

# Tema 1 - Escalas, condición y estabilidad

## Curso de Física Computacional

M. en C. Gustavo Contreras Mayén

Facultad de Ciencias - UNAM

2 de septiembre de 2017



1. Precisión de la máquina

2. Tipos de errores

3. Conceptos importantes

4. Algoritmo de Horner

# 1. Precisión de la máquina

## 1.1 Precisión en punto flotante

## 1.2 Definición de la precisión de la máquina

## 1.3 Calculando la precisión de la máquina

# 2. Tipos de errores

# 3. Conceptos importantes

# 4. Algoritmo de Horner

Una preocupación importante en el cálculo computacional para la representación en punto flotante que se utiliza para almacenar números es que se cuenta con una precisión limitada.

# Precisión de la máquina

Como hemos visto, para una máquina de 32 bits, los números de precisión simple son buenos para 6 – 7 decimales, mientras que los de precisión doble son buenos para 15 – 16 lugares.

# Precisión de la máquina

Como hemos visto, para una máquina de 32 bits, los números de precisión simple son buenos para 6 – 7 decimales, mientras que los de precisión doble son buenos para 15 – 16 lugares.

Para ver cómo la precisión limitada afecta a los cálculos, consideremos la suma en la computadora de dos palabras de precisión simple:

# Precisión de la máquina

Como hemos visto, para una máquina de 32 bits, los números de precisión simple son buenos para 6 – 7 decimales, mientras que los de precisión doble son buenos para 15 – 16 lugares.

Para ver cómo la precisión limitada afecta a los cálculos, consideremos la suma en la computadora de dos palabras de precisión simple:

$$7 + 1.0 \times 10^{-7}$$

# Representación de los valores

La computadora extrae estos números de la memoria y los almacena en cadenas de bits

$$\begin{aligned} 7 &= 0\ 10000010\ 1110\ 0000\ 0000\ 0000\ 0000\ 000 \\ 10^{-7} &= 0\ 01100000\ 1101\ 0110\ 1011\ 1111\ 1001\ 010 \end{aligned}$$



# Suma de los dos números

Debido a que los exponentes son diferentes, no es correcto sumar las mantisas, y así el exponente del número más pequeño se hace más grande mientras que disminuye progresivamente la mantisa desplazando bits a la derecha (insertando ceros) hasta que ambos números tengan el mismo exponente.

$$10^{-7} = 0\ 01100000\ 1101\ 0110\ 1011\ 1111\ 1001\ 010$$

$$\begin{aligned}
 10^{-7} &= 0\ 01100000\ 1101\ 0110\ 1011\ 1111\ 1001\ 010 \\
 &= 0\ 01100001\ 0110\ 1011\ 0101\ 1111\ 1100\ 101\ (0)
 \end{aligned}$$

$$\begin{aligned}
 10^{-7} &= 0\ 01100000\ 1101\ 0110\ 1011\ 1111\ 1001\ 010 \\
 &= 0\ 01100001\ 0110\ 1011\ 0101\ 1111\ 1100\ 101\ (0) \\
 &= 0\ 01100010\ 0011\ 0101\ 1010\ 1111\ 1110\ 010\ (10)
 \end{aligned}$$

$$\begin{aligned}
10^{-7} &= 0\ 01100000\ 1101\ 0110\ 1011\ 1111\ 1001\ 010 \\
&= 0\ 01100001\ 0110\ 1011\ 0101\ 1111\ 1100\ 101\ (0) \\
&= 0\ 01100010\ 0011\ 0101\ 1010\ 1111\ 1110\ 010\ (10) \\
&\dots
\end{aligned}$$

$$\begin{aligned}
10^{-7} &= 0\ 01100000\ 1101\ 0110\ 1011\ 1111\ 1001\ 010 \\
&= 0\ 01100001\ 0110\ 1011\ 0101\ 1111\ 1100\ 101\ (0) \\
&= 0\ 01100010\ 0011\ 0101\ 1010\ 1111\ 1110\ 010\ (10) \\
&\dots \\
&= 0\ 10000010\ 0000\ 0000\ 0000\ 0000\ 00000\ 000\ (0001101\dots) \\
&\rightarrow 7 + 1.0 \times 10^{-7} = 7
\end{aligned}$$

# Dígitos significativos

Debido a que no hay espacio para almacenar los últimos dígitos, éstos se pierden, y después de todo este largo proceso, la suma sólo devuelve 7 como respuesta (tenemos un error de truncamiento)

# Dígitos significativos

Debido a que no hay espacio para almacenar los últimos dígitos, éstos se pierden, y después de todo este largo proceso, la suma sólo devuelve 7 como respuesta (tenemos un error de truncamiento)

En otras palabras, debido a que un equipo de 32 bits almacena sólo 6 o 7 decimales, ignora efectivamente cualquier cambio más allá del sexto decimal.



