

Math 55I

2/18/21

• HW #2 due 2/25 @ 9PM  
Th

• Read 3.1-3.3 for Tuesday

Find  $x$  such that  $\underbrace{f(x) = 0}$   
non-linear

• Look over 2.1 - 2.3

## Condition and Stability: Evaluate $f(x)$

Condition number of  $f$  at  $x$

$$K = \max_{|x-x^*|} \frac{\text{Rel}(f(x))}{\text{Rel}(x)} \approx \max_x \left| \frac{x \cdot f'(x)}{f(x)} \right|$$

$\Rightarrow K = O(1)$ , well-conditioned  
 $K \gg 1$ , III-conditioned

Ex:  $f(x) = \sqrt{x+1} - \sqrt{x}$ , uniformly well-conditioned

$$K \leq \frac{1}{2}$$



$$\sqrt{x+1} - \sqrt{x} = \frac{1}{\sqrt{x+1} + \sqrt{x}}$$

$$\left( \frac{\sqrt{x+1} + \sqrt{x}}{\sqrt{x+1} + \sqrt{x}} \right)$$

$$(1) \quad f(x) = \sqrt{x+1} - \sqrt{x} \quad \text{4-digit arithmetic}$$

$$f(1984) = \sqrt{1985} - \sqrt{1984} \quad \text{UNSTABLE}$$

Subtractive cancellation  $\approx [0.4455 \times 10^2] - [0.4454 \times 10^2]$   
 $= 0.0001 \times 10^2 = 0.1000 \times 10^{-1}$

*1-digit*

MATLAB:  $\sqrt{1985} - \sqrt{1984} = 0.112239\dots \times 10^{-1}$

$$(2) \quad f(x) = \frac{1}{\sqrt{x+1} + \sqrt{x}}$$

$$f(1984) = \frac{1}{\sqrt{1985} + \sqrt{1984}} \approx \frac{1}{0.4455 \times 10^2 + 0.4454 \times 10^2}$$

$$= \frac{1}{0.8909 \times 10^2} \approx \underbrace{0.1122 \times 10^{-1}}_{4\text{-digit}} \quad \text{STABLE}$$

Evaluate  $f(x)$



Well-conditioned

Ill-conditioned

Algorithm

**STABLE**

**UNSTABLE**

$$\text{Ex: } ax^2 + bx + c = 0 \Rightarrow x_{+-} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$a = 0.2$$

$$b = -47.91$$

$$c = 6$$

MATLAB  
⇒

$$x_+ = \underline{239.4247\dots}$$

$$x_- = \underline{0.1253\dots}$$

4-digit arithmetic

$$x_{+-} \approx \frac{47.91 \pm \sqrt{(-47.91)^2 - 4(0.2) \cdot 6}}{2(0.2)}$$

$$\approx \frac{47.91 \pm \sqrt{2290}}{0.4}$$

$$\approx \frac{47.91 \pm 47.85}{0.4}$$

Subtractive  
cancellation

$$= \left\{ \begin{array}{l} \overbrace{239.4}^{4\text{-digits}} \\ \overbrace{0.1500}^{1\text{-digit}} \end{array} \right\}$$

(+)

(-)

Fix?

$$\frac{-b - \sqrt{b^2 - 4ac}}{2a} \cdot \frac{-b + \sqrt{b^2 - 4ac}}{-b + \sqrt{b^2 - 4ac}} = \boxed{\frac{2c}{-b + \sqrt{b^2 - 4ac}}}$$

$$x_- \approx \frac{2 \cdot 6}{47.91 - \sqrt{(-47.91)^2 - 4 \cdot (0.2) \cdot 6}}$$

$$\stackrel{=}{\underbrace{4\text{-digit}}_{\text{full 4-digit accuracy}}} 0.\underbrace{1253}_{\text{accuracy!}}$$

# Floating Point Number Systems : fl(x)

$x \in \mathbb{R}$ , base  $\beta$  (typically 2), n

$$fl(x) = \pm (0.b_1 b_2 b_3 \dots b_n)_\beta \cdot \beta^e$$

(\*)  $0 \leq b_i \leq \beta - 1$        $\begin{cases} \beta = 2 & \{0, 1\} \\ \beta = 10 & \{0, 1, 2, \dots, 9\} \end{cases}$

