

Questão: Como armazenar números em caixas?

Representação finita => erros de arredondamento

$$423_{10} = 4 \cdot 100 + 2 \cdot 10 + 3 \cdot 1 = 4 \cdot 10^2 + 2 \cdot 10^1 + 3 \cdot 10^0$$

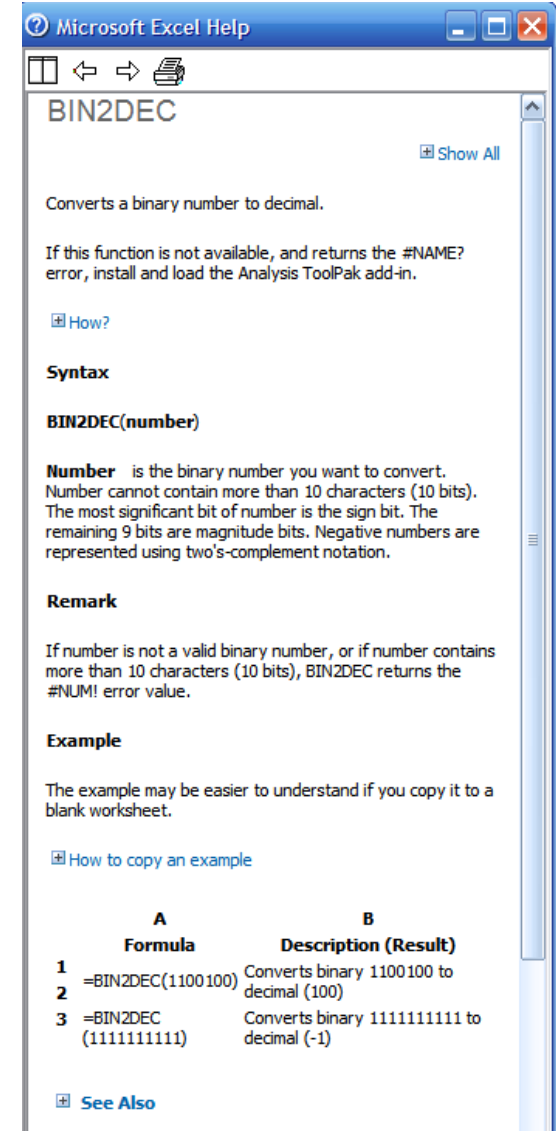
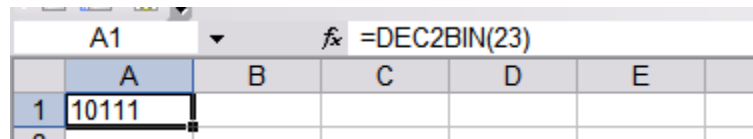
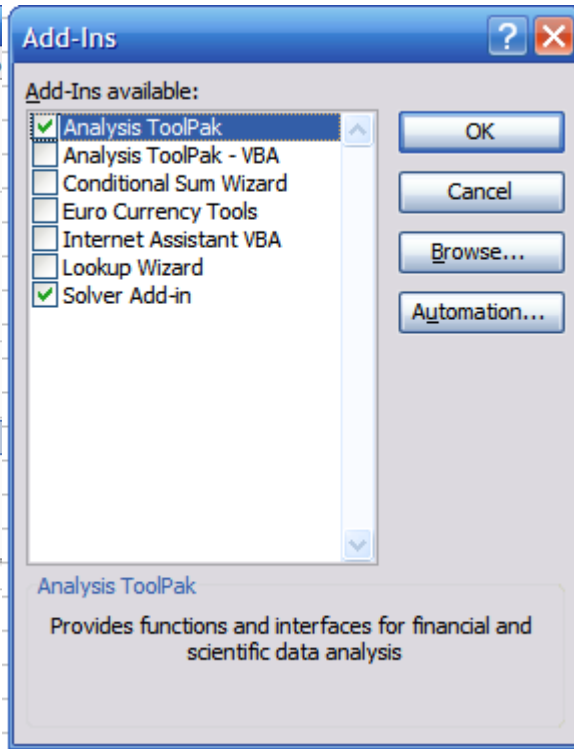
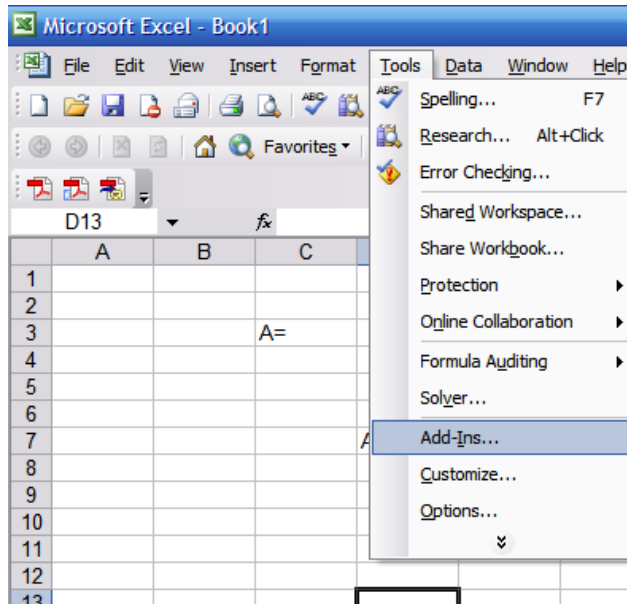
$$35.23_{10} = 3 \cdot 10 + 5 \cdot 1 + 2 \cdot 0.1 + 3 \cdot 0.01 = 3 \cdot 10^1 + 5 \cdot 10^0 + 2 \cdot 10^{-1} + 3 \cdot 10^{-2}$$

```
>> dec2bin(87)
```

```
ans= 1010111
```

