



UNIVERSIDADE FEDERAL DE SANTA CATARINA - UFSC
CENTRO TECNOLÓGICO - CTC
DEPARTAMENTO DE INFORMÁTICA E ESTATÍSTICA - INE

Calculo Numérico em Computadores:

Capítulo 1

Sistemas de Numeração e Erros numéricos

Autores: Prof. Sérgio Peters
Acad. Andréa Vergara da Silva

e-mail: sergio.peters@ufsc.br

Florianópolis, 2013.

1. SISTEMAS DE NUMERAÇÃO

1.1 - INTRODUÇÃO

Atualmente o sistema padronizado de representação de quantidades para o uso e a comunicação entre as pessoas é o sistema decimal. Entretanto, para facilitar a representação física, a definição das operações aritméticas e a comunicação entre as máquinas digitais, é necessário fazer uso de outros sistemas de representação.

Como premissa básica, conceitua-se **número** como a representação simbólica de determinada quantidade matemática e **base** de um sistema de numeração a quantidade de símbolos distintos utilizados nesta representação. Desta forma, um número real qualquer X na base β pode ser algebricamente representado através de:

$$X_{\beta} = (a_1 a_2 \dots a_k , a_{k+1} a_{k+2} \dots a_{k+n})_{\beta} \quad (1)$$

onde β é a base, $a_i \in \{0,1,2,\dots,\beta-1\}$, $i = 1,2,\dots,k+n$, k é o comprimento da parte inteira e n da parte fracionária do número, com $k,n \in \mathbb{N}$.

Ex. 1: $X=(309,57)_{10}$

Para fins de uso algébrico X pode também ser representado na **forma fatorada** equivalente:

$$X_{\beta} = \sum_{i=1}^k a_i \beta^{k-i} + \sum_{j=1}^n a_{k+j} \beta^{-j} \quad (2)$$

Ex. 2: $X=(309,57)_{10} = 3 \cdot 10^2 + 0 \cdot 10^1 + 9 \cdot 10^0 + 5 \cdot 10^{-1} + 7 \cdot 10^{-2}$

A seguir, serão abordados alguns sistemas de numeração e as formas de representação de números de amplo uso nas máquinas digitais.

1.2 - SISTEMA DECIMAL ($\beta = 10$)

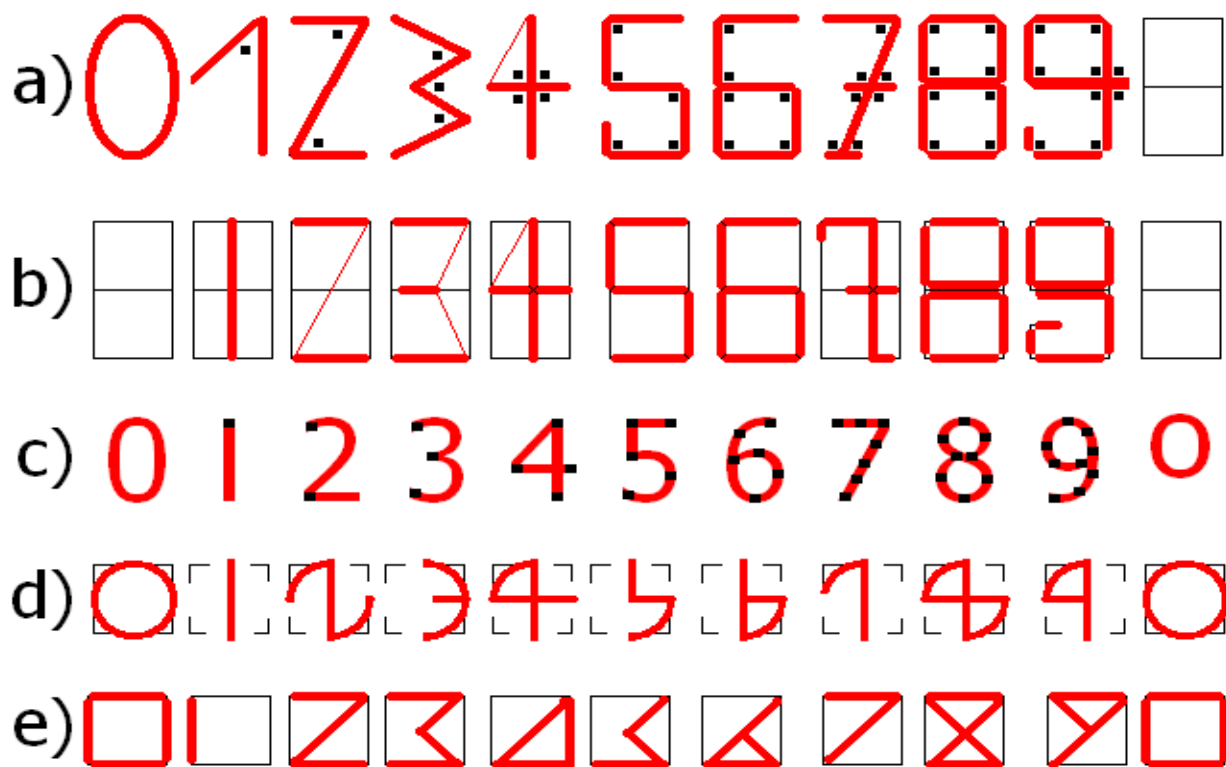
O sistema decimal de numeração, adotado pela maioria dos países, foi desenvolvido pelos astrônomos hindus por volta do século V e divulgado ao mundo islâmico em 825 no livro do matemático Alkharizmi e definitivamente adotado no ocidente no século XVI. Sua aceitação como padrão deve-se a algumas de suas características tais como:

i). Utiliza dez símbolos, dígitos (digitus = dedo em latim) ou algarismos (currupela lógica de Alkharizmi). Tais símbolos atualmente são representados por: 0, 1, 2, 3, 4, 5, 6, 7, 8 e 9, conhecidos por algarismos arábicos, que são derivados da versão ainda hoje usada no mundo muçulmano:

٩, ٨, ٧, ٦, ٥, ٤, ٣, ٢, ١, ٠

Historicamente estes números tem diversas hipóteses sobre sua origem::

- a) Número de ângulos existentes no desenho de cada algarismo;
- b) Número de traços contidos no desenho de cada algarismo;
- c) Número de pontos de cada algarismo;
- d) Número de diâmetros e arcos de circunferência contidos no desenho de cada algarismo;
- e) Figuras desenhadas a partir dos traços de um quadrado e suas diagonais.



ii). Faz uso do zero. O zero, aceito com muita relutância, é o indicador da ausência de certas potências da base na representação de um número na forma fatorada;

iii). Adota o princípio da posicionalidade. No sistema posicional o valor de cada símbolo é relativo, isto é, depende da sua posição no número.

Ex. 3: Nos números

a) $(574)_{10} = 5 \times 10^2 + 7 \times 10^1 + 4 \times 10^0$

b) $(348)_{10} = 3 \times 10^2 + 4 \times 10^1 + 8 \times 10^0$

c) $(432,5)_{10} = 4 \times 10^2 + 3 \times 10^1 + 2 \times 10^0 + 5 \times 10^{-1}$

O símbolo 4 representa, respectivamente, quatro unidades, quatro dezenas e quatro centenas.

Note que nenhum dígito interfere na posição do outro, eles são inteiramente independentes entre si.

Utilizando-se das duas últimas características do sistema decimal, a seguir serão estabelecidos outros sistemas de numeração, para facilitar a comunicação homem-máquina.

1.3 - SISTEMA BINÁRIO ($\beta = 2$)

Fazendo uso apenas dos símbolos 0 e 1, também chamados de *bits* (abreviatura de *binary digits*); do zero e da posicionalidade, gera-se um novo sistema de numeração cuja correspondência com o decimal será:

Decimal	0	1	2	3	4	5	6	7	8	9	10	.	19	.
Binário	0	1	10	11	100	101	110	111	1000	1001	1010	.	10011	.

Utilizando-se da notação fatorada, tem-se por exemplo,

Ex. 4:

$$(10011)_2 = (1.2^4 + 0.2^3 + 0.2^2 + 1.2^1 + 1.2^0)_{10} = (19)_{10}$$

Obs.: A forma fatorada do número binário (base 2) está representada na base 10.

- **Vantagens do Sistema Binário em Relação ao Sistema Decimal**

(i). Simplicidade de representação física, bastam 2 estados distintos de uma máquina digital para representar os dígitos da base: 0 = -, off
1 = +, on

Obs.: No futuro poderá se chegar a distinção de bits através dos átomos, que é a concepção do computador quântico (*quantum-bits*, ou *qubits*).

(ii). Simplicidade na definição de operações aritméticas fundamentais:

Ex. 5: Adição:

$$+ : \mathfrak{X} \times \mathfrak{X} \rightarrow \mathfrak{X}$$

$$(x, y) \rightarrow x + y \text{ definida por:}$$

Em $\beta = 10$ necessita-se de 100 combinações dos possíveis valores de x e y para se definir a função adição.

Em $\beta = 2$ tem-se apenas 4 combinações:

$$0 + 0 = 0$$

$$0 + 1 = 1$$

$$1 + 0 = 1$$

$$1 + 1 = 10$$

- **Desvantagens do Sistema Binário**

(i). Necessidade de registros longos para armazenamento de números.

Ex. 6: $(597)_{10} = (1001010101)_2$

Observa-se que foi necessário um registro com capacidade de armazenamento de dez símbolos binários para representar a grandeza decimal $(597)_{10}$ de apenas três dígitos decimais.

1.4 - SISTEMA HEXADECIMAL ($\beta = 16$)

Símbolos representativos: 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E e F. Onde A, B, C, D, E e F representam as quantidades decimais 10, 11, 12, 13, 14 e 15, respectivamente.

Este também é um sistema posicional.

$$\begin{aligned} \text{Ex. 7: } (1A0, C)_{16} &= 1 \times 16^2 + A \times 16^1 + 0 \times 16^0 + C \times 16^{-1} \\ &= (1 \times 16^2 + 10 \times 16^1 + 0 \times 16^0 + 12 \times 16^{-1})_{10} \\ &= (256 + 160 + 0 + 12/16 = 416,75)_{10} \end{aligned}$$

- **Vantagem do Sistema Hexadecimal**

(i). Número reduzido de símbolos para representar grandes quantidades, por isso é um sistema de numeração interessante para visualização e armazenamento de dados. Os registros binários internos de uma máquina digital são convertidos de forma direta para Hexadecimal quando são necessárias visualizações externas, requisitadas pelo usuário.

