

## - LIMITE DE UMA SEQUÊNCIA DE REAIS OU COMPLEXOS

1.2

\* DEFINIÇÃO: SGA

$$\{x_n\}_{n=1}^{\infty}$$

UMA SEQUÊNCIA INFINITA DE REAIS OU COMPLEXOS.

ESTA SEQUÊNCIA TEM O LIMITE  $x$  (i.e., CONVERGE PARA  $x$ ), SE, PARA QUALQUER  $\epsilon > 0$ , EXISTE UM NÚMERO POSITIVO INTEIRO  $N(\epsilon)$ , TAL QUE  $|x_n - x| < \epsilon$ , SEMPRE QUE  $n > N(\epsilon)$ .

\* NOTAÇÃO:  $\lim_{n \rightarrow \infty} x_n = x$  OU  $x_n \rightarrow x$  QUANDO  $n \rightarrow \infty$

## - CONVERGÊNCIA E CONTINUIDADE

\* TEOREMA: SE  $f$  É DEFINIDA EM  $X \subset \mathbb{R}$  E  $x_0 \in X$ , ENTÃO AS SEGUINTEs AFIRMAÇÕES SÃO EQUIVALENTES:

A)  $f$  É CONTÍNUA EM  $x_0$

B) SE  $\{x_n\}_{n=1}^{\infty}$  É QUALQUER SEQUÊNCIA EM  $X$  CONVERGENTE PARA  $x_0$ , ENTÃO  
$$\lim_{n \rightarrow \infty} f(x_n) = f(x_0)$$

OBSERVAÇÃO: ASSUMIREMOS QUE TODAS AS FUNÇÕES CONSIDERADAS EM ANÁLISE NUMÉRICA SÊTAM CONTÍNUAS.

## - DERIVADA

1.3

\* DEFINIÇÃO: SEJA  $f$  UMA FUNÇÃO DEFINIDA NUM INTERVALO ABERTO CONTENDO  $x_0$ .

$f$  É DIFERENCIÁVEL EM  $x_0$ , SE EXISTE

$$f'(x_0) = \lim_{x \rightarrow x_0} \frac{f(x) - f(x_0)}{x - x_0}$$

$f'(x_0)$  É CHAMADA DE DERIVADA DE  $f$  EM  $x_0$ .

$f$  É DITA DIFERENCIÁVEL EM  $X \subset \mathbb{R}$  SE TEM DERIVADA PARA CADA  $x \in X$ .

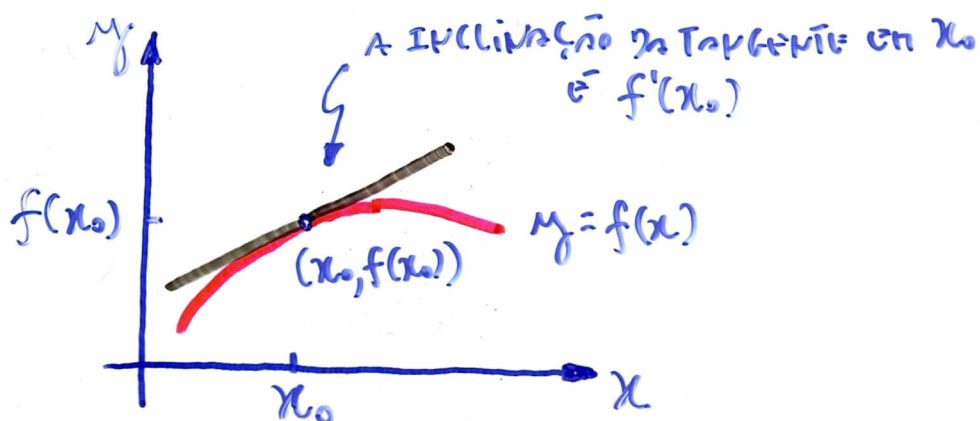
\* TEOREMA: SE  $f$  É DIFERENCIÁVEL EM  $x_0$ , ENTÃO  $f$  É CONTÍNUA EM  $x_0$ .

\* NOTAÇÃO:  $C^n(X)$  É CONJUNTO DE TODAS AS FUNÇÕES QUE TÊM  $n$  DERIVADAS CONTÍNUAS EM  $X \subset \mathbb{R}$

$C^\infty(X)$  É CONJUNTO DAS FUNÇÕES QUE TÊM DERIVADAS DE TODAS AS ORDENS EM  $X \subset \mathbb{R}$

## - TEOREMAS DE ROLLE, VALOR MÉDIO E VALOR EXTREMO 1.4

OBSERVAÇÃO: A DERIVADA DE  $f$  EM  $x_0$  É O VALOR DA INCLINAÇÃO DA LINHA TANGENTE À CURVA REPRESENTATIVA DE  $f$ , NO PONTO  $(x_0, f(x_0))$

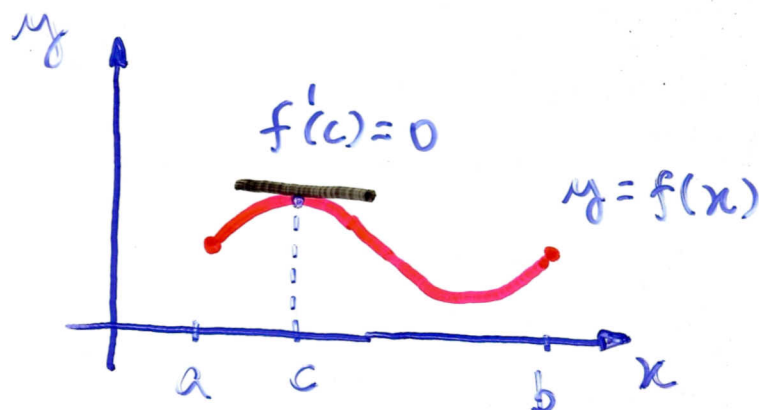


\* TEOREMA DE ROLLE:

SEJA  $f \in C[a, b]$  E DIFERENCIÁVEL EM  $(a, b)$ .

SE  $f(a) = f(b)$ , EXISTE UM NÚMERO  $c$  EM  $(a, b)$  PARA O QUAL  $f'(c) = 0$ .

