

Aritmética del Computador



Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

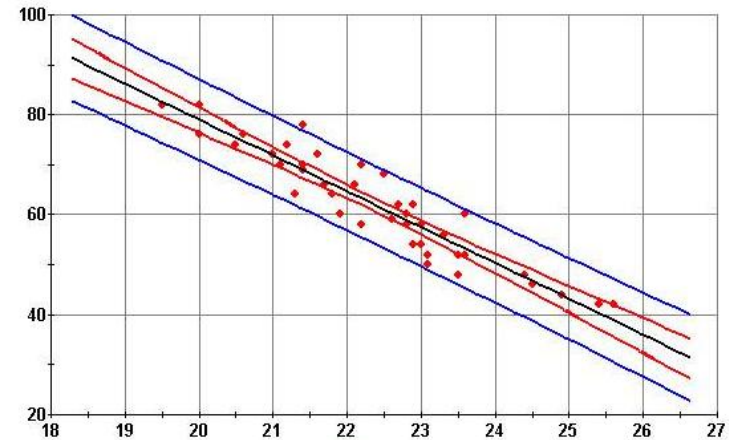
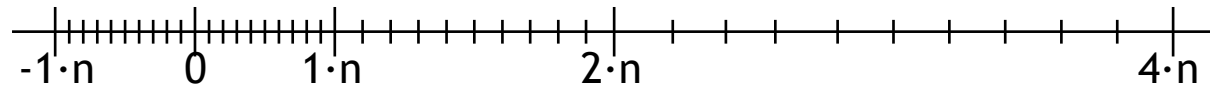
Posibles fuentes de Error

1. Datos de entrada

- Experimental
- Cálculos previos

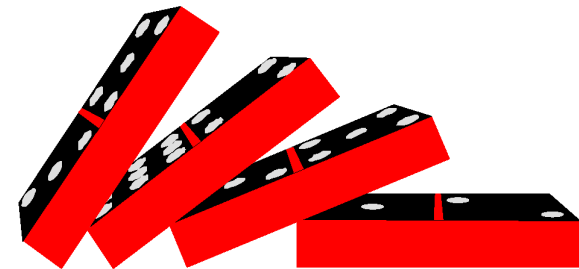
2. Representación de los números :

- Redondeo
- Desbordamiento



3. Cálculos:

- Acumulación de errores de redondeo
- Anulación catastrófica
- Desbordamiento



Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

Estabilidad del algoritmo

dt T

UCM

Representación binaria y decimal

Los **números reales** pueden representarse convenientemente mediante una recta que se extiende hacia el infinito en ambas direcciones. Los **enteros** son los números 0, 1, -1, 2, -2, 3, -3,...

Los **números racionales** son aquéllos que consisten en el cociente de dos enteros, como $1/2$, $2/3$, $6/3$; algunos de éstos, por ejemplo $6/3$, son enteros. Las expansiones de números reales no enteros pueden ser **finitas** o **no terminar**. Por ejemplo, en el caso de $11/2$ ambas expansiones terminan:

$$11/2 = (5.5)_{10} = 5 \times 10^0 + 5 \times 10^{-1}$$

$$11/2 = (101.1)_2 = 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 1 \times 2^{-1}$$

Sin embargo, el número $1/10$, el cual obviamente tiene la representación decimal finita $(0.1)_{10}$, no tiene una representación binaria finita, sino una expansión que no termina:

$$1/10 = (0.0001100110011...)_{2} = 1 \times 2^{-4} + 1 \times 2^{-5} + 0 \times 2^{-6} + 0 \times 2^{-7} + 1 \times 2^{-8} + 1 \times 2^{-9} + 0 \times 2^{-10} + \dots$$

Nótese que esta representación, aunque no termina, es **repetitiva** o **periódica**. La fracción $1/3$ tiene una expansión interminable tanto en binario como en decimal:

$$1/3 = (0.333...)_{10} = (0.010101...)_{2}$$

Los **números racionales** siempre tienen expansiones **finitas** o **repetitivas**. Por ejemplo, $1/7$:



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

irracionales siempre tienen expansiones **interminables** y **no repetitivas**.

Representación de punto fijo (primeros computadores)

Para la mayoría de propósitos de cómputo numérico, los **números reales**, sean racionales o irracionales, se almacenan aproximadamente usando la representación binaria del número. Hay dos métodos posibles, llamados **punto fijo** y **punto flotante**.

En la representación de **punto fijo**, la palabra puede verse como dividida en tres campos:

1. Un campo de bit para el signo del número
2. Un campo de bits para la representación binaria del número anterior al punto binario
3. Un campo de bits para la representación binaria posterior al punto binario.

Por ejemplo, en una palabra de **32 bits** con anchos de campo de 15 y 16 bits respectivamente, el número $(11/2)_{10} = (101.1)_2$ se almacenaría como:

0	000000000000000101	100000000000000000
---	--------------------	--------------------

mientras el número $(1/10)_{10} = (0.0001100110011\dots)_2$ se almacenaría aproximadamente como:

0	000000000000000000	0001100110011001
---	--------------------	------------------

Cartagena99

CLASES PARTICULARES, TUTORIAS TECNICAS ONLINE
LLAMA O ENVIA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

Representación en Punto Flotante (I)

Los computadores representan los números reales utilizando una notación semejante a la conocida notación científica normalizada:

$$\begin{array}{cc}
 \textit{signo} & \textit{exponente} \\
 \underbrace{\quad} & \underbrace{\quad} \\
 + 0.602 \cdot 10^{-22} \\
 \underbrace{\quad} & \underbrace{\quad} \\
 \textit{mantisa} & \textit{base} \\
 \text{Notación científica} \\
 \text{Normalizada}
 \end{array}$$



$$\begin{array}{cc}
 \textit{signo} & \textit{exponente} \\
 \underbrace{\quad} & \underbrace{\quad} \\
 + 0.101110 \cdot 2^{-1100} \\
 \underbrace{\quad} & \underbrace{\quad} \\
 \textit{mantisa} & \textit{base} \\
 \text{Representación en punto} \\
 \text{flotante normalizada}
 \end{array}$$