$$87 = 1 \times 2^6 + 1 \times 2^4 + 1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0$$

No Excel há que introduzir o add-in para incluir a instrução de conversão



Sistemas de codificação

Decimal	Hexadecimal	Binário
0	0	0000
1	1	0001
2	2	0010
3	3	0011
4	4	0100
5	5	0101
6	6	0110
7	7	0111
8	8	1000
9	9	1001
10	A	1010
11	B	1011
12	C	1100
13	D	1101
14	E	1110
15	F	1111

Problema de Arredondamentos

Exemplo simples: aplicar uma sucessão de operações e depois a sucessão inversa de operações

Aplicar a raiz quadrada a um número, N vezes. Depois aplicar N vezes o quadrado. Conseguirei recuperar o número inicial?

Fazer em excel...

Exemplo 2: funções

Fazer no EXCEL:

`sin(PI())`

Quanto dará?

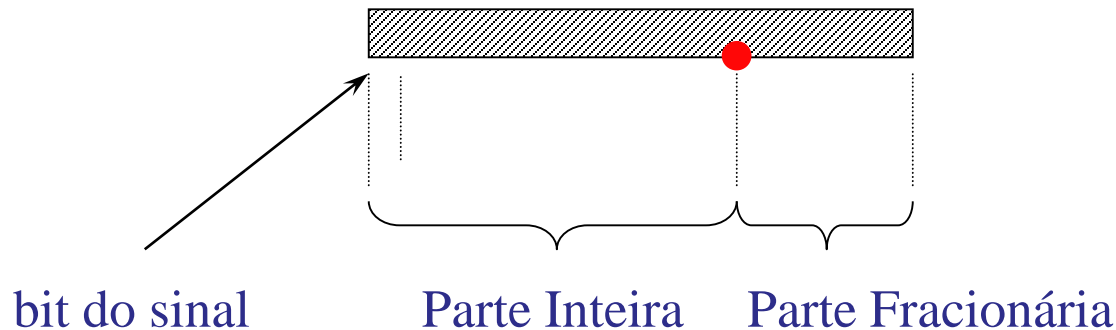
Fazer no EXCEL:

`cos(PI())`

Quanto dará?

