

# Análise de Erros em Aritmética de Ponto Flutuante

---

Márcio Antônio de Andrade Bortoloti

Cálculo Numérico

Departamento de Ciências Exatas e Tecnológicas - DCET

Universidade Estadual do Sudoeste da Bahia

# Sumário

## Análise de Erros

Truncamento e Arredondamento

Truncamento e Arredondamento

Erros Absoluto e Relativo

Operações em Aritmética de Ponto Flutuante

## Análise de Erros nas Operações de Ponto Flutuante

Adição

Subtração

Multiplicação

Divisão

## Efeitos Numéricos

Cancelamento Subtrativo

# Análise de Erros

---

### Definição

Se  $x \in \mathcal{F}(10, t, m, M)$  então ele pode ser representado como

$$x = f_x \times 10^e + g_x \times 10^{e-t},$$

onde  $0.1 \leq f_x < 1$  e  $0 \leq g_x < 1$ .

### Exemplo:

Seja  $x = 234.57$  e  $t = 4$ . logo

$$\begin{aligned} x &= 234.57 \\ &= 0.23457 \times 10^3 \\ &= (0.2345 + 0.00007) \times 10^3 \\ &= 0.2345 \times 10^3 + 0.00007 \times 10^3 \\ &= 0.2345 \times 10^3 + 0.7 \times 10^{-1} \end{aligned}$$

## Definição de Truncamento

Seja  $\mathcal{F}(10, t, m, M)$  uma máquina e  $x$  um número que em geral não pode ser representado em  $\mathcal{F}$  de forma exata. Quando isso ocorre, devemos utilizar uma aproximação  $\bar{x}$  para  $x$ . Assim, se  $x$  é tal que

$$x = f_x \times 10^e + g_x \times 10^{e-t}, \text{ para } m \leq e \leq M,$$

onde  $0.1 \leq f_x < 1$  e  $0 \leq g_x < 1$  então a operação de truncamento gera uma aproximação  $\bar{x}$ , de  $x$ , da forma

$$\bar{x} = f_x \times 10^e.$$

# Truncamento e Arredondamento

## Definição de Arredondamento

Seja  $\mathcal{F}(10, t, m, M)$ . No caso de obtermos uma aproximação,  $\bar{x}$ , de um número  $x = f_x \times 10^e + g_x \times 10^{e-t}$ , usando arredondamento, teremos que analisar  $g_x$  de forma que

$$\bar{x} = \begin{cases} f_x \times 10^e & \text{se } g_x < 1/2 \\ f_x \times 10^e + 10^{e-t} & \text{se } g_x \geq 1/2 \end{cases}$$

## Exemplo:

Considere uma máquina  $\mathcal{F}(10, 3, -5, 5)$ . Vamos representar  $x = 45.8787$  em  $\mathcal{F}$ . De fato,

$$x = 45.8787 = 0.458 \times 10^2 + 0.787 \times 10^{-1}$$

Fazendo o arredondamento

$$\bar{x} = 0.458 \times 10^2 + 10^{-1} = 0.459 \times 10^2 = 0.459$$

## Definição

Seja  $x \in \mathbb{R}$  e  $\bar{x}$  sua aproximação. O erro absoluto, cometido na representação de  $x$  por  $\bar{x}$  é definido por

$$EA_x = x - \bar{x}.$$

## Exemplo

O erro absoluto cometido na aproximação de  $\pi$  por  $\bar{\pi} = 3.14$  é

$$|EA_\pi| = |\pi - \bar{\pi}| = |\pi - 3.14| \leq 0.01.$$

## Erro Relativo

Considere dois números  $x = 1991.67$  e  $y = 3.67$ . Se aproximarmos  $x$  e  $y$  por  $\bar{x} = 1991.7$  e  $\bar{y} = 3.7$  teremos

$$|EA_x| = |EA_y| = 0.03.$$

No entanto, os dois números estão aproximados da “mesma forma” ?



## Erro Relativo

Considere dois números  $x = 1991.67$  e  $y = 3.67$ . Se aproximarmos  $x$  e  $y$  por  $\bar{x} = 1991.7$  e  $\bar{y} = 3.7$  teremos

$$|EA_x| = |EA_y| = 0.03.$$

No entanto, os dois números estão aproximados da “mesma forma” ? Qual aproximação está mais precisa ?

