



Profesor: Gonzalo Hernández.

Auxiliar: Gonzalo Ríos.

Fecha: 15 de Abril

Pauta Control 1

1. Codificación binaria floating point tipo inicial

- (a) Como el número real máximo representable en una codificación binaria siempre es $\frac{2^m-1}{2^m} \times 2^{(2^e-1)}$, donde m es la cantidad de bits de la mantisa, y e la cantidad de bits del exponente, esto equivale a $2^{2^e-1} - 2^{(2^e-1)-m}$, queda un sistema:

$$2^{63} - 2^{54} = 2^{2^e-1} - 2^{(2^e-1)-m} \implies \begin{matrix} 2^e - 1 = 63 \\ (2^e - 1) - m = 54 \end{matrix} \implies \begin{matrix} e = 6 \\ m = 9 \end{matrix}$$

Sumando el bit para el signo de la mantisa y el bit del signo del exponente, la codificación necesita $2 + 6 + 9 = 17$ bits.

- (b) **Como son 17 bits, entonces**

i. Numero en codificación binaria: $a = a_1 a_2 a_3 \dots a_{17} = 11011101110111011$

ii. Número real:

A. Signo Mantisa: —

B. Signo Exponente: —

C. Exponente: $0 \times 2^5 + 1 \times 2^4 + 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 = 29$

D. Mantisa: $1 \times 2^{-1} + 1 \times 2^{-2} + 0 \times 2^{-3} + 1 \times 2^{-4} + 1 \times 2^{-5} + 1 \times 2^{-6} + 0 \times 2^{-7} + 1 \times 2^{-8} + 1 \times 2^{-9} = \frac{443}{512} = 0.865234375$

E. Numero: $a = -\frac{443}{512} \times 2^{-29} = -443 \times 2^{-38} = -1.6116246115416288376 \times 10^{-9}$

- (c) **Polinomio de Taylor**

i. $p_4(x) = 1 + x^2 + \frac{1}{2}x^4$

ii. Se debe calcular con la aritmética finita de 5 cifras con redondeo, entonces $\frac{0.0625}{2} = 0.03125$, cumple la aritmética. Luego $p_4(0.5) = 1 + (0.5)^2 + \frac{1}{2}(0.5)^4 = 1 + 0.25 + \frac{0.0625}{2} = 1.25 + 0.03125 = 1.28125$, reduciendo a 5 cifras con redondeo, queda $p_4(0.5) = 1.2813$

iii. En representación de punto flotante quedan:

A. $\boxed{f(0.5)}$ $1.2840 = 0.642 \times 2^1$

| | | | | | | | | |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| 2^{-1} | 2^{-2} | 2^{-3} | 2^{-4} | 2^{-5} | 2^{-6} | 2^{-7} | 2^{-8} | 2^{-9} |
| 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |

0.642 0.142 0.017 0.001375

$f(0.5) = 00000001101001000$

B. $\boxed{p_4(0.5)}$ $1.2813 = 0.64065 \times 2^1$

| | | | | | | | | |
|----------|----------|----------|----------|----------|----------|----------|----------|----------|
| 2^{-1} | 2^{-2} | 2^{-3} | 2^{-4} | 2^{-5} | 2^{-6} | 2^{-7} | 2^{-8} | 2^{-9} |
| 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |

0.64065 0.14065 0.01565 0.000025

$p_4(0.5) = 00000001101001000$

C. $E_{rel} = \frac{|f(0.5) - p_4(0.5)|}{|f(0.5)|} = \frac{0.0027}{1.2840} = 2.1028037383177570093 \times 10^{-3}$

La representación de $f(0.5)$ y de $p_4(0.5)$ en esta decodificación es exactamente la misma, por lo que quiere decir que la aproximación, en términos del punto flotante, es muy buena.

2. Serie, suma parcial exacta y suma parcial aproximada

- (a) $E_{abs} = |S - Z_n| = |S - Z_n + S_n - S_n| = |S - S_n + S_n - Z_n| \leq |S - S_n| + |S_n - Z_n|$

El primer término corresponde al error de truncación que tiende a cero cuando $n \rightarrow \infty$ ya que la suma parcial tiende a la serie. El segundo término corresponde al error de redondeo, que aumenta cuando $n \rightarrow \infty$, ya que al aumentar el número de operaciones, se cometen más errores de redondeo.