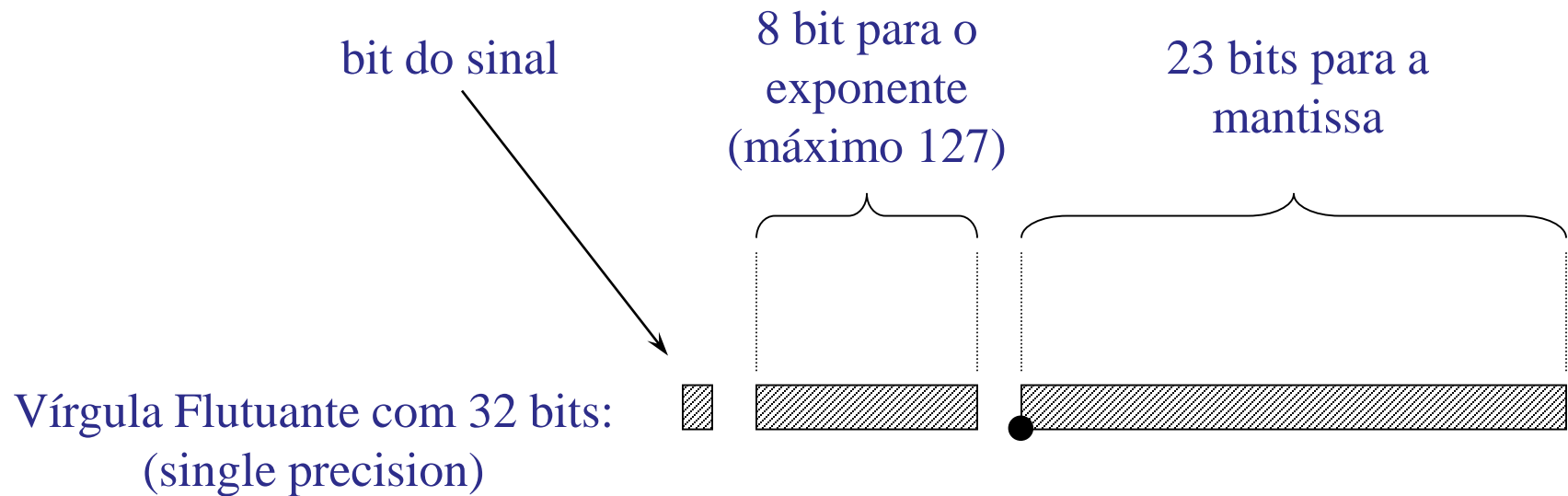
E NO MATLAB DARÁ O MESMO?

Representação em Vírgula Fixa



Exemplo: 51.23, 0.23, 0.99 estão todos representados com 2 dígitos na parte fracionária.

Representação em Vírgula Flutuante



$$x = (-1)^s \times b^e \times (0.a_1a_2...a_t)$$

Representação em vírgula flutuante normalizada (**binário**):

$$x = (1 + f) \times 2^e \quad 0 \leq f < 1$$

Os números em dupla precisão dedicam 52 bits para a fracção e 11 bits para o expoente.

Variação mínima na precisão: $2^{-52} = \text{eps}$

Name	Binary	Decimal
eps	$2^{(-52)}$	2.2204e-16
realmin	$2^{(-1022)}$	2.2251e-308
realmax	$(2-\text{eps}) \times 2^{1023}$	1.7977e+308

Ver help do matlab com
floating point precision

Cleve's Corner

Floating points

IEEE Standard unifies arithmetic model

by Cleve Moler

If you look carefully at the definition of fundamental arithmetic operations like addition and multiplication, you soon encounter the mathematical abstraction known as the *real numbers*. But actual computation with real numbers is not very practical because it involves limits and infinities. Instead, MATLAB and most other technical computing environments use *floating-point* arithmetic, which involves a finite set of numbers with finite precision. This leads to phenomena like *roundoff error*, *underflow*, and *overflow*. Most of the time, MATLAB can be effectively used without worrying

where f is the fraction or mantissa and e is the exponent. The fraction must satisfy

$$0 \leq f < 1$$

and must be representable in binary using at most 52 bits. In other words, $2^{52}f$ must be an integer in the interval

$$0 \leq 2^{52}f < 2^{53}$$

The exponent must be an integer in the interval

$$-1022 \leq e \leq 1022$$

What is the output?

$$a = 4/3$$

$$b = a - 1$$

$$c = b + b + b$$

$$e = 1 - c$$

Conversão pode levar a séries infinitas

A very important example occurs with the simple MATLAB statement

```
t = 0.1
```

The value stored in `t` is not exactly 0.1 because expressing the decimal fraction $1/10$ in binary requires an infinite series.

