

# Tema 1 - Más sobre errores de truncamiento

Curso de Física Computacional

M. en C. Gustavo Contreras Mayén

## 1 Números de punto flotante

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F
- 3 Más ejercicios.

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F
- 3 Más ejercicios.
- 4 Error por corte/redondeo

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F
- 3 Más ejercicios.
- 4 Error por corte/redondeo
- 5 Errores de truncamiento

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F
- 3 Más ejercicios.
- 4 Error por corte/redondeo
- 5 Errores de truncamiento
- 6 Acumulación del error por redondeo

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F
- 3 Más ejercicios.
- 4 Error por corte/redondeo
- 5 Errores de truncamiento
- 6 Acumulación del error por redondeo



# Números de punto flotante

Un conjunto  $F$  de números de punto flotante está caracterizado por los siguientes parámetros:

- 1 La base del sistema  $\beta$ .
- 2 El número de dígitos  $n$  en la mantisa.
- 3 Un exponente  $m \leq e \leq M$ .

donde  $\beta$ ,  $n$ ,  $m$ ,  $e$  son enteros.

Cada número en el sistema de punto flotante  $F$  tiene la forma

$$\pm (.d_1 d_2 \dots d_n)_\beta \beta^e$$

donde  $d_i = 0, 1, \dots, \beta - 1$ ,  $i = 1, 2, \dots, n$ ,  
 $(.d_1 d_2 \dots d_n)_\beta$  es una  $\beta$ -fracción llamada *mantisa*,  $e$   
es un entero llamado *exponente*.

El sistema F de punto flotante está *normalizado* si  $d_1 \neq 0$ , en general, todos los números flotantes se normalizan con la excepción del cero, en el cual  $d_1 = d_2 = \dots = d_n = 0$ .

El conjunto F es discreto y finito. Su cardinalidad (i.e. el número de elementos que lo constituyen, está dada por

$$2(\beta - 1)\beta^{n-1}(M - m + 1) + 1$$

Como consecuencia de la finitud de F para representar al conjunto de números reales  $\mathbb{R}$ , existirá una infinidad de números en  $\mathbb{R}$  que no pueden representarse en forma exacta en F.

Sea  $x$  un número real denotemos por  $fl(x)$  el número en  $F$  que es más cercano a  $x$ . La diferencia entre  $x$  y  $fl(x)$  se llama *error de redondeo*, éste depende de la magnitud de  $x$  y es por lo tanto medido relativo a  $x$

$$\delta(x) = \frac{fl(x) - x}{x}, \quad x \neq 0$$

luego  $|\delta(x)|$  es el error relativo introducido en la representación de  $x$  en el sistema de punto flotante  $F$ ; de la ecuación anterior obtenemos que

$$fl(x) = x(1 + \delta(x))$$

Hay dos formas generalmente usadas para convertir un número real  $x$  a un  $n - \beta$  número flotante  $fl(x)$ : *redondeado* o *truncando*.

Cuando se redondea,  $fl(x)$  se elige como el número de punto flotante normalizado más cercano a  $x$ , si hay empate se usa alguna regla especial, por ejemplo, se toma el de la derecha.

Si se trunca,  $fl(x)$  se escoge como el número flotante normalizado más cercano entre 0 y  $x$ , esto es, se toman algunas d's y se desprecian otras.

A continuación encontraremos una cota para  $\delta(x)$  que sea independiente de  $x$ . En el sistema de números de punto flotante F, el número cuyo valor absoluto es el más pequeño está dado por

$$+ (.100 \dots 0)_{\beta} \beta^m = \beta^{m-1}$$

el sucesor inmediato de  $\beta^{m-1}$  se encuentra sumando a éste el número

$$+ (.000 \dots 1)_{\beta} \beta^m = \beta^{m-n}.$$

Se concluye entonces, que la distancia entre dos números consecutivos en el intervalo  $[\beta^{m-1}, \beta^m]$  es  $\beta^{m-n}$ .