# Precisión de la máquina

La pérdida de precisión anterior se categoriza definiendo la precisión de la máquina  $\epsilon_m$  como el máximo número positivo que en la computadora, se puede sumar al número almacenado como 1 sin cambiarlo:

$$1_c + \epsilon_m = 1_c \quad (1)$$

Por consiguiente, un número arbitrario  $x$  puede considerarse relacionado con su representación en punto flotante  $x_c$  por

$$x_c = x(1 \pm \epsilon) \quad |\epsilon| \leq \epsilon_m$$

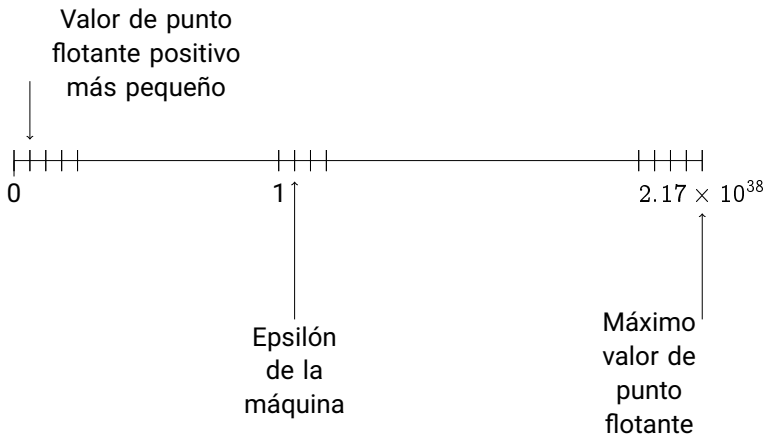
donde no se conoce el valor real de  $\epsilon$ .

En otras palabras, con excepción de las potencias de 2 que están representadas exactamente, debemos suponer que todos los números de precisión simple contienen un error en el sexto decimal y que todos los dobles tienen un error en el decimoquinto lugar.

Y como siempre ocurre con los errores, debemos suponer que no sabemos cuál es el error, porque si lo supiéramos: **¡entonces lo eliminaríamos!**

En consecuencia, los argumentos que planteamos con respecto a los errores son siempre aproximados, y eso es lo mejor que podemos hacer.

# El epsilon de la máquina



# Calculando el epsilon de la máquina

Una tarea inicial que debemos de atender, es el cálculo de la precisión de la máquina, a la que llamaremos **epsilon de la máquina**.

# Calculando el epsilon de la máquina

Una tarea inicial que debemos de atender, es el cálculo de la precisión de la máquina, a la que llamaremos **epsilon de la máquina**.

Veamos cómo calcular con `python` ese valor:

# Punto de partida

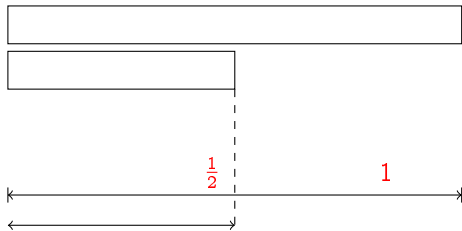
Tomamos el valor de 1





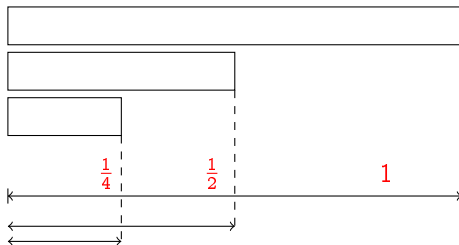
# Punto de partida

Dividimos a la mitad esta unidad



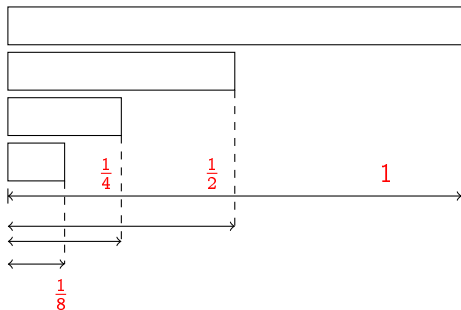
# Punto de partida

Volvemos a dividir a la mitad



# Punto de partida

Seguimos dividiendo a la mitad ...