Ex. 8: $(1101\ 0110)_2 = (D6)_{16} = (214)_{10}$

oito bits

2 símbolos hexadecimais

Obs.: Até a década de 70 as máquinas digitais se utilizavam do sistema de numeração octal, de base

8, para visualizar os registros binários internos.

1.5 - CONVERSÕES ENTRE SISTEMAS DE NUMERAÇÃO

As conversões entre bases são necessárias para que se possa melhor entender algumas das causas dos erros existentes nas representações digitais de quantidades, pois o homem utiliza o sistema decimal e os computadores as convertem para bases binária, hexadecimal, ou outra.

1.5.1 - CONVERSÃO DE BASE β PARA BASE 10

Nestes casos as conversões são obtidas escrevendo o respectivo número na sua forma fatorada, representada na base decimal.

$$(a_1a_2a_3a_4a_5)_\beta = (a_1 \cdot \beta^2 + a_2 \cdot \beta^1 + a_3 \cdot \beta^0 + a_4 \cdot \beta^{-1} + a_5 \cdot \beta^{-2})_{10}$$

$$\text{Ex. 9: } (101,1)_2 = 1 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 + 1 \cdot 2^{-1} = (5,5)_{10}$$

$$(1A,B)_{16} = 1 \cdot 16^1 + A \cdot 16^0 + B \cdot 16^{-1} = (16,6875)_{10}$$

1.5.2 - CONVERSÃO DE BASE 10 PARA BASE β

$$(\overline{17},5)_{10} = ((i),(ii))_\beta$$

Procedimento de conversão:

i) na parte inteira do número: divide-se sucessivamente a parte inteira do número decimal pela base β , e constrói-se o novo número escrevendo o último quociente e os restos obtidos nas divisões, para separar as potências sucessivas de β componentes da parte inteira. Desta forma agrupam-se as diferentes potências da base β .

$$\text{Ex. 10: } (19)_{10} = ()_2 \Rightarrow 19 \overline{)2}$$

$$1 \ 9 \overline{)2}$$

$$1 \ 4 \overline{)2}$$

$$0 \ 2 \overline{)2}$$

$$0 \ 1 \rightarrow (10011)_2 = 1 \cdot 2^4 + 0 \cdot 2^3 + 0 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0$$

$$(527)_{10} = ()_{16} \Rightarrow 527 \overline{)16}$$

$$15 \ 32 \overline{)16}$$

$$0 \ 2 \rightarrow (20F)_{16} = 2 \cdot 16^2 + 0 \cdot 16^1 + F \cdot 16^0$$

ii) na parte fracionária do número: multiplica-se sucessivamente a parte fracionária do número decimal pela base β , e constrói-se o número escrevendo os inteiros resultantes de cada produto. A parte fracionária restante é novamente multiplicada por β até que o produto final seja um inteiro ou, a quantidade limite de dígitos na representação seja atingida. Desta forma constróem-se as frações sucessivas de β .

$$\text{Ex. 11: } (0,03125)_{10} = ()_2$$

$$\begin{array}{r} 0,03125 \\ \times 2 \\ \hline 0,06250 \end{array}$$

$$\begin{array}{r} 0,06250 \\ \times 2 \\ \hline 0,12500 \end{array}$$

$$\begin{array}{r} 0,125 \\ \times 2 \\ \hline 0,250 \end{array}$$

$$\begin{array}{r} 0,25 \\ \times 2 \\ \hline 0,50 \end{array}$$

$$\begin{array}{r} 0,5 \\ \times 2 \\ \hline 1,0 \end{array}$$

$$= (0,00001)_2$$

$$\text{Ex. 12: } (0,1)_{10} = ()_2$$

$$(119)_{10}$$

$$(120)_{10}$$

Represente os números na forma fatorada e converta para a base decimal:

$$\text{i) } (3021)_{F!} = 3 \times 4! + 0 \times 3! + 2 \times 2! + 1 \times 1! = (77)_{10}$$

$$\text{ii) } (321,123)_{F!} =$$

$$\text{iii) } (0,02)_{F!} =$$

$$\text{iv) } (0,113)_{F!} =$$

Note que nos exercícios (ii), (iii) e (iv) tem-se representações exatas de números racionais, que na base decimal são dízimas periódicas.

1.5.3 - CONVERSÕES DIRETAS ENTRE BINÁRIO E HEXADECIMAL:

Estas conversões são importantes para se entender os mecanismos de operacionalização de máquinas digitais que implementam as operações aritméticas em base binária e visualizar as representações em base hexadecimal.

Sabemos que um dígito hexadecimal corresponde a quatro dígitos binários, pois

$$16^1 = 2^4 \text{ e note que } (15)_{10} = (F)_{16} = (1111)_2$$

Então, fazemos a conversão direta associando a cada um dígito hexadecimal quatro dígitos binários. Para tal, agrupamos os dígitos binários em grupos de quatro a partir da posição da vírgula, para a direita e para a esquerda. Caso seja necessário, completa-se o grupo de quatro *bits* com zeros não significativos.

$$\text{Ex. 13: } (A1,B)_{16} = (\quad)_2$$

Como:

$$(A)_{16} = (1010)_2$$

$$(1)_{16} = (0001)_2$$

$$(B)_{16} = (1011)_2$$

temos:

$$(A1,B)_{16} = (1010 \ 0001, 1011)_2$$

$$\text{Ex. 14: } (10001, 01001 \ 1001 \ 1001 \dots)_2 = (\quad)_{16}$$

Agrupando-se os dígitos em grupos de quatro e completando com zeros:

$$\Rightarrow (\underline{0001} \ \underline{0001}, \ \underline{0100} \ \underline{1100} \ \underline{1100} \ \underline{1100} \dots)_2$$

$$\begin{matrix} 1 & 1 & 4 & C & C & C \end{matrix} \Rightarrow (11,4CCC\dots)_{16}$$

$$\text{Ex. 15: } (110001,01) = (\quad)_{16}$$

$$(\underline{0011} \ \underline{0001}, \ \underline{0100})_2 = (31,4)_{16}$$

$$\begin{matrix} 3 & 1 & 4 \end{matrix}$$

$$\text{Pois, } (00110001,0100)_2 =$$

$$\begin{array}{c} \underline{0 \times 2^7 + 0 \times 2^6 + 1 \times 2^5 + 1 \times 2^4} + \underline{0 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0} + \underline{0 \times 2^{-1} + 1 \times 2^{-2} + 0 \times 2^{-3} + 0 \times 2^{-4}} \\ \swarrow \quad \searrow \quad \quad \quad \swarrow \quad \searrow \quad \quad \quad \swarrow \quad \searrow \\ (0 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0) (2^4)^1 + (0 \times 2^3 + 0 \times 2^2 + 0 \times 2^1 + 1 \times 2^0) (2^4)^0 + (0 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 0 \times 2^0) (2^4)^{-1} \\ \swarrow \quad \searrow \quad \quad \quad \swarrow \quad \searrow \quad \quad \quad \swarrow \quad \searrow \\ \underline{3} \quad \times \quad \underline{16^1} \quad \quad \quad \underline{1} \quad \times \quad \underline{16^0} \quad \quad \quad \underline{4} \quad \times \quad \underline{16^{-1}} \end{array}$$

logo

$$3 \times 16^1 + 1 \times 16 + 4 \times 16 = (31,4)_{16}$$

Obs.: Nas conversões diretas entre as bases binária e hexadecimal (ou entre as bases binária e octal), não há perda de dígitos (arredondamento). Mas nas conversões de base decimal para base

binária, ou para base hexadecimal, ou para base octal, podemos perder dígitos significativos. Por exemplo:

$$(17,3)_{10} = (11,4CCC \dots)_{16} = (\quad)_{10} \\ = 1 \cdot 16^1 + 1 \cdot 16^0 + 4 \cdot 16^{-1} + 12 \cdot 16^{-2} + 12 \cdot 16^{-3} + \dots$$

Assim,

$$(17,3)_{10} = (11,4CCCC\dots)_{16} \cong (11,4CC)_{16} \quad (5 \text{ significativos}) \\ \cong (17,296875)_{10}$$

Exercícios:

1.4 - Converter os números para as bases na ordem indicada:

- a) $(10111,1101)_2 = (\quad)_{16} = (\quad)_{10}$
- b) $(BD,0E)_{16} = (\quad)_{10} = (\quad)_2$
- c) $(41,1)_{10} = (\quad)_2 = (\quad)_{16}$

Obs.: Verifique se houve perda de dígitos significativos em alguma das conversões, considerando um número limitado de dígitos representáveis.

2. REPRESENTAÇÃO DIGITAL DE NÚMEROS:

2.1. - INTRODUÇÃO

De uma maneira geral, nos sistemas computacionais, um número $X \in \Re$ é representado na forma de notação em ponto flutuante, de maneira a racionalizar o armazenamento digital.

Se utilizássemos um armazenamento em **ponto fixo** (vírgula fixa) seria necessário um número de posições (dígitos) no mínimo igual a variação dos limites dos expoentes. Por exemplo, para se obter a representação de uma calculadora científica comum com limites positivos entre $1,0 \cdot 10^{-99}$ e $9,999999999 \cdot 10^{+99}$ seria necessário:

(i). Entre $1,0 \cdot 10^{-99}$ e 1 seriam necessárias 99 posições:

$$1,0 \cdot 10^{-99} = 0,000 \dots 00001 \rightarrow 99 \text{ dígitos após a vírgula}$$

(ii). Entre 1 e $9,999999999 \cdot 10^{+99}$ seriam necessárias 100 posições:

$$9,999999999 \cdot 10^{+99} = 9999999999000 \dots 0000, \rightarrow 100 \text{ dígitos inteiros}$$

(iii). Seria necessário mais uma posição para o sinal (s), para as representações de negativos, totalizando 200 posições em cada registro:

s		
	100 posições para a parte inteira	99 posições para a parte fracionária

Por outro lado, em uma representação em **Ponto Flutuante**, esta calculadora científica funciona com pouco mais de dez dígitos, incluindo posições reservadas ao expoente.

Então, em uma representação em Ponto Flutuante, onde a vírgula flutua segundo um certo padrão, temos a seguinte representação genérica na base β :

$$X = \pm [d_1/\beta + d_2/\beta^2 + d_3/\beta^3 + \dots + d_t/\beta^t] \cdot \beta^{\text{exp}}$$

ou

$$X = \pm (0, d_1 d_2 d_3 \dots d_t)_{\beta} \cdot \beta^{\text{exp}}$$

onde

d_i = números inteiros contidos em $0 \leq d_i \leq (\beta - 1)$ ($i = 1, 2, \dots, t$) que constituem a mantissa.