En forma similar se demuestra que la distancia entre dos números consecutivos que pertenezcan a cualquier intervalo de la forma

$[\beta^j, \beta^{j+1}] j = m, \dots, M - 1$  está dada por

$$\beta^{j+1-n}$$

Hemos demostrado un aspecto singular de la distribución de los números del sistema de punto flotante F:

- *estos no están igualmente espaciados a través de todo su rango, sino únicamente cuando se encuentran entre potencias sucesivas de la base  $\beta$ .*

Sunpongamos que  $x \in [\beta^j, \beta^{j+1}]$  para alguna  $j = m, \dots, M - 1$ , si el número flotante  $fl(x)$  que representa a  $x$  es seleccionado por redondeo, entonces de acuerdo con el resultado anterior, el error introducido es a lo más  $(1/2)\beta^{j+1-n}$ , si  $fl(x)$  se selecciona por truncamiento el error es a lo más  $\beta^{j+1-n}$ . Lo anterior nos da la *medida del error de redondeo absoluto*

$$|fl(x) - x| \leq \begin{cases} \frac{1}{2}\beta^{j+1-n} & \text{redondeo,} \\ \beta^{j+1-n} & \text{truncamiento} \end{cases}$$

El *error de redondeo relativo*  $|\delta(x)|$  se obtiene dividiendo al error de redondeo absoluto por  $|x|$ .

Dado que  $0 < \beta^j \leq |x|$ , se tiene

$$|\delta(x)| \leq \begin{cases} \frac{1}{2}\beta^{1-n} & \text{redondeo,} \\ \beta^{1-n} & \text{truncamiento} \end{cases}$$



El  $\epsilon$  de la máquina se define como

$$\epsilon = \begin{cases} \frac{1}{2}\beta^{1-n} & \text{redondeo,} \\ \beta^{1-n} & \text{truncamiento} \end{cases}$$

De lo anterior, tenemos que se cumple  $|\delta(x)| \leq \epsilon$  para toda  $x$ .

La cota para  $\delta(x)$  independiente de  $x$  es el  $\epsilon$  de la máquina.

Dado  $x \in \mathbb{R}$ , su flotante  $fl(x)$  se define como

$$fl(x) = x(1 + \delta), \quad |\delta| \leq \epsilon$$

La exactitud de la aritmética de punto flotante está entonces caracterizada por el  $\epsilon$  de la máquina.

Supongamos que estamos en el intervalo  $[1, \beta]$ . Queremos calcular  $1 \oplus \epsilon$ , donde  $\oplus$  denota la suma entre números que pertenecen a  $F$

$$1 \oplus \epsilon = fl(1 + \epsilon)$$

Si el error introducido en la representación de  $1 + \epsilon$  es por redondeo, entonces éste queda localizado en la mitad del intervalo  $[1, 1 + \epsilon^{1-n}]$ .

Ya que la distancia entre dos números consecutivos de  $F$  que pertenecen al intervalo  $[1, \beta]$  es  $\beta^{1-n}$ , el número de punto flotante  $1 \oplus \epsilon$  es tomado como  $1 + \beta^{1-n} > 1$ .

Sin embargo, si  $0 < \epsilon_1 < \epsilon$  obtenemos por un procedimiento análogo al anterior que  $1 \oplus \epsilon_1 = 1$ .

Un resultado similar es encontrado cuando el error introducido en la representación de  $1 \oplus \epsilon$  es por truncamiento.

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F
- 3 Más ejercicios.
- 4 Error por corte/redondeo
- 5 Errores de truncamiento
- 6 Acumulación del error por redondeo

# Modelo de aritmética en F

La aritmética en el sistema numérico de punto flotante F permite aproximar a la del sistema de números reales  $\mathbb{R}$ .

Como notación emplearemos  $\oplus, \ominus, \otimes, \oslash$  para indicar las aproximaciones a las operaciones aritméticas  $+, -, \times, /$  de  $\mathbb{R}$ :

$$x \oplus y = fl(x + y)$$

$$x \ominus y = fl(x - y)$$

$$x \otimes y = fl(x \times y)$$

$$x \oslash y = fl(x / y)$$

El modelo que asumiremos para la aritmética en F es el siguiente:

$$fl(x \text{ op } y) = (x \text{ op } y)(1+\delta), \quad |\delta| \leq \epsilon, \quad \text{op} = +, -, *, /$$

Para efectuar operaciones en forma manual en este modelo aritmético, por cada operación  $+$ ,  $-$ ,  $\times$ ,  $/$  encontrada, hágala en aritmética exacta, normalice el resultado, trunque o redondee de acuerdo al número de dígitos permitido.