- (b) $E_{abs} = |S - Z_n| \leq |S - S_n| + |S_n - Z_n| \leq \int_n^\infty \frac{1}{x^2} dx + C \cdot n = -\frac{1}{x} \Big|_n^\infty + C \cdot n = \frac{1}{n} + C \cdot n = g(n)$

$$(c) \min g(n) = \frac{1}{n} + C \cdot n \iff \frac{dg}{dn}(\bar{n}) = 0 \wedge \frac{d^2g}{dn^2}(\bar{n}) > 0$$

$$\frac{dg}{dn} = -\frac{1}{n^2} + C = 0 \implies \bar{n}_1 = \frac{1}{\sqrt{C}} \wedge \bar{n}_2 = -\frac{1}{\sqrt{C}}$$

$$\frac{d^2g}{dn^2} = 2\frac{1}{n^3} > 0 \implies \boxed{\bar{n} = \frac{1}{\sqrt{C}}}$$

$$n = \frac{1}{\sqrt{4 \times 10^{-16}}} = 5 \times 10^7$$

$$E_{abs} \leq \frac{1}{n} + C \cdot n = \sqrt{C} + \frac{C}{\sqrt{C}} = 2\sqrt{C} = 2\sqrt{4 \times 10^{-16}} = 4 \times 10^{-8}$$

Lo que implica que tendrá 7 dígitos decimales correctos.

$$(d) \phi(x, y) = g(x) \times g(y) = \left(\frac{1}{x} + cx\right) \left(\frac{1}{y} + cy\right) = \frac{1}{xy} + c^2xy + c\left(\frac{x}{y} + \frac{y}{x}\right)$$

$$\varepsilon\phi = \frac{x}{\phi(x,y)} \frac{\partial\phi}{\partial x} \varepsilon_x + \frac{y}{\phi(x,y)} \frac{\partial\phi}{\partial y} \varepsilon_y$$

$$\frac{\partial\phi}{\partial x} = \frac{-1}{x^2y} + c^2y + c\left(\frac{1}{y} - \frac{y}{x^2}\right)$$

$$x \frac{\partial\phi}{\partial x} = c^2xy + c\left(\frac{x}{y} - \frac{y}{x}\right) - \frac{1}{xy} = \frac{c^2x^2y^2 + c(x^2 - y^2) - 1}{xy}$$

$$\phi(x, y) = \frac{c^2x^2y^2 + c(x^2 + y^2) + 1}{xy}$$

$$\begin{aligned} \frac{x}{\phi(x,y)} \frac{\partial\phi}{\partial x} &= \frac{c^2x^2y^2 + c(x^2 - y^2) - 1}{c^2x^2y^2 + c(x^2 + y^2) + 1} = \frac{c^2x^2y^2 + c(x^2 - y^2) - 1 + (2-2) + (2y^2 - 2y^2)}{c^2x^2y^2 + c(x^2 + y^2) + 1} = \frac{c^2x^2y^2 + c(x^2 + y^2) + 1 - 2 - 2y^2}{c^2x^2y^2 + c(x^2 + y^2) + 1} \\ &= 1 - \frac{2 + 2y^2}{c^2x^2y^2 + c(x^2 + y^2) + 1} = 1 - 2 \frac{1 + y^2}{c^2x^2y^2 + c(x^2 + y^2) + 1} = 1 - \frac{2}{c^2x^2y^2 + c(x^2 + y^2) + 1} - 2 \frac{y^2}{c^2x^2y^2 + c(x^2 + y^2) + 1} = \\ &= 1 - \frac{2}{c^2x^2y^2 + c(x^2 + y^2) + 1} - 2 \frac{1}{c^2x^2 + c\left(\frac{x^2 + y^2}{y^2}\right) + \frac{1}{y^2}} \end{aligned}$$

entonces $\frac{2}{c^2x^2y^2 + c(x^2 + y^2) + 1}$ es acotada, ya que $c > 0$ lo que implica que $c^2x^2y^2 + c(x^2 + y^2) + 1 >$