En todo número en punto flotante se distinguen cuatro componentes:

- ★ **Signo:** indica el signo del número (0=positivo, 1=negativo)
- ★ **Mantisa:** contiene la magnitud del número (en binario puro)
- ★ **Exponente:** contiene el valor de la potencia de la base (notación exceso)
- ★ **Base:** queda implícita y es común a todos los números (habitual base 2)

Otras alternativas menos extendidas para la representaciones de los números reales:

Cartagena99

CLASES PARTICULARES, TUTORIAS TECNICAS ONLINE
LLAMA O ENVIA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

Representación en Punto Flotante (II)

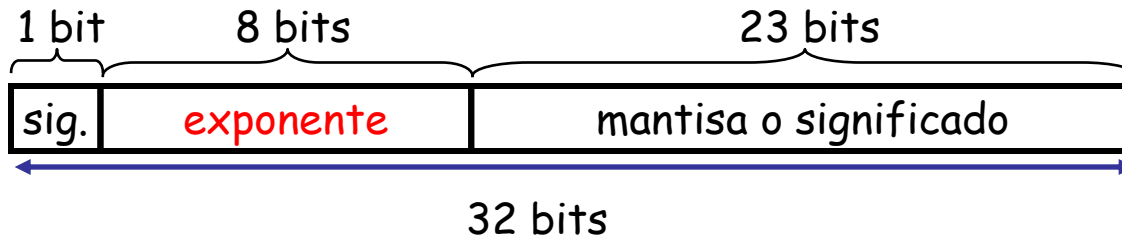
Un mismo número puede tener **varias representaciones** ($0.110 \cdot 2^5 = 110 \cdot 2^2 = 0.0110 \cdot 2^6$)
Los números suelen estar normalizados. Un número está normalizado en base dos si tiene la forma:

$$0.1 \text{ xxxxxxxx} \cdot 2^{\text{xxx}}$$

$$1. \text{ xxxxxxxx} \cdot 2^{\text{xxx}-1}$$

Tienen siempre un 1 a la izquierda: se puede almacenar de forma implícita (**bit fantasma**).

Ejemplo: Formato de punto flotante de 32 bits



Si el número real decimal se puede representar con **exponente** y **mantisa** ocupando los bits asignados, entonces el número real tendrá una representación en punto flotante exacta denominándose **número de máquina**.

El rango de valores representable por cada uno de los campos es:

- **Exponente** (8 bits en exceso a la mitad del exp. max: $254/2=127$): indica el **rango** que el computador entiende como números reales.

$$\text{Exp}_{\text{min}}: (00000001)_2 = (1)_{10} \rightarrow \text{emir} = (\text{exponente} - 127)_{10} = (1 - 127)_{10} = (-126)_{10}$$

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

entre 100, 0 y 1, 11, 1 (1 y 2⁻²³)

El formato simple IEEE



32 bits

Valores del exponente ($a_1 \dots a_8$)	Valore numérico representado
$(00000000)_2 = (0)_{10}$	$\pm(0.b_1 b_2 b_3 \dots b_{23})_2 \cdot 2^{-126}$ \longrightarrow Número subnormal
$(00000001)_2 = (1)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{23})_2 \cdot 2^{(1-127)=-126}$ N_{min}
$(00000010)_2 = (2)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{23})_2 \cdot 2^{-125}$
$(00000011)_2 = (3)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{23})_2 \cdot 2^{-124}$
↓	↓
$(01111111)_2 = (127)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{23})_2 \cdot 2^{(127-127)=0}$
$(10000000)_2 = (128)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{23})_2 \cdot 2^{(128-127)=1}$
↓	↓
$(11111100)_2 = (252)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{23})_2 \cdot 2^{125}$
$(11111101)_2 = (253)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{23})_2 \cdot 2^{126}$
$(11111110)_2 = (254)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{23})_2 \cdot 2^{(254-127)=+127}$ N_{max}

254 \rightarrow exceso (254/2)=127

Épsilon de máquina:
 $\epsilon = 2^{-23} \approx 1.2 \times 10^{-7}$



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

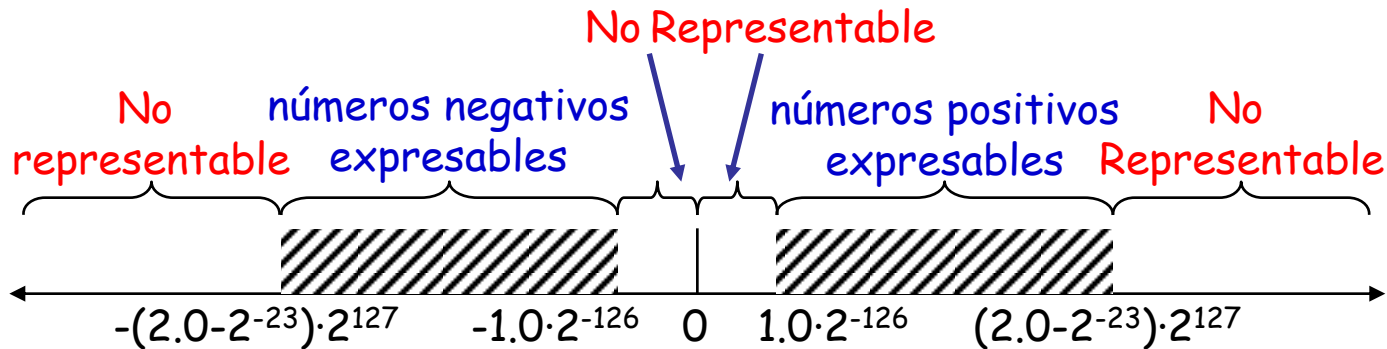
ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

Representación en Punto Flotante (IV)

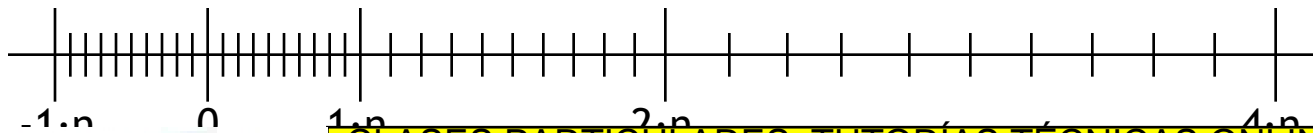
Un número x arbitrario se representará de la forma:

$$(x)_2 = (-1)^{(sm)_2} \times (1.m_1\dots m_{23}) \times 2^{(-1)^{(se)_2} (e_7\dots e_1)_2 - (127)_{10}}$$

✓ **Rango Representable:** en coma flotante con mantisa normalizada nos encontramos con un hueco alrededor del cero de números reales que no se pueden representar. El cero tampoco se podrá representar salvo que se le asigne una representación específica.



