

*Programador Universitario*  
*Licenciatura en Informática*  
*Ingeniería en Informática*

## **METODOS NUMERICOS I**

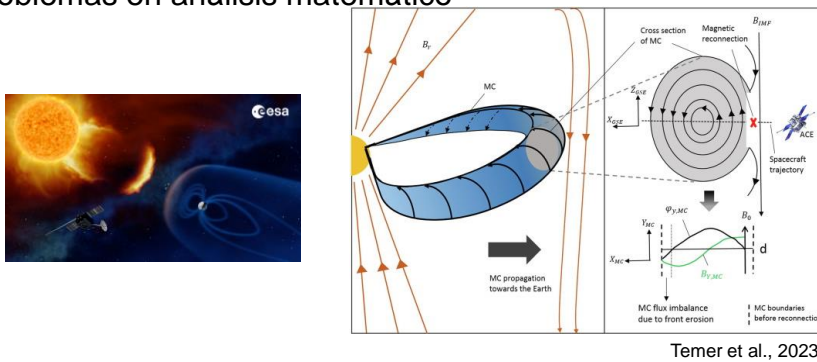
### **(P10)**

Métodos Numéricos I

1

### **CALCULO NUMERICO**

Es la rama de las matemáticas encargada de diseñar algoritmos para simular aproximaciones de solución a problemas en análisis matemático



Temer et al., 2023

Métodos Numéricos I

## **METODOS NUMERICOS I**

### **TEMAS**

- Unidad I:** Teoría de errores
- Unidad II:** Solución de Ecuaciones no Lineales
- Unidad III:** Solución Numérica de Sistemas Lineales
- Unidad IV:** Interpolación
- Unidad V:** Integración numérica

Métodos Numéricos I

### **Unidad I - Teoría de errores**

- Conceptos básicos
- Fuentes de error
- Error de truncamiento
- Representación en punto flotante.
- Error de representación
- Aritmética de números reales. Propagación del error
- Métodos de estimación del error

Métodos Numéricos I

## **Unidad I - *Teoría de errores***

Porque nos interesa como programadores?

### **RESOLUCION DE UN PROBLEMA NUMERICAMENTE**

- i) Formulación del problema
- ii) Elección del método
- iii) Programación y codificación
- iv) Análisis de los resultados

## **Formulación del problema**

El proceso de análisis del mundo real para interpretar los aspectos de un problema y expresarlo en términos precisos se denomina **abstracción**.

Abstraer un problema del mundo real y simplificar su expresión, tratando de encontrar los aspectos principales que se pueden resolver, los datos que se van a procesar y el contexto del problema se denomina **modelización** (mental)

Formular un **Modelo Matemático** es:

Definir las ecuaciones matemáticas del modelo y los parámetros que intervienen

Métodos Numéricos I

Un **modelo matemático** es una ecuación que muestra las características de un sistema físico:

$$y = f(x, \text{parámetros}, \text{fn.fuerza})$$

$y$ : vble. dep., muestra el comportamiento o estado del sistema

$x$ : vbles. indep., determinan el comportamiento del sistema (tiempo, espacio)

*parámetros*: constantes que muestran la composición o propiedades del sistema

*fn. fuerza*: influencias externas sobre el sistema

Métodos Numéricos I

### Ejemplos de Modelos Matemáticos

- Distancia recorrida por un objeto
- Velocidad de descenso de un objeto
- Concentración de una sustancia en una reacción química
- Variación en la demanda de un producto
- Evolución de una enfermedad según tratamiento
- Dinámica de poblaciones
- Desarrollo económico
- Movimientos sísmicos

Etc..

Métodos Numéricos I

-el porcentaje de riesgo  $R$  de tener un accidente al manejar un auto, en función de la cantidad de alcohol  $x$  en sangre:

$$R(x) = 6(1,013)^x$$

Diagram illustrating the components of the equation  $R(x) = 6(1,013)^x$ :

- Parámetros** (Parameters) points to the constants  $6$  and  $1,013$ .
- Vble dep** (Dependent variable) points to  $x$ .
- vble ind.** (Independent variable) points to  $x$ .

Si  $R = 20\%$  no deben conducir vehículos (por ley)  
 O sea que  $x = 95$ , es el mayor valor permitido

Métodos Numéricos I

## CLASIFICACION DE LOS ERRORES

- i) Errores del modelo o del problema
  - ii) Errores de cálculo y programación
  - iii) Errores en los datos
  - iii) Errores de truncamiento
  - iv) Errores de redondeo

Métodos Numéricos I

### *i) Errores del modelo o del problema*

Son los **más graves**, si modelamos mal, por ejemplo hacemos hipótesis incorrectas o demasiadas simplificaciones para llegar al modelo matemático, entonces **no resolvemos** el problema real planteado

Ej.; Ley de Malthus para dinámica poblacional

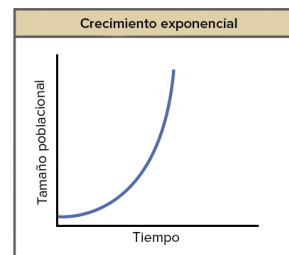
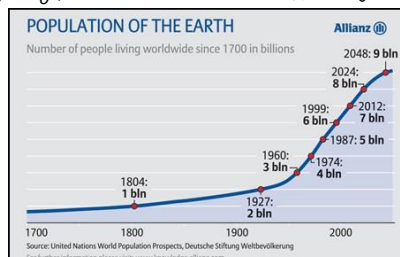
$$\frac{dP(t)}{dt} = rP(t)$$

población en el año t  
Parámetro

Dado  $P(0)=P_0$ , la solución es:  $P(t) = P_0 e^{rt}$

Pero

...