$0 \forall x, y$ Ahora si $\frac{1}{c^2x^2 + c\left(\frac{x^2 + y^2}{y^2}\right) + \frac{1}{y^2}}$ no fuera acotada, entonces $\exists x, y$ tal que $c^2x^2 + c\left(\frac{x^2 + y^2}{y^2}\right) +$

$\frac{1}{y^2} = 0$, pero $c^2x^2 \geq 0$, $c\left(\frac{x^2 + y^2}{y^2}\right) \geq 0$ y $\frac{1}{y^2} \geq 0$, entonces la única opción sería que las tres fueran igual a cero, y eligiendo $x = 0$ para que se cumpla la primera igualdad, la segunda se convierte en c , que por enunciado es > 0 . Esto implica que $\frac{1}{c^2x^2 + c\left(\frac{x^2 + y^2}{y^2}\right) + \frac{1}{y^2}}$ es acotada, y la

suma de expresiones acotadas es acotada, luego se puede acotar $\frac{x}{\phi(x,y)} \frac{\partial\phi}{\partial x}$. De forma análoga se acota $\frac{y}{\phi(x,y)} \frac{\partial\phi}{\partial y}$.

3. SEL

(a) Métodos Directos

$$i. A^0 = \begin{bmatrix} 4 & 2 & 0 \\ 2 & 4 & 2 \\ 0 & 2 & 4 \end{bmatrix}$$

$$A^1 = \begin{bmatrix} 1 & 0 & 0 \\ -\frac{2}{4} & 1 & 0 \\ -\frac{0}{4} & 0 & 1 \end{bmatrix} \begin{bmatrix} 4 & 2 & 0 \\ 2 & 4 & 2 \\ 0 & 2 & 4 \end{bmatrix} = \begin{bmatrix} 4 & 2 & 0 \\ 0 & 3 & 2 \\ 0 & 2 & 4 \end{bmatrix}$$

$$A^2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & -\frac{2}{3} & 1 \end{bmatrix} \begin{bmatrix} 4 & 2 & 0 \\ 0 & 3 & 2 \\ 0 & 2 & 4 \end{bmatrix} = \begin{bmatrix} 4 & 2 & 0 \\ 0 & 3 & 2 \\ 0 & 0 & \frac{8}{3} \end{bmatrix}$$

$$U = \begin{bmatrix} 4 & 2 & 0 \\ 0 & 3 & 2 \\ 0 & 0 & \frac{8}{3} \end{bmatrix} \quad L = \begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ 0 & \frac{2}{3} & 1 \end{bmatrix}$$

Definiendo $U\vec{x} = \vec{y}$ el sistema queda

$$\begin{bmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 1 & 0 \\ 0 & \frac{2}{3} & 1 \end{bmatrix} \vec{y} = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \implies \vec{y} = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{2}{3} \end{bmatrix}$$

Resolviendo $U\vec{x} = \vec{y}$

$$\begin{bmatrix} 4 & 2 & 0 \\ 0 & 3 & 2 \\ 0 & 0 & \frac{8}{3} \end{bmatrix} \vec{x} = \begin{bmatrix} 1 \\ \frac{1}{2} \\ \frac{2}{3} \end{bmatrix} \implies \vec{x} = \begin{bmatrix} \frac{1}{4} \\ 0 \\ \frac{1}{4} \end{bmatrix}$$

ii. Es definida positiva si sus subdeterminantes son positivas

$$|4| > 0$$

$$\begin{vmatrix} 4 & 2 \\ 0 & 3 \end{vmatrix} = 12 > 0$$

$$\begin{vmatrix} 4 & 2 & 0 \\ 0 & 3 & 2 \\ 0 & 0 & \frac{8}{3} \end{vmatrix} = 32 > 0 \implies \text{es definida positiva}$$