# Calculemos el epsilon con python

```
1 t = 1.0;
2
3 while 1+ t != 1:
4     eps = t
5     t = t/2
6
7 print ('el epsilon de la maquina es: ',
        eps)
```

# Calculemos el epsilon con python

```
1 t = 1.0;
2
3 while 1+ t != 1:
4     eps = t
5     t = t/2
6
7 print ('el epsilon de la maquina es: ',
        eps)
```

En mi equipo obtengo el siguiente resultado:

El cero de la maquina es:

$2.220446049250313e - 16$

# 1. Precisión de la máquina

## 2. Tipos de errores

2.1 Error absoluto verdadero

2.2 Error relativo verdadero

2.3 Error relativo aproximado

## 3. Conceptos importantes

## 4. Algoritmo de Horner

# Nuevos conceptos

Una vez que se ha establecido la clasificación del error, definiremos los conceptos de:

- Error absoluto verdadero.
- Error relativo verdadero.
- Error relativo aproximado.

todos ellos como una suma o consecuencia de los errores de redondeo y truncamiento.

Supóngase que  $\hat{p}$  es una aproximación a  $p$ .

El error absoluto verdadero se define con la siguiente expresión:

$$E_v = |p - \hat{p}|$$

Esta definición de error, lo cuantifica en términos brutos.



No obstante, una medida que puede describir con mayor detalle o proporción el error, es aquella que lo expresa en términos porcentuales.

Para ello se emplea el error verdadero relativo.

Supóngase que  $\hat{p}$  es una aproximación a  $p$ . El error relativo verdadero se calcula con la siguiente expresión:

$$e_v = \frac{|p - \hat{p}|}{p}$$

El resultado suele expresarse en términos porcentuales.

# Error relativo aproximado

El error relativo aproximado, mide el error de un método numérico, determinando el error de la iteración actual respecto el error surgido en la iteración anterior:

$$e_a = \frac{|\hat{x}_i - \hat{x}_{i-1}|}{\hat{x}_i}$$

Donde  $x_i$  es la aproximación actual a  $x$  y  $x_{i-1}$  es la aproximación anterior a  $x$ .

En métodos numéricos suele establecerse una tolerancia porcentual como criterio de paro, tal que el error relativo aproximado de un método, no exceda dicha tolerancia.

$$e_a < t$$

donde  $t$ , es tolerancia fijada de antemano.

A menor tolerancia se tiene mayor precisión en la aproximación al valor verdadero, sin embargo esto implica un aumento en el número de iteraciones requeridas para detener el método.

# 1. Precisión de la máquina

## 2. Tipos de errores

## 3. Conceptos importantes

- 3.1 Contaminación en los cálculos

- 3.2 Condición

- 3.3 Estabilidad

- 3.4 Eficiencia

## 4. Algoritmo de Horner

# Contaminación en los cálculos.

Un error en un cálculo numérico “contamina” las sucesivas evaluaciones.

Esta propagación puede describirse en términos de dos conceptos relacionados: los de estabilidad y condición.

La condición de una función  $f(x)$  mide la sensibilidad de los valores de  $f(x)$  a pequeños cambios de  $x$ , se define como:

$$C = \left| \frac{E_{rel}(f(x))}{E_{rel}(x)} \right|$$



Del teorema del valor medio en cálculo, podemos expresar

$$\begin{aligned} f(x_T) - f(x_A) &\approx f'(x_t)(x_T - x_A) \rightarrow E_{rel}(f(x)) \approx \\ &\approx \frac{f'(x_T)}{f(x_T)}(x_T - x_A) \end{aligned}$$

luego

$$C \approx \left| x_T \frac{f'(x_T)}{f(x_T)} \right|$$

Se utilizará ésta definición como definición de condición para funciones  $f(x)$  de una variable real.

