

Representación de Números en Coma Flotante

Elementos de Programación y Cálculo Numérico - 2025

Material de Estudio - Optativo

1. Introducción

En el desarrollo de métodos numéricos para la resolución de problemas de ingeniería, uno de los aspectos fundamentales que debe comprenderse es cómo los sistemas computacionales representan y manipulan los números reales. Esta comprensión es crucial porque determina tanto la precisión como la estabilidad de los algoritmos implementados.

Los números reales, tal como los concebimos matemáticamente, forman un conjunto infinito y denso. Sin embargo, los sistemas digitales deben trabajar con representaciones finitas y discretas. Esta limitación fundamental introduce inevitablemente errores y restricciones que deben ser considerados cuidadosamente en el diseño e implementación de algoritmos numéricos.

La representación en coma flotante constituye el estándar universal para el manejo de números reales en sistemas computacionales modernos. Esta representación, formalizada por la norma IEEE 754 en 1985 y adoptada ampliamente desde entonces, permite trabajar con un amplio rango de magnitudes manteniendo una precisión relativa aproximadamente constante, característica que resulta especialmente útil para aplicaciones científicas y de ingeniería.

La adopción del estándar IEEE 754 representó un hito crucial en la computación científica, pues antes de su establecimiento, diferentes fabricantes de computadoras utilizaban representaciones propias e incompatibles. Esta falta de estandarización generó problemas serios, incluyendo casos catastróficos como el del cohete Ariane 5 en 1996, donde una conversión incorrecta de formato numérico causó la destrucción del vehículo espacial con pérdidas superiores a los 500 millones de dólares. La uniformidad introducida por IEEE 754 garantizó la portabilidad de software y la reproducibilidad de resultados entre diferentes plataformas, aspectos fundamentales para la confiabilidad de los cálculos científicos y de ingeniería.

2. Fundamentos de la Representación Digital

2.1. Conceptos Básicos: Bit y Byte

Antes de abordar la representación en coma flotante, es necesario establecer los conceptos fundamentales de la representación digital de información.

Un **bit** (contracción de "binary digit") constituye la unidad básica de información en sistemas digitales. Un bit puede almacenar únicamente dos valores: 0 o 1. Esta representación binaria se debe a que los sistemas digitales utilizan circuitos electrónicos que pueden encontrarse en dos estados distintos y estables.

Un **byte** está compuesto por 8 bits consecutivos y representa la unidad básica de direccionamiento en la mayoría de sistemas computacionales modernos. Con 8 bits es posible representar $2^8 = 256$ valores diferentes, típicamente los enteros del 0 al 255 en representación sin signo, o del -128 al 127 en representación con signo en complemento a dos.

La organización de bits en patrones específicos permite representar diferentes tipos de datos. Para números enteros, la representación es directa utilizando el sistema binario. Sin embargo, para números reales, la situación se torna considerablemente más compleja debido a la necesidad de representar tanto la magnitud como la escala del número.

2.2. Motivación para la Representación en Coma Flotante

La representación en coma flotante surge de la necesidad de manejar números que abarcan rangos de magnitudes muy amplios, característica común en problemas científicos y de ingeniería. Consideremos los siguientes valores típicos que pueden aparecer en aplicaciones forestales:

- Constante de Avogadro: $6,022 \times 10^{23}$ moléculas/mol
- Radio de un átomo de carbono: $7,0 \times 10^{-11}$ metros
- Densidad de la madera: ~ 500 kg/m³
- Concentración de CO₂ atmosférico: $4,1 \times 10^{-4}$ (fracción molar)

Una representación de punto fijo, donde se asigna un número fijo de dígitos para la parte entera y otro para la parte decimal, resultaría inadecuada para manejar simultáneamente esta diversidad de escalas. La representación en coma flotante resuelve este problema separando la información sobre la magnitud (exponente) de la información sobre la precisión (mantisa).

Esta separación permite que la precisión relativa se mantenga aproximadamente constante a través de diferentes órdenes de magnitud. En otras palabras, tanto el número $1,234 \times 10^{23}$ como el número $1,234 \times 10^{-11}$ pueden representarse con la misma precisión relativa, lo cual es fundamental para la estabilidad de los algoritmos numéricos.

3. Estructura de la Representación en Coma Flotante

3.1. Formulación General

Un número en representación de coma flotante se expresa mediante la siguiente forma general:

$$x = (-1)^s \times m \times b^e \quad (1)$$

donde cada componente tiene un significado específico:

- s es el **bit de signo**: determina si el número es positivo ($s = 0$) o negativo ($s = 1$)
- m es la **mantisa** o significando: representa los dígitos significativos del número
- b es la **base** del sistema: para sistemas binarios, $b = 2$
- e es el **exponente**: determina la escala o magnitud del número

Esta representación es análoga a la notación científica decimal que se utiliza comúnmente en ciencias, donde un número como $3,14159 \times 10^2$ se descompone en mantisa (3,14159), base (10) y exponente (2).

3.2. Formatos Estándar IEEE 754

El estándar IEEE 754, adoptado universalmente por la industria computacional, define varios formatos de representación. Los dos más comunes son la precisión simple (32 bits) y la precisión doble (64 bits).

3.2.1. Precisión Simple (32 bits)

En el formato de precisión simple, los 32 bits disponibles se distribuyen según el siguiente esquema:

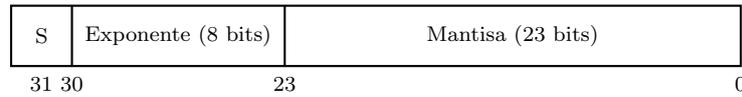


Figura 1: Distribución de bits en formato de precisión simple (32 bits)

La distribución específica es:

- **Bit de signo:** 1 bit (posición 31)
- **Campo del exponente:** 8 bits (posiciones 30 a 23)
- **Campo de la mantisa:** 23 bits (posiciones 22 a 0)

3.2.2. Precisión Doble (64 bits)

El formato de precisión doble utiliza 64 bits con la siguiente distribución:

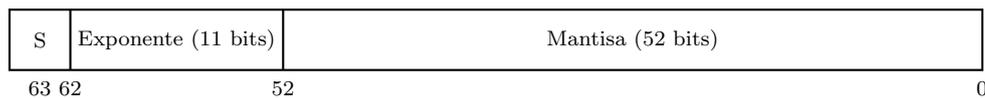


Figura 2: Distribución de bits en formato de precisión doble (64 bits)

La distribución correspondiente es:

- **Bit de signo:** 1 bit (posición 63)
- **Campo del exponente:** 11 bits (posiciones 62 a 52)
- **Campo de la mantisa:** 52 bits (posiciones 51 a 0)

3.3. Cálculo del Exponente y la Mantisa

Para comprender completamente la representación en coma flotante, es fundamental entender cómo se calculan y almacenan tanto el exponente como la mantisa. Sin embargo, antes de abordar estos conceptos, es necesario revisar brevemente la conversión de números binarios a decimales.

3.3.1. Repaso: Conversión de Binario a Decimal

La representación binaria utiliza potencias de 2 para expresar números. Para un número entero binario, cada posición representa una potencia de 2, comenzando desde 2^0 en la posición más a la derecha.

Ejemplo - Número entero: El número binario 1101_2 se convierte a decimal como:

$$1101_2 = 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 \quad (2)$$

$$= 1 \times 8 + 1 \times 4 + 0 \times 2 + 1 \times 1 \quad (3)$$

$$= 8 + 4 + 0 + 1 = 13_{10} \quad (4)$$

Ejemplo - Número con parte fraccionaria: Para números con parte fraccionaria, las posiciones a la derecha del punto binario representan potencias negativas de 2:

$$1101,101_2 = 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} \quad (5)$$

$$= 8 + 4 + 0 + 1 + 0,5 + 0 + 0,125 \quad (6)$$

$$= 13,625_{10} \quad (7)$$

Ejemplo - Número fraccionario puro: El número binario $0,1011_2$ se convierte como:

$$0,1011_2 = 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} + 1 \times 2^{-4} \quad (8)$$

$$= 0,5 + 0 + 0,125 + 0,0625 \quad (9)$$

$$= 0,6875_{10} \quad (10)$$

3.3.2. Cálculo del Exponente con Sesgo

El exponente no se almacena directamente, sino que se utiliza una representación con sesgo (bias). Este sesgo permite representar tanto exponentes positivos como negativos utilizando únicamente valores positivos en el campo del exponente.

Para precisión simple, el sesgo es 127, mientras que para precisión doble es 1023. El exponente real se calcula como:

$$e = e_{almacenado} - sesgo \quad (11)$$

donde $e_{almacenado}$ es el valor binario almacenado en el campo del exponente.

El rango de exponentes válidos es:

- **Precisión simple:** $-126 \leq e_{real} \leq 127$ (los valores extremos 0 y 255 se reservan para casos especiales)
- **Precisión doble:** $-1022 \leq e_{real} \leq 1023$ (los valores extremos 0 y 2047 se reservan para casos especiales)

3.3.3. Normalización y Cálculo de la Mantisa

La normalización es un proceso fundamental que maximiza la precisión de la representación. Un número normalizado tiene la forma:

$$1.f_1f_2f_3\dots f_n \times 2^e \quad (12)$$

donde f_1, f_2, \dots, f_n son los bits fraccionarios que se almacenan en el campo de la mantisa.

El aspecto crucial de la normalización es que el bit más significativo de la mantisa siempre es 1 para números normalizados. Dado que este bit es predecible, no se almacena explícitamente, lo que se conoce como el "bit implícito" o "bit escondido". Esto permite ganar un bit adicional de precisión.

La mantisa completa se calcula como:

$$m = 1 + \sum_{i=1}^n f_i \times 2^{-i} \quad (13)$$

donde f_i son los bits almacenados en el campo de la mantisa, y n es el número de bits de la mantisa (23 para precisión simple, 52 para precisión doble).

3.3.4. Ejemplo Detallado de Cálculo

Consideremos la representación del número decimal 12,375 en precisión simple:

Paso 1: Conversión a binario $12,375_{10} = 1100,011_2$

Paso 2: Normalización $1100,011_2 = 1,100011_2 \times 2^3$

Paso 3: Identificación de componentes

- Signo: positivo, por lo tanto $s = 0$
- Exponente real: $e = 3$
- Exponente almacenado: $e_{almacenado} = e + sesgo = 3 + 127 = 130 = 10000010_2$
- Mantisa fraccionaria: 100011 (se completa con ceros: 100011000000000000000000)

Paso 4: Representación final

Signo	Exponente	Mantisa
0	10000010	100011000000000000000000

4. Motivación y Ventajas de la Normalización

La normalización no es simplemente una convención arbitraria, sino que proporciona ventajas fundamentales para la representación numérica:

4.1. Maximización de la Precisión

Al garantizar que el bit más significativo sea siempre 1, la normalización elimina los ceros iniciales innecesarios en la mantisa. Esto permite que todos los bits disponibles se utilicen para almacenar información significativa, maximizando la precisión de la representación.

4.2. Unicidad de Representación

Sin normalización, un mismo número podría representarse de múltiples formas. Por ejemplo, el número 1.5 podría representarse como $1,5 \times 2^0$, $0,75 \times 2^1$, o $3,0 \times 2^{-1}$. La normalización garantiza que cada número tenga una única representación válida, lo que simplifica las operaciones aritméticas y las comparaciones.

4.3. Aprovechamiento del Bit Implícito

El bit implícito permite ganar efectivamente un bit adicional de precisión sin incrementar el espacio de almacenamiento. En precisión simple, aunque solo se almacenan 23 bits de mantisa, la precisión efectiva es de 24 bits gracias al bit implícito.

5. Características de la Representación Discreta

5.1. Distribución No Equiespaciada

Una característica fundamental de la representación en coma flotante es que los números representables no están uniformemente distribuidos en la recta real. Esta característica contrasta marcadamente con la representación de números enteros, donde los valores están equiespaciados.

La densidad de números representables en coma flotante sigue un patrón específico: entre cualquier par de potencias consecutivas de 2, existe exactamente el mismo número de valores representables. Específicamente, en el intervalo $[2^k, 2^{k+1})$, existen exactamente 2^{n+1} números representables, donde n es el número de bits de la mantisa.

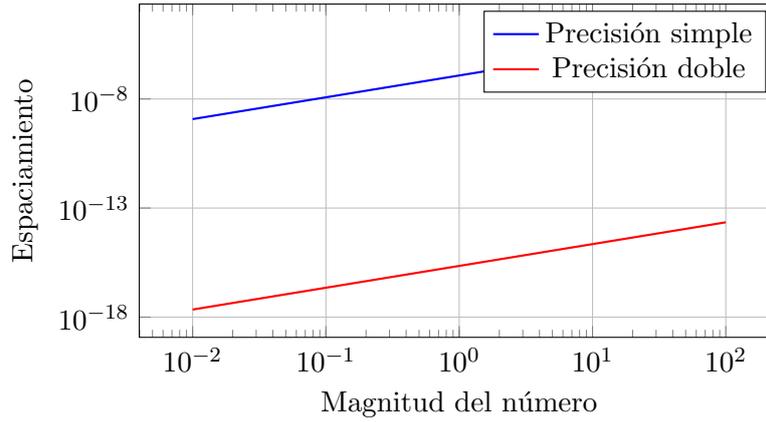


Figura 3: Espaciamiento entre números consecutivos en función de la magnitud

Esta distribución no uniforme tiene implicaciones importantes para el análisis numérico. El espaciamiento absoluto entre números consecutivos crece proporcionalmente con la magnitud del número, pero el espaciamiento relativo se mantiene aproximadamente constante.

5.2. Épsilon de Máquina

El epsilon de máquina (ϵ_m) es un parámetro fundamental que caracteriza la precisión del sistema de coma flotante. Se define como el número positivo más pequeño tal que:

$$fl(1 + \epsilon_m) > fl(1) \quad (14)$$

donde $fl(\cdot)$ denota la representación en coma flotante de un número. Matemáticamente, el epsilon de máquina se calcula como:

$$\epsilon_m = 2^{-(n+1)} \quad (15)$$

donde n es el número de bits de la mantisa.

Para los formatos estándar:

$$\epsilon_m = \begin{cases} 2^{-24} \approx 5,96 \times 10^{-8} & \text{(precisión simple)} \\ 2^{-53} \approx 1,11 \times 10^{-16} & \text{(precisión doble)} \end{cases} \quad (16)$$

El epsilon de máquina establece un límite fundamental para la precisión relativa que puede obtenerse en cualquier cálculo. Ningún algoritmo numérico puede producir resultados con una precisión relativa mejor que ϵ_m , independientemente de su sofisticación.

6. Análisis Detallado de Limitaciones y Errores

6.1. Error de Representación y Truncamiento

El error de representación surge cuando un número real no puede expresarse exactamente en el formato de coma flotante disponible. Este error es inherente y inevitable en la representación digital de números reales.

Para cualquier número real x , su representación en coma flotante $fl(x)$ satisface:

$$fl(x) = x(1 + \delta) \quad (17)$$

donde $|\delta| \leq \frac{\epsilon_m}{2}$ representa el error relativo de representación.
 El error absoluto de representación está acotado por:

$$|x - fl(x)| \leq \frac{\epsilon_m}{2}|x| \quad (18)$$

Esta cota demuestra que el error absoluto crece proporcionalmente con la magnitud del número, mientras que el error relativo se mantiene acotado por $\epsilon_m/2$.

6.1.1. Ejemplo de Error de Representación

Consideremos la representación del número decimal 0,1 en precisión simple. Este número, aparentemente simple, no puede representarse exactamente en binario:

$$0,1_{10} = 0,0001\overline{1001}_2$$

La secuencia 1001 se repite infinitamente, pero solo podemos almacenar un número finito de bits. Esto resulta en un error de representación que, aunque pequeño, puede acumularse en cálculos repetitivos.

6.2. Errores de Redondeo en Operaciones Aritméticas

Cada operación aritmética en coma flotante introduce potencialmente un error de redondeo. Para las operaciones básicas (suma, resta, multiplicación, división), el estándar IEEE 754 garantiza que el resultado es el número representable más cercano al resultado matemático exacto.

6.2.1. Suma y Resta

Para la suma de dos números a y b , el resultado calculado satisface:

$$fl(a + b) = (a + b)(1 + \delta) \quad (19)$$

donde $|\delta| \leq \epsilon_m$.

Sin embargo, cuando a y b tienen magnitudes muy diferentes, puede ocurrir que la suma se comporte como si el número menor no existiera. Esto sucede cuando:

$$|b| < \epsilon_m|a| \quad (20)$$

En este caso, $fl(a + b) = a$, y la información contenida en b se pierde completamente.

6.2.2. Multiplicación y División

Para la multiplicación, el error relativo también está acotado por ϵ_m :

$$fl(a \times b) = (a \times b)(1 + \delta) \quad (21)$$

La división presenta características similares, pero debe considerarse adicionalmente el riesgo de división por cero o por números muy pequeños que pueden llevar a overflow.

6.3. Cancelación Catastrófica

La cancelación catastrófica representa uno de los fenómenos más problemáticos en la aritmética de coma flotante. Ocurre cuando se restan dos números muy próximos, resultando en una pérdida significativa de dígitos significativos.

Consideremos dos números a y b tales que $a \approx b$ y $a > b$. Si ambos números contienen errores de representación:

$$fl(a) = a(1 + \delta_a) \quad (22)$$

$$fl(b) = b(1 + \delta_b) \quad (23)$$

Entonces la diferencia calculada es:

$$fl(a) - fl(b) = a(1 + \delta_a) - b(1 + \delta_b) \quad (24)$$

$$= (a - b) + a\delta_a - b\delta_b \quad (25)$$

El error relativo en la diferencia es:

$$\frac{|fl(a) - fl(b) - (a - b)|}{|a - b|} = \left| \frac{a\delta_a - b\delta_b}{a - b} \right| \quad (26)$$

Cuando $a \approx b$, este error relativo puede ser enormemente amplificado, incluso para valores pequeños de δ_a y δ_b . Para ilustrar este fenómeno, se presenta un análisis detallado en el Anexo A.

6.4. Acumulación de Errores

En secuencias largas de operaciones, los errores individuales se propagan y pueden acumularse. Para n operaciones sucesivas, cada una con error relativo acotado por ϵ_m , el error acumulado puede crecer hasta aproximadamente $n\epsilon_m$ en el peor caso.

Esta acumulación es particularmente problemática en:

- Sumas largas con muchos términos
- Algoritmos iterativos con muchas iteraciones
- Evaluación de polinomios de alto grado
- Integración numérica con pasos muy pequeños

6.5. Underflow y Overflow

El underflow y overflow representan situaciones extremas donde los límites del formato de representación son excedidos.

6.5.1. Underflow

El underflow ocurre cuando el resultado de una operación es tan pequeño que no puede representarse como un número normalizado. En este caso, el resultado puede:

- Convertirse en cero (flush-to-zero)
- Representarse como un número desnormalizado (gradual underflow)

Los números desnormalizados tienen exponente mínimo y no tienen el bit implícito, lo que resulta en menor precisión pero permite una transición gradual hacia cero.

6.5.2. Overflow

El overflow ocurre cuando el resultado excede el número más grande representable. En este caso, el resultado se convierte en infinito positivo o negativo, según corresponda.

Los rangos aproximados para evitar overflow son:

$$\text{Precisión simple: } |x| < 3,40 \times 10^{38} \quad (27)$$

$$\text{Precisión doble: } |x| < 1,79 \times 10^{308} \quad (28)$$

7. Números Especiales en IEEE 754

El estándar IEEE 754 define valores especiales para representar situaciones excepcionales que pueden surgir durante los cálculos.

7.1. Representación del Cero

Existen dos representaciones distintas de cero: $+0$ y -0 , que se diferencian únicamente en el bit de signo. Aunque matemáticamente equivalentes, pueden tener comportamientos ligeramente diferentes en ciertas operaciones, como la división por cero.

7.2. Infinito

Los valores $+\infty$ y $-\infty$ se representan con exponente máximo (todos los bits en 1) y mantisa completamente nula. Estos valores resultan de operaciones de overflow o divisiones por cero con operandos finitos.

7.3. NaN (Not a Number)

Los valores NaN se utilizan para representar resultados de operaciones matemáticamente indefinidas, tales como:

- $0/0$
- $\infty - \infty$
- $\sqrt{-1}$ (en aritmética real)
- Cualquier operación que involucre un NaN como operando

Los NaN se propagan a través de los cálculos: cualquier operación aritmética que involucre un NaN produce como resultado un NaN.

8. Implicaciones para Métodos Numéricos

8.1. Estabilidad Numérica

Un algoritmo es numéricamente estable si pequeñas perturbaciones en los datos de entrada (del orden de ϵ_m) producen pequeñas perturbaciones en el resultado final. La estabilidad numérica es una propiedad crucial que debe evaluarse cuidadosamente en el diseño de algoritmos.

La estabilidad puede analizarse mediante el concepto de amplificación de errores. Si un algoritmo amplifica los errores de entrada por un factor mayor que $1/\epsilon_m$, el resultado puede carecer completamente de precisión.

8.2. Estrategias de Implementación Robusta

Para desarrollar implementaciones numéricamente robustas, deben considerarse las siguientes estrategias:

1. **Reformulación algebraica:** Evitar cancelación catastrófica mediante manipulación algebraica de las expresiones.
2. **Escalado:** Normalizar los datos para evitar overflow y underflow.
3. **Pivoteo:** En algoritmos matriciales, utilizar estrategias de pivoteo para mejorar la estabilidad.
4. **Verificación de casos especiales:** Implementar verificaciones explícitas para divisiones por cero, overflow, y underflow.
5. **Uso de precisión extendida:** En cálculos críticos, utilizar precisión doble o extendida para cálculos intermedios.

9. Conclusiones

La representación en coma flotante constituye un compromiso fundamental entre el rango de magnitudes representables y la precisión disponible dentro de las limitaciones de almacenamiento finito. La comprensión profunda de sus características, limitaciones y fuentes de error es esencial para el desarrollo exitoso de métodos numéricos en ingeniería.

Las limitaciones inherentes de esta representación no deben verse como obstáculos insalvables, sino como restricciones que deben incorporarse conscientemente en el diseño algorítmico. Con una comprensión adecuada de estos principios, es posible desarrollar implementaciones numéricamente estables y confiables que produzcan resultados precisos dentro de las limitaciones fundamentales de la aritmética de coma flotante.

La consideración cuidadosa de estos aspectos es particularmente crucial en aplicaciones de ingeniería forestal, donde los modelos numéricos deben manejar datos que abarcan múltiples escalas temporales y espaciales, desde procesos celulares hasta dinámicas de ecosistemas completos.