OBSERVAÇÃO: PODE EXISTIR MAIS DE UM PONTO COM TAL PROPRIEDADE

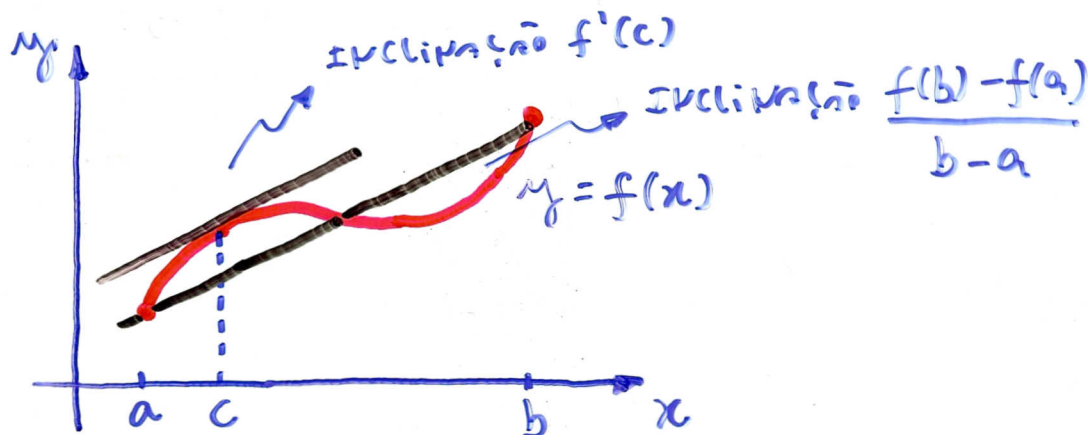


GENERALIZAÇÃO: SE  $f \in C[a, b]$  E  $f$  É DIFERENCIÁVEL n Vezes EM  $(a, b)$ , E ADÉMIS SE  $f(x)=0$  NOS  $n+1$  VALORES DISTINTOS  $x_0, \dots, x_n$  EM  $[a, b]$ , ENTÃO EXISTE  $c$  EM  $(a, b)$  TAL QUE  $f^{(n)}(c)=0$

# \* TEOREMA DO VALOR MÉDIO (DERIVADAS):

Se  $f \in C[a, b]$  e  $f$  é diferenciável em  $(a, b)$ , existe um número  $c$  no intervalo  $(a, b)$  tal que

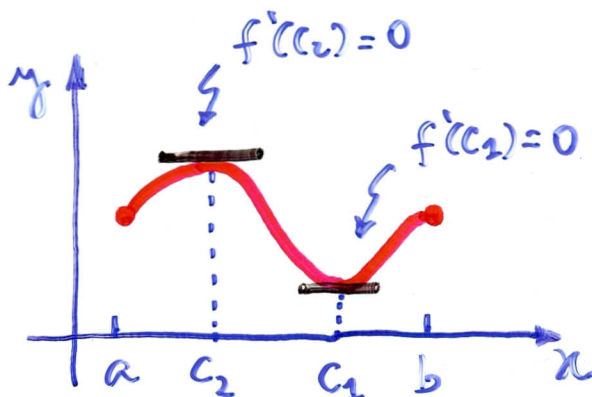
$$f'(c) = \frac{f(b) - f(a)}{b - a}$$



## \* TEOREMA DO VALOR EXTREMO:

Se  $f \in C[a, b]$ , então existem  $c_1$  e  $c_2$  em  $[a, b]$  tais que  $f(c_1) \leq f(x) \leq f(c_2)$  para todo  $x \in [a, b]$ .

Ademais, se  $f$  é diferenciável em  $(a, b)$ , os números  $c_1$  e  $c_2$  podem existir tanto nos extremos de  $[a, b]$  como onde  $f' = 0$ .





\* DEFINIÇÃO: INTEGRAL DE RIEMANN

A INTEGRAL DE RIEMANN DE  $f$  EM  $[a, b]$  É O SEGUINTE LIMITE, DESDE QUE ELE EXISTA:

$$\int_a^b f(x) dx = \lim_{\max \Delta x_i \rightarrow 0} \sum_{i=1}^n f(\xi_i) \Delta x_i$$

ONDE  $a = x_0 \leq x_1 \leq \dots \leq x_n = b$ , E ONDE

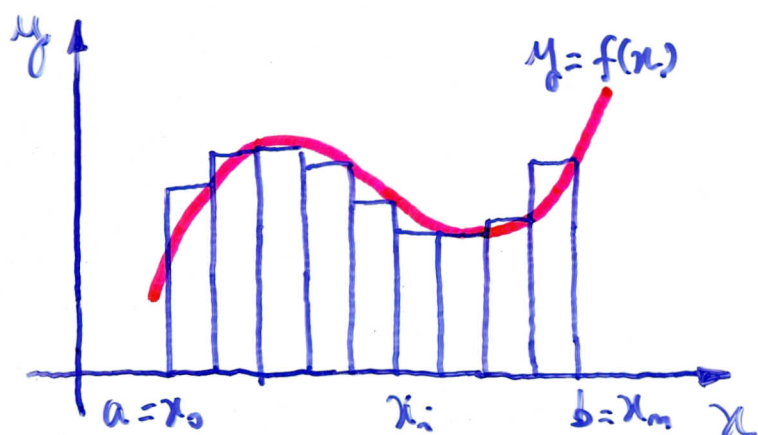
$\Delta x_i = x_i - x_{i-1}$  PARA CADA  $i = 1, 2, \dots, n$ , COM

$\xi_i$  ESCOLHIDO ARBITRARIAMENTE EM  $[x_{i-1}, x_i]$

OBSERVAÇÕES: i) Toda FUNÇÃO CONTÍNUA EM  $[a, b]$  ADMITE INTEGRAL DE RIEMANN NESSE INTERVALO.

ii) PARA MAIOR CONVENIÊNCIA COMPUTACIONAL, PODEMOS ESCOLHER OS PONTOS  $x_i$  ESPACADOS EM INTERVALOS IGUAIS DE  $[a, b]$ , E ESCOLHER  $\xi_i = x_i$  PARA CADA  $i = 1, 2, \dots, n$

NESSE CASO, 
$$\int_a^b f(x) dx = \lim_{n \rightarrow \infty} \frac{b-a}{n} \sum_{i=1}^n f(x_i)$$



$$x_i = a + \left(\frac{b-a}{n}\right)i$$

$$\therefore \Delta x_i = \frac{b-a}{n}$$

$$\lim_{\max \Delta x_i \rightarrow 0} \Delta x_i = \lim_{n \rightarrow \infty} \left(\frac{b-a}{n}\right)$$

