

## **Tema 1.**

**Computación con precisión finita. Representación en coma flotante. Normas vectoriales y matriciales.**

## **Bibliografía:**

**“Matrix Computations”. G.Golub & C. Van Loan. Baltimore ; London : Johns Hopkins University Press, 1996 (u otra edición del libro)**

## **Lecturas recomendadas:**

**“Matrix Computations”. G.Golub & C. Van Loan.**

**Capítulo 1. Puntos 1.1, 1.2 y 1.3**

**Capítulo 2. Puntos 2.1, 2.2, 2.3, 2.4 y 2.7**

# Computación en precisión finita

- Las operaciones aritméticas realizadas en un computador suelen estar afectadas por errores de redondeo
- Solo se dispone de una cantidad limitada de memoria para almacenar números
- Los sistemas de representación en coma flotante solo permiten representar un subconjunto finito de elementos del conjunto de los números reales
- Existen números consecutivos y números máximo y mínimo representables
- Los algoritmos van a trabajar con datos aproximados
- Esto va a dar lugar a que tengan que ser analizados y diseñados con cuidado para evitar errores.

## Sistemas de representación en coma flotante

- Se caracterizan por cuatro números enteros:

Base :  $\beta$ . Precisión :  $t$ . Rango de exponentes :  $[L, U]$

$F \equiv$  Conjunto de números en coma flotante

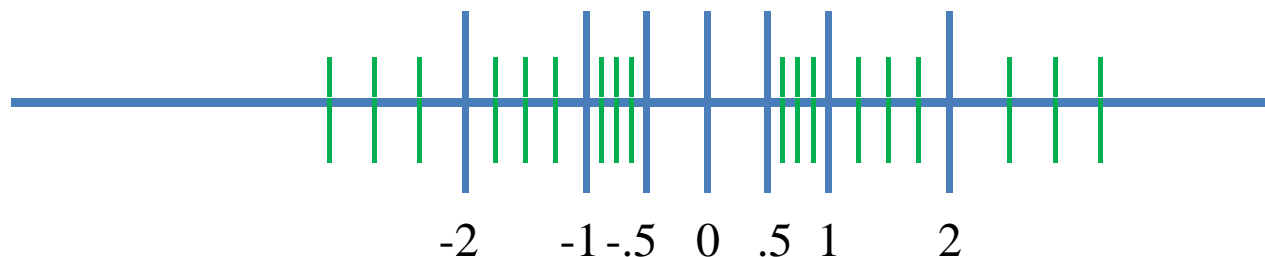
$$F = \{f \in \mathbb{R} : f = \pm .d_1 d_2 \dots d_t \times \beta^e, \text{ con } 0 \leq d_i < \beta, d_1 \neq 0, L \leq e \leq U\}$$

- Números extremos representables

Para  $f \in F, f \neq 0$ , se tiene

$$m \leq |f| \leq M, \text{ con } m = \beta^{L-1} \text{ y } M = \beta^U (1 - \beta^{-t});$$

- Ejemplo:  $\beta = 2; t = 3; L = 0; U = 2$ .



# Modelo aritmético de las operaciones en coma flotante

- \* Sea  $G = \{x \in \mathfrak{R} : m \leq |x| \leq M\} \cup \{0\}$

Se define el operador en coma flotante  $fl: G \rightarrow F$

$$fl(x) = \begin{cases} c \in F \text{ más próximo a } x, \text{ si se usa aritmética redondeada} \\ c \in F \text{ más próximo a } x \text{ que verifique } |c| \leq |x|, \text{ si se usa aritmética truncada} \end{cases}$$

- \* Redondeo unidad

$$u = \begin{cases} \frac{1}{2} \beta^{1-t}, \text{ si se usa aritmética redondeada} \\ \beta^{1-t} \text{ si se usa aritmética truncada} \end{cases}$$

- \* Errores debidos al uso del operador en coma flotante

$$fl(x) = x(1 + \varepsilon), \text{ con } |\varepsilon| \leq u$$

# Sistemas de representación en coma flotante en el Matlab

Doble precisión :  $\beta = 2$ ;  $t = 52$ ;  $L = -1023$ ;  $U = 1024$ .

