

Numeri in virgole mobili

Note Title

2026-02-18

"numeri di macchina"
"floating point"

Representazione scientifica in base β

Teo: fissata una base $\beta > 1$, ogni reale $\neq 0$ si scrive come

$$x = \pm \beta^p \left(\sum_{k=0}^{\infty} c_k \beta^{-k} \right)$$

$$[-764,8888... \rightarrow -10^2 \cdot 7,64888...]$$

c_k cifre in $\{0, 1, \dots, \beta-1\}$

Per avere una rappresentazione unica, supponiamo

$c_0 \neq 0$, e che c_k non sono tutti uguali a $\beta-1$ da un certo punto in poi

$c_0 \neq 0$ per avere unicità: $\underbrace{10^2 \cdot 7,6488...}_p = 10^3 \cdot \cancel{0,76488...}$
rappresentazione normalizzata $= 10^4 \cdot \cancel{0,076488...}$

c_k non tutti uguali a $\beta-1$ da un certo punto:

~~$0,49999... = 0,5$~~

$\underbrace{-10}_{\text{segno}} \cdot \underbrace{7,6488...}_{\text{mantissa}} \cdot \underbrace{2}_{\text{esponente}}$

Definizione: i numeri di macchina normalizzati sono i numeri della forma

$$c_k \in \{0, 1, \dots, \beta-1\}$$

$$x = \pm \beta^p \sum_{k=0}^{t-1} C_k \beta^{-k} \quad p \in \{m, m+1, \dots, M\} \quad \boxed{C_0 \neq 0}$$

L'insieme dei numeri fatti in questo modo si indica con $F(\beta, t, m, M)$

Esempio: scriviamo tutti i numeri di $F(10, 3, -2, 3)$

Partiamo a scrivere i numeri positivi con esponente $p = -2$

$$x = 10^{-2} \cdot (C_0 10^0 + C_1 10^{-1} + C_2 10^{-2}) = 0.0C_0C_1C_2$$

| $\overset{w}{\boxed{0.0100}}, 0.0101, 0.0102, \dots, 0.0999$

Poi con esponente $p = -1$: $x = 10^{-1} (C_0 10^0 + C_1 10^{-1} + C_2 10^{-2})$:

| $0.100, 0.101, 0.102, \dots, 0.999$

| $p=0$: $1.00, 1.01, 1.02, \dots, 9.99$

| $p=1$: $10.0, 10.1, 10.2, \dots, 99.9$

| $p=2$: $100, 101, 102, \dots, 999$

| $p=3$: $\boxed{1000}, \boxed{1010}, \boxed{1020}, \dots, \boxed{9990}^{\Omega}$

più tutti i loro negativi; questi sono tutti e soli gli elementi di $F(10, 3, -2, 3)$

• Numeri che non stanno in questo insieme andranno approssimati, ad es. $1001 \rightarrow 1000$

• I numeri non sono equispaziali:



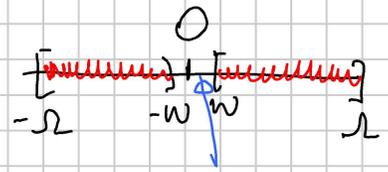
• esiste un numero γ più piccolo, $w = \beta^m \cdot 1.000\dots 0$

e un numero più grande $\Omega = \beta^M \cdot \underbrace{(\beta-1)(\beta-1)\dots(\beta-1)}_{\substack{\downarrow \\ t \text{ cifre uguali a } \beta-1}}$

• per passare da un numero di macchina v al suo successivo,

mi basta aggiungere 1 all'ultima cifra:
il successivo di

$$x = \beta^p (c_0 \beta^0 + c_1 \beta^{-1} + \dots + c_{t-1} \beta^{-(t-1)})$$



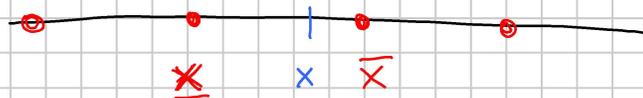
$$\tilde{x} = x + \beta^p \cdot \beta^{-(t-1)} = x + \beta^{p-t+1}$$

Teo: dato un numero reale $x \in [-\Omega, -w] \cup [w, \Omega]$,
esiste un numero di macchine \tilde{x} tale che

$$\left| \frac{\tilde{x} - x}{x} \right| < \beta^{-(t-1)}$$

errore relativo

Dim: possiamo assumere $x > 0$, l'altro caso è uguale, ma con
dei meno davanti. Possiamo assumere che x non sia
un numero di macchine.