# - GENERALIZAÇÃO DO TEOREMA DO VALOR MÉDIO (INTEGRAIS)

## \* TEOREMA DO VALOR MÉDIO Ponderado Para INTEGRAIS:

SUPONDO QUE EXISTA  $f \in C[a,b]$ , QUE EXISTA A INTEGRAL DE RIEMANN DE  $g$  EM  $[a,b]$ , E QUE  $g(x)$  NÃO TENHA DE SINAL EM  $[a,b]$ , ENTÃO EXISTE UM NÚMERO  $c$  EM  $(a,b)$  TAL QUE

$$\int_a^b f(x)g(x)dx = f(c) \int_a^b g(x)dx$$

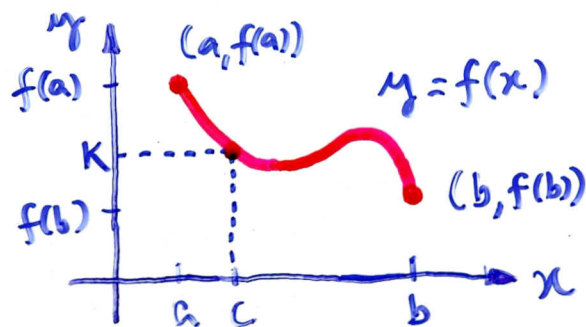
OBSERVAÇÃO: Quando  $g(x)=1$ , ESTE TEOREMA SE TORNA O TRADICIONAL TEOREMA DO VALOR MÉDIO PARA INTEGRAIS, QUE DÁ O VALOR MÉDIO DE  $f$  NO INTERVALO  $[a,b]$  COMO SENDO

$$f(c) = \frac{1}{b-a} \int_a^b f(x)dx$$

(NESTE CASO,  $c$  É O PONTO - INDETERMINADO - TAL QUE O VALOR DA FUNÇÃO ALI É O SEU VALOR MÉDIO NO INTERVALO  $[a,b]$ )

## \* TEOREMA DO VALOR INTERMEDIÁRIO:

SE  $f \in C[a,b]$  E  $K$  É QUALQUER NÚMERO ENTRE  $f(a)$  E  $f(b)$ , EXISTE UM NÚMERO  $c$  EM  $(a,b)$  PARA O QUAL  $f(c)=K$ .



# - Polinômios de Taylor

\* TEOREMA DE TAYLOR:

SUPONDO QUE  $f \in C^n[a, b]$ , QUE  $f^{(n+1)}$  EXISTA EM  $[a, b]$   
 E QUE  $x_0 \in [a, b]$ , EXISTE UM NÚMERO  $\xi(x)$  ENTRE  $x_0$  E  $x$ ,  
 PARA TODO  $x \in [a, b]$ , TAL QUE

$$f(x) = P_n(x) + R_n(x)$$

ONDE

$$P_n(x) = f(x_0) + f'(x_0)(x-x_0) + \frac{f''(x_0)}{2!}(x-x_0)^2 + \dots + \frac{f^{(n)}(x_0)}{n!}(x-x_0)^n$$

$$= \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x-x_0)^k$$

$$\text{E } R_n(x) = \frac{f^{(n+1)}(\xi(x))}{(n+1)!} (x-x_0)^{n+1}$$

$P_n(x)$  É CHAMADO POLINÔMIO DE TAYLOR DE  $n$ -ÉSIMO  
 GRAU PARA  $f(x)$ , CENTRADO EM  $x_0$

$R_n(x)$  É CHAMADO RESTO O ERRO DE TRUNCAMENTO  
 ASSOCIADO A  $P_n(x)$

A SÉRIE INFINITA OBTIDA TOMANDO-SE  $\lim_{n \rightarrow \infty} P_n(x)$  É  
 CHAMADA SÉRIE DE TAYLOR PARA  $f(x)$ , CENTRADA EM  $x_0$

QUANDO  $x_0 = 0$ , O POLINÔMIO E A SÉRIE DE TAYLOR TAMBÉM  
 SÃO CHAMADOS COMO DE MACLAURIN



## 2. ERROS DE ARREDONDAMENTO E ARITMÉTICA COMPUTACIONAL

### - ARITMÉTICA COMPUTACIONAL

⇒ ARITMÉTICA DE DÍGITOS FINITOS

i.e., CADA NÚMERO REPRESENTÁVEL TEM APENAS  
UM NÚMERO FIXO E FINITO DE DÍGITOS

⇒ APENAS UM SUBCONJUNTO DOS NÚMEROS RACIONAIS  
PODE SER REPRESENTADO EXATAMENTE

⇒ O CÁLCULO COM OS DEMAIS NÚMEROS LEVA A ERROS  
DE ARREDONDAMENTO

### - REPRESENTAÇÃO DE NÚMEROS

⇒ DEPENDE DA BASE DISPONÍVEL NA MÁQUINA EM USO,  
E DO NÚMERO MÁXIMO DE DÍGITOS USADO NA REPRESENTAÇÃO  
COMPUTADORES OPERAM NORMALMENTE EM BASE BINÁRIA

⇒ DADOS DE ENTRADA SÃO ENVIADOS AO COMPUTADOR  
EM SISTEMA DECIMAL E CONVERTIDOS AO SISTEMA  
BINÁRIO.

AS OPERAÇÕES SÃO EFETUADAS NO SISTEMA BINÁRIO  
E OS RESULTADOS SÃO CONVERTIDOS AO SISTEMA  
DECIMAL PARA TRANSMISSÃO AO USUÁRIO



# - CONVERSÃO ENTRE OS SISTEMAS BINÁRIO E DECIMAL

110

## \* CONVERSÃO DE NÚMEROS INTEIROS

REPRESENTAÇÃO NUMA BASE  $\beta$ :

$$\text{NÚMERO } N = (a_j a_{j-1} \dots a_2 a_1 a_0)_\beta$$

$$= a_j \beta^j + a_{j-1} \beta^{j-1} + \dots + a_2 \beta^2 + a_1 \beta^1 + a_0 \beta^0$$

$$\text{ONDE } 0 \leq a_k \leq \beta - 1, \quad k = 0, \dots, j$$

$$\text{e.g., } (10111)_2 = (1 \times 2^4) + (0 \times 2^3) + (1 \times 2^2) + (1 \times 2^1) + (1 \times 2^0)$$

$$(347)_{10} = (3 \times 10^2) + (4 \times 10^1) + (7 \times 10^0)$$

### 1. CONVERSÃO DO SISTEMA BINÁRIO PARA O DECIMAL

$$\begin{aligned} \text{e.g., } (10111)_2 &= 1 \times 2^4 + 0 \times 2^3 + 1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 \\ &= (1 \times 2^4) + (1 \times 2^2) + (0 \times 2^3) + (1 \times 2^1) + 1 \end{aligned}$$

