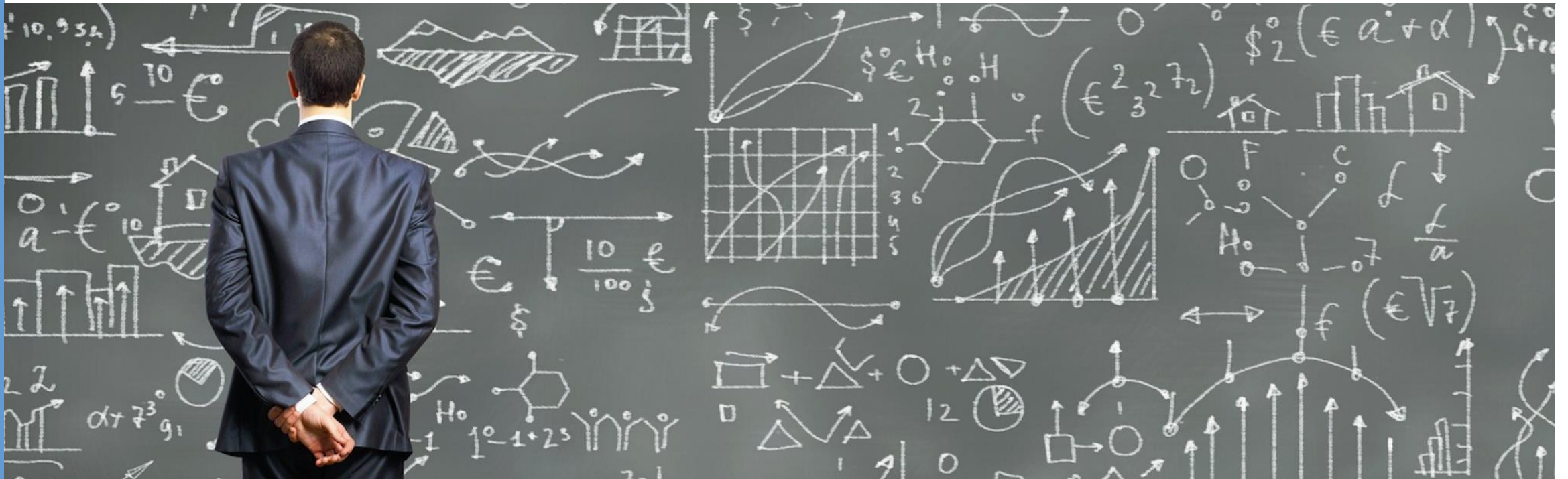




UNIVERSIDADE
FEDERAL DO CEARÁ

Métodos Numéricos



Unidade I: Ponto Flutuante e Erros



Representação dos Números

- Os números empregados no cálculo computacional podem ser de dois tipos: **números inteiros** e **números em “ponto flutuante”** (números reais da matemática, por exemplo $3.56 \rightarrow 0.356 \times 10^1$).
- Os computadores atuais **representam os números internamente no formato binário**, como uma sequência de 0s e 1s.
- Apesar dessa representação ser **conveniente para as máquinas** é **antinatural para os seres humanos**, cujo sistema de numeração é o decimal.
 - Obs. No passado o nosso sistema de numeração já foi também na base 12 (ex. contar nas falanges dos dedos) na base 60 (ex. sistema horário).



Representação em Ponto Flutuante

- A representação de números reais mais utilizada em máquinas é a do **ponto flutuante**.
- Esse número tem três partes: o sinal, a parte fracionária (mantissa) e o expoente:

$$m = \pm , d_1 d_2 d_3 \dots d_t \times \beta^e$$

- Sendo:
 - d_i 's : dígitos da parte fracionária, $d_1 \neq 0$, $0 \leq d_i \leq \beta - 1$.
 - β : base (em geral 2, 10 ou 16).
 - t : número de dígitos na mantissa.
 - e : expoente inteiro.



Exemplos

- $x=34,2$ (decimal)
- $\beta=10$
- $t=4$

$$x = 0,3420 \times 10^2$$

- $x=0,1$ (decimal)
- $\beta=2$
- $t=9$

$$x = 0,110011001 \times 2^{-3}$$

Obs: $0,1_{10} = 0,0001100110011 \dots_2$



Representação em Ponto Flutuante

- Nas máquinas digitais, **um dígito binário é denominado BIT** (do inglês, binary digit).
- Um grupo de **oito bits** corresponde a **1 byte**.
- Dessa forma, percebemos que a representação dos números binários num computador é feita com um **número finito de bits**.
- A esse tamanho finito de bits é dado o nome **palavra de computador**.
- O tamanho da palavra do computador **depende de características internas** à arquitetura do mesmo.



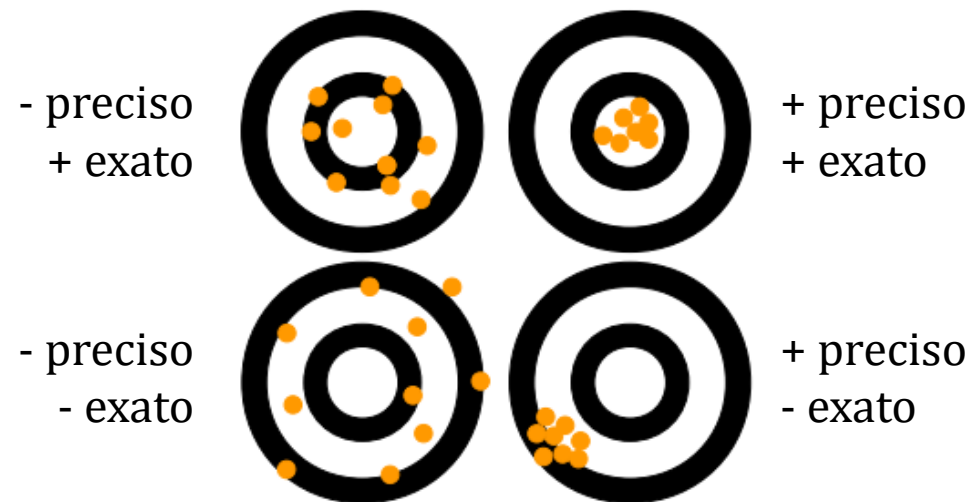
Teoria dos Erros

- Na busca de uma solução do modelo matemático por meio do cálculo numérico **temos o surgimento de erros através de diversas fontes.**
- Além disso, **toda medida experimental apresenta uma incerteza** e desta forma a solução do problema pode ser influenciada pela mesma.
- Como consequência, **métodos numéricos podem chegar a resultados distantes do que se esperaria** ou mesmo fornecer respostas sem **nenhuma relação com a solução** do problema original.

Teoria dos Erros

■ Precisão x Exatidão

- Conceitos erroneamente tratados como sinônimos no cotidiano.
 - **Exatidão:** Grau de concordância entre o resultado de uma medição e um valor verdadeiro mensurado.
 - **Precisão:** Grau de concordância entre resultados de medição obtidos sob as mesmas condições.





Teoria dos Erros

- Agora analisamos uma aplicação prática desses conceitos em números.
 - Exatidão
 - É governada pelos **erros no método numérico empregado**. Assim, se os números $\pi_1 = 3,1416304958$ e $\pi_2 = 3,1415809485$ almejam ambos a representar o número $\pi = 3,141592654 \dots$, o número π_2 possui maior exatidão que π_1 , embora ambos possuam a mesma precisão.
 - Precisão
 - A precisão de um número é governada pelo **número de dígitos empregados** na representação e na álgebra. Assim, a constante π será representada com maior precisão utilizando 8 bytes do que utilizando 4 bytes para armazenar o número.



Teoria dos Erros

- Para medir a acurácia de um número podemos utilizar como ferramentas de cálculo o conceito de erros absoluto e relativo.
 - **Erro Absoluto:** Diferença entre o valor exato de um número e o seu valor aproximado.

$$EA_x = \Delta x = |x_{exato} - x_{aprox}| = |x - \bar{x}|$$

- **Erro Relativo:** Trata-se do erro absoluto dividido pelo valor verdadeiro.

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



Teoria dos Erros

- Dado um número x já na forma normalizada que não possua representação exata no sistema $F[b, n, e_{min}, e_{max}]$. Pode-se escrever x como:

$$x = (0, d_1 d_2 \dots d_n) \times b^e + g_x \times b^{e-n}, \text{ com } 0 \leq g_x < 1$$

- Onde g_x é a parcela que não pode ser incluída em sua representação.
- Dessa forma, existem duas formas de realizarmos a aproximação.
 - Truncamento
 - Arredondamento



Teoria dos Erros

- Truncamento

- Consistem em simplesmente ignorar o g_x . Assim,

$$\bar{x} = (\mathbf{0}, d_1 d_2 \dots d_n) \times b^e$$

- Arredondamento

- No arredondamento, executa-se a seguinte operação,

$$\bar{x} = \begin{cases} (\mathbf{0}, d_1 d_2 \dots d_n) \times b^e, & \text{se } |g_x| < 1/2 \\ (\mathbf{0}, d_1 d_2 \dots (d_n + 1)) \times b^e, & \text{se } |g_x| \geq 1/2 \end{cases}$$



Teoria dos Erros

- Dado um sistema em ponto flutuante fictício com 4 algarismos na mantissa e base 10.
- Ex. $x = 0,037 \times 10^4$, $y = 0,1272 \times 10^2$, calcule $x+y$.

No procedimento da adição em ponto flutuante devemos alinhar as casas decimais de ambos os números igualando os expoentes ao maior expoente presente na soma.



Teoria dos Erros

- Dado uma sequência de operações algébricas é importante observar **como o erro se propaga** ao longo destas operações consecutivas.
- Em um sistema de representação em ponto flutuante qualquer, **a soma de dois números exatos fornecerá um resultado exato?**



Teoria dos Erros

- $x_1 = 0,3491 \times 10^4$
- $x_2 = 0,2345 \times 10^0$

- $(x_2 + x_1) - x_1$

$$= (0,2345 \times 10^0 + 0,3491 \times 10^4) - 0,3491 \times 10^4$$

$$= 0,3491 \times 10^4 - 0,3491 \times 10^4$$

$$= 0,0000$$

- $x_1 = 0,3491 \times 10^4$
- $x_2 = 0,2345 \times 10^0$

- $x_2 + (x_1 - x_1)$

$$= 0,2345 \times 10^0 + (0,3491 \times 10^4 - 0,3491 \times 10^4)$$

$$= 0,2345 \times 10^0 - 0,0000$$

$$= 0,2345$$



Análise dos Erros

- Adição
 - Erro Absoluto

$$x + y = (\bar{x} + EA_x) + (\bar{y} + EA_y) = (\bar{x} + \bar{y}) + (EA_x + EA_y)$$

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



Análise dos Erros

- Adição
 - Erro Relativo

$$ER_{x+y} = \frac{EA_{x+y}}{x+y} = \frac{EA_x}{x} \cdot \frac{x}{x+y} + \frac{EA_y}{y} \cdot \frac{y}{x+y}$$

$$ER_{x+y} = ER_x \cdot \frac{x}{x+y} + ER_y \cdot \frac{y}{x+y}$$

$$ER_{x+y} = \frac{ER_x}{1 + \frac{y}{x}} + \frac{ER_y}{1 + \frac{x}{y}}$$



Análise dos Erros

- Adição
 - Erro Relativo
 - $x \gg y$

$$ER_{x+y} = \frac{ER_x}{1 + \frac{y}{x}} + \frac{ER_y}{1 + \frac{x}{y}}$$

$$ER_{x+y} \approx ER_x$$



Análise dos Erros

- Adição
 - Erro Relativo
 - $x \ll y$

$$ER_{x+y} = \frac{ER_x}{1 + \frac{y}{x}} + \frac{ER_y}{1 + \frac{x}{y}}$$

$$ER_{x+y} \approx ER_y$$



Análise dos Erros

- Adição
 - Erro Relativo
 - $x \approx y$

$$ER_{x+y} = \frac{ER_x}{1 + \frac{y}{x}} + \frac{ER_y}{1 + \frac{x}{y}}$$

$$ER_{x+y} \approx \frac{1}{2} (ER_x + ER_y)$$



Análise dos Erros

- Adição
 - Erro Relativo

$$ER_{x+y} \approx \text{Max} (ER_x, ER_y)$$



Análise dos Erros

- Subtração
 - Erro Absoluto

$$x - y = (\bar{x} + EA_x) - (\bar{y} + EA_y) = (\bar{x} + \bar{y}) - (EA_x + EA_y)$$

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



Análise dos Erros

- Subtração
 - Erro Relativo

$$ER_{x-y} = \left| \frac{EA_{x-y}}{x-y} \right| = \left| \frac{EA_x}{x} \cdot \frac{x}{x-y} + \frac{EA_y}{y} \cdot \frac{y}{x-y} \right|$$

$$ER_{x-y} = \left| ER_x \cdot \frac{x}{x-y} + ER_y \cdot \frac{y}{x-y} \right|$$

$$ER_{x-y} = \left| \frac{ER_x}{1 - \frac{y}{x}} + \frac{ER_y}{1 - \frac{x}{y}} \right|$$



Análise dos Erros

- Adição
 - Erro Relativo
 - $x \gg y$

$$ER_{x-y} = \left| \frac{ER_x}{1 - \frac{y}{x}} + \frac{ER_y}{1 - \frac{x}{y}} \right|$$

$$ER_{x-y} \approx ER_x$$



Análise dos Erros

- Adição
 - Erro Relativo
 - $x \ll y$

$$ER_{x-y} = \left| \frac{ER_x}{1 - \frac{y}{x}} + \frac{ER_y}{1 - \frac{x}{y}} \right|$$

$$ER_{x-y} \approx ER_y$$



Análise dos Erros

- Adição
 - Erro Relativo
 - $x \approx y$

$$ER_{x-y} = \left| \frac{ER_x}{1 - \frac{y}{x}} + \frac{ER_y}{1 - \frac{x}{y}} \right|$$

$$\text{Se } 1 - \frac{y}{x} \ll 1 \text{ e } 1 - \frac{x}{y} \ll 1$$

$$ER_{x+y} \gg (ER_x + ER_y)$$

Este resultado mostra claramente como o erro relativo pode se tornar muito grande quando $X \approx Y$. Isto ocorre porque a subtração de dois números muito próximos entre si resulta em um número cuja representação ocorre nos últimos dígitos da mantissa, resultando em um grande erro de arredondamento.



Análise dos Erros

- Desafio
 - Encontre as expressões de erros para a multiplicação e para a divisão.



Exercícios

- Seja um sistema que opera em aritmética de ponto flutuante de $t = 4$, na base 10, calcule os erros absolutos e relativos por truncamento e arredondamento:
 - a) 123,456
 - b) $374,3 + 3,345$
 - c) $124,34 + 0,1234$
 - d) $22,12 \times 0,123$
 - e) $0,3212 \times 12,32$