

- (d) $(0.02)_{F!} = \frac{0}{2!} + \frac{2}{3!} = (\frac{2}{6})_{10} = (\frac{1}{3})_{10} = (\mathbf{0.333333\dots})_{10}$
(e) $(0.113)_{F!} = \frac{1}{2!} + \frac{1}{3!} + \frac{3}{4!} = (\frac{19}{24})_{10} = (\mathbf{0.791\bar{6}}\dots)_{10}$
(f) $(321.123)_{F!} = 3 \cdot 3! + 2 \cdot 2! + 1 \cdot 1! + \frac{1}{2!} + \frac{2}{3!} + \frac{3}{4!} = (\frac{575}{24})_{10} = (\mathbf{23.958\bar{3}})_{10}$

- 4) (a) $(10111.1101)_2 = (0001\ 0111.1101)_2 = (17.D)_{16} = 1 \cdot 16^1 + 7 \cdot 16^0 + 13 \cdot 16^{-1} = (\frac{381}{16})_{10} = (\mathbf{23.8125})_{10}$
(b) $(BD.0E)_{16} = (1011\ 1101.0000\ 1110)_2 = 1 \cdot 2^7 + 1 \cdot 2^5 + 1 \cdot 2^4 + 1 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^0 + 1 \cdot 2^{-5} + 1 \cdot 2^{-6} + 1 \cdot 2^{-7} = (\frac{24199}{128})_{10} = (\mathbf{189.0546875})_{10}$
(c) $(41.1)_{10} = (?)_2 = (?)_{16}$

$$(41)_{10} = (?)_2 \longrightarrow \begin{array}{c} 41 \begin{array}{|l} 2 \\ 1 \end{array} \begin{array}{|l} 20 \\ 0 \end{array} \begin{array}{|l} 2 \\ 10 \\ 0 \end{array} \begin{array}{|l} 2 \\ 5 \\ 1 \end{array} \begin{array}{|l} 2 \\ 2 \\ 0 \end{array} \begin{array}{|l} 2 \\ 1 \\ 1 \end{array} \begin{array}{|l} 2 \\ 0 \end{array} \end{array}$$

←

$$(41)_{10} = (101001)_2 = 1 \cdot 2^5 + 1 \cdot 2^3 + 1 \cdot 2^0$$

$$(0.1)_{10} = (?)_2 \longrightarrow$$

$$\begin{array}{cccccccccccc} 0.1 & 0.2 & 0.4 & 0.8 & 0.6 & 0.2 & 0.4 & 0.8 & 0.6 & 0.2 & & \\ \times 2 & \times 2 & \times 2 & \times 2 & \times 2 & \times 2 & \times 2 & \times 2 & \times 2 & \times 2 & \times 2 & \dots \\ \mathbf{0.2} & \mathbf{0.4} & \mathbf{0.8} & \mathbf{1.6} & \mathbf{1.2} & \mathbf{0.4} & \mathbf{0.8} & \mathbf{1.6} & \mathbf{1.2} & \mathbf{0.4} & & \end{array}$$

$$(0.1)_{10} = (0.00011001100110011\dots)_2$$

Logo, $(41.1)_{10} = (101001.\overline{00011})_2 = (0010\ 1001.0001\ \overline{1001})_2 = (\mathbf{29.19})_{16}$. Haverá perda de dígitos significativos.

5) $F(\beta, t, I, S) \longrightarrow F(2, 3, -3, +3) \longrightarrow \boxed{s_1 \mid d_1 \mid d_2 \mid d_3 \mid s_2 \mid e_1 \mid e_2 \mid e_3}$

(a) $NM = (\beta - 1) \cdot \beta^{t-1} = (2 - 1) \cdot 2^{3-1} = \mathbf{4}$

(b) $NE = S - I + 1 = 3 - (-3) + 1 = \mathbf{7}$

(c) $NR = 2 \cdot NM \cdot NE + 1 = \mathbf{57}$

(d) $\boxed{0 \mid 1 \mid 0 \mid 0 \mid 1 \mid 1 \mid 0 \mid 0} \longrightarrow m.p. = (0.1)_2 \cdot (2^{-3})_{10} = (2^{-1})_{10} \cdot (2^{-3})_{10} = (0.0625)_{10}$

Logo, a região de *underflow* é $\{x \in \mathbb{R} \mid -(\mathbf{0.0625})_{10} < x < (\mathbf{0.0625})_{10}\}$.

