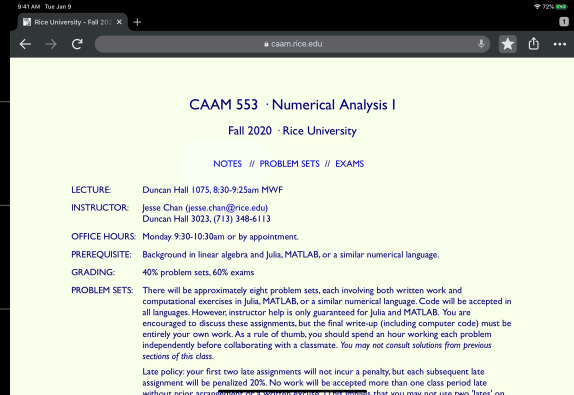


## Class overview:

[www.caam.rice.edu/~caam553/](http://www.caam.rice.edu/~caam553/)

- homework, exams, + notes will be posted here.
- Canvas for homework submission + gradebook.



- What should you learn?
  - Analysis of algorithms for problems in continuous (vs discrete) mathematics.
  - Tools for research: analysis + implementation of algorithms
  - **Qualifying exam preparation**
    - ↳ Make exams like quals.
    - ↳ How is **practice**.

## Floating point number systems

- What you should learn from this lecture:

- how real numbers are represented on computers
- when "catastrophic" rounding errors occur

- Naive ideas for representations

- $3.14159 \dots \times 10^e$

⇒ store digits and a scale factor

- What are some issues w/ this

⇒ **fixed** precision numbers

$$a \times b = \underline{56028}$$

$$\begin{array}{ccc} 123 & 456 & \\ \uparrow\uparrow\uparrow & \uparrow\uparrow\uparrow & \end{array}$$

$$5.60 \times 10^4$$

$$= 1.23 \times 10^2 \quad 4.56 \times 10^2 = 3 \text{ digits} \\ + \text{ exponent}$$

"Floating" vs "fixed" point

- aims for **relative** vs **absolute** precision

$$X = \pm \left( d_0 \beta^0 + d_1 \beta^{-1} + \dots + d_{p-1} \beta^{-(p-1)} \right) \beta^e$$

↑ sign      ↓ coefficients      base<sup>power</sup>      p = precision

- everything is an integer

- $\beta =$  base (radix)
- $p =$  precision

- $e =$  exponent (usually with  $e_{\min} \leq e \leq e_{\max}$ )

**Example:** suppose  $\beta=2, p=3$

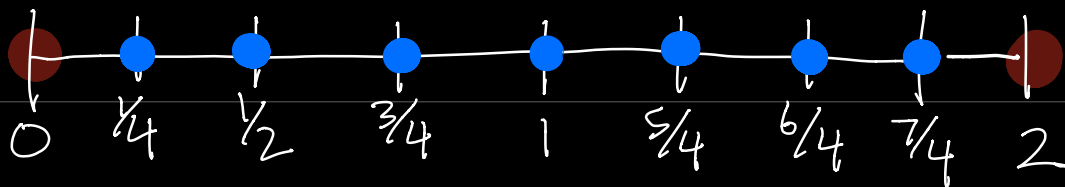
$$X = \pm (d_0 2^0 + d_1 2^{-1} + d_2 2^{-2}) 2^e$$
$$= \pm \left( d_0 + d_1 \frac{1}{2} + d_2 \frac{1}{4} \right) 2^e$$

Suppose  $e=0$ : **what numbers can we represent?**

$$X = \pm \left( d_0 + d_1 \frac{1}{2} + d_2 \frac{1}{4} \right) 2^e$$

$$\begin{array}{l}
 e \neq 0 \\
 \left\{ \begin{array}{l}
 x = \frac{1}{4} \Rightarrow d = (d_0, d_1, d_2) = (0, 0, 1) \\
 x = \frac{1}{2} \Rightarrow d = (0, 1, 0) \\
 x = \frac{3}{4} \Rightarrow d = (0, 1, 1) \\
 x = 1 \Rightarrow d = (1, 0, 0) \\
 x = \frac{5}{4} \Rightarrow d = (1, 0, 1) \\
 x = \frac{6}{4} \Rightarrow d = (1, 1, 0) \\
 x = \frac{7}{4} \Rightarrow d = (1, 1, 1)
 \end{array} \right.
 \end{array}$$

$x \in (0, 2) \Rightarrow 0, 2$  not representable.



