

$$2^0 \quad 2^0 + \varepsilon_m \quad 2^0 + 2\varepsilon_m$$

Single
Double

$$\varepsilon_m \approx 10^{-8}$$

$$\varepsilon_m \approx 10^{-16}$$

Def: - Let x denote a ^{float} representable number then $\text{fl}(x)$ is the associated float.

$$- \frac{|x - \text{fl}(x)|}{|x|} \leq \varepsilon_m$$

Floating pt arithmetic

Q: How many real # are there?

Uncountably infinite.

Q: Can the machine represent all of them?

No.

The largest number that you can represent on

The largest number that you can represent on the machine is 1.79×10^{308}

This number is called the overflow (OFL).

The smallest # that you can represent on the machine is 2.23×10^{-308} . This is called the underflow (UFL).

On the machine, numbers are represented by truncated binary series of the form

$$x = (a_0 2^0 + a_1 2^{-1} + a_2 2^{-2} + \dots + a_p 2^{-(p-1)}) 2^e$$

Where p is the precision, $\{a_i\}$ is the exponential. (common factor in base 2)

Def: The representation of a real # on the machine is called floating pt representation.

For $x \in \mathbb{R}$ that is represented w/a float, $fl(x)$ denotes its floating pt representation.

Types of precision

	P	Machine Epsilon ϵ_m
single	24	$\epsilon_m = 2^{-23} \approx 10^{-6}$
double	53	$\epsilon_m = 2^{-52} \approx 10^{-16}$

Warm-up

$$\text{Let } \epsilon_m = 10^{-2}$$

(i) set $x = \frac{1}{3}$. What is $\tilde{x} = \text{fl}(x)$?

$$\text{What is } \frac{|x - \tilde{x}|}{|x|}?$$

(ii) set $x = 1234$. What is $\tilde{x} = \text{fl}(x)$?

$$\text{What is } \frac{|x - \tilde{x}|}{|x|}?$$

Soln: (i) $\tilde{x} = \text{fl}(\frac{1}{3}) = 0.33$

$$\text{absolute error } |x - \tilde{x}| = 0.003333\bar{3}$$

$$\text{relative error } \frac{|x - \tilde{x}|}{|x|} = 10^{-2}$$

$$(ii) \tilde{x} = 1200$$

$$\text{absolute error } |x - \tilde{x}| = 1234 - 1200 = 34 > 10^{-2}$$

$$\text{relative error } \frac{|x - \tilde{x}|}{|x|} = \frac{34}{1234} \sim 10^{-2}$$

$$0.\underbrace{027}_{2.7 \times 10^{-2}}$$

Def: The most accurately we can approximate a real # that is representable w/a float is machine

that is representable w/a float is machine epsilon ε_m .

i.e. let $x \in \mathbb{R}$ be representable by a float, let $\tilde{x} = \text{fl}(x)$
then $\frac{|x - \tilde{x}|}{|x|} \leq \varepsilon_m$

The biggest float $\sim 10^{308}$. This # is called
the overflow #. (OFL). The smallest float is

$\sim 10^{-308}$. This is called the underflow #. (UFL)

The ordering of everything you can represent on the computer is

$$0 < \text{UFL} < \varepsilon_m < \text{OFL}$$

Exceptional values

computer name	How to make it
Inf	Infinity
NaN	Not a #

or go over OFL

$\%/\%$, $0 \cdot \text{Inf}$, Inf/Inf