

# Computação Científica

SEMANA 1

ERROS EM  
REPRESENTAÇÕES  
NUMÉRICAS E

ARITMÉTICA EM  
PONTO FLUTUANTE

PARTE I

ERROS E ERROS DE  
ARREDONDAMENTO

# Computação Científica

prof. Marco Villaça

# EXATIDÃO E PRECISÃO

## INSTRUMENTAÇÃO

---

**Exatidão (acurácia):** Refere-se a quão próximo um valor medido está do valor real.

**Precisão (repetibilidade):** Refere-se a quão próximo um valor medido está de outro valor medido.

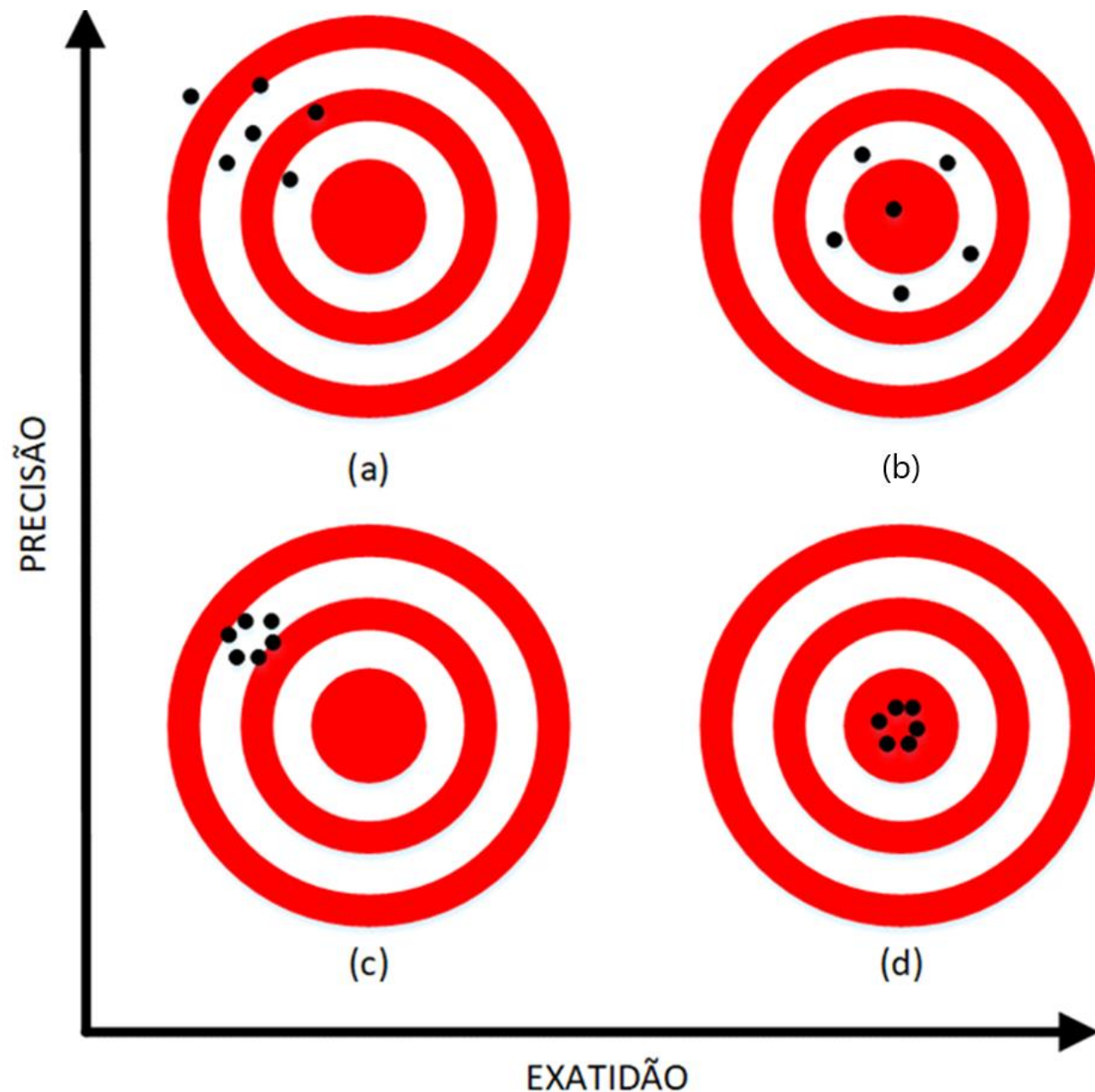
**Inexatidão:** Sistemático desvio do real.

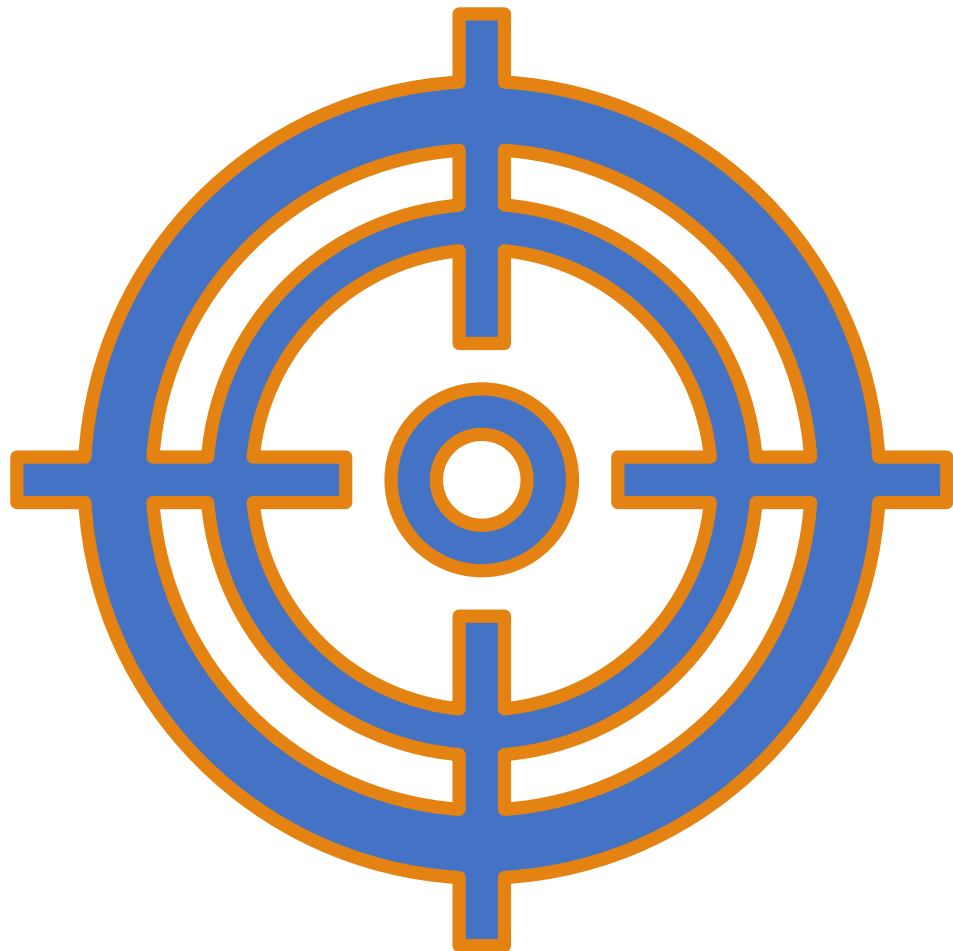
**Imprecisão ou incerteza (dispersão):** magnitude do espalhamento

