

Conversione:

$$(0.1)_{10} = (?)_2$$

$$0.1 \times 2 = 0.2 \quad \text{periodo}$$

$$\rightarrow 0.2 \times 2 = 0.4$$

$$0.4 \times 2 = 0.8$$

$$0.8 \times 2 = 1.6$$

$$\rightarrow 0.6 \times 2 = 1.2$$

allora

$$(0.1)_{10} = (0.00011)_2 =$$

$$= (1.1001 \times 2^{-4})_2$$

sviluppo finito in base 10 ma ∞ in base 2!

INSIEME DEI NUMERI DI MACCHINA

L'insieme dei numeri d'Macchine
in base β e $t+1$ cifre significative
e range per l'esponente (M_1, M_2) è
costituito dai seguenti elementi

$$\left\{ x \in \mathbb{R} : x = \pm \sum_{i=0}^t d_i \beta^{-i} \times \beta^p \right\}$$

dove $d_i \in \mathbb{N}$, $0 \leq d_i \leq \beta - 1$

per $i = 0, 1, \dots, t$, $d_0 \neq 0$

e $p \in \mathbb{Z}$ con $M_1 + 1 \leq p \leq M_2 - 1$

E_{SS} è dimostrato con

$$\mathbb{F}(\beta, t, M_1, M_2)$$

Osservazioni:

(1) $0 \notin \mathbb{F}$

(2) inferenzialmente i numeri machine hanno la seguente forma:

$$X = \pm \underbrace{d_0.d_1d_2\dots d_t}_{\text{cifre}} \times \beta^p$$

↑
base

↑
esponente

ESEMPIO

$$3.14 \in \mathbb{F}(10, 2, -1, 2) =$$

$$= \{ \pm d_0.d_1d_2 \times 10^p, p \in \{0, 1\} \}$$

$$3.14 \notin \mathbb{F}(10, \textcolor{red}{1}, -1, 2)$$

$$3.14 \notin \mathbb{F}(10, 2, \textcolor{red}{0}, 2)$$

più piccolo numero macchina strett. positivo:

Matlab \rightarrow $\text{realmin} := 1.00 \dots 0 \times \beta^{M_1+1} = \beta^{M_1+1}$

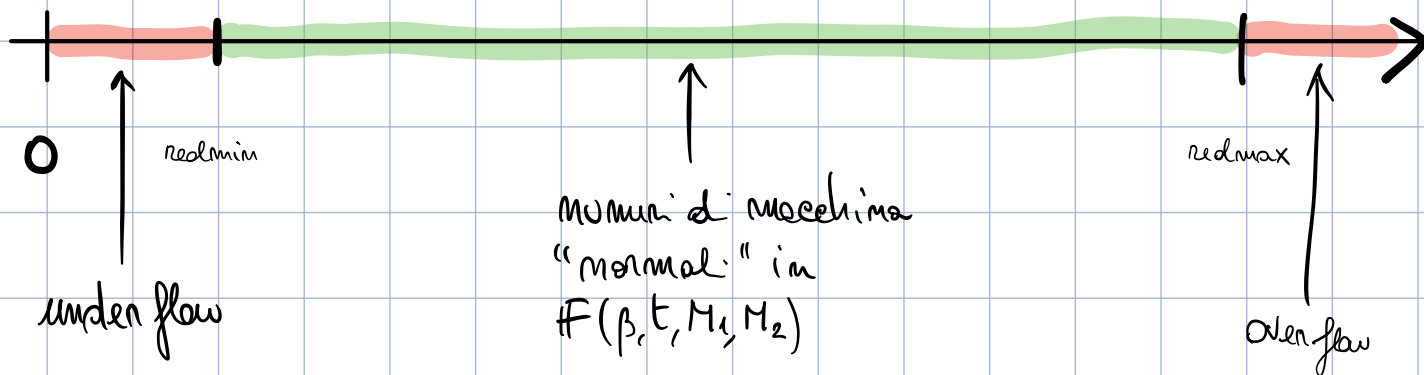
più grande numero di macchine:

$$\text{realmax} := \beta^{-1}, \beta^{-1} \beta^{-1} \dots \beta^{-1} \times \beta^{M_2-1}$$

aggiungo e sottraggo $0.00 \dots 01 \times \beta^{M_2-1}$

$$\begin{array}{r}
 \left| \begin{array}{ccccccc} \beta^{-1} & \beta^{-1} & \beta^{-1} & \dots & \beta^{-1} & \beta^{-1} & \times \beta^{M_2-1} \\ 0 & 0 & 0 & \dots & 0 & 1 & \times \beta^{M_2-1} \\ 0 & 0 & 0 & \dots & 0 & 1 & \times \beta^{M_2-1} \end{array} \right. + \\
 \left| \begin{array}{ccccccc} 0 & 0 & 0 & \dots & 0 & 1 & \times \beta^{M_2-1} \\ 0 & 0 & 0 & \dots & 0 & 1 & \times \beta^{M_2-1} \end{array} \right. - \\
 \hline
 1 \left| \begin{array}{ccccccc} 0 & 0 & 0 & \dots & 0 & 0 & \times \beta^{M_2-1} \\ 0 & 0 & 0 & \dots & 0 & 0 & \times \beta^{M_2-1-t} \end{array} \right. - \\
 1 \left| \begin{array}{ccccccc} 0 & 0 & 0 & \dots & 0 & 0 & \times \beta^{M_2-1-t} \end{array} \right. =
 \end{array}$$

$$\beta^{M_2} - \beta^{M_2-1-t} = \beta^{M_2-1} \left(\beta - \beta^{-t} \right)$$



Come vengono gestiti underflow e overflow?

Definizione "graduale" dell'underflow:

$$\text{numeri denormali} = \left\{ \pm 0.d_1d_2\dots d_t \times \beta^{M_1+1}, d_i \neq 0 \right. \\ \left. \uparrow \text{per almeno un } i=1,2,\dots,t \right\}$$

Viene "rilassata" la condizione di normalizzazione ($d_0 \neq 0$)

più piccolo numero denormale strett. positivo:

$$\underbrace{0.0\dots 01}_{t \text{ zeri}} \times \beta^{M_1+1} = \beta^{M_1+1-t}$$

ESEMPIO: $F(2,2,-2,3) =$

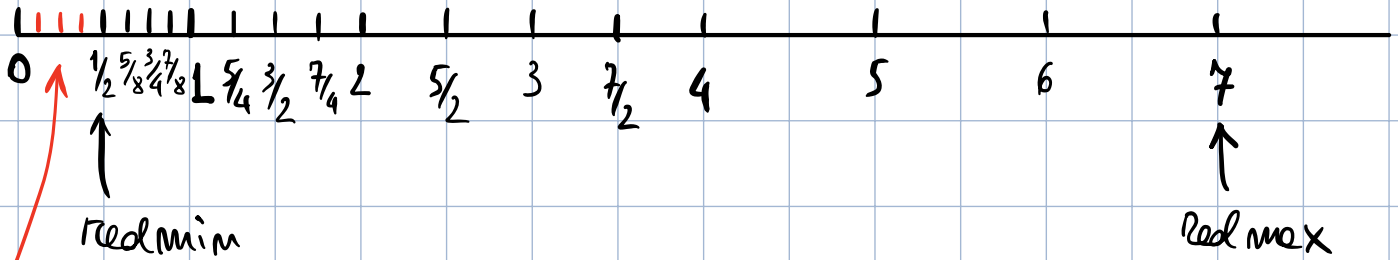
$$= \left\{ \pm 1.d_1d_2 \times 2^p, p = -1, 0, 1, 2 \right\}$$

Sono 32 numeri (16 positivi). I positivi sono:

$p = -1$	1.00×2^{-1}	1.01×2^{-1}	1.10×2^{-1}	1.11×2^{-1}
$p = 0$	1.00×2^0	1.01×2^0	1.10×2^0	1.11×2^0

$p=1$	1.00×2^1	1.01×2^1	1.10×2^1	1.11×2^1
$p=2$	1.00×2^2	1.01×2^2	1.10×2^2	1.11×2^2

Rappresentazioni sull'asse reale:



Osserv.: maggiore densità per numeri piccoli

Osserv.: **denormali** = $\left\{ \pm 0.d_1d_2 \times 2^{-1} \right\} =$
 $= \left\{ \pm \frac{1}{8}, \pm \frac{1}{4}, \pm \frac{3}{8} \right\}$

Distanze tra numeri macchina consecutivi:

Sia $x = d_0.d_1d_2 \dots d_t \times \beta^p \in \mathbb{F}(\beta, t, n_1, n_2)$.