✓ La cantidad de números representables es 2^{32} . El espacio de representación no es uniforme: Mayor densidad cerca del 0 !!



CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

Cartagena99

siguiente número de punto flotante)

IEEE 754 (I)

✓ Hasta la década de los 90 cada computador utilizaba su propio formato en punto flotante (numero de bits para exponente y mantisa, como se realiza el redondeo, acciones tomadas en condiciones excepcionales) :

- ★ En algunas máquinas (Crays, IBM System 370) era necesario utilizar ciertos trucos para obtener resultados correctos.

Ejemplo: En algunos computadores, algunos números próximos a cero sólo se consideraban distintos de cero en comparaciones y sumas (se trataban como ceros en multiplicaciones y divisiones). Consecuencia: Para poder utilizar de forma segura una variable como denominador debía multiplicarse por 1.0 antes de comparar con cero!!

- ★ Era imposible escribir programas portables que produjesen los mismos resultados en máquinas diferentes (Sólo el pentágono, AT&T... Podían permitírsele)

✓ **John Palmer** (Intel) promueve la creación de un estándar para punto flotante (1976) antes de comenzar el diseño del coprocesador matemático del i8086/8 y del i432. Primeras reuniones comenzaron en **1977**

✓ El primer borrador fue elaborado por **William Kahan** (Intel), **Jerome Coonen** (Berkeley) v **Harold Stone** (Berkeley): **documento K-C-S**. Inspirado en el trabajo de Kahan en el

Cartagena99

CLASES PARTICULARES, TUTORIAS TECNICAS ONLINE
LLAMA O ENVIA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

Precisión simple (32 bits): Exponente Exceso 127

Exponente (8)	Mantisa (23)	Representa
-127 (combinación binaria 0)	0	$(-1)^{(sm)_2} \cdot (0.0000\dots 0)_2 \cdot 2^{\underbrace{(00000000)_2}_{2^{(-127)_{10}}}} = (\pm 0)_{10}$
-127 (combinación binaria 0)	$\neq 0$	$(-1)^{(sm)_2} \cdot (0.m_{23}\dots m_1)_2 \cdot 2^{(-126)_{10}}$
[-126, 127] ($[1, 254]_{10}$)	\forall	$(-1)^{(sm)_2} \cdot (1.m_{23}\dots m_1)_2 \cdot 2^{(e_8\dots e_1)_2 - (127)_{10}}$
128 (combinación binaria 255)	0	$(-1)^{(sm)_2} \cdot (0.0000\dots 0)_2 \cdot 2^{\underbrace{(11111111)_2}_{2^{(255-127)_{10}} = 2^{(128)_{10}}}} = (\pm \infty)_{10}$
128 (combinación binaria 255)	\neq	$(-1)^{(sm)_2} \cdot (0.m_{23}\dots m_1)_2 \cdot 2^{(128)_{10}} = (NaN)_{10}$

Precisión doble (64 bits): Exponente Exceso 1023

Exponente (11)	Mantisa (52)	Representa
-1023 (comb. binaria 0)	0	$(-1)^{(sm)_2} \cdot (0.0000\dots 0)_2 \cdot 2^{\underbrace{(000000000000)_2}_{2^{(-1023)_{10}}}} = \pm(0)_{10}$
-1023	$\neq 0$	$(-1)^{(sm)_2} \cdot (0.m_{52}\dots m_1)_2 \cdot 2^{(-1022)_2}$
[-1022, 1023] ($[1, 2046]_{10}$)	\forall	$(-1)^{(sm)_2} \cdot (1.m_{52}\dots m_1)_2 \cdot 2^{(e_{11}\dots e_1)_2 - (1023)_{10}}$

CLASES PARTICULARES, TUTORIAS TECNICAS ONLINE
 LLAMA O ENVIA WHATSAPP: 689 45 44 70

 ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70



Codificaciones con significado especial (IV)

e=11111111
(Combinación binaria 255)

$m=0 \rightarrow +\infty \text{ o } -\infty$
(23 bits = 0)

Representan cualquier valor de la región de overflow. El estándar especifica aritmética en el infinito. Ejemplo: $1/\infty=0$, $\arctan(\infty)=\pi/2$.

$m>0 \rightarrow \text{NaN (resultado sin sentido)}$

Se obtiene n como resultado de operaciones inválidas. El estándar especifica que cuando el argumento de una operación es un NaN, el resultado es NaN.

Ejemplo: $0 \times \infty$, $0/0$ y $\infty - \infty$. Estas operaciones no tienen sentido matemático.

e=00000000
(Combinación binaria 0)

$m=0000\dots0 \rightarrow 0 \text{ decimal}$

El cero se representa como número no normalizado.

$m>0 \rightarrow \text{Números desnormalizados}$

Número sin normalizar cuyo bit implícito se

CLASES PARTICULARES, TUTORIAS TECNICAS ONLINE
LLAMA O ENVIA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

Cartagena99

números subnormales).

