

Erros de representação

Na construção de um equipamento computacional, uma questão importante a ser considerada em sua arquitetura é a forma que será adotada para representar os dados numéricos.

Basicamente, **na memória** de um equipamento, **cada número** é armazenado em uma posição que **consiste de um sinal que identifica se o número é positivo ou negativo** e **um número fixo e limitado de dígitos significativos**.

Erros de representação

Sistema de ponto flutuante

Um número no sistema de ponto flutuante é caracterizado por uma base **b**, um número de dígitos significativos **n** e um expoente **exp**.

Dizemos que um número real **nr** está representado no sistema de ponto flutuante se for possível escrevê-lo da seguinte maneira:

$$nr = m \times b^{\text{exp}}$$

Onde, **m** é a mantissa do número,

Neste sistema de ponto flutuante, as seguintes condições devem ser verificadas:

$$m = \pm 0.d_1 d_2 \cdots d_n \quad n \in N$$

Sendo n o número máximo do sistema de dígitos na mantissa d_1, d_2, \dots, d_n , dígitos significativos da mantissa, do sistema de representação, com o primeiro dígito satisfazendo a condição $1 \leq d_1 \leq (b-1)$ e os demais dígitos da seguinte maneira:

$$0 \leq d_i \leq (b-1); i = 2, 3, \dots, n.$$

Erros de representação

O expoente exp varia da seguinte maneira:

$$\text{exp}_{\min} \leq \text{exp} \leq \text{exp}_{\max}$$

Sendo, $\text{exp}_{\min} \leq 0$ e $\text{exp}_{\max} \geq 1$ com exp_{\min} e exp_{\max} inteiros.

A união de todos os números em ponto flutuante, juntamente com a representação do zero, constitui o sistema de ponto flutuante normalizado, que indicamos por **SPF** ($b, n, \text{exp}_{\min}, \text{exp}_{\max}$).

No sistema o zero é representado da seguinte maneira:

$$\text{Zero: } (0.\underbrace{0000\dots 0}_{n \text{ vezes}})b^{\text{exp}(\min)}$$

Considerando o SPF dado na fórmula genérica $(b, n, \text{exp}_{\min}, \text{exp}_{\max})$, temos:

a) O menor positivo exatamente representável, não nulo, é o real formado pela menor mantissa multiplicada pela base elevada ao menor expoente, isto é:

$$\text{menor: } (0.\underbrace{1000\dots 0}_{(n-1) \text{ vezes}})b^{\text{exp}(\min)}$$

Erros de representação

b) O maior positivo exatamente representável é o real formado pela maior mantissa multiplicada pela base elevada ao maior expoente, isto é:

$$\text{maior: } (0 \cdot \underbrace{[b-1] [b-1] \dots [b-1]}_{n \text{ vezes}}) b^{\text{exp}(\text{max})}$$

c) O número máximo de mantissas positivas possíveis é dado por:

$$\text{Mantissas}_+ = (b-1) b^{n-1}$$

d) O número máximo de expoentes possíveis é dado por:

$$\text{exp}_{\text{possíveis}} = \text{exp}_{\text{max}} - \text{exp}_{\text{min}} + 1$$

e) O número de elementos positivos representáveis é dado pelo produto entre o número máximo de mantissas pelo máximo de expoentes, isto é:

$$\text{NR}_+ = \text{Mantissas}_+ \times \text{exp}_{\text{possíveis}}$$

Se considerarmos que dado um número real $nr \in \text{SPF}$ temos que $-nr \in \text{SPF}$ e a representação do zero, podemos concluir que o número total de elementos exatamente representáveis NR_t é dado por:

$$\text{NR}_t = 2 \times \text{NR}_+ + 1$$

Exemplo

Considere o SPF $(b, n, \text{expmin}, \text{expmax}) = \text{SPF}(3, 2, -1, 2)$, isto é, de base 3, 2 dígitos de mantissa, menor expoente igual a -1 e maior expoente 2. Para este sistema temos:

a) O menor exatamente representável:

$$0.10 \times 3^{-1} = (1 \times 3^{-1} + 0 \times 3^{-2}) \times 3^{-1} = \frac{1}{9}$$

b) O maior exatamente representável:

$$0.22 \times 3^2 = (2 \times 3^{-1} + 2 \times 3^{-2}) \times 3^2 = 8$$

c) A quantidade de reais positivos exatamente representáveis:

Temos que a quantidade de reais positivos exatamente representáveis é dada pelo produto entre todas as mantissas possíveis de 2 dígitos, formadas com os dígitos da base 3, isto é, 0.10, 0.11, 0.12, 0.20, 0.21, 0.22, e todas as possibilidades de expoentes, que no caso são -1, 0, 1, 2.

Desta forma, os 24 positivos exatamente representáveis estão listados a seguir

Exemplo

$$\text{exp} = -1: 0.10 \times 3^{-1} = \frac{1}{9}$$

$$\text{exp} = 0: 0.10 \times 3^0 = \frac{1}{3}$$

$$\text{exp} = 1: 0.10 \times 3^1 = 1$$

$$\text{exp} = 2: 0.10 \times 3^2 = 3$$

Resolver para 0.11, 0.2, 0.21 e 0.22 ?

OBS: O menor real positivo representável é $1/9$
o maior positivo representável é o real 8.

$$\text{exp} = -1: 0.12 \times 3^{-1} = 5/27$$

$$\text{exp} = 0: 0.12 \times 3^0 = 5/9$$

$$\text{exp} = 1: 0.12 \times 3^1 = 5/3$$

$$\text{exp} = 2: 0.12 \times 3^2 = 5$$

Exemplo 2

Considere o SPF $(b, n, \text{expmin}, \text{expmax}) = \text{SPF}(3, 2, -1, 2)$, isto é, de base 3, 2 dígitos de mantissa, menor expoente igual a -1 e maior expoente 2. Para este sistema temos:

a) O menor exatamente representável:

$$x = (0.10)_3 \times 3^{-1} = (1 \times 3^{-1} + 0 \times 3^{-2}) \times 3^{-1} = \frac{1}{9} \quad \text{e} \quad y = 5 = (0.12)_3 \times 3^2$$

São exatamente representáveis, no entanto,

$$(x+y) = (0.00010)_3 \times 3^2 + (0.12)_3 \times 3^2 = (0.1201)_3 \times 3^2$$

Não é exatamente representável em SPF, uma vez que a mantissa é de 2 dígitos.

OBS: Pode ocorrer de outras propriedades consagradas no conjunto dos números reais não serem verdadeiras, no sentido da exatidão da representação, no sistema de ponto flutuante normalizado, como as propriedades comutativa e associativa.

Exemplo 3

Dados $x, y, z \in \text{Reais}$ e os Sistema de Ponto Flutuante normalizado SPF (3,2,-1,2) temos:

$$\text{Se } x = \frac{5}{3} = (0.12)_3 \times 3^1, y = \frac{7}{27} = (0.21)_3 \times 3^{-1} \text{ e}$$
$$z = \frac{8}{9} = (0.22)_3 \times 3^0$$

Temos:

e

$$x + (y + z) = 0.22 \times 3^1$$

$$(x + y) + z = 0.21 \times 3^1$$

Logo,

$$X+(y+z) \neq (x+y) +z$$