Métodos Numéricos I

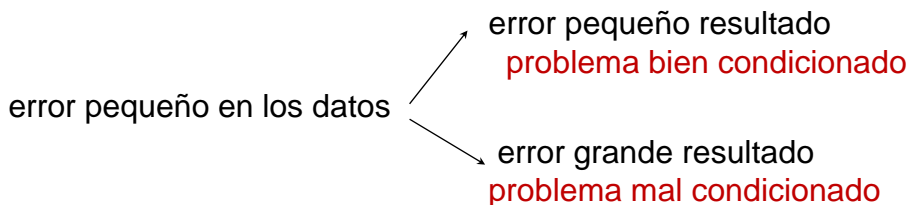


## ii) Errores de cálculo y programación

Causados por el **factor humano**, antes errores de calculo, ahora los de programación, para evitarlos es fundamental probar adecuadamente los programas

## lii) Errores en los datos

En un problema físico los **datos** son **empíricos**, por ello tienen **errores experimentales**, asociados al instrumento de medición. El error en los resultados no puede ser menor que el error en los datos. Relacionado a estos errores esta el concepto de **condicionamiento**



problema bien  
condicionado

$$\begin{cases} 3x - y = 1 \\ -x + 3y = 13 \end{cases} \rightarrow \begin{matrix} X= \\ 2 \\ y=5 \end{matrix}$$

$$\begin{cases} 3,1x - 0,9y = 1 \\ -1,1x + 2,9y = 13 \end{cases} \rightarrow \begin{matrix} X=1,825 \\ Y=5,175 \end{matrix}$$

problema mal condicionado

$$\begin{cases} 1x + 2y = 10 \\ 1,05x + 2y = 10,4 \end{cases} \rightarrow \begin{matrix} X= \\ 8 \\ y=1 \end{matrix}$$

$$\begin{cases} 1x + 2y = 10 \\ 1,1x + 2y = 10,4 \end{cases} \rightarrow \begin{matrix} X= \\ 4 \\ y=3 \end{matrix}$$



## Errores Sistemáticos

Son los que tienen siempre *aproximadamente el mismo tamaño y signo*, es decir que el error tiene una causa constante, son siempre por exceso o por defecto

## Errores Accidentales

Son los relacionados a *factores aleatorios* en la medición

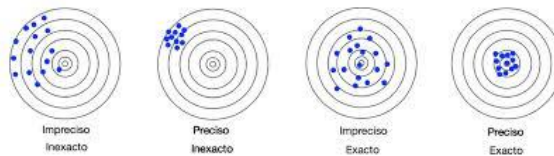


## Precisión y Exactitud

Una medición es **exacta** cuando su valor es muy cercano al valor verdadero.

Una medición es **precisa** cuando el instrumento que empleamos nos proporciona *muchas cifras decimales*, pero la medida puede no ser exacta, por ej. si estamos cometiendo un error sistemático al usar el instrumento.

Una medida para ser exacta debe ser precisa, pero no todas las medidas precisas son exactas.



Métodos Numéricos I

Los errores vistos en *i)* , *ii)* y *iii)* son independientes del instrumento con el que se resuelva el problema.

Una diferencia fundamental entre el tratamiento matemático de un problema y el numérico está en las limitaciones que tiene el cálculo numérico, tanto para representar procesos como en la representación de los números. Esto genera dos tipos de errores:

**Errores de truncamiento y Errores de redondeo**

## ERROR de TRUNCAMIENTO

Se origina al substituir procesos infinitos por procesos finitos

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$$

$$\int_a^b f(x)dx \sim \frac{h}{2} \sum_{i=1}^n (f(x_{i-1}) + f(x_i))$$

$$\frac{dy}{dx} \cong \frac{\Delta y}{\Delta x}$$

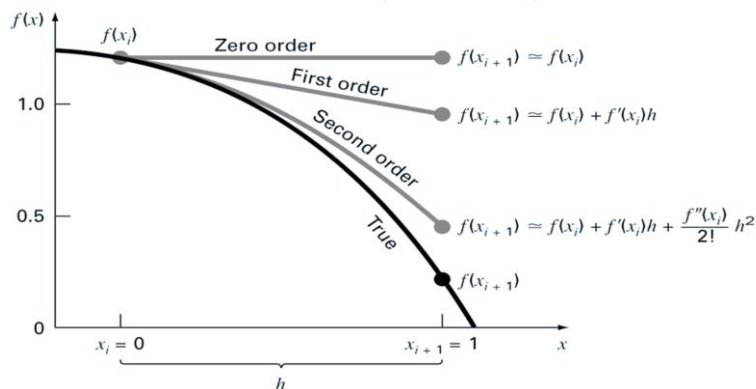
Los errores de truncamiento causan inexactitud de los resultados.

