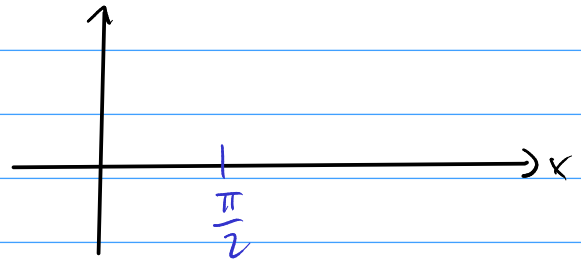


$$K \approx \left| \frac{x \cdot f'(x)}{f(x)} \right|$$

$$f(x) = \sin x \rightarrow \frac{x \cdot \cos(x)}{\sin(x)}$$

$$\text{biggest near } 0 \rightarrow \frac{x \cdot 1}{x} \approx 1$$



Stability and Accuracy

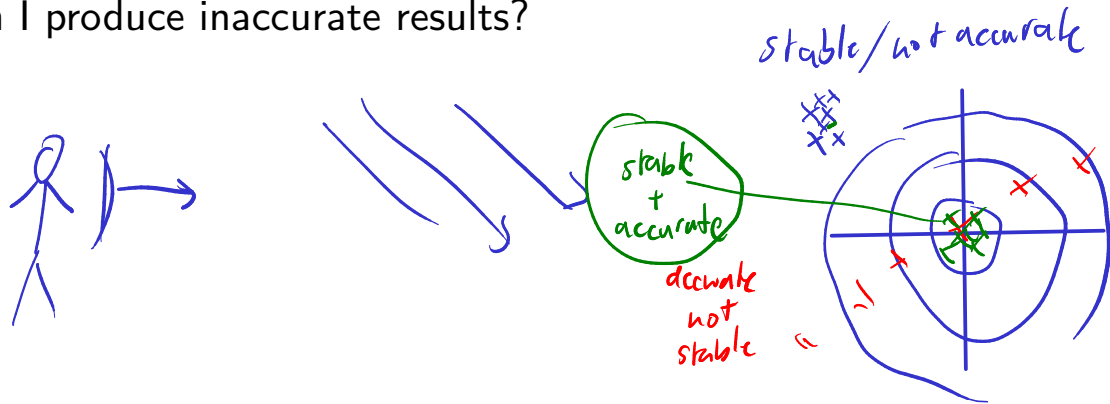
- When is a method **stable**?

If its sensitivity to variation in input is no (or not much) greater than that of the underlying problem.

- When is a method **accurate**?

Closeness of method output to the actual answer for completely accurate input

- How can I produce inaccurate results?



1.2 Floating Point

Wanted: Real Numbers... in a computer

- Computers can represent *integers*, using bits:

$$23 = 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 = (10111)_2$$

How would we represent fractions?

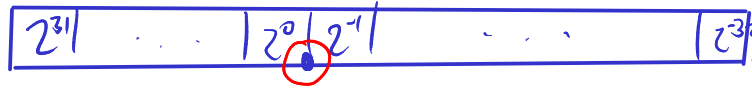
↑	↑	↑	↑	↑
16	8	4	2	1
1	0	1	1	1

$$23.625 = 1 \cdot 2^4 + 0 \cdot 2^3 + (1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0) + \underbrace{1}_{\substack{\uparrow \\ 0.5 \\ \uparrow \\ .125}} \cdot 2^{-1} + \underbrace{0}_{\substack{\uparrow \\ 0.25 \\ \uparrow \\ .125}} \cdot 2^{-2} + \underbrace{1}_{\substack{\uparrow \\ .125 \\ \uparrow \\ 0}} \cdot 2^{-3}$$

Fixed-point arithmetic

Fixed-Point Numbers

- Suppose we use units of 64 bits, with 32 bits for exponents ≥ 0 and 32 bits for exponents < 0 . What numbers can we represent?



biggest: $2^{32} - 2^{-32}$

smallest: 2^{-32}

- How many 'digits' of relative accuracy (think relative rounding error) are available for the smallest vs. the largest number?

Largest: 64 bits: $2^{64} \approx 10^{19} \rightsquigarrow 19$ digits
Smallest: 0 bits $\rightsquigarrow 0$ digits

Floating Point numbers

- Convert $13 = (1101)_2$ into floating point representation.

$$13 = (1.101)_2 \cdot 2^3$$

$$\text{e.g. } 10^{15}$$

- What pieces do you need to store an FP number?

$$\begin{array}{l} \begin{array}{l} \underline{\underline{1}} \\ \uparrow \\ \text{not} \\ \text{stored} \end{array} (1.101)_2 \rightarrow \text{"significand" stored: } -1023 + \underline{1026} \\ \begin{array}{l} \text{stored} \\ 3 \\ +/- \end{array} \rightarrow \text{"exponent"} \\ \rightarrow \text{sign} \end{array}$$

$$12.25 =$$

$$1100.01$$

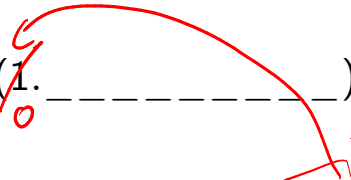
$$1.10001 \cdot 2^3$$

$$.125 = 2^{-3}$$

In-class activity: Floating Point

Unrepresentable numbers?

- Can you think of a somewhat central number that we cannot represent as

$$x = (1.\text{-----})_2 \cdot 2^{-p}?$$


"special exponent" : **-1023**

Demo: Picking apart a floating point number