In fact,

$$1/10 = 1/2^4 + 1/2^5 + 0/2^6 + 0/2^7 + 1/2^8 + 1/2^9 + 0/2^{10} + 0/2^{11} + 1/2^{12} + \dots$$

After the first term, the sequence of coefficients 1, 0, 0, 1 is repeated infinitely often. The floating-point number nearest 0.1 is obtained by rounding this series to 53 terms, including rounding the last four coefficients to binary 1010. Grouping the resulting terms together four at a time expresses the approximation as a base 16, or *hexadecimal*, series. So the resulting value of `t` is actually

$$t = (1 + 9/16 + 9/16^2 + 9/16^3 + \dots + 9/16^{12} + 10/16^{13}) \cdot 2^{-4}$$

Cleve Moler .

Ou seja:

Há números racionais que só admitem uma representação aproximada no computador.

Mas de notar que o mesmo já acontecia na representação decimal de $1/3$...

De forma mais geral ainda: qualquer número irracional tem de ser representado de forma aproximada num computador...

Acidentes causados por tratamento incorrecto dos erros de arredondamento

Ver a página na internet <http://ta.twi.tudelft.nl/nw/users/vuik/wi211/disasters.html#patriot>

Guerra do Golfo (1991):

24 bits $\Rightarrow 2^{24}=16777216$ números

Os relógios internos do bateria Patriot tinham registos de 24 bits. Além disso utilizavam a unidade da décima do segundo. A conversão para binário produz erros de truncatura da ordem de 0.00000000000000000000000011001100... ou seja 0.000000095 decimal. Estes erros podiam-se acumular ao longo do tempo, de tal maneira que se a bateria estivesse em funcionamento há 100 horas, os erros acumulados poderiam ser da ordem de $0.000000095 \times 100 \times 60 \times 60 \times 10 = 0.34$. Como um missil Scud viaja 1,676 metros por segundo, a imprecisão gerada pode ser superior ao raio de correcção/detecção do antimissil.



*The range gate's prediction of where the Scud will next appear is a function of the Scud's known velocity and the time of the last radar detection. Velocity is a real number that can be expressed as a whole number and a decimal (e.g., 3750.2563...miles per hour). Time is kept continuously by the system's internal clock in tenths of seconds but is expressed as an integer or whole number (e.g., 32, 33, 34...). The longer the system has been running, the larger the number representing time. To predict where the Scud will next appear, both time and velocity must be expressed as real numbers. Because of the way the Patriot computer performs its calculations and the fact that **its registers are only 24 bits long**, the conversion of time from an integer to a real number cannot be any more precise than 24 bits. **This conversion results in a loss of precision causing a less accurate time calculation.** The effect of this inaccuracy on the range gate's calculation is directly proportional to the target's velocity and the length of the system has been running. Consequently, performing the conversion after the Patriot has been running continuously for extended periods causes the range gate to shift away from the center of the target, making it less likely that the target, in this case a Scud, will be successfully intercepted.*

The Vancouver Stock Exchange

In 1982 the Vancouver Stock Exchange instituted a new index initialized to a value of 1000.000. The index was updated after each transaction. Twenty two months later it had fallen to 520. The cause was that the updated value was truncated rather than rounded. The rounded calculation gave a value of 1098.892.

Sources

The Wall Street Journal November 8, 1983, p.37.

The Toronto Star, November 19, 1983.

B.D. McCullough and H.D. Vinod, Journal of Economic Literature, Vol XXXVII (June 1999), pp. 633-665.