Métodos Numéricos I

## SERIE DE TAYLOR

El teorema de Taylor nos dice que toda función suave se puede aproximar con un polinomio. Esto se hace usando la Serie de Taylor

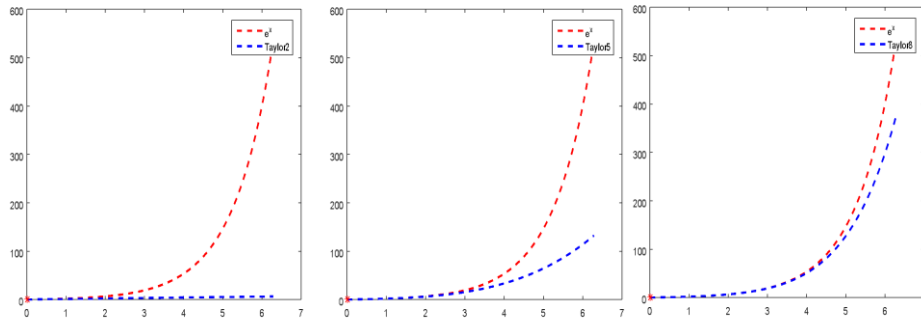
$$f(x_i + h) = f(x_i) + hf'(x_i) + \frac{h^2}{2!} f''(x_i) + \frac{h^3}{3!} f'''(x_i) + \dots + R_n.$$



Métodos Numéricos I

## Desarrollo en serie de Taylor

**Ejemplo:  $f(x) = e^x$**



Métodos Numéricos I

## ERROR

Sean  $x$  y  $\bar{x}$  un número y su valor aproximado

$$e = x - \bar{x} \quad \text{error absoluto}$$

$$e_r = \frac{x - \bar{x}}{x} \quad \text{error relativo si } x \neq 0 \quad (\text{error por unidad de medida})$$

$$e_r \% = e_r \cdot 100 \quad \text{error relativo porcentual}$$

$$|x - \bar{x}| \leq \varepsilon \quad \text{cota de error} \quad \Longrightarrow \quad x = \bar{x} \pm \varepsilon$$

Métodos Numéricos I

22

*Ejemplo : Cota de error*

$$|x - \bar{x}| \leq \varepsilon \quad \longrightarrow \quad x = \bar{x} \pm \varepsilon$$

Se quiere estimar  $11 / 3$

$$x_{real} = 11 / 3 = 3.666666....$$

$$x_{aprox} = 3.666$$

$$Ea = |3.666666... - 3.666| = 0.000666... < 0.001$$

El error absoluto  $\varepsilon$  es menor que una milésima

## **ERROR**

Sean  $x$  y  $\bar{x} \in \mathbb{R}^n$

$$e = \|x - \bar{x}\| \quad \text{error absoluto}$$

$$e_r = \frac{\|x - \bar{x}\|}{\|x\|} \quad \text{error relativo, si } x \neq 0$$

Si trabajamos con vectores o matrices, se realiza el cálculo del error utilizando el concepto de norma

## ERROR DE REDONDEO

Se debe a que una máquina sólo puede representar cantidades con un número finito de dígitos.

Podemos distinguir dos tipos:

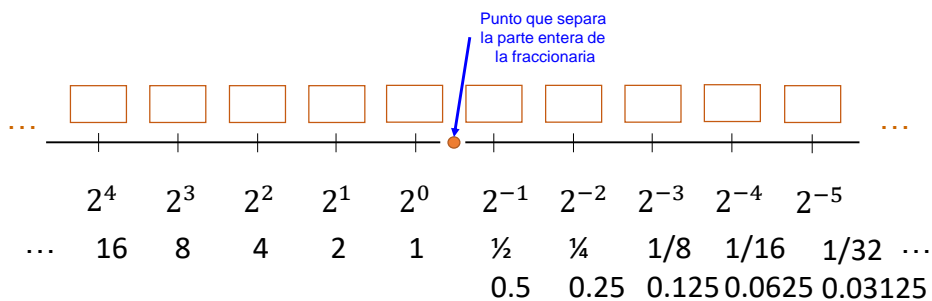
- a) error de representación
- b) error debido a los cálculos

Métodos Numéricos I

25

## SISTEMA DE NUMERACIÓN POSICIONAL

En un sistema de numeración posicional, el valor de cada símbolo o cifra depende tanto del **símbolo** como de su **posición**



Métodos Numéricos I

26

## SISTEMA DE NUMERACIÓN POSICIONAL

**Ejemplo:** Convertir  $(25)_{10}$  a binario

25		Resultado	Resto
25 / 2	=	12	1
12 / 2	=	6	0
6 / 2	=	3	0
3 / 2	=	1	1
1 / 2	=	0	1



**Resultado:**  $(25)_{10} = (11001)_2$

Métodos Numéricos I

27

## SISTEMA DE NUMERACIÓN POSICIONAL

**Ejemplo:** Convertir  $(0.71875)_{10}$  a binario

0.71875		Parte entera
$0.71875 * 2 = 1.4375$		1
$0.43750 * 2 = 0.8750$		0
$0.87500 * 2 = 1.7500$		1
$0.75000 * 2 = 1.5000$		1
$0.50000 * 2 = 1.0000$		1