Colocando 2 em evidência:

$$\begin{aligned} (10111)_2 &= 2 \times (1 + 2 \times (1 + 2 \times (0 + 2 \times 1))) + 1 \\ &= 22 + 1 = 23 \end{aligned}$$

$$\therefore (10111)_2 = (23)_{10}$$

NO CASO GERAL: PARA OBTIVER A REPRESENTAÇÃO DECIMAL,  
 DEPARTANDO POR  $b_0$ , DO NÚMERO BINÁRIO  
 $(a_j a_{j-1} \dots a_2 a_1 a_0)_2$ , FAZEMOS

$$b_j = a_j$$

$$b_{j-1} = a_{j-1} + 2b_j$$

$$b_{j-2} = a_{j-2} + 2b_{j-1}$$

$$\vdots$$

$$b_1 = a_1 + 2b_2$$

$$b_0 = a_0 + 2b_1$$

NO EXEMPLO: PARA  $(10111)_2$

$$b_4 = a_4 = 1$$

$$b_3 = a_3 + 2b_4 = 0 + 2 \times 1 = 2$$

$$b_2 = a_2 + 2b_3 = 1 + 2 \times 2 = 5$$

$$b_1 = a_1 + 2b_2 = 1 + 2 \times 5 = 11$$

$$b_0 = a_0 + 2b_1 = 1 + 2 \times 11 = 23$$

## 2. CONVERSÃO DO SISTEMA DECIMAL PARA O BINÁRIO

e.g.,  $N_0 = 347$

Seja  $(a_j a_{j-1} \dots a_2 a_1 a_0)_2$  a representação binária procurada

$$N_0 = 347 = a_j 2^j + a_{j-1} 2^{j-1} + \dots + a_2 2^2 + a_1 2 + a_0$$

$$\therefore 2 \times 173 + 1 = 2 \times (a_j 2^{j-1} + a_{j-1} 2^{j-2} + \dots + a_2 2 + a_1) + a_0$$

$$\therefore a_0 = 1$$

Repetindo o processo para  $N_1 = 173$ :

$$N_1 = 173 = a_j 2^{j-1} + a_{j-1} 2^{j-2} + \dots + a_2 2 + a_1$$

$$\therefore 2 \times 86 + 1 = 2 \times (a_j 2^{j-2} + a_{j-1} 2^{j-3} + \dots + a_2) + a_1$$

$$\therefore a_1 = 1$$

Continuando, obtemos

$$N_2 = 86 = 2 \times 43 + 0 \Rightarrow a_2 = 0$$

$$N_3 = 43 = 2 \times 21 + 1 \Rightarrow a_3 = 1$$

$$N_4 = 21 = 2 \times 10 + 1 \Rightarrow a_4 = 1$$

$$N_5 = 10 = 2 \times 5 + 0 \Rightarrow a_5 = 0$$

$$N_6 = 5 = 2 \times 2 + 1 \Rightarrow a_6 = 1$$

$$N_7 = 2 = 2 \times 1 + 0 \Rightarrow a_7 = 0$$

$$N_8 = 1 = 2 \times 0 + 1 \Rightarrow a_8 = 1$$

$$\therefore 347 = (101011011)_2$$

NO CASO GERAL: PARA OBTER A REPRESENTAÇÃO 1.13  
BINÁRIA, DEVEMOS POR  $(a_j a_{j-1} \dots a_0)_2$ ,  
DO NÚMERO DECIMAL  $N$ , IMPLEMENTAMOS  
O SEGUINTE ALGORITMO, QUE OBTÉM, A  
CADA  $k$ , O DÍGITO BINÁRIO  $a_k$ .

PASSO 0:  $k=0$

$$N_k = N$$

PASSO 1: OBTENHA  $q_k \in \mathbb{Z}_k$  TAIS QUE

$$N_k = 2 \times q_k + r_k$$

$$\text{FAÇA } a_k = r_k$$

PASSO 2: SE  $q_k = 0$ , PARE

$$\text{SE NÃO, FAÇA } N_{k+1} = q_k$$

$$\text{FAÇA } k = k + 1$$

E VOLTE AO PASSO 1



## \* CONVERSÃO DE NÚMEROS FRACTIONÁRIOS

1.14

### 1. CONVERSÃO DO SISTEMA DECIMAL PARA O BINÁRIO

UM NÚMERO REAL ENTRE 0 E 1 PODE TER UMA REPRESENTAÇÃO FINITA NO SISTEMA DECIMAL E INFINITA NO BINÁRIO

=> USAREMOS A REPRESENTAÇÃO BINÁRIA

$$(0.d_1d_2\dots)_2 = 0 \times 2^0 + d_1 \times 2^{-1} + d_2 \times 2^{-2} + \dots$$

e.g.,  $N = (0.125)_{10}$

MULTIPLICAMOS OS DOIS LADOS POR 2:

$$2 \times 0.125 = 0.25 = 0 + 0.25 = d_1 + d_1 2^{-1} + \dots$$

=>  $d_1$  REPRESENTA A PARTE INTEIRA DE  $2 \times 0.125$ , QUE É ZERO.

AS DEMAIS PARCELAS REPRESENTAM A PARTE FRACTIONÁRIA, QUE É 0.25

REPETIMOS O PROCESSO PARA 0.25:

$$0.5 = d_2 + d_2 2^{-1} + d_2 2^{-2} + \dots \Rightarrow d_2 = 0$$

NA MESMA FORMA, PARA 0.5:

$$1.0 = d_3 + d_3 2^{-1} + d_3 2^{-2} + \dots \Rightarrow d_3 = 1$$

COMO A PARTE FRACTIONÁRIA É ZERO, O PROCESSO PARA

$$\text{É } (0.125)_{10} = (0.001)_2$$

No caso Geral: Para obter a representação 1.15  
Binária, denotada por  $(0.d_1d_2\dots d_j\dots)_2$ ,  
do número decimal  $R$  entre 0 e 1,  
implementamos o algoritmo seguinte  
que obtém, a cada  $k$ , o dígito  
binário  $d_k$ .

Passo 0:  $k=1$   
 $R_1 = R$

Passo 1: Calcule  $2R_k$   
Se  $2R_k \geq 1$ , Faça  $d_k = 1$   
Senão, Faça  $d_k = 0$

Passo 2: Faça  $R_{k+1} = 2R_k - d_k$   
Se  $R_{k+1} = 0$ , Pare.  
Senão, Faça  $k = k+1$   
e volte ao Passo 1

OBSERVAÇÕES:

- i) O algoritmo pode não parar após um número finito de passos, quando o número tem representação binária infinita.
- ii) O fato de alguns números terem representação binária infinita pode acarretar erros aparentemente inexplicáveis nos cálculos efetuados em computadores.

ex., Para  $\pi = (0.1)_{10}$

$$k=1 \quad \pi_1=0.1, 2\pi_1=0.2 < 1 \therefore d_1=0, \pi_2=0.2$$

$$k=2 \quad 2\pi_2=0.4 < 1 \therefore d_2=0, \pi_3=0.4$$

$$k=3 \quad 2\pi_3=0.8 < 1 \therefore d_3=0, \pi_4=0.8$$

$$k=4 \quad 2\pi_4=1.6 > 1 \therefore d_4=1, \pi_5=0.6$$

$$k=5 \quad 2\pi_5=1.2 > 1 \therefore d_5=1, \pi_6=0.2 = \underline{\pi_1}$$

$\Rightarrow$  OS RESULTADOS PARA  $k$  DE 2 A 5 SE REPETEM

$$\text{i.e., } \pi_2 = \pi_6 = \pi_{10} = \dots$$

$$\therefore (0.1)_{10} = (0.000110011\overline{0011}\dots)_2$$

ONDE O TRAÇO INDICA A PARTE QUE SE REPETE

$\Rightarrow$  O NÚMERO 0.1 TEM REPRESENTAÇÃO

INFINITA NA BASE 2.



