

17/09/24

ⁿ ARITMETICA di RACCHINA

$$x = 123456.789$$

$\uparrow \quad \uparrow \quad \uparrow$
 $10^3 \quad 10^0 \quad 10^{-1}$

notazione posizionale

$$x = \sum_{k=-m}^n a_k \cdot \beta^k$$

β - base (10)

a_k sono le cifre : $a_k \in \{0, 1, \dots, \beta-1\}$

Notazione scientifica o esponenziale

$$x = 0.123456789 \cdot 10^6 \quad \text{FLOATING POINT}$$

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

$$s \in \{0, 1\}$$

$$a_i \in \{0, 1, \dots, \beta-1\}, a_1 \neq 0$$

l'ci fpe delle mantissa

Mantissa $m = a_1 a_2 \dots a_t$

B = base di rappresentazione

e = esponente della base (ordine di grandezza)

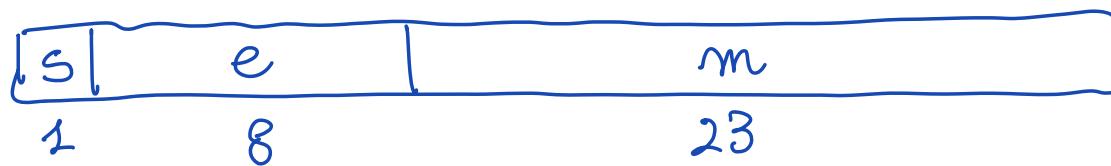
t = n° cifre della mantissa (precisione)

IEEE semplice precisione 4 Byte

doppia precisione 8 Byte

$B = 2$

4 Byte (semplice)



a 1.4^0 con 23 bit per la mantissa, ho in realtà

$t = 24$ (n° cifre in base 2)

$\sim 6 \circ 7$ cifre in base 10

$\beta = 2$, 8 Byte (doppia precisione)

| s | e | m |
|---|----|----|
| 1 | 11 | 52 |

$a_1 \neq 0 \Rightarrow a_1 = 1$, $t = 53$

$-1021 \leq e \leq 1024$ (2^e $\beta = 2$)

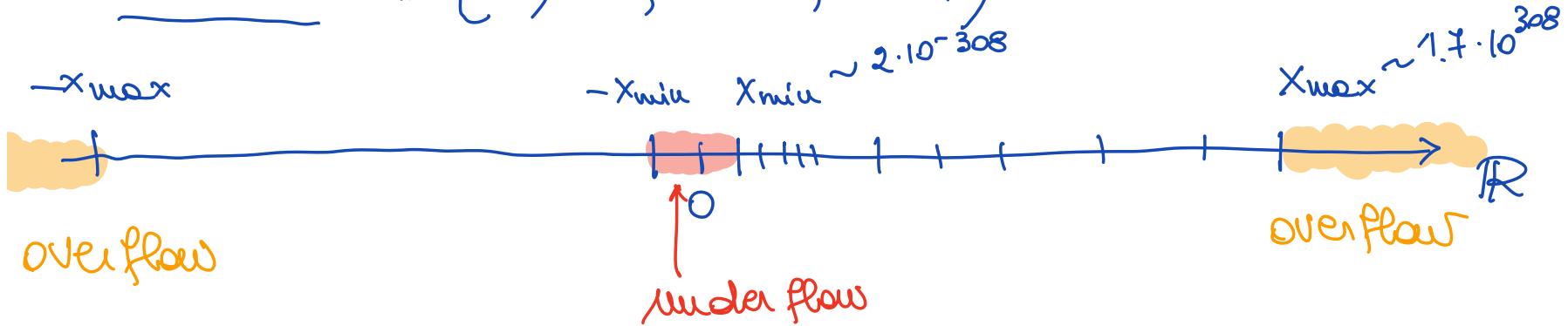
Sistemi floating point

$F(\beta, t, L, V) = \{x = (-1)^s \cdot 0.a_1 a_2 \dots a_t \cdot \beta^e, \text{ con:}$
 $s \in \{0, 1\}, a_i \in \{0, \dots, \beta - 1\}, a_1 \neq 0, L \leq e \leq V$

con $L < 0$, $U > 0 \} U \{ 0 \}$

MatLab

$$\mathbb{F}(2, 53, -1021, 1024)$$



\mathbb{F} è simmetrico rispetto a 0

è limitato ($\text{near } x = \pm \infty$)

è frutto (ha un n° finito di elementi)

i numeri fp non sono equispaziati, ma sono + densi verso lo 0 e meno densi man mano che cresce l'ordine di grandezza.