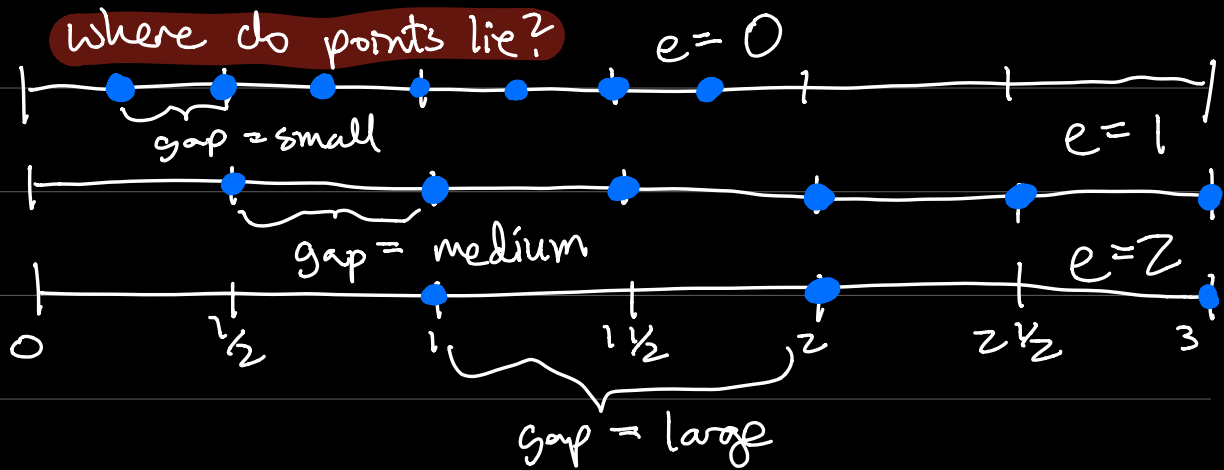
Suppose now  $e=1$ : what numbers are representable?

$$x = \pm \left( d_0 + d_1 \frac{1}{2} + d_2 \frac{1}{4} \right) 2$$

$$\Rightarrow x = \left\{ \frac{1}{4}, \frac{1}{2}, \frac{3}{4}, 1, \frac{5}{4}, \frac{6}{4}, \frac{7}{4} \right\} \times 2$$

•  $\Rightarrow x = \{\frac{1}{2}, 1, \frac{6}{4}, 2, 2\frac{1}{2}, 3, 3\frac{1}{2}\}$  for  $e=1$

•  $x = \{1, 2, 3, 4, 5, 6, 7\}$  for  $e=2$



### Observations:

- representations of numbers is not unique. Can fix by assuming  $d_0=1$  (referred to as normalization)

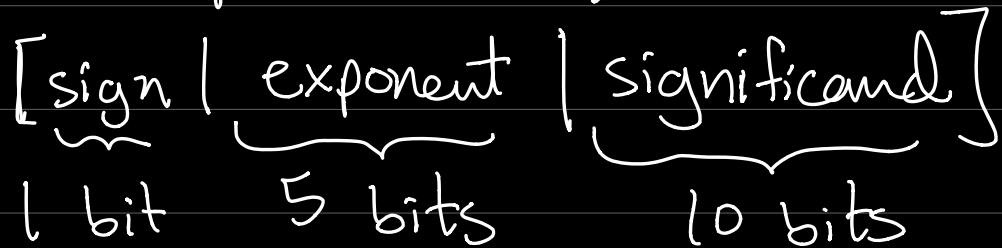
- **Main idea**: spacing between numbers wider for larger numbers  
**precision depends on  $e$** .

- radix point  $\Rightarrow$   $1\ 2\ 3\ .\ 4\ 5\ 6$   
moves or "floats" stored separately in fixed pt. format

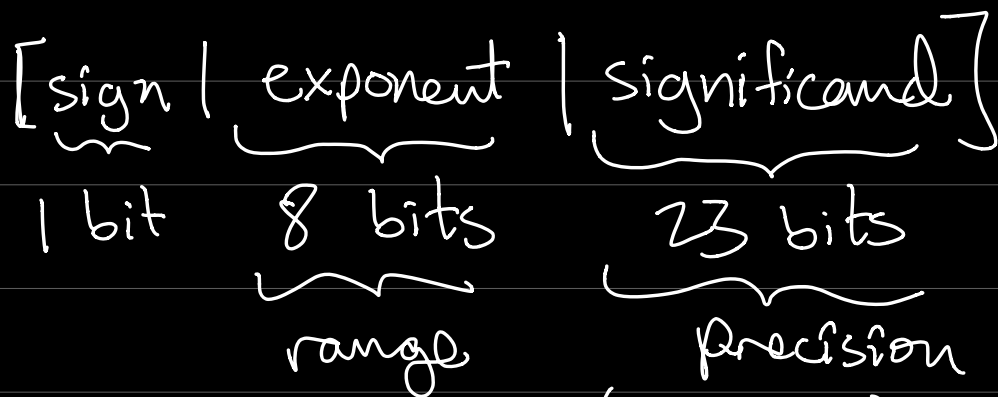
IEEE - Inst. of Electrical + Electronic Engineers