UCM

2046 → exceso (2046/2)=1023

Valores del exponente ($a_1 \dots a_{11}$)	Valore numérico representado
$(00000000000)_2 = (0)_{10}$	$\pm(0.b_1 b_2 b_3 \dots b_{52})_2 \cdot 2^{-1022}$ → Número subnormal
$(00000000001)_2 = (1)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{52})_2 \cdot 2^{(1-1023)=-1022}$ N_{\min}
$(00000000010)_2 = (2)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{52})_2 \cdot 2^{-1021}$
$(00000000011)_2 = (3)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{52})_2 \cdot 2^{-1020}$
↓	↓
$(01111111111)_2 = (1023)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{52})_2 \cdot 2^0$
$(10000000000)_2 = (1024)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{52})_2 \cdot 2^1$
↓	↓
$(11111111100)_2 = (2044)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{52})_2 \cdot 2^{1021}$
$(11111111101)_2 = (2045)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{52})_2 \cdot 2^{1022}$
$(11111111110)_2 = (2046)_{10}$	$\pm(1.b_1 b_2 b_3 \dots b_{52})_2 \cdot 2^{(2046-1023)=+1023}$ N_{\max}

Épsilon de máquina:
 $\epsilon = 2^{-52} \approx 2.2 \times 10^{-16}$

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

IEEE 754 (IV).Ejemplos

$$(x)_{10} = (-1)^{(sm)_2} \times (1 + \text{fracción}) \times 2^{(\text{Exponente}-\text{Exceso})_{10}}$$

Ejemplo 1: ¿Cuál es el valor decimal de:

1 0111100 110000000000000000000000?

- ✓ El bit de signo es **1** : **número negativo**
- ✓ El exponente contiene **0111100** = $(124)_{10}$
- ✓ La mantisa es **0.11000...** = $(0.75)_{10}$

El valor es : $(-1)^1 \times (1 + 0.75) \times (2^{(124-127)}) = -1.75 \times (2^{-3}) = -0.21875$

Ejemplo 2: ¿Cuál es el valor decimal de:

0 1000001 010000000000000000000000?

- ✓ El bit de signo es **0** : **número positivo**
- ✓ El exponente contiene **1000001** = $(129)_{10}$

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

Cartagena99

IEEE 754 (V).Ejemplos

Ejemplo 3: ¿Cuál es la representación en simple precisión de : 347.625 ?

1- Convertir a binario: $347.625 = 101011011.101$

2- Normalizar el número (mover el punto decimal hasta que haya un solo 1 a la izquierda) $101011011.101 = 1.01011011101 \times (2^8)$

3- mantisa: 01011011101

4- exponente: 8 en exceso 127 $\rightarrow 8 = (\text{exponente} - \text{exceso } 127) \rightarrow \text{exponente} = (127 + 8) = (135)_{10} = 10000111$

5- El número es positivo: bit de signo 0

Resultado : 0 10000111 010110111010000000000000

Ejemplo 4: Representar el número decimal 5/3 en precisión simple.

$(5/3)_{10} = 1.6666... \times 10^0$

Mantisa: 1.66666...

\Rightarrow Parte entera: $1 \times 2^0 = 1$

\Rightarrow Parte decimal: $0.666... \times 2 = 1.333...$

$0.333... \times 2 = 0.666...$ (Es cíclico. Infinitos decimales)

Por tanto $(0.66...)_10 = (0.101010...)_2$

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

0 0111111 10101010101010101010101

Errores en la representación de números reales

Número real que no sea un número máquina

Se puede deber a:

1. Que el exponente, una vez normalizado, esté fuera del rango admitido (demasiado grande o pequeño).
2. Que la mantisa, una vez normalizada, tenga más de 23 bits (más bits de los que puede almacenar).

¿Cómo resolverlo?

1. Si el exponente se sale del rango admitido se produce el fenómeno de **desbordamiento**, bien por *exceso* (*overflow*) bien por *defecto* (*underflow*).

Desbordamiento por exceso: toma, p.e., el valor $\pm\infty$ (según sea el número positivo o negativo), y sigue trabajando siempre y cuando tengan sentido las operaciones siguientes (p.e. $1/\infty=0$). Normalmente es el resultado de un error de cálculo que no ha considerado el programador.

Desbordamiento por defecto: se asigna, por ejemplo, el valor cero y continua el cálculo reemplazando el resultado por cero o se desnormaliza el número para poder tener en número máquina próximo en el entorno cero.

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

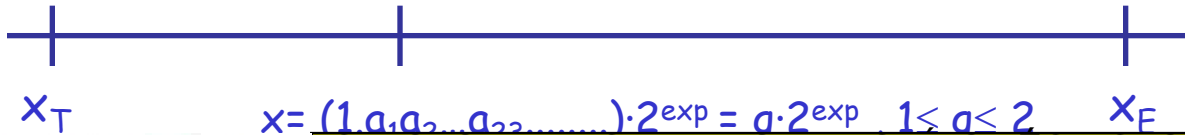
Error de Redondeo Unitario (I)

Representación en punto flotante de $x \rightarrow fl(x) = (1.a_1a_2\dots a_{23}) \cdot 2^{exp}$ (suponemos 23 bits de mantisa)
 2^{-23}

Error de redondeo unitario:

Depende del nº de bits de la mantisa
 épsilon (*eps*) o precisión del computador

- TRUNCAMIENTO: $fl(x) = x_T = (1.a_1a_2\dots a_{23}) \cdot 2^{exp}$
 Se descartan el resto de bits
- EXCESO: $fl(x) = x_E = ((1.a_1a_2\dots a_{23}) + 2^{-23}) \cdot 2^{exp}$
 El último bit se incrementa en uno
- REDONDEO: $fl(x) = \begin{cases} x_T \\ x_E \end{cases}$ Se usa el más cercano



Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP: 689 45 44 70

Error de Redondeo Unitario (II)

Formato n Bits $eps=2^{-n}$

"Axioma de la Representación en punto Flotante"

$$fl(x) = x_{\text{redondeado}} = x \cdot (1 + \delta); \quad |\delta| \leq eps$$

Matlab:

Precisión del sistema: Número de bits de la mantisa (incluido el bit escondido)

$$\downarrow$$

$$eps = 2^{-(p-1)} = 2^{-52} = 2.22044 \text{ e} - 016$$

Si está en efecto el modo de redondeo al más cercano:

$$|\delta| \leq \frac{1}{2} eps = 2^{-p}$$

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

Error de Redondeo Unitario (III)

