



ESCUELA POLITÉCNICA NACIONAL
FACULTAD DE INGENIERÍA DE SISTEMAS
INGENIERÍA CIENCIAS DE LA COMPUTACIÓN

PERÍODO ACADÉMICO: 2025-A

ASIGNATURA: ICCD412 Métodos Numéricos

GRUPO: GR2

TIPO DE INSTRUMENTO: Práctica 02

FECHA DE ENTREGA LÍMITE: 04/05/2025

ALUMNO: Freire Ismael

TEMA

Representación numérica 64 bits

OBJETIVOS

- Comprender la representación de números de punto flotante aplicando el Estándar IEEE 754 de 64 bits en la computadora mediante un ejercicio práctico.

MARCO TEÓRICO

Dentro de la computación, ha sido un gran desafío el lograr representar adecuadamente la precisión y exactitud en los diferentes cálculos que se realicen. Por eso, surgió el Estándar IEEE 754 en 1985 por la misma institución del cual lleva su nombre, IEEE (Instituto de Ingenieros Eléctricos y Electrónicos, por sus siglas en inglés). Este estándar revisado por la misma entidad en el año 2008 se estableció como la manera de representación y manipulación de números de punto flotante a nivel mundial en sistemas digitales, esto con la finalidad de garantizar una uniformidad en los sistemas computacionales.

El estándar define varios formatos precisos para dicha representación de números, que se diferencian entre ellos por la cantidad de bits utilizados. Existen tres formatos los cuales mantienen una estructura similar en cuanto a los bits para el signo, exponente y mantisa, pero cuyas diferencias se detallarán más adelante. Primero, se debe tener claro que el signo indica si el número a representar es positivo o negativo, el exponente se usa para escalar el número y la mantisa corresponde a la parte fraccionaria de este, esta es la que tiene un solo dígito distinto de cero en su parte entera [1].

A continuación, se indican los formatos:

1. Precisión simple (32 bits) con 1 bit para el signo, 8 bits para el exponente y 23 bits para la mantisa.
2. Precisión doble (64 bits) con 1 bit para el signo, 11 bits para el exponente y 52 bits para la mantisa.
3. Precisión quad (128 bits) con 1 bit para el signo, 15 bits para el exponente y 112 bits para la mantisa. Este formato no se trata en el resto de la práctica.

Por como se ha mostrado, el estándar maneja distintos formatos los cuales agregan mayor exactitud al número flotante, junto con un aumento en la cantidad de bits que van utilizando. Como tal, este estándar tiene aplicaciones en los lenguajes de programación, algoritmos y hasta cálculos científicos. Por ello, su implementación asegura resultados confiables y consistentes.

DESARROLLO

1. **Pasar el número $263,3_{10}$ al formato IEEE 754 de 64 bits, una vez en binario pasarlo a decimal y calcular el error relativo a 3 cifras de redondeo.**

| | | | |
|---|-----------------|--|-------------------|
| Transformar $263,3_{10} \Rightarrow$ IEEE 754 en 64 bits. | | | |
| $263,3 \div 2 = 131,5$ | $\rightarrow 1$ | | |
| $131,5 \div 2 = 65,75$ | $\rightarrow 1$ | | |
| $65,75 \div 2 = 32,875$ | $\rightarrow 1$ | | $263 = 100000111$ |
| $32,875 \div 2 = 16,4375$ | $\rightarrow 0$ | | |
| $16,4375 \div 2 = 8,21875$ | $\rightarrow 0$ | | |
| $8,21875 \div 2 = 4,109375$ | $\rightarrow 0$ | | |
| $4,109375 \div 2 = 2,0546875$ | $\rightarrow 0$ | | |
| $2,0546875 \div 2 = 1,02734375$ | $\rightarrow 0$ | | |
| $1,02734375 \div 2 = 0,513671875$ | $\rightarrow 1$ | | |
| Parte fraccionaria: | | | |
| $0,3 \times 2 = 0,6$ | $\rightarrow 0$ | | |
| $0,6 \times 2 = 1,2$ | $\rightarrow 1$ | | |
| $0,2 \times 2 = 0,4$ | $\rightarrow 0$ | | |
| $0,4 \times 2 = 0,8$ | $\rightarrow 0$ | | $0,3 = 010011$ |
| $0,8 \times 2 = 1,6$ | $\rightarrow 1$ | | |
| $0,6 \times 2 = 1,2$ | $\rightarrow 1$ | | |
| $0,2 \times 2 = 0,4$ | $\rightarrow 0$ | | |
| $0,4 \times 2 = 0,8$ | $\rightarrow 0$ | | |
| $0,8 \times 2 = 1,6$ | $\rightarrow 1$ | | |
| $0,6 \times 2 = 1,2$ | $\rightarrow 1$ | | |
| $0,2 \times 2 = 0,4$ | $\rightarrow 0$ | | |

