

① Representação de números no computador

- \* bit :  $\rightarrow 0$  ou  $1$  (com ou sem corrente)
  - $\rightarrow$  menor unidade de informação
  - $\rightarrow$  números binários

- \* Representação dos números:

$$\begin{aligned}(1328)_{10} &= 1 \times 10^3 + 3 \times 10^2 + 2 \times 10^1 + 8 \times 10^0 \\ &= 1000 + 300 + 20 + 8\end{aligned}$$

$\Rightarrow$  em qualquer base:

$$(a_3 a_2 a_1 a_0)_\beta = (a_3 \beta^3 + a_2 \beta^2 + a_1 \beta^1 + a_0 \beta^0)_{10}$$

onde  $0 \leq a_i \leq \beta - 1$

$$\begin{aligned}\Rightarrow \text{base } 2: (10010)_2 &= (1 \times 2^4 + 0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 0 \times 2^0)_{10} \\ &= (16 + 2)_{10} = (18)_{10}\end{aligned}$$

$$\bullet (0)_2 = (0 \times 2^0)_{10} = (0)_{10}$$

$$\bullet (1)_2 = (1 \times 2^0)_{10} = (1)_{10}$$

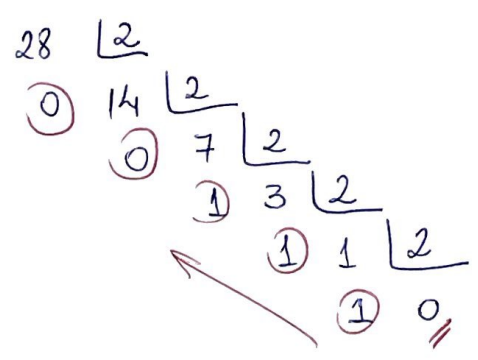
$$\bullet (10)_2 = (1 \times 2^1 + 0 \times 2^0)_{10} = (2)_{10}$$

$$\bullet (11)_2 = (1 \times 2^1 + 1 \times 2^0)_{10} = (3)_{10}$$

$$\bullet (100)_2 = (1 \times 2^2 + 0 \times 2^1 + 0 \times 2^0)_{10} = (4)_{10}$$

•  $(28)_{10} = ?$

$= (11100)_2$



•  $x = 2 \rightarrow (10)_2 \rightarrow 2 \text{ bits}$

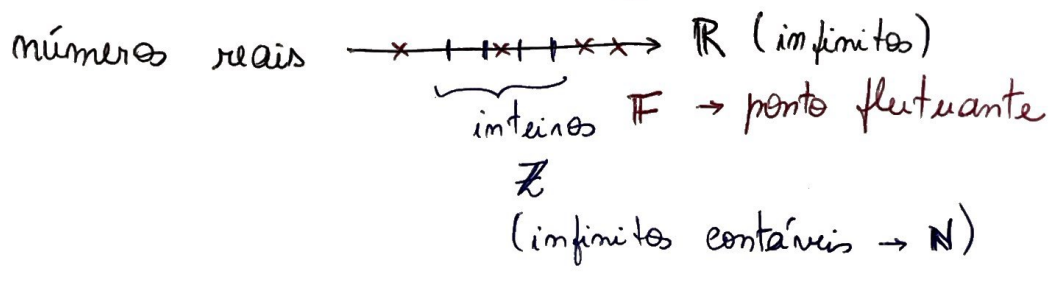
•  $x = 0,25 \quad ?$

•  $x = 1/3 = 0,333 \quad ?$

•  $x = \sqrt{2} \quad ?$

} não tem solução exata no computador

↳ irracional,  $\neq \frac{A}{B}$



$\Rightarrow F \rightarrow$  ponto flutuante  $\rightarrow$  número finito de casas decimais  
↳ arredondamento

$x = 1/3 = 0,3333$

$x = 2/3 = 0,6667$

$\Rightarrow$  posição do número decimal não é fixa

$$x = (-1)^s (0.\underbrace{a_1 a_2 \dots a_t}_m) \beta^e = (-1)^s m \beta^{e-t}$$

•  $\beta = \text{base} \geq 2$  (inteiro)