$$\text{Es: } F(2, 3, -2, 3) \\ \beta \quad t \quad L \quad U$$

$$x = (-1)^s \cdot 0.Q_1Q_2Q_3 \cdot 2^e$$

per simmetria prendo solo $s=0$

$$x_{\min} = \begin{cases} \text{base 2} \\ 0.100 \cdot 2^{-2} \end{cases} \text{ base 10} = 2^{-1} \cdot 2^{-2} = \frac{1}{8} = \frac{4}{32}$$

$$0.101 \cdot 2^{-2} = (2^{-1} + 2^{-3}) \cdot 2^{-2} = \left(\frac{1}{2} + \frac{1}{8}\right) \cdot \frac{1}{4} = \frac{5}{32}$$

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

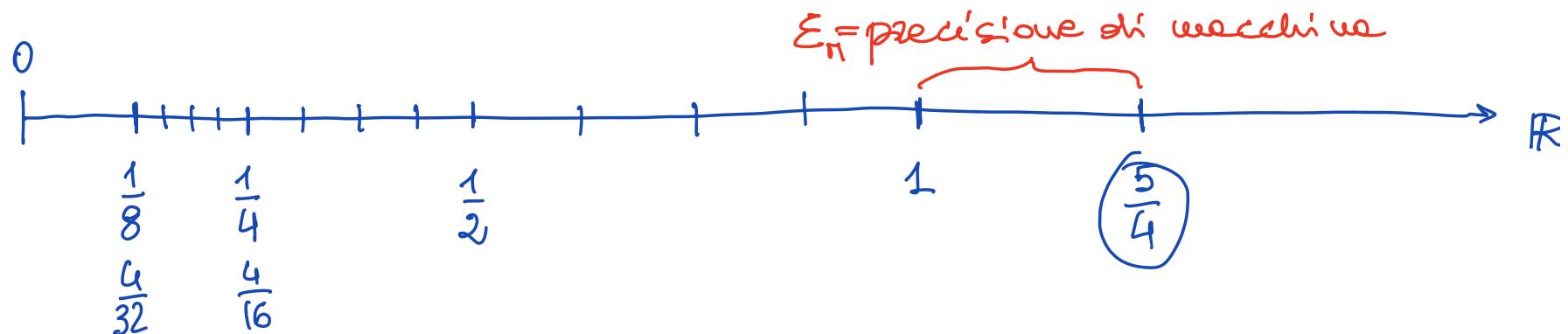
$$0.111 \cdot 2^{-2} = \dots - \dots = \frac{7}{32}$$

$$\begin{cases} 0.100 \cdot 2^{-1} \\ 0.101 \cdot 2^{-1} \\ 0.110 \cdot 2^{-1} \\ 0.111 \cdot 2^{-1} \end{cases} = 2^{-1} \cdot 2^{-1} = \frac{1}{4} \quad \frac{1}{4} = \frac{8}{32} = \frac{4}{16}$$

$$0.101 \cdot 2^{-1} = (2^{-1} + 2^{-3}) \cdot 2^{-1} = \left(\frac{1}{2} + \frac{1}{8}\right) \cdot \frac{1}{2} = \frac{5}{16}$$

$$0.110 \cdot 2^{-1} = \frac{6}{16}$$

$$0.111 \cdot 2^{-1} = \frac{7}{16}$$



$E_n = \text{distanza tra il } n^{\circ} 1 \text{ e il più piccolo numero floating point maggiore di } 1$

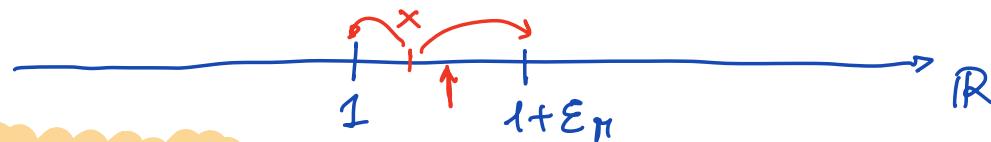
$$\frac{5}{4} = 0.101 \cdot 2^1$$

$$1 = 0.100 \cdot 2^1$$

$$E_n = \frac{0.001 \cdot 2^1}{B^{-t} \cdot 2^1} = B^{-t} \cdot B^1 = B^{1-t}$$

$E_n = B^{1-t}$

24/09



$$x \in \mathbb{R} : 1 < x < 1 + \varepsilon_H$$

Se $x \notin \mathbb{F} \Rightarrow x$ sarà approssimato dal n° fp più vicino a x

$fl_t(x)$ = rappresentante
di x nelle
macchine

$$\text{Si dimostra } |x - fl_t(x)| \leq \frac{1}{2} \varepsilon_H$$

Se $x \notin (1, 1 + \varepsilon_H)$ allora si ha

errore relativo
tra x e il
uso float

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

$$\mu = \frac{1}{2} \varepsilon_H$$

mità di
roundoff

MATLAB (double prec) $\beta=2, t=53 \Rightarrow \varepsilon_H \sim 2.22 \cdot 10^{-16}$

$$\mu \sim 1.11 \cdot 10^{-16}$$

Propriété dei codici in \mathbb{F}

1) Non vale più la proprietà associativa di + e .