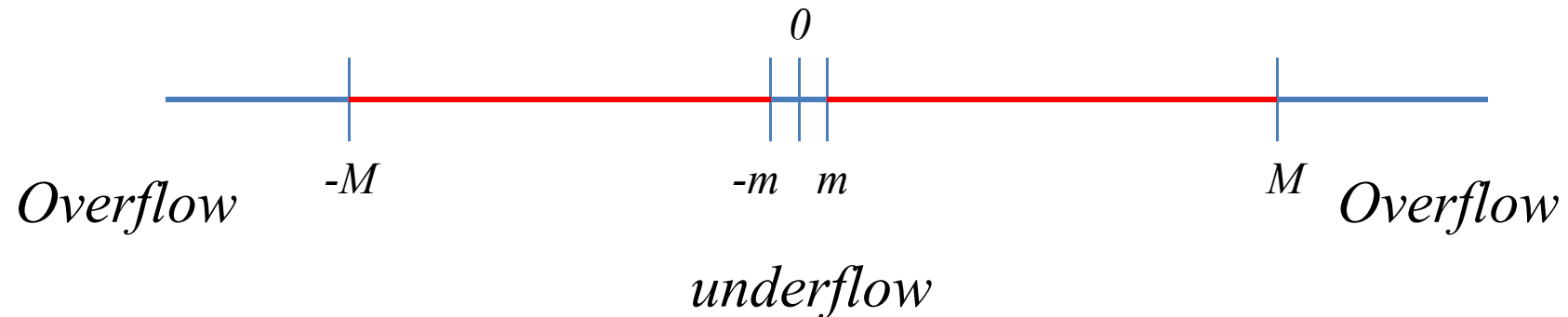
Simple precisión :  $\beta = 2$ ;  $t = 23$ ;  $L = -127$ ;  $U = 128$ .

- Help/Math Constants

	Simple precisión	Doble precisión
Redondeo unidad: <b>eps</b>	1.1920929e-007	2.220446049250313e-016
Real máximo: <b>realmax</b>	3.4028235e+038	1.797693134862316e+308
Real mínimo : <b>realmin</b>	1.1754944e-038	2.225073858507201e-308
Infinito: <b>Inf</b>	1/0 (Overflow)	1/0 (Overflow)
Not-a-Number: <b>NaN</b>	0/0 , Inf-Inf (Indeterminación)	0/0, Inf-Inf (Indeterminación)

# Operaciones que producen inestabilidad

Representaciones fuera de rango: *Overflow* y *underflow*



*Overflow*: puede presentarse al dividir números grandes por números pequeños

*Cancelaciones catastróficas*: puede presentarse al obtener números pequeños a partir de operaciones con números grandes, especialmente restas

## Normas Vectoriales

$f : \mathfrak{R}^n \rightarrow \mathfrak{R}$  es una norma vectorial si :

$$1. f(x) \geq 0 \text{ y } f(x) = 0 \leftrightarrow x = 0, \forall x \in \mathfrak{R}^n$$

$$2. f(x + y) \leq f(x) + f(y), \forall x, y \in \mathfrak{R}^n$$

$$3. f(\alpha x) = |\alpha| f(x), \forall \alpha \in \mathfrak{R}, \forall x \in \mathfrak{R}^n$$

Ejemplos: p-normas vectoriales

$$x \in \mathfrak{R}^n$$

$$\|x\|_p = (|x_1|^p + |x_2|^p + \dots + |x_n|^p)^{1/p}$$

$$\|x\|_1 = (|x_1| + |x_2| + \dots + |x_n|)$$

$$\|x\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2} = \sqrt{x^T x}$$

$$\|x\|_\infty = \max_{1 \leq i \leq n} (|x_i|)$$



# Algunas propiedades de las normas vectoriales

- Desigualdad de Holder:  $|x^T y| \leq \|x\|_p \|y\|_q$ , si  $\frac{1}{p} + \frac{1}{q} = 1$
- Desigualdad de Cauchy-Schwartz:  $|x^T y| \leq \|x\|_2 \|y\|_2$
- Equivalencia de todas las normas vectoriales en  $\mathfrak{R}^n$

Si  $\|\bullet\|_\alpha$  y  $\|\bullet\|_\beta$  son normas en  $\mathfrak{R}^n$ , existen constantes positivas  $c_1$  y  $c_2$  tales que  $c_1 \|x\|_\alpha \leq \|x\|_\beta \leq c_2 \|x\|_\alpha$ ,  $\forall x \in \mathfrak{R}^n$

Ejemplos: Si  $x \in \mathfrak{R}^n$  :

$$\|x\|_2 \leq \|x\|_1 \leq \sqrt{n} \|x\|_2$$
$$\|x\|_\infty \leq \|x\|_2 \leq \sqrt{n} \|x\|_\infty$$
$$\|x\|_\infty \leq \|x\|_1 \leq n \|x\|_\infty$$

