

# **Algoritmos Numéricos**

## Uma abordagem direta

BERNARDO SUNDERHUS

29 de maio de 2016



# Sumário

<b>1</b>	<b>Definições em computação numérica</b>	<b>3</b>
1.1	Representação numérica . . . . .	3
1.1.1	Números inteiros e positivos . . . . .	3
1.1.2	Números em geral . . . . .	4
1.2	Bases importantes para computação . . . . .	4
1.3	Algumas trocas de bases . . . . .	4
1.3.1	Números inteiros . . . . .	4
1.3.2	Números fracionários . . . . .	5
1.4	Aritmética de ponto flutuante . . . . .	6
1.4.1	Um sistema de ponto flutuante pode ser representado por: . . . . .	7
1.4.2	Erros . . . . .	8
<b>2</b>	<b>Resolução de Sistemas Lineares</b>	<b>11</b>
2.1	Métodos Diretos . . . . .	11



# Prefácio

## Structure of book

Add short description about each chapter in this book.

## Acknowledgements

- Um agradecimento especial ao professor Thomas W. Rauber que disponibilizou suas anotações pessoais.



# 1

## Definições em computação numérica

### 1.1 Representação numérica

É possível representar um número sendo um sistema contendo um polinómio e uma base

$$\begin{aligned} \text{número} &= \text{senal} \times (a_j, a_{j-1}, \dots, \bullet, a_{-1}, a_{-2}, \dots, a_{-k})_B \\ \text{senal} &\in \{-1, 1\} \\ B &\in \mathbb{N}, B \geq 2 \end{aligned}$$

Do que foi estabelecido a cima podemos gerar

$$\begin{aligned} \text{número} &= \text{senal} \times a_j B^j + a_{j-1} B^{j-1} + \dots + a_0 B^0 + a_{-1} B^{-1} + \dots + a_{-k} B^{-k} \\ &= \text{senal} \times \sum_{i=-k}^j a_i B^i \end{aligned}$$

#### 1.1.1 Números inteiros e positivos

$$k = 0, \text{ senal} = +1$$

Exemplo com base  $B = 10$

$$\begin{aligned} (1\ 2\ 3) &= 1 \times 10^2 + 2 \times 10^1 + 3 \times 10^0 \\ (4\ 5\ 6) &= 4 \times 10^2 + 5 \times 10^1 + 6 \times 10^0 \\ (1\ 0\ 0\ 1) &= 1 \times 10^3 + 0 \times 10^2 + 0 \times 10^1 + 1 \times 10^0 \end{aligned}$$

### 1.1.2 Números em geral

$$-(1\ 2\ 3\ .\ 5)_{10} = -(1 \times 10^2 + 2 \times 10^1 + 3 \times 10^0 + 4 \times 10^{-1} + 5 \times 10^{-2})$$

$$(A\ .\ B\ C)_{16}^1 = 10 \times 16^0 + 11 \times 16^{-1} + 12 \times 16^{-2}$$

## 1.2 Bases importantes para computação

Base	Name	Digits
2	binary	0,1
8	octal	0,...,7
10	decimal	0,...,9
16	hexadecimal	0,...9,A,...,F

## 1.3 Algumas trocas de bases

$$(1\ 0\ 1)_2 = (5)_{10}$$

$$(A\ B)_{16} = (1\ 7\ 1)_{10}$$

### 1.3.1 Números inteiros

número  $\in \mathbb{Z}$

**Decimal  $\longrightarrow$  binário**

Dado um número decimal  $N$  e sendo  $b$  o equivalente binário que queremos

**while**  $N \neq 0$  **do**

$b \leftarrow b + N \% 2$ <sup>2</sup>

$N \leftarrow N \setminus 2$ <sup>3</sup>

**end while**

---

<sup>1</sup>Em hexadecimal os números vão de 0,...9,A,...F

<sup>2</sup>resto da divisão

<sup>3</sup>a divisão exata ou divisão inteira é aquela que assume valores inteiros como resposta. Caso a divisão não gere um número inteiro, esse valor extra é descartado da resposta e é chamada de resto da divisão. ex:  $3/2 = 2$ , *resto* = 1