## 2. CONVERSÃO DO SISTEMA BINÁRIO PARA O DECIMAL

⇒ SEGUIMOS UM ALGORITMO SEMELHANTE AO ANTERIOR

SEJA A REPRESENTAÇÃO DECIMAL

$$(0.b_1b_2\dots b_j)_{10} = b_110^{-1} + b_210^{-2} + \dots + b_j10^{-j}$$

Ex.:  $R = (0.000111)_2 = b_110^{-1} + b_210^{-2} + \dots + b_j10^{-j}$

MULTIPLICAMOS OS DOIS LADOS POR 10 [NA BASE BINÁRIA, I.E.,  $(1010)_2$ ]:

$$\begin{aligned} (0.000111)_2 \times (1010)_2 &= (1.11 \times 2^{-4})_2 \times (1.01 \times 2^3)_2 = \\ &= (1.11 \times 1.01)_2 \times 2^{-1} = (1.00011)_2 = b_1 + b_210^{-1} + \dots + b_j10^{-j+1} \end{aligned}$$

⇒  $b_1$  REPRESENTA A PARTE INTEIRA DO RESULTADO,

$$\underline{b_1 = 1}$$

REPETIMOS O PROCESSO PARA A PARTE FRACTIONÁRIA,

$$R_1 = (0.00011)_2 :$$

$$\begin{aligned} (0.00011)_2 \times (1010)_2 &= (1.1 \times 1.01) \times 2^{-1} = (0.1111)_2 = \\ &= b_2 + b_310^{-1} + \dots + b_j10^{-j+2} \end{aligned}$$

$$\Rightarrow \underline{b_2 = 0} \quad \& \quad R_2 = 0.1111$$

CONTINUAMOS O PROCESSO, CHEGAMOS A

$$(0.000111)_2 = (0.209375)_{10}$$



No caso geral: Para obter a representação 1.18  
Decimal, Denotada por  $(0.b_1b_2\dots b_k)_{10}$ ,  
Do número Binário  $R$  entre 0 e 1,  
Implementamos o Algoritmo seguinte,  
que obtém, a cada  $k$ , o dígito  
Decimal  $b_k$ .

Passo 0:  $k=1$   
 $R_1 = R$

Passo 1: Calcule  $W_k = (1010)_2 \times R_k$   
Seja  $z_k$  a parte inteira de  
 $W_k$   
 $b_k$  é a conversão de  $z_k$  para  
a base 10

Passo 2: Faça  $R_{k+1} = W_k - z_k$   
Se  $R_{k+1} = 0$ , Pare.  
Senão, Faça  $k = k+1$   
e volte ao passo 1.

# - DIFICULDADES DA COMPUTAÇÃO COM PRECISÃO FINITA

## \* REPRESENTAÇÃO EM PONTO FLUTUANTE:

$\Rightarrow$  UM NÚMERO REAL É REPRESENTADO COMO

$$\pm (0.d_1d_2\dots d_k) \times \beta^l$$

ONDE:  $\beta \equiv$  BASE EM QUE A MÁQUINA OPERA

$k \equiv$  NÚMERO DE DÍGITOS NA MANTISSA  
OU PARTE FRACTIONÁRIA

$l \equiv$  EXPONENTE INTEIRO

$0 \leq d_i \leq (\beta - 1)$ , PARA  $i = 1, \dots, k$   
COM  $d_1 \neq 0$

ex.: REPRESENTAÇÃO DECIMAL EM PONTO FLUTUANTE:

$$\pm (0.d_1d_2\dots d_n) \times 10^m \quad 0 \leq d_i \leq 9, d_1 \neq 0$$

$\Rightarrow$  DADO UM NÚMERO REAL POSITIVO NO FORMATO

$$y = 0.d_1d_2\dots d_nd_{n+1}\dots \times 10^m \quad (\text{DECIMAL NORMALIZADO})$$

A SUA REPRESENTAÇÃO EM PONTO FLUTUANTE É OBTIDA LIMITANDO-SE A MANTISSA DE  $y$  A  $k$  DÍGITOS DECIMAIS.

POSSIBILIDADES: i) TRUNCAMENTO

ii) ARREDONDAMENTO

## i) TRUNCAMENTO

1.20

$\Rightarrow$  Cortamos os dígitos  $d_{k+1}d_{k+2}\dots$ , Para obter

$$fl(y) = 0.d_1d_2\dots d_k \times 10^m$$

## ii) ARREDONDAMENTO

$\Rightarrow$  Adicionamos  $5 \times 10^{m-(k+1)}$  A  $y$ , e truncamos

o resultado, Para obter

$$fl(y) = 0.d_1d_2\dots d_k \times 10^m$$

$\Rightarrow$  Adicionamos 1 a  $d_k$ , se  $d_{k+1} \geq 5$  (Arredondamos Para cima)

Se não, cortamos todos os dígitos seguintes  
aos  $k$  primeiros (Arredondamos Para Baixo)

e.g.,  $\pi = 3.14159265\dots$  TEM INFINITOS DÍGITOS

Notação decimal normalizada:  $\pi = 0.314159265\dots \times 10^1$

Notação em ponto flutuante com 5 dígitos:

- TRUNCAMENTO:  $fl(\pi) = 0.31415 \times 10^1 = 3.1415$

- ARREDONDAMENTO:

$$\pi + 5 \times 10^{1-(5+1)} = \pi + 5 \times 10^{-5} =$$

$$= 3.14159265\dots + 0.00005 =$$

$$= 3.14164265$$

$$\therefore fl(\pi) = 0.31416 \times 10^1 = 3.1416$$



- Erros De Arredondamento e Truncamento

≡ Erros Resultantes Da Substituição Do Número Real Por Sua Notação Em Ponto Flutuante

\* Definição: Erro Absoluto e Erro Relativo

Se  $p^*$  é uma aproximação para  $p$ ,

O Erro Absoluto é  $|p - p^*|$

O Erro Relativo é  $\frac{|p - p^*|}{|p|}$ , com  $p \neq 0$

ou  $\frac{|p - p^*|}{|p^*|}$ , com  $p^* \neq 0$

Observação: Como uma medida de precisão, o Erro Relativo é mais significativo, pois leva em conta a magnitude dos valores

\* Definição: Algarismos Significativos

O número  $p^*$  se aproxima do valor  $p$  com  $t$  Algarismos Significativos, se  $t$  é o maior valor inteiro não negativo para o qual

$$\frac{|p - p^*|}{|p|} \leq 5 \times 10^{-t}$$

e.g.,  $p = 0.5$   $\Rightarrow \frac{|p - p^*|}{|p|} = 0.0004 \leq 5 \times 10^{-4}$   
 $p^* = 0.5002$

$\Rightarrow t = 4$  é o maior inteiro não-negativo para o qual  $|p - p^*|/|p| \leq 5 \times 10^{-t}$

$\Rightarrow p^*$  aproxima  $p$  com 4 Algarismos Significativos



## - LIMITES PARA OS ERROS DE ARREDONDAMENTO E TRUNCAMENTO 1.22

SEJA  $y$  UM NÚMERO REAL A SER EXPRESSO EM POTÊNCIA DE PONTO FLUTUANTE DECIMAL COM  $k$  DÍGITOS.