Ejemplo: Error al representar $(2/3)_{10}$ en formato simple IEEE754

Calcular el error de truncado absoluto y relativo
Cual es la representación final

- Se convierte en binario:

$$(2/3)_{10} = (0.1010...)_{2} = (1.010101...)_{2} \times 2^{-1}$$

- Los números de maquinas más cercanos son:

$$x_{-} = (1.0101...010)_{2} \times 2^{-1} \longleftarrow \text{Obtenido truncando}$$

$$x_{+} = (1.0101...011)_{2} \times 2^{-1} \longleftarrow \text{Obtenido redondeando}$$

- Determinamos el más cercano calculando $(x-x_{-})$ y $(x_{+}-x)$

$$(x-x_{-}) = (0.1010...)_{2} \times 2^{-24} = (2/3)_{10} \times 2^{-24}$$

$$(x_{+}-x) = (x_{+}-x_{-}) - (x-x_{-}) = 2^{-24} - (2/3)_{10} \times 2^{-24} = (1/3) \times 2^{-24} \longleftarrow f(x) = x_{+}$$

- Error de truncado absoluto:

$$|f(x)-x| = (1/3) \times 2^{-24}$$

- Error de truncado relativo:

$$|f(x)-x|/|x| = (1/3 \times 2^{-24}) / (2/3)_{10} = 2^{-25}$$

- Representación final:

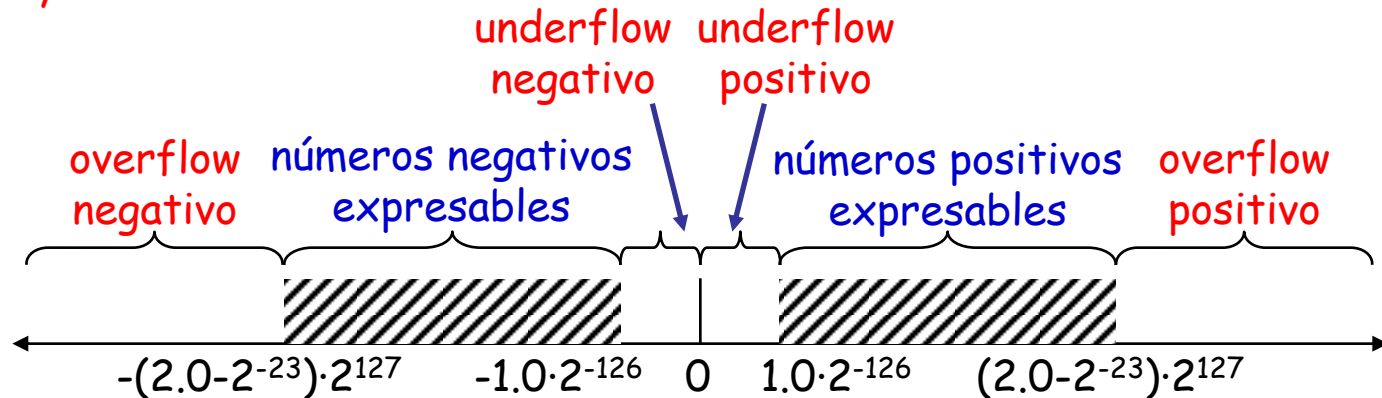
Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

Error de Desbordamiento (I)

- ✓ Solo se puede representar un rango limitado de valores
 1. Principalmente dependen del nº de bits asignados al exponente
- ✓ **Overflow y Underflow :**



- ✓ Proteger las conversiones de tipos utilizando los límites representables (`realmax`, `intmax`...)

```
realmax=1.79769313486231e+308
intmax= 2147483647
```

- ✓ Utilizar códigos especiales para marcar que se ha producido desbordamiento, ***+-inf***

```
If a>intmax then a=intmax
```

```
If a~=0 then c=b/a
```

```
else c=b/eps
```

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

`d = Inf` ←

Error de Desbordamiento (II)

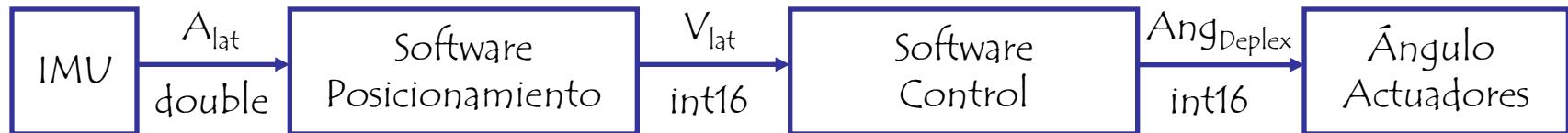
Cuidado con las Conversiones (punto flotante \rightarrow entero)



French Guyana, June 4, 1996
\$800 million embedded software failure

Guayana Francesa, 4 Junio 1996. Un Cohete **Ariane V** de la Agencia Espacial Europea estalla a los 40 segundos del despegue. Era el primer lanzamiento del Ariane V, después de 10 años de desarrollo (7 billones de \$). El cohete y su carga estaban valorados en 800 millones de \$.

El error fue debido a un desbordamiento en el software de Posicionamiento. Cuando se presentó la operación inválida, el programa cerró todo el sistema de guía y el cohete se autodestruyó.



