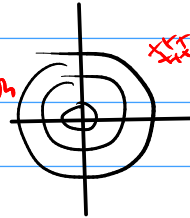


Stability  $\Rightarrow$  accuracy?

$\rightarrow$  "backward stable"  
measure output perturbation  
to the true answer



UFL?

$\rightarrow$  smallest normal (not subnormal)  
number.

## Floating Point and Rounding Error

$$(1.\underline{0}\underline{0}\underline{0}\underline{1})_2 \cdot 2^{p \leftarrow 0}$$

What is the relative error produced by working with floating point numbers?

- What is smallest floating point number  $> 1$ ? Assume 4 bits in the significand.

$$(1.\underline{0}\underline{0}\underline{0}\underline{1})_2 \cdot 2^{p \leftarrow 0}$$

- What's the smallest FP number  $> 1024$  in that same system?

$$(1.\underline{0}\underline{0}\underline{0}\underline{1})_2 \cdot 2^{10}$$

- Can we give that number a name?

machine precision  $\epsilon_{mach}$

the smallest number so that  $fl(1 + \epsilon_{mach}) > 1$

- What does this say about the relative error incurred in floating point calculations?

$$(1.\underline{0}\underline{0}\underline{0}\underline{1})_2 = fl\left(1 + (1.\underline{0}\underline{0}\underline{0}\underline{0}\underline{1})_2\right)$$

- What's that same number for double-precision floating point? (52 bits in the significand)

$$2^{-53} \quad \text{or} \quad 2^{-52} \quad (\text{dep. on rounding})$$

$$fl(a+b) = \tilde{x} \quad a+b = x$$

$$\uparrow$$
$$= x \cdot (1 + \epsilon_{mach})$$

rel. error from FP roundiy

$$\frac{|\tilde{x} - x|}{|x|} = \frac{|x(1 + \epsilon_{mach}) - x|}{|x|} = \epsilon_{mach}$$

## Demo: Floating Point and the Harmonic Series

## Implementing Arithmetic

- How is floating point addition implemented?  
Consider adding  $a = (1.101)_2 \cdot 2^1$  and  $b = (1.001)_2 \cdot 2^{-1}$  in a system with three bits in the significand.

$$\begin{array}{r} a = (1.101)_2 \cdot 2^1 \\ + \\ b = (0.01001)_2 \cdot 2^1 \\ \hline \underline{1.11101} \cdot 2^1 \end{array}$$

## Problems with FP Addition

- What happens if you subtract two numbers of very similar magnitude?  
As an example, consider  $a = (1.1011)_2 \cdot 2^0$  and  $b = (1.1010)_2 \cdot 2^0$ .

$$\rightarrow a = (1.1011)_2 \cdot 2^0$$

$$\rightarrow b = (1.1010)_2 \cdot 2^0$$

---

$$a - b = 0.0001 \cdot 2^0$$

$$\rightarrow = \underline{1.0000} \cdot 2^{-4}$$

## Demo: Catastrophic Cancellation

## 2 Systems of Linear Equations



## 2.1 Theory: Conditioning



## Solving a Linear System

Given:

- $m \times n$  matrix  $A$
- $m$ -vector  $b$

$$Ax = b$$

- What are we looking for here, and when are we allowed to ask the question?

- If  $b \in \text{span}(\text{columns}(A)) \rightarrow$  even if the answer is not unique
- For unique answer, need  $A$  to be invertible.

## Matrix Norms

- What norms would we apply to matrices?

Need:

$$\underbrace{\|Ax\|}_v \leq \underbrace{\|A\|}_? \cdot \underbrace{\|x\|}_v$$

"submultiplicativity"

Given  $\|v\| \leftarrow$  a vector norm

$$\|A\| := \max_{\|x\|=1} \|Ax\|$$

$\leftarrow$  matrix norm

$$\|x\|_1 = |x_1| + |x_2| + \dots + |x_n|$$

$$n \times 1 \rightarrow \underbrace{\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix}}_A \quad \left\| A \cdot \underbrace{\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}}_{\uparrow} \right\|_1 = \left\| \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \right\|_1$$

If you compute the norm of a  $n \times 1$  matrix, you obtain the vector norm of that column vector.

$$1 \times n \quad (\alpha_1 \dots \alpha_n) \cdot \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix} \quad \dots \text{move subtle.}$$

$$\|x\| = 1$$

$$\|A\|_1 = \max_{\text{col } j} \underbrace{\sum_{\text{row } i} |A_{ij}|}$$

$$\|A\|_\infty = \max_{\text{row } i} \sum_{\text{col } j} |A_{ij}|$$

**Demo:** Matrix norms

**In-class activity:** Matrix norms