**Resulta:**  $(0.71875)_{10} = (0.10111)_2$

$(0.71875)_{10}$  tiene  
representación exacta en el  
sistema binario

**Ejemplo:** Convertir  $(0.6)_{10}$  a binario

0.6		Parte entera
$0.6000 * 2 = 1.2000$		1
$0.2000 * 2 = 0.4000$		0
$0.4000 * 2 = 0.8000$		0
$0.8000 * 2 = 1.6000$		1
$0.6000 * 2 = 1.2000$		1
$0.2000 * 2 = 0.4000$		0



**Resulta:**  $(0.6)_{10} = (0.100110011 \dots)_2$

$(0.6)_{10}$  **NO** tiene  
representación exacta en el  
sistema binario

Métodos Numéricos I

28

## **REPRESENTACION INTERNA**

Existen dos maneras de representar los números:

**1. Punto fijo:** Los números se representan con un número fijo de cifras decimales. 6.358, 0.013  
(introducida inicio de los '80 ; la mayoría de los chips DSP de bajo costo la usan pues no requiere FPU)

**2. Punto flotante:** Los números se representan con un número fijo de dígitos significativas 0.636E01, 0.135E-01

**Dígito Significativo:** Dado un número  $x$ , es cualquier dígito, excepto los ceros a la izquierda del primer dígito diferente de cero y que solo sirven para fijar la posición del punto decimal Ej. **1360**, **1.360**; **0.001360**; tienen cuatro dígitos significativos.

Métodos Numéricos I

29

## **REPRESENTACION EN PUNTO FLOTANTE**

Sea un número  $x$ , representado en punto flotante en una base  $b$

$$x = (\text{sign } x) (.a_1 a_2 a_3 \dots a_t)_b \times b^e$$

$a_i$  : dígitos en el sistema de base  $b$  ,  $a_1 \neq 0$  o  $a_i = 0 \quad \forall i$

$t$ : número de dígitos de la mantisa (determina la **precisión**)

$.a_1 a_2 a_3 \dots a_t$  : mantisa normalizada, por ser  $a_1 > 0$

$e$ : exponente o característica (determina el **rango**)

Métodos Numéricos I

30

Ej:  $x = -9.25_{10}$   $m = 8$   $e = 4$

Representación en binario:

$$\begin{array}{l} 9 = 1001 \\ 0.25 = 0.01 \end{array} \longrightarrow 9.25_{10} = 1001.01_2$$

Normalizamos:  $0.100101 \times 10^{0100}$

Representación en punto flotante:

$$\begin{array}{ccc} 1 & 0100 & 10010100 \\ \text{sgn} & e & m \end{array}$$

Métodos Numéricos I

31

Ej:  $x = 11.125_{10}$   $m = 8$   $e = 4$

Representación en binario:

$$\begin{array}{l} 11 = 1011 \\ 0.125 = 0.001 \end{array} \longrightarrow 11.125_{10} = 1011.001_2$$

Normalizamos:  $0.1011001 \times 10^{0100}$

Representación en punto flotante:

$$\begin{array}{ccc} 0 & 0100 & 10110010 \\ \text{sgn} & e & m \end{array}$$

Métodos Numéricos I

32

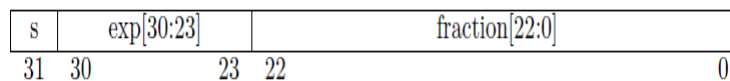


Estándar **IEEE 754** (85') se estableció para facilitar la portabilidad de los programas de un procesador a otro.

Define el formato para precisión simple de 32 bits y para precisión doble de 64 bits.

## REPRESENTACION EN PUNTO FLOTANTE

### IEEE Standard 754

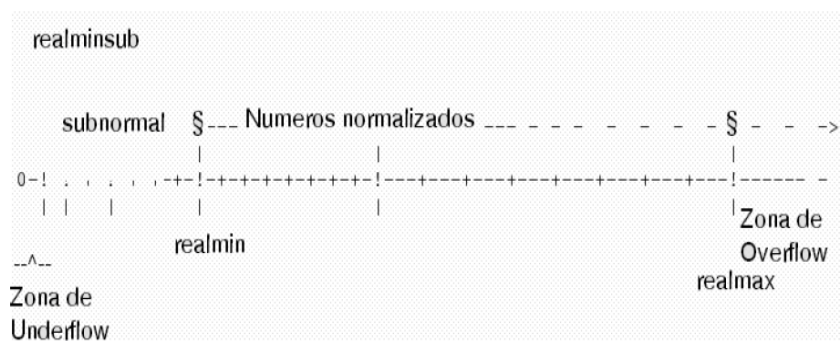


32 bits {

- 1 bit para signo número
- 8 bits para exponente ( $10^{-38}$ ,  $10^{38}$ )
- 23 bits para mantisa ( ~7 dígitos )

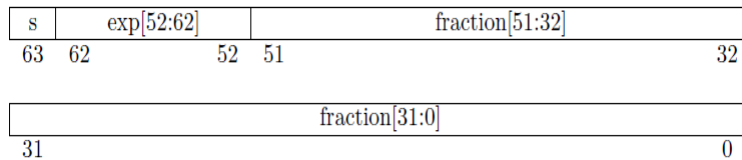
**Overflow:** Resultado del cálculo mayor que el número más grande que se puede representar

**Underflow:** Resultado del cálculo menor que el número más pequeño (no nulo) que se puede representar. Se considera 0 al valor