$\boxed{0 \mid 1 \mid 1 \mid 1 \mid 0 \mid 1 \mid 0 \mid 0} \longrightarrow M.P. = (0.111)_2 \cdot (2^3)_{10} = (0.875)_{10} \cdot (2^3)_{10} = (7)_{10}$

Logo, a região de *overflow* é $\{x \in \mathbb{R} \mid x < -(\mathbf{7})_{10} \cup x > (\mathbf{7})_{10}\}$.

(e) A precisão decimal equivalente é aproximadamente $\log_{10}(2^3) = 3 \cdot \log_{10}(2) \approx \mathbf{0.90308998699}$.

(f) $F(2, 3, 0, 6)$ (polarização $p = +3$)

6) $F(\beta, t, I, S) \longrightarrow F(2, 3, 0, +7) \longrightarrow \boxed{s_1 \mid d_1 \mid d_2 \mid d_3 \mid 0 \mid e_1 \mid e_2 \mid e_3}$

(a) $NM = (\beta - 1) \cdot \beta^{t-1} = (2 - 1) \cdot 2^{3-1} = \mathbf{4}$

(b) $NE = S - I + 1 = 7 - 0 + 1 = \mathbf{8}$

(c) $NR = 2 \cdot NM \cdot NE + 1 = \mathbf{65}$

(d) $\boxed{0 \mid 1 \mid 0 \mid 0 \mid 0 \mid 0 \mid 0 \mid 0} \longrightarrow m.p. = (0.1)_2 \cdot (2^0)_{10} = (2^{-1})_{10} \cdot (1)_{10} = (\mathbf{0.5})_{10}$

Logo, a região de *underflow* é $\{x \in \mathbb{R} \mid -(\mathbf{0.5})_{10} < x < (\mathbf{0.5})_{10}\}$.

$\boxed{0 \mid 1 \mid 1 \mid 1 \mid 0 \mid 1 \mid 1 \mid 1} \longrightarrow m.p. = (0.111)_2 \cdot (2^7)_{10} = (0.875)_{10} \cdot (2^7)_{10} = (\mathbf{112})_{10}$

Logo, a região de *overflow* é $\{x \in \mathbb{R} \mid x < -(\mathbf{112})_{10} \cup x > (\mathbf{112})_{10}\}$.

(e) A precisão decimal equivalente é aproximadamente $\log_{10}(2^3) = 3 \cdot \log_{10}(2) \approx \mathbf{0.90308998699}$.

- 7) O formato de ponto flutuante de precisão dupla tem 1 bit de sinal, 11 bits para expoente e 52 bits explícitos para decimais (o 1º bit é implícito por conta do armazenamento do número de acordo com o padrão IEEE 754).

Também foi decidido neste padrão **facilitar** o armazenamento do expoente. Assim, um conceito chamado *exponent bias*, calculado por $2^{k-1} - 1$ onde k é o número de bits do expoente, foi criado com esta função. Este *bias* é essencial para a codificação do expoente como um valor **unsigned**. Para obter o expoente verdadeiro, subtrai-se o *bias* do expoente representado.

Assim sendo, $bias = 2^{11-1} - 1 = 1023$. $e_{min} = 1 - 1023 = -1022$ e $e_{max} = 2046 - 1023 = 1023$. (Os expoentes 0 e 2047 são reservados.)

Portanto, o menor número passível de ser representado com precisão dupla é $2^{-1022} \approx \mathbf{2.225 \cdot 10^{-308}}$, e o maior número é 2^{1023} com todos os bits da mantissa ativados, ou seja, $2^{1023} \cdot (2 - 2^{-52}) \approx \mathbf{1.797 \cdot 10^{308}}$. Às custas da perda de precisão, existem números que vão abaixo da fronteira de *underflow*, chamados de **subnormais**.


```

17, 18) from functools import reduce
from math import factorial as fact
from numpy import float32 as single

n = 4
x = -.111
desired = 0.894938748929031

for y in [single(x), x]:
    obtained = reduce(lambda w, i: w + (y**i)/fact(i), range(n), 0)
    print("{:.70f}%".format((desired - obtained) / desired))

```

- O decimal desejado foi retirado do [Wolfram Alpha](#), para comparação com os resultados calculados pelo programa em variáveis de precisão simples e dupla, respectivamente. O erro da variável de precisão simples é maior pela quantidade reduzida de bits disponíveis para armazenar a parte decimal do resultado.
- Novamente, a função *reduce* é utilizada para simular o cálculo da série de Maclaurin com os argumentos fornecidos.
- Saída do programa:

```

0.0000069152325516942810432695874778286082573686144314706325531005859375%
0.0000069138016857893659460184906939694826633058255538344383239746093750%

```