# EXATIDÃO E PRECISÃO INSTRUMENTAÇÃO

Na figura ao lado, os buracos no alvo correspondem as medidas realizadas e a “mosca” o valor real:

- (a) Inexato e impreciso;
- (b) Exato e impreciso;
- (c) Inexato e preciso
- (d) Exato e Preciso





# EXATIDÃO E PRECISÃO COMPUTAÇÃO CIENTÍFICA

---

**Precisão:** refere-se ao número de dígitos utilizados nas operações aritméticas básicas.

**Exatidão:** refere-se ao erro absoluto ou relativo de uma solução numérica.

# Definições de erro

---

- Erro real ou absoluto

$$E_t = \text{valor real} - \text{valor aproximado}$$

- O problema desta definição de erro é que ela não leva em conta a magnitude do valor em exame.

- Erro relativo real

$$\varepsilon_t = \frac{\text{valor real} - \text{valor aproximado}}{\text{valor real}}$$

- O erro relativo pode também ser multiplicado por 100 para ser expresso percentualmente

$$\varepsilon_t = \frac{\text{valor real} - \text{valor aproximado}}{\text{valor real}} \times 100\%$$

# Definições de erro

---

- A informação do valor real raramente está disponível.
  - Para métodos numéricos, o valor real só será conhecido quando se trabalha com funções que podem ser resolvidas analiticamente
- Quando não se conhece o valor real a priori, a alternativa é normalizar o erro usando a melhor estimativa disponível do valor real:

$$\varepsilon_a = \frac{\textit{erro aproximado}}{\textit{valor aproximado}} \times 100\%$$

ou seja,  $\varepsilon_a$  é o **erro aproximado relativo**.

# Definições de erro

---

- Em processos computacionais iterativos:

$$\varepsilon_a = \frac{\text{valor aproximado atual} - \text{valor aproximado anterior}}{\text{valor aproximado atual}} \times 100\%$$

O erro pode ser positivo ou negativo. Geralmente, se está interessado se o valor absoluto do erro relativo percentual é menor que uma tolerância pré-definida  $\varepsilon_s$ , sendo a computação repetida até:

$$|\varepsilon_a| = \varepsilon_s$$

relação conhecida como *critério de parada*.



# Definições de erro

---

- Se um número é correto para  $n$  algarismos significativos, é evidente que o erro absoluto não pode ser superior a metade de uma unidade no  $n$ -ésimo lugar.
- Por exemplo, se o número 4,629 for correto para quatro algarismos, o erro absoluto não é superior a

$$0,001 \times (1/2) = 0,0005$$

# Definições de erro

---

- Scarborough (1930, p. 7) afirma que se o erro relativo de qualquer número não é maior que  $1/(2 \times 10^n)$ , o número está com certeza correto para  $n$  algarismos significativos.
- Assim, se o critério abaixo for obedecido, assume-se que o resultado alcançado é correto para pelo menos  $n$  algarismos significativos.

$$\varepsilon_s = 0,5 \times 10^{2-n} \%$$

# EXEMPLO 1

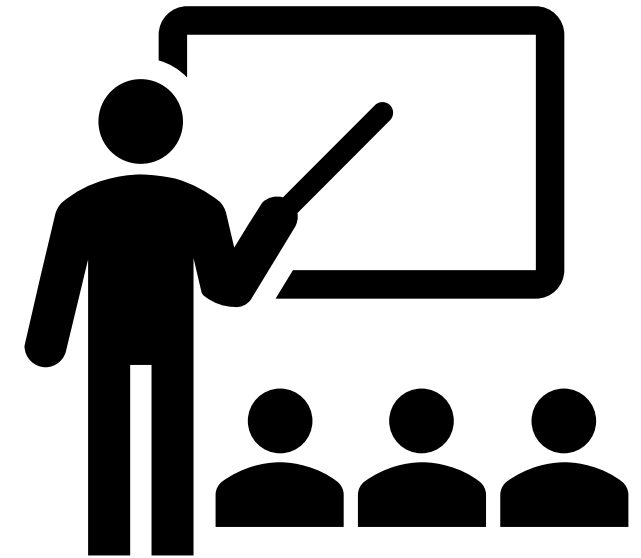
As funções matemáticas frequentemente podem ser representadas por séries de potências (séries de Maclaurin).

É o caso da função  $e^x$ :

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} \dots + \frac{x^n}{n!}$$

Determine  $e^{0,5}$  com pelo menos 3 algarismos significativos de certeza, utilizando a série acima.

Considere o valor real de  $e^{0,5}$  igual a 1,64872127.



# EXEMPLO 1 - SOLUÇÃO

$$e^{0,5} = 1,64872127$$
$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} \dots + \frac{x^n}{n!}$$

- Como foi estabelecido um critério de erro que garanta que o resultado está correto para pelo menos 3 algarismos significativos de resulta:

$$\varepsilon_s = 0,5 \times 10^{2-n} = 0,5 \times 10^{-1} = 0,05 \%$$

- Então, o número de termos calculados deve ser o suficiente para que o  $|\varepsilon_a|$  caia abaixo deste nível.

# EXEMPLO 1 - SOLUÇÃO

$$e^{0,5} = 1,64872127$$
$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} \dots + \frac{x^n}{n!}$$

- Se forem utilizados 3 termos para avaliar  $e^{0,5}$ , resulta:

$$e^x = 1 + 0,5 + \frac{0,5^2}{2} = 1,625$$

- Com 4 termos fica:

$$e^x = 1,625 + \frac{0,6^3}{6} = 1,64583333$$

- O erros aproximado e real relativo serão, respectivamente

$$\varepsilon_a = \frac{1,64583333 - 1,625}{1,64583333} \times 100 = 1,27\%$$

$$\varepsilon_t = \frac{1,64872127 - 1,64583333}{1,64872127} \times 100 = 0,75\%$$

# EXEMPLO 1 - SOLUÇÃO

$$e^{0,5} = 1,64872127$$
$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} \dots + \frac{x^n}{n!}$$

- Pelo critério de parada a computação continuará até:

$$|\varepsilon_a| < 0,05\%$$

- A tabela abaixo resume o processo de computação.

<b>Termo</b>	<b>Resultado</b>	<b><math>\varepsilon_t</math> %</b>	<b><math>\varepsilon_a</math> %</b>
1	1	39,3	-
2	1,5	9,02	33,3
3	1,625	1,44	7,69
4	1,64583333	0,175	1,27
5	1,64843750	0,00172	0,158
6	1,64869791	0,00142	0,0158

## EXEMPLO 2

Escrever uma função Scilab que calcule uma função matemática por séries de potências.

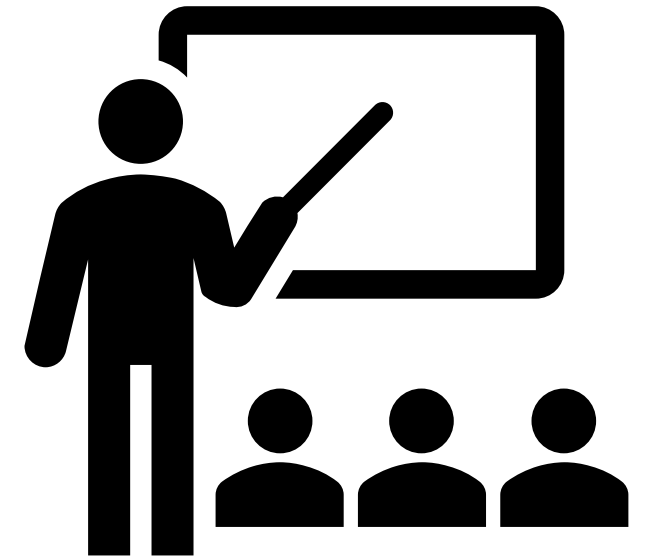
Teste seu código com a função  $\cos x$ , com  $x$  real, utilizando a seguinte série de potências:

$$f(x) = \cos(x) = 1 - \frac{x^2}{2!} + \frac{x^4}{4!} - \frac{x^6}{6!} + \dots + \frac{x^n}{n!}$$

A condição de parada deve garantir pelo menos 4 algarismos significativos.

Lembre que o valor de  $x$  deve ser em radianos.

Durante o processo de cálculo, compute o resultado e os erros relativos percentuais real e aproximado a cada passo.



# EXEMPLO 2

---

```
function [fx, term]=series(fun, x, vreal, n_sig)
// Cálculo de funções com serie de potências
// function [fx,term]=series(funcao, x, n_sig)
// onde fx é o valor da função e x
//   term é o número de termos usados para o erro especificado
//   fun é a função de entrada literal em x e n
//   x é o número np qual se deseja obter o valor da função
//   vreal é o valor real da função em x
//   n é o numero do termo
//   n_sig é o numero de alg. significativos da resposta - opcional
// Exemplo de chamada:
// exec('path\series.sci',-1)
// fun = '(-1)^(n)*x^(2*n)/factorial(2*n)'
// [fx,term]=series(fun, %pi/6)
```



# EXEMPLO 2

---

```
if argn(2) < 4 then
```

```
    n_sig = 3;
```

```
end
```

```
es = 0.5*10^(2-n_sig); // condição de parada relativa %
```

```
fx_old = 0; fx = 0; n = 0;
```

```
printf("Termo fx      et %%      erro %%\n");
```

# EXEMPLO 2

---

```
while 1 do
    fx = fx + evstr(fun);
    erro = abs((fx - fx_old)/fx)*100;
    errot = abs((vreal - fx)/vreal)*100;
    printf("%-5d %-10f %-10f %-10f\n",n+1,fx,errot,erro);
    if(erro <= es) then
        break;
    end
    fx_old = fx;    n = n + 1;
end
```

O método “divisão e média”, um antigo método para a aproximação da raiz quadrada de um número positivo  $a$ , pode ser formulado por:

$$x_{i+1} = \frac{x_i + a/x_i}{2}$$

Com  $i = 0, 1, 2, \dots, n$  e condição de parada que garanta 4 algarismos significativos, escreva um script Scilab para implementar este método, onde:

Considere o valor inicial estimado para a raiz de  $a$  igual a 1.

## Exercício 1



# ERROS DE ARREDONDAMENTO

---

- Números como o  $\pi$  e  $\sqrt{2}$  não podem ser expressos por um número fixo de algarismos significativos. Logo eles não podem ser representados com exatidão por computadores, pois os mesmos apresentam uma palavra de tamanho fixo.
- A discrepância introduzida pela omissão de algarismos significativos é chamada de erro de arredondamento.

# ERROS DE ARREDONDAMENTO

---

- Se  $\pi = 3,141592653\dots$  deve ser armazenado em um sistema decimal com apenas 7 algarismos significativos, surgem duas possibilidades:

- **Chopping** (corte):

$$\pi = 3,141592 \rightarrow \text{erro: } E_t = 0,00000065$$

- **Rounding** (arredondamento):

$$\pi = 3,141593 \rightarrow \text{erro: } E_t = -0,00000035$$

- Os números seguintes são arredondados aplicando o critério mais conhecido de arredondamento:

$$3,6543 = 3,65; \quad 0,49781 = 0,498$$

$$22,65 = 22,6; \quad 1,735 = 1,74$$

# Erros de arredondamento

## Representação em ponto flutuante

---

- Similar a notação científica, na representação em ponto flutuante um número real é expresso por:

$$\pm s \cdot b^e$$

onde  $s$  = significando

$b$  = base do sistema numérico

$e$  = expoente

- Uma forma normalizada elimina os zeros imediatamente a direita do ponto decimal :

$$0.0456 = 0,456 \times 10^{-1}$$

# Erros de arredondamento

## Implicações da representação em ponto flutuante

---

- Suponha um computador hipotético de base 10 e palavra de 5 dígitos, onde 1 dígito é reservado para o sinal, dois para o expoente e dois para o significando. Um dígito do expoente é reservado para o seu sinal:
- A representação do número será
- $s_1 d_1 . d_2 \times 10^{s_2 d_3}$   
onde
  - $s_1$  e  $s_2$  são sinais
  - $d_1$  e  $d_2$  a magnitude do significando (sendo  $d_1 \neq 0$ )
  - $d_3$  a magnitude do expoente

# Erros de arredondamento

## Implicações da representação em ponto flutuante

---

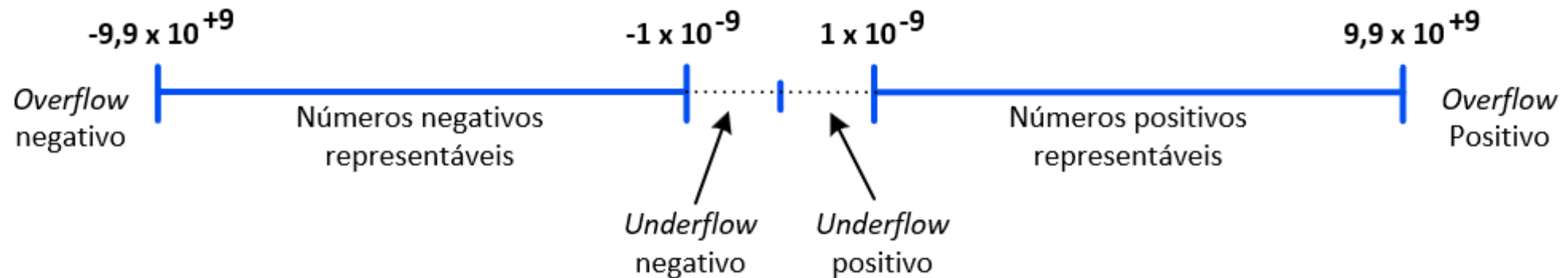
- Nesse computador:
  - O maior valor absoluto representado é
$$9,9 \times 10^9$$
  - O menor valor absoluto representado é
$$1,0 \times 10^{-9}$$



# Erros de arredondamento

## Implicações da representação em ponto flutuante

- O Diagrama abaixo mostra todos os possíveis números representados pelo computador hipotético.



# Erros de arredondamento

## Implicações da representação em ponto flutuante

---

- Números positivos ou negativos muito grandes saem da faixa de representação causando um erro por *overflow*.
- Os números muito pequenos, que não podem ser representados, usualmente são convertidos em zero.

# Erros de arredondamento

## Implicações da representação em ponto flutuante

---

- A utilização de apenas dois dígitos para o significando afeta dramaticamente a exatidão:
  - O número 0,01845 seria armazenado como  $0,18 \times 10^{-1}$  gerando um erro relativo de

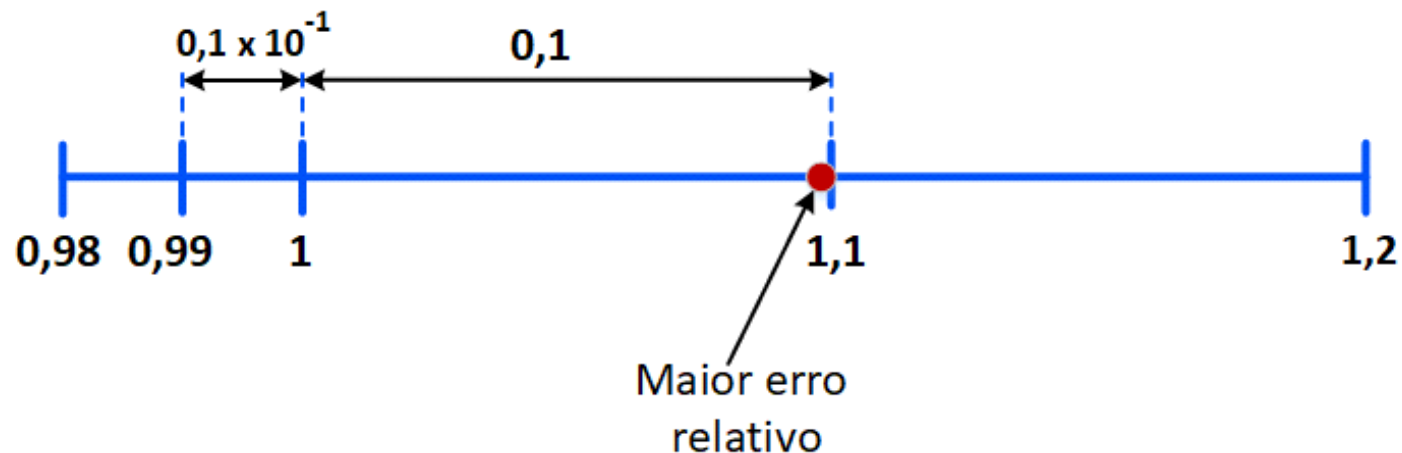
$$e_t = \frac{0,01845 - 0,018}{0,01845} \times 100 = 2,44 \%$$

- Embora seja possível armazenar exatamente um número como 0,01845, expandindo-se os dígitos do significando, quantidades com infinitos dígitos sempre devem ser aproximadas, como no caso do  $\pi$ .
- Logicamente quanto maior o número de dígitos do significando, melhor será essa aproximação. Entretanto, tais números **sempre** apresentarão um erro de arredondamento.

# Erros de arredondamento

## Implicações da representação em ponto flutuante

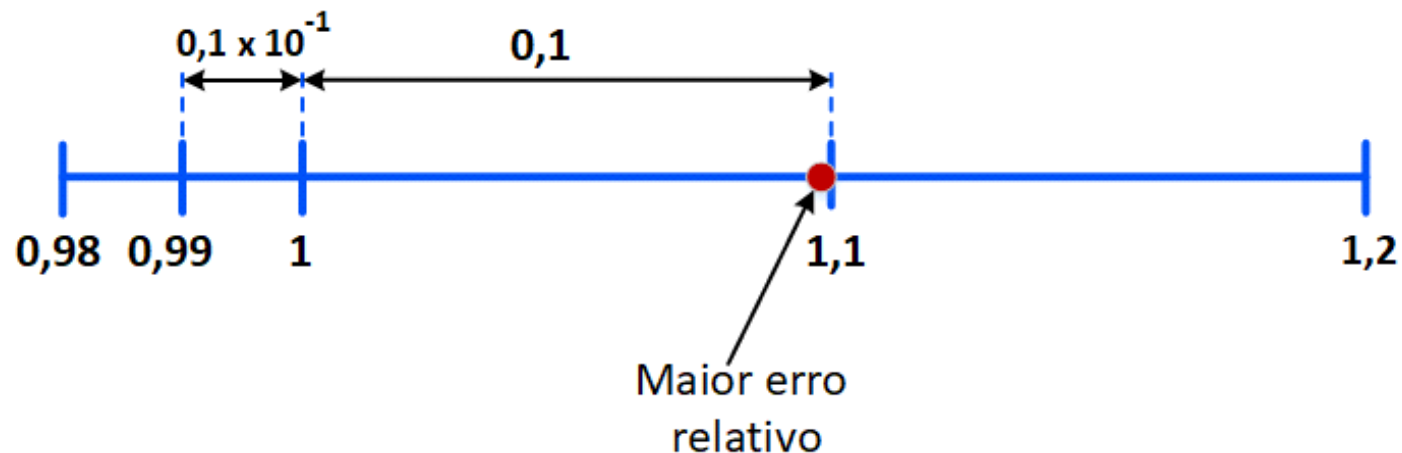
- Outro efeito mais sutil da representação em ponto flutuante é mostrado pelo diagrama abaixo.
- Observa-se que o intervalo entre números incrementa quando se move entre faixas de magnitude:
  - Entre 0,1 e 1, o intervalo é de  $0,01 = 0,1 \times 10^{-1}$
  - Entre 1 e 10, o intervalo é de  $0,1 = 0,1 \times 10^0$



# Erros de arredondamento

## Implicações da representação em ponto flutuante

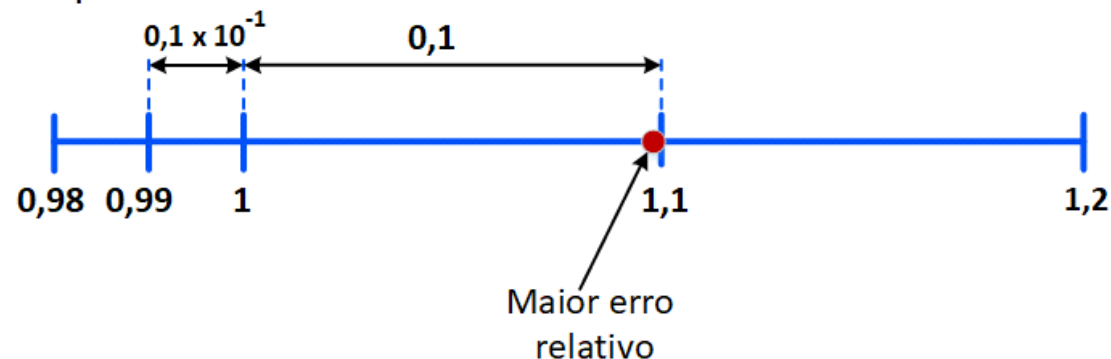
- Isso permite que a representação em ponto flutuante preserve o número de algarismos significativos.
- Em termos de erro de arredondamento, significa que dentro de uma faixa, o erro absoluto será limitado ao tamanho do intervalo.



# Erros de arredondamento

## Implicações da representação em ponto flutuante

- O maior erro relativo, por sua vez, ocorrerá para os números situados imediatamente abaixo do limite superior do primeiro de uma série de intervalos igualmente espaçados. Utilizando-se o corte, erro relativo será, portanto, limitado pela diferença entre 1 e o próximo número que pode ser representado.
- Este valor é chamado de *épsilon da máquina* ( $\epsilon$ ) ou precisão da máquina, que em nosso computador hipotético vale 0,1.
- Foi visto que para o número 0,01845 o erro relativo vale 0,024, um número menor que  $\epsilon$ .



ERROS EM  
REPRESENTAÇÕES  
NUMÉRICAS E

ARITMÉTICA EM  
PONTO FLUTUANTE

PARTE I

ERROS E ERROS DE  
ARREDONDAMENTO

(CONTINUAÇÃO)

# Computação Científica

prof. Marco Villaça

# IEEE754

---

- É o padrão de **base 2** utilizado atualmente para representar números em ponto flutuante em computadores :
  - Na forma normalizada, o bit à esquerda do ponto sempre será 1, ou seja, ele é implícito, não precisando ser armazenado.
  - Assim, qualquer número, diferente de zero será representado na forma normalizada por

$$(-1)^s \times (1.f) \times 2^E$$

onde  $s$  é o bit de sinal,  $f$  é a mantissa (parte fracionária do significando) e  $E$  o expoente representado em “excesso de”.

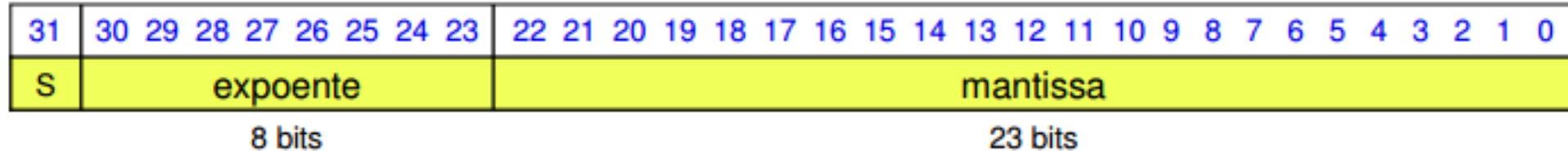


# IEEE754

## Precisão simples

---

- O padrão IEEE 754 de precisão simples:



- Bit de sinal (S): 1bit
  - $S = 0 \rightarrow$  número positivo
  - $S = 1 \rightarrow$  número negativo
- Expoente: 8 bits
- Significando: 24 bits (23 armazenados explicitamente)

# IEEE754

## Precisão simples

---

- Expoente  $E$  em excesso de 127 ( $2^{(n-1)} - 1$ ):  
$$E = \text{expoente real} + 127$$
- Para representar números, utiliza-se  $E$  entre 1 e 254, resultando em expoentes reais entre -126 e +127.

# IEEE754

## Precisão simples

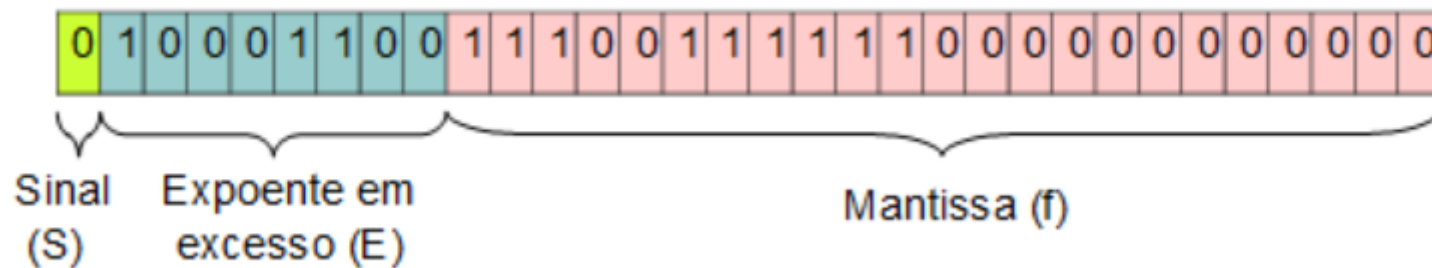
- O número decimal 15612 é representado por:

$$15612 = 11110011111100_2 = 1,1110011111100_2 \times 2^{13}$$

Mantissa (f) = 111001111110000000000000

Expoente (E) =  $(13 + 127)_{10} = 10001100$

Sinal = 0



# IEEE754

## Precisão simples

---

- Os expoentes 0 (todos os bits em 0) e 255 (todos os bits em 1) são reservados para denotar os valores especiais:
    - Zero:  $E = 0$ ,  $S = 0$  e  $f = 0$
    - Números (subnormais):  $E = 0, f \neq 0$ .
    - $+\infty$  :  $E = 255$ ,  $S = 0$  e  $f = 0$
    - $-\infty$  :  $E = 255$ ,  $S = 1$  e  $f = 0$
    - NaN (Not a Number):  $E = 255$ ,  $S = 0$  e  $f \neq 0$
- Resultados de operações indeterminadas (  $0/0$ ,  $\sqrt{-1}$ ,  $\infty - \infty$ , etc.) ou exceções.

# IEEE754

## Precisão simples

---

- Números desnormalizados ou subnormais:
  - Possuem um bit 0 implícito à esquerda do ponto binário, permitindo a representação de números menores do que os possíveis na forma normalizada com redução da perda de significância quando ocorre um underflow.
  - Qualquer número diferente de zero será representado na forma desnormalizada por

$$(-1)^s \times (0.f) \times 2^{-126}$$

- Perdem gradualmente a precisão à medida que diminuem (bits da esquerda da mantissa se tornam zero).

# IEEE754

## Precisão simples

---

- No padrão IEEE 754 de precisão simples:
  - Como a mantissa tem 23 bits, o *epsilon* ou precisão da máquina é:

$$\epsilon = 2^{-23} = 1,192 \times 10^{-7}$$

- O maior número normalizado em valor absoluto que pode ser representado é:

$$\text{Maior valor} = 1,111111 \dots \times 2^{127} = (2 - \epsilon) \times 2^{127} = 3,403 \times 10^{38}$$

- O menor número normalizado em valor absoluto :

$$\text{Menor valor} = 1,00000 \dots \times 2^{-126} = 1,175 \times 10^{-38}$$

# IEEE754

## Precisão simples

---

- No padrão IEEE 754 de precisão simples:
  - O menor número subnormal em valor absoluto :  
Menor valor  $= \varepsilon \times 2^{-126} = 2^{-23} \times 2^{-126} = 1,401 \times 10^{-45}$

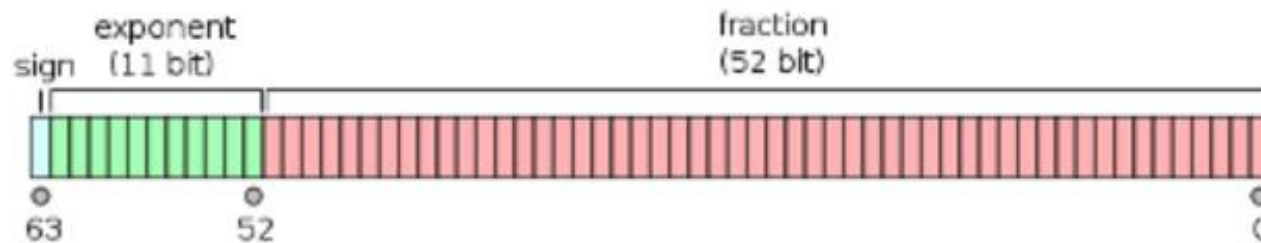


# IEEE754

## Dupla precisão

---

- O padrão IEEE 754 de dupla precisão:



- Utiliza 52 bits para a mantissa (M);
- Expoente E em excesso de 1023:  
$$E = \text{expoente real} + 1023$$
- Para representar números, utiliza-se E entre 1 e 2046, resultando em expoentes reais entre -1022 e 1023.
- Casos especiais com  $E = 0$  e  $E = 2047$



# IEEE754

## Dupla precisão

---

- No padrão IEEE 754 de precisão dupla:
  - Como a mantissa tem 52 bits, o *epsilon* ou precisão da máquina é:

$$\epsilon = 2^{-52} = 2,22 \times 10^{-16}$$

- O maior número que pode ser representado é:

$$\begin{aligned} \text{Maior valor} &= 1,111111 \dots \times 2^{1023} = (2 - \epsilon) \times 2^{1023} = \\ &+1,80 \times 10^{308} \end{aligned}$$

- O menor número é:

$$\text{Menor valor} = 1,00000 \dots \times 2^{-1022} = 2,23 \times 10^{-308}$$

# IEEE754

## Dupla precisão

---

- No padrão IEEE 754 de precisão dupla:
  - O menor número subnormal em valor absoluto :  
Menor valor =  $\epsilon \times 2^{-1022} = 2^{-52} \times 2^{-1022} = 4,94 \times 10^{-324}$

# IEEE754

## Algarismos significativos

---

- No padrão IEEE 754 de precisão simples:
  - Há 23 bits denotando a mantissa mais o bit implícito, totalizando 24 bits significativos. Em dígitos decimais, tem-se aproximadamente:

$$\text{Número de dígitos} = \frac{24 \cdot \log(2)}{\log(10)} \approx 7$$

- Similarmente, no padrão IEEE 754 de dupla precisão:

$$\text{Número de dígitos} = \frac{53 \cdot \log(2)}{\log(10)} \approx 15$$

- Em muitos casos, o uso da precisão dupla pode atenuar os efeitos dos erros de arredondamento, às custas de consumo de memória e tempo de execução, consumo que pode ser insignificante para pequenos programas, mas considerável para grandes aplicações.

# IEEE754

## Precisão estendida

---

- Conforme foi visto, O padrão IEEE 754 especifica dois formatos básicos de ponto flutuante: precisão simples e dupla.
- O formato de precisão simples tem um significando com precisão de 24 bits e ocupa 32 bits no total.
- O formato de dupla precisão tem um significando com precisão de 53 bits e ocupa 64 bits.
- Especifica, também, duas classes de formatos de ponto flutuante estendidos: precisão simples estendida e dupla estendida. O padrão não prescreve a precisão e o tamanho exatos desses formatos, mas especifica a precisão e o tamanho mínimos. Por exemplo, um formato IEEE precisão estendida dupla deve ter um significando com precisão de pelo menos 64 bits e ocupar pelo menos 79 bits no total.

	Precisão Simples	Dupla Precisão	Precisão estendida
Número de bits	32	64	80
Extensão do expoente em bits	8	11	15
Offset do expoente	127	1023	16383
Extensão da mantissa em bits	23	52	63
Precisão em bits	24	53	64
Precisão em dígitos decimais	7	16	19
Épsilon da máquina	$1,192 \times 10^{-7}$	$2,22 \times 10^{-16}$	$1,08 \times 10^{-19}$
Maior número em valor absoluto	$3,40 \times 10^{38}$	$1,80 \times 10^{308}$	$1,19 \times 10^{4932}$
Menor número normalizado	$1,18 \times 10^{-38}$	$2,23 \times 10^{-308}$	$3,36 \times 10^{-4932}$
Menor número subnormal	$1,401 \times 10^{-45}$	$4,94 \times 10^{-324}$	$3,65 \times 10^{-4951}$

Representações  
padrão IEEE 754  
incluindo a  
representação de  
precisão estendida  
de 80 bits x86

Para computadores, o épsilon da máquina pode ser pensado como o menor número que quando adicionado a 1 resulta um número maior que 1. Um algoritmo baseado nesta ideia pode ser desenvolvido da seguinte maneira:

- **Passo 1:** Faça  $\varepsilon = 1$ ;
- **Passo 2:** Se  $(1 + \varepsilon) \leq 1$ , então vá para o passo 5
- **Passo 3:** Faça  $\varepsilon = \varepsilon / 2$
- **Passo 4:** Retorne ao passo 2
- **Passo 5:** Imprima  $2 * \varepsilon$

Escreva um script Scilab para determinar o épsilon da máquina e valide o resultado comparando-o valor calculado com o valor da constante %eps.

## Exercício 2



# Aritmética de ponto flutuante

## Adição/subtração

---

- Devido a necessidade de igualar os expoentes para alinhar o ponto decimal, a adição e a subtração são especialmente suscetíveis a perda potencial de algarismos significativos. A adição/subtração de ponto flutuante segue 4 etapas:
  - Alinhamento da mantissa se os expoentes dos números forem diferentes.
  - Adição / subtração das mantissas alinhadas.
  - Normalização do resultado se não normalizado.
  - Arredondamento do resultado

# Aritmética de ponto flutuante

## Adição/subtração

- Para simplificar a discussão, suponha que se deseja adicionar 2 números reais usando um computador decimal com uma mantissa normalizada de 4 dígitos ( $0.d_1d_2d_3d_4 \times 10^{\text{exp}}$ )
- Para adicionar  $0,2556 \times 10^1$  e  $0,4481 \times 10^{-1}$ , o computador ajusta o menor número para que seu expoente coincida com o expoente de maior número:

	0	,	2	5	5	6			×	1	0	1
+	0	,	0	0	4	4	8	1	×	1	0	1
	0	,	2	6	0	0	8	1	×	1	0	1

- O resultado da adição é  $0,2600 \times 10^1$  (corte) ou  $0,2601 \times 10^1$  (arredondamento).



# Aritmética de ponto flutuante

## Adição/subtração

---

- Observa-se que para ajustar o expoente do menor número, seus dois últimos dígitos foram deslocados para a direita.
- Isso é permitido porque a mantissa de um número pode ser estendida pela adição de guard bits para preservar a exatidão durante o arredondamento.
- Assim, os guard bits permitem que um número de ponto flutuante possa ser expandido durante as operações internas antes de ser armazenado em 4 bits.
- Em resumo, no exemplo, a falta de precisão (o número de dígitos ou bits mantidos ao final de um cálculo) afeta a exatidão da computação, pois o resultado foi truncado para 4 algarismos significativos

# Aritmética de ponto flutuante

## Adição/subtração

---

- Para subtrair 36,36 de 26,41, se deve normalizar os operandos:

	0	,	3	6	3	6	×	1	0	<sup>2</sup>
-	0	,	2	6	4	1	×	1	0	<sup>2</sup>
	0	,	0	9	9	5	×	1	0	<sup>2</sup>

- Normalizando o resultado, obtém-se  $0,9950 \times 10^1$ . Observe que zero colocado no fim da mantissa não é significativo, ele serve apenas para preenchê-la.

# Aritmética de ponto flutuante

## Regra 1

---

- Esse resultado permite formular uma primeira regra importante da aritmética de ponto flutuante:

*Sempre que se subtrair dois números de mesmo sinal ou adicionar dois números com sinais diferentes, a precisão do resultado pode ser menor que a precisão disponibilizada pelo formato de ponto flutuante.*

# Aritmética de ponto flutuante

## Multiplicação e divisão

---

- Na multiplicação, somam-se os expoentes e multiplicam-se as mantissas.
- Duas mantissas com  $n$  dígitos produzem um resultado com  $2n$  dígitos. Assim, o resultado intermediário deve ser armazenado em um registrador com o dobro de extensão. Por exemplo:

$$\begin{array}{r} 0,3641 \times 10^3 \\ \times \quad 0,2686 \times 10^{-1} \\ \hline 0,09779726 \times 10^2 \end{array}$$

- Normalizando o resultado:  $0,97797260 \times 10^1$
- Cortando:  $0,9779 \times 10^1$
- O procedimento de divisão é similar, porém as mantissas são divididas e os expoentes subtraídos.

# Aritmética de ponto flutuante

## Regra 2

---

- Por si só, multiplicação e divisão não produzem resultados ruins. No entanto, eles tendem a incrementar qualquer erro que já exista em um valor.
- Por exemplo, multiplicando-se 0,453 por dois, quando se deveria estar multiplicando 0,454 por dois, o resultado é ainda menos preciso. Isso traz uma segunda regra importante quando se trabalha com aritmética de ponto flutuante:

*Ao realizar uma cadeia de cálculos envolvendo adição, subtração, multiplicação e divisão, deve se executar primeiro as operações de multiplicação e divisão.*

# Aritmética de ponto flutuante

## Regra 3

---

- Entretanto, a multiplicação e divisão também apresentam problemas.
- Ao multiplicar dois números muito grandes ou ao dividir um número grande por um número pequeno pode ocorrer um overflow.
- Por outro lado, a multiplicação de dois números pequenos e a divisão de um número pequeno por um número grande podem levar a um underflow. Isso nos leva a formular uma terceira regra que a ser levada em consideração:

*Procurar organizar os cálculos de forma que os números grandes sejam multiplicados por números pequenos e que a divisão ocorra entre números da mesma ordem de grandeza.*

# Aritmética de ponto flutuante

## Grandes computações

---

- A perda de precisão durante um único cálculo geralmente não é de grande preocupação, a menos que se esteja muito preocupado com a exatidão dos cálculos.
- No entanto, quando for calculado um valor que é o resultado de uma sequência de operações de ponto flutuante, mesmo que o erro de arredondamento de cada operação seja pequeno, o efeito cumulativo no curso de toda a computação pode ser significativo.

# Aritmética de ponto flutuante

## Grandes computações

---

- Seja o programa em linguagem C abaixo:

```
#include <stdio.h>
int main()
{
    float s = 0;
    int i;
    for (i=1;i<=10000;i++)
        s= s+0.0001;
    printf("Resultado = %f",s);
}
```

que ao ser executado, resulta em:

Resultado = 1.00054



# Aritmética de ponto flutuante

## Grandes computações

---

- Poderia ser esperado o resultado igual 1. Entretanto, embora 0,0001 seja um número redondo na base 10, ele não é redondo na base 2:

$$0,0001_{10} = (0.0000000000000011010001101000 \dots)_2$$

- Assim, usando a representação IEEE 754 de precisão simples (`float` em C), as somas sucessivas levam a um resultado ligeiramente diferente de 1.

# Aritmética de ponto flutuante

## Adicionando um número pequeno e um grande

---

- Suponha que, utilizando nosso computador hipotético de 4 dígitos, adiciona-se 0,0010 a 4000:

$$\begin{array}{r} 0,4000 \quad \times 10^4 \\ + 0,0000001 \times 10^4 \\ \hline 0,4000001 \times 10^4 \end{array}$$

- Efetuando o corte, o resultado final é  $0,4000 \times 10^4$ , parecendo que a adição não foi realizada!
- Este tipo de erro pode ocorrer na computação de séries infinitas, onde os termos iniciais são frequentemente grandes quando comparados aos termos posteriores.

Uma maneira de minimizar esse tipo de erro é somar os termos da série em ordem inversa, do menor para o maior.

- Use aritmética de 3 dígitos com truncamento para calcular a soma de

$$\sum_{i=1}^{10} \frac{1}{i^2}$$

primeiro com

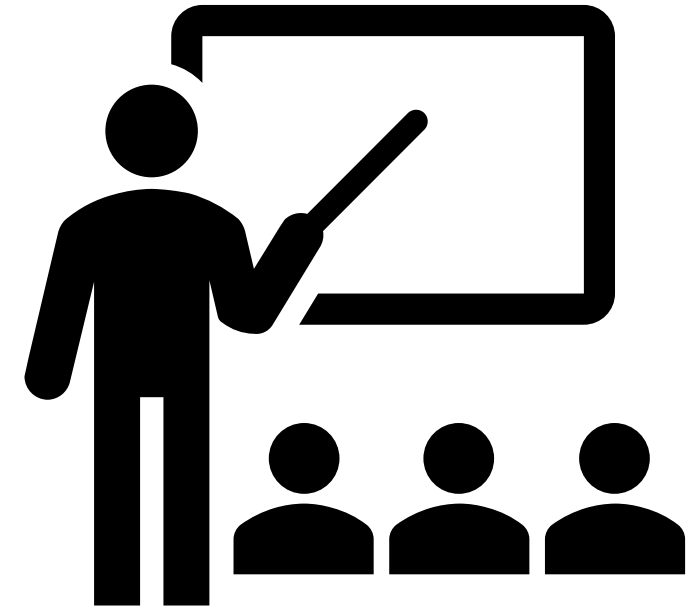
$$\frac{1}{1} + \frac{1}{4} + \dots + \frac{1}{100}$$

e depois com

$$\frac{1}{100} + \frac{1}{81} + \dots + \frac{1}{1}$$

Qual método é mais preciso. Por quê?

- Obtenha uma aproximação melhor com o auxílio de um computador (Resposta: 1,5497677).



## EXEMPLO 3

0,	1	0	0	x	1	0	1	1,000
+ 0,	0	2	5	x	1	0	1	0,250
= 0,	1	2	5	x	1	0	1	
+ 0,	0	1	1	x	1	0	1	0,111
= 0,	1	3	6	x	1	0	1	
+ 0,	0	0	6	x	1	0	1	0,0625
= 0,	1	4	2	x	1	0	1	
+ 0,	0	0	4	x	1	0	1	0,040
= 0,	1	4	6	x	1	0	1	
+ 0,	0	0	2	x	1	0	1	0,0277
= 0,	1	4	8	x	1	0	1	
+ 0,	0	0	2	x	1	0	1	0,0204
= 0,	1	5	0	x	1	0	1	
+ 0,	0	0	1	x	1	0	1	0,0156
= 0,	1	5	1	x	1	0	1	
+ 0,	0	0	1	x	1	0	1	0,0123
= 0,	1	5	2	x	1	0	1	
+ 0,	0	0	1	x	1	0	1	0,010
= 0,	1	5	3	x	1	0	1	<b>1,53</b>

## EXEMPLO 3

Soma dos termos da série do maior para menor

---

$$\sum_{i=1}^{10} \frac{1}{i^2} = \frac{1}{1^2} + \frac{1}{2^2} + \frac{1}{3^2} + \dots$$

- Erro verdadeiro:

$$\varepsilon_a = \frac{1,5497677 - 1,53}{1,5497677} \times 100 = 1,27 \%$$

0,	1	0	0	x	1	0	-1	0,01
+ 0,	1	2	3	x	1	0	-1	0,0123
= 0,	2	2	3	x	1	0	-1	
+ 0,	1	5	6	x	1	0	-1	0,0156
= 0,	3	7	9	x	1	0	-1	
+ 0,	2	0	4	x	1	0	-1	0,0204
= 0,	5	8	3	x	1	0	-1	
+ 0,	2	7	7	x	1	0	-1	0,0277
= 0,	8	6	0	x	1	0	-1	
+ 0,	4	0	0	x	1	0	-1	0,04
= 0,	1	2	6	x	1	0	0	1,260 x 10 <sup>-1</sup>
+ 0,	0	6	2	x	1	0	0	0,0625
= 0,	1	8	8	x	1	0	0	
+ 0,	1	1	1	x	1	0	0	0,111
= 0,	2	9	9	x	1	0	0	
+ 0,	2	5	0	x	1	0	0	0,25
= 0,	0	5	4	x	1	0	1	0,549 x 10 <sup>0</sup>
+ 0,	1	0	0	x	1	0	1	1,000
= 0,	1	5	4	x	1	0	1	1,54

Ajuste

Ajuste

## EXEMPLO 3

Soma dos termos da série do menor para o maior

$$\sum_{i=10}^1 \frac{1}{i^2} = \frac{1}{10^2} + \frac{1}{9^2} + \frac{1}{8^2} + \dots$$

- Erro verdadeiro:

$$\epsilon_a = \frac{1,5497677 - 1,54}{1,5497677} \times 100 = 0,63\%$$

# Aritmética de ponto flutuante

## Regra 4

---

- Um fato muito importante que deve ser conhecido quando se trabalha com aritmética de ponto flutuante é:  
*A ordem da avaliação pode afetar a exatidão do resultado.*

- A série infinita

$$f(n) = \sum_{i=1}^n \frac{1}{i^4}$$

converge para o valor de  $f(n) = \pi^4/90$  quando  $n$  se aproxima do infinito.

- Escreva um programa em C **em precisão simples** para calcular  $f(n)$  para  $n = 10000$  computando a soma de  $i = 1$  a  $10000$  usando incrementos de 1
- Repita o cálculo mas na ordem inversa, isto é de  $i = 10000$  a  $1$ , usando incrementos de  $-1$ .
- Em cada caso estime o erro relativo percentual real. Explique os resultados.

Resposta: Crescente –  $f(n) = 1,082322$ ,  $\varepsilon_t = 0,000103$  %

Decrescente –  $f(n) = 1,082323$ ,  $\varepsilon_t = 0,000004$  %

## Exercício 3



Seja o seguinte padrão de representação de um número flutuante de um computador de 16 bits:

**Bit de sinal (S) :** 1 bit

- S = 0, número positivo
- S = 1, número negativo

**Expoente:** 4 bits, em excesso de 7:

- expoente = expoente real + 7
- Para representar números utiliza-se o expoente entre 1 e 14. resultando em expoentes reais de -6 a 7

**Mantissa:** 12 bits (11 armazenados explicitamente)

**Representação:**  $(-1)^{\text{sinal}} \times 1.\text{mantissa} \times 2^{\text{expoente}}$

## Exercício 4





Pede-se

- a) O  $\epsilon$  da máquina
- b) O menor número que pode ser representado em valor absoluto;
- c) O maior número que pode ser representado em valor absoluto;
- d) O número de algarismos significativos da representação

## Exercício 4



# Bibliografia e crédito das figuras

---



CHAPRA, Steven. **Applied numerical methods with MATLAB for engineers and scientists.** McGrawHill, 2012.



CHAPRA, Steven e CANALE, Raymond. **Numerical methods for engineers.** McGrawHill, 2010.



SCARBOROUGH, James. **Numerical Mathematical Analysis.** London: Oxford Press, 1930.