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F
- 3 Más ejercicios.
- 4 Error por corte/redondeo
- 5 Errores de truncamiento
- 6 Acumulación del error por redondeo

# Problema 1

Considera la siguiente suma finita

$$S_N^{(1)} = \sum_{n=1}^{2N} (-1)^n \frac{n}{n+1} \quad (1)$$

Si sumamos de manera separada los valores impares y los pares de  $x$ , tendremos dos sumas:

$$S_N^{(2)} = - \sum_{n=1}^N \frac{2n-1}{2n} + \sum_{n=1}^N \frac{2n}{2n+1} \quad (2)$$



## Tercera suma

Podemos eliminar la diferencia mediante una combinación entre las dos sumas, quedando de la siguiente manera

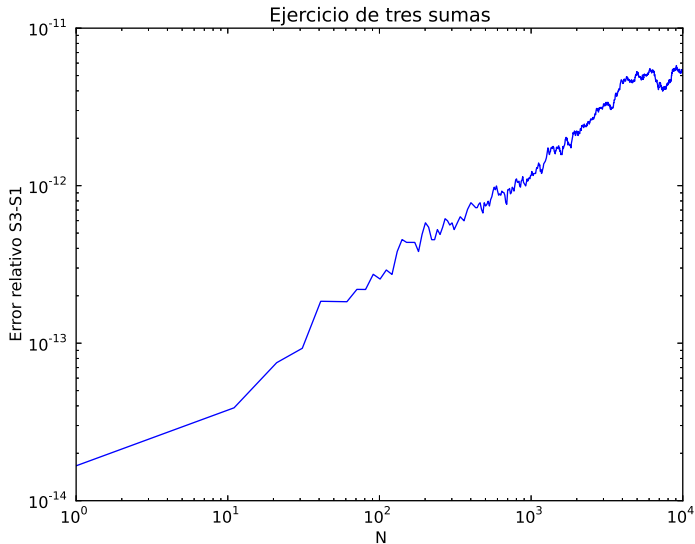
$$S_N^{(3)} = \sum_{n=1}^N \frac{1}{2n(2n+1)} \quad (3)$$

Sabemos que aunque el valor de las tres sumas  $S_N^{(1)}$ ,  $S_N^{(2)}$ ,  $S_N^{(3)}$ , es el mismo, pero el resultado numérico puede ser diferente.

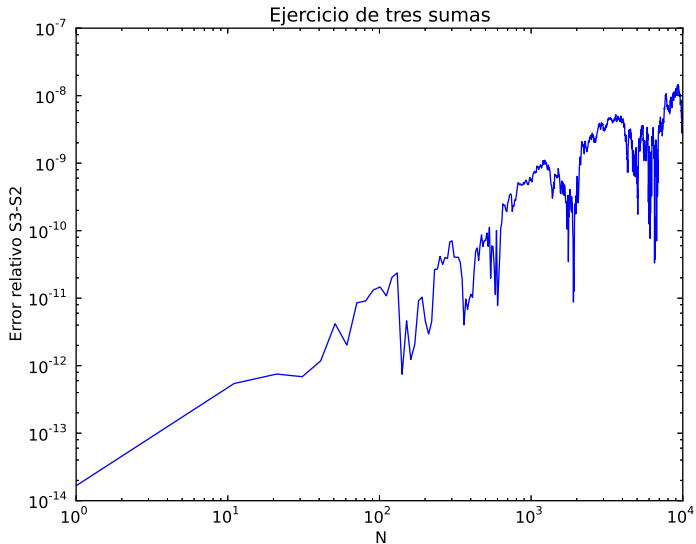
# Ejercicio a resolver

- 1 Escribe un programa que calcule  $S_N^{(1)}$ ,  $S_N^{(2)}$ ,  $S_N^{(3)}$ .
- 2 Supongamos que  $S_N^{(3)}$  es el valor exacto de la suma. Grafica el error relativo contra el número de términos en la suma (tip: usa una escala log-log). Comienza con  $N = 1$  hasta  $N = 1000000$ . Describe la gráfica.
- 3 Identifica en tu gráfica una región en donde la tendencia es casi lineal, ¿qué representa ésta sección con respecto al error?