Obs.: É necessário algum tipo de normalização para padronização da mantissa, no caso adota-se $d_1 \neq 0$.

exp = expoente de β , assume valores limites I (Inferior) e S (Superior) onde $I \leq \text{exp} \leq S$.

t = número de dígitos significativos do sistema de representação, é chamado de precisão da máquina.

Ex. 16: Representar em ponto flutuante:

a) $(3,501)_{10} = (3/10 + 5/10^2 + 0/10^3 + 1/10^4) \cdot 10^1 = 0,3501 \times 10^1$

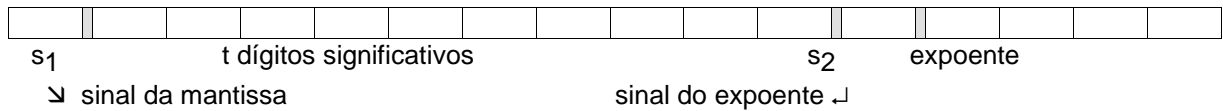
b) $(101,011)_2 = (1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} + 0 \times 2^{-4} + 1 \times 2^{-5} + 1 \times 2^{-6}) \cdot 2^3 = (0,101011)_2 \cdot 2^3$

Vamos agora exemplificar a representação digital de números em computador mostrando três exemplos práticos:

2.2. - PADRÃO 16 BITS

Vamos mostrar o sistema de Representação em Ponto Flutuante em uma máquina binária ($\beta = 2$), com $t = 10$ dígitos na mantissa e expoentes limitados entre $I = -15$ e $S = +15$ ($15_{10} = 1111_2$), de modo que simbolicamente temos: $F(\beta, t, I, S) = F(2, 10, -15, 15)_{10}$. Esta é a representação clássica da *máquina de 16 bits*, que é detalhada aqui por motivos históricos.

Representação esquemática de dígitos significativos binários, onde cada *bit* é alocado em um registro (célula quadrada):



Convenciona-se que:

- Se $s_1 = 0 \rightarrow$ número positivo.
- Se $s_1 = 1 \rightarrow$ número negativo.
- s_2 idem.

No registro total tem-se:

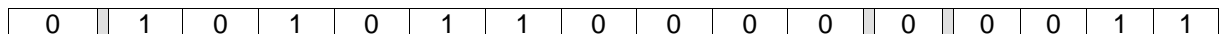
- 1 *bit* para sinal da mantissa .
- 10 *bits* para armazenar os dígitos significativos da mantissa ($t=10$).
- 1 *bit* para sinal do expoente.
- 4 *bits* para o módulo do expoente

Totalizando 16 *bits* neste registro.

Devemos notar também que os dígitos significativos são armazenados no padrão de normalização com $d_1 \neq 0$, conforme estabelecido anteriormente.

Ex. 17: Representar $+0,101011 \cdot 2^3$ na máquina de 16 *bits* estabelecida anteriormente.

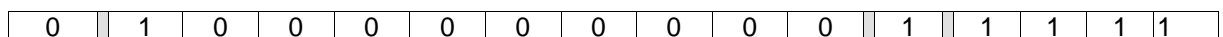
Convertendo-se o expoente: $(3)_{10} = (0011)_2$, tem-se



Limites da Representação em ponto Flutuante

Estes limites de representação serão exemplificados através da máquina de 16 bits:

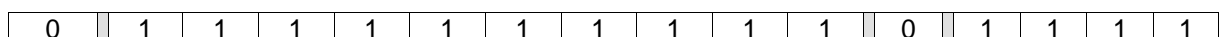
a). Menor positivo representável (m.p.):



└ Lembre-se de que toda representação na máquina de 16 bits usa normalização com padrão $d_1 \neq 0$.

$$m.p. = +(0,1)_2 \cdot 2^{-15} = (2^{-1} \cdot 2^{-15})_{10} = (2^{-16})_{10} = (0,0000152587)_{10}$$

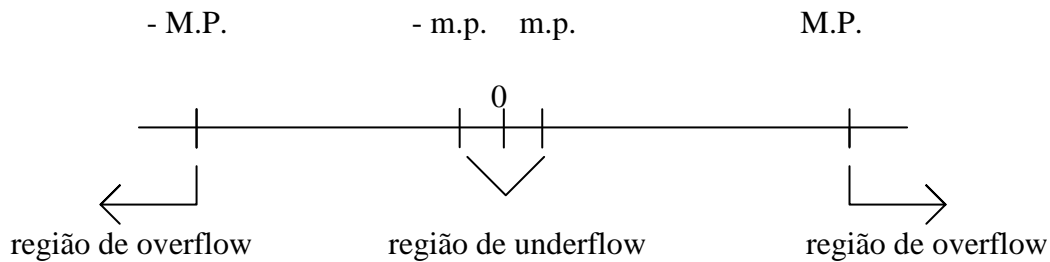
b) Maior positivo representável (M.P.):



$$M.P. = +(0,1111111111)_2 \cdot 2^{15} = (2^{-1} + 2^{-2} + 2^{-3} + \dots + 2^{-10}) \cdot 2^{15} = (32736)_{10} \cong (1 \cdot 2^{15})$$

Obs.: Os limites de representação dos números negativos são simétricos aos limites positivos apresentados.

Na reta real temos a seguinte representação para $F(2, 10, -15, +15)$:



Obs.:

- região de underflow: $\{x \in \mathbb{R} \mid -mp < x < mp\}$
- região de overflow: $\{x \in \mathbb{R} \mid x < -MP \text{ e } x > MP\}$

c) Representação do zero:

É obtida com mantissa nula e o menor expoente representável (l).

Ex. 18: Representar o zero em $F(2,10,-15,+15)$.

0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

expoente mínimo

Devemos lembrar que este é o único número escrito não normalizado, pois sua mantissa é zero.

No exemplo 20, a seguir, pode-se visualizar o que poderia acontecer se o expoente do zero fosse diferente do limite inferior l.

Ex. 19: Simular a operação de adição: $0,000135 + 0$ na máquina $F(10,4,-10,+10)$

\downarrow \downarrow
a b

1º) Considerando a representação do zero com expoente nulo ($b = 0,0000 \cdot 10^0$):

$$a = 0,1350 \cdot 10^{-3}$$

$$b = 0,0000 \cdot 10^0$$

$$a = 0,000135 \cdot 10^0$$

$$b = 0,0000 \cdot 10^0$$

$$a + b = 0,0001 \cdot 10^0 = 0,1000 \cdot 10^{-3}$$

2º) Considerando a representação do zero com expoente mínimo l ($b = 0,0000 \cdot 10^{-10}$):

$$a = 0,1350 \cdot 10^{-3}$$

$$b = 0,0000 \cdot 10^{-10}$$

$$a = 0,1350 \cdot 10^{-3}$$

$$b = 0,0000 \cdot 10^{-3}$$

$$a + b = 0,1350 \cdot 10^{-3} = a$$

Neste segundo caso o zero representado pela máquina digital representa corretamente o elemento neutro da operação de adição.

d). Número máximo de elementos representáveis:

Podemos notar que a distribuição de números representáveis em ponto flutuante é discreta (somente alguns valores são representáveis), enquanto a distribuição de valores na parte Real é contínua (qualquer valor é representável).

Ex. 20: Representar os dois primeiros números positivos do sistema $F(2, 10, -15, +15)$

1º positivo -

0	1 0 . . .	0	1	1 1 1 1
---	-----------	---	---	---------

= (0,0000152587)₁₀

2º positivo -

0	1 0 . . .	0 1	1	1 1 1 1
---	-----------	-----	---	---------

= (0,0000152885)₁₀

Ex. 21: Caso uma operação aritmética gere o número (0,00001527)₁₀, como ele será representado?

Como o valor acima não tem representação binária exata, ele será representado pelo valor discreto mais próximo, no caso (0,0000152587)₁₀ que é o menor positivo representável (mp).

Pode-se notar que a distribuição de números representáveis de $F(\beta, t, l, S)$ não é uniforme em \mathfrak{R} , e que para cada potência da base β existe uma quantidade fixa de números representáveis dada por:

$$NC = (\beta - 1) \cdot \beta^{t-1}$$

Ex. 22: Em $F(2, 3, -1, +2)$ temos as seguintes representações possíveis:

a) mantissas possíveis:

0,100
0,101
0,110
0,111

b) expoentes possíveis:

2⁻¹
2⁰
2⁺¹
2⁺²

A combinação de quatro possibilidades de mantissas em cada potência da base $((\beta - 1) \cdot \beta^{t-1} = 4 \text{ para } \beta = 2 \text{ e } t = 3)$, com as quatro possibilidades de expoentes $(S - l + 1 = 4 \text{ para } S = 2 \text{ e } l = -1)$ define o número total de positivos representáveis (NP = 16).

Desta forma o número total de elementos representáveis em uma máquina genérica $F(\beta, t, l, S)$ é dado por:

$$NF(\beta, t, l, S) = 2 \cdot (S - l + 1) \cdot (\beta - 1) \cdot \beta^{t-1} + 1$$

incluindo os positivos, negativos e o zero.

Ex. 23: Em $F(2, 10, -15, +15)$ (máquina de 16 bits) temos:

$$NF = 2 \cdot (2 - 1) \cdot (15 - (-15) + 1) \cdot 2^{10-1} = 31745 \text{ elementos incluindo os positivos, negativos e o zero.}$$

Ex. 24: Em $F(10, 10, -99, +99)$ (calculadora científica comum) temos:

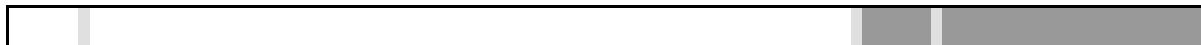
$$NF = 2 \cdot (10 - 1) \cdot (99 - (-99) + 1) \cdot 10^{10-1} + 1 = 3,582 \cdot 10^{12} \text{ elementos}$$

Esta representação da máquina padrão de 16 bits evoluiu, juntamente com os computadores, e atingiu uma forma mais otimizada de representação, incluindo a polarização dos expoentes, mais flexibilidade na normalização da mantissa, dentre outras. Mais tarde surgiu o padrão **IEEE 754** (1985), que é amplamente utilizado no armazenamento de variáveis (vide seção 4.1).

e). Polarização na Representação em Ponto Flutuante

Polarização (ou Excesso) é um valor acrescentado (em Excesso) a todos expoentes de um sistema de representação em ponto flutuante com o objetivo de tornar todos os expoentes positivos e ampliar a representação do expoente superior (S). Naturalmente todas as operações aritméticas devem considerar esta polarização introduzida.