Entonces los números de condición serán

$$C = \left| x \frac{f'(x)}{f(x)} \right|$$

- 1 Para un  $x$  dado  $0 < C(x) < 1$  se dirá que el problema está bien condicionado, y cuanto menor sea  $C$ , mejor condicionado.

Se utilizará ésta definición como definición de condición para funciones  $f(x)$  de una variable real.

Entonces los números de condición serán

$$C = \left| x \frac{f'(x)}{f(x)} \right|$$

- 1 Para un  $x$  dado  $0 < C(x) < 1$  se dirá que el problema está bien condicionado, y cuanto menor sea  $C$ , mejor condicionado.
- 2 Si  $C(x) > 1$ , el problema estará mal condicionado.

Se utilizará ésta definición como definición de condición para funciones  $f(x)$  de una variable real.

Entonces los números de condición serán

$$C = \left| x \frac{f'(x)}{f(x)} \right|$$

- 1 Para un  $x$  dado  $0 < C(x) < 1$  se dirá que el problema está bien condicionado, y cuanto menor sea  $C$ , mejor condicionado.
- 2 Si  $C(x) > 1$ , el problema estará mal condicionado.
- 3 Si  $C(x) = 1$ , el error relativo se mantiene.

Las siguientes funciones están bien condicionadas?

❶  $f(x) = \sqrt{x}$        $C(x) = ?$

❷  $g(x) = x^2 - 1$        $C(x) = ?$

La estabilidad de un algoritmo describe la sensibilidad de un método numérico específico respecto a los inevitables errores de redondeo cometidos durante su ejecución en aritmética de precisión finita.

Consideremos la siguiente función:

$$f(x) = \sqrt{x+1} - \sqrt{x}$$

Su número de condición es:

$$C(x) = \left| x \frac{f'(x)}{f(x)} \right| = \frac{x}{2\sqrt{x}\sqrt{x+1}}$$

Vemos que  $C(x) < \frac{1}{2}$  para  $x > 0$ , por lo que la función está bien condicionada pero ...

El algoritmo para calcular  $x$  de tal forma que se vayan realizando las operaciones, es:

- 1 obtener  $x$



El algoritmo para calcular  $x$  de tal forma que se vayan realizando las operaciones, es:

① obtener  $x$

②  $y = x + 1$

El algoritmo para calcular  $x$  de tal forma que se vayan realizando las operaciones, es:

- 1 obtener  $x$
- 2  $y = x + 1$
- 3  $f = \text{sqrt}(y)$

El algoritmo para calcular  $x$  de tal forma que se vayan realizando las operaciones, es:

① obtener  $x$

②  $y = x + 1$

③  $f = \text{sqrt}(y)$

④  $g = \text{sqrt}(x)$

El algoritmo para calcular  $x$  de tal forma que se vayan realizando las operaciones, es:

① obtener  $x$

②  $y = x + 1$

③  $f = \text{sqrt}(y)$

④  $g = \text{sqrt}(x)$

⑤  $h = f - g$

El algoritmo para calcular  $x$  de tal forma que se vayan realizando las operaciones, es:

① obtener  $x$

②  $y = x + 1$

③  $f = \text{sqrt}(y)$

④  $g = \text{sqrt}(x)$

⑤  $h = f - g$

El algoritmo para calcular  $x$  de tal forma que se vayan realizando las operaciones, es:

- 1 obtener  $x$
- 2  $y = x + 1$
- 3  $f = \text{sqrt}(y)$
- 4  $g = \text{sqrt}(x)$
- 5  $h = f - g$

es inestable para  $x$  grandes, dado el paso 5, por lo que debemos de re-estructurar la función.

*Debemos evitar que todo algoritmo sea inestable.*  
Si existieran varios métodos para evaluar una misma función, entonces conviene utilizar aquel que sea más eficiente, es decir, más rápido.

*Debemos evitar que todo algoritmo sea inestable.*  
Si existieran varios métodos para evaluar una misma función, entonces conviene utilizar aquel que sea más eficiente, es decir, más rápido.

Hay que aprovechar al máximo los recursos: hardware, software, algoritmos, para resolver problemas más complejos y no para resolver peor problemas simples.



Por ejemplo, para calcular  $x ** 4$  para un  $x$  dado, no es buena idea calcular  $x ** 4.0$  (exponente en punto flotante).