Para responder a pergunta vamos usar a seguinte definição:

# Erro Relativo

## Definição

O Erro Relativo,  $ER_x$ , cometido na aproximação de  $x$  por  $\bar{x}$  é definido como

$$ER_x = \frac{EA_x}{\bar{x}} = \frac{x - \bar{x}}{\bar{x}}$$

## Voltando ao exemplo ...

Se  $x = 1991.67$  e  $y = 3.67$  as aproximações  $\bar{x} = 1991.7$  e  $\bar{y} = 3.7$  cometem erros relativos da ordem de

$$|ER_x| = \frac{|EA_x|}{|\bar{x}|} = \frac{0.03}{1991.7} = 1.506250941 \times 10^{-5}.$$

$$|ER_y| = \frac{|EA_y|}{|\bar{y}|} = \frac{0.03}{3.7} = 0.810810810 \times 10^{-2}.$$

### Teorema

Sejam  $x \in \mathbb{R}$  e  $\mathcal{F}(10, t, m, M)$  uma máquina. Os erros absoluto e relativo cometidos na aproximação de  $x$  por  $\bar{x}$ , utilizando truncamento, são da ordem de

$$|EA_x| = |x - \bar{x}| < 10^{e-t} \quad \text{e} \quad |ER_x| = \frac{|EA_x|}{|\bar{x}|} < 10^{-t+1}.$$

### Prova:

Note que

$$x = f_x \times 10^e + g_x \times 10^{e-t},$$

onde  $0.1 \leq f_x < 1$  e  $0 \leq g_x < 1$ .

Usando o truncamento, tem-se

$$\bar{x} = f_x \times 10^e.$$

Logo

$$\begin{aligned}|EA_x| &= |x - \bar{x}| \\&= |f_x \times 10^e + g_x \times 10^{e-t} - f_x \times 10^e| \\&= |g_x| \times 10^{e-t} \\&< 10^{e-t} \quad (|g_x| < 1)\end{aligned}$$

Agora, o erro relativo ...

$$\begin{aligned}|ER_x| &= \frac{|EA_x|}{|\bar{x}|} \\&= \frac{|g_x| \times 10^{e-t}}{|f_x| \times 10^e} \\&< \frac{10^{e-t}}{0.1 \times 10^e} \\&< 10^{-t+1}\end{aligned}$$

## Teorema

Sejam  $x \in \mathbb{R}$  e  $\mathcal{F}(10, t, m, M)$  uma máquina. Os erros absoluto e relativo cometidos na aproximação de  $x$  por  $\bar{x}$ , utilizando arredondamento, são da ordem de

$$|EA_x| \leq 0.5 \times 10^{e-t} \quad \text{e} \quad |ER_x| \leq 0.5 \times 10^{-t+1}.$$

## Prova:

Note que

$$x = f_x \times 10^e + g_x \times 10^{e-t},$$

onde  $0.1 \leq f_x < 1$  e  $0 \leq g_x < 1$  e

$$\bar{x} = \begin{cases} f_x \times 10^e & \text{se } g_x < 1/2 \\ f_x \times 10^e + 10^{e-t} & \text{se } g_x \geq 1/2 \end{cases}$$

Se  $g_x < 1/2$  então

$$\begin{aligned}|EA_x| &= |x - \bar{x}| = |g_x| \times 10^{e-t} \\ &< \frac{1}{2} \times 10^{e-t}\end{aligned}$$

E também

$$\begin{aligned}|ER_x| &= \frac{|EA_x|}{|\bar{x}|} \\ &= \frac{|g_x| \times 10^{e-t}}{|f_x| \times 10^e} \\ &< \frac{0.5 \times 10^{e-t}}{0.1 \times 10^e} \\ &< \frac{1}{2} \times 10^{-t+1}\end{aligned}$$

Se  $g_x \geq 1/2$  então

$$\begin{aligned}|EA_x| &= |x - \bar{x}| \\&= |(f_x \times 10^e + g_x \times 10^{e-t}) - (f_x \times 10^e + 10^{e-t})| \\&= |g_x \times 10^{e-t} - 10^{e-t}| \\&= |g_x - 1| \times 10^{e-t} \\&\leq \frac{1}{2} \times 10^{e-t}\end{aligned}$$