Ex. 25: Na máquina de 16 *bits* F(2, 10, -15, 15) podemos usar uma polarização $p = +15$.



$$p = +15 = + (1111)_2$$

$$I + p = - (15)_{10} + p = - (1111)_2 + p = - (1111)_2 + (1111)_2 = (00000)_2$$

$$S + p = + (15)_{10} + p = + (1111)_2 + p = + (1111)_2 + (1111)_2 = (11110)_2$$

5 bits

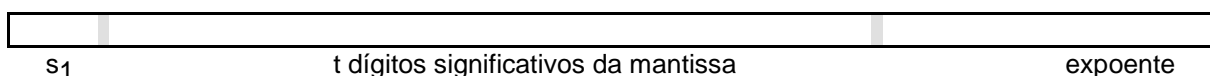
Como I e S têm agora o mesmo sinal, (+), podemos usar todos os registros binários reservados ao expoente, inclusive a posição do sinal, para representar o expoente polarizado (sem o sinal). Assim, os limites polarizados do expoente são:

$$I = (00000)_2 = (0)_{10}$$

$$S = (11110)_2 = (30)_{10}$$

Podemos aqui aproveitar melhor os 5 *bits* reservados ao expoente tomando o maior valor possível, adotando $S = (11111)_2 = (31)_{10}$.

Na forma polarizada qualquer número v representado nesta máquina deverá seguir a forma abaixo:



$$v = (-1)^S \cdot (0, \text{mantissa})_2 \cdot 2^{\text{exp} - 15}$$

Ex. 26: Na representação

0	1	1	0	1	0	0	0	0	0	0	0	1	0	0	1	0
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

tem-se: $s = 0$, $m = 110100000$ e $\text{exp} = (10010)_2 = (18)_{10}$

$$v = (-1)^0 \cdot (0, 1101000000)_2 \cdot 2^{18 - 15} = +(110,1)_2$$

Outra otimização adotada no padrão IEEE 754 foi o não armazenamento do primeiro bit da mantissa, que é sempre unitário, ou seja, usa-se a representação implícita do primeiro bit e abre-se uma posição binária para armazenamento de um novo bit:

$$v = (-1)^S \cdot (1, \text{mantissa})_2 \cdot 2^{\text{exp} - 15}$$

que veremos em detalhes na seção 4.1.

Exercícios:

2.1 - Na máquina F(2, 3, -3, +3) com $d_1 \neq 0$ (não polarizada) calcule:

- O número de elementos representáveis;
- Esquematize a representação de todos os elementos positivos na base 2;
- Defina as regiões de underflow e overflow;
- Estime a precisão decimal equivalente;

- e) Proponha uma transformação da máquina F apresentada em uma máquina com polarização que utilize os limites dos 3 *bits* totais reservados ao sinal.

2.2 - Na máquina $F(2,3,0,7)$ com $d_1 \neq 0$ e polarização $p = +3$ calcule:

- a) O número de elementos representáveis;
- b) Esquematize a representação de todos os elementos positivos na base 2;
- c) Defina as regiões de underflow e overflow;
- d) Estime a precisão decimal equivalente;

2.3). Precisão versus Exatidão

- **PRECISÃO**: é um conceito objetivo que estabelece a quantidade de algarismos significativos que representam um número. A precisão de uma máquina digital é definida como o número de dígitos t da mantissa na base β , e a precisão decimal “ d ” equivalente pode ser definida baseada na equivalência entre as variações dos dígitos menos significativos em cada base, da seguinte forma:

$$\begin{aligned}10^{1-d} &= \beta^{1-t} \\ \log(10^{1-d}) &= \log(\beta^{1-t}) \\ 1-d &= (1-t) \log \beta \\ d &= 1 + (t-1) \log \beta\end{aligned}$$

Ex. 27: Calcule a precisão decimal equivalente da máquina de 16 *bits* F(2, 10, -15, 15)

$$\begin{aligned}\beta &= 2 \quad t = 10 \\ d &= 1 + (10 - 1) \log 2 \\ d &= 1 + 9 \cdot (\log 2) \\ d &= 3,71 \\ d &= 3 - 4 \text{ dígitos (ou seja, pode representar entre 3 e 4 dígitos)}\end{aligned}$$

Ex. 28: Considere uma máquina F(2, 27, -20, 20)

$$\begin{aligned}2^{-26} &= 10^{-d+1} \Rightarrow \log_2 2^{-26} = \log_{10} 10^{-d+1} \\ -d + 1 &= -26 \log_2 \Rightarrow d = 1 + 26 \log_2 \approx 8,8\end{aligned}$$

Assim, esta última máquina tem entre 8 e 9 dígitos significativos equivalentes. Isto não significa que todas as frações decimais de 8 dígitos possam ser representadas precisamente em 27 *bits*, visto que a representação é discreta (com espaços vazios entre dois números consecutivos) e nem todas as frações decimais têm representação binária finita. Isto significa que todas as representações binárias da máquina estão “corretas” para 8 dígitos significativos na base 10, ou seja, apresentam decimais equivalentes com pelo menos 8 dígitos corretos.

- **EXATIDÃO**: conceito relacionado com a forma que melhor representa uma grandeza numérica, ou seja, uma representação é mais exata quando tem o menor desvio (erro) em relação ao valor exato.

Ex. 29: Representar o π (3,1415926535) por:

- (a) 3,14 → precisão de três dígitos
- (b) 3,151 → precisão de quatro dígitos
- (c) 3,1416 → precisão de cinco dígitos

Obs.: Note que se fossemos classificar o números acima quanto a exatidão teríamos o seguinte:

- a) é mais exato que (b) (ou seja, (a) está mais próximo de π do que (b));
- b) é menos exato que (c);
- c) é mais exato que (a).

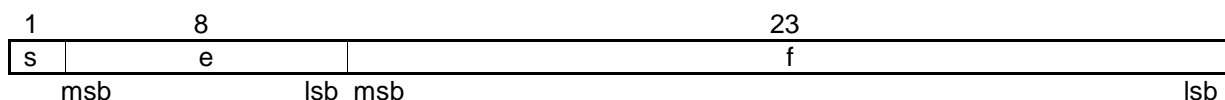
3 – REPRESENTAÇÃO NUMÉRICA SEGUNDO O PADRÃO IEEE 754 (1985)

Este padrão é utilizado em linguagens comerciais como o Pascal e C, e será apresentado a seguir na sua forma esquemática para representação em ponto flutuante e na forma de variáveis inteiras.

4.1). VARIÁVEL SINGLE DO PASCAL (OU FLOAT DO C):

Padrão: 4 bytes ou 32 bits (precisão de 7 a 8 dígitos significativos equivalentes).

Neste padrão um número real v pode ser representado por:



onde

$s = 0 \Rightarrow v$ positivo e $s = 1 \Rightarrow v$ negativo

e = expoente

f = mantissa

polarização = $(127)_{10} = 2^7 - 1 = (01111111)_2$

msb = bit mais significativo e lsb = bit menos significativo

Um número v armazenado no registro acima é interpretado da seguinte forma:

- Se $0 < e < 255$, então $v = (-1)^s \cdot 2^{(e-127)} \cdot (1,f)_2$
- Se $e = 0$ e $f \neq 0$, então $v = (-1)^s \cdot 2^{-126} \cdot (0,f)_2$
- Se $e = 0$ e $f = 0$, então $v = (-1)^s \cdot 2^{-126} \cdot (0,0) = (-1)^s \cdot 0$ (zero)
- Se $e = 255$, então v pertence a região de overflow.

Obs.:

(i). A representação destes registros binários em computadores digitais é feita em grupos de bytes (8 bits) escritos de forma invertida (de traz para frente) em relação ao esquema binário apresentado acima. Neste exemplo tem-se quatro bytes, onde cada byte é composto por dois registros hexadecimais (8 bits).

A seguir apresenta-se um exemplo desta representação para a fração $1/10$, que representada na variável SINGLE do Pascal gera:

$x = 0.10000000149$ - (representação decimal, note o erro de arredondamento)

(ii). Se na janela 'watch' do Pascal, onde se pode visualizar as variáveis na base hexadecimal escrevendo 'x,m', tem-se a seguinte representação hexadecimal de uma variável SINGLE:

$x,m = CD\ CC\ CC\ 3D$ - (representação hexadecimal no computador)

Esta deve ser interpretada na forma de bytes em ordem invertida, para compor o registro binário correspondente. Para obter este registro procede-se da seguinte forma:

a). Invertamos os bytes (grupos de dois hexadecimais):

$3D\ CC\ CC\ CD$

b). Efetua-se a conversão direta para base binária:

c). Distribui-se os bits no registro SINGLE:

d). Interpretando os 32 bits acima pode-se converter o registro para decimal, conforme segue:

Os registros sublinhados representam os arredondamentos gerados na representação SINGLE de $1/10$.

e). Para representar $(0,10)_{10}$ na base binária no padrão IEEE de 32 bits temos duas possibilidades:

e.1). Converter $(0,10)_{10}$ para binário e adequar o resultado ao padrão 32 bits:

$$(0,10)_{10} = (0,0001100110011001100110011001100110011001100110011...)_{2}$$
$$(0,10)_{10} = 2^{-4} \cdot (1,100110011001100110011001100110011001100110011...)_{2}$$

Numero positivo $s=0$.

$$\begin{aligned}
 (0,10)_{10} &= (-1)^0 \cdot 2^{-4} \cdot (1,10011001100110011001100 \mid 110011001100110011...)_{2} \\
 &\quad \rightarrow \text{ARREDONDAMENTO na mantissa f} \\
 &\quad (0,110011001100110011...)_{2} > (0,5)_{10} \\
 &\quad \quad \quad +1 \\
 (0,10)_{10} &\cong (-1)^0 \cdot 2^{-4} \cdot (1,10011001100110011001100)_{2} \\
 \hline
 (0,10)_{10} &\cong (-1)^0 \cdot 2^{-4} \cdot (1,10011001100110011001101)_{2} \\
 (e-127)_{10} &= (-4)_{10} \Rightarrow e = (123)_{10} = (01111011)_{2}
 \end{aligned}$$