EXPRESSAMOS  $y$  COMO

$$y = f \times 10^m + g \times 10^{m-k}$$

ONDE  $0.1 \leq f < 1$

$$0 \leq g < 1$$

ex. g.,  $y = 234.57$  COM  $k = 4$  DÍGITOS:

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

$$\therefore f = 0.2345 \text{ e } g = 0.7$$

$\Rightarrow$  O FATOR  $g \times 10^{m-k}$  REPRESENTA A PARTE DE  $y$  QUE NÃO PODE SER INCORPORADA À MANTISSA COM  $k$  DÍGITOS

$\Rightarrow$  O ERRO CORTADO DEPENDE DA APROXIMAÇÃO UTILIZADA: POR TRUNCAMENTO OU ARREDONDAMENTO.

### i) TRUNCAMENTO

$\Rightarrow$  O FATOR  $g \times 10^{m-k}$  É DESPREZADO E  $y^* = f \times 10^m$

$$\text{ERRO ABSOLUTO: } E_A = |y - y^*| = |g| \times 10^{m-k} < 10^{m-k}$$

$$\text{POIS } |g| < 1$$

$$\text{ERRO RELATIVO: } E_R = \frac{|E_A|}{|y^*|} = \frac{|g| \times 10^{m-k}}{|f| \times 10^m} < \frac{10^{m-k}}{0.1 \times 10^m} = 10^{-k+1}$$

$$\text{POIS } f \geq 0.1$$

## ii) ARREDONDAMENTO

$\Rightarrow$  O Fator  $g \times 10^{m-k}$  é levado em conta, modificando  $f \times 10^m$

$$y^* = \begin{cases} f \times 10^m, & \text{se } |g| < \frac{1}{2} \\ f \times 10^m + 10^{m-k}, & \text{se } |g| \geq \frac{1}{2} \end{cases}$$

\* Se  $|g| < \frac{1}{2}$ :

ERRO ABSOLUTO:  $E_A = |y - y^*| = |g| \times 10^{m-k} < \frac{1}{2} \times 10^{m-k}$

ERRO RELATIVO:  $E_R = \frac{|E_A|}{|y^*|} = \frac{|g| \times 10^{m-k}}{|f| \times 10^m} < \frac{0.5 \times 10^{m-k}}{0.1 \times 10^m} = \frac{1}{2} \times 10^{-k+1}$

\* Se  $|g| \geq \frac{1}{2}$ :

ERRO ABSOLUTO:  $E_A = |y - y^*| = |(f \times 10^m + g \times 10^{m-k}) - (f \times 10^m + 10^{m-k})|$   
 $= |g \times 10^{m-k} - 10^{m-k}| = |(g-1)| \times 10^{m-k} \leq$   
 $\leq \frac{1}{2} \times 10^{m-k}$

ERRO RELATIVO:  $E_R = \frac{|E_A|}{|y^*|} \leq \frac{\frac{1}{2} \times 10^{m-k}}{|f \times 10^m + 10^{m-k}|} < \frac{\frac{1}{2} \times 10^{m-k}}{|f| \times 10^m} <$   
 $< \frac{\frac{1}{2} \times 10^{m-k}}{0.1 \times 10^m} = \frac{1}{2} \times 10^{-k+1}$

$\Rightarrow$  em ambos os casos,  $|E_A| \leq \frac{1}{2} \times 10^{m-k}$  e  $|E_R| < \frac{1}{2} \times 10^{-k+1}$

$\Rightarrow$  O ARREDONDAMENTO ACORTA ERROS NUMÉRICOS, AO CUSTO DE UM TEMPO MAIOR DE EXECUÇÃO.

OPERAÇÕES: Dados  $x$  e  $y$  reais

$$\text{Adição: } x \oplus y = fl(fl(x) + fl(y))$$

$$\text{Subtração: } x \ominus y = fl(fl(x) - fl(y))$$

$$\text{Multiplicação: } x \otimes y = fl(fl(x) \times fl(y))$$

$$\text{Divisão: } x \oslash y = fl(fl(x) \div fl(y))$$

$\Rightarrow$  OPERAÇÕES ARITMÉTICAS EXATAS com os valores de  $x$  e  $y$  em ponto flutuante e CONVERSÃO DO RESULTADO EXATO PARA Ponto Flutuante

OBSERVAÇÕES: i) numa sequência de operações, o erro em uma delas se propaga ao longo das operações subsequentes.

ii) o erro total em uma operação é composto pelo erro das parcelas ou fatores e pelo erro no resultado da operação.

- ERROS ABSOLUTO E RELATIVO NAS OPERAÇÕES ARITMÉTICAS:

Supomos Parcelas ou Fatores com erro

$$\text{e.g., } x = x^* + E_A(x)$$

$$y = y^* + E_A(y)$$

Supomos que o erro final é arredondando ou truncando, mas NÃO CONSIDERAMOS ESSE PASSO nas fórmulas a seguir.