Allora y , numero macchina successivo, è dato da

$$y = x + 0.0 \dots 01 \times \beta^p \Rightarrow$$

la distanza di x da y è $d = 0.0 \dots 01 \times \beta^p = \beta^{p-t}$

Conclusione: se $x < y \in F(\beta, t, M_1, M_2)$ sono

numeri macchine consecutivi e $x, y \in [\beta^p, \beta^{p+1}]$,

allora la loro distanza è

$$y - x = \beta^{p-t}$$

Esempio: distanza di $5/4$ dal numero macchina consecutivo in $F(2, 2, -2, 3)$:

$$5/4 \in [1, 2] = [2^0, 2^1] \Rightarrow$$

$$\text{distanza} = 2^{0-2} = 1/4$$

Standard IEEE 754 single/double precision

Single prec.: il numero è memorizzato sotto

forma di stringa a 32 bit:

	S	q	f
# bit:	1	8	23
Cosa:	segno	esponente	mantissa

Sono i numeri $(-1)^S \times 1.f \times 2^P$, dove

positivo $\times S=0$
negativo $\times S=1$

$$S = 0, 1, P = q - \underbrace{\text{"bias"}}_?$$

$$q = (00000000)_2, (00000001)_2, \dots,$$

$$(11111110)_2, (11111111)_2 =$$

$$= \cancel{0}, 1, \dots, 254, \cancel{255}$$

sottruggo 127 per distribuirli
meglio attorno a zero

$$P = q - \underset{\substack{\uparrow \\ \text{bias}}}{127} = -126, \dots, 127$$

Si tratta dell'insieme $\mathbb{F}(2, 23, -127, 128)$

Esempi:

$$\text{redmin} = \beta^{11+1} = 2^{-126} \approx 1.2 \times 10^{-38}$$

$$\begin{aligned} \text{redmax} &= \beta^{11-1}(\beta - \beta^{-t}) = 2^{127}(2 - 2^{-23}) \approx \\ &\approx 2^{128} \approx 1.7 \times 10^{38} \end{aligned}$$

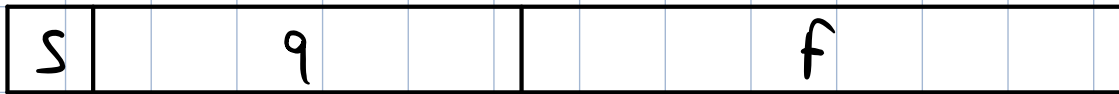
$$\begin{aligned} \text{più piccolo denormale} \\ \text{positivo} &= \beta^{11+1-t} = 2^{-149} \approx 1.4 \times 10^{-45} \end{aligned}$$

Le stringhe $q = 00000000, 11111111$ non contribuiscono all'esponente perché rivestono un ruolo speciale:

q	mantissa	uso
00000000	0	zero macchina
00000000	$\neq 0$	denormali
11111111	0	$\pm \text{Inf}$
11111111	$\neq 0$	NaN

ESEMPLI (Matlab): $1/0 = \text{Inf}$, $\text{Inf} + \text{Inf} = \text{Inf}$, $\text{Inf} - \text{Inf} = \text{NaN}$

Double prec.: (default in Matlab)



#25: 1 11 52

$$q = 0, 1, \dots, 2046, 2047$$

$$\text{bias} : 1023$$

$$p = q - \text{"bias"} = -1022, \dots, 1023$$

$$\text{si Trovate di } F(2, 52, -1023, 1024)$$

Determinare realmin , realmax e

" realmin normale".