| Today | Announcements |
|---|---|
| − $\nabla P$ | − HW1 out |
| − num. LA | − Exam1et 0 ongoing |
|  ↳ norms for matrices | − Quiz deadlines todg |
|  ↳ conditiong |  ↳ next Wed. |
|  ↳ solving | |

# 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.

$$a = (1.101)_2 \cdot 2^1$$

$$b = 0.01001)_2 \cdot 2^1$$

round $\quad \left(\begin{array}{l} 1.11101 \quad \cdot 2^1 \\ 1.111 \quad \cdot 2^1 \end{array}\right.$

# Problems with FP Addition

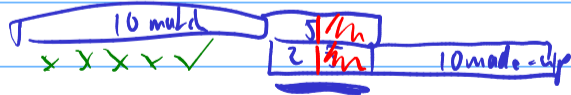What happens if you subtract two numbers of very similar magnitude?
As an example, consider $a = (1.1011)_2 \cdot 2^0$ and $b = (1.1010)_2 \cdot 2^0$.

$$a = 1.1011$$
$$b = 1.1010$$
$$a-b = 0.0001 \cdot 2^0$$
$$1.????\cdot 2^{-4}$$

**Demo:** Catastrophic Cancellation

Suppose $a, b$ are known to $12$ digits
and stored in DP. $\rightarrow$ $52$ bits in the stored fraction

Their $10$ leading digits match.



$\hookrightarrow$ $5$ digits left because of cancellation
$\hookrightarrow$ $2$ left because of rel. error

$$a \cdot 2^c \cdot b \cdot 2^f$$

# Supplementary Material

- Josh Haberman, Floating Point Demystified, Part 1
- David Goldberg, What every computer programmer should know about floating point

# Outline

# Solving a Linear System

Given:
- $m \times n$ matrix $A$
- $m$-vector $\mathbf{b}$

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

$$A \vec{x} = \vec{b}$$

$\leadsto$ lin. comb. of the col. of $A$ to yield $\vec{b}$

$\leadsto$ $m = n$ for now

$\leadsto$ sol. may not exist.

Next: Want to talk about conditioning of this operation. Need to measure

# Matrix Norms

What norms would we apply to matrices?



submultiplicativity

$$\|A x\| \leq \|A\| \|x\| \qquad \text{for all } x \neq 0$$

defining now

$$\frac{\|A x\|}{\|x\|} \leq \|A\|$$

$$= \max_{\|y\|=1} \|A y\|$$

$$\|A\| := \max_{x \neq 0} \frac{\|A x\|}{\|x\|}$$

$$= \max_{x \neq 0} \left\| A \frac{x}{\|x\|} \right\|$$

· norm 1

# Matrix Norm Properties

What is $\|A\|_1$? $\|A\|_\infty$?

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

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

How do matrix and vector norms relate for $n \times 1$ matrices?

$$\left\|\begin{array}{c} n\times1 \\ \tilde{x} \end{array}\right\|_{\|\cdot\|=1} \quad \leq \quad \|Ax\| \quad \underbrace{\quad}_{} \text{"they agree"}$$

**Demo:** Matrix norms

# Properties of Matrix Norms

Matrix norms inherit the vector norm properties:

- $\|A\| > 0 \Leftrightarrow A \neq \mathbf{0}$.
- $\|\gamma A\| = |\gamma| \, \|A\|$ for all scalars $\gamma$.
- Obeys triangle inequality $\|A + B\| \leqslant \|A\| + \|B\|$

But also some more properties that stem from our definition:

$$\|A x\| \leq \|A\| \, \|x\|$$

$$\|A B\| \leq \|A\| \, \|B\|$$

# Conditioning

What is the condition number of solving a linear system $A\mathbf{x} = \mathbf{b}$?

output ↓    input ↓

$\Delta \mathbf{b}$

$\to \Delta \mathbf{x}$

$$\frac{\text{rel. error in output}}{\text{rel. error in input}} = \frac{\|\Delta x\| / \|x\|}{\|\Delta b\| / \|b\|} = \frac{\|\Delta x\| \|b\|}{\|\Delta b\| \|x\|}$$

$$= \frac{\|A^{-1}\Delta b\| \|Ax\|}{\|\Delta b\| \|x\|} \underset{\text{sub}}{\leq} \cdot \|A^{-1}\| \|A\| \cdot \frac{\|\Delta b\| \|x\|}{\|\Delta b\| \cdot \|x\|}$$

↳ shows an upper bound

↳ need to show that bound is sharp
  ↳ find an example that reaches the bound

# Conditioning of Linear Systems: Observations

Showed $\kappa(\text{Solve } A\mathbf{x} = \mathbf{b}) \leq \|A^{-1}\| \, \|A\|$.
I.e. found an *upper bound* on the condition number. With a little bit of fiddling, it's not too hard to find examples that achieve this bound, i.e. that it is *sharp*.

So we've found the *condition number of linear system solving*, also called the condition number of the matrix $A$:

$$\text{cond}(A) = \kappa(A) = \|A\| \, \|A^{-1}\|.$$

# Conditioning of Linear Systems: More properties

▶ cond is relative to a given norm. So, to be precise, use

$$\text{cond}_2 \quad \text{or} \quad \text{cond}_\infty .$$

▶ If $A^{-1}$ does not exist: $\text{cond}(A) = \infty$ by convention.

What is $\kappa(A^{-1})$?

$$\kappa(A) \qquad \|A\underline{v}\| \leq \|A\| \|\underline{v}\|^{?}$$

What is the condition number of matrix-vector multiplication?

$$x \rightarrow Ax$$
$$x \rightarrow A^{-1}y = x \qquad \rightarrow \quad y = Ax$$

**Demo:** Condition number visualized
**Demo:** Conditioning of 2x2 Matrices

# Residual Vector

What is the residual vector of solving the linear system

$$\mathbf{b} = A\mathbf{x}?$$

$$1 = \|I\| = \|AA^{-1}\| \leq \|A\|\,\|A^{-1}\|$$