•  $m = \text{mantissa}$  (inteiro)  $0 \leq a_i \leq \beta - 1$

•  $t = \text{n}^\circ$  de dígitos da mantissa

•  $e = \text{expoente}$  (inteiro),  $s = \text{ sinal } \begin{cases} 0 & \oplus \\ 1 & \ominus \end{cases}$

ex)  $\cdot 250 = (-1)^0 0,25 \times 10^3 = (-1)^0 \textcircled{25} \times 10^{\textcircled{3}}$

$\cdot -0,33 = (-1)^1 0,33 \times 10^0 = (-1)^1 \textcircled{33} \times 10^{\textcircled{-2}}$

$\cdot 0,001 = (-1)^0 0,1 \times 10^{-2} = (-1)^0 \textcircled{1} \times 10^{\textcircled{-3}}$

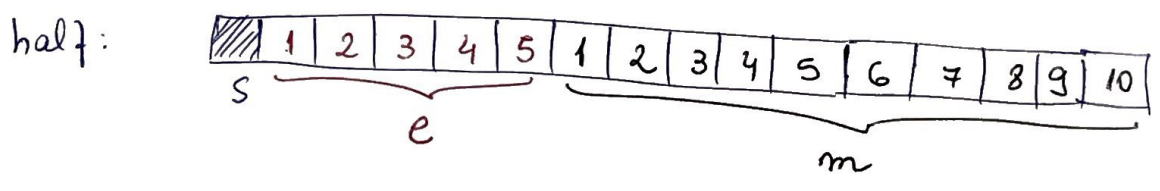
obs: se  $a_1 \neq 0$ , representação é única.

\* Padrão IEEE 754

$(-1)^s \textcircled{1} \cdot m \times 2^{e-\textcircled{b}}$   $E = e - b$

$\swarrow$  implícito  $\nwarrow$  bias

- byte = 8 bits
- half = 16 bits (2 bytes)
- single = 32 bits (4 bytes)
- double = 64 bits (8 bytes)



Cases especiais:

{	0	00000	000000000000	$\rightarrow (0)_{10}$
	0	11111	"	$\rightarrow +\infty$
	1	11111	"	$\rightarrow -\infty$

obs:  $e = \begin{cases} 00000 \\ 11111 \end{cases}$  são casos especiais

$(1)_{10} = (1)_2 \Rightarrow 0 \ 01111 \ 000 \ 000 \ 0000$   
 $= 1 \times 2^0 = (-1)^0 1,0 \times 2^0 \quad s=0, m=0, E=0$

⇒ half: 16 bits: 1(s), 5(e), 10(m)

$$F(2, 11, -14, 15)$$

- $e_{\max} = (11110)_2 = (30)_{10}$

- $e_{\min} = (00001)_2 = (1)_{10}$

números pequenos  
↓

- $e = [1, 30] \Rightarrow b = 15 \Rightarrow E = [-14, 15]$

↑  
números grandes

- $|x|_{\max} = 0 \ 11110 \ 1111111111$

$$= (-1)^0 (1, 1111111111)_2 \times 2^{(11110)_2 - (15)_{10}}$$

$$= [1 \times 2^0 + (1 - 2^{-10})] \times 2^{15}$$

$$= (2 - 2^{-10}) 2^{15} = 65504 = 6,55 \times 10^4$$

- $|x|_{\min} = 0 \ 00001 \ 0000000000$

$$= (-1)^0 1,0 \times 2^{(00001)_2 - (15)_{10}}$$

$$= 1 \times 2^{-14} = 1/16384 = 6,10 \times 10^{-5}$$

- precisão: 10 (mantissa) + 1 = 11 bits

$$\underbrace{(1111111111)_2}_{11 \text{ bits}} = \underbrace{(2047)_{10}}$$

11 bits

a partir de 2048, precisa de 12 bits, logo, todos os números com 3 dígitos (até 999) são representados corretamente.

$$p_{10} = 3$$

⇒ single: 32 bits: 1(s), 8(e), 23(m)

(5)

$\mathbb{F}(2, 24, -126, 127)$

- $e_{\max} = (11111110)_2 = (254)_{10}$
- $e_{\min} = (00000001)_2 = (1)_{10}$
- $e = [1, 254] \Rightarrow b = 127 \Rightarrow E = [-\underline{126}, \underline{127}]$
- $|x|_{\max} = (2 - 2^{-23}) \times 2^{\underline{127}} = 3,4028235 \times 10^{38}$
- $|x|_{\min} = 1 \times 2^{-\underline{126}} = 1,1754944 \times 10^{-38}$
- $p_{10} = 7$  (24 bits em  $\beta=10$ )

