



Sinal	Expoente	Mantissa	Valor do número
0	0000...0000	0000...0000	+0
0	0000...0000	0000...0001 1111...1111	$0, M \times 2^{-126}$
0	0000...0001 1111...1110	XXXX...XXXX	$1, XXXX...XXXX \times 2^{(e-b)}$
0	1111...1111	0000...0000	+ infinito
0	1111...1111	0000...0001 1111...1111	NaN
1	0000...0000	0000...0000	-0
1	0000...0000	0000...0001 1111...1111	$-0, M \times 2^{-126}$
1	0000...0001 1111...1110	XXXX...XXXX	$-1, XXXX...XXXX \times 2^{(e-b)}$
1	1111...1111	0000...0000	- infinito
1	1111...1111	0000...0001 1111...1111	NaN

### Soma em ponto flutuante

Somar  $9,999 \times 10^1$  com  $1,610 \times 10^{-1}$ .

Suponha que pode-se utilizar 4 dígitos para o número e 2 dígitos para expoente.

1. Alinha o ponto decimal do número que tem o menor expoente para igualar ao expoente do número que tem maior expoente.

$$1,610 \times 10^{-1} = 0,1610 \times 10^0 = 0,016 \times 10^1 \text{ (só pode usar 4 dígitos)}$$

Porque alinha menor com maior ? Pode perder dígitos, então melhor perder dígitos menos significativos. Maior chance de resultado ser normalizado.

2.

3.  $1.101 \times 2^{-6},$

4.  $1.001001 \times 2^6$

5.  $.000000000001101 \times 2^6$

6.  $1.001001 \times 2^6$

7. -----

8.  $1.001001000001101 \times 2^6$

### 9. Soma as mantissas

10.

11.  $9,999$

12.  $0,016$

13. -----

14.  $10,015 \times 10^1$

15. Coloca em notação científica normalizada

$$10,015 \times 10^1 = 1,0015 \times 10^2$$

16. Ajusta para número de dígitos possível com arredondamento

$$1,0015 \times 10^2 = 1,002 \times 10^2$$

17. Verifica se houve overflow ou underflow (expoente maior ou menor do que é possível ser representado)

Exemplo de adição de números em ponto flutuante em binário:

Somar  $0,5_{10}$  com  $-0,4375_{10}$

1. Converte para a base 2.

$$0,5 = 0,1 \times 2^0 = 1,000 \times 2^{-1}$$

$$-0,4375 = -0,01110 \times 2^0 = -1,110 \times 2^{-2}$$

2. Desloca para a direita o número com menor expoente. O número de deslocamentos é igual a diferença entre os expoentes. Se os expoentes são iguais, não há deslocamento.

$$-1,110 \times 2^{-2} = -0,1110 \times 2^{-1} = -0,111 \times 2^{-1}$$

3. Soma os números

$$1,000 \times 2^{-1} + (-0,111 \times 2^{-1}) = 0,001 \times 2^{-1}$$

4. Normaliza a soma se necessário

(ex. resultado igual a 10. \_\_ ou 11. \_\_ deverá sofrer um deslocamento para a direita e um incremento no expoente. Podemos também ter um resultado de 0. \_\_\_\_ se somarmos um número positivo com negativo. Deverá sofrer um deslocamento para a esquerda e um decremento do expoente.

$$0,001 \times 2^{-1} = 0,010 \times 2^{-2} = 0,100 \times 2^{-3} = 1,000 \times 2^{-4}$$

5. Para padrão IEEE 754 precisão simples, ocorre overflow se expoente maior que 127 ou underflow se expoente menor que -126 então será gerada uma exceção

6. Trunca resultado para número de bits que pode ser utilizado.

$$1,000 \times 2^{-4}.$$