

## CAPÍTULO 1

# Aritmética de Punto Flotante

### 1. Números de Punto Flotante

Un conjunto  $F$  de números de *punto flotante* está caracterizado por los siguientes parámetros:

- (1) La base del sistema  $\beta$ .
- (2) El número de dígitos  $n$  en la mantisa.
- (3) Un exponente  $m \leq e \leq M$ .

donde  $\beta, n, m, e$  son enteros.

Cada número en el sistema de punto flotante  $F$  tiene la forma

$$\pm(.d_1 d_2 \cdots d_n)_\beta \beta^e,$$

donde  $d_i = 0, 1, \dots, \beta - 1$ ,  $i = 1, 2, \dots, n$ ,  $(.d_1 d_2 \cdots d_n)_\beta$  es una  $\beta$ -fracción llamada *mantisa*,  $e$  es un entero llamado el *exponente*.

El sistema  $F$  de punto flotante está *normalizado* si  $d_1 \neq 0$ . en general, todos los números flotantes se normalizan con la excepción del cero, en el cual  $d_1 = d_2 = \cdots = d_n = 0$ .

El conjunto  $F$  es discreto y finito. Su cardinalidad, el número de elementos que lo constituyen, está dada por

$$2(\beta - 1)\beta^{n-1}(M - m + 1) + 1.$$

Como una consecuencia de la finitud de  $F$  para representar al conjunto de números reales  $\mathbb{R}$ , existirá una infinidad de números en  $\mathbb{R}$  que no pueden representarse en forma exacta en  $F$ .

Sea  $x$  un número real denotemos por  $\text{fl}(x)$  el número en  $F$  que es más *cercano* a  $x$ . La diferencia entre  $x$  y  $\text{fl}(x)$  se llama *error de redondeo*, éste depende de la magnitud de  $x$  y es por lo tanto medido relativo a  $x$

$$(1.1) \quad \delta(x) = \frac{\text{fl}(x) - x}{x}, \quad x \neq 0,$$

luego  $|\delta(x)|$  es el error relativo introducido en la representación de  $x$  en el sistema de punto flotante  $F$ . De (1.1) obtenemos

$$\text{fl}(x) = x(1 + \delta(x)).$$

Hay dos formas generalmente usadas para convertir un número real  $x$  a un  $n$ - $\beta$  número flotante  $\text{fl}(x)$ : *redondeado* o *truncado*. Cuando se redondea,  $\text{fl}(x)$  se elige como el número de punto flotante normalizado más cercano a  $x$ , si hay empate se usa alguna regla especial, por ejemplo, se toma el de la derecha. Si se trunca,  $\text{fl}(x)$  se escoge como el número flotante normalizado más cercano entre 0 y  $x$ , esto es, se toman algunas  $d$ 's y se desprecian otras.

A continuación encontraremos una cota para  $\delta(x)$  que sea independiente de  $x$ . En el sistema de números de punto flotante  $F$ , el número cuyo valor absoluto es el más pequeño está dado por

$$+ (.100 \dots 0)_\beta \beta^m = \beta^{m-1},$$

el sucesor inmediato de  $\beta^{m-1}$  se encuentra sumándole a éste el número

$$+ (.000 \dots 1)_\beta \beta^m = \beta^{m-n}.$$

Se concluye entonces, que la distancia entre dos números consecutivos en el intervalo  $[\beta^{m-1}, \beta^m]$  es  $\beta^{m-n}$ .

En forma similar se demuestra que la distancia entre dos números consecutivos que pertenezcan a cualquier intervalo de la forma  $[\beta^j, \beta^{j+1}]$   $j = m, \dots, M-1$  está dada por

$$(1.2) \quad \beta^{j+1-n}.$$

Hemos demostrado un aspecto singular de la distribución de los números del sistema de punto flotante  $F$ :

- *estos no están igualmente espaciados a través de todo su rango, sino únicamente cuando se encuentran entre potencias sucesivas de la base  $\beta$ .*