# Error relativo entre $S_3$ y $S_1$



# Error relativo entre $S_3$ y $S_2$



## Problema 2

Aunque tengamos el apoyo de una buena computadora, el cálculo de la suma de una serie requiere reflexión y cuidado.

Considera la serie:

$$S^{(u)} = \sum_{n=1}^N \frac{1}{n}$$

que será una suma finita mientras  $N$  sea finito.

Cuando hacemos la suma de manera analítica, no importa si se hace de manera ascendente: desde  $n = 1, 2, 3, \dots, N - 1, N$ , o descendente: desde  $n = N, N - 1, N - 2, \dots, 3, 2, 1$

$$S^{(d)} = \sum_{n=N}^1 \frac{1}{n}$$

Sin embargo, debido a los errores por redondeo, cuando calculamos de manera analítica, el valor de las sumas no es el mismo,  $S^{(u)} \neq S^{(d)}$

- 1 Escribe un programa que calcule  $S^{(u)}$  y  $S^{(d)}$  como función de  $N$ .
- 2 Grafica (log-log) la diferencia relativa entre la suma relativa contra  $N$ .
- 3 Identifica en tu gráfica una región en donde la tendencia es casi lineal, ¿qué representa ésta sección con respecto al error?

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F
- 3 Más ejercicios.
- 4 Error por corte/redondeo
- 5 Errores de truncamiento
- 6 Acumulación del error por redondeo

Volvamos a nuestro sistema decimal tradicional. Supongamos ahora que los números se pueden representar de la siguiente manera:

$$fl(x) = \pm(0.d_1d_2d_3 \dots d_td_{t+1}d_{t+2} \dots) \times 10^e$$

Si la precisión elegida es  $t$ , entonces "recortar" el número definido arriba, pues no podemos representar los  $d_i$  para  $i > t$ .



En consecuencia, tenemos dos alternativas básicas para efectuar dicho recorte:

- ① **Corte:** Ignorar los dígitos  $d_i$  cuando  $i > t$

En consecuencia, tenemos dos alternativas básicas para efectuar dicho recorte:

- 1 **Corte:** Ignorar los dígitos  $d_i$  cuando  $i > t$
- 2 **Redondeo:** Sumar 1 a  $d_t$  si  $d_{t+1} \geq \frac{10}{2}$  e ignorar los restantes  $d_i$  para  $i > t + 1$ , o aplicar corte si  $d_{t+1} < \frac{10}{2}$

Esto nos permite obtener una cota del error absoluto para ambos casos:

$$e_A = \begin{cases} 10^{-t} \times 10^e & \text{para corte} \\ \frac{1}{2} 10^{-t} \times 10^e & \text{para redondeo} \end{cases}$$

Y como definimos el error absoluto, también podemos definir un límite para el error relativo, que será:

**Corte:**

$$e_r \leq \frac{10^{-t} \times 10^e}{0.1 \times 10^e} = 10^{1-t}$$

**Redondeo:**

$$e_r \leq \frac{1}{2} \frac{10^{-t} \times 10^e}{0.1 \times 10^e} = \frac{1}{2} 10^{1-t}$$

Al valor  $10^{1-t}$  lo identificaremos con la letra  $\mu$ , y resulta ser importante porque nos da una idea del error relativo que cometemos al utilizar una representación de coma flotante. Suele denominarse como **unidad de máquina o unidad de redondeo**. El negativo del exponente de  $\mu$  suele llamarse también *cantidad de dígitos significativos*.

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F
- 3 Más ejercicios.
- 4 Error por corte/redondeo
- 5 Errores de truncamiento**
- 6 Acumulación del error por redondeo

# Errores de truncamiento

Sabemos que este error surge de aproximar procesos continuos mediante procedimientos discretos o de procesos "infinitos" mediante procedimientos "finitos".