⇒ double: 64 bits: 1(s), 11(e), 52(m)

$\mathbb{F}(2, 53, -1022, 1023)$

- $e_{\max} = (11111111110)_2 = (2046)_{10}$
- $e_{\min} = (00000000001)_2 = (1)_{10}$
- $e = [1, 2046] \rightarrow b = 1023 \rightarrow E = [-1022, 1023]$
- $|x|_{\max} = (2 - 2^{-52}) 2^{1023} = 1,797693 \dots \times 10^{308}$
- $|x|_{\min} = 1 \times 2^{-1022} = 2,2250 \dots \times 10^{-308}$
- $p_{10} = 15$  (53 bits em  $\beta=10$ )

$\mathbb{F}(\beta, t, E_{\min}, E_{\max})$

$E_{\min} \leq E \leq E_{\max}$

obs: Quaterni:  $\mathbb{F}(2, 53, -\underline{1021}, \underline{1024}) \Rightarrow (-1)^s \underline{0.1_m} 2^{e-b}$



## ① Representação de números no computador

⇒ ponto flutuante:  $F(\beta, t, \underbrace{E_{\min}, E_{\max}}_{\text{decimal}})$ ,  $E_{\min} \leq E \leq E_{\max}$

\* Padrão IEEE 754:  $x = (-1)^s \underbrace{(1)}_{\text{implícito}} m \cdot 2^{e - \underbrace{b}_{\text{viés}}}$

$$\begin{cases} \beta = \text{base} = 2 \\ t = \text{precisão} = m + 1 \quad (m = \text{memória}) \quad t = \text{mantissa} \\ \quad \quad \quad \text{binária} \quad \quad \quad \text{"real"} \quad \quad \quad \text{real} \\ E = \text{exponente} = \underbrace{e}_{\text{memória}} - b \end{cases}$$

• half:  $F(2, 11, -14, 15)$   $\left\{ \begin{array}{l} 16 \text{ bits: } 1(s), 5(e), 10(m) \\ b = 15 \end{array} \right.$

• single:  $F(2, 24, -126, 127)$   $\left\{ \begin{array}{l} 32 \text{ bits: } 1(s), 8(e), 23(m) \\ b = 127 \end{array} \right.$

• double:  $F(2, 53, -1022, 1023)$   $\left\{ \begin{array}{l} 64 \text{ bits: } 1(s), 11(e), 52(m) \\ b = 1023 \end{array} \right.$

ex:  $\left\{ \begin{array}{l} x = (1)_{10} \\ x = (1)_2 \end{array} \right. \quad x = (-1)^0 1,0 \cdot 2^0 \quad s = 0, E = 0, m = 0$   
 $e = E + b = 0 + 15 = (15)_{10}$   
 $e = 1111$

$x_{\text{half}} = 0 \ 01111 \ 0000000000$

obs: Quaterni: double:  $F(2, 53, -1021, 1024)$

pois  $x = (-1)^s 0, \underline{1m} 2^{e-b}$

$|x|_{\min} = 2^{E_{\min}}, \quad |x|_{\max} = (2 - 2^{-\overset{t-1}{\downarrow} m}) 2^{E_{\max}}$

\* precisão em decimal:

• half: 11 bits em  $\beta=2 \Rightarrow (1111111111)_2 = (2047)_{10}$   
 $p_{10} = 3$  (até 999).