[illegible]

e.2). Converter $(0,10)_{10}$ diretamente para o padrão 32 bits, considerando $0 < e < 255$:

$$(0,10)_{10} = (-1)^s \cdot 2^{(e-127)} \cdot (1,f)_2$$

(e.2.1). $s = 0 \Rightarrow$ para números positivos $+(0,10)_{10}$

(e.2.2). Considerando que ainda temos 2 incógnitas, 'e' e 'f', vamos considerar o valor Mínimo de 'f', $f = 0$, e calcular 'e':

$$(0,10)_{10} = (-1)^0 \cdot 2^{(-127)} \cdot (1,0)_2$$

No caso, determina-se um valor de 'e' maior que o verdadeiro, no caso, 'e' = (123,6781...) ₁₀ e toma-se o seu menor inteiro, 'e' = (123) ₁₀ (DESCONSIDERAR A PARTE FRACIONARIA)

(e.2.3). Agora determinamos o valor de 'f' com os valores de 's' e 'e' obtidos acima:

$$m_p = 1,4012985 \cdot 10^{-45}$$

$$\text{iii). Maior positivo (MP):} \quad \begin{cases} s = 0 \\ e = (11111110)_2 = (254)_{10} \\ f = 11111111111111111111111111111111 \end{cases}$$

[illegible]

$$MP = (-1)^0 \cdot 2^{254-127} \cdot (1,111111111111111111111111)_2$$

$$MP = 2^{127} \cdot (1,999999988)_{10}$$

$$MP = (3,4028235 \cdot 10^{38})_{10}$$

4.2) Variável DOUBLE do Pascal (ou DOUBLE do C):

Padrão: 8 bytes ou 64 bits (precisão de 16 a 17 dígitos significativos equivalentes).

onde

polarização = $(1023)_{10} = 2^{10} - 1 = (0111111111)_2$

Um número v armazenado no registro acima é interpretado da seguinte forma:

- Se $0 < e < 2047$, então $v = (-1)^s \cdot 2^{(e-1023)} \cdot (1, f)_2$
- Se $e = 0$ e $f \neq 0$, então $v = (-1)^s \cdot 2^{-1022} \cdot (0, f)_2$
- Se $e = 0$ e $f = 0$, então $v = (-1)^s \cdot 2^{-1022} \cdot (0,) = (-1)^s \cdot 0$ (zero)
- Se $e = 2047$, então v pertence a região de overflow.

4.3). VARIÁVEL EXTENDED DO PASCAL (OU LONG DOUBLE DO C):

Padrão: 10 bytes ou 80 bits (precisão de 19 a 20 dígitos significativos equivalentes).

Diagram illustrating the structure of a 64-bit word. The word is divided into four fields: 's' (1 bit, msb), 'e' (15 bits), 'i' (1 bit, lsb), and 'f' (63 bits, msb). The 's' field is labeled 'msb' and the 'i' field is labeled 'lsb'.

onde

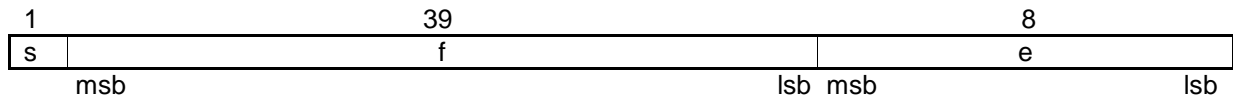
$$\text{polarização} = (16383)_{10} = 2^{14} - 1 = (01111111111111)_{2}$$

Um número v armazenado no registro acima é interpretado da seguinte forma:

- Se $0 < e < 32767$, então $v = (-1)^s \cdot 2^{(e-16383)} \cdot (i, f)_2$ (onde i pode assumir 0 ou 1)
(se $e = 0 \Rightarrow i = 1$)
- Se $e = 32767$ e $f = 0$, então v pertence a região de overflow.

4.4) Variável REAL (padrão BORLAND):

Padrão: 6 bytes ou 48 bits (precisão de 12 a 13 dígitos significativos equivalentes).



onde

$$\text{polarização} = (129)_{10} = 2^7 + 1 = (10000001)_2$$

Um número v armazenado no registro acima é interpretado da seguinte forma:

- Se $0 < e \leq 255$, então $v = (-1)^s \cdot 2^{(e-129)} \cdot (1,f)_2$
- Se $e = 0$, então $v = 0$, (independe de f)

Ex. 41: Simulação do algoritmo para avaliação da precisão decimal equivalente de uma variável.

```

p1=1
repita
    p1=p1/2
    p2=1+p1
até p2=1
    
```

Devemos ressaltar que no caso da operação soma os expoentes devem estar alinhados pelo maior expoente para que a soma possa ocorrer, independentemente da normalização da representação em ponto flutuante (vide cap. 3). Assim,

$$i). \text{ Representação da unidade: } 1 = \begin{cases} e = (01111111)_2 = (127)_{10} \\ f = 00000000000000000000000000000000 \end{cases}$$

0	0 1 1 1 1 1 1 1	0 0
(00 00 80 3F - representação em hexadecimal)		

$$1 = (-1)^0 \cdot 2^{(127-127)} \cdot (1,00000000000000000000000000000000)_2$$

$$1 = 2^0 \cdot (1,00000000000000000000000000000000)_2$$

ii). Representação do menor número (P1) que pode ser somado a unidade (note que no momento da soma os expoentes devem estar alinhados para que a soma possa ocorrer no processador aritmético e portanto o expoente de p1 deve ser o mesmo da representação da unidade):

$$P1 = \begin{cases} e = (01111111)_2 = (127)_{10} & - \text{ expoente igual a unidade (item i)} \\ f = 00000000000000000000000000000001 & - \text{ menor mantissa possível} \end{cases}$$

$$P1 = (-1)^0 \cdot 2^{(127-127)} \cdot (0,00000000000000000000000000000001)_2 \quad - \quad \text{este valor de P1 é conseguido após 23 divisões binárias.}$$

Então,

$$1 = 2^0 \cdot (1,00000000000000000000000000000000)_2$$

$$P1 = 2^0 \cdot (0,00000000000000000000000000000001)_2$$

$$1+P1 = 2^0 \cdot (1,00000000000000000000000000000001)_2$$

$$\text{Logo, } P1 = 2^0 \cdot 2^{-23} = 2^{-23} = 1,1921 \cdot 10^{-7}$$

Neste caso a precisão decimal equivalente está entre 7 e 8 dígitos significativos.

$$P1 = \begin{cases} e = (01101000)_2 = (104)_{10} \\ f = 000000000000000000000000 \end{cases}$$

Após a 23ª divisão de $P1=1$ por 2, $P1$ representará um número desprezível na soma frente a unidade, pois $1 + P1 = 1$.

4.1). a). Avalie as regiões de underflow e overflow de cada uma das variáveis apresentadas acima e faça uma verificação em um compilador Pascal ou C;
b). Avalie a precisão decimal equivalente de cada variável, através da fórmula de equivalência;

21

4.5). REPRESENTAÇÃO DE VARIÁVEIS INTEIRAS

Os formatos para representação de variáveis do tipo inteiras podem seguir diversos padrões:

4.5.1). Shortint: são tipos inteiros limitados aos valores -128_{10} e $+127_{10}$, sua representação interna é feita em forma de *byte* (8 *bits*) com *bit* de sinal. Valores negativos são armazenados em forma de complemento de dois.

$$\begin{aligned}\text{Ex. 42: } +0_{10} &= 00000000_2 = 00_{16} \\ +127_{10} &= 01111111_2 = 7F_{16} \\ -128_{10} &= 2^8 - 128 = 128 = 10000000_2 = 80_{16} \\ -5_{10} &= 2^8 - 5 = 123 = 11111011_2 = FC_{16}\end{aligned}$$

Obs.: Veja como é feito o complemento de dois no exemplo 32.

Ex. 43: Se for efetuada a operação de adição entre $+127_{10}$ e -127_{10} , em forma de complemento e dois, tem-se:

$$\begin{array}{r} \left\{ \begin{array}{l} +127 \\ -127 \end{array} \right\} \begin{array}{l} 01111111 \\ 10000001 \end{array} \\ \hline 00000000 \end{array}$$

4.5.2). byte: são tipos inteiros sem bit de sinal (não permite armazenar negativos). Estão limitados aos valores 0 e 255.

$$\begin{aligned}\text{Ex. 44: } 0_{10} &= 0000\ 0000_2 = 00_{16} \\ 255_{10} &= 1111\ 1111_2 = FF_{16}\end{aligned}$$

4.5.3). Integer: são tipos inteiros limitados à faixa entre -32768 e $+32767$, correspondendo ao armazenamento como 2 *bytes* com *bit* de sinal ("negativos" são armazenados em forma de complemento de dois). **Integer** equivale ao '**Integer**' da linguagem Pascal.

$$\begin{aligned}\text{Ex. 45: } \\ \text{Zero} \quad 0_{10} &= 0000\ 0000\ 0000\ 0000 = 0000_{16}\end{aligned}$$

$$\text{Maior Positivo } +32767_{10} = 0111\ 1111\ 1111\ 1111_2 = 7FFF_{16}$$

$$\begin{aligned}-32767_{10} &= -0111\ 1111\ 1111\ 1111_2 \\ &\quad 1000\ 0000\ 0000\ 0000_2 \rightarrow \text{complemento de um de } -32767_{10} \\ &\quad \quad \quad +1 \\ &\quad 1000\ 0000\ 0000\ 0001_2 \rightarrow \text{complemento de dois de } -32767_{10} \\ &= 8001_{16}\end{aligned}$$

$$\begin{aligned}\text{Menor Negativo } -32768_{10} &= -1000\ 0000\ 0000\ 0000_2 \\ &\quad 0111\ 1111\ 1111\ 1111_2 \rightarrow \text{complemento de um de } -32768_{10} \\ &\quad \quad \quad +1 \rightarrow \text{soma 1} \\ &\quad 1000\ 0000\ 0000\ 0000_2 \rightarrow \text{complemento de dois de } -32768_{10} \\ &= 8000_{16}\end{aligned}$$

Agora observe o que acontece se for executada a soma da 'unidade' com o 'maior positivo' da variável **Integer** 32767_{10} , ou seja, estamos calculando um número na região de **overflow**:

$$\begin{aligned}\text{Unidade} \quad 1_{10} &= 0000\ 0000\ 0000\ 0001_2 = 0001_{16} \\ &+ \\ \text{Maior Positivo } +32767_{10} &= 0111\ 1111\ 1111\ 1111_2 = 7FFF_{16}\end{aligned}$$

$$1_{10} + 32767_{10} = 1000\ 0000\ 0000\ 0000_2 = 8000_{16} \rightarrow \text{Complemento de 2 (inicia com 1)}$$

Precisamos obter o complemento de 2 novamente para obter o número 'negativo' armazenado:

$$\begin{array}{rcl}
 0111\ 1111\ 1111\ 1111_2 & \rightarrow & \text{complemento de 1 de } 1000\ 0000\ 0000\ 0000_2 \\
 \hline
 & + 1 & \rightarrow \text{soma 1} \\
 -1000\ 0000\ 0000\ 0000_2 & \rightarrow & \text{complemento de 2 de } 1000\ 0000\ 0000\ 0000_2 \\
 = -8000_{16} \\
 = -32768_{10}
 \end{array}$$

Então, cuidado podemos achar que $1_{10} + 32767_{10} = +32768_{10}$, mas estamos armazenando -32768₁₀.