E também

$$\begin{aligned}|ER_x| &= \frac{|EA_x|}{|\bar{x}|} \leq \frac{1/2 \times 10^{e-t}}{|f_x \times 10^e + 10^{e-t}|} \\&< \frac{1/2 \times 10^{e-t}}{|f_x| \times 10^e} < \frac{1/2 \times 10^{e-t}}{0.1 \times 10^e} < \frac{1}{2} \times 10^{-t+1}\end{aligned}$$

## Operações em Aritmética de Ponto Flutuante

- O arredondamento não é muito utilizado, pois mesmo acarretando erros menores, ele aumenta o tempo de execução de um programa.
- Mesmo que  $x$  e  $y$  estejam representados de forma exata, a soma  $x + y$ , por exemplo, também gera erros numéricos.
- **Exemplo:** Sejam  $x = 0.234 \times 10^5$  e  $y = 0.567 \times 10^2$  em uma máquina  $\mathcal{F}(10, 3, -5, 5)$ . Então

$$\begin{aligned}x + y &= 0.234 \times 10^5 + 0.567 \times 10^2 \\&= 0.234 \times 10^5 + 0.000567 \times 10^5 \\&= (0.234 + 0.000567) \times 10^5 \\&= 0.234567 \times 10^5 \\&= 0.234 \times 10^5 \quad (\text{se truncarmos}) \\&= 0.235 \times 10^5 \quad (\text{se arredondarmos})\end{aligned}$$



# **Análise de Erros nas Operações de Ponto Flutuante**

---

## Erros nas Operações Aritméticas de Ponto Flutuante

Já vimos que o erro relativo no resultado de uma operação (supondo que as parcelas são representadas exatamente) será

$$|ER_{op}| < 10^{-t+1} \quad (\text{Para o caso de Truncamento})$$

$$|ER_{op}| < \frac{1}{2} \times 10^{-t+1} \quad (\text{Para o caso de Arredondamento})$$

### Exemplo:

Considere uma máquina com  $t = 4$  e dois números  $x = 0.9765 \times 10^4$  e  $y = 0.2456 \times 10^2$ . Logo

$$\begin{aligned} x + y &= 0.9765 \times 10^4 + 0.2456 \times 10^2 \\ &= 0.9765 \times 10^4 + 0.002456 \times 10^4 = 0.978956 \times 10^4 \\ &= \begin{cases} 0.9789 \times 10^4 & \text{Se for usado truncamento} \\ 0.9790 \times 10^4 & \text{Se for usado arredondamento} \end{cases} \end{aligned}$$

# Erros nas Operações Aritméticas de Ponto Flutuante

## Análise de Erros - Adição

Sejam  $x$  e  $y$  tais que

$$x = \bar{x} + EA_x$$

$$y = \bar{y} + EA_y$$

Logo

$$\begin{aligned}x + y &= (\bar{x} + EA_x) + (\bar{y} + EA_y) \\ &= (\bar{x} + \bar{y}) + (EA_x + EA_y)\end{aligned}$$

Assim

$$\boxed{EA_{x+y} = EA_x + EA_y}.$$

**Atenção:** Aqui discutimos como o erro na representação de cada parcela da soma influencia o resultado final da operação.

## Erros nas Operações Aritméticas de Ponto Flutuante

Por outro lado, temos

$$\begin{aligned} ER_{x+y} &= \frac{EA_{x+y}}{\bar{x} + \bar{y}} = \frac{EA_x + EA_y}{\bar{x} + \bar{y}} = \frac{EA_x}{\bar{x} + \bar{y}} + \frac{EA_y}{\bar{x} + \bar{y}} \\ &= \frac{EA_x}{\bar{x}} \left( \frac{\bar{x}}{\bar{x} + \bar{y}} \right) + \frac{EA_y}{\bar{y}} \left( \frac{\bar{y}}{\bar{x} + \bar{y}} \right) \\ &= ER_x \left( \frac{\bar{x}}{\bar{x} + \bar{y}} \right) + ER_y \left( \frac{\bar{y}}{\bar{x} + \bar{y}} \right) \end{aligned}$$

## Análise de Erros - Subtração

É fácil ver que

$$EA_{x-y} = EA_x - EA_y$$

$$\begin{aligned} ER_{x-y} &= \frac{EA_{x-y}}{\bar{x} - \bar{y}} \\ &= \frac{EA_x - EA_y}{\bar{x} - \bar{y}} \\ &= ER_x\left(\frac{\bar{x}}{\bar{x} - \bar{y}}\right) - ER_y\left(\frac{\bar{y}}{\bar{x} - \bar{y}}\right) \end{aligned}$$