Entonces su factorizacion de Cholesky es

$$\begin{bmatrix} 4 & 2 & 0 \\ 2 & 4 & 2 \\ 0 & 2 & 4 \end{bmatrix} = \begin{bmatrix} 2 & 0 & 0 \\ 1 & \sqrt{3} & 0 \\ 0 & \frac{2}{3}\sqrt{3} & \frac{2}{3}\sqrt{2}\sqrt{3} \end{bmatrix} \begin{bmatrix} 2 & 1 & 0 \\ 0 & \sqrt{3} & \frac{2}{3}\sqrt{3} \\ 0 & 0 & \frac{2}{3}\sqrt{2}\sqrt{3} \end{bmatrix}$$

(b) **Métodos Iterativos**

i. Gauss-Seidel

$$M = -(D + L)^{-1}U = -\begin{bmatrix} 4 & 0 & 0 \\ 2 & 4 & 0 \\ 0 & 2 & 4 \end{bmatrix}^{-1} \begin{bmatrix} 0 & 2 & 0 \\ 0 & 0 & 2 \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & -\frac{1}{2} & 0 \\ 0 & \frac{1}{4} & -\frac{1}{2} \\ 0 & -\frac{1}{8} & \frac{1}{4} \end{bmatrix}$$

$$\vec{N} = (D + L)^{-1}\vec{b} = \begin{bmatrix} 4 & 0 & 0 \\ 2 & 4 & 0 \\ 0 & 2 & 4 \end{bmatrix}^{-1} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} \frac{1}{4} \\ \frac{1}{8} \\ \frac{3}{16} \end{bmatrix}$$

$$\vec{x}^{k+1} = M\vec{x}^k + \vec{N} = \begin{bmatrix} 0 & -\frac{1}{2} & 0 \\ 0 & \frac{1}{4} & -\frac{1}{2} \\ 0 & -\frac{1}{8} & \frac{1}{4} \end{bmatrix} \begin{bmatrix} x_1^k \\ x_2^k \\ x_3^k \end{bmatrix} + \begin{bmatrix} \frac{1}{4} \\ \frac{1}{8} \\ \frac{3}{16} \end{bmatrix} = \begin{bmatrix} -\frac{1}{2}x_2^k + \frac{1}{4} \\ \frac{1}{4}x_2^k - \frac{1}{2}x_3^k + \frac{1}{8} \\ -\frac{1}{8}x_2^k + \frac{1}{4}x_3^k + \frac{3}{16} \end{bmatrix}$$

$$\text{A. } \vec{x}^0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{B. } \vec{x}^1 = \begin{bmatrix} \frac{1}{4} \\ \frac{1}{8} \\ \frac{3}{16} \end{bmatrix}$$

$$\text{C. } \vec{x}^2 = \begin{bmatrix} \frac{3}{16} \\ \frac{1}{16} \\ \frac{7}{32} \end{bmatrix}$$

$$\text{D. } \vec{x}^3 = \begin{bmatrix} \frac{7}{32} \\ \frac{1}{16} \\ \frac{15}{64} \end{bmatrix}$$

$$\text{E. } \vec{x}^4 = \begin{bmatrix} \frac{15}{64} \\ \frac{1}{64} \\ \frac{31}{128} \end{bmatrix} = \begin{bmatrix} 0.234375 \\ 0.015625 \\ 0.2421875 \end{bmatrix}$$

$$\text{F. } \vec{x}^5 = \begin{bmatrix} \frac{31}{128} \\ \frac{1}{128} \\ \frac{63}{256} \end{bmatrix} = \begin{bmatrix} 0.2421875 \\ 0.0078125 \\ 0.24609375 \end{bmatrix}$$

Solo 4 iteraciones eran necesarias, es decir, hasta \vec{x}^4 .

ii. $\rho(T_G) = \max_{i=1, \dots, n} |\lambda_i|$ donde λ_i es un valor propio de T_G

$$\det(T_G - I\lambda) = 0$$

$$\begin{vmatrix} -\lambda & -\frac{1}{2} & 0 \\ 0 & \frac{1}{4} - \lambda & -\frac{1}{2} \\ 0 & -\frac{1}{8} & \frac{1}{4} - \lambda \end{vmatrix} = \frac{1}{2}\lambda^2 - \lambda^3 = 0 \implies \lambda_1 = \lambda_2 = 0 \wedge \lambda_3 = \frac{1}{2} \implies \rho(T_G) = \frac{1}{2}$$