(\*)  $b_1 \neq 0$

$\pm$  : sign

$(0.b_1 b_2 \dots b_n)_\beta$  : mantissa

e : exponent

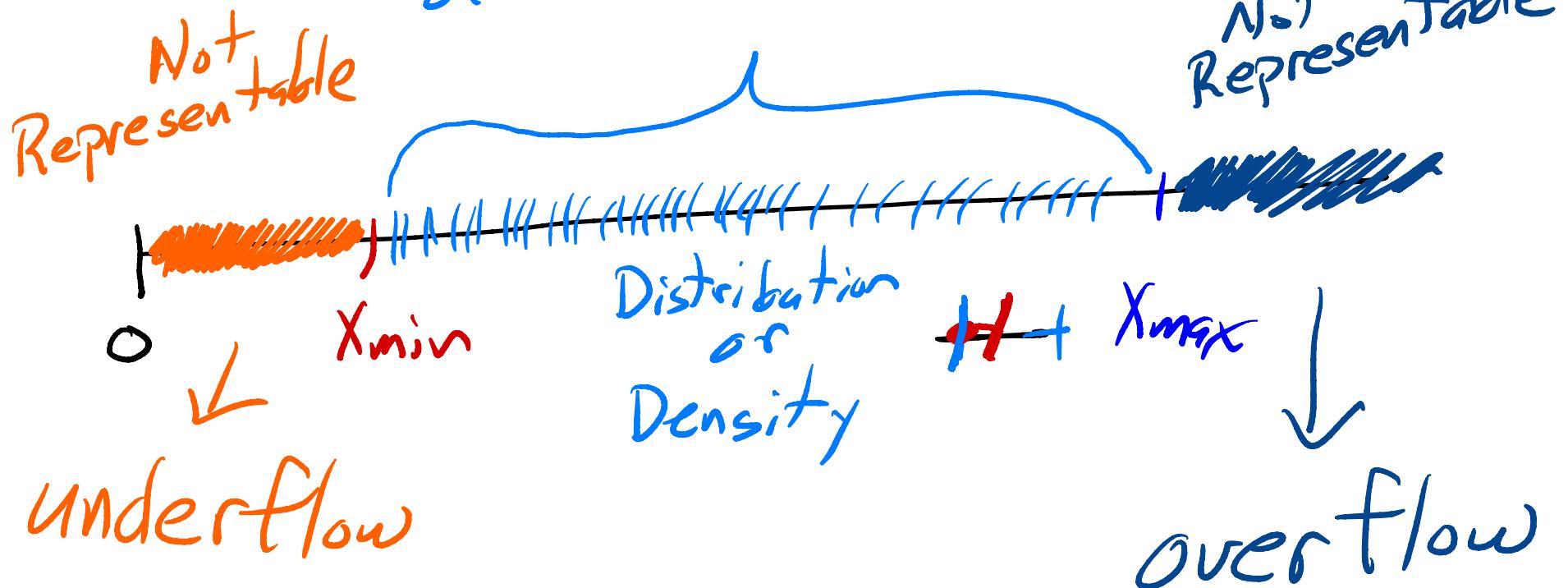
$$(-m \leq e \leq m)$$

$$\text{Ex: } \underline{n = 4} \quad \beta = 2 \quad \underline{m = 3}$$

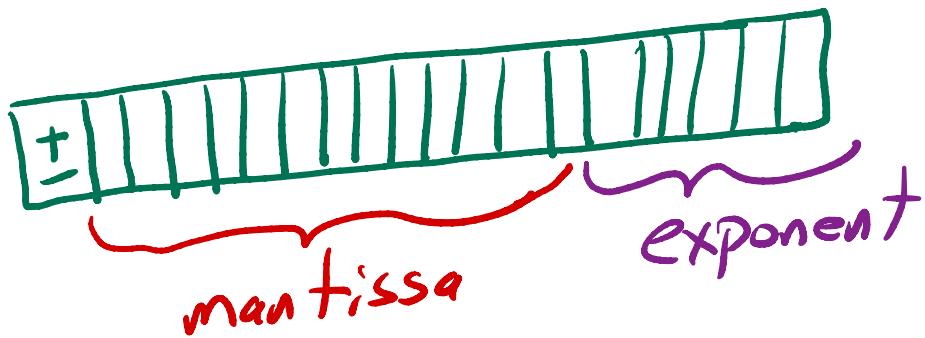
$$\underline{x_{\min}} = (0.1000)_2 \times 2^{-3} = \underline{2^{-4}} = (0.0625)_{10}$$

$$\underline{x_{\max}} = (0.1111)_2 \times 2^3 = (111.1)_2 = (7.5)_{10}$$

Only a FINITE number  
of numbers are representable



# IEEE 754 Standard $F1(x)$



single-precision (32-bits)	(1, 23, 8)
double-precision (64-bits)	(1, 52, 11)
quadruple-precision (128-bits)	(1, 112, 4)