

Arrottondamento

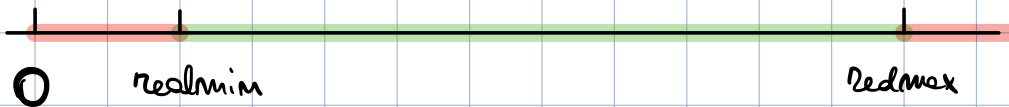
$$x \in \mathbb{R} \mapsto fl(x) \in \mathbb{F} \cup \{0, \pm Inf, \text{"denormali"}\}$$

Si legge "float di x " o
"arrottondamento di x "

numeri
machine
normali

gestione
overflow

gestione
underflow

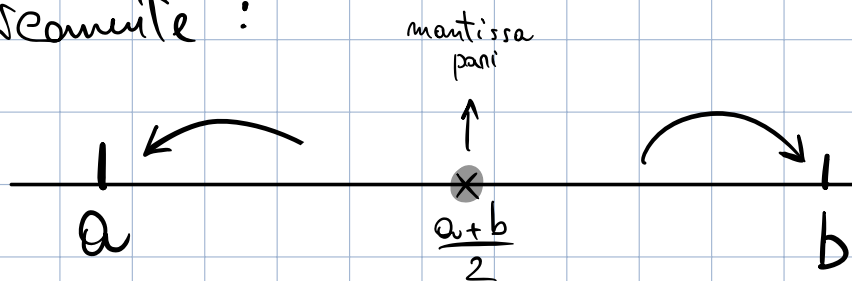


Supponiamo che l'arrottondamento non
coinvolga underflow o overflow,
ovvero $|x| \in [\text{redmin}, \text{redmax}]$

L'arrottondamento segue la regola

"RNTE": "Round To Nearest
, Ties to Even"

Graficamente:



dove $a < b \in \mathbb{F}$ consecutivi

Formalmente:

Siano $x \in [a, b]$, $a < b \in \mathbb{F}(\beta, t, n_1, n_2)$
consecutivi, β pari.

- se $x \in [a, \frac{a+b}{2})$, allora $fl(x) = a$
- se $x \in (\frac{a+b}{2}, b]$, allora $fl(x) = b$
- se $x = \frac{a+b}{2}$ e a ha mantissa pari,
allora $fl(x) = a$, altrimenti $fl(x) = b$.

Massimo errore relativo commesso
nell'arrotondamento:

$$x \in [\text{realmin}, \text{realmax}] \mapsto fl(x) \in \mathbb{F}$$

(Per semplicità, supponiamo $x > 0$)

errore relativo

$$\frac{|x - f(x)|}{|x|} \leq \frac{\frac{1}{2}(b-a)}{a} \leq \frac{\frac{1}{2}\beta^{p-t}}{\beta^p} = \frac{1}{2}\beta^{-t}$$

$a < b \in [\beta^p, \beta^{p+1}]$ numeri macchina consecutivi

Definizione : il numero β^{-t} viene detto "precisione di macchina" e indicato con eps (anche in Matlab).

Dunque il massimo errore relativo commesso nell'arrotondamento RNTF è la metà della precisione di macchina:

$$\frac{1}{2} \beta^{-t}$$

ESEMPLI:

IEEE single prec. : $\frac{1}{2} \beta^{-t} = 2^{-24} \approx 10^{-8}$

IEEE double prec. : $\frac{1}{2} \beta^{-t} = 2^{-53} \approx 10^{-16}$

cifre
corrette

Conditionamento dell'aritmetica di macchina

Consideriamo \tilde{x} approx di x ($x, \tilde{x} \in \mathbb{R}$) e definiamo l'errore relativo con segno:

$$\varepsilon_x := \frac{\tilde{x} - x}{x}$$

N.B. Sia $x > 0$. Se $\varepsilon_x > 0$, l'approx è per eccesso ($\tilde{x} > x$) altrimenti è per difetto.

Si ha:

$$\begin{aligned} x \varepsilon_x &= \tilde{x} - x \Leftrightarrow x + x \varepsilon_x = \tilde{x} \Leftrightarrow \\ \Leftrightarrow \tilde{x} &= x(1 + \varepsilon_x) \end{aligned}$$

Somma/sottrazione

Siano $x, y \in \mathbb{R}$. Consideriamo

$$(X, Y) \mapsto X + Y \in \mathbb{R} \leftarrow \text{"aritmetica esatta"}$$

$$(X, Y) \mapsto fl(fl(X) + fl(Y)) \in \mathbb{F}$$

↑
"aritmetica di macchine"
o "aritmetica finita"

Abbiamo:

$$fl(X) = X(1 + \varepsilon_X)$$

$$fl(Y) = Y(1 + \varepsilon_Y)$$

$$fl(fl(X) + fl(Y)) = (X + Y)(1 + \varepsilon_{X+Y})$$

Vogliamo mettere in relazione

$\varepsilon_X, \varepsilon_Y$ (errori rel. in input) con

ε_{X+Y} (errore rel. in output)

