

Lec. 1 - MAT 226A - Bob Ojij - 9/27

- Dahlquist book has everything + depth
Numerical computations require that you rethink & reformulate problems.

Ex: Compute eigenvalues of a matrix $A \in \mathbb{C}^{n \times n}$

1. compute characteristic polynomial

$$p(\lambda) = \det(A - \lambda I)$$

2. Find the roots of $p(\lambda)$ to get the eigenvalues.

Horrible algorithm to implement numerically!

1. calc $\det(A)$ → standard way is expansion by minors

Let C_n = cost for $n \times n$ matrix

$$C_n = n \cdot C_{n-1}, \quad c_1 = 3 \text{ flops} \Rightarrow C_n = 3^{\frac{n}{2}} n! \text{ flops}$$
$$\Rightarrow C_n = \mathcal{O}(n!)$$

Reformulation: $A = LU$ where L, U are triangular

$\mathcal{O}(n^3)$ work to compute LU factorization

then $\det(A) = \det(L) \det(U)$
 $\overbrace{\det(U)}$ work $\Rightarrow \mathcal{O}(n^4)$ total work

2. Find roots of high-deg poly $p(\lambda) \rightarrow$ poorly conditioned
(very sensitive to round-off errors)

↳ use eigen solver

* Numerically generate $A = QTQ^*$ (Schur decomposition)
where Q unitary, T triangular & eigenvalues are diag. elmts.
Same eigenvals as A .

Floating point #5

Supp. I can store 4 decimal digits, which ones to store?

Naive: Pick first four left of decimal point

Can store 0, smallest = 1, largest = 9999

Pick first four to the right: 0.xxxx

from 0.0001 to 0.9999

To expand range, use scientific notation

$$x = \underbrace{r}_{\substack{1 \text{ digit here} \\ 3 \text{ digits here}}} \cdot 10^n$$

can store #s on 10 orders of magnitude

$$\frac{1}{10} \leq r < 1 \text{ or } r=0$$

(Most) computers store numbers in binary:

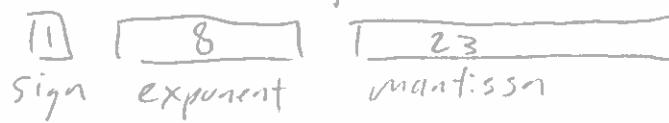
$$x = \underbrace{\pm q}_{\substack{0. \dots \text{ w/ } q=0 \text{ or } \frac{1}{2} \leq q < 1}} \times 2^m$$

↑ know leading bit is 1

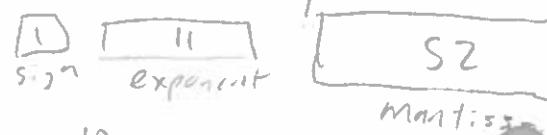
$$\text{Instead } x = \underbrace{\pm (1+q)}_{\substack{1 \text{ bit for sign} \\ \text{for mantissa}}} \times 2^m$$

↑ bits for exponent/
characteristic

IEEE standard: 32-bit precision



64-bit double precision:



(Characteristic in double precision ranges from $c=0$ to $c=\sum_{j=0}^{10} 2^j = 2^{11}-1 = 2047$)

$$\text{Then } m = c - 1023 \Rightarrow$$

$$x = \pm (1+q) \times 2^{c-1023}$$

Double precision: $x = \pm (1+q) \times 2^{c-1023}$

Dont actually use $m=1024$ or $m=1023 \rightarrow NaN$, underflow
smallest positive # representable: $c=1, q=0$ flags!

$$\Rightarrow x = 2^{-1022} \approx 2 \times 10^{-308}$$

largest positive # representable: $x = 2^{1024} \approx 2 \times 10^{308}$
 $q = 1/2 + 1/4 + 1/8 + \dots \approx 1$

Same # of #'s btwn 2^{-1022} and 2^{-1021} as btwn 2^{1022} and 2^{1023}
but relative size of the gaps is the same.

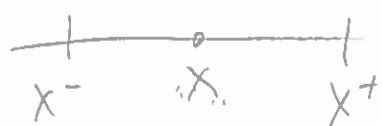
Store 0 as $c=0, q=0$.

How much precision in decimals?

$$\log_{10} 2^{52} \approx 15.7, \log_{10} 2^{53} \approx 15.95$$

\Rightarrow 15-16 decimal digits of accuracy.

What is the error of representing real #'s as floating point #'s?

 $x \in \mathbb{R}, x, x^+ \in \mathbb{F}$ closest float to x .