Supongamos que  $x \in [\beta^j, \beta^{j+1}]$  para alguna  $j = m, \dots, M-1$ , si el número flotante  $\text{fl}(x)$  que representa a  $x$  es seleccionado por redondeo, entonces de acuerdo a (1.2) el error introducido es a lo más  $(1/2)\beta^{j+1-n}$ , si  $\text{fl}(x)$  se selecciona por truncamiento el error es a lo más  $\beta^{j+1-n}$ . Lo anterior nos da la *medida del error de redondeo absoluto*

$$(1.3) \quad |\text{fl}(x) - x| \leq \begin{cases} \frac{1}{2}\beta^{j+1-n} & \text{redondeo,} \\ \beta^{j+1-n} & \text{truncamiento.} \end{cases}$$

El *error de redondeo relativo*  $|\delta(x)|$  se obtiene dividiendo al error de redondeo absoluto por  $|x|$ . Ya que  $0 < \beta^j \leq |x|$  obtenemos

$$|\delta(x)| \leq \begin{cases} \frac{1}{2}\beta^{1-n} & \text{redondeo,} \\ \beta^{1-n} & \text{truncamiento.} \end{cases}$$

El *épsilon de la máquina* se define como

$$\varepsilon = \begin{cases} \frac{1}{2}\beta^{1-n} & \text{redondeo,} \\ \beta^{1-n} & \text{truncamiento.} \end{cases}$$

Consecuente con lo anterior se cumple que  $|\delta(x)| \leq \varepsilon$  para toda  $x$ . La cota para  $\delta(x)$  independiente de la  $x$  es el *épsilon de la máquina*.

Dado  $x \in \mathbb{R}$  su flotante  $\text{fl}(x)$  se define como

$$\text{fl}(x) = x(1 + \delta), \quad |\delta| \leq \varepsilon.$$

La exactitud de la aritmética de punto flotante está entonces caracterizada por el *épsilon de la máquina*.

Supongamos que estamos en el intervalo  $[1, \beta]$ . Queremos calcular  $1 \oplus \varepsilon$ , donde  $\oplus$  denota la suma entre números que pertenecen a  $F$

$$1 \oplus \varepsilon = \text{fl}(1 + \varepsilon).$$

Si el error introducido en la representación de  $1 + \varepsilon$  es por redondeo, entonces éste queda localizado en la mitad del intervalo  $[1, 1 + \beta^{1-n}]$ . Ya que la distancia entre dos números consecutivos de  $F$  que pertenecen al intervalo  $[1, \beta]$  es  $\beta^{1-n}$  el número de punto flotante  $1 \oplus \varepsilon$  es tomado como  $1 + \beta^{1-n} > 1$ . Sin embargo, si  $0 < \varepsilon_1 < \varepsilon$  obtenemos por un procedimiento análogo al anterior que  $1 \oplus \varepsilon_1 = 1$ . Un resultado similar es encontrado cuando el error introducido en la representación de  $1 \oplus \varepsilon$  es por truncamiento.

---

El épsilon de la máquina se caracteriza por ser el menor número positivo de  $F$  que satisface

$$1 \oplus \varepsilon > 1.$$


---

## 2. Modelo de Aritmética en $F$

La aritmética en el sistema numérico de punto flotante  $F$  permite aproximar a la del sistema de números reales  $\mathbb{R}$ . Como notación emplearemos  $\oplus, \ominus, \otimes, \oslash$  para indicar las aproximaciones a las operaciones aritméticas  $+, -, \times, /$  de  $\mathbb{R}$ :

$$x \oplus y = \text{fl}(x + y),$$

$$x \ominus y = \text{fl}(x - y),$$

$$x \otimes y = \text{fl}(x \times y),$$

$$x \oslash y = \text{fl}(x/y),$$

El modelo que asumiremos para la aritmética en  $F$  es el siguiente

$$\text{fl}(x \text{ op } y) = (x \text{ op } y)(1 + \delta), \quad |\delta| \leq \varepsilon, \quad \text{op} = +, -, *, /.$$

Para efectuar operaciones en forma manual en este modelo aritmético, por cada operación  $+, -, \times, /$  encontrada hágala en aritmética exacta, normalice el resultado, trunque o redondee de acuerdo al número de dígitos permitidos.