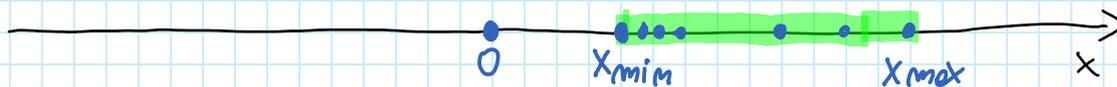


$$F(\beta, t, L, U) = \{0\} \cup \left\{ x \in \mathbb{R} \mid x = \pm \beta^p \sum_{k=1}^t d_k \beta^{-k} \right\}$$

con  $L \leq p \leq U$   
 $d_k \in \{0, \dots, \beta-1\}$ ,  $d_1 \geq 1$ .



VEDIAMO COME RAPPRESENTARE UN NUMERO  $x \in \mathbb{R}$   
 NELL'INSIEME  $F$

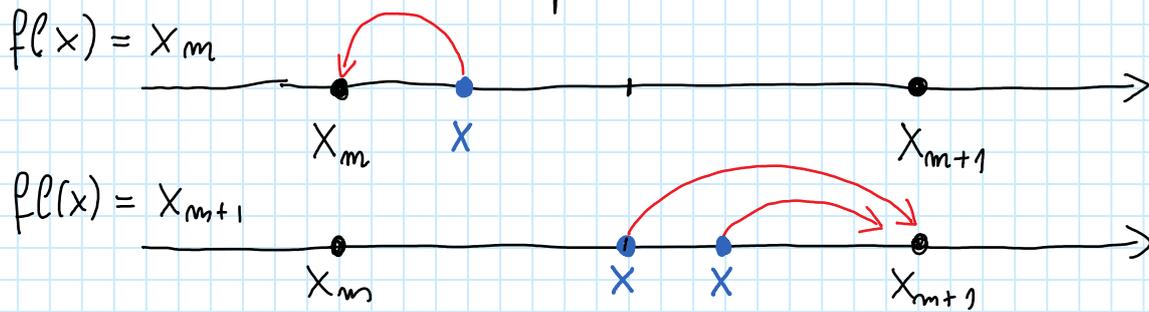
Per semplicità, consideriamo  $x \geq 0$ .

Si hanno questi casi:

- 1)  $x > X_{max}$  : viene generato un ERRORE di overflow. Il calcolo si arresta.
- 2)  $0 < x < X_{min}$  : viene generata una segnalazione di UNDERFLOW. Il calcolo, di solito, prosegue.
- 3)  $X_{min} \leq x \leq X_{max}$  : in questo caso  $x$  ha una rappresentazione, indicata con  $fl(x)$ , all'interno di  $F$ . Ho due possibilità:
  - 3.1)  $x \in \mathbb{R}$  è anche un numero di macchina, ossia di  $F(\beta, t, L, U)$ . In questo caso, è rappresentato esattamente (cioè, senza errore) in  $F$ :  

$$fl(x) = x$$
  - 3.2)  $x \notin F(\beta, t, L, U)$ . Il rappresentante di  $x$  nell'insieme  $F$  viene ottenuto per arrotondamento verso il numero di  $F$ .

più vicino a  $x$ :



dove  $X_m$  e  $X_{m+1}$  sono due numeri di macchina in  $\mathbb{F}(\beta, t, L, U)$  consecutivi.

Nel processo appena descritto viene commesso un errore  $|fl(x) - x|$  detto **errore di arrotondamento**

**PROP.** Sia  $x \in \mathbb{R}$  che NON dia né underflow né overflow. Allora, vale la seguente maggiorazione

$$\left| \frac{fl(x) - x}{x} \right| \leq \frac{1}{2} \beta^{1-t}$$

usando  $\mathbb{F}(\beta, t, L, U)$ .

Dim. Sia  $x \in \mathbb{R}$  della forma (per semplicità,  $x > 0$ )

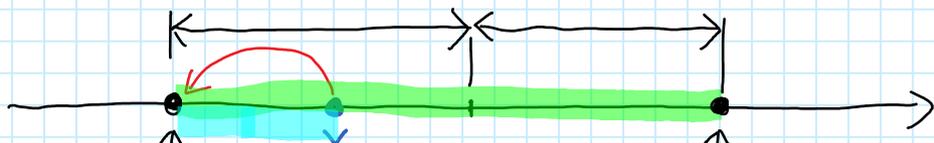
$$x = \beta^p \sum_{k=1}^{+\infty} d_k \beta^{-k}$$