4.5.4). Word: armazenamento de 2 bytes sem *bit* de sinal limitado a faixa entre 0 e 65535.

Ex. 46: $0_{10} = 0000_{16}$
 $65535_{10} = 1111\ 1111\ 1111\ 1111_2 = \text{FFFF}_{16}$
 $65536_{10} = 0000_{16} = 0_{10}$

Obs.: Note que a representação de números inteiros acima do limite superior acarreta uma grande perda de significação, cuja representação volta ao zero (vide ex. 45).

4.5.5). Longint: corresponde ao "double word" (4 bytes) com bit de sinal. Limita-se entre -2147483648 (-2^{31}) e +2147483647 ($2^{31} - 1$). **Longint** equivale ao 'Int' da linguagem C/C++.

Ex. 47: $0_{10} = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_2 = 00\ 00\ 00\ 00_{16}$

$+2147473647_{10} = 0111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111_2 = 7F\ FF\ FF\ FF_{16}$

$-2147483647_{10} = -0111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111_2$
 $1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_2 \rightarrow \text{complemento de 1}$
 \hline
 $1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0001_2 = 80000001_{16} \rightarrow \text{complemento de 2}$

$-2147483648_{10} = -1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_2$
 $0111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111_2 \rightarrow \text{complemento de 1}$
 \hline
 $1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_2 = 80000000_{16} \rightarrow \text{complemento de 2}$

Observe novamente que se for executada a soma da 'unidade' com o 'maior positivo' da variável **Longint** 2147473647₁₀, estaremos calculando um número na região de overflow da variável **Longint**.

Unidade $1_{10} = 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0001_2 = 00\ 00\ 00\ 01_{16}$
 $+$
 Maior Positivo $+2147473647_{10} = 0111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111_2 = 7F\ FF\ FF\ FF_{16}$
 $1_{10} + 2147473647_{10} = 1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_2 = 80\ 00\ 00\ 00_{16} \rightarrow$

É um complemento de 2 (inicia com 1)

Precisamos obter o complemento de 2 novamente para obter o número 'negativo' armazenado:

$$\begin{array}{rcl}
 1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_2 & & \\
 0111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111\ 1111_2 & \rightarrow & \text{complemento de 1} \\
 \hline
 & + 1 & \rightarrow \text{soma 1} \\
 -1000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000\ 0000_2 & \rightarrow & \text{complemento de 2} \\
 = -80\ 00\ 00\ 00\ 00\ 00\ 00\ 00_{16} \\
 = -2147483648_{10}
 \end{array}$$

Então, cuidado podemos achar que $1_{10} + 2147483647_{10} = +2147483648_{10}$, mas estamos armazenando -2147483648₁₀.

Obs.: Note que o armazenamento de inteiros negativos é sempre feito em forma de complemento de dois. Isto é uma vantagem do ponto de vista de operações aritméticas, pois o carregamento na unidade aritmética já está na sua forma final com expoentes iguais (alinhados) para efetuar as operações de adição (ou subtração). Note que isto só é possível, pois os expoentes equivalentes na notação de Tipos Inteiros são todos iguais a zero.

Por exemplo, tem-se a seguinte nomenclatura em C:

Type	Bytes	Bits	Range	
short int	2	16	-32,768 -> +32,767	(16kb)
unsigned short int	2	16	0 -> +65,535	(32Kb)
unsigned int	4	16	0 -> +4,294,967,295	(4Gb)
int	4	32	-2,147,483,648 -> +2,147,483,647	(2Gb)
long int	4	32	-2,147,483,648 -> +2,147,483,647	(2Gb)
signed char	1	8	-128 -> +127	
unsigned char	1	8	0 -> +255	
float	4	32		
double	8	64		
long double	12	96		

4). TIPOS DE ERROS EXISTENTES EM MÁQUINAS DIGITAIS

É muito importante conhecer as possibilidades de erros na representação numérica em máquinas digitais e entender as suas causas para se poder estabelecer a confiabilidade de um software.

Todo estudo apresentado neste capítulo é necessário para que se possa entender as causas de cada tipo de erro existente em máquinas digitais.

Pode-se classificar os erros nos seguintes tipos principais:

5.1) Erros Inerentes:

São aqueles existentes nos dados de entrada de um software numérico. Decorre, por exemplo, de medições experimentais, de outras simulações numéricas, ...

5.2) Erros de truncamento:

Ocorrem quando quebramos um processo matematicamente infinito, tornando-o finito, por incapacidade de execução ou armazenamento.

A seguir serão apresentados exemplos de fontes de erro de truncamento:

Ex. 48: Veja a seguinte expansão em série infinita para e^x :

$$e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!} + \dots$$

Sabe-se que não é possível usar infinitos termos para avaliar uma função, então é necessário estabelecer um limite para o número de parcelas utilizadas. Esta limitação nas parcelas gera um erro de truncamento na série, que corresponde ao somatório dos termos abandonados.

Obs.: A representação de $f(x)$ em série será mostrada de forma ilustrativa a seguir,

É possível representar, de forma exata, uma função $f(x)$ em um ponto qualquer x_0+x a partir de sua representação em x_0 , através de expansão em Séries de Taylor, dada genericamente por:

$$f(x_0 + x) = f(x_0) + f'(x_0) \cdot \frac{x}{1!} + f''(x_0) \cdot \frac{x^2}{2!} + f'''(x_0) \cdot \frac{x^3}{3!} + \dots + f^{(n)}(x_0) \cdot \frac{x^n}{n!} + \dots$$

Expandindo a função, por exemplo, em torno de $x_0 = 0$ (quando $x_0=0$, temos o caso particular da série Maclaurin), tem-se:

$$\begin{aligned}
f(x) &= e^x \\
f(x_0 = 0) &= e^0 = 1 \\
f'(x_0 = 0) &= e^0 = 1 \\
f''(x_0 = 0) &= e^0 = 1 \\
&\vdots \\
f^n(x_0 = 0) &= e^0 = 1
\end{aligned}$$

Gerando então,

$$f(0+x) = e^{x+0} = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!} + O(x^{(n+1)}) = \sum_{i=0}^n \frac{x^i}{i!} + O(x^{(n+1)})$$

Se aproximamos e^x usando até o termo de ordem 'n', estamos desprezando os termos

$$O(x^{(n+1)}) = \frac{x^{n+1}}{(n+1)!} + \frac{x^{n+2}}{(n+2)!} + \dots$$

$$\text{então } e^x \cong 1 + \frac{x}{1!} + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!}$$

Assim, o termo $O(x^{(n+1)})$ caracteriza o Erro de Truncamento da aproximação, lembrando que x, nesse caso, representa o incremento de x entre o ponto inicial x_0 e o novo ponto x_0+x .

Trabalho:

Dado $e^1 = 2.718281828459045235360287471352662497757247093699959574966967\dots$

- Calcule o Erro exato cometido quando se aproxima e^x , via série de Taylor, usando até o termo de ordem 'n=7', para x=1, em uma variável IEEE de 32 bits (tipo float);
- Calcule os erros Estimados de Arredondamento, Truncamento e Total, quando não dispomos do valor exato.

Ps.:

- O valor exato "Estimado", com pouco Arredondamento, pode ser calculado aproximando e^x , via série de Taylor, usando até o termo de ordem 'n' em uma variável IEEE de 64 bits (tipo double), cujos arredondamentos são muito inferiores ($\text{Erro}_{\text{Arred}} = |VA_{32\text{bits},n} - VE_{64\text{bits},n}|$);
- O valor exato "Estimado", com pouco Truncamento, pode ser calculado aproximando e^x via série de Taylor, usando até o termo de ordem 'n₂=2*n' com variáveis IEEE de 64 bits (tipo double), para não ter influência dos arredondamentos ($\text{Erro}_{\text{Truncam}} = |VA_{64\text{bits},n} - VE_{64\text{bits},2*n}|$);
- O erro Total, quando não dispomos do valor exato, pode ser obtido por:
 $\text{Erro}_{\text{Total}} = |VA_{32\text{bits},n} - VE_{64\text{bits},2*n}|$, que considera variações no arredondamento das variáveis e no truncamento do processo de aproximação em série.

Ex. 49: Aproximações numéricas de limites de funções,

Por definição $f'(x)$ é dada por,

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

Porém, se este limite exato não puder ser obtido, pode-se promover uma aproximação numérica deste, tomando o incremento h como finito e promovendo sucessivos refinamentos. Assim, pode-se obter uma sequência de aproximações sucessivas de $f'(x)$, com incremento cada vez menor, mas não se pode chegar ao incremento nulo ($h \rightarrow 0$). Então, também se quebra o processo matemático, de refinamentos sucessivos, gerando um erro de truncamento do processo, que era matematicamente infinito, tornando-o finito.