OBS.: ERROS com SINAL!



i) Adição:

$$\begin{aligned} x + y &= (x^* + E_A(x)) + (y^* + E_A(y)) = \\ &= (x^* + y^*) + (E_A(x) + E_A(y)) \end{aligned}$$

$$\Rightarrow \text{ERRO ABSOLUTO: } E_A(x+y) = E_A(x) + E_A(y)$$

$$\text{ERRO RELATIVO: } E_R(x+y) = \frac{E_A(x+y)}{x^* + y^*} = \frac{E_A(x) + E_A(y)}{x^* + y^*}$$

$$= \frac{E_A(x)}{x^*} \left( \frac{x^*}{x^* + y^*} \right) + \frac{E_A(y)}{y^*} \left( \frac{y^*}{x^* + y^*} \right) =$$

$$= E_R(x) \left( \frac{x^*}{x^* + y^*} \right) + E_R(y) \left( \frac{y^*}{x^* + y^*} \right)$$

ii) Subtração:

ANALOGAMENTE, OBTENEMOS

$$\text{ERRO ABSOLUTO: } E_A(x-y) = E_A(x) - E_A(y)$$

ERRO RELATIVO:

$$E_R(x-y) = E_R(x) \left( \frac{x^*}{x^* - y^*} \right) - E_R(y) \left( \frac{y^*}{x^* - y^*} \right)$$

### iii) Multiplicação:

$$xy = (x^* + E_A(x))(y^* + E_A(y)) =$$

$$= x^*y^* + x^*E_A(y) + y^*E_A(x) + E_A(x)E_A(y)$$

Considerando que  $E_A(x)E_A(y)$  se é muito pequeno,  
Podemos desprezê-lo na expressão acima

Assim,

$$\text{ERRO ABSOLUTO: } E_A(xy) \approx x^*E_A(y) + y^*E_A(x)$$

$$\text{ERRO RELATIVO: } E_R(xy) \approx \frac{x^*E_A(y) + y^*E_A(x)}{x^*y^*} =$$

$$= \frac{E_A(x)}{x^*} + \frac{E_A(y)}{y^*} = E_R(x) + E_R(y)$$

### iv) Divisão:

$$\frac{x}{y} = \frac{x^* + E_A(x)}{y^* + E_A(y)} = \frac{x^* + E_A(x)}{y^*} \left( \frac{1}{1 + \frac{E_A(y)}{y^*}} \right) \approx$$

$$\approx \left( \frac{x^* + E_A(x)}{y^*} \right) \left( 1 - \frac{E_A(y)}{y^*} \right) \approx \frac{x^*}{y^*} + \frac{E_A(x)}{y^*} - \frac{x^*E_A(y)}{y^{*2}}$$

Assim,

$$\text{ERRO ABSOLUTO: } E_A(x/y) \approx \frac{E_A(x)}{y^*} - \frac{x^*E_A(y)}{y^{*2}} = \frac{y^*E_A(x) - x^*E_A(y)}{y^{*2}}$$

$$\text{ERRO RELATIVO: } E_R(x/y) \approx \left( \frac{y^*E_A(x) - x^*E_A(y)}{y^{*2}} \right) \frac{y^*}{x^*} =$$

$$= \frac{E_A(x)}{x^*} - \frac{E_A(y)}{y^*} = E_R(x) - E_R(y)$$

OBSERVAÇÃO: Dos resultados acima, podemos concluir 127

QUE: i) A SUBTRAÇÃO DE DOIS NÚMEROS APROXIMADAMENTE IGUAIS PODE LEVAR A UM GRANDE AUMENTO NO ERRO RELATIVO.

ii) MULTIPLICAÇÃO POR UM NÚMERO DE GRANDE MAGNITUDE, OU DIVISÃO POR UM NÚMERO DE PEQUENA MAGNITUDE, PODE ALCANÇAR UM GRANDE AUMENTO NO ERRO ABSOLUTO.

- ARTIFÍCIOS PARA MINIMIZAR ERROS NA ARITMÉTICA COMPUTACIONAL

i) EVITAR A SUBTRAÇÃO DE VALORES PRÓXIMOS

e.g., Equação do 2º Grau:  $ax^2 + bx + c = 0$

$$\text{Raízes: } x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}, \quad x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}$$

$\Rightarrow$  Quando  $b^2 \gg 4ac$ ,  $x_2$  envolve subtração de valores muito próximos

$\Rightarrow$  Fórmula alternativa: Racionalizando o numerador

$$\begin{aligned} x_2 &= \frac{-b + \sqrt{b^2 - 4ac}}{2a} \times \left( \frac{-b - \sqrt{b^2 - 4ac}}{-b - \sqrt{b^2 - 4ac}} \right) = \frac{b^2 - (b^2 - 4ac)}{2a(-b - \sqrt{b^2 - 4ac})} \\ &= -2c / (b + \sqrt{b^2 - 4ac}) \end{aligned}$$

ii) Reduzir o número de termos, no cálculo de polinômios

$$\begin{aligned} \text{e.g., } f(x) &= x^3 - 6.1x^2 + 3.2x + 1.5 \quad (4 \text{ produtos, 3 adições}) \\ &= [(x - 6.1)x + 3.2]x + 1.5 \quad (2 \text{ produtos, 3 adições}) \end{aligned}$$

Menos operações  $\Rightarrow$  menor probabilidade de erros



1 BIT INDICADOR DE SINAL (+ ou -): s

Parte Fracionária de 52 Bits (Mantissa): f

$$(-1)^s 2^{c-1023} (1+f)$$

ex., número de página: 01000000001110111001000000

$S = 0 \Rightarrow$  número positivo

$$C = 10000000011 = 1 \cdot 2^{10} + 0 \cdot 2^9 + \dots + 0 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0$$

$$= 1027$$

$\Rightarrow$  Parte exponencial:  $2^{1027-1023} = 2^4$

$$f = 1 \cdot \left(\frac{1}{2}\right)^1 + 1 \cdot \left(\frac{1}{2}\right)^3 + 1 \cdot \left(\frac{1}{2}\right)^4 + 1 \cdot \left(\frac{1}{2}\right)^5 + 1 \cdot \left(\frac{1}{2}\right)^8 + 1 \cdot \left(\frac{1}{2}\right)^{12}$$

$\Rightarrow$  MARTISSA: 0.72290039

$\Rightarrow$  O número de máquina representa o decimal

$$(-1)^5 2^{6-1023} (1+f) = (-1)^0 \cdot 2^4 \cdot (1.72250038)$$

$$= 27.56640625$$

NÚMERO DADO: 0 100...011 10111001000100...0

$\underbrace{\hspace{1cm}}_s \quad \underbrace{\hspace{2cm}}_c \quad \underbrace{\hspace{3cm}}_f$

Nº IMEDIATAMENTE  
INFERIOR: 0 100...011 10111001000011...1

$\underbrace{\hspace{1cm}}_s \quad \underbrace{\hspace{2cm}}_c \quad \underbrace{\hspace{3cm}}_f$

Nº IMEDIATAMENTE  
SUPERIOR: 0 100...011 1011100100010...1

$\underbrace{\hspace{1cm}}_s \quad \underbrace{\hspace{2cm}}_c \quad \underbrace{\hspace{3cm}}_f$

⇒ O NÚMERO DE MÁQUINA ORIGINAL REPRESENTA MEIO  
DOS NÚMEROS REAIS QUE ESTÃO ENTRE ELE E OS SEUS  
DOIS VIZINHOS MAIS PRÓXIMOS EM LINGUAGEM DE MÁQUINA

⇒ ELE REPRESENTA UM INTERVALO

- UNDERFLOW E OVERFLOW

\* O MEHOR NÚMERO POSITIVO QUE SE PODE REPRESENTAR  
NO PADRÃO IEEE TEM TODOS OS DÍGITOS ZERO, EXCETO  
O ÚLTIMO À DIREITA, QUE É 1

$$\Rightarrow 2^{-1023} \cdot (1 + 2^{-52}) \simeq 10^{-308}$$

\* O MAIOR NÚMERO POSITIVO TEM UM ZERO INICIAL  
SEGUINDO POR 63 DÍGITOS 1

$$\Rightarrow 2^{1024} \cdot (2 - 2^{-52}) \simeq 10^{308}$$

Se nos cálculos computacionais aparecem números menores do que  $\approx 10^{-308}$ , estes são ajustados para zero (UNDERFLOW).