Allora x sarà compreso tra due numeri di macchine
successivi: $\underline{x} < x < \bar{x}$

$$\left. \begin{array}{l} x = 10^{-6} \quad \tilde{x} = 10^{-2} \quad \text{errore: } \frac{10^{-2} - 10^{-6}}{10^{-6}} = \frac{10^{-6} \cdot (10^{4-6} - 1)}{10^{-6}} = 9999 \\ x = 10^{-6} \quad \tilde{x} = 0 \quad \text{errore: } \frac{0 - 10^{-6}}{10^{-6}} = -1 \end{array} \right\}$$

Errore relativo = rapporto tra l'errore assoluto $\tilde{x} - x$
e la quantità che voglio misurare x .

↳ (Tra l'altro, riusciamo anche a scrivere \underline{x} facilmente: è
quello che si ottiene troncando le cifre di x

$$x = \beta^p (c_0 + c_1 \beta^{-1} + c_2 \beta^{-2} + c_3 \beta^{-3} + \dots)$$

\tilde{x} si ottiene troncando la somma dopo t addendi.

\bar{x} è il num. di macchina successivo di x ,

$$\text{cioè, } \underline{x} + \beta^{p+1-t}$$



Sia prendendo \underline{x} (troncamento), sia \bar{x} (approssimazione per eccesso),

ho che

$$|\tilde{x} - x| \leq \bar{x} - \underline{x} = \beta^{p+1-t}$$

Devo stimare anche il denominatore della frazione:

$$x = \beta^p (c_0 \beta^0 + c_1 \beta^{-1} + c_2 \beta^{-2} + \dots) \geq \beta^p$$

↑
almeno 1

$$\frac{|\tilde{x} - x|}{x} \leq \frac{\beta^{p+1-t}}{\beta^p} = \beta^{-(t-1)} \quad \square$$

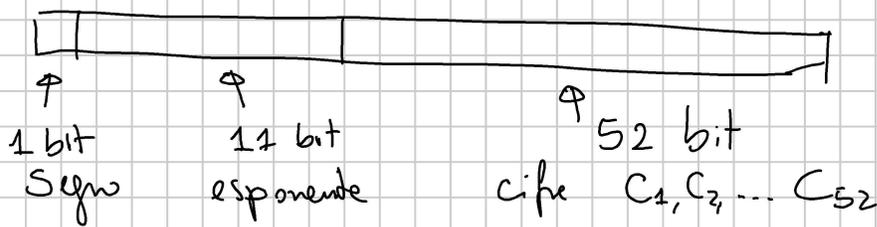
La costante $u = \beta^{-(t-1)}$ dipende solo dal sistema di numeri di macchine che sto usando, ed è detta precisione di macchine

Nel nostro esempio $\mathbb{F}(10, 3, -2, 3)$, $u = 10^{-2}$

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

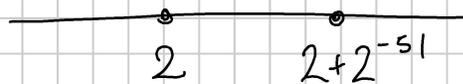
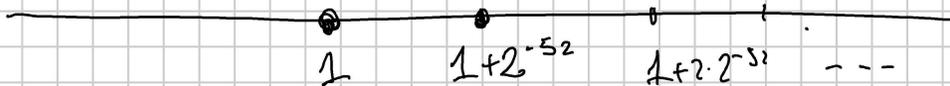
double / float64 / binary64 (IEEE)

$$\mathbb{F} \left(\underset{\beta}{2}, \underset{t}{53}, \underset{m}{-1022}, \underset{M}{1023} \right)$$



C_0 non si memorizza perché è sempre 1.

$$u = \beta^{-(p+1)} = 2^{-52} \approx 2.2 \cdot 10^{-16}$$



$$w \approx 10^{-308} \quad \Omega \approx 10^{308}$$

single/float/float32 : 32 bit, $u \approx 10^{-8}$

Dettagli: oltre a questi numeri, vogliamo rappresentare:

- 0
- ∞ (inf), $-\infty$, $-\emptyset$ ed es. per $\frac{1}{\infty} = 0$.
- NaN, "not a number", $0/0$, $\infty - \infty$
- "numeri denormalizzati" $\in (-w, w)$.

⚠ Somme, prodotti, ecc. di numeri di macchina in generale non sono numeri di macchina \rightarrow serve approssimare anche i risultati.

Si indica con un'operazione cerchiata, ad esempio

$$a \oplus b, a \ominus b, a \odot b, a \oslash b$$

quando effettuiamo un'operazione tra numeri di macchina e