Como se faz uma soma?

```
>> x= 1.1*2^30 + 9*2^(-30)
```

```
x = 1.1811e+009
```

Os dois números são convertidos em números com igual expoente, e depois as mantissas são somadas.

```
>> y=1.1*2^30
```

```
y = 1.1811e+009
```

```
>> x==y
```

```
ans = 1
```

Por isso podem surgir arredondamentos importantes. Porém hoje a precisão é grande, chegando-se a utilizar 52 bits na parte da mantissa. Por isso neste exemplo não se nota o efeito do arredondamento.

Mas o mesmo não sucede em pequenos processadores (microcontroladores) utilizados em múltiplos equipamentos de electrónica.

Mais uma consequência de haver arredondamentos:

instruções `if x==y` são problemáticas sempre que x e y não forem inteiros.

Cancelamento Catastrófico

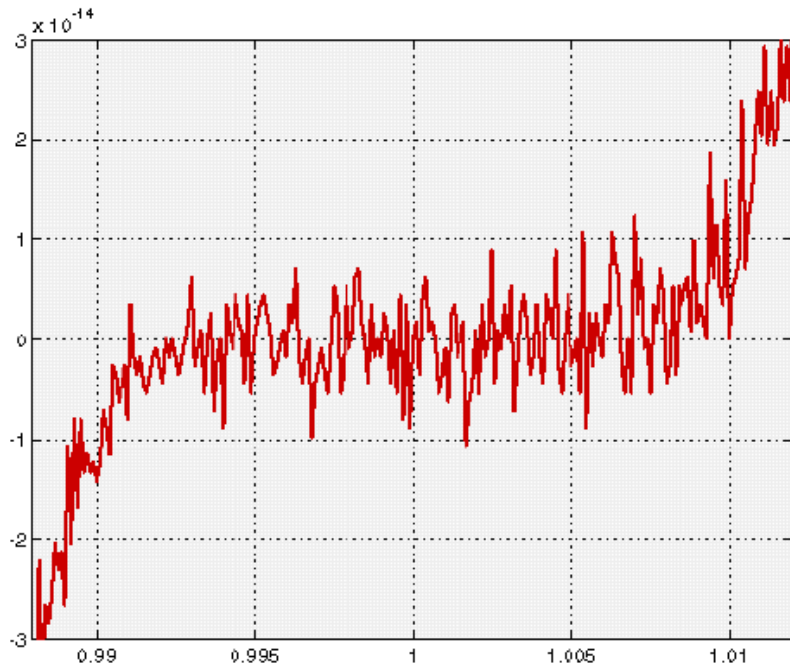
Se x e y forem números grandes, e os subtraímos podemos gerar erros de arredondamentos grandes em termos de erro relativo no resultado final:

$$f(x, y) = x - y$$

pode ser grande se os
números forem
grandes

$$\delta_{x-y} = \frac{\varepsilon_{x-y}}{\text{Diferença}} = \frac{\varepsilon_{x-y}}{x-y}$$

é pequeno quando há
cancelamento



But the resulting plot doesn't look anything like a polynomial. It isn't smooth. You are seeing roundoff error in action. The y -axis scale factor is tiny, 10^{-14} . The tiny values of y are being computed by taking sums and differences of numbers as large as $35 \cdot 1.012^4$. There is severe subtractive cancellation. The example was contrived by using the Symbolic Toolbox to expand $(x - 1)^7$ and carefully choosing the range for the x -axis to be near $x = 1$. If the values of y are computed instead by

$$y = (x-1).^7;$$

then a smooth (but very flat) plot results. ■

Our final example plots a seventh degree polynomial.

```
x = 0.988:.0001:1.012;
y = x.^7-7*x.^6+21*x.^5-35*x.^4+35*x.^3-...
    21*x.^2+7*x-1;
plot(x,y)
```

*Cleve Moler is chairm.
and co-founder of
The MathWorks.
His e-mail address is
moler@mathworks.com*

Como analisar facilmente o computador a cometer erros de arredondamento?

```
>> 1-0.9999  
ans = 1.0000e-004
```

```
>> format long  
>> 1-0.9999  
ans = 9.9999999999998899e-005
```

```
>> 1-0.99999  
ans = 9.9999999999954490e-006
```

```
>> 1-0.999999  
ans = 1.0000000000028756e-006
```

não altera a precisão com que
são feitos os cálculos
(só a apresentação)

Default: short

Como estudar de uma forma analítica o **condicionamento** de um algoritmo?