## Normas Matriciales

$f : \mathfrak{R}^{m \times n} \rightarrow \mathfrak{R}$  es una norma matricial si :

1.  $f(x) \geq 0$  y  $f(x) = 0 \leftrightarrow x = 0, \forall x \in \mathfrak{R}^{m \times n}$

2.  $f(x + y) \leq f(x) + f(y), \forall x, y \in \mathfrak{R}^{m \times n}$

3.  $f(\alpha x) = |\alpha| f(x), \forall \alpha \in \mathfrak{R}, \forall x \in \mathfrak{R}^{m \times n}$

## Ejemplos de Normas Matriciales

$$A \in \mathfrak{R}^{m \times n}, x \in \mathfrak{R}^n$$

Norma de Frobenius

$$\|A\|_F = \left( \sum_{i=1}^m \sum_{j=1}^n a_{ij}^2 \right)^{1/2} = (tr(A^T A))^{1/2} \quad (* \quad tr(A) = \sum_{i=1}^{\min(m,n)} a_{ii} \quad *)$$

p - norma matricial

$$\|A\|_p = \sup_{x \in \mathfrak{R}^n} \frac{\|Ax\|_p}{\|x\|_p} = \sup_{\|x\|_p=1} \|Ax\|_p$$

## Errores absoluto y relativo al aproximar vectores y matrices

$$\text{Si } x \in \mathfrak{R}^n \text{ se aproxima por } x^* \in \mathfrak{R}^n \quad E_{ab}(x) = \|x - x^*\| \quad E_r(x) = \frac{\|x - x^*\|}{\|x\|}$$

$$\text{Si } A \in \mathfrak{R}^{m \times n} \text{ se aproxima por } A^* \in \mathfrak{R}^{m \times n} \quad E_{ab}(A) = \|A - A^*\| \quad E_r(A) = \frac{\|A - A^*\|}{\|A\|}$$

## Distancia entre vectores y entre matrices

Si  $x, y \in \mathfrak{R}^n$

$$d(x, y) = \|x - y\|$$

Si  $A, B \in \mathfrak{R}^{m \times n}$

$$d(A, B) = \|A - B\|$$

## Propiedades de las p-normas

$$1. \quad \|A.B\|_p \leq \|A\|_p \cdot \|B\|_p$$

$$2. \quad \|Ax\|_p \leq \|A\|_p \cdot \|x\|_p$$

# Algunas propiedades de las normas matriciales

Si  $A \in \mathbb{R}^{m \times n}$  :

$$i) \quad \|A\|_2 \leq \|A\|_F \leq \sqrt{n} \|A\|_2$$

$$ii) \quad \max_{i,j} |a_{ij}| \leq \|A\|_2 \leq \sqrt{n} \max_{i,j} |a_{ij}|$$

$$iii) \quad \|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^m |a_{ij}|$$

$$iv) \quad \|A\|_\infty = \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}|$$

$$v) \quad \frac{1}{\sqrt{n}} \|A\|_\infty \leq \|A\|_2 \leq \sqrt{m} \|A\|_\infty$$

$$vi) \quad \frac{1}{\sqrt{m}} \|A\|_1 \leq \|A\|_2 \leq \sqrt{n} \|A\|_1$$

$$vii) \quad \|A\|_2 \leq \sqrt{\|A\|_1 \|A\|_\infty}$$

**Todas las normas matriciales en  $\mathbb{R}^{m \times n}$  son equivalentes**

Límite de una sucesión de matrices,  $\{A^{(k)}\}$ :

$$\{A^{(k)}\} \text{ converge si } \lim_{k \rightarrow \infty} \|A^{(k)} - A\| = 0$$

# Ejercicios propuestos

- 1. Probar que las transformaciones ortogonales conservan la 2-norma vectorial:

Si  $v \in \mathbb{R}^n$ , y  $Q \in \mathbb{R}^{n \times n}$  es ortogonal ( $Q^T Q = I$ ),  
se verifica que :  $\|v\|_2 = \|Qv\|_2$

- 2. Probar que las transformaciones ortogonales conservan la norma de Frobenius y la 2-norma matricial:

Sea  $A \in \mathbb{R}^{m \times n}$ . Si  $P \in \mathbb{R}^{m \times m}$  y  $Q \in \mathbb{R}^{n \times n}$  son matrices ortogonales,  
se verifica que :  $i) \|A\|_F = \|PAQ\|_F$        $ii) \|A\|_2 = \|PAQ\|_2$