exemplo:

```

 $N = 23_{10}$ 
 $N = N \setminus 2 = 11$  ( $N \% 2 = 1 \longrightarrow b_0 = 1$ )
 $N = N \setminus 2 = 5$  ( $N \% 2 = 1 \longrightarrow b_1 = 1$ )
 $N = N \setminus 2 = 2$  ( $N \% 2 = 0 \longrightarrow b_2 = 0$ )
 $N = N \setminus 2 = 1$  ( $N \% 2 = 1 \longrightarrow b_3 = 1$ )
 $N = N \setminus 2 = 0$  ( $N \% 2 = 0 \longrightarrow b_4 = 0$ )
return  $b = (10111)_2$ 

```

### Binário $\longrightarrow$ decimal

Dado um texto/vetor binário  $\beta$  e  $j$  o comprimento de  $\beta$

```

 $N \leftarrow \beta(j)$ 
for  $i = j - 1, \dots, 0$  do
   $N \leftarrow 2 \times N + \beta(i)$ 
end for

```

exemplo:

```

 $\beta = (10111)_2 \rightarrow j = 4$ 
 $N_4 \leftarrow 1$ 
 $N_3 \leftarrow 0 + 2 \times 1 \rightarrow 2$ 
 $N_2 \leftarrow 1 + 2 \times 2 \rightarrow 5$ 
 $N_1 \leftarrow 1 + 2 \times 5 \rightarrow 11$ 
 $N_0 \leftarrow 1 + 2 \times 11 \rightarrow 23$ 
return  $N_0 = N = 23_{10}$ 

```

### 1.3.2 Números fracionários

Os métodos a baixo servem para computar somente a parte fracionária, eles assumem que o entrada do algoritmo só tem valor fracionário.

### Decimal $\longrightarrow$ binário

Dado um número natural finito  $N$  e um número maximo de iterações  $i$

```

 $b \leftarrow "0."$ 
 $aux \leftarrow 0$ 
while  $N \neq 0$  and  $aux < i$  do
   $N2 \leftarrow N \times 2$ 
   $digit \leftarrow N2 \setminus 1$ 
   $b \leftarrow b + digit$ 
   $N \leftarrow N2 - digit$ 
   $aux \leftarrow aux + 1$ 

```

```

end while
return b

```

#### Binária $\longrightarrow$ decimal

```

Dado uma string/vetor binário b e j o tamanho de b.
N  $\leftarrow$  b[1]
for i = 2, ..., j do
    N = 2 * N + b[i]
end for
return N

```

## 1.4 Aritmética de ponto flutuante

$$x = sinal \times (. d_1, \dots, d_t) \times B^e$$

Representação do subconjunto de  $\mathbb{R}$ .

Tal representação é impossível usando um número  $t$  finito de dígitos fracionários, o que gera erro na representação e nos cálculos.

$$\begin{aligned}
 Sinal &\in \{-1, 1\} \\
 d_1, \dots, d_t &: \text{mantissa de } t \text{ dígitos} \\
 d_j &\in \{0, \dots, B-1\} \quad j = 1, \dots, t \\
 d_1 &\neq 0
 \end{aligned}$$

Exemplo:

$$\begin{aligned}
 B = 10, t = 3, e &\in \{-5, 5\} \\
 \longrightarrow x &= \pm 0 . d_1 d_2 d_3 \times 10^e
 \end{aligned}$$

Maior e menor número representável:

$$\begin{aligned}
 menor &= 0.1000 \times 10^5 \rightarrow 0.000001 \\
 Maior &= 0.999 \times 10^5 \rightarrow 99900
 \end{aligned}$$

Possíveis casos para  $x$

Seja  $G = \{x \in \mathbb{R} | m \leq |x| \leq M \text{ e } x \in \mathbb{R}\}$

1.  $x \in G$

a) Representação exata (sem erros)

Ex:  $1.23_{10} \longrightarrow 0.123 \times 10^1$

b) Representação aproximada

Ex:  $456.789_{10} \longrightarrow 0.456789 \times 10^3 \xrightarrow{t=3} 0.457 \times 10^3$  arredondado.

2.  $|x| < m$

Ex:  $0.0000000345 = 0.345 \times 10^7 \xrightarrow{e \in \{-5,5\}} \text{Não representável, "under flow"}$ .