Objectivo: Calcular xy/z

Algoritmo 1

1) $x \cdot y$

2) res/z

Algoritmo 2

1) x / z

2) $\text{res} * y$

Algoritmo 1

$$1) \ x \ y \quad \Rightarrow \quad xy + \varepsilon_{xy}$$

$$2) \ \text{res}/z \Rightarrow (xy + \varepsilon_{xy})/z + \varepsilon_{xy/z}$$

$$\frac{xy}{z} \Rightarrow \frac{xy}{z} + \frac{\varepsilon_{xy}}{z} + \varepsilon_{xy/z}$$

Algoritmo 2

$$1) \ x/z \quad \Rightarrow \quad x/z + \varepsilon_{x/z}$$

$$2) \ \text{res} * y \quad \Rightarrow \quad (x/z + \varepsilon_{x/z}) \times y + \varepsilon_{xy/z}$$

$$\frac{xy}{z} \Rightarrow \frac{xy}{z} + y\varepsilon_{x/z} + \varepsilon_{xy/z}$$

Se x e y forem muito diferentes em magnitude, a magnitude do erro intermédio não diverge usando o primeiro algoritmo. Se x/z tiverem magnitudes semelhantes, pode-se utilizar o segundo algoritmo pois o erro intermédio também não diverge.

Dicas para evitar grandes erros

- 1) Devemos evitar que os cálculos intermédios tenham grandes números (por exemplo, usem expressões com singularidades) pois os erros nos cálculos intermédios podem tornar-se grandes, ainda que o resultado final possa não o vir a ser.

Por exemplo: para calcular $x y / z$ devemos:

- a) calcular primeiro xy e depois $/z$ se x e y forem muito diferentes em magnitude.
- b) devemos efectuar primeiro x/z ou y/z se estes números forem de magnitudes semelhantes.

Assim, numericamente certas equações não são equivalentes porque conduzem a erros muito diferentes

Exemplo:

$$f_1(x) = \sqrt{x}(\sqrt{x+1} - \sqrt{x}), \quad f_2(x) = \frac{\sqrt{x}}{\sqrt{x+1} + \sqrt{x}}$$

Na primeira expressão, podemos ter erros importantes se $x=10^{15}$ pois o erro será dessa ordem.

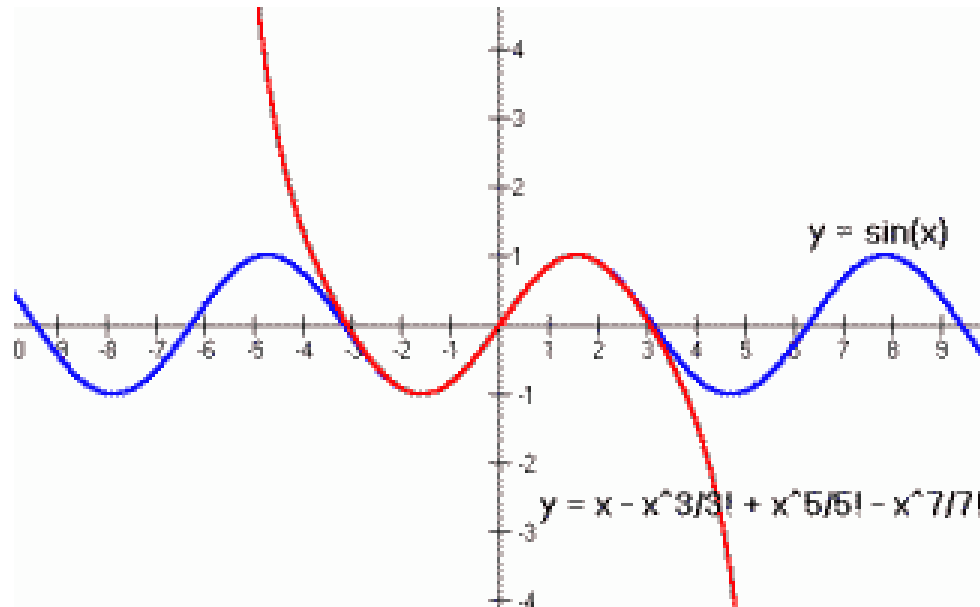
$$f_1(x) = \frac{1 - \cos x}{x^2}, \quad f_2(x) = \frac{\sin^2 x}{x^2(1 + \cos x)}$$

