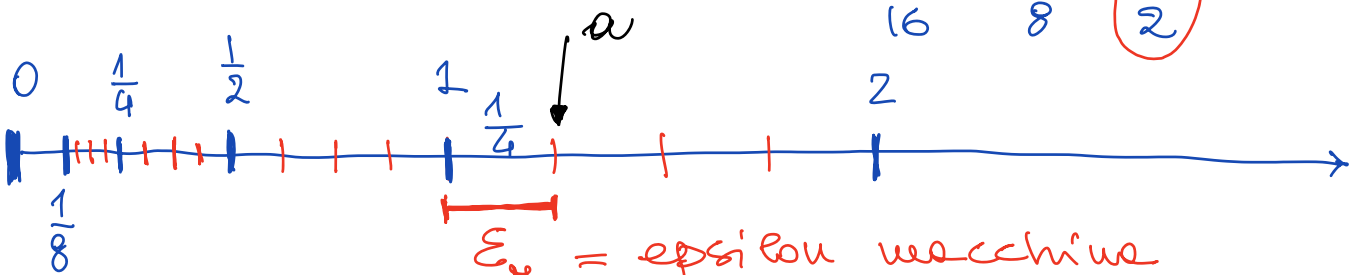
$$\frac{6}{16}$$

$$\frac{7}{16}$$

---


$$0.100 \cdot 2^0$$

$$\frac{8}{16} = \frac{4}{8} = \frac{1}{2}$$



$\epsilon_H = \text{epsilon macchina}$

precisione di macchina

$\epsilon$  è la distanza tra 1 e il + numero di  $\mathbb{F}$  maggiore di 1.

26/09/2023

$$1 = 0.100 \cdot \beta^1$$

$$a = 0.101 \cdot \beta^1$$

$$\varepsilon_M = a - 1 = 0.001 \cdot \beta^1 = \beta^{-t} \cdot \beta^1 = \beta^{1-t}$$

↑  
posiz. t

$$\boxed{\varepsilon_M = \beta^{1-t}}$$

nell'esempio avevo  $\beta=2$  e  $t=3 \Rightarrow \varepsilon_M = 2^{-2} = \frac{1}{4}$

in semplice precisione  $\varepsilon_M = 2^{1-24} \sim 10^{-7}$

in doppia precisione  $\varepsilon_M = 2^{1-53} = 2^{-52} \sim 2.22 \cdot 10^{-16}$

(eps)

$$x \in \mathbb{R}, \quad \boxed{1 < x < 1 + \varepsilon_M} \Rightarrow x \notin \mathbb{F}$$

$fl_t(x) = 1$



A horizontal number line with tick marks at 1 and 1+ε<sub>M</sub>. A point x is marked between 1 and 1+ε<sub>M</sub>. A red arrow points from x down to the tick mark at 1.

$fl_t(x)$  è l'arrotondamento di  $x$  in  $\mathbb{F}$

$$|x - fl_t(x)| \leq \frac{1}{2} \varepsilon_M$$

$u = \frac{1}{2} \varepsilon_M$   
unità di arrotondamento  
(roundoff unit)

Più in generale :

$$|x - fl_t(x)| \leq \frac{1}{2} \epsilon_H \cdot |x|$$

se  $x \neq 0$

$$\frac{|x - fl_t(x)|}{|x|} \leq \frac{1}{2} \epsilon_H$$

errore di arrotondamento relativo

Proprietà dell'aritmetica classica  
che cadono in  $\mathbb{F}$

1)  $\nexists !$  lo zero come elemento neutro  
della somma

ciò  $\exists x \in \mathbb{F}$  con  $x \neq 0$  :  $1+x=1$

es:  $x = \beta^{-2t} = 2^{-106} \in \mathbb{F}$   
( $t=53$ )

$$1 = 0.1 \dots \dots \overset{\boxed{t}}{0} \cdot \beta^1$$

$$x = 0.0 \dots \dots \dots 0 \cdot \beta^{-2t+1}$$

allineando gli esponenti ho:

$$1 = 0.1 \dots \dots \dots 0 \Big| \cdot \beta^1$$

$$x = 0.0 \dots \dots \dots 0 \Big| \dots \dots \dots \boxed{1} \cdot \beta^1$$

$$1+x = 0.1 \dots \dots \dots 0 \cdot \beta^1$$

si perde nel fare la somma

t-sime cifre

2t-sime cifre

2) In  $\mathbb{F}$  vale la prop associativa di + e -

$$\text{prop ass: } (a_1 + a_2) + a_3 = a_1 + (a_2 + a_3)$$

$$(a_1 \cdot a_2) \cdot a_3 = a_1 \cdot (a_2 \cdot a_3)$$

(es)

$$a_1 = 10^{308}$$

$$a_2 = 10^{308}$$

$$a_3 = -4 \cdot 10^{307}$$

$\mathbb{F}$  doppia precisione

$$x_{\max} \sim 1.7 \cdot 10^{308}$$

$$(a_1 + a_2) + a_3 = \underbrace{2 \cdot 10^{308}}_{> x_{\max}} - 4 \cdot 10^{307} = \text{NaN}$$

NaN

$$a_1 + (a_2 + a_3) = 10^{308} + \underbrace{(10^{308} - 4 \cdot 10^{307})}_{< x_{\max}}$$

Es 2

$$\frac{1 + (-1+x)}{x} \stackrel{x \neq 0}{=} \frac{(1-1)+x}{x} = 1 \quad \text{ne certo}$$

$\neq$

## Propagazione degli errori di arrotondamento

(+)

$$x, y \in \mathbb{R} \quad \bar{x} = \text{fl}_t(x) \quad \bar{y} = \text{fl}_t(y)$$

in  $\mathbb{R}$  ho  $x+y$   
in  $\mathbb{F}$  ho  $\text{fl}_t(\bar{x} + \bar{y})$

$$\frac{|\text{fl}_t(\bar{x} + \bar{y}) - (x+y)|}{|x+y|}$$

se  $\underline{x = 1 + 10^{-15}}$       $\underline{y = -1}$

la somma è un'operazione  
potenzialmente  
instabile

si riesce a dim

$$\leq \left( \frac{|x| + |y|}{|x+y|} + 1 \right) u$$

$$\sim \frac{2}{10^{-15}} + 1$$

$$\sim 2 \cdot 10^{15}$$

$$\leq 2 \cdot 10^{15} \cdot 10^{-16}$$

$$\sim \underline{\underline{2 \cdot 10^{-1}}}$$

$$\textcircled{\bullet} \quad \frac{|\text{fl}_t(\bar{x} \cdot \bar{y}) - (x \cdot y)|}{|x \cdot y|} \leq 3u$$

↑

Il prodotto è un'op. stabile,

e piccoli errori sui dati corrispondono  
piccoli errori nella soluzione