## **REPRESENTACION EN PUNTO FLOTANTE**

### **Doble Precision**



$$64 \text{ bits} \left\{ \begin{array}{l} 1 \text{ bit para signo número} \\ 11 \text{ bits para exponente ( } 10^{-308}, 10^{308} \text{)} \\ 53 \text{ bits para mantisa ( } \sim 15 \text{ dig.)} \end{array} \right.$$

Métodos Numéricos I

37

### **ERROR DE REPRESENTACION**

Como la mantisa contiene n dígitos en la base b, todo número más largo debe cortarse

Ej:

$$7/3 = 2.3333333333\dots$$

$$1/6 = 0.1666666666\dots$$

$$\pi = 3.141596265358\dots$$

También puede ocurrir que haya números con representación exacta en una base pero no en otra

$$\text{Ej. } 4/5 = (0.8)_{10}, 3/5 = (0.6)_{10}, 1.6_{10}$$

Métodos Numéricos I

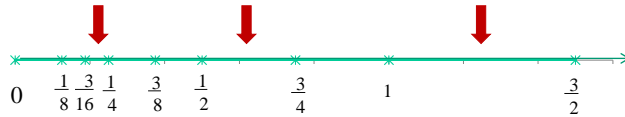
38

## ERROR DE REPRESENTACION

Los valores de **b**, **t** y **e** determinan que valores reales se pueden representar exactamente en una computadora

Ej : **b** = 2 , **t** = 2 , **e** = 2

Observamos que no podemos representar todos los números reales, la distancia entre los números representados (gap) es proporcional a la magnitud del número



Métodos Numéricos I

39

## ERROR DE REPRESENTACION

Un número  $x$  que no tiene representación exacta se denomina **fl(x)**.

Los dos métodos mas comunes para determinar la mantisa son :

- a) **Redondeo**: Se suma  $0.5 \times 10^{n-(k+1)}$  y luego se corta los dígitos  $k + 1$  en adelante.
- b) **Corte**: Cortar los dígitos  $k + 1$  en adelante

Esto es equivalente a la conocida “regla” para redondear:

- a) **Redondeo**: se elige como  $fl(x)$  el número de punto flotante normalizado más cercano a  $x$
- b) **Corte**: se elige como  $fl(x)$  el número de punto flotante normalizado más cercano entre 0 y  $x$

Métodos Numéricos I

40

## ERROR DE REPRESENTACIÓN

### Ejemplo:

Se desea almacenar el número irracional  $\pi = 3.14159265 \dots$  usando 5 dígitos para la mantisa.

$$\pi = 0.314159265 \dots \times 10^1$$

**Corte:**  $fl(\pi)_{CORTE} = 0.31415 \times 10^1$

**Redondeo:**  $fl(\pi)_{REDONDEO} = 0.314159265 \dots \times 10^1 + 5.0 \times 10^{1-(5+1)} =$

$$+ \begin{array}{r} 0.314159265 \dots \times 10^1 \\ 0.000005000 \dots \times 10^1 \\ \hline \end{array}$$

$$0.314164265 \dots \times 10^1$$

$$fl(\pi)_{REDONDEO} = 0.31416 \times 10^1$$

Métodos Numéricos I

41

## ERROR DE REPRESENTACION

Se define  $fl(x) = x(1+\varepsilon)$

$\varepsilon$  : error de redondeo,

Si redondeamos:

$$|\varepsilon| < \frac{1}{2}b^{1-t}$$

Si cortamos:

$$|\varepsilon| \leq b^{1-t}$$

Métodos Numéricos I

42

Ejemplo: calcular los errores absoluto y relativo en los casos planteados a continuación:

$x$	$x^*$	$e$	$er$
$0.3000 \cdot 10^1$	$0.3100 \cdot 10^1$		
$0.3000 \cdot 10^{-3}$	$0.3100 \cdot 10^{-3}$		
$0.3000 \cdot 10^4$	$0.3100 \cdot 10^4$		

Qué se concluye de estos resultados?

### **OPERACIONES ARITMETICAS EN PUNTO FLOTANTE**

Ej:  $b = 10, t = 3$

$$X = 0.164 \cdot 10^3, \quad y = 0.280 \cdot 10^3$$

$$Z = x+y = 0.167 \cdot 10^3$$

Operaciones en punto flotante:  $\hat{+}, \hat{=}, \hat{*}, \hat{/}$

$$x \hat{op} y = fl(x op y) = (x op y)(1+\epsilon)$$

### OPERACIONES ARITMETICAS EN PUNTO FLOTANTE

- Si sumamos dos números de magnitud muy diferente, puede ocurrir que el mas chico no se considere

Ej, si  $m = 5$  y tenemos los valores:

$$x = 8647300; y = 12 \longrightarrow \text{fl}(x) = 0.86473 \times 10^7; \text{fl}(y) = 0.12 \times 10^2$$

es decir que  $x$  e  $y$  se pueden representar exactamente.