# Erros nas Operações Aritméticas de Ponto Flutuante

## Análise de Erros - Multiplicação

Notemos que

$$\begin{aligned}xy &= (\bar{x} + EA_x)(\bar{y} + EA_y) = \bar{x} \bar{y} + \bar{x}EA_y + \bar{y}EA_x + EA_xEA_y \\ &\approx \bar{x} \bar{y} + \bar{x}EA_y + \bar{y}EA_x \quad \text{pois} \quad EA_xEA_y \rightarrow 0\end{aligned}$$

Logo

$$EA_{xy} \approx \bar{x}EA_y + \bar{y}EA_x$$

E o Erro Relativo ...

$$\begin{aligned}ER_{xy} &= \frac{EA_{xy}}{\bar{x} \bar{y}} \\ &\approx \frac{\bar{x}EA_y + \bar{y}EA_x}{\bar{x} \bar{y}} = \frac{EA_x}{\bar{x}} + \frac{EA_y}{\bar{y}} = ER_x + ER_y\end{aligned}$$

Logo

$$ER_{xy} = ER_x + ER_y$$

# Erros nas Operações Aritméticas de Ponto Flutuante

## Análise de Erros - Divisão

$$\begin{aligned}\frac{x}{y} &= \frac{\bar{x} + EA_x}{\bar{y} + EA_y} \\ &= \frac{\bar{x} + EA_x}{\bar{y}(1 + \frac{EA_y}{\bar{y}})} \\ &= \frac{\bar{x} + EA_x}{\bar{y}} \left( \frac{1}{1 + \frac{EA_y}{\bar{y}}} \right)\end{aligned}$$

Agora note que

$$\begin{aligned}\frac{1}{1 + \frac{EA_y}{\bar{y}}} &= \sum_{n=0}^{\infty} (-1)^n \left( \frac{EA_y}{\bar{y}} \right)^n \quad (\text{Note que } EA_y \rightarrow 0) \\ &\approx 1 - \frac{EA_y}{\bar{y}}\end{aligned}$$

# Erros nas Operações Aritméticas de Ponto Flutuante

## Análise de Erros - Divisão

Logo

$$\begin{aligned}\frac{x}{y} &\approx \frac{\bar{x} + EA_x}{\bar{y}} \left(1 - \frac{EA_y}{\bar{y}}\right) \\ &\approx \frac{\bar{x}}{\bar{y}} + \frac{EA_x}{\bar{y}} - \frac{\bar{x}EA_y}{\bar{y}^2} - \frac{EA_x EA_y}{\bar{y}^2} \\ &\approx \frac{\bar{x}}{\bar{y}} + \frac{EA_x}{\bar{y}} - \frac{\bar{x}EA_y}{\bar{y}^2} \quad (EA_x EA_y \rightarrow 0)\end{aligned}$$

Logo

$$EA_{x/y} \approx \frac{EA_x}{\bar{y}} - \frac{\bar{x}EA_y}{\bar{y}^2} = \frac{\bar{y}EA_x - \bar{x}EA_y}{\bar{y}^2}$$

Assim,

$$EA_{x/y} \approx \frac{\bar{y}EA_x - \bar{x}EA_y}{\bar{y}^2}$$



# Erros nas Operações Aritméticas de Ponto Flutuante

## Análise de Erros - Divisão

E finalmente,

$$\begin{aligned} ER_{x/y} &= \frac{EA_{x/y}}{\bar{x}/\bar{y}} \\ &\approx \left( \frac{\bar{y}EA_x - \bar{x}EA_y}{\bar{y}^2} \right) \frac{\bar{y}}{\bar{x}} \\ &\approx \frac{EA_x}{\bar{x}} - \frac{EA_y}{\bar{y}} \\ &\approx ER_x - ER_y \end{aligned}$$

Logo

$$\boxed{ER_{x/y} \approx ER_x - ER_y}$$

# Erros nas Operações Aritméticas de Ponto Flutuante

## Exemplo

Sejam  $x, y, z, t$  representados exatamente. Qual o erro relativo total cometido no cálculo de

$$u = (x + y)z - t \quad ?$$

$$\begin{aligned} ER_{x+y} &= ER_x\left(\frac{\bar{x}}{\bar{x} + \bar{y}}\right) + ER_y\left(\frac{\bar{y}}{\bar{x} + \bar{y}}\right) + RA \\ &= 0 + 0 + RA \end{aligned}$$

Assim,

$$\begin{aligned} ER_{(x+y)z} &= ER_{x+y} + ER_z + RA \\ &= RA + 0 + RA \\ &= 2RA \end{aligned}$$

---

$RA$  é o erro relativo cometido na aproximação após a operação.