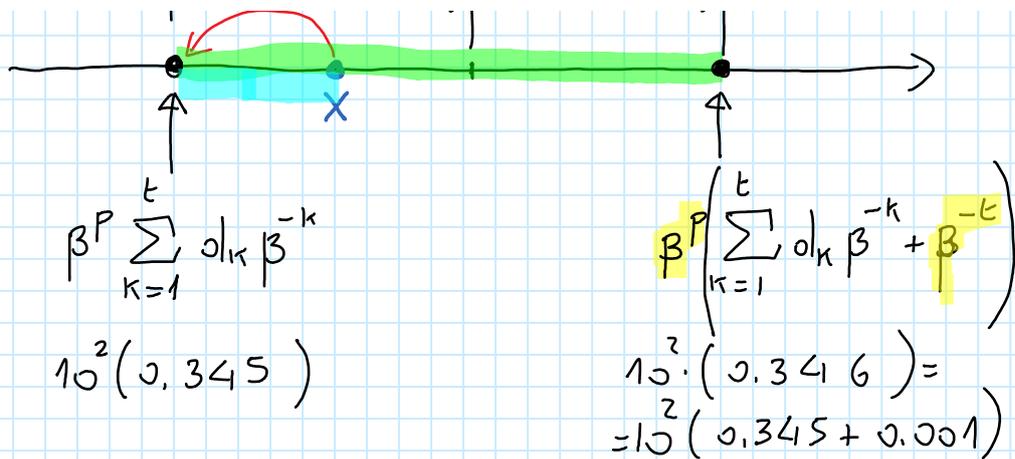
Dato che NON ci sono under o over flow,  $x$  è nel range dei numeri macchina. Ho due casi.

1)  $x \in \mathbb{F}(\beta, t, L, U)$ . Allora è  $fl(x) = x$

$$\left| \frac{fl(x) - x}{x} \right| = \left| \frac{0}{x} \right| = |0| = 0 < \frac{1}{2} \beta^{1-t}$$

2)  $x \notin \mathbb{F}(\beta, t, L, U)$ . Allora  $x$  cade tra due numeri di macchina consecutivi





Quindi, la distanza tra i due numeri macchina è  $\beta^P \cdot \beta^{-t} = \beta^{P-t}$ . Ne segue che

$$|f(x) - x| \leq \frac{1}{2} \beta^{P-t}$$

Perciò ho

$$\left| \frac{f(x) - x}{x} \right| = \frac{|f(x) - x|}{|x|} \leq \frac{\frac{1}{2} \beta^{P-t}}{x} \leq \frac{\frac{1}{2} \beta^{P-t}}{\beta^{P-1}} = \frac{1}{2} \beta^{1-t}$$

$$x = \beta^P \cdot (d_1 \beta^{-1} + d_2 \beta^{-2} + \dots) \geq \beta^P \cdot 1 \cdot \beta^{-1} = \beta^{P-1}$$

( $\bar{e}$ , o potrebbe essere, molto pessimistica).

DEF. Sia  $IF(\beta, t, L, U)$ . Il numero positivo

$$\text{eps} = \frac{1}{2} \beta^{1-t}$$

è detto PRECISIONE DI MACCHINA del sistema floating point  $IF$ .

La precisione di macchina è il più piccolo numero positivo tale che

$$1 + \text{eps} > 1$$

nel sistema floating-point  $IF(\beta, t, L, U)$ .

Infatti,  $1 \in IF(\beta, t, L, U)$  perché

$$1 = \beta^0 \cdot (1 \cdot \beta^{-1} + 0 \cdot \beta^{-2} + \dots + 0 \cdot \beta^{-t})$$

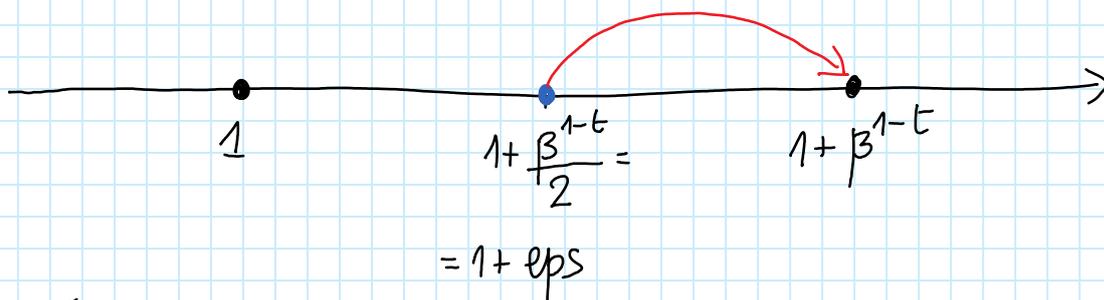
$$1 = \beta^p \cdot \left( \overset{1}{\beta^{-1}} + \overset{0}{\beta^{-2}} + \dots + \overset{0}{\beta^{-t}} \right)$$