Si calculamos

$$x + y = 8647312 ,$$

$$\text{fl}(x+y) = 0.86473 \times 10^7 = \text{fl}(x) = x$$

Es decir que si en un calculo se suman primero los términos más pequeños se pierden menos cifras significativas que si se empieza sumando los términos de mayor valor

**Resolver:** si  $b=10$  ,  $m = 3$

$$\text{i) } (121 - 0.327) - 119 \quad \text{ii) } (121 - 119) - 0.327$$

Métodos Numéricos I

45

### OPERACIONES ARITMETICAS EN PUNTO FLOTANTE

-Comparación de números de punto flotante

Nunca se deben comparar con el operador de igualdad directamente:  $a = b$  sino  $\text{abs}(a-b) < \text{eps}$

O si es un numero muy pequeño, no con 0, sino con una cota  $\text{abs}(x) < \text{eps}$

Métodos Numéricos I

46

## OPERACIONES ARITMETICAS EN PUNTO FLOTANTE

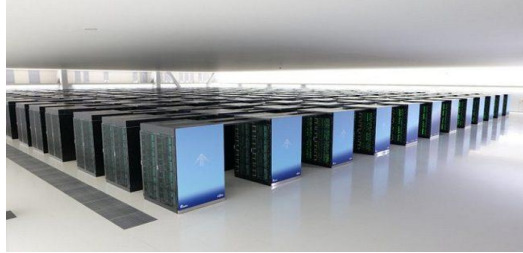
### FLOPS

(Floating point Operations Per Second)

Top500.org

2024

Lawrence Livermore National  
Laboratory - California



Rank	System	Vendor	Total Cores	Rmax (PFlops)	Rpeak (PFlops)	Power (kW)
1	El Capitan	HPE	11,039,616	1,742.00	2,746.38	29,580
	EX255a, AMD 4th Gen EPYC 24C 1.8GHz, AMD Instinct MI300A, Slingshot-11, TOSS					

Métodos Numéricos I

47

## Error de Representacion

**Unidad de redondeo:**

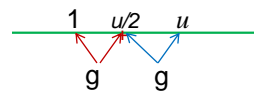
Se define como el menor valor  $u$  tal que

$$1 + u > 1$$

Esto significa que no puede representarse ningún numero entre 1 y  $1+u$

Dado un numero  $g / 1 + g$ ,

- $0 < g < u/2$ , se redondea a 1,
- si  $u/2 < g < 1$  se redondea a  $1+u$



Ej: utilizando Phyton

```
In [8]: print(np.finfo(np.float32).eps)
1.19209e-07

In [9]: print(np.finfo(float).eps)
2.22044604925e-16
```

Métodos Numéricos I

48



## ERROR DE REDONDEO ACUMULADO

Es la suma de todos los errores efectuados durante el cálculo

**Ej:**  $b = 10$ ,  $t = 8$

$$a = 0.23371258 \times 10^{-4}$$

$$b = 0.33678429 \times 10^2$$

$$c = -0.33677811 \times 10^2$$

Calculamos  $a \hat{+} b \hat{+} c$  de dos formas diferentes:

i)  $a \hat{+} (b \hat{+} c) = 0.64137126 \times 10^{-3}$

ii)  $(a \hat{+} b) \hat{+} c = 0.64100000 \times 10^{-3}$

Resultado exacto:  $a + b + c = 0.641371258 \times 10^{-3}$

Métodos Numéricos I

49

## ERROR DE REDONDEO ACUMULADO

Cada operación de punto flotante que realiza la computadora tiene asociado un error de redondeo; este error se acumula.

Analicemos el ej. anterior: dado  $a$ ,  $b$ ,  $c$  calculamos  $z = a + b + c$

$$a = 0.23371258 \times 10^{-4}$$

$$b = 0.33678429 \times 10^2$$

$$c = -0.33677811 \times 10^2$$

Se realiza el cálculo de dos maneras

i.  $a \hat{+} (b \hat{+} c)$

ii)  $(a \hat{+} b) \hat{+} c$

$$\begin{array}{rcl}
 & a \hat{+} (b \hat{+} c) & (a \hat{+} b) \hat{+} c \\
 \begin{array}{r}
 0.33678429 \times 10^2 \\
 - 0.33677811 \times 10^2 \\
 \hline
 0.00000618 \times 10^2 \\
 + 0.02337126 \times 10^{-3} \\
 \hline
 0.61800000 \times 10^{-3} \\
 \hline
 0.64137126 \times 10^{-3}
 \end{array}
 & \left. \begin{array}{l} \\ \\ \\ \\ \end{array} \right\} (b \hat{+} c) & + \left. \begin{array}{l}
 0.00000023371258 \times 10^2 \\
 0.33678429000000 \times 10^2 \\
 \hline
 0.33678452371258 \times 10^2 \\
 - 0.33677811 \times 10^2 \\
 \hline
 0.00000641 \times 10^2 \equiv 0.641 \times 10^{-3}
 \end{array} \right\} (a \hat{+} b) \\
 & \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} a \hat{+} (b \hat{+} c) & - \left. \begin{array}{l} \\ \\ \end{array} \right\} (a \hat{+} b) \hat{+} c
 \end{array}$$

Si  $a + b + c = 0.641371258 \times 10^{-3}$  Qué concluye ?

Métodos Numéricos I

50

## **ERROR DE REDONDEO ACUMULADO**

### En resumen:

Los resultados de las operaciones en la computadora tendrán en general errores debido a los errores de los operandos y al redondeo o truncamiento que ocurre al efectuar estas operaciones.