$$\bar{\omega} = \frac{2}{1 + \sqrt{1 - \rho(T_G)}} = \frac{2}{1 + \sqrt{1 - \frac{1}{2}}} = \frac{2}{\frac{1}{2}\sqrt{2} + 1} = 4 - 2\sqrt{2} = 2(2 - \sqrt{2}) = 1.171572875$$

$$2538099024$$

iii. Sor

$$\vec{x}^{k+1} = (1 - \bar{\omega})\vec{x}^k + \bar{\omega}(M\vec{x}^k + \vec{N}) = (1 - 1.1715) \begin{bmatrix} x_1^k \\ x_2^k \\ x_3^k \end{bmatrix} + 1.1715 \begin{bmatrix} 0 & -\frac{1}{2} & 0 \\ 0 & \frac{1}{4} & -\frac{1}{2} \\ 0 & -\frac{1}{8} & \frac{1}{4} \end{bmatrix} \begin{bmatrix} x_1^k \\ x_2^k \\ x_3^k \end{bmatrix} + \begin{bmatrix} \frac{1}{4} \\ \frac{1}{8} \\ \frac{3}{16} \end{bmatrix} =$$

$$-0.1715 \begin{bmatrix} x_1^k \\ x_2^k \\ x_3^k \end{bmatrix} + 1.1715 \begin{bmatrix} -\frac{1}{2}x_2^k + \frac{1}{4} \\ \frac{1}{4}x_2^k - \frac{1}{2}x_3^k + \frac{1}{8} \\ -\frac{1}{8}x_2^k + \frac{1}{4}x_3^k + \frac{3}{16} \end{bmatrix} = \begin{bmatrix} -0.1715x_1^k - 0.58575x_2^k + 0.292875 \\ 0.121375x_2^k - 0.58575x_3^k + 0.1464375 \\ -0.1464375x_2^k + 0.121375x_3^k + 0.21965625 \end{bmatrix}$$

$$\text{A. } \vec{x}^0 = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\text{B. } \vec{x}^1 = (1 - 1.1715) \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + 1.1715 \begin{bmatrix} 0 & -\frac{1}{2} & 0 \\ 0 & \frac{1}{4} & -\frac{1}{2} \\ 0 & -\frac{1}{8} & \frac{1}{4} \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + \begin{bmatrix} \frac{1}{4} \\ \frac{1}{8} \\ \frac{3}{16} \end{bmatrix} = 1.1715 \begin{bmatrix} \frac{1}{4} \\ \frac{1}{8} \\ \frac{3}{16} \end{bmatrix} =$$

$$\begin{bmatrix} 0.292875 \\ 0.1464375 \\ 0.21965625 \end{bmatrix}$$

$$\text{C. } \vec{x}^2 = (1 - 1.1715) \begin{bmatrix} 0.292875 \\ 0.1464375 \\ 0.21965625 \end{bmatrix} + 1.1715 \begin{bmatrix} \begin{bmatrix} 0 & -\frac{1}{2} & 0 \\ 0 & \frac{1}{4} & -\frac{1}{2} \\ 0 & -\frac{1}{8} & \frac{1}{4} \end{bmatrix} \begin{bmatrix} 0.292875 \\ 0.1464375 \\ 0.21965625 \end{bmatrix} + \begin{bmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{3}{16} \end{bmatrix} \end{bmatrix} =$$

$$\text{D. } \vec{x}^3 = (1 - 1.1715) \begin{bmatrix} 0.156871171875 \\ 0.035547703125 \\ 0.2248730859375 \end{bmatrix} + 1.1715 \begin{bmatrix} \begin{bmatrix} 0 & -\frac{1}{2} & 0 \\ 0 & \frac{1}{4} & -\frac{1}{2} \\ 0 & -\frac{1}{8} & \frac{1}{4} \end{bmatrix} \begin{bmatrix} 0.156871171875 \\ 0.035547703125 \\ 0.2248730859375 \end{bmatrix} + \begin{bmatrix} \frac{1}{4} \\ \frac{1}{4} \\ \frac{3}{16} \end{bmatrix} \end{bmatrix} =$$

Hasta acá es necesario, es decir, 3 iteraciones