

ARITMÉTICA DE PONTO FLUTUANTE

o UM COMPUTADOR UTILIZA UMA QUANTIDADE FINITA DE BITS PARA REPRESENTAR NÚMEROS REAIS, LOGO, ELE SÓ PODE REPRESENTAR UM SUBSET FINITO DOS REAIS. ISSO TRAZ DUAS DIFICULDADES, OS NÚMEROS NÃO PODEM SER ARBITRARIAMENTE GRANDES E HÁ UM ESPAÇAMENTO ENTRE OS REPRESENTADOS PELO COMPUTADOR.

o ATUALMENTE, O PROBLEMA DE OVERFLOW E UNDERFLOW NÃO É UM PROBLEMA COMUM, JÁ QUE OS COMPUTADORES ARMAZENAM NÚMEROS SUFICIENTEMENTES GRANDES PARA A MAIORIA DOS PROBLEMAS

o PORÉM, O PROBLEMA DO GAP ENTRE OS NÚMEROS AINDA EXISTE, COMO, POR EXEMPLO, NO IEEE ARITMÉTICA DE PRECISÃO DOBRADA, O INTERVALO $[1, 2]$ É REPRESENTADO POR

$$1, 1 + 2^{-52}, 1 + 2 \cdot 2^{-52}, 1 + 3 \cdot 2^{-52}, \dots, 2$$

O GAP ENTRE OS NÚMEROS ADJACENTES NÃO SÃO MAIORES QUE 2^{-52} . É UM GAP TÃO PEQUENO QUE, EM ALGORITMOS ESTÁVEIS (PRÓXIMO DOCUMENTO), ELAS PODEM SER IGNORADOS, MAS ALGORITMOS INSTÁVEIS PODEM CAUSAR IMPRECIÇÕES

o VAMOS DEFINIR UM CONJUNTO DE NÚMEROS DE PONTO FLUTUANTE IDEAL PARA TRABALHARMOS

DEFINIÇÃO

SEJA $F \subset \mathbb{R}$ ONDE:

$$F = \left\{ \pm \left(\frac{m}{\beta^t} \right) \beta^e / e \in \mathbb{Z}, \beta \in \mathbb{N} \geq 2, t \in \mathbb{N} \geq 1, 1 \leq m \leq \beta^t \right\} \cup \{0\}$$

ONDE β É A BASE OU RADIX, t É A PRECISÃO, e É O EXPOENTE, $\pm \left(\frac{m}{\beta^t} \right)$ É A MANTISSA DE $x \in F$

o O SISTEMA É IDEAL POR NÃO TER UNDERFLOW NEM OVERFLOW, LOGO, É UM SET INFINITO E CONTÁVEL, ALÉM DE SER SEMELHANTE A SI MESMO: $F = \beta F$

o m ESTÁ SEMPRE EM BASE β (PARA OS CÁLCULOS PODEMOS DESCONSIDERAR ISSO)

o t INDICA QUANTOS DÍGITOS (NA BASE DADA) A MINHA MANTÍSSA POSSUI, OU SEJA, QUANTO MAIOR É t , MAIS NÚMEROS EU POSSO ARMAZENAR NO MEU CONJUNTO F (A QUANTIDADE DE DÍGITOS É ANALISADA NA BASE β)

o POR QUE DIVIDO POR β^t ? POIS, AO FAZER ISSO, EU FAÇO COM QUE O NÚMERO INTEIRO QUE TRABALHO, EM BASE β , TENHA EXATAMENTE t DÍGITOS

↳ POSSO ESCREVER 15, EM BASE 10, COMO:

$$0,15 \cdot 10^2, 1,5 \cdot 10, 0,00015 \cdot 10^5, \dots$$

DIVIDIR POR 10^t ME DEIXARÁ COM UM NÚMERO DE CASAS RELEVANTE PARA A REPRESENTAÇÃO

$$\Rightarrow \frac{15}{10^2} = \boxed{0,15}$$

$$t=2$$

DEFINIÇÃO ALTERNATIVA

PODEMOS DEFINIR, TAMBÉM, F COMO:

$$F = \{ 0, d_1 d_2 \dots d_t \cdot \beta^e / e \in \mathbb{Z} \}$$

$$\text{com } 0 \leq d_j < \beta$$

o UMA DEFINIÇÃO EQUIVALENTE, MAS DE MAIS FÁCIL VISUALIZAÇÃO. A PARTIR DISSO, PODEMOS DEFINIR CLARAMENTE O MENOR E MAIOR VALOR QUE F CONTÉM (CASO $e \in I$ COM I SENDO UM SUBCONJUNTO FINITO DE \mathbb{Z})

↳ SEJA $t=3, e \in [-5, 5], \beta=10$

O MENOR VALOR ABSOLUTO E NÃO NULO DO CONJUNTO F :

$$\text{menor} = 0,001 \cdot 10^{-5}$$

O MAIOR //

$$\text{maior} = 0,999 \cdot 10^5$$

◦ QUANDO TENTAMOS REPRESENTAR UM NÚMERO QUE A PRECISÃO NÃO SUPORTA (EXEMPLO, $0,3145 \cdot 10^2$, MAS $t=3$), PODEMOS ARREDONDAR OU TRUNCAR O NÚMERO

EX:

$$t=3, e \in [-10, 10]$$

$$\rightarrow 14,38 = 0, \underbrace{1438}_{4>3} \cdot 10^2$$

▷ ARREDONDAR:

PEGAMOS OS n DÍGITOS RELEVANTES DO NÚMERO ($n > t$), ANALISAMOS ENTÃO O $(t+1)$ -ÉSIMO DÍGITO (TUDO ISSO NA BASE β), ENTÃO, SE O $(t+1)$ -DÍGITO FOR MAIOR QUE $\frac{\beta}{2}$, EU ELIMINO DO $(t+1)$ -ÉSIMO DÍGITO EM DIANTE E SOMO 1 AO t -ÉSIMO DÍGITO.

$$0,1438 \cdot 10^2 \rightarrow \boxed{8 > 5} \rightarrow \boxed{0,144 \cdot 10^2}$$

SE O $(t+1)$ -ÉSIMO DÍGITO É IGUAL A $\frac{\beta}{2}$, PASSAMOS AO $(t+2)$ -ÉSIMO E ASSIM POR DIANTE, SE TODOS SÃO IGUAIS A $\frac{\beta}{2}$, TANTO FAZ A DECISÃO

▷ TRUNCAMENTO

ELIMINAMOS DO $(t+1)$ -ÉSIMO DÍGITO EM DIANTE

$$0,1438 \cdot 10^2 \rightarrow \underline{0,143 \cdot 10^2}$$

EPSILON DA MÁQUINA

◦ É UM NÚMERO QUE RESUME A RESOLUÇÃO DE F

DEFINIÇÃO

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

◦ ESSE NÚMERO É METADE DA DISTÂNCIA ENTRE 1 E O PRÓXIMO MAIOR NÚMERO REPRESENTÁVEL DE PONTO FLUTUANTE

◦ OU SEJA, TECNICAMENTE: $1 + \epsilon_{\text{ma}} > 1$

◦ MAS, POR QUE A DISTÂNCIA ENTRE 2 NÚMEROS REPRESENTÁVEIS É β^{1-t} ?

↳ Ao representar um número x em base β , eu tenho que

① $d_1 \dots d_t$ com $0 \leq d_j < \beta$, eu deixo fazer a menor alteração possível no meu número, para isso, eu simplesmente somo β , logo:

② Tenho t dígitos significativos

$$x = d_1 \cdot \beta^0 + d_2 \cdot \beta^{-1} + d_3 \cdot \beta^{-2} + \dots + d_t \cdot \beta^{1-t}$$

Logo, a menor alteração possível para o número é acrescentar 1 unidade a d_t , logo:

$$d_1 \beta^0 + \dots + (d_t + 1) \beta^{1-t} = d_1 \beta^0 + \dots + d_t \beta^{1-t} + \beta^{1-t} = \frac{x}{\beta^t} + \beta^{1-t}$$

• Se somarmos um número menor que ϵ_{ma} , pela precisão limitada que temos, teríamos $1 + \alpha = 1$ ($\alpha < \epsilon_{ma}$).

• Temos então a propriedade

$$\forall x \in \mathbb{R}, \exists x' \in F / |x - x'| \leq \epsilon_{mac} \cdot |x|$$

DEFINIÇÃO

Seja F o conjunto de ponto flutuante, temos:

$$fl: \mathbb{R} \rightarrow F$$

onde $fl(x)$ dá a aproximação mais próxima de x em F .

• Com essa definição, podemos reescrever a última desigualdade:

$$\forall x \in \mathbb{R}, \exists \epsilon / |\epsilon| \leq \epsilon_{mac}; fl(x) = x(1 + \epsilon) \quad \textcircled{E}$$

• Ou seja, a diferença entre um real e sua aproximação é sempre menor ou igual ao ϵ_{mac}

ARITMÉTICA DE PONTO FLUTUANTE

• As operações $+$, $-$, \times , \div são feitas em reais, no computador, ou seja, em F , vamos representar por \oplus , \ominus , \otimes , \oslash .

• UM COMPUTADOR ENTÃO É CONSTRUÍDO SEGUINDO O SEGUINTE PRINCÍPIO. SEJAM $x, y \in F$ E $*$ DENOTA UMA DAS OPERAÇÕES $(+, \times, -, \div)$ E \odot SEU EQUIVALENTE EM F , ENTÃO:

$$x \odot y = fl(x * y)$$

TEOREMA ①

(AXIOMA FUNDAMENTAL DA ARITMÉTICA COM PONTO FLUTUANTE)

$$\forall x, y \in F, \exists \epsilon \text{ com } |\epsilon| \leq \epsilon_{\text{mac}} / x \odot y = (x * y)(1 + \epsilon)$$

• OU SEJA, TODA A OPERAÇÃO FEITA EM F TEM UM ERRO DE, NO MÁXIMO, ϵ_{mac}

→ AGORA PODEMOS DEFINIR ϵ_{mac} COMO O MENOR NÚMERO QUE

① E ② VALEM.