Errores de redondeo invalidan leyes básicas de la aritmética tal como la ley asociativa

$$(x + y) + z \neq x + (y + z).$$

Si en un método los errores crecen mucho hablamos de método **mal condicionado o inestable**

Ej: Cálculo de  $e^1$  utilizando la serie:

$$e^x = \sum_{i=0}^n \frac{x^i}{i!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \dots + \frac{x^n}{n!}$$

Si consideramos una precisión de 5 cifras

Término	Valor	Error
1	1	1.71828
2	2	0.71828
3	2.5	0.21828
4	2.66667	0.05161
5	2.70833	0.00995
6	2.71667	0.00161
7	2.71806	0.00022
.	.	.

## ESTABILIDAD

Si pequeños cambios en los datos producen pequeños cambios en los resultados diremos que es **estable**, caso contrario es **inestable**, en algunos casos la estabilidad depende del conjunto de datos, entonces se dice:  
**condicionalmente estable**

**Inestabilidad inherente** es aquella propia del problema o sistema.

**Inestabilidad inducida** es la que se produce por usar un método equivocado para resolver un determinado problema.

Ej: Polinomio de Wilkinson

$$p_n(x) = (x - 1)(x - 2) \dots (x - n) = \prod_{i=1}^{20} (x - i) = x^{20} - 210 x^{19} + \dots$$

El valor de  $a_{19} = -210$  si le restamos  $2^{-23}$ ;  $a_{19} = -210.0000001192$

Las nuevas raíces son:

1.00000	2.00000	3.00000	4.00000	5.00000
6.00001	6.99970	8.00727	8.91725	20.84691
10.09527 ± 0.64350i	11.79363 ± 1.65233i	13.99236 ± 2.51883i	16.73074 ± 2.81262i	19.50244 ± 1.94033i

## **Inestabilidad inherente**

### **ERROR DE SIGNIFICACIÓN**

Se dice que el número  $p^*$  aproxima a  $p$  hasta  $t$  dígitos significativos, si  $t$  es el mayor entero no negativo para el cual :

$$\frac{|p - p^*|}{|p|} \leq 0.5 \times 10^{1-t}$$

#### **Ejemplo:**

Sea  $p = \sqrt{3}$ . ¿En cuantos dígitos significativos aproxima  $p^* = 1.7$  a  $p$  y  $p^* = 1.73$  ?

Si se toma  $p^* = 1.7$ :

$t = 2$  es el mayor entero no negativo que verifica la desigualdad:

$$\frac{|\sqrt{3} - 1.7|}{\sqrt{3}} = 0.0185 \leq 0.5 \times 10^{1-2}$$

Si se toma  $p^* = 1.73$ :

$t = 3$  es el mayor el entero no negativo que verifica la desigualdad:

$$\frac{|\sqrt{3} - 1.73|}{\sqrt{3}} = 0.001184 \leq 0.5 \times 10^{1-3}$$

Métodos Numéricos I

55

### **ERROR DE SIGNIFICACIÓN**

Este error se produce cuando al operar se produce una pérdida de dígitos significativos

#### **Ejemplo:**

Sea  $x_1 = 0.84456 \cdot 10^0$  e  $y_1 = 0.84444 \cdot 10^0$  aproximaciones de  $x$  e  $y$ , en 4 cifras significativas. Queremos calcular  $z = x - y$

$$\begin{aligned} \text{Calculamos : } z_1 &= x_1 - y_1 \\ z_1 &= 0.82457 \cdot 10^0 - 0.82444 \cdot 10^0 \\ z_1 &= 0.00013 \cdot 10^0 \\ z_1 &= 0.13000 \cdot 10^{-3} \end{aligned}$$

¿En cuantos dígitos significativos aproxima  $z_1$  a  $z$  ?

Solo en 1 dígito pues el 2 proviene de la 5ta cifra, que estaba afectada de error. Esto ocurre cuando se restan cantidades muy próximas entre si

Métodos Numéricos I

56

## **PROPAGACION DE ERRORES**

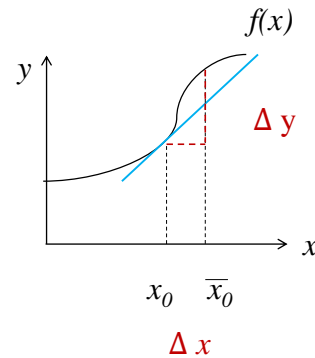
Cuando usamos métodos numéricos , el error del resultado será la suma de los errores en el desarrollo del mismo

### ***Funciones de una variable***

$$y = f(x)$$

$$\Delta y = y(x_0) - y(\bar{x}_0)$$

$$\Delta y \approx y'(\bar{x}_0) \Delta x$$



Métodos Numéricos I

57

## **FORMULA GRAL DE PROPAGACION DEL ERROR**

Supongamos conocer  $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n \sim x_1, x_2, \dots, x_n$  y sea  $y$  función de estas vbles.

