

Lecture 13: Floating Point Arithmetic

$$[1, 2] = 1, 1 + 2^{-52}, 1 + 2 \times 2^{-52}, \dots, 2$$

$$\text{largest number} = 1.79 \times 10^{308}$$

$$\text{smallest number} = 2.23 \times 10^{-308}$$

$$[2, 4] = 2[1, 2] = 2, 2 + 2^{-51}, 2 + 2 \times 2^{-51}, \dots, 4$$

$$\epsilon_{\text{machine}} = 2^{-52} \approx 2.22 \times 10^{-16}$$

Floating point numbers F

* For all $x \in \mathbb{R}$, there exists $x' \in F$ such that
$$\frac{|x - x'|}{|x|} \leq \epsilon_{\text{machine}} \rightarrow \text{relative error is always } \epsilon_{\text{machine}}$$

* Let $S: \mathbb{R} \rightarrow F$ gives the closest floating point approximation to a real number (rounding).

* For all $x \in \mathbb{R}$, there exists ϵ with $|\epsilon| \leq \epsilon_{\text{machine}}$ such that
$$S(x) = x(1 + \epsilon).$$

Floating point theorem

$\oplus, \ominus, \otimes, \oslash \rightarrow$ floating point operation

For all $x, y \in F$, there exists ϵ with $|\epsilon| < \epsilon_{\text{machine}}$ such that

$$x \otimes y = (x * y)(1 + \epsilon)$$