Se aparecem números maiores do que  $\approx 10^{308}$ , se caracteriza OVERFLOW, o que geralmente leva a falhas no processamento.



### 3. ALGORITMOS E CONVERGÊNCIA

1.31

\* ALGORITMO  $\equiv$  PROCEDIMENTO QUE DESCREVE, SEM AMBIGÜIDADES, UMA SEQUÊNCIA FINITA DE PASSOS A SEREM EXECUTADOS, NUMA ORDEM ESPECIFICADA, COM O OBJETIVO DE RESOLVER UM PROBLEMA OU OBTER UMA SOLUÇÃO APROXIMADA PARA ELE.

#### DESCRIÇÃO DOS ALGORITMOS: PSEUDOCÓDIGO

\* PSEUDOCÓDIGO  $\Rightarrow$  ESPECIFICA A FORMA DOS DADOS DE ENTRADA E DOS RESULTADOS A SEREM PRODUZIDOS INCORPORA UMA PRESCRIÇÃO INDEPENDENTE PARA A INTERRUPÇÃO DOS PROCEDIMENTOS

CONVENÇÕES: i) PONTO (.): INDICA O FIM DE UM PASSO

ii) PONTO-E-VÍRGULA (;): SEPARA TAREFAS EM UM MESMO PASSO

iii) PARÁGRAFO: INDICA GRUPO DE INSTRUÇÕES A SEREM TRATADAS COMO UMA ÚNICA ENTIDADE

iv) CÁLCULOS EM LOOP:

SÃO CONTROLADOS POR CONTAGEM:

e.g., PARA  $i = 1, 2, \dots, n$

FAÇA ----

OU POR CONDIÇÕES: e.g., ENQUANTO  $i \leq n$

v) EXECUÇÕES CONDIÇIONAIS: SE... ENTÃO....  
SE NÃO ---

vi) COMENTÁRIOS: SÃO COLOCADOS ENTRE PARÊNTESES

EXEMPLO: CALCULAR  $\sum_{i=1}^N x_i = x_1 + x_2 + \dots + x_N$

ONDE  $x_1, x_2, \dots, x_N$  SÃO DADOS

ALGORITMO:

ENTRADA  $N, x_1, x_2, \dots, x_N$

SAÍDA  $SUM = \sum_{i=1}^N x_i$

Passo 1 Faça  $SUM = 0$ . (Inicializa o acumulador)

Passo 2 PARA  $i = 1, 2, \dots, N$

Faça  $SUM = SUM + x_i$  (adiciona o termo seguinte)

Passo 3 Saída (SUM);

PARA.

\* ALGORITMO ESTÁVEL  $\Rightarrow$  PEQUENAS ALTERAÇÕES NOS DADOS INICIAIS PRODUZEM PEQUENAS ALTERAÇÕES NO RESULTADO FINAL

\* ALGORITMO CONDICIONALMENTE ESTÁVEL

$\Rightarrow$  ESTÁVEL APENAS PARA CERTOS VALORES DOS DADOS INICIAIS

A ESTABILIDADE DE UM ALGORITMO ESTÁ RELACIONADA AO CRESCIMENTO DO ERRO DE ARREDONDAMENTO

\* DEFINIÇÃO: SEJA  $E_0 > 0$  UM ERRO DE ARREDONDAMENTO INICIAL.

SEJA  $E_n$  A MAGNITUDE DO ERRO APÓS  $n$  OPERAÇÕES SUCESSIVAS.

SE  $E_n \approx C \cdot n \cdot E_0$ , ONDE  $C$  É UMA CONSTANTE INDEPENDENTE DE  $n$ , O CRESCIMENTO DO ERRO É DITO LINEAR.

SE  $E_n \approx C^n E_0$  PARA ALGUM  $C > 1$ , O CRESCIMENTO DO ERRO É DITO EXPONENCIAL.

ALGORITMOS COM CRESCIMENTO LINEAR DO ERRO SÃO ESTÁVEIS.

ALGORITMOS COM CRESCIMENTO EXPONENCIAL DO ERRO SÃO INSTÁVEIS.



## - CONVERGÊNCIA DOS ALGORITMOS

1.34

\* DEFINIÇÃO: Seja  $\{\beta_m\}_{m=1}^{\infty}$  uma sequência convergente para 0.

Seja  $\{\alpha_m\}_{m=1}^{\infty}$  uma sequência convergente para o número  $\alpha$ .

Se existe uma constante positiva  $K$  tal que

$$|\alpha_m - \alpha| \leq K|\beta_m|$$

para um  $n$  grande, então dizemos que

$\{\alpha_m\}_{m=1}^{\infty}$  converge para  $\alpha$  com taxa de

convergência  $O(\beta_m)$ , i.e.,  $\alpha_m = \alpha + O(\beta_m)$

OBSERVAÇÃO: Em geral, utilizamos a sequência  $\beta_m = \frac{1}{n^p}$ ,

para um  $p > 0$  qualquer, e estamos

interessados no maior valor de  $p$  para o qual

$$\alpha_m = \alpha + O\left(\frac{1}{n^p}\right)$$

e.g., se  $\alpha_m = \frac{n+1}{n^2}$  e  $\hat{\alpha}_m = \frac{n+3}{n^3}$

temos

$$\alpha_m = 0 + O\left(\frac{1}{n}\right)$$

$$\hat{\alpha}_m = 0 + O\left(\frac{1}{n^2}\right)$$

De maneira semelhante, podemos definir a Taxa

de convergência de uma função:

\* Definição: Sejam  $\lim_{h \rightarrow 0} G(h) = 0$  e  $\lim_{h \rightarrow 0} F(h) = L$

Se existe uma constante  $K > 0$  para a qual

$$|F(h) - L| \leq K |G(h)|$$

Para  $h$  suficientemente pequeno,

então escrevemos

$$F(h) = L + O(G(h))$$

Observação: As funções utilizadas para comparação

geralmente têm a forma  $G(h) = h^p$ , com  $p > 0$ ,

e estamos interessados no maior valor  $p$

para o qual  $F(h) = L + O(h^p)$