Sea  $\bar{x} = (\bar{x}_1, \bar{x}_2, \dots, \bar{x}_n)$ ,  $x = (x_1, x_2, \dots, x_n)$

$$\Delta x_i = \bar{x}_i - x_i, \quad \Delta y = y(\bar{x}) - y(x)$$

Veamos el siguiente Teorema: Dados  $\bar{x}$ ,  $x$  e  $y(x)$  entonces,

$$\Delta y = \sum_{i=1}^n \frac{\partial y}{\partial x_i}(\bar{x}) \Delta x_i, \quad \text{por lo tanto}$$

$$|\Delta y| \leq \sum_{i=1}^n \left| \frac{\partial y}{\partial x_i}(\bar{x}) \right| |\Delta x_i|$$

Métodos Numéricos I

58

Ej: Calcular el error absoluto y relativo que se cometería al calcular  $y = x^3$  si  $x = 3 \pm 0.1$

$$x = 3, \Delta x = 0.1$$

$$\Delta y \approx (d(y)/dx)\Delta x = (3x^2)\Delta x$$

$$\Delta y \approx (3 \times 9)(0.1)$$

$$\Delta y \approx 2.7 = ea$$

$$yap = 3^3 = 27$$

$$yv = yap \pm ea = 27 \pm 2.7$$

$$er = 2.7/27 = 0.10 \%$$

Ej: Una corriente pasa a través de una resistencia de 10 Ohmios, este valor tiene una precisión de 5%, la corriente es de 2 A y fue medida con una aproximación de  $\pm 0.1$  A.

A) Hallar el valor aproximado del voltaje ( $v=i*r$ ).

B) Hallar el error absoluto y relativo

$$i = 2, r = 10$$

$$\Delta i = 0.1, \Delta r = 5\%(10) = 0.5$$

$$v = i * r$$

$$\Delta v \approx \left| \frac{\partial v}{\partial i} \right| \Delta i + \left| \frac{\partial v}{\partial r} \right| \Delta r$$

$$\Delta v \approx r \Delta i + i \Delta r$$

$$\Delta v \approx (10)(0.1) + (2)(0.5)$$

$$\Delta v \approx 2$$

$$v = i * r = 2 * 10 = 20$$

$$v = 20 \pm 2$$

$$vr = 2/20 = 10 \%$$

Ej: Se tiene un triángulo rectángulo cuya altura  $h \sim 3$  cm y la base,  $b \sim 4$  cm. Si se quiere calcular el área con un error no mayor al 10 %. Qué errores se pueden tener en los valores de  $h$  y  $b$  ?

$$h = 3$$

$$b = 4$$

$$area = (b * h) / 2 = 6$$

$$\xi_a^* = 0.1(6) = 0.6$$

$$\xi_a \approx \xi_b * h + \xi_h * b$$

Suponiendo que cada variable contribuye en igual proporción

$$\xi_b^* = \frac{\xi_a^*}{2h} = \frac{0.6}{2(3)} = 0.1$$

$$\xi_h^* = \frac{\xi_a^*}{2b} = \frac{0.6}{2(4)} = 0.075$$

Ej: Calcular el error absoluto y relativo que se cometería al calcular  $z = a^3 + b^2$  si  $a = 2.02 \pm 0.01$  y  $b = 0.60 \pm 0.01$

$$a = 2.02, \quad b = 0.60$$

$$\Delta a = 0.01, \quad \Delta b = 0.01$$

$$z = a^3 + b^2$$

$$\Delta z = \left| \frac{\partial z}{\partial a} \right| \Delta a + \left| \frac{\partial z}{\partial b} \right| \Delta b$$

$$\Delta z = ?$$

$$\left| \frac{\partial z}{\partial a} \right| = ?$$

$$\left| \frac{\partial z}{\partial b} \right| = ?$$

$$z_{ap} = ?$$

$$z = z_{ap} \pm \Delta z = ?$$

La pérdida de precisión debida al error de redondeo puede ser resuelta, reescribiendo los cálculos o con una reformulación del problema; si es posible.

También es conveniente buscar métodos que reduzcan el número de operaciones Por ej: Método de Horner para evaluar un polinomio en un punto

Analizando las magnitudes de los números que intervienen en los cálculos

....

## **METODOS DE ESTIMACION DEL ERROR**

- Doble precisión
- Análisis regresivo del error
- Aritmética de intervalo
- Enfoque estadístico



## MÉTODOS DE ESTIMACION DEL ERROR

- **Enfoque estadístico:** en este método se adopta un modelo estocástico de la propagación del error de redondeo, los errores locales se tratan como si fueran variables aleatorias y se asume que tienen una distribución normal entre sus valores extremos. Se pueden calcular así la desviación estándar, la varianza y estimaciones del error de redondeo acumulado.
- Aun cuando requiere un mayor análisis matemático y tiempo adicional de calculo da muy buenas estimaciones del error , siendo el método más usado actualmente

Métodos Numéricos I

65

## Elementos de Estadística

La medida se repite  $n$  veces y se obtienen los valores  $x_1, x_2, x_3, x_4, \dots, x_i, \dots, x_n$ .

**Valor medio**  $\bar{x} = \frac{\sum x_i}{n}$

**Desviación**  $D_i = x_i - \bar{x}$

**Desviación media**  $\bar{D} = \frac{\sum |D_i|}{n}$

**Desviación standard**  $\sigma = \sqrt{\frac{\sum D_i^2}{n - 1}}$

**Error cuadrático medio**  $\Delta x = \varepsilon = \frac{\sigma}{\sqrt{n}} = \sqrt{\frac{\sum D_i^2}{n (n - 1)}}$