es: doppie prec , $x_{\max} \sim 1.8 \cdot 10^{308}$

$$a = 10^{308}$$

$$b = 10^{308}$$

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

$$\underbrace{(a+b)+c}_{2 \cdot 10^{308} > x_{\max}}$$

Nan o Inf

\neq

$$\underbrace{a+(b+c)}_{1.6 \cdot 10^{308} < x_{\max}}_{6 \cdot 10^{307}}$$

Ese $x = 10^{-14}$

$$\frac{1 + (-1+x)}{x} = 0,9992$$

in weetlab
↑

$$\frac{(-1)+x}{x} = 1$$

in weetlab

2) lo zero, se fosse come elemento neutro delle somme,
non è unico

$$0 : x + 0 = 0 + x = x \quad \forall x \in \mathbb{R}$$

In \mathbb{F} , dato $x \in \mathbb{F}$ \exists molti numeri vescrivere $y \in \mathbb{F}$

$$\text{t.c. } x+y = x$$

Es $\mathbb{F}(10, 4, -5, 6)$ $\varepsilon_h = \beta^{t-t} = 10^{-3}$

$$x = 1 \quad y = 10^{-4} \quad ? x+y$$

$$\begin{array}{rcl} x=1 & = 0.1000 \cdot 10^1 & \\ y=10^{-4} & = 0.1000 \cdot 10^{-3} & \\ \hline & & 0.1000 \end{array} \quad \begin{array}{l} \cdot 10^1 \\ \cancel{1} \cdot 10^1 \\ \hline \end{array}$$

$$\Rightarrow x+y = 0.1000 \cdot 10^1 = 1 = x$$

però nel solvare
il si riscontra in un

regresso

Propagazione degli errori di arrotondamento

$$x \rightarrow \text{fl}_t(x)$$

1) Somme

dati $x, y \in \mathbb{R}$? $x+y$, in realtà calcolo $\bar{x}+\bar{y}$
 $\bar{x} = \text{fl}_t(x)$ $\bar{y} = \text{fl}_t(y)$ addì calcolo
addì calcolo
 $\text{fl}_t(\bar{x}+\bar{y})$

la somma di 2 numeri \mathbb{F} non è stessa che \mathbb{F}

es: $\mathbb{F}(10, 3, -2, 3)$

$$\begin{array}{r} \bar{x} = 0.123 \cdot 10^0 = 0.123 \\ \bar{y} = 0.456 \cdot 10^{-2} = 0.00456 \\ \hline \bar{x}+\bar{y} = 0.12756 \cdot 10^0 \\ \text{fl}_t(\bar{x}+\bar{y}) = 0.128 \cdot 10^0 \end{array}$$

Che errore commetto nel prendere come soluzione

$\text{fl}_t(\bar{x}+\bar{y})$ al posto di $(x+y)$

risultato
numerico

risultato
esatto

$c_{x,y}$

$$\left| \frac{(x+y) - f_{\text{ex}}(\bar{x}+\bar{y})}{|x+y|} \right| \leq \left(\frac{|x|+|y|}{|x+y|} + 1 \right) \cdot u$$

$u = \text{unità di roundoff}$ (adp $u \sim 1.11 \cdot 10^{-16}$)

Se $x \sim -y \Rightarrow |x+y| \sim 0, |x|+|y| > 0$

$\frac{|x|+|y|}{|x+y|}$ può diventare molto grande

(es)

$$x = 1 + 10^{-15}$$

$$y = -1$$

$$|x|+|y| = 2 + 10^{-15}$$

$$\frac{|x|+|y|}{|x+y|} \sim \frac{2}{10^{-15}} \sim 2 \cdot 10^{15}$$

$$|x+y| = 10^{-15}$$

$$\Rightarrow C_{xy} \cdot u \approx 2 \cdot 10^{15} \cdot 1.11 \cdot 10^{-16} \approx 2 \cdot 10^{-1}$$

A fronte di una unità di roundoff di 10^{-16} in una sola somma arrivo ad un errore nello zcl di circa 10^1 .

La **somma** è un'operazione **instabile**
perché a piccoli errori sui dati
grandi errori sulla soluzione.
possono covi' spandersi

Gli errori alle + sono detti errori di cancellazione

2) prodotto

$$x, y \in \mathbb{R}, \bar{x} = f_{lt}(x), \bar{y} = f_{lt}(y) \quad ? x \cdot y$$

trovo $f_{lt}(\bar{x} \cdot \bar{y})$ (numerico) $x, y \neq 0$

$$\text{err} = \frac{|(x \cdot y) - f_{lt}(\bar{x} \cdot \bar{y})|}{|x \cdot y|} \leq C \cdot u$$

$C \sim 3$ indip de x e y

il progetto è un'op. storica.