La velocidad horizontal del cohete con respecto a la plataforma se convertía a un

Interpretación
opuesta

Valor $\rightarrow +5.1^\circ$
maxint $\rightarrow +5^\circ$

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

representable con 16 bits) y la conversión

$\rightarrow -\text{maxint} \rightarrow -5^\circ$

Acumulación de Errores de Redondeo (I)

Ejemplo : Suma en punto flotante ($99.99 + 0.161 = 100.151$). Suponemos que:

- ✓ La mantisa puede ser de **4** dígitos decimales
- ✓ El exponente puede tener **2** dígitos

- Extraer signos, exponentes y magnitudes.
- Tratar operandos especiales (por ejemplo, alguno de ellos a cero)
- Desplazar la mantisa del número más pequeño a la derecha $|e_1 - e_2|$ bits
- Fijar el exponente del resultado al máximo de los exponentes

1. Alineamiento

$$\begin{array}{r} 9.999 \cdot 10^1 \\ 1.610 \cdot 10^{-1} \end{array}$$

$$\begin{array}{r} 9.999 \cdot 10^1 \\ + 0.0161 \cdot 10^1 \\ \hline \end{array}$$

$$10.015 \times 10^1$$

2. Operación

- Si la operación es suma y los signos son

3. Normalización

$$\begin{array}{r} 1.0015 \cdot 10^2 \\ 1.002 \cdot 10^2 \end{array}$$

4. Redondeo

- Redondear el resultado y renormalizar la mantisa si es necesario (IEE754 ultimo par si hay duda)
- Corregir el exponente en función de los desplazamientos realizados sobre la mantisa.
- Detectar overflow o underflow del exponente

Cartagena99

CLASES PARTICULARES, TUTORIAS TECNICAS ONLINE
LLAMA O ENVIA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

5. Re-normalización: A veces es necesario volver a normalizar ($9.9999 \rightarrow 10.000 \rightarrow \dots$)

UCM

Acumulación de Errores de Redondeo (II)

Hay **SITUACIONES EXTREMAS** en las que no se puede evitar considerables errores de redondeo: ERRORES DEBIDO A LA PERDIDA DE DIGITOS SIGNIFICATIVOS

Ejemplo 1: Suma de cantidades grandes a números pequeños

$$(1.5 \cdot 10^{38}) + (1.0 \cdot 10^0) = 1.5 \cdot 10^{38}$$

Alineamiento : es necesario desplazar el punto decimal 38 posiciones a la izquierda



Con 16 dígitos decimales de mantisa ese desplazamiento supone convertir al segundo sumando en 0!

No hay Propiedad Conmutativa !!!!!

$$x = 1.5 \cdot 10^{38}$$

$$(x + y) + 1 \longrightarrow 1$$

$(x + y) + 1 \neq x + (y + 1)$
 CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
 LLAMA O ENVÍA WHATSAPP: 689 45 44 70

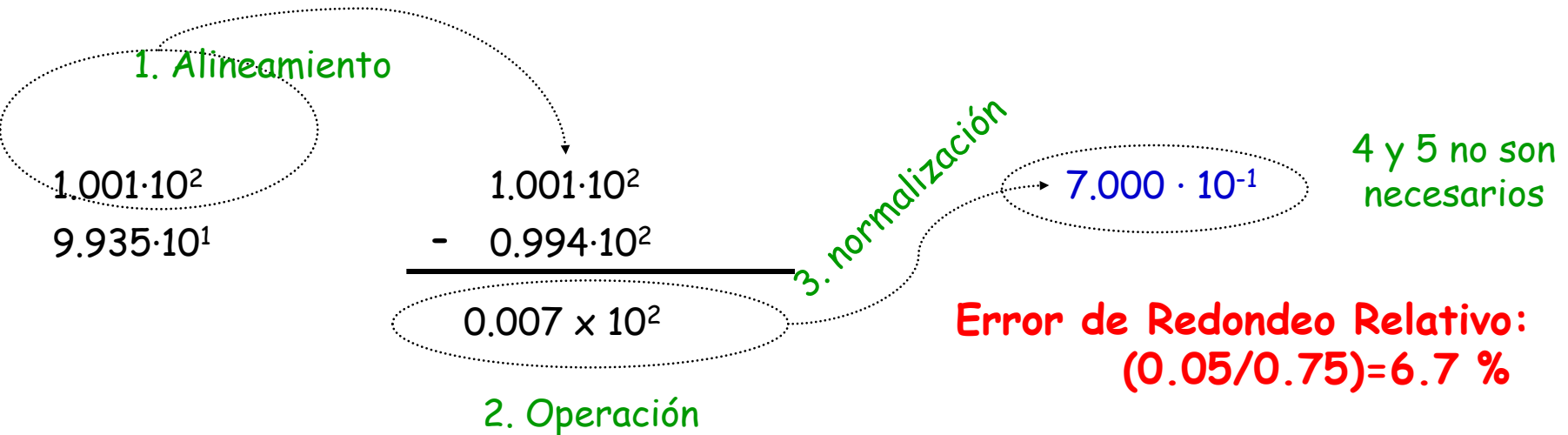
ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
 CALL OR WHATSAPP:689 45 44 70

Cartagena99

Acumulación de Errores de Redondeo (III)

Ejemplo 2 : Sustracción de cantidades casi iguales ($100.1 - 99.35 = 0.75$).

- ✓ La mantisa puede ser de **4** dígitos decimales
- ✓ El exponente puede tener **2** dígitos



Problema: En el Alineamiento se pierde 1 dígito significativo

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

Acumulación de Errores de Redondeo (IV)

Serie 360 de IBM (1968)

Dígitos de Protección (*guard digits*) :

Para minimizar este problema, el estándar IEEE 754 especifica que los cálculos intermedios deben realizarse utilizando 2 dígitos de protección.

Ejemplo 2 revisado

1. Alineamiento

$$\begin{array}{r} 1.0010 \cdot 10^2 \\ 9.9350 \cdot 10^1 \end{array}$$

$$\begin{array}{r} 1.0010 \cdot 10^2 \\ - 0.9935 \cdot 10^2 \\ \hline \end{array}$$

$$0.0075 \times 10^2$$

2. Operación

3. normalización

$$7.5000 \cdot 10^{-1}$$

4 y 5 no son necesarios

Resultado Correcto

Si utilizamos 2 dígitos de guarda para realizar las subtracciones el error de redondeo relativo será siempre menor que $\frac{1}{2^{20}}$.