Figura 1: P1 Transformación a binario

| | |
|------------------------------------|--------------------------------|
| P2 Escribir en notación científica | |
| $263,3 =$ | $100000111,010011$ |
| $263,3 =$ | $1,00000111010011 \times 10^8$ |

Figura 2: P2 Notación Científica

Resolución con el estándar IEEE 754 64 bits:

→ Transformación al valor decimal

$$\begin{aligned}
 x &= (-1)^s & 2^{(E-1023)} & & (1+f) \\
 &= (-1)^0 & & & f = \sum_{i=1}^{52} (f_i \times \frac{1}{2^i}) \\
 &= 1 & = 2^{(1031-1023)} & & \\
 & & = 2^8 & & \\
 & & = 256 & & \\
 & & & & f = 2^{-6} + 2^{-7} + 2^{-8} + 2^{-10} + 2^{-13} + 2^{-14} + 2^{-17} + 2^{-18} \\
 & & & & + 2^{-21} + 2^{-22} + 2^{-25} + 2^{-26} + 2^{-29} + 2^{-30} + 2^{-33} + 2^{-34} \\
 & & & & + 2^{-37} + 2^{-38} + 2^{-41} + 2^{-42} + 2^{-45} + 2^{-46} + 2^{-49} \\
 & & & & + 2^{-50} = 0,028515625 \\
 & & & & (1 + 0,028515625) = 1,028515625 \\
 x &= 1 \times 256 \times 1,028515625 \\
 &= 263,3
 \end{aligned}$$

Figura 5: P4 Transformación del número a decimal

Cálculo del error relativo:

$$\begin{aligned}
 \text{error rel} &= \left| \frac{p - p^*}{p} \right| \\
 &= \left| \frac{263,3 - 263,3}{263,3} \right| \\
 &= 0
 \end{aligned}$$

Figura 6: Cálculo error relativo

CONCLUSIONES

En el ejercicio, se obtuvo un error relativo igual a cero lo que sugiere que se tuvo una representación binaria exacta con los bits disponibles en la mantisa. Aun así, se debe

tener en cuenta que generalmente el error está cerca de ser muy pequeño, pero no necesariamente cero. Por ello, hay que tener precaución a la hora de aplicar el estándar y de ser posible validar el resultado obtenido con el cálculo del error.

Por otro lado, gracias a que el método de doble precisión posee una mayor cantidad de bits en su mantisa permite que se pueda manejar un amplio rango numérico a diferencia del formato de 32 bits que se analizó en clase. Es por ello por lo que es perfecto para la gran mayoría de aplicaciones y usos en áreas como la computación.

RECOMENDACIONES

Teniendo en cuenta que la cantidad de bits hace que la representación se vuelva mas exacta, esto no quiere decir que sea infalible. Por lo que siempre es importante tener presente la existencia de números que no necesariamente tienen una representación binaria exacta, lo cuales pueden generar errores que deben ser evaluados para comprobar que el método sea el adecuado y evitar problemas con los cálculos.

REFERENCIAS

- [1] H. A. V. Martínez, “Representación de números en binario,” pp. 8–11, 2008.