In aritmetica finite:

$$(X + Y)(1 + \varepsilon_{X+Y}) = X(1 + \varepsilon_X) + Y(1 + \varepsilon_Y)$$

$$\begin{aligned}
 \cancel{(x+y)} + \cancel{(x+y)} \varepsilon_{x+y} &= \\
 &= \cancel{x} + x \varepsilon_x + \cancel{y} + y \varepsilon_y
 \end{aligned}$$

$$(x+y) \varepsilon_{x+y} = x \varepsilon_x + y \varepsilon_y \quad \stackrel{x+y \neq 0}{\Leftrightarrow}$$

$$\varepsilon_{x+y} = \frac{x \varepsilon_x + y \varepsilon_y}{x+y}$$

$$|\varepsilon_{x+y}| = \frac{|x \varepsilon_x + y \varepsilon_y|}{|x+y|} \stackrel{\text{disug. triang.}}{\leq}$$

$$\leq \frac{|x| |\varepsilon_x| + |y| |\varepsilon_y|}{|x+y|} \leq$$

maggiore $|\varepsilon_x|$ e $|\varepsilon_y|$
con $\max\{|\varepsilon_x|, |\varepsilon_y|\}$
e metto in evidenza

$$\leq \frac{|x| + |y|}{|x+y|} \max\{|\varepsilon_x|, |\varepsilon_y|\}$$

Risultato:

numero di
evoluzione
 K

$$|\varepsilon_{x+y}| \leq \frac{|x| + |y|}{|x+y|} \max\{|\varepsilon_x|, |\varepsilon_y|\}$$

err. rel. in output

err. rel. in input

Osservazioni

(1) disug. Triang. $\Rightarrow K \geq 1$

(2) K grande se $|x+y|$ piccolo

Somma / sottrazione:

- ben condizionata se $K \approx 1$
- mal condizionata se $K \gg 1$

Osservazione

x, y hanno segno concorde \Rightarrow

$$\Rightarrow |x+y| = |x| + |y| \Rightarrow$$

$$\Rightarrow K = 1$$

Dunque:

La somma può essere mal condizionata solo se gli addendi hanno segno opposto e lo è se $y \approx -x$;
È tanto più mal condizionata quanto più y è vicino a $-x$.

Le perdite di precisione dovute alla somma $x+y$ quando

$y \approx -x$ è detto errore di cancellazione (di cifre...)

" x, y , vicini a essere uno l'opposto dell'altro"

Moltiplicazione

Siano $x, y \in \mathbb{R}$. Consideriamo:

$$(x, y) \mapsto xy \in \mathbb{R} \quad \text{"esatta"}$$

$$(x, y) \mapsto fl(fl(x) fl(y)) \quad \text{"di macchina"}$$

Abbiamo:

$$\cancel{(x, y)} (1 + \varepsilon_{xy}) = \cancel{x} (1 + \varepsilon_x) \cancel{y} (1 + \varepsilon_y)$$

e dunque

$$\cancel{1 + \varepsilon_{xy}} = (1 + \varepsilon_x)(1 + \varepsilon_y) = \cancel{1 + \varepsilon_x + \varepsilon_y + \varepsilon_x \varepsilon_y}$$

trascurabile rispetto
a $\varepsilon_x, \varepsilon_y$, se
 $|\varepsilon_x|, |\varepsilon_y|$ sono molto
piccoli (ipotesi
naturale)

Dunque, trascurando l'ultimo
addendo:

$$\varepsilon_{xy} \approx \varepsilon_x + \varepsilon_y$$

Infine

$$\begin{aligned} |\varepsilon_{xy}| &\approx |\varepsilon_x + \varepsilon_y| \stackrel{\text{disug. } \Delta}{\leq} \\ &\leq |\varepsilon_x| + |\varepsilon_y| \leq \\ &\leq 2 \max \{|\varepsilon_x|, |\varepsilon_y|\} \end{aligned}$$

allora

$$|\varepsilon_{xy}| \leq 2 \max \{|\varepsilon_x|, |\varepsilon_y|\}$$

↑ err. rel. output ↑ Numero di cond. K ↑ err. rel. input

l'operazione è sempre ben condizionata

Divisione

Siano $x, y \in \mathbb{R}$, $y \neq 0$. Consideriamo:

$$(x, y) \mapsto \frac{x}{y} \in \mathbb{R} \quad \text{"esatta"}$$

$$(x, y) \mapsto \text{fl}\left(\frac{\text{fl}(x)}{\text{fl}(y)}\right) \quad \text{"di macchina"}$$

$$\cancel{\frac{x}{y}} (1 + \varepsilon_{x/y}) = \frac{\cancel{x} (1 + \varepsilon_x)}{\cancel{y} (1 + \varepsilon_y)}$$

$$1 + \varepsilon_{x/y} = \frac{1 + \varepsilon_x}{1 + \varepsilon_y} \approx \dots$$

Richiamo (serie geometrica)



$$1 + x + x^2 + x^3 + \dots = \sum_{k=0}^{\infty} x^k = \frac{1}{1-x}, \quad |x| < 1$$

per x molto piccolo,
 x^2, x^3, \dots sono trascurabili
rispetto a x . Per cui

$$\frac{1}{1-x} \approx 1 \pm x \quad \text{per } x \text{ piccolo}$$

$$\begin{aligned}
 \dots &\approx (1 + \varepsilon_x)(1 - \varepsilon_y) = \\
 &= 1 + \varepsilon_x - \varepsilon_y - \varepsilon_x \varepsilon_y \approx \\
 &\approx 1 + \varepsilon_x - \varepsilon_y \quad (\text{per } |\varepsilon_x|, |\varepsilon_y| \\
 &\quad \text{piccoli}) ;
 \end{aligned}$$

Altrimenti

$$\cancel{1 + \varepsilon_{x/y}} \approx \cancel{1 + \varepsilon_x - \varepsilon_y}$$

$$\begin{aligned}
 |\varepsilon_{x/y}| &\leq |\varepsilon_x| + |\varepsilon_y| \leq \\
 &\leq \underset{\substack{\uparrow \\ K}}{2} \max \{|\varepsilon_x|, |\varepsilon_y|\}
 \end{aligned}$$

Quindi l'operazione è sempre
ben condizionata

Ricapitoliamo :

Somme $X+Y$	mod cond. se $Y \approx -X$
moltiplicazione	ben cond.
divisione	ben cond.