Ex. 50: Aproximações de derivadas

Para avaliar numericamente $f'(x_0)$ a partir de três pontos vizinhos de $f(x)$:

$$\begin{array}{ccc} f(x_0 - h) & f(x_0) & f(x_0 + h) \\ \hline x_0 - h & x_0 & x_0 + h \end{array}$$

Pode-se subtrair $f(x_0 - h)$ de $f(x_0 + h)$:

$$f(x_0 + h) = f(x_0) + f'(x_0)h + f''(x_0)\frac{h^2}{2!} + f'''(x_0)\frac{h^3}{3!} + f^{iv}(x_0)\frac{h^4}{4!} + \dots + f^n(x_0)\frac{h^n}{n!}$$

-

$$f(x_0 - h) = f(x_0) - f'(x_0)h + f''(x_0)\frac{h^2}{2!} - f'''(x_0)\frac{h^3}{3!} + f^{iv}(x_0)\frac{h^4}{4!} - \dots + f^n(x_0)\frac{h^n}{n!}$$

$$f(x_0 + h) - f(x_0 - h) = f'(x_0)(2h) + 2.f'''(x_0)\frac{h^3}{3!} + 2.f^{iv}(x_0)\frac{h^4}{4!} + \dots$$

$$f'(x_0) = \frac{f(x_0 + h) - f(x_0 - h)}{2h} - \frac{1}{h}.f'''(x_0)\frac{h^3}{3!} - \frac{1}{h}.f^{v}(x_0)\frac{h^5}{5!} + \dots$$

$$f'(x_0) = \frac{f(x_0 + h) - f(x_0 - h)}{2h} - O(h^2)$$

onde o termo de segunda ordem $O(h^2)$ representa o somatório de todos os termos decorrentes da aproximação em série, e é definido como Erro de Truncamento da aproximação.

$$O(h^2) = f'''(x_0)\frac{h^2}{3!} + f^{v}(x_0)\frac{h^4}{5!} + \dots$$

Então, desprezando o termo $O(h^2)$, assumindo assim um erro de truncamento de segunda ordem, tem-se uma aproximação para $f'(x_0)$ dada por,

$$e \quad f'(x_0) \cong \frac{f(x_0 + h) - f(x_0 - h)}{2h}$$

5.3). Erros de Arredondamento

Ocorrem quando são desprezados os últimos dígitos que, ou não são fisicamente significativos na representação numérica, ou estão além da capacidade de armazenamento na máquina digital.

5.3.1). Arredondamento manual:

Ex. 51: Representar os seguintes números com quatro dígitos significativos:

- 69,348 = 69,35 → parcela descartada é maior que 5 ⇒ +1 no dígito anterior.
 69,34433 = 69,34 → parcela descartada é menor que 500 ⇒ dígito anterior inalterado.
 69,335 = 69,34 → parcela descartada é igual a 5 e dígito anterior ímpar ⇒ +1 no dígito anterior.
 69,345 = 69,34 → parcela descartada é igual a 5 e dígito anterior par ⇒ dígito anterior inalterado.

No exemplo anterior o arredondamento foi feito de forma ponderada, baseado em critérios estatísticos para descartar parte dos dígitos do número. Pode-se seguir o seguinte raciocínio:

Analisando estatisticamente um conjunto de números que terão parcelas descartadas. Podemos admitir que, em uma distribuição normal de erros, 50% dos valores a descartar são maiores que 5 e que 50% são menores, então

- Se todos os dígitos a descartar forem simplesmente cancelados, sem nenhum critério de compensação, o conjunto inicial de números perde parte de seu significado, gerando um erro de arredondamento global.

- Se, por outro lado, procurar-se distribuir o erro de arredondamento entre os números do conjunto escolhido, pode-se minimizar os efeitos globais dos erros de arredondamento sobre este conjunto. Assim, promove-se uma atualização do dígito anterior ao descartado nos 50% dos casos cuja parcela descartada é maior que 5, adicionando-se uma unidade. Nos demais casos simplesmente descarta-se a parcela indesejada, sem nenhuma atualização no dígito anterior. Desta forma distribui-se, estatisticamente, a parcela perdida entre os elementos do conjunto.

No caso específico de parcelas descartadas iguais a 5, tem-se um impasse que deve também ser decidido estatisticamente. Nestes casos também dividem-se as possibilidades em dois grupos, um com parcela anterior par e outro com parcela anterior ímpar. No primeiro grupo simplesmente descartam-se os dígitos indesejados, sem nenhuma atualização no dígito anterior, e segundo grupo, com parcela anterior ímpar, atualiza-se este dígito adicionando-lhe uma unidade. Também nesta situação procurou-se distribuir estatisticamente o erro devido as parcelas descartadas.

O arredondamento também pode ser feito por cancelamento puro, onde a parte indesejada do número é simplesmente cancelada, independente do seu valor, assumindo um erro de arredondamento global para valores menores em todos os elementos de um conjunto de números.

5.3.2). Arredondamento em máquinas digitais:

Neste caso o arredondamento pode ocorrer nas seguintes situações básicas:

(i). Armazenamento de racionais ilimitados:

Ex. 52: Representar a fração $(1/3)_{10}$ em F(10,6,-99,+99) (normalização: vírgula antes do 1º dígito).

$$(1/3)_{10} = (0, \underbrace{33333333}_{t=6} \dots 10^0)_{10}$$

$$(1/3)_{10} \cong (0,333333.10^0)_{10} \Rightarrow \text{Representação arredondada de } (1/3)_{10}.$$

Ex. 53: Representar a fração decimal $(1/10)_{10}$ na máquina binária F(2,10,-15,+15).

$$\left(\frac{1}{10}\right)_{10} = \left(\frac{1}{1010}\right)_2 = (0,00011001100110011\dots)_2 \Rightarrow \text{dízima periódica em base binária}$$

$$(1/10)_{10} = (0, \underbrace{11001100110011\dots}_{t=10})_2 \cdot 2^{-3}$$

$$\left(\frac{1}{10}\right)_{10} \cong (0,1100110011)_2 \cdot 2^{-3}$$

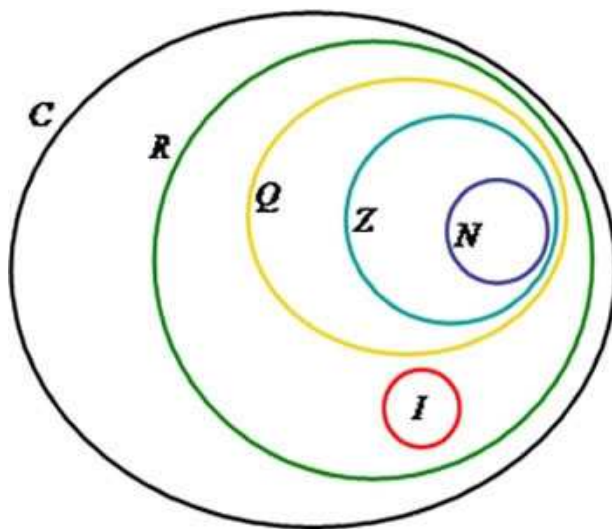
Obs.: Note que uma fração decimal exata ($1/10 = 0,1$) quando armazenada em uma máquina binária se transforma em uma fração binária periódica, que deve ser aproximada devido à limitação do registro em ponto flutuante utilizado no armazenamento.

(ii). Armazenamento de Irracionais:

O conjunto dos números irracionais compreende todas as representações através de dízimas

não periódicas e infinitas.

Conjuntos:



N - Naturais
Z - Inteiros
Q - Racionais
I - Irracionais
R - Reais
C - Complexos

Obs.: Existem representações que generalizam todos os números como complexos, de modo que um complexo com parte imaginária nula se torna um real.

Ex. 54: Representar em F(10,6,-99,+99),

$$(a). \pi \Rightarrow 3.141592653589\dots$$

$$\pi \Rightarrow (0,3141592653589\dots)_{10} \cdot 10^{+1}$$

$$\pi \equiv (0,314159)_{10} \cdot 10^{+1}$$

$$(b) \sqrt{2} \Rightarrow 1.414213562373\dots$$

$$\sqrt{2} \Rightarrow (0,1414213562373\dots)_{10} \cdot 10^{+1}$$

$$\sqrt{2} \equiv (0,141421)_{10} \cdot 10^{+1}$$

(iii). Abrangência limitada da notação em ponto flutuante:

Ex. 55: Efetue a soma de $a = 0,0135$ e $b = 10,51$ em F(10,4,-10,+10) e $g = 0$.

Representação em ponto flutuante:

$$a = 0,1350 \cdot 10^{-1}$$

$$b = 0,1051 \cdot 10^2$$

Vamos implementar a soma, de forma simplificada, usando alinhamento pelo maior expoente

$$a = 0,1350 \cdot 10^{-1}$$

$$b = 0,1051 \cdot 10^2$$

$$a = 0,000135 \cdot 10^2$$

+

$$b = 0,1051 \cdot 10^2$$

$$a + b = 0,1052 \cdot 10^2$$

Ex. 56: Efetue a soma de $a = (10,01)_2$ e $b = (0,0101)_2$ em F(2,4,-15,+15) e $g = 0$.

Representação em ponto flutuante:

$$a = (0,1001)_2 \cdot 2^2$$

$$b = (0,1010)_2 \cdot 2^{-1}$$

Vamos novamente implementar a soma, de forma simplificada,

$$a = (0,1001)_2 \cdot 2^2$$

$$b = (0,0001010)_2 \cdot 2^2$$

$$a + b = (0,1010)_2 \cdot 2^2$$

Obs.: Estes fatos ocorrem, geralmente, na soma de números com potências muito diferentes. Neste caso o número de menor potência pode perder significação, total ou parcial, frente ao número de maior potência. Ou seja, devido a faixa limitada de abrangência dos registradores em ponto flutuante, o número menor perde dígitos significativos quando comparado com o número maior.

(iv). Mudança de base para armazenamento e operações aritméticas

Sabe-se que a representação de números em base binária é amplamente utilizada em máquinas digitais (computadores), devido as suas vantagens no armazenamento e implementação de operações aritméticas.

O que ocorre na prática é que a interface entre o usuário e os computadores deve ser feita em base decimal, para que a representação de grandezas físicas seja naturalmente entendida pelos usuários dos computadores.

Então toda grandeza física é expressa inicialmente em base decimal, e o seu efetivo armazenamento nos computadores é feito em base binária, por isso é necessária uma conversão entre as bases decimal e binária e vice-versa.

Ex. 57: Representar $(0,1)_{10}$ em $F(2,10,-15,+15)$.

$$(0,1)_{10} = (0,00011001100110011...)_{2} \Rightarrow \text{fração decimal exata gerou dízima periódica binária}$$

$$(0,1)_{10} = 2^{-4} \cdot (1,100110011\underline{0011}...)_{2} \rightarrow \text{representação exata, anterior ao arredondamento}$$

$$(0,1)_{10} \cong 2^{-4} \cdot (1,100110011...)_{2} \cong (0,099975585)_{10} \rightarrow \text{representação aproximada, pós arredondamento}$$

Nesse caso, o Erro de Arredondamento pode ser calculado, pois temos o Valor Exato (VE) original e o Valor Aproximado (VA), pós arredondamento:

$$VE = (0,1)_{10}$$

$$VA = 2^{-4} \cdot (1,100110011...)_{2} = (0,099975585)_{10}$$

$$\text{Erro Relativo \%} = \left| \frac{VA - VE}{VE} \right| \cdot 100\% = -0,02441\%$$

Consequências:

Os erros de arredondamento podem causar:

(a). Perda de significação:

Esta é uma consequência de erros de arredondamento, que gera perda, total ou parcial, de dígitos significativos.