## Erros nas Operações Aritméticas de Ponto Flutuante

Para simplificar a notação, faremos  $m = (x + y)z$ .

$$\begin{aligned}ER_{m-t} &= ER_m\left(\frac{\bar{m}}{\bar{m} - \bar{t}}\right) - ER_t\left(\frac{\bar{t}}{\bar{m} - \bar{t}}\right) + RA \\&= ER_m\left(\frac{\bar{m}}{\bar{m} - \bar{t}}\right) - 0 + RA \\&= 2RA\left(\frac{\bar{m}}{\bar{m} - \bar{t}}\right) + RA\end{aligned}$$

Logo

$$\begin{aligned}|ER_{m-t}| &= \left|2RA\left(\frac{\bar{m}}{\bar{m} - \bar{t}}\right) + RA\right| \\&\leq 2|RA|\left|\frac{\bar{m}}{\bar{m} - \bar{t}}\right| + |RA| \\&\leq 10^{-t+1}\left(\left|\frac{\bar{m}}{\bar{m} - \bar{t}}\right| + \frac{1}{2}\right)\end{aligned}$$

## Efeitos Numéricos

---

## Efeitos Numéricos - Cancelamento Subtrativo

O erro relativo na subtração é dado por

$$ER_{x-y} = ER_x\left(\frac{\bar{x}}{\bar{x} - \bar{y}}\right) - ER_y\left(\frac{\bar{y}}{\bar{x} - \bar{y}}\right)$$

### Exemplo

Considere  $\bar{x} = 0.2357 \times 10^3$ ,  $\bar{y} = 0.2353 \times 10^3$  e  $\bar{r} = 0.4537 \times 10^3$ . Vamos analisar  $w = (x - y)r$  em uma máquina com quatro algarismos significativos. Temos  $\bar{x} - \bar{y} = 0.0004 \times 10^3$ . Note que

$$\begin{aligned} |ER_{x-y}| &= \left| ER_x\left(\frac{\bar{x}}{\bar{x} - \bar{y}}\right) - ER_y\left(\frac{\bar{y}}{\bar{x} - \bar{y}}\right) \right| \\ &< \left( \left| \frac{\bar{x}}{\bar{x} - \bar{y}} \right| + \left| \frac{\bar{y}}{\bar{x} - \bar{y}} \right| \right) \times \frac{1}{2} \times 10^{-3} \quad (t = 4) \\ &< \left( \frac{0.2357 \times 10^3 + 0.2353 \times 10^3}{0.0004 \times 10^3} \right) \times \frac{1}{2} \times 10^{-3} \\ &< 0.5888 \approx 59\% \end{aligned}$$

Agora, note que

$$\begin{aligned} |ER_w| &= |ER_{x-y} + ER_r| \\ &\leq |ER_{x-y}| + |ER_r| \\ &\leq 0.59 + \frac{1}{2} \times 10^{-3} \\ &\leq 0.5905 \approx 59\% \end{aligned}$$

## Efeitos Numéricos - Cancelamento Subtrativo

### Exemplo

Calcular, usando uma máquina com 10 dígitos,

$$\sqrt{9876} - \sqrt{9875}.$$

Temos que

$$\sqrt{9876} = 0.9937806599 \times 10^2 \quad \text{e} \quad \sqrt{9875} = 0.9937303457 \times 10^2$$

Portanto

$$\begin{aligned}\sqrt{9876} - \sqrt{9875} &= 0.0000503142 \times 10^2 \\ &= 0.5031420000 \times 10^{-4}\end{aligned}$$

Podemos obter um resultado mais preciso? Neste caso a resposta é sim! Basta notar que

Logo

$$\sqrt{x} - \sqrt{y} = \frac{x - y}{\sqrt{x} + \sqrt{y}}$$

$$\sqrt{9876} - \sqrt{9875} = \frac{1}{\sqrt{9876} + \sqrt{9875}} = 0.5031418679 \times 10^{-4}$$