Como ejemplo suele tomarse la diferenciación numérica como forma de aproximar el cálculo de una derivada en un punto (o su equivalente, la integración numérica), en tanto que para el caso de discretización, el ejemplo más usual es la utilización de métodos iterativos para resolver sistemas de ecuaciones lineales.

En general, el error de truncamiento está asociado al uso de la serie de Taylor para aproximar funciones, de modo que estimar una cota del error no conlleva una dificultad mayor. Sin embargo, en él suelen interactuar el error inherente y/o el de redondeo, con lo que muchas veces su influencia no es bien advertida o es muy reducida.

Veamos un ejemplo clásico: Supongamos que queremos calcular una aproximación de  $f'(x_0)$  para una función continua, pues no es posible obtener la derivada en forma analítica o resulta muy difícil. Por lo tanto, usaremos un entorno del punto  $x_0$  para calcular  $f'(x_0)$  utilizando solamente  $f(x)$ .



Para ello nos valdremos de la serie de Taylor. En efecto, para cualquier punto distante  $h$  de  $x_0$  tendremos:

$$\begin{aligned} f(x_0 + h) = & f(x_0) + f'(x_0)h + f''(x_0)\frac{h^2}{2} + \\ & + f'''(x_0)\frac{h^3}{6} + f^4(x_0)\frac{h^4}{24} + \dots \end{aligned}$$

Para ello nos valdremos de la serie de Taylor. En efecto, para cualquier punto distante  $h$  de  $x_0$  tendremos:

$$\begin{aligned} f(x_0 + h) = & f(x_0) + f'(x_0)h + f''(x_0)\frac{h^2}{2} + \\ & + f'''(x_0)\frac{h^3}{6} + f^4(x_0)\frac{h^4}{24} + \dots \end{aligned}$$

despejamos  $f'(x_0)$ , por tanto

$$\begin{aligned} f'(x_0) = & \frac{f(x_0 + h) - f(x_0)}{h} + \\ - & \left[ f''(x_0)\frac{h^2}{2} + f'''(x_0)\frac{h^3}{6} + f^4(x_0)\frac{h^4}{24} + \dots \right] \end{aligned}$$

Si el algoritmo que proponemos para aproximar  $f'(x_0)$  es

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

El error que se comete en la aproximación viene dado por:

$$\begin{aligned} & \left[ f'(x_0) - \frac{f(x_0 + h) - f(x_0)}{h} \right] = \\ & = \left[ f''(x_0) \frac{h^2}{2} + f'''(x_0) \frac{h^3}{6} + f^{(4)}(x_0) \frac{h^4}{24} + \dots \right] \end{aligned}$$

El término de la derecha es el denominado error de truncamiento, pues es lo que se truncó a la serie de Taylor para aproximar el valor buscado.

Este error suele asociarse también con la convergencia (o la velocidad de convergencia), que suele representarse como  $O(n)$  (generalmente, como  $O(h^n)$ , siendo  $n$  el parámetro que determina la velocidad o la convergencia.

En nuestro ejemplo, y dado que  $h$  generalmente es menor a 1, podemos decir que la aproximación es del tipo:

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

En donde el error que se comete es proporcional a  $h$ .

Se verifica que además están los términos con  $h^2$ ,  $h^3$ , etc. pero como  $h < 1$  se tiene que  $h^2 \ll h$ ,  $h^3 \ll h^2$ , etc. por lo que la influencia de éstos es mucho menos y despreciable.

Supongamos por un momento que todas las derivadas  $f^i(x_0) = 0$  para  $i \geq 3$ . Entonces, tenemos que:

$$\left[ f'(x_0) - \frac{f(x_0 + h) - f(x_0)}{h} \right] = \frac{h}{2} |f''(\xi)|$$

con  $\xi \in [x, x + h]$

por lo que, si conociéramos  $f''(\xi)$  se podría acotar el error que se está cometiendo por despreciar el término  $h/2 f''(x_0)$

Como ejercicio, apliquemos el algoritmo para obtener la derivada en  $x_0 = 0.45$ , es decir,  $f'(0.45)$  de la función  $f(x) = \sin(2\pi x)$ .