Esta perda de dígitos significativos pode ocorrer nos seguintes casos:

(a1). Soma de parcelas de grandezas muito diferentes:

Vide exemplos 55 e 56 apresentados anteriormente.

(a2). Subtração de parcelas de grandezas muito próximas:

Ex. 58: Efetuar $a - b$ com $a = 0,1351$ e $b = 0,1369$ em $F(10,4,-10,+10)$ e $g = 0$.

Efetuando a subtração de forma simplificada tem-se:

$$\begin{array}{r} a = 0,1351 \\ - b = - 0,1369 \\ \hline a - b = - 0,0018 = - 0,1800 \cdot 10^{-2} \end{array}$$

Note que o resultado final não sofreu arredondamentos, mas perdeu dígitos significativos, pois as parcelas a e b tem quatro dígitos significativos e a subtração $a - b$ tem apenas dois dígitos significativos.

Obs.: A expressão de Baskara, para a solução exata da equação de segundo grau, muitas vezes, aparece expressa de forma alternativa para minimizar perdas de significação:

Para $ax^2 + bx + c = 0$, tem-se as seguintes raízes:

$$x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \quad (\text{Formula de Baskara})$$

$$x_{1,2} = \frac{-2c}{b \pm \sqrt{b^2 - 4ac}} \quad (\text{Fórmula obtida da racionalização do numerador da expressão anterior})$$

Pode-se observar que as duas formas apresentadas para a solução podem apresentar perdas de significação, quando as parcelas b e $\sqrt{b^2 - 4ac}$ forem de magnitudes próximas e estiverem sujeitas a uma operação de subtração.

Assim, recomenda-se utilizar as expressões propostas acima, escolhendo o sinal do radicando de modo que as parcelas b e $\sqrt{b^2 - 4ac}$ fiquem sujeitas a operação de adição nas duas parcelas, em uma expressão obtendo x_1 e em outra obtendo x_2 .

Exercícios:

5.1). Achar as duas raízes de $x^2 + 62,10x + 1 = 0$, utilizando o operador aritmético $F(10,4,-99,+99)$ e $g = 0$ (quatro dígitos significativos nas operações).

- Use a fórmula de Baskara normal;
- Use a fórmula de Baskara racionalizada;
- Avalie os erros relativos nas duas formas de avaliação das raízes, sabendo que os seus valores exatos são $x_1 = - 0,01610$ e $x_2 = - 62,08$.

(a3). Nas operações de divisão:

Em geral, nas operações de divisão entre duas parcelas quaisquer, normalmente, são gerados resultados com um número de dígitos maior que o permitido na representação em ponto flutuante.

Ex. 59: Efetuar a / b , com $a = 1332$ e $b = 0,9876$, no operador $F(10, 4,-99,+99)$ e $g = 0$.

Efetuando esta operação em uma calculadora de 10 dígitos, tem-se,
 $a / b = 1348,72418$

Porém, se esta operação está sujeita a apenas 4 dígitos significativos ($F10,4,-99,+99$), o resultado será,
 $a / b = 1348$

Causando desta forma uma perda de dígitos significativos.

(b). Instabilidade Numérica:

A acumulação sucessiva de erros de arredondamento pode conduzir um algoritmo de repetição a resultados absurdos, por exemplo,

A avaliação sucessiva de $x = f(x)$ com $f(x) = (N+1)x - 1 \Rightarrow$ é uma constante $\forall x=1/N$ (na ausência de arredondamentos)

Verifique esta afirmação, implementando o exercício 4.2 em computador.

Existem também outras formas de instabilidade, como aquelas associadas ao modelo matemático, por exemplo,

Ex. 60: Dada a função
$$f(x) = \frac{27985}{9,1 - x^2}$$

Avalie $f(x)$ em $x = 3$ e em $x = 3,00001$ utilizando uma calculadora científica com representação de 10 dígitos.

$$\begin{aligned} f(3) &= 279850 \\ f(3.00001) &= 280018,0108 \end{aligned}$$

Note que uma variação de 0,0003333% na variável independente x gera uma variação de 0,06% no resultado final de $f(x)$, ou seja, uma variação no resultado de cerca de 180 vezes superior. Isto caracteriza uma instabilidade intrínseca do modelo matemático em relação aos seus dados de entrada.

Ex. 61: Avaliar $f(x) = 1 - \cos(x)$ para $x = 10^{-4}$ (radianos), na mesma calculadora científica (10 dígitos).

$$f(10^{-4}) = 4,9 \cdot 10^{-9}$$

Se for utilizada uma forma alternativa para o modelo da função:

$$g(x) = (1 - \cos(x)) \cdot \left(\frac{1 + \cos(x)}{1 + \cos(x)} \right) = \left(\frac{\sin^2(x)}{1 + \cos(x)} \right)$$

tem-se então,

$$g(10^{-4}) = 5 \cdot 10^{-5}$$

Note que devido a perda de dígitos significativos por arredondamento, duas representações idênticas da mesma função geram respostas diferentes. Neste caso, também tem-se um exemplo de instabilidade do modelo, porém neste exemplo específico, se tem a possibilidade de reformular a sua representação, gerando uma forma com menor perda de significação ($g(x)$).

Considerações finais:

(i). A avaliação de erros numéricos em máquinas digitais pode ser feita caso se tenha disponibilidade de uma estimativa do valor exato para o resultado desejado.

Desta forma pode-se avaliar o Erro Numérico através das seguintes formas:

- Erro absoluto = | Valor obtido - Valor exato |

- Erro relativo = $\frac{\text{Erro absoluto}}{\text{Valor exato}}$

- Erro percentual = Erro relativo . 100%

A representação dos erros numéricos em forma de Erro relativo (ou percentual) é mais realística (vide exercício 4.3)

(ii). A estimativa de valores exatos, ou mais próximos do exato, para os resultados de algoritmos numéricos, podem ser obtidas das seguinte formas:

- **Estimar o valor mais exato como sendo aquele obtido com o algoritmo numérico operando em dupla precisão, ou seja, reavalia-se o valor obtido numericamente, mas agora com precisão superior, na expectativa de que os resultados sejam mais exatos.**

- Tentar estimar o valor mais exato através de simulação do algoritmo numérico em Sistemas de Computação Algébrica (SCA). Nestes sistemas pode-se recorrer a simulações com precisão ilimitada e alguns casos especiais é possível se proceder a simulação exata do algoritmo.

Obs.: Existem outras alternativas, como por exemplo, proceder simulações numéricas utilizando matemática intervalar, de forma que se possa limitar o erro existente a um intervalo aritmético. Outra possibilidade é o tratamento do erro como uma variável de comportamento estatístico, deste modo pode-se prever os limites do erro, segundo um tratamento estatístico de variáveis.

Exercícios:

5.2). Implemente o algoritmo abaixo em um compilador com aritmética numérica para P1 e P2 tipo real (Pascal, C, Fortran, ...)

```
p1:= 1;
i:= FALSE;
j:= 1;
Repita
  p1:= p1/2;
  p2:= p1 + 1;
  j=j+1;
  SE p2 <= 1 então
    Escreva ('Para minha precisao unitária', p1, ' eh 0 ', 'em ', j, ' repeticoes');
    i:= TRUE;
  fim SE
Até (i=TRUE);
```

Avalie a posição do 1º dígito significativo de P1, que indica a sua precisão relativa ao número 1 (unidade). Verifique que P1 vai diminuindo até que fique tão pequeno que não altera mais o valor de P2. Avalie a influência dos erros de arredondamento no resultados.

5.3) Considerando o compilador Turbo Pascal e suas cinco variáveis do tipo real:

TIPO	FAIXA	DÍGITOS	BYTES
real	2.9e-39...1.7e38	11-12	6
single	1.5e-45..3.4e38	7-8	4
double	5.0e-324...1.7e308	15-16	8
extended	3.4e-4932..1.1e4932	19-20	10
comp	-9.2e18..9.2e18	19-20	8

Execute o seguinte algoritmo e avalie seus resultados em um processador numérico:

Variáveis reais e, f, g, h, x, y;

```
h:= 1/2;
x:= 2/3 - h;
y:= 3/5 - h;
e:= (x + x + x) - h;
f:= (y + y + y + y + y) - h;
```



```

g:= e/f;
Imprima ('h = ',h, ' ',h:5:15); {imprima em notação científica e também na notação decimal}
Imprima ('x = ',x, ' ',x:5:15);
Imprima ('y = ',y, ' ',y:5:15);
Imprima ('e = ',e, ' ',e:5:15);
Imprima ('f = ',f, ' ',f:5:15);
Imprima ('g = ',g, ' ',g:5:15);

```

Experimente trocar o tipo das variáveis, no início do programa, execute o programa novamente e compare os resultados, com os tipos REAL, SINGLE, DOUBLE, EXTENDED E COMP.

5.4). Dadas algumas estimativas do valor exato e de valores aproximados numericamente em um algoritmo, avalie o erro absoluto, relativo e percentual, existente nas seguintes situações:

- Valor aproximado = 1102,345 e Valor exato = 1100,9.
- Valor aproximado = 0,01245 e Valor exato = 0,0119.
- Verifique que o erro absoluto obtido, segundo as várias formas de avaliação, pode não refletir a realidade.

Conclusões:

O perfeito entendimento das causas dos erros numéricos gerados em máquinas digitais é um fator decisivo para que o analista numérico possa fazer a escolha computacionalmente mais eficiente, no momento de programar a resolução de um modelo matemático. Para resolver um determinado modelo o analista numérico deve se preocupar com:

- a escolha de métodos de resolução com menor número de operações aritméticas envolvidas;
- a escolha de um algoritmo com menor número de recursividade;
- a escolha de um compilador, para implementar o algoritmo, que represente as variáveis envolvidas com precisão suficientemente grande;

No final de todo o processo o analista numérico deve ser capaz de:

- obter uma solução de custo mínimo, ou seja, com menor demanda de memória e menor tempo de CPU;
- dimensionar o grau de confiabilidade dos resultados obtidos como solução do modelo matemático.

Pois obter resultados em um programa de computador é muito fácil, qualquer programa compilado gera resultados numéricos, mas daí a se dizer que estes resultados são as respostas para a solução de um modelo matemático, tem-se um longo caminho. É necessário que se faça um minucioso processo de análise numérica dos resultados obtidos, para depois atribuí-los à solução de um modelo matemático.