La mejor idea consiste en economizar el cálculo en dos pasos:

$$x^2 = x * x$$

$$x^4 = x^2 * x^2$$

y no un producto  $x^4 = x * x * x * x$

# Ejemplo: Evaluación de polinomios

Supongamos que queremos evaluar el polinomio:

$$P(x) = 2 + 4x - 5x^2 + 2x^3 - 6x^4 + 8x^5 + 10x^6$$

Contando con que cada potencia de exponente  $k$  entero como  $k - 1$  productos, tendríamos que el total de productos para evaluar en forma directa es:

$$1 + 2 + 3 + 4 + 5 + 6 = 21$$

Además de seis sumas.

Una mejora en el algoritmo , es calcular primero las potencias de forma sucesiva:

$$x^2 = x * x$$

$$x^3 = x * x^2$$

$$x^4 = x * x^3$$

$$x^5 = x * x^4$$

$$x^6 = x * x^5$$

De tal forma que se añade un producto por cada potencia, para un total de productos:

$$1 + 2 + 2 + 2 + 2 + 2 = 11$$

Con el polinomio

$$P(x) = 2 + 4x - 5x^2 + 2x^3 - 6x^4 + 8x^5 + 10x^6$$

podemos mejorar el algoritmo de la siguiente manera

$$P(x) = 2 + x \{4 + x (-5 + x [2 + x (-6 + x \{8 + x * 10\})])\}$$

Para evaluar un polinomio de grado  $n$  en el que ninguno de los coeficientes es cero, se necesitan

Para evaluar un polinomio de grado  $n$  en el que ninguno de los coeficientes es cero, se necesitan

$\frac{n(n+1)}{2}$  Productos para el primer método

$2n - 1$  para el segundo métodos

$n$  para el tercero

Antes de escribir una línea de código, hay que revisar la manera en que podemos optimizar la solución del problema.

1. Precisión de la máquina
2. Tipos de errores
3. Conceptos importantes
4. Algoritmo de Horner
  - 4.1 Descripción del algoritmo
  - 4.2 Pseudocódigo
  - 4.3 Ejercicio
  - 4.4 Extendiendo la respuesta
  - 4.5 El código con `python`

# Algoritmo de Horner

Dado el polinomio

$$P(x) = a_0 + a_1x + \dots + a_nx^n \quad a_n \neq 0$$

La evaluación de  $P(x)$  para cierto valor de  $x = z$  se puede realizar en  $n$  pasos mediante:

$$\begin{aligned} b_n &= a_n \\ b_{n-1} &= a_{n-1} + z * b_n \\ b_{n-2} &= a_{n-2} + z * b_{n-1} \\ &\vdots \\ b_0 &= a_0 + z * b_1 \end{aligned}$$



# Pseudocódigo

$$\textcircled{1} \ b_n = a_n$$

# Pseudocódigo

- 1  $b_n = a_n$
- 2 repetir mientras  $n > 0$

# Pseudocódigo

- 1  $b_n = a_n$
- 2 repetir mientras  $n > 0$
- 3  $n = n - 1$

# Pseudocódigo

- 1  $b_n = a_n$
- 2 repetir mientras  $n > 0$
- 3  $n = n - 1$
- 4  $b = a_n + z * b$

# Pseudocódigo

- 1  $b_n = a_n$
- 2 repetir mientras  $n > 0$
- 3  $n = n - 1$
- 4  $b = a_n + z * b$
- 5 volver al paso 2

# Pseudocódigo

- 1  $b_n = a_n$
- 2 repetir mientras  $n > 0$
- 3  $n = n - 1$
- 4  $b = a_n + z * b$
- 5 volver al paso 2
- 6  $p(z) = b$

# Completa la siguiente tabla

$$P(x) = 2 + 4x - 5x^2 + 2x^3 - 6x^4 + 8x^5 + 10x^6$$

Evalúa el polinomio  $P(x)$  en:

$x$	$P(x)$
-1.5	
-0.65	
0.1	
1.4	
2.87	

# Resultado

$$P(x) = 2 + 4x - 5x^2 + 2x^3 - 6x^4 + 8x^5 + 10x^6$$

El polinomio  $P(x)$  evaluado es:

-1.5	0.781 25
-0.65	-4.506 83
0.1	2.351 49
1.4	98.559 68
2.87	6758.702 45



# Extendiendo la respuesta al problema

¿Cómo resolver el problema usando una función?  
¿mostrando una tabla con formato de salida?  
usando una gráfica que muestre  $P(x)$  y un  
conjunto de datos evaluados? ¿Evaluar el error  
relativo?