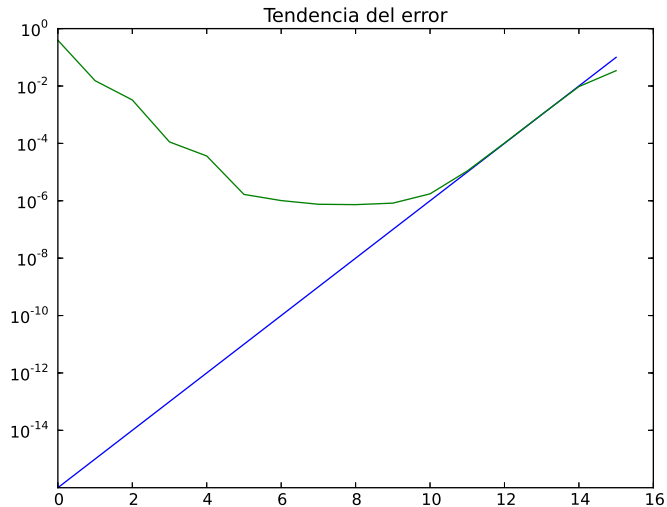
Considera el valor exacto de la derivada  $f'(0.45) = 2\pi \cos(2\pi * 0.45) = -5.97566$ , toma el valor de  $h = 0.1$

# Construye una tabla

$h$	$f'(x_0)$	Error
$10^{-1}$		
$10^{-2}$		
$10^{-3}$		
$10^{-4}$		
$10^{-5}$		
$10^{-6}$		
$\dots$		
$10^{-16}$		



# Gráfica de la tendencia del error



$h$	$f'(x_0)$	Error
1e-01	-6.180340	3.425226e-02
1e-02	-6.032711	9.547183e-03
1e-03	-5.981725	1.014908e-03
1e-04	-5.976274	1.027353e-04
1e-05	-5.975725	1.093152e-05
1e-06	-5.975670	1.745271e-06
1e-07	-5.975665	8.269814e-07
1e-08	-5.975664	7.339002e-07
1e-09	-5.975665	7.561951e-07
1e-10	-5.975666	1.016302e-06

$h$	$f'(x_0)$	Error
1e-11	-5.975670	1.666570e-05
1e-12	-5.975442	3.642056e-05
1e-13	-5.976331	1.122121e-04
1e-14	-5.995204	3.270657e-03
1e-15	-5.884182	1.530843e+02
1e-16	-8.326673	3.934315e+01

Si analizamos en detalle, vemos que la tendencia del error de truncamiento es lineal (en escala logarítmica) pero para  $h < 10^{-8}$  el error aumenta y no sigue una ley determinada. Este "empeoramiento" de la aproximación se debe a la incidencia del error de redondeo, es decir, la unidad de máquina pasa a ser más importante que el error de truncamiento.

Es por eso que no siempre el utilizar una "mejor precisión" ayuda a mejorar los resultados finales. En este tipo de problemas, es conveniente que el error que domine los cálculos sea el de truncamiento o dediscretización.

# Contenido

- 1 Números de punto flotante
- 2 Modelo de aritmética en F
- 3 Más ejercicios.
- 4 Error por corte/redondeo
- 5 Errores de truncamiento
- 6 Acumulación del error por redondeo

Desde que se creó la primera computadora, la acumulación del error de redondeo ha sido uno de los "dolores de cabeza" de los especialistas, como se puede ver en esta frase:

"La extraordinaria rapidez de las actuales computadoras significa que en un problema típico se realizan millones de operaciones con coma (punto) flotante. Esto quiere decir que la acumulación de errores de redondeo puede ser desastrosa".

En muchas ocasiones la inestabilidad está dada por la incidencia de unos pocos errores de redondeo y no por la acumulación de millones de ellos.

Un ejemplo en ese sentido está dado por el algoritmo del ejemplo inicial, en el cual el error está dado por el redondeo de  $y_{n-1}$ , que se propaga a medida que el valor es cada vez más chico.

Calcula el valor de  $e$  para  $n$  suficientemente grandes, a partir de la su definición:

$$f(n) = \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n$$



Completa la tabla:

$n$	$f(n)$	$ \exp - f(n) $
$10^1$		
$10^2$		
$10^3$		
$\dots$		
$10^{14}$		
$10^{15}$		

Discute tus resultados!