$\uparrow$   $\uparrow$   $\uparrow$   $\uparrow$   
 $p$   $d_1$   $d_2$   $d_t$

Il numero floating point successivo al numero 1 è

$$\beta^p \cdot \left( 1 \cdot \beta^{-1} + 0 \cdot \beta^{-2} + \dots + 0 \cdot \beta^{-(t-1)} + 1 \cdot \beta^{-t} \right) =$$

$$= \beta^p \left( \beta^{-1} + \beta^{-t} \right) = 1 + \beta^{p-t}$$



Parisi,  $\epsilon$

- )  $0 < \epsilon < \text{eps} \Rightarrow 1 + \epsilon = 1$  in  $\mathbb{F}$
- )  $\text{eps} \leq \epsilon \Rightarrow 1 + \epsilon = 1 + \beta^{p-t}$  in  $\mathbb{F}$ .

ESERCIZIO Cosa succede se consideriamo il numero

$$x = \beta^p \sum_{k=1}^t d_k \beta^{-k} ?$$

Ossia, quale è il più piccolo numero positivo che devo sommare ad  $x$  per ottenere il successivo numero macchina?

## OPERAZIONI IN $\mathbb{F}(\beta, t, L, U)$

In modo APPROSSIMATO (la realtà è più complicata!) definiamo le operazioni elementari in  $\mathbb{F}(\beta, t, L, U)$  in questo modo:

SOMMA IN  $\mathbb{F}$   $x \oplus y = fl(x + y)$

CONTRAZIONE IN  $\mathbb{F}$   $\forall a, b \in \mathbb{D} \quad (a - b)$

operazione in  $\mathbb{R}$

SOMMA IN IF  $x \oplus y = fl(\hat{x+y})$

SOTTRAZIONE IN IF  $x \ominus y = fl(\hat{x-y})$

MOLTIPLICAZIONE IN IF  $x \odot y = fl(\hat{x \cdot y})$

DIVISIONE IN IF  $x \oslash y = fl(\hat{x/y})$

dove  $x, y \in \mathbb{F}(\beta, t, L, U)$ .

- Però, ogni operazione in IF ha due passi:
- eseguo la corrispondente operazione in  $\mathbb{R}$ ;
  - rappresento il risultato in IF.

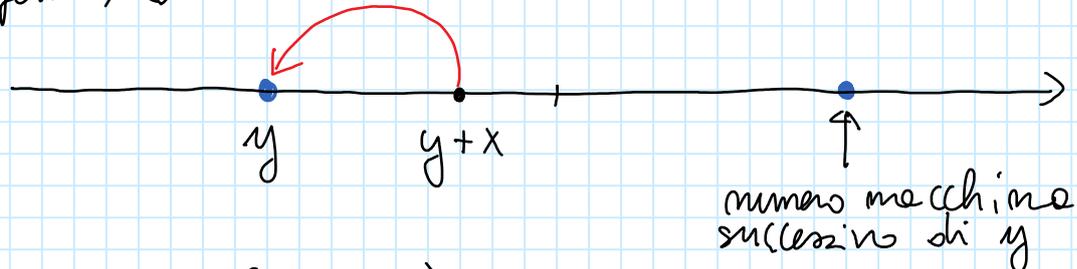
In modo analogo si definiscono le altre operazioni che si possono ritenere elementari, quali, ad esempio, la radice quadrata.

**OSSERVAZIONE** Le normali regole dell'algebra possono NON essere più valide quando operiamo con la aritmetica di macchina.

Ad esempio, se  $x > 0$ ,  $x \in \mathbb{F}(\beta, t, L, U)$  ed  $y \in \mathbb{F}(\beta, t, L, U)$  NON è più detto che valga la relazione (valide in  $\mathbb{R}$ )

$$y \oplus x > y$$

Infatti, e



$$y \oplus x = fl(y+x) = y.$$

**ESEMPIO** Consideriamo  $\mathbb{F}(10, 1, -1, 2)$

Sia  $x = 0.1$ ,  $y = 0.08$ , numeri di  $\mathbb{F}$ .