3.  $|x| > M$

Ex:  $875000000 = 0.875 \times 10^9 \xrightarrow{e \in \{-5,5\}} \text{Não representável "over flow"}$ .

#### 1.4.1 Um sistema de ponto flutuante pode ser representado por:

$F(B, t, l, u) = F(\text{Base}, \text{número de dígitos}, \text{expoente menor}, \text{expoente maior})$

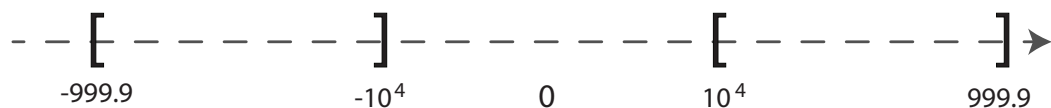
Ex:  $F(10, 4, -3, 3)$

$menor = 0.1 \times 10^3 \longrightarrow 10^{-4}$

$maior = 0.9999 \times 10^3 \longrightarrow 999.9$

região de under flow:  $(-10^4, 0) \cup (0, 10^{-4})$

região de over flow:  $(-\infty, -999.9) \cup (999.9, \infty)$



\*IEEE 754 :  $F(2, 23, -127, 128)$

### 1.4.2 Erros

#### Erros absolutos e relativos

$x$  : Valor exato

$\bar{x}$  : Valor aproximado de  $x$

*Def* : Erro absoluto  $EA_x = x - \bar{x}$

$$|EA_x| \geq 0$$

*Def* : Erro relativo  $ER_x = \frac{EA_x}{\bar{x}}$

Exemplo:

$$x \in (2112.8, 2113) \quad \bar{x} = 2112.9$$

$$y \in (5.2, 5.3) \quad \bar{y} = 5.3$$

$$|EA_x| < 0.1$$

$$|ER_x| < 4.7 \times 10^{-5}$$

$$|EA_y| < 0.1$$

$$|ER_y| < 0.02$$

Representação de  $x$  por  $\bar{x}$  é relativamente mais precisa do que  $y$  por  $\bar{y}$

#### Erro de arredondamento e truncamento

$B = 10, t$  dígitos

$$x = 0.d_1d_2\dots d_t \times 10^e$$

$$= 0.d_1d_2\dots d_{t-1} \times 10^e$$

$$+ 0.000\dots 0d_t \times 10^e$$

$$= 0.d_1d_2\dots d_{t-1} \times 10^e + 0.d_t \times 10^{e-t}$$

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

Exemplo:

$$t = 4, x = 234.57$$

$$= 0.2345 \times 10^3 + 0.7 \times 10^{-1}$$

$$= f_x \times 10^3 + g_x \times 10^{-1}$$

## 1. Truncamento

$$\begin{aligned}
g_x &\leftarrow 0, \bar{x} \leftarrow f_x \times 10^e \\
|EA_x| &< 10^{e-t} \\
|ER_x| &< 10^{-t+1}
\end{aligned}$$

## 2. Arredondamento

a)  $|g_x| < \frac{1}{2}$

$$\begin{aligned}
g_x &\leftarrow 0, \bar{x} = f_x \times 10^e, \\
|EA_x| &< \frac{1}{2} \times 10^{e-t} \\
|ER_x| &< \frac{1}{2} \times 10^{-t+1}
\end{aligned}$$

b)  $|g_x| \geq \frac{1}{2}$

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

$$\longrightarrow |EA_x| \leq \frac{1}{2} \times 10^{e-t}, |ER_x| < \frac{1}{2} \times 10^{-t+1}$$

