

$$(19.625)_{10} = (10011.101)_2 = (1.0011101)_2 \times 2^4$$

## Machine floating point number.

Not all real number can be exactly represented as a machine floating-point number

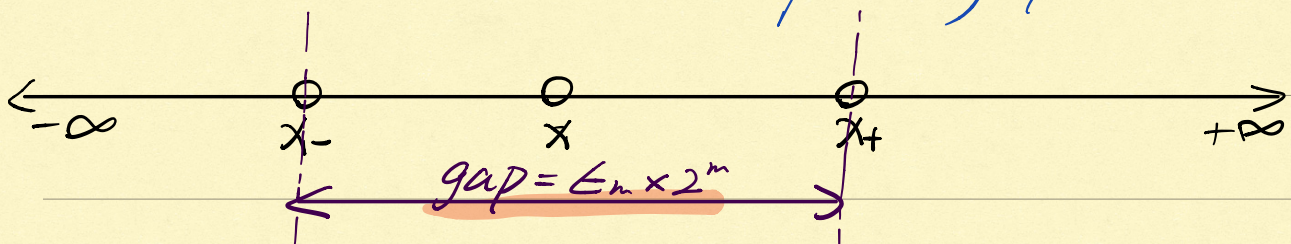
$$x = \pm 1. b_1 b_2 \dots b_n \times 2^m$$

*Data bits*

*can't represented.*

$x$  will be approximated by either  $x_-$  or  $x_+$ ,

the nearest two machine floating point numbers



$$x_- = 1. b_1 b_2 b_3 \dots b_n \times 2^m$$

$$x_+ = x_- + \underbrace{0.00 \dots 01}_{2^{-n}} \times 2^m$$

$$\underline{2^{-n}} \times 2^m = \underline{\epsilon_m} \times 2^m$$

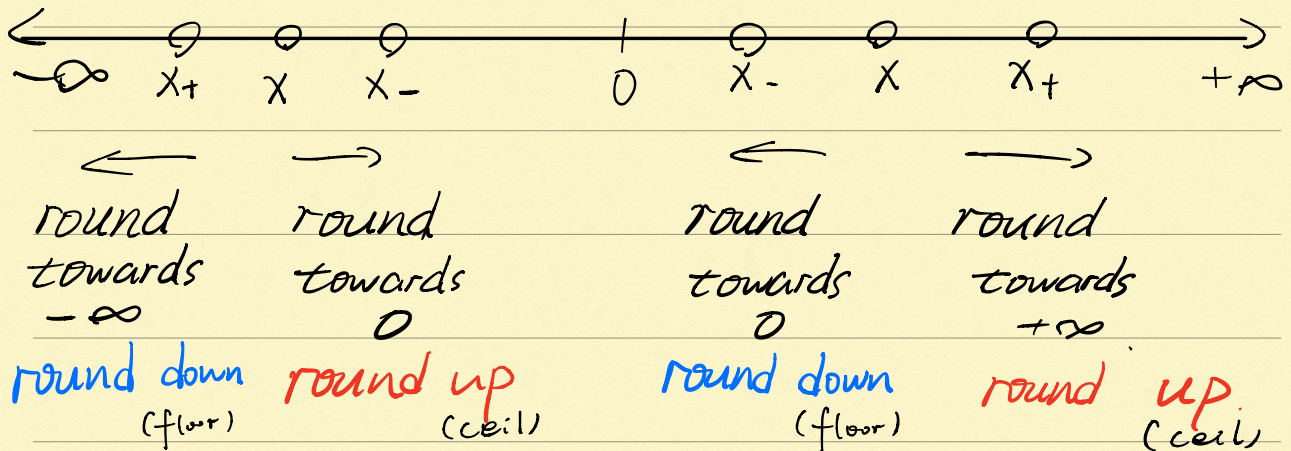
$$|(x_+) - (x_-)| = \epsilon_m \times 2^m$$

## Rounding

The process of replacing  $x$  by a nearby



machine number is called rounding, and the error involved is called roundoff error.



①. Round by chopping  $\uparrow f(x) = x_+ / x_-$

②. Round to nearest

Rounding errors  
(round off)

Absolute error:  $|f(x) - x| \leq \epsilon_m \times 2^m$

Relative error:  $\frac{|f(x) - x|}{|x|} \leq \frac{\epsilon_m \times 2^m}{1. b_1 b_2 \dots b_n \times 2^n}$

$$\Rightarrow \frac{|f(x) - x|}{|x|} \leq \epsilon_m \quad = \frac{\epsilon_m}{1. b_1 b_2 \dots b_n} \leq \epsilon_m$$



IEEE Single precision  $\frac{|f(x) - x|}{|x|} \leq 2^{-23} \approx 1.2 \times 10^{-7}$

Double precision  $\frac{|f(x) - x|}{|x|} \leq 2^{-52} \approx 2.2 \times 10^{-16}$

Gap:  $\delta < \underline{\text{gap}}$ ,  $f(x + \delta) = f(x)$

$\epsilon_m \times 2^m$

## Arithmetic

Possible:  $(x + y) + z \neq x + (y + z)$

$$z(x + y) \neq zx + zy$$

$$x + y = x$$

really small

Arithmetic:

①. compute the exact result.

②. Round.

$$x + y = f(x + y)$$

$$x \times y = f(x \times y)$$



## Addition and subtraction.

- ①. bring both numbers onto a common exponent.
- ②. Do "grade-school" operation.
- ③. Round result.

Example:  $x = \pm 1.b_1b_2b_3b_4 \times 2^m$

$$a = (1.1011)_2 \times 2^1$$

$$b = (1.1010)_2 \times 2^1$$

$$c = a - b = (0.0001)_2 \times 2^1$$

$$c = (1.?????)_2 \times 2^{-5}$$

lose of precision

machine fills them with its best guess, (not good).

Catastrophic  
Cancellation

Example: ① 16 decimal digits of accuracy:

$$f(a) = 3004.45$$

$$f(b) = 3004.46$$

$$b-a = \underbrace{0000.01}_{\text{lose 5 digits}} \underbrace{\text{-----}}_{\text{Still has 11 digits left.}}$$

lose 5 digits

Still has 11 digits left.

(significant)

②. 5 decimal accurate digits.

$$f(x) = \sqrt{x^2 + 1} - 1$$

$$f(10^{-3}) = \sqrt{10^{-6} + 1} - 1$$

$$f(10^{-6} + 1) = 1 \Rightarrow f(10^{-3}) = 0.$$

③.  $f(x) = \sqrt{x^2 + 1} - 1$

rewrite  $f(x) = \frac{x^2}{\sqrt{x^2 + 1} + 1} \Rightarrow f(10^{-3}) = \frac{10^{-6}}{2}$