Allora è

$$x \oplus y = fl(0.1 + 0.08) = fl(0.18) = 0.2$$

$$x \oplus y = fl(0.1 + 0.08) = fl(\overset{\text{"X"}}{\circlearrowleft} 0.18) = 0.2$$

perché  $0.1 < 0.18 < 0.2$

e  $0.2$  è più vicino a  $0.18$ .

$$x \ominus y = fl(0.1 - 0.08) = fl(\overset{\text{"IF"}}{\circlearrowleft} 0.02) = 0.02$$

$$x \odot y = fl(0.1 \cdot 0.08) = fl(0.008) \underset{\uparrow}{=} 0$$

UNDERFLOW

$$x \oslash y = fl(0.1 / 0.08) = fl(1.25) = 1.$$

Nelle varie operazioni, tranne le  $2^{\pm}$ , ho degli errori!

ESEMPIO Consideriamo  $\mathbb{F}(10, 3, -2, 2)$  e i tre numeri  
meccanici

$$x = 0.123$$

$$y = 45.6$$

$$z = -45.5$$

Calcoliamo  $x + y + z$  in due modi:

$$(x \oplus y) \oplus z \quad \text{or} \quad x \oplus (y \oplus z)$$

In  $\mathbb{R}$  i due modi sono equivalenti. E in  $\mathbb{F}$ ?

$$\begin{aligned} x \oplus y &= fl(x + y) = fl(0.123 + 45.6) = \\ &= fl(\circlearrowleft 45.723) = 45.7 \quad (\text{ho "perso" le cifre 2, 3}). \end{aligned}$$

$$(x \oplus y) \oplus z = fl(45.7 + (-45.5)) = fl(0.2) = 0.200$$

Per l'altra possibilità ho

$$y \oplus z = fl(45.6 + (-45.5)) = fl(0.1) = 0.100$$

$$x \oplus (y \oplus z) = fl(0.123 + 0.100) = fl(0.223) = 0.223$$

Per tanto,

→ NON vale associatività delle somme

→ ho  $x+y+z = 0.123 + 45.6 + (-45.5) = 0.223$   
 per cui la seconda è corretto ma la prima non lo è!

**DEF. (STABILITÀ DI UN ALGORITMO)** Un algoritmo, che funziona come una successione di operazioni meccaniche elementari, si dice STABILE se e solo se NON AMPLIFICA gli errori nelle varie operazioni. In caso contrario, l'algoritmo si dice INSTABILE.

Chiaramente, ci interessano algoritmi stabili!

**ESEMPIO** Supponiamo di valutare  $f(x) = \sqrt{x+1} - \sqrt{x}$  per  $x=49$  usando  $\mathbb{F}(10, 1, -1, 2)$ .

In aritmetica esatte è

$$f(49) = \sqrt{50} - \sqrt{49} = 0.0710678...$$

In  $\mathbb{F}$ , abbiamo invece

$$\sqrt{50} = fl(\sqrt{50}) = fl(7.071...) = 7$$

$$\sqrt{49} = fl(\sqrt{49}) = fl(7) = 7$$

e quindi, in  $\mathbb{F}$ , è

$$\sqrt{50} \ominus \sqrt{49} = 7 - 7 = 0$$

Se cambio le modalità di valutare  $f$ , ossia CAMBIO l'ALGORITMO, nel seguente modo

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

$$= \frac{\cancel{\sqrt{x+1}} - \cancel{\sqrt{x}}}{\sqrt{x+1} + \sqrt{x}} = \frac{1}{\sqrt{x+1} + \sqrt{x}}$$

Ho, in  $\mathbb{F}$ ,

$$\sqrt{49+1} = \sqrt{50} = 7$$

$$\sqrt{49} = 7$$

$$\sqrt{50} \oplus \sqrt{49} = 7 \oplus 7 = fl(14) = 10$$

$$\frac{1}{\sqrt{50} \oplus \sqrt{49}} = fl(1/10) = 0.1$$

Quindi, per lo STESSO PROBLEMA, cambiando algoritmo è possibile contenere meglio gli errori nelle varie operazioni (ovvero, avere un algoritmo più stabile).

ESEMPIO Riscrivere un algoritmo stabile per valutare

$$a) f(x) = \frac{\sin(x)}{x}$$

$$b) \frac{1 - \cos(x)}{x^2}$$

per  $x \approx 0$ .