Cartagena99

CLASES PARTICULARES, TUTORIAS TECNICAS ONLINE
LLAMA O ENVIA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

Acumulación de Errores de Redondeo (V)

Supongamos un computador decimal de **5** dígitos (**epsilon del computador** = 10^{-5}) en el que se utiliza formato extendido para cálculos intermedios y dos números maquina x e y :

$$x = 3.1426 \cdot 10^2 \quad y = 9.2577 \cdot 10^4$$

OPERACIÓN	RESULTADO NORMALIZADO	REDONDEO	ERROR RELATIVO
Multiplicación	$2.909324802 \cdot 10^7$	$2.9093 \cdot 10^7$	$8.5 \cdot 10^{-6}$
Suma	$9.289126 \cdot 10^4$	$9.2891 \cdot 10^4$	$2.3 \cdot 10^{-6}$
Resta	$-9.2262740 \cdot 10^4$	$-9.2263 \cdot 10^4$	$2.8 \cdot 10^{-6}$
División	$3.394579647 \cdot 10^{-3}$	$3.3946 \cdot 10^{-3}$	$6.0 \cdot 10^{-6}$

- ✓ x e y números máquina (es decir, son representables de forma exacta)
- ✓ \odot operación aritmética básica (+, -, \times , \div)

Axioma fundamental de la aritmética en punto flotante:

$$\text{Flotante } (x \odot y) = (x \odot y) (1 + \delta) \quad |\delta| \leq \text{eps}$$

Cartagena99

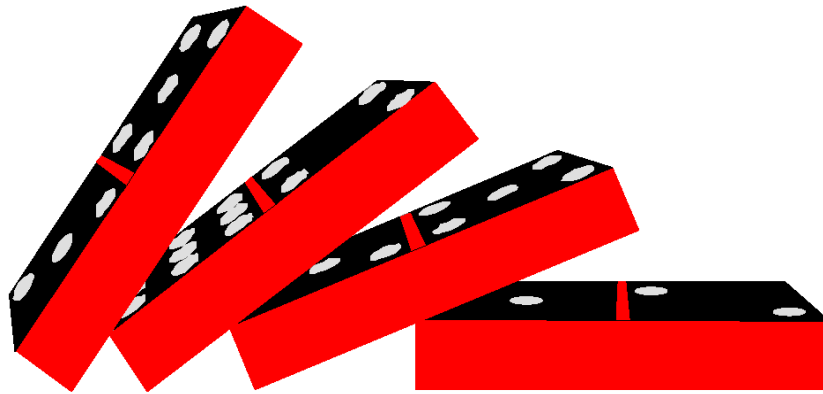
CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

Acumulación de Errores de Redondeo (VI)

Cómo se acumulan los errores al encadenar varias operaciones?

A pesar de que el error de redondeo al realizar una operación aritmética básica pueda ser pequeño, **al encadenar operaciones** dichos **errores pueden magnificarse** (acumularse) considerablemente.



Ejemplo: x , y , z números máquina, se desea realizar la operación $(x \cdot (y+z))$. Siendo $|\delta_1|$, $|\delta_2|$ y $|\delta| \leq \text{eps}$.

$$\text{Flotante}(x \cdot (y + z)) = (x \cdot \text{Flotante}(y + z)) (1 + \delta_1) \quad |\delta_1| \leq \text{eps}$$

$$= (x \cdot (y + z)) (1 + \delta_2) (1 + \delta) \quad |\delta_2| \leq \text{eps}, |\delta| \leq \text{eps}$$

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

aproximadamente igual a $n \cdot \text{eps}$.

Acumulación de Errores de Redondeo(VII). Ejemplo

Guerra del Golfo, 25 Febrero 1991. Una batería de antimisiles Patriot en Dharan (Arabia Saudita), no consiguió interceptar un misil Scud iraquí, que impactó contra un barracón de la armada americana, muriendo 28 soldados.

La predicción de la "puerta de alcance" (*range gate*) se calcula en función de:

- ✓ Velocidad del Scud (dato conocido : 1676 m/s)
- ✓ Tiempo transcurrido desde que el Scud fue detectado por el radar por última vez



- ✓ Para las mediciones de tiempo, el sistema utiliza un reloj interno con una precisión de décimas de segundo. El tiempo transcurrido desde el *boom* se almacenaba en un registro

Cartagena99

CLASES PARTICULARES, TUTORIAS TECNICAS ONLINE
LLAMA O ENVIA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

Acumulación de Errores de Redondeo(VIII). Ejemplo

El tiempo en segundo se almacena en un registro de 24 bits. Aunque el reloj interno daba décimas de segundo, para hacer los cálculos, el número entero con las décimas de segundo transcurridas se convertía a punto flotante, multiplicándose el valor por 0.1 para obtener el tiempo en segundos. Es decir, **el reloj interno de la computadora que controlaba el sistema de defensa almacenaba el tiempo como un valor entero en unidades de décimas de segundo, y el programa de la computadora lo convertía a un valor de punto flotante en unidades de segundos, redondeando consiguientemente la expansión.**

¿Qué pasa con el 0.1?



¿Cuál es la representación en simple precisión de: $(0.10)_{10}$?

- 1- Convertir a binario: $0.10 = 0.0001100110011001100110011001100 \dots$
- 2- Normalizar el número : $1.100110011001100110011001100 \dots \times 2^{-4}$
- 3- Truncar el número: 23 bits de mantisa!!! $1.10011001100110011001100 \times 2^{-4}$

0.10 no se puede escribir como una suma finita de potencias de 2 !!!

0.10 se representa en precisión simple por el número 0.09999999403953552

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

end

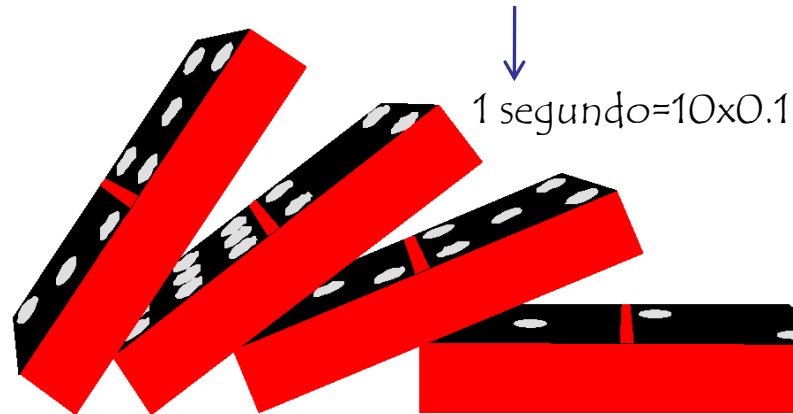
Acumulación de Errores de Redondeo(IX). Ejemplo

Batería Antimisiles Patriot: El error de redondeo cometido al representar 0.1 es \approx de 5.9×10^{-9} .