## Floating point standards

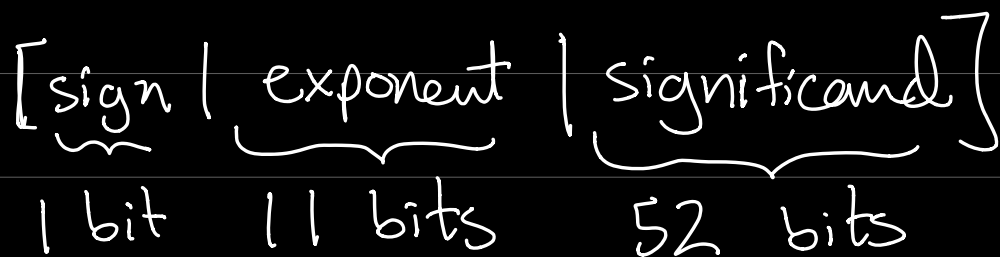
- half :  $p=11$  (16 bit)



- single :  $p=24$  (32 bit)



- double :  $p=53$  (64 bit)



- quad:  $p = 113$  (128 bit)

Note: cost of operations do not scale linearly wrt. number of bits!

### Machine Epsilon

- Binary  $\Rightarrow \beta = 2$ , double prec.

For  $e = 0$ ,  $x \in [0, 2]$  and

$$x = d_0 2^0 + d_1 2^{-1} + \dots + d_{52} 2^{-52}$$

$\epsilon_{\text{mach}} =$  smallest rel. diff b/w  
2 floating pt  
numbers

## Rounding in floating point systems

- Floating point systems designed to ensure that the **relative floating point representation error is bounded.**

- Let  $fl(x)$  denote the floating point representation of  $x$ .

$$\Rightarrow \frac{|x - fl(x)|}{|x|} \leq \epsilon_{mach}$$

Pf: Assume  $x \in \mathbb{R}$  & within representable range of floating point system

$\Rightarrow x$  is b/w 2 floating pt numbers

$$a \leq x \leq b, |b-a| \leq \epsilon_{mach}$$

$$fl(x) = x(1+\delta) \Rightarrow \frac{|x - fl(x)|}{|x|} \leq \epsilon_{mach} (|\delta| \leq \epsilon_{mach})$$



- Floating point arithmetic is not exact,  
 & requires truncating information.

Let  $\mathbb{F}$  = set of all floating point numbers

Let  $x, y \in \mathbb{F}$ .

$$\Rightarrow x = \pm (x_0 \beta^0 + x_1 \beta^1 + \dots + x_{p-1} \beta^{-(p-1)}) \beta^{e_x}$$

$$+ y = \pm (y_0 \beta^0 + y_1 \beta^1 + \dots + y_{p-1} \beta^{-(p-1)}) \beta^{e_y}$$

Assume  $e_y = e_x - 1$

$$\beta^{e_y} = \beta^{e_x-1}$$

$$y = \pm (y_0 \beta^{-1} + y_1 \beta^{-2} + \dots + y_{p-1} \beta^{-(p-2)}) \beta^{e_x}$$

$$x+y = \pm (x_0 \beta^0 + (x_1 + y_0) \beta^{-1} + \dots) \beta^{e_x}$$

Main issue:  $\mathbb{F}$  (set of floating point numbers) is not closed, e.g.,

$\Rightarrow$  for  $x, y \in \mathbb{F}$ ,  $x+y \notin \mathbb{F}$

Solution: design operations st.

$$fl(x+y) = (x+y)(1+\delta)$$

$$fl(x-y) = (x-y)(1+\delta)$$

$$fl(x*y) = (x*y)(1+\delta)$$

$$fl(x/y) = (x/y)(1+\delta)$$

$$\text{w/ } |\delta| \leq \epsilon_{mach}$$

- Can still have issues ...

## Rounding and catastrophic cancellation

let  $x, y \in \mathbb{R} \Rightarrow$  consider  $\text{fl}(x-y)$

$$\begin{aligned} \text{let } \hat{x} &= (1+\delta_x)x \\ \hat{y} &= (1+\delta_y)y \Rightarrow \text{fl}(x-y) = \text{fl}(\hat{x}-\hat{y}) \\ &= (x-y + \delta_x x - \delta_y y) \\ |\delta| &\leq \epsilon_{\text{mach}} \end{aligned}$$

$$\frac{|\text{fl}(\hat{x}-\hat{y}) - (x-y)|}{|x-y|} = \frac{|\delta_x x - \delta_y y|}{|x-y|}$$

$$\leq \frac{|\delta_x||x| + |\delta_y||y|}{|x-y|} \leq$$

$$\leq \underbrace{\max\{\delta_x, \delta_y\}}_{\leq \epsilon_{\text{mach}}} \underbrace{\frac{|x| + |y|}{|x-y|}}$$

Error can be large if  $x \approx y$  or if  $x, y \gg 1$ !!