

Aritmética de Ponto Flutuante

Um número real $x \in \mathbb{R}$ pode ser escrito como

$$x = \pm \left(\frac{d_0}{b^0} + \frac{d_1}{b^1} + \dots + \frac{d_t}{b^t} + \dots \right) \times \underbrace{b^e}_{\substack{\text{expoente} \\ \text{Base}}}$$

\rightarrow dígitos dependem da base

Ex: $\bullet c \in \mathbb{Z}$

$$\frac{1}{3} = \left(\frac{3}{10^0} + \frac{3}{10^1} + \frac{3}{10^2} + \dots + \frac{3}{10^t} + \dots \right) \times 10^{-1}$$

$$= 3.3333 \dots \times \underbrace{10^{-1}}_{\text{Base}} = 0.3333 \dots = 0.\bar{3}$$

• Sistema de Ponto Flutuante

$\left\{ \begin{array}{l} \beta \rightarrow \text{base} \\ t \rightarrow \text{precisão (nº de dígitos significativos)} \\ l \rightarrow \text{Limitante inferior do expoente } e \\ u \rightarrow \text{Limitante superior do expoente } e. \end{array} \right.$

$$fp(x) = \overset{\text{sign}}{\uparrow} \text{sgn}(x) \times \underbrace{(0.\tilde{d}_0\tilde{d}_1 \dots \tilde{d}_{t-1})}_{\text{mantissa}} \times \beta^e, e \in [l, u]$$

Para base decimal ($\beta = 10$) $\neq 0$

$$fl(x) = \text{sgn}(x) \cdot (0.\underbrace{d_0 d_1 \dots d_{t-1}}_{\text{mantissa finita}}) \times 10^e$$

em que t é a precisão, $0 \leq d_j \leq 9$, $d_j \in \mathbb{Z}$, $j = 0, \dots, t-1$,

$$d_0 \neq 0, \text{ e } e \in [l, u] \cap \mathbb{Z}$$

2º Se $x \in \mathbb{D}^c$, temos infinitos dígitos de x
em qualquer base β .

Como guardar em um computador com t dígitos de precisão o número ($\beta=10$)

$$x = \pm (0.d_0 d_1 \dots d_{t-1} d_t d_{t+1} \dots) \times 10^e$$

nas cabem todos os dígitos na máquina

Temos duas estratégias:

- Arredondamento
- Truncamento

(t digits)

• Arredondamento

$$fl_A(x) = \begin{cases} \pm (0.d_0d_1 \dots d_{t-1}) \times 10^e, & \text{se } d_t < 5 \\ \pm (0.d_0d_1 \dots d_{t-1} + \frac{1}{10^{t-1}}) \times 10^e, & \text{se } d_t \geq 5 \end{cases} \quad (P/2)$$

• Truncamento

$$fl_T(x) = \pm (0.d_0d_1 \dots d_{t-1}) \times 10^e$$

Ex: $\beta = 10$, $\underline{t = 3}$

x	Trunc	Arred
5.672	5.67	5.67
-5.672	-5.67	-5.67
5.677	5.67	5.68
5.692	5.69	5.69
-5.695	5.69	-5.70

$$\underline{Ex}: x = \frac{8}{3} = 2.6666\dots$$

(t=4)

$$fl_A(x) = 2.667 \times 10^0$$

$$= 0.2667 \times 10^1$$

$$fl_+(x) = 2.666 \times 10^0$$

$$= 0.2666 \times 10^1$$

normalized

• 64 bits (double-precision)

$$\boxed{\beta = 2}$$

1 bit - sign

11 bits - exponent ~

52 bits - mantissa (digits) ~ 16 digits
decimals

$\underline{Ex:}$
 $\beta = 10$, $\underline{t = 3}$ e $\underline{l = -5}$, $\underline{u = 5}$.
 $d_0 \neq 0$ $e \in [-5, 5]$

maior número representável (em módulo)

$$M = 0.999 \times 10^5 = 99900$$

menor número representável (em módulo)

$$m = 0.100 \times 10^{-5} = 0.000001$$

Teorema: Seja $x \rightarrow fl(x) = g \times \beta^e$, com
 $x \neq 0$, g é a mantissa normalizada.

Então

$$E_A: |x - fl(x)| \leq \begin{cases} \beta^{1-t} \times \beta^e, & p/\text{truncamento} \\ \frac{1}{2} \beta^{1-t} \times \beta^e, & p/\text{arredondamento} \end{cases}$$

$$E_R: \frac{|x - fl(x)|}{|x|} \leq \begin{cases} \beta^{1-t}, & p/\text{truncamento} \\ \frac{1}{2} \beta^{1-t}, & p/\text{arredondamento} \end{cases}$$

Erro nas Operações do P. flutuante

Ex: $p=10$, $t=4$ e $[e, v]$ suficientes

$$x = 0.9377 \times 10^4 \quad \text{e} \quad y = 0.1272 \times 10^2$$

$$\underline{x+y} \Rightarrow fl(\overset{\downarrow}{fl(x)} + \overset{\downarrow}{fl(y)})$$

$$f(x+y) = f\left(0.937 \times 10^4 + 0.001272 \times 10^4\right)$$

$$= f\left(0.\underbrace{9382}_{\text{red circle}}72 \times 10^4\right)$$

$$= \begin{cases} 0.9382 \times 10^4 \text{ (T)} \\ 0.9383 \times 10^4 \text{ (A)} \end{cases}$$

$$fp(x * y) = ?$$

TAKEDA