**Propagação de erros em sequências de operações aritméticas**

## 1. Propagação dos erros dos operandos

## a) Adição

Ex:  $F(10, 4, -3, 3)$

$X = 1.3749213$  Valor exato

$\bar{X} = 0.1374 \times 10^1 \quad EA_x = X - \bar{X} = 0.0009213\dots$

$Y = 573.27482$

$\bar{Y} = 0.5732 \times 10^3 \quad EA_y = Y - \bar{Y} = 0.0007482$

$$X + Y = \overline{X} + \overline{Y} + EA_X + EA_Y$$

$\oplus$ : Operação de adição no sistema de ponto flutuante

$$\overline{X} \oplus \overline{Y} (EA_x + EA_y \text{ desconhecidos})$$

$$= \overline{\overline{X} + \overline{Y}}$$

$$\frac{0.0013 \times 10^3}{0.5732 \times 10^3} = 0.5745 \times 10^3 (\text{somente quatro casas na mantissa})$$

$$\overline{X} + \overline{Y} = 1.374 + 573.2 = 574.574$$

$$\overline{\overline{X} + \overline{Y}} = 574.5$$

$$EA_{\overline{x}+\overline{y}} = \overline{X} + \overline{Y} - \overline{\overline{X} + \overline{Y}} = 0.074$$

$$X + Y = 574.6497413 \text{ Vs}$$

b) Operações:

Operação	Erro abs.	Erro rel.
Adição	$EA_{x+y} = EA_x + EA_y$	$ER_{x+y} = ER_x \times \frac{\bar{x}}{\bar{x}+\bar{y}} + ER_y \times \frac{\bar{y}}{\bar{x}+\bar{y}}$
Subtração	$EA_{x-y} = EA_x - EA_y$	$ER_{x-y} = ER_x \times \frac{\bar{x}}{\bar{x}-\bar{y}} - ER_y \times \frac{\bar{y}}{\bar{x}-\bar{y}}$
Multiplicação	$EA_{xy} \approx \bar{x} \times EA_y + \bar{y} \times EA_x$	$ER_{xy} \approx ER_x + ER_y$
Divisão	$EA_{\frac{x}{y}} \approx \frac{\bar{y}EA_x - \bar{x}EA_y}{\bar{y}^2}$	$ER_{\frac{x}{y}} \approx ER_x - ER_y$

## 2

# Resolução de Sistemas Lineares

Sistema Linear: m equações lineares em n variáveis

Variável:  $x_j$   $j = 1, \dots, n$

Equação:  $i \rightarrow a_{i1}x_1 + a_{i2}x_2 + \dots + a_{in}x_n = b_i$

Em cálculo numérico levaremos em conta na maioria o caso  $m = n$

Em forma matricial temos:  $Ax = b$

com

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ \vdots & \vdots & \vdots & \vdots \\ a_{n1} & a_{n2} & \cdots & a_{nn} \end{bmatrix} \quad x = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad b = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix}$$

Problema em questão: Dado A e b determinar x.

## 2.1 Métodos Diretos

### Regra de Cramer

$$x_i = \frac{\det(A_i)}{\det(A)}$$

det: Determinante

$A_i$ : Substituição da coluna i por b em A

Esse método se torna extremamente complexo para  $n \gg 2$ , pois se houver  $n+1$  determinantes de ordem n, teremos:  $(n+1) \times n! \times (n-1)$  multiplicações. Com  $n = 20$  um computador que faça 100 milhões de multiplicações por segundo levaria 300.000 anos para computar.

## Inversa

$Ax = b \longrightarrow A^{-1}Ax = A^{-1}b \longrightarrow x = A^{-1}b$  caso  $A^{-1}$  exista.

Problemas: ineficiente e instável (graças ao erro de arredondamento existente nas inversões, que se acumulam).

## Resolução de SL Triangulares

Dado a matrix A triangular superior n x n com elementos diagonais  $a_{ii}$  não nulos.

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ 0 & a_{22} & \cdots & a_{2n} \\ 0 & 0 & \ddots & \vdots \\ 0 & 0 & \cdots & a_{nn} \end{bmatrix}$$

$$x_n = b_n / a_{nn}$$

**for**  $k = n - 1, \dots, 1$  **do**

    soma =  $b_k$

**for**  $j = k + 1, \dots, n$  **do**

        soma = soma -  $a_{kj}x_j$

**end for**

$$x_k = \frac{\text{soma}}{a_{kk}}$$

**end for**

Esforço computacional:

n divisões

$$\sum_{j=1}^{n-1} j = n(n-1)/2 \text{ adições}$$

$$\sum_{j=1}^{n-1} j = n(n-1)/2 \text{ multiplicações}$$

Gerando uma complexidade  $O(n^2)$  (a mesma do bubblesort)