00011001100110011001100

La batería llevaba unas 100 horas en funcionamiento antes de producirse el fallo. **El error cometido al representar 0.1 SE HABIA ACUMULADO**

$$5.9 \times 10^{-9} \times (100 \times 60 \times 60 \times 10) = 0.0212 \text{ segundos}$$



En ese intervalo, un misil *Scud* recorre **35.60 m**. Este error solo se tenía en cuenta en algunas de las partes del software **(LOS ERRORES NO SE CANCELABAN!!!!)**

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

Software Problem Led to System Failure at Dhahran, Saudi Arabia

Anulación Catastrófica (II)

Regla Práctica: Se debe evitar situaciones en las que se restan cantidades casi iguales

Ejemplo 1

$$y \leftarrow \sqrt{x^2 + 1} - 1$$

Implica una cancelación por sustracción y pérdida de dígitos significativos para valores **pequeños de x**

¿Cómo evitarlo?



Solución: utilizar una forma alternativa de asignación

Multiplicando y Dividiendo por el Conjugado

$$y \leftarrow \frac{(\sqrt{x^2 + 1} - 1)(\sqrt{x^2 + 1} + 1)}{(\sqrt{x^2 + 1} + 1)} = \frac{x^2}{\sqrt{x^2 + 1} + 1}$$

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

Cartagena99

Anulación Catastrófica (IV)

Ejemplo 3

$y \leftarrow x - \sin(x)$ **Implica una cancelación para valores pequeños de x**

¿Cómo evitarlo? Desarrollo en serie de Taylor para $\sin(x)$

$$\begin{aligned} y &= x - \sin(x) \\ &= x - \left(x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} + \dots \right) \\ &= \frac{x^3}{3!} - \frac{x^5}{5!} + \frac{x^7}{7!} + \dots \end{aligned}$$

Si x está cercano a cero, se puede utilizar una serie truncada

$$y \leftarrow \frac{x^3}{6} \left\{ 1 - \frac{x^2}{20} \left[1 - \frac{x^2}{42} \left(1 - \frac{x^2}{72} \right) \right] \right\}$$

Cartagena99

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

uno con su propio rango

Cálculos Estables e Inestables (I)

Proceso/Algoritmo Numérico Inestable (Informal) :

Los pequeños errores que se producen en alguna de sus etapas se agrandan en etapas posteriores y degradan seriamente la exactitud del resultado final, a pesar de que en aritmética exacta el algoritmo sea correcto.

Algoritmo Estable (*backward stability*) :

Si $\text{alg}(x)$ representa la versión discreta de un algoritmo $\text{ALG}(x)$, en el que se incluye el efecto del error de redondeo, se dice que $\text{alg}(x)$ es estable si para todo los datos de entrada x existe un pequeño δx tal que $\text{alg}(x) = \text{ALG}(x + \delta x)$

backward error

EJEMPLO 1: sucesión de números reales

$$\begin{cases} x_0 = 1 \\ x_1 = \frac{1}{3} \\ \dots \\ x_{13} = \dots \\ \dots \\ x_4 = \dots \end{cases}$$



$$x'_n = \left(\frac{1}{3}\right)^n = \frac{1}{3} x'_{n-1}$$

Por inducción, se puede

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP: 689 45 44 70

Cartagena99

Estimar el error de redondeo de x_n .

Cálculos Estables e Inestables (II)

# iter	Valor Inestable	Valor Estable	Error Absoluto	Error Relativo
	x_n	x'_n	$ x_n - x'_n $	$\frac{ x_n - x'_n }{x_n}$
	1.0000e+000	1.0000e+000	0	0
	3.3333e-001	3.3333e-001	0	0
1	1.1111e-001	1.1111e-001	1.6653e-016	1.4988e-015
2	3.7037e-002	3.7037e-002	7.7716e-016	2.0983e-014
3	1.2346e-002	1.2346e-002	3.1641e-015	2.5629e-013
4	4.1152e-003	4.1152e-003	1.2675e-014	3.0800e-012
5	1.3717e-003	1.3717e-003	5.0707e-014	3.6966e-011
6	4.5725e-004	4.5725e-004	2.0283e-013	4.4359e-010
7	1.5242e-004	1.5242e-004	8.1133e-013	5.3231e-009
8	5.0805e-005	5.0805e-005	3.2453e-012	6.3878e-008
9	1.6935e-005	1.6935e-005	1.2981e-011	7.6653e-007
10	5.6450e-006	5.6450e-006	5.1925e-011	9.1984e-006
11	1.8815e-006	1.8817e-006	2.0770e-010	1.1038e-004
	6.2639e-007	6.2723e-007	8.3080e-010	1.3246e-003
	2.0575e-007	2.0908e-007	3.3232e-009	1.5895e-002
	5.6399e-008	6.9692e-008	1.3293e-008	1.9074e-001
	-2.9941e-008	2.3231e-008	5.3171e-008	2.2889e+000

• Error absoluto : se multiplica por $\approx 13/3$ en cada iteración, en 10 iteraciones:

$$(13/3)^{10} 10^{-16} \approx 10^{-10}$$

• x_n es una **sucesión** monótona **decreciente**



El error Relativo crece considerablemente

Si, en cambio, la solución

CLASES PARTICULARES, TUTORÍAS TÉCNICAS ONLINE
LLAMA O ENVÍA WHATSAPP: 689 45 44 70

ONLINE PRIVATE LESSONS FOR SCIENCE STUDENTS
CALL OR WHATSAPP:689 45 44 70

Cartagena99

3.4031e-006

8.6039e-010

3.4030e-006

3.9551e-003