f_1 deve ser utilizada para calcular a expressão em $x \sim \pi$;
neste caso a expressão f_2 tem uma singularidade e diverge.

No entanto, f_2 deve ser utilizada para calcular a expressão em $x \sim 0$;
nesse caso a expressão f_1 tem uma singularidade e diverge.

TPC: Usar o MATLAB para representar a função
entre -2π e 2π .

Como é que o computador calcula o seno de um número

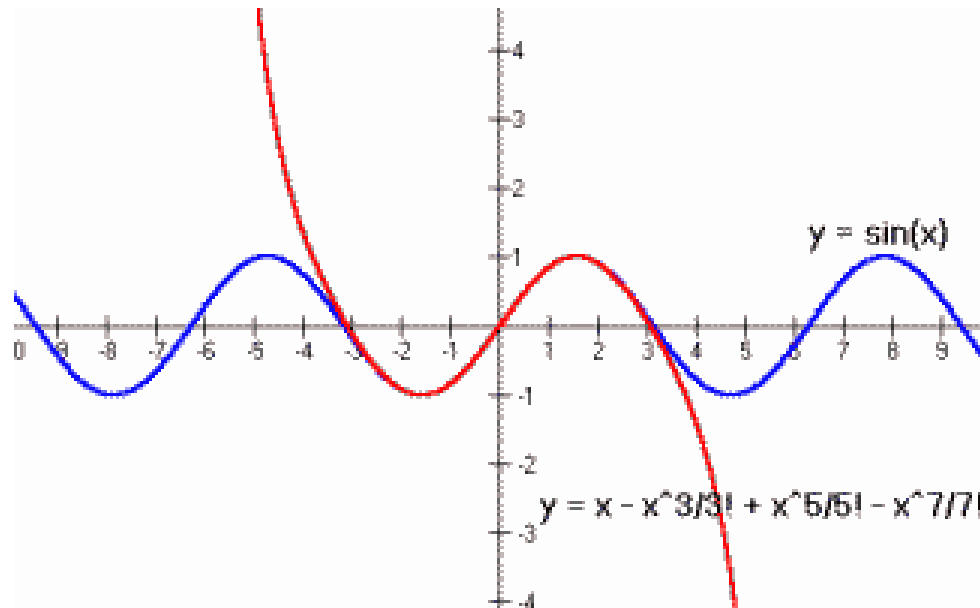


Questão: Quanto vale $\sin(2\pi)$?

E $\sin(20\pi)$? E $\sin(2000\pi)$?

TPC: Usar o Excel para fazer o gráfico de $\sin(2n\pi)$ entre $n=0$ e 300.