• single: 24 bits em  $\beta=2 \Rightarrow (111 \dots 11)_2 = 16.777.215$   
 $p_{10} = 7$  (até 9.999.999)

• double: 53 bits em  $\beta=2 \Rightarrow (111 \dots 11)_2 = ?$   
 $p_{10} = 15$

---


$$\begin{cases} * |x| > |x|_{\max} \Rightarrow \text{overflow } (\pm\infty) \\ * |x| < |x|_{\min} \Rightarrow \text{underflow } (zero) \end{cases}$$


---

\* exemplo memória: 1 milhão de dados

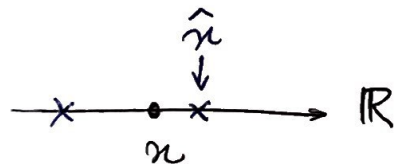
- half:  $10^6 \times 2 \text{ bytes} \approx 2 \text{ Mb}$
- single:  $10^6 \times 4 \text{ bytes} \approx 4 \text{ Mb}$
- double:  $10^6 \times 8 \text{ bytes} \approx 8 \text{ Mb}$

↳ usado pela RAM ou armazenamento binário

↳ ASCII → caracteres que representam números →  
 $\sim 3 \times$  arquivo binário.

(1.1)

Erros em ponto flutuante:



•  $EA_n = |x - \hat{n}| \rightarrow$  erro absoluto, minimizado na escolha de  $\hat{n}$  (arredondamento)

•  $ER_n = \frac{|x - \hat{n}|}{|x|} \rightarrow$  erro relativo,  $\sim$  constante em ponto flutuante

ex)  $\beta = 10, p_{10} = 4$

⇒ ponto fixo:

$$a) \begin{cases} n = 3507,6 \\ \hat{n} = 3508 \end{cases} \left\{ \begin{array}{l} EA_n = 0,4 \\ ER_n = \frac{0,4}{3507,6} \simeq 1,1 \times 10^{-4} \end{array} \right.$$

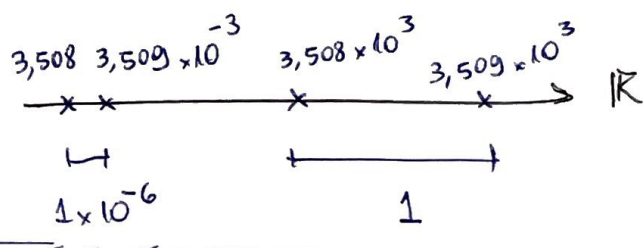
$$b) \begin{cases} n = 0,0035076 \\ \hat{n} = 0,004 \end{cases} \left\{ \begin{array}{l} EA_n = 0,0004924 \\ ER_n \simeq 0,14 \end{array} \right.$$

⇒ ponto flutuante:

$$a) \begin{cases} n = 3507,6 = 3,5076 \times 10^3 \\ \hat{n} = 3,508 \times 10^3 \end{cases} \left\{ \begin{array}{l} EA_n = 0,4 \\ ER_n \simeq 1,1 \times 10^{-4} \end{array} \right.$$

$$b) \begin{cases} n = 0,0035076 = 3,5076 \times 10^{-3} \\ \hat{n} = 3,508 \times 10^{-3} \end{cases} \left\{ \begin{array}{l} EA_n = 0,0004 \times 10^{-3} \\ ER_n = \frac{0,0004 \times 10^{-3}}{3,5076 \times 10^{-3}} \simeq 1,1 \times 10^{-4} \end{array} \right.$$

obs: quanto maior o valor absoluto, maior o erro absoluto e o espaçamento na reta  $\mathbb{R}$



\* operações aritméticas com ponto flutuante

a) associativa:  $(2 + 3) + 4 = 2 + (3 + 4)$

$(2 \times 3) \times 4 = 2 \times (3 \times 4)$

$\times \mathbb{R}$



b) comutativa:  $x + y = y + x$

$xy = yx$

✓ F

c) distributiva:  $2(1+3) = 2 \cdot 1 + 2 \cdot 3$  X F

erros de arredondamento a cada operação.

ex) 
$$\left\{ \begin{array}{l} (23,4 + 5,18) + 3,05 = 31,7 \\ 23,4 + (5,18 + 3,05) = 31,6 \end{array} \right.$$

ex) 
$$\left\{ \begin{array}{l} 3,18 (5,05 + 11,4) = 52,5 \\ 3,18 \cdot 5,05 + 3,18 \cdot 11,4 = 52,4 \end{array} \right.$$

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

se  $x = 1 \times 10^{-15}$  (próximo da precisão de máquina)

$$\frac{(1+x)-1}{x} = 1,1102... \quad \text{erro de } 11\% \quad !$$

obs: - trabalhe com números na ordem de 1 !

- Cap 1 Quaternioni: