

Introducción teórica

Teorema:

Para cada $x \in \mathbb{R}$, (de signo s) existe un par $(m \in \mathbb{R}, \text{exp} \in \mathbb{Z})$ que verifica: $x = (-1)^s \cdot 1.m \cdot 2^{\text{exp}}$.

Demostración:

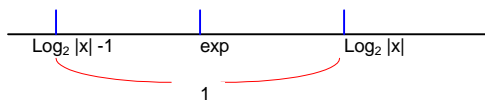
Vamos a calcular m y exp ; y así al mismo tiempo:

- Demostraremos que existen, es decir, que el teorema es cierto.
- Proporcionaremos un método para calcularlos.

a) Cálculo de exp .

$$|x| = 1.m \cdot 2^{\text{exp}} \Rightarrow \begin{cases} |x| \geq 1.0 \cdot 2^{\text{exp}} = 2^{\text{exp}} \Rightarrow \text{exp} \leq \log_2 |x| \\ |x| < 1.9 \cdot 2^{\text{exp}} < 2 \cdot 2^{\text{exp}} = 2^{\text{exp}+1} \Rightarrow \text{exp}+1 > \log_2 |x| \Rightarrow \text{exp} > \log_2 |x| - 1 \end{cases} \Rightarrow (\log_2 |x|) - 1 < \text{exp} \leq \log_2 |x|$$

Hay que tener en cuenta que en cualquier intervalo $(r-1, r]$, con $r \in \mathbb{R}$, hay sólo un entero de \mathbb{Z} .



Así pues, exp es el mayor entero de \mathbb{Z} que está entre los reales $\log_2 |x| - 1$ y $\log_2 |x|$

Conclusión: la forma de encontrar exp es:

1º: Calcular $\log_2 |x| = \frac{\ln |x|}{\ln 2}$. Lo más probable es que sea un real no entero.

2º: El exp buscado es el mayor de los enteros menores o iguales que $\log_2 |x|$, es decir: $\text{exp} = \text{floor}(\log_2 |x|)$

b) Cálculo de m :

$$|x| = 1.m \cdot 2^{\text{exp}} \Rightarrow 1.m = \frac{|x|}{2^{\text{exp}}} \Rightarrow m = \text{FRAC}\left(\frac{|x|}{2^{\text{exp}}}\right)$$

Corolarios:

1º.-

$$|x| = 1.m \cdot 2^{\text{exp}} \Rightarrow \text{En formato IEEE 754 :}$$

$$1.a) \text{ Campo exponente} = \text{exp} + \text{Exceso} = \text{exp} + (2^{n-1} - 1) = \text{floor}(\log_2 |x|) + (2^{n-1} - 1)$$

$$1.b) \text{ Campo mantisa} = m = \text{FRAC}\left(\frac{|x|}{2^{\text{exp}}}\right)$$

2º.-

$$|x| = 1.m \cdot 2^{\text{exp}} \Rightarrow |x| = (1 + 0.m) \cdot 2^{\text{exp}} = 2^{\text{exp}} + 0.m \cdot 2^{\text{exp}} \Rightarrow \begin{cases} |x| = 2^{\text{exp}} + 0.m \cdot 2^{\text{exp}} \geq 2^{\text{exp}} \\ |x| = 2^{\text{exp}} + 0.m \cdot 2^{\text{exp}} < 2^{\text{exp}} + 1.0 \cdot 2^{\text{exp}} = 2 \cdot 2^{\text{exp}} = 2^{\text{exp}+1} \end{cases} \Rightarrow$$

$$2^{\text{exp}} \leq |x| < 2^{\text{exp}+1}$$

⇓

La mayor potencia de 2 que se puede restar de $|x|$ es 2^{exp} .

⇓

El primer bit significativo de $|x| = b_{n-1} b_{n-2} \dots b_1 b_0 . b_{-1} b_{-2} b_{-3} \dots$ es b_{exp} .

Algoritmo para un número fraccionario de decimal a binario:

Recordemos que en la igualdad:

$$X = 0.D_{-1} D_{-2} D_{-3} \dots = 0.b_{-1} b_{-2} b_{-3} \dots$$

tenemos dados los dígitos D_i y se trata de encontrar los dígitos b_i que cumplen dicha igualdad. Cuando el número decimal tiene muchos dígitos fraccionarios y carece de parte entera, es preferible, por su rapidez, este algoritmo:

```
Acumulador := X ;  
i := -1 ;  
REPEAT  
  IF Acumulador * 2 > 1  
  THEN  
    bi := 1 ;  
    i := i - 1 ;  
    Acumulador := Acumulador - 1 ;  
  ELSE  
    bi := 0 ;  
  END  
UNTIL Acumulador = 0
```

Su justificación puede encontrarse, por ejemplo, en el capítulo 2 del libro de ETC I (Gestión).



Ejercicios.



Convierta el número decimal 1227 a complemento a dos de 16 bits.

Solución:

Al tratarse de un número positivo, nos basta con calcular su módulo, pues coinciden número y módulo.

1227	÷ 2 = 613.5	⇒ b ₀ = 1
613	÷ 2 = 306.5	⇒ b ₀ = 1
306	÷ 2 = 153.0	⇒ b ₀ = 0
153	÷ 2 = 76.5	⇒ b ₀ = 1
76	÷ 2 = 38.0	⇒ b ₀ = 0
38	÷ 2 = 19.0	⇒ b ₀ = 0
19	÷ 2 = 9.5	⇒ b ₀ = 1
9	÷ 2 = 4.5	⇒ b ₀ = 1
4	÷ 2 = 2.0	⇒ b ₀ = 0
2	÷ 2 = 1.0	⇒ b ₀ = 0
1	÷ 2 = 0.5	⇒ b ₀ = 1
0	÷ 2 = 0.0	⇒ b ₀ = 0

&1227 = % 0000 0100 1100 1011 = \$04CB



Convierta el número decimal -1227 a complemento a dos de 16 bits.

Solución:

Al tratarse de un número negativo, primero convertimos su módulo de decimal a binario y después los complementamos a dos.

- 1º Cálculo del módulo. Es un problema ya resuelto anteriormente: &1227 = % 0000 0100 1100 1011
- 2º Complemento a dos:
 - 2º.1 Inversión de todos los bits: %1111 1011 0011 0100
 - 2º.2 Suma de 1 en el bit menos significativo %1111 1011 0011 0101 = \$FB35



Convierta el número decimal -1227 a complemento a dos de 20 bits.

Solución:

Si aprovechamos el resultado de un problema resuelto anteriormente, -1227 en complemento a dos de 16 bits es

%1111 1011 0011 0101

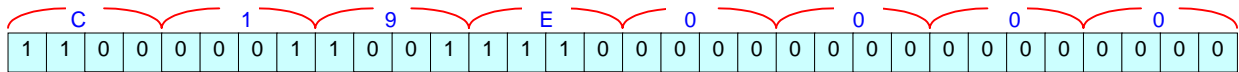
Para almacenarlo en un registro de 20 bits debemos hacer una extensión de signo:

%1111 1111 1011 0011 0101 = \$FBB35

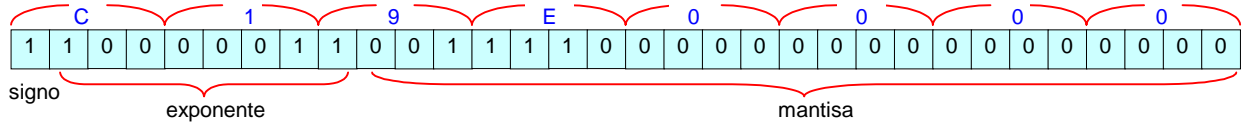
Encuentre el número decimal cuya representación en formato IEEE 754 en coma flotante de 32 bits, compactada en hexadecimal es C19E0000.

Solución:

Paso 1º: Obtención de la información almacenada en el registro a partir de la expresión compactada en hexadecimal



Paso 2º: Identificación de los diferentes campos presentes en el registro



Paso 3º: Cálculo de los diferentes componentes de la representación en punto flotante

Bit de signo = 1 \Rightarrow El número es negativo

Mantisa = %1.001111

Exponente auténtico = Exponente representado - Exceso

$$\begin{array}{r}
 \%10000011 \\
 - \%11111111 \\
 \hline
 \%00000100 = 4
 \end{array}$$

Paso 4º: Cambio de representación del número desde punto flotante a punto fijo

$$\text{Número} = -(\%1.001111) \cdot (2^4) = -1.234375 \cdot 16 = -19.75$$

Obtenga la representación del número 1540 en formato normalizado IEEE 754 para coma flotante de 32 bits.

Solución:

Paso 1º: Cálculo del campo exponente

Según el apartado a) del primer corolario:

$$\text{exp} = \text{floor}(\text{Log}_2 |x|) = \text{floor}\left(\frac{\text{Ln} |x|}{\text{Ln} 2}\right) = \text{floor}\left(\frac{\text{Ln} |1540|}{\text{Ln} 2}\right) = \text{floor}(10.58871464) = 10$$

$$\text{Campo exponente} = \text{exp} + \text{Exceso} = 10 + (2^7 - 1) = 10 + 127 = 137 = \%10001001$$

Paso 2º: Cálculo del campo mantisa

Según el apartado b) del primer corolario:

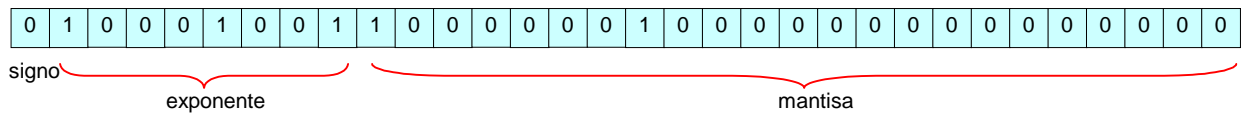
$$m = \text{FRAC}\left(\frac{|x|}{2^{\text{exp}}}\right) = \text{FRAC}\left(\frac{|1540|}{2^{10}}\right) = \text{FRAC}(1.050390625) = 0.050390625$$

$$\text{Mantisa } 1.050390625 = \% 1.10000001$$

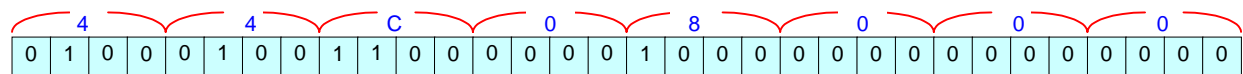
Paso 3º: Cálculo del campo de signo

$$1540 \geq 0 \quad \Rightarrow \quad \text{Bit de signo} = 0$$

Paso 4º: Almacenamiento en un registro de 32 bits de todos los campos anteriormente calculados



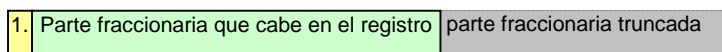
Paso 5º: Expresión compactada en hexadecimal de la información almacenada en el registro



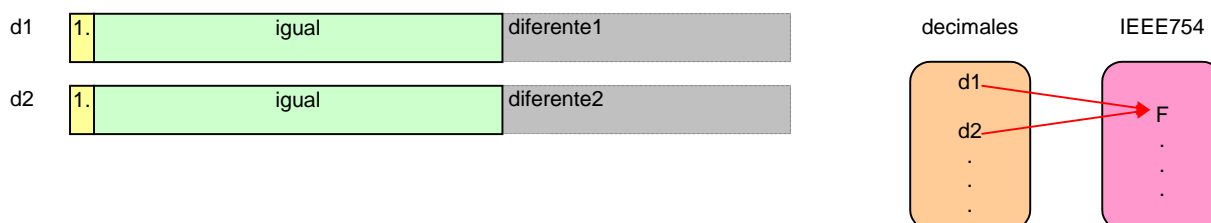


Solución:

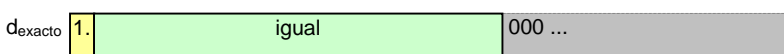
Debido a la limitación de la longitud del registro donde se almacena el número IEEE754 (en este caso es 16 bits, pero siempre existirá algún límite), puede suceder que al convertir de decimal a IEEE754 haya que truncar la mantisa del número decimal.



Por este motivo, existe la posibilidad de que diferentes números decimales tengan la misma representación IEEE754:



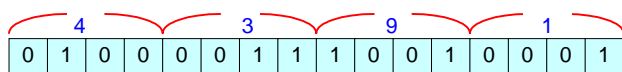
En este caso, el número en IEEE754 se corresponde exactamente sólo con aquel decimal cuya parte truncada es nula:



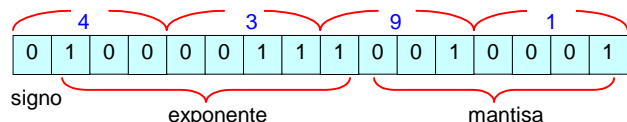
Cuando se quiere conocer el error (en valor absoluto) cometido al representar un decimal d_i con el número F en IEEE754, por definición se calcula el valor absoluto de su resta: $Error\ absolute = |d_i - F|$ (en esencia, toda comparación consiste en una substracción). Por supuesto, las comparaciones sólo se pueden hacer entre entidades homogéneas, por ello ambos números deben estar en el mismo formato. Como las operaciones aritméticas se realizan de forma más directa en decimal que en IEEE754, convertimos el IEEE754 a decimal; y luego calculamos el valor absoluto de la resta. Observe el lector que al hacer esto, estamos comparando el número d_i con el número d_{exacto}

Conversión del IEEE754 a decimal

Paso 1º: Obtención de la información almacenada en el registro a partir de la expresión compactada en hexadecimal



Paso 2º: Identificación de los diferentes campos presentes en el registro



Paso 3º: Cálculo de los diferentes componentes de la representación en punto flotante

Bit de signo = 0 \Rightarrow El número es positivo

Mantisa = %1.0010001

Exponente auténtico = Exponente representado - Exceso

$$\begin{array}{r} \%10000111 \\ - \%11111111 \\ \hline \%00001000 = 8 \end{array}$$

Paso 4º: Cambio de representación del número desde punto flotante a punto fijo

Número = (%1.0010001) * (2⁸) = %100100010 = &290

Comparación

$$\text{Error absoluto} = |d_i - d_{\text{exacto}}| = |291.072 - 290| = 1.072$$