Como é que o computador calcula o seno de um número



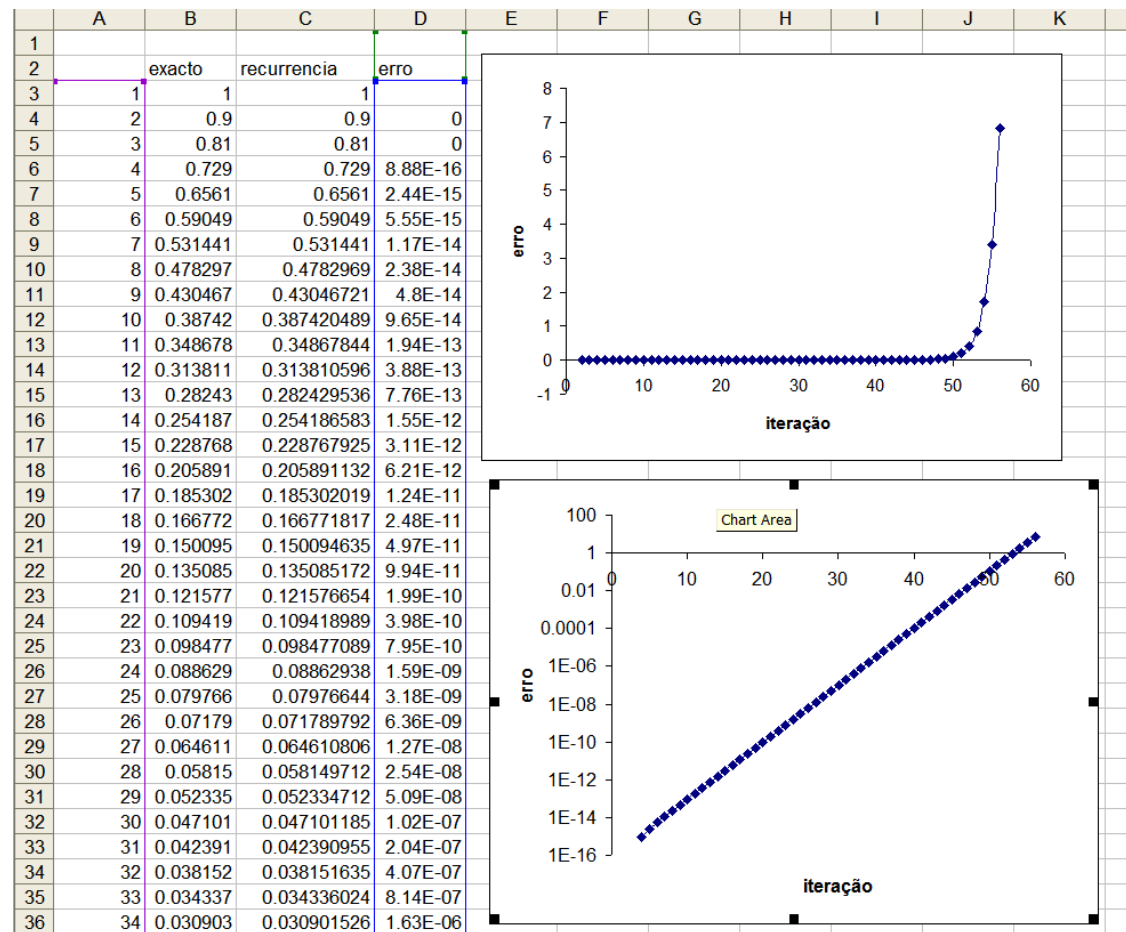
O uso de séries de Taylor não é apropriado, pelo menos em todo o domínio. Os algoritmos usam diversas técnicas (mapeamentos, aproximações com polinómios de Chebyshev).

Exemplo: cálculos exactos versus numéricos

Sequência: $f_{n+1} = 2.9 f_n - 1.8 f_{n-1}$

Solução exacta: $f_n = (0.9)^{n-1}$

Ainda que seja desejável conhecer soluções exactas, isso nem sempre é possível, sendo essencial compreender em que medida a natureza dos resultados numéricos não são um artefacto do método numérico.



Eficiência Computacional (em termos de número de operações)

$$p_4(x) = a_1x^4 + a_2x^3 + a_3x^2 + a_4x + a_5$$

$$p_{4n}(x) = (((a_1x + a_2)x + a_3)x + a_4)x + a_5$$

A multiplicações são operações que requerem mais operações elementares que as adições. Por isso, uma computação eficiente deve evitá-las...
No primeiro caso gastamos 4 multiplicações para calcular a_1x^4 , tantas quantas as multiplicações usadas no segundo método.