# Extendiendo la respuesta al problema

Ya contamos con las herramientas necesarias para extender la respuesta al problema, en nuestro código podemos agregar funciones, ajustar formatos de salida en los resultados, graficar el polinomio y los puntos (o un conjunto diferente de puntos), y obtener el error relativo. Basta con que le dediquemos un poco más de tiempo.

# Pasos a resolver

- 1 Conviene definir una función que resuelva la evaluación del método de Horner.

# Pasos a resolver

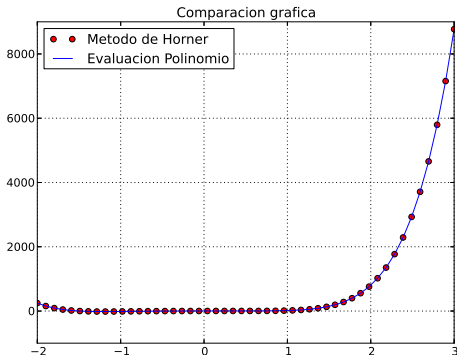
- 1 Conviene definir una función que resuelva la evaluación del método de Horner.
- 2 Para obtener el error relativo, se debe de evaluar el polinomio y considerar que los valores obtenidos, son los valores exactos.

# Pasos a resolver

- 1 Conviene definir una función que resuelva la evaluación del método de Horner.
- 2 Para obtener el error relativo, se debe de evaluar el polinomio y considerar que los valores obtenidos, son los valores exactos.
- 3 Comparamos los resultados mediante una gráfica que represente los dos resultados de la evaluación.

# Resultado gráfico más completo

En la gráfica se muestra un conjunto de datos que se evalúan y posteriormente con la función (en línea continua) se compara, podemos ver que los resultados son prácticamente los mismo.



# Estructuramos el código en python .

Definimos mediante dos listas:

- 1 Los valores donde queremos evaluar el polinomio  $P(x)$ .

```
1 # Valores de x0 para evaluar P(x0)
2 x0=[-1.5, -0.65, 0.1, 1.4, 2.87]
3
4 # Coeficientes a de P(x)
5 a=[2, 4, -5, 2, -6, 8, 10]
```

# Estructuramos el código en python .

Definimos mediante dos listas:

- 1 Los valores donde queremos evaluar el polinomio  $P(x)$ .
- 2 Los coeficientes de  $P(x)$ .

```
1 # Valores de x0 para evaluar P(x0)
2 x0=[-1.5, -0.65, 0.1, 1.4, 2.87]
3
4 # Coeficientes a de P(x)
5 a=[2, 4, -5, 2, -6, 8, 10]
```



# Definimos una función que resuelva por Horner.

```
1 # Metodo de Horner
2
3 def P_Horner(x):
4     P_Hor=0
5     for n in range(len(a)-1,-1,-1):
6         P_Hor=a[n]+P_Hor*x
7     return P_Hor
```

Definimos una función que evalúe directamente  $P(x)$ .

```
1 # Evaluacion directa
2
3 def P_Directo(x):
4     return 2+4*x-5*x**2+2*x**3-6*x**4+8*x**5+10*x**6
```

# Calculamos el error relativo.

```
1 # Calculo de error relativo
2
3 def Err_Rel(p,p_): return (p-p_)/p*100
```

# Mostramos el error relativo de los puntos a evaluar.

```
1 # Evaluacion de valores de P(x0)
2
3 for i in range(len(x0)):
4     print ("P(%.2f) =" %x0[i], P_Horner(x0[i]
        ), "; Error rel. =", Err_Rel(
            P_Directo(x0[i]), P_Horner(x0[i])))
```

# Comparamos los resultados con una gráfica.

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3
4 x=np.linspace(-2.,3.)
5
6 plt.plot(x,P_Horner(x),'ro', label='Metodo de
   Horner')
7
8 plt.plot(x,P_Directo(x), label='Evaluacion
   Polinomio')
9
10 plt.title('Comparacion grafica')
11 plt.legend(loc='upper left')
12
13 plt.grid(True)
14 plt.show()
```