$fl(x)$ is the floating point representation of x

$$|fl(x) - x| \leq |x^+ - x^-| \quad \text{or if we're rounding.}$$

$$|fl(x) - x| \leq \frac{1}{2} |x^+ - x^-|$$

Assuming same m . $|x^+ - x^-| = 0.0 \dots l_k \times 2^m = 2^{m-k}$
k-bit mantissa

Relative error: $\frac{|fl(x) - x|}{|x|} \leq \frac{2^{m-k}}{|x|} = \frac{2^{m-k}}{\underbrace{1.b_1 b_2 \dots b_k b_{k+1} \dots}_{\text{independent on size of } x} \times 2^m} \leq 2^{-k}$

(independent on size of x ... only depends on # bits > 1)

* Relative error in representing real # as floating point # depends on # bits in the mantissa.

For double precision, $k=52 \Rightarrow 2^{-52} \approx 2.22 \times 10^{-16}$
 call this $\epsilon_{\text{machine}}$ $\xrightarrow{\text{relative spacing between floating point #s}}$

Experiment: $f1(1+10^{-16}) = 1$ if $\epsilon_{\text{machine}} > 10^{-16}$

Thm: $\forall x \in \mathbb{R}$ in some range ($\pm 2 \times 10^{-308}$ to $\pm 2 \times 10^{308}$), $\exists x' \in \mathbb{F}$ s.t. $|x - x'| \leq \epsilon_{\text{machine}} |x|$.

There exists an ϵ w/ $|\epsilon| < \epsilon_{\text{machine}}$
 s.t. $f1(x) = x(1+\epsilon)$

Thm: Let \odot represent $(+, -, \times, \div)$ fl pt. operation
 $\forall x, y \in \mathbb{F}, \exists \epsilon$ with $|\epsilon| < \epsilon_{\text{machine}}$ s.t.
 $x \odot y = x \underset{\text{real operation}}{\hat{\odot}} y (1+\epsilon)$.

Condition # Ex 3) Polynomial Eval.

$$p(x) = \sum_{k=0}^n a_k x^k \rightarrow \text{How sensitive to coefficient perturbations?}$$

$$\frac{\partial p}{\partial a_i} = x^i$$

$$\text{Then } K(a_i) = \left| \frac{\partial p}{\partial a_i} \right| \cdot \frac{|a_i|}{|p(x)|} = \frac{a_i x^i}{|p(x)|}$$

One way to measure condition # wrt. all coeffs

$$K = \frac{\sum_{k=0}^n |a_k x^k|}{\left| \sum_{k=0}^n a_k x^k \right|} \rightarrow \text{poorly-conditioned near roots}$$

↳ look at abs. cond. #!
 $\geq - \sum_{k=0}^n |x^k|$

Lec 3 - MAT226A - 10/2

Need to be careful subtracting nearly equal numbers!

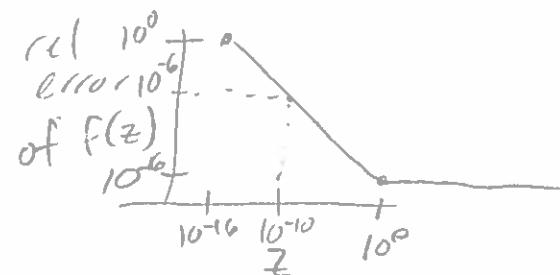
Ex: Evaluate the function $f(z) = e^z - 1$ for z small

$$f(z) = 1 + z + \frac{z^2}{2!} + \dots - 1 = z + \frac{z^2}{2!} + \frac{z^3}{3!} + \dots$$

\Rightarrow For z below ϵ_{mach} , $f'(e^z) = 1$

$$\text{So Rel error } \approx \frac{|1 - e^z|}{|e^z|} = 1$$

What to do? Reformulate the function



1) Use Taylor series to replace f → only good for small z .

2) Do some algebra:
$$\begin{aligned} f(z) &= e^{z/2} (e^{z/2} - e^{-z/2}) \\ &= 2e^{z/2} \sinh(z/2) \end{aligned}$$

Conditioning & Stability

- Conditioning is related to how sensitive a problem is to perturbation.
- Stability is related to the error in an algorithm to solve the problem.

Condition # of a problem: Let $f: X \rightarrow Y$ where X, Y normed vector spaces. Want to know how changes in input affect changes in output.

- Solve $A\underline{x} = \underline{b}$ where $A \in \mathbb{C}^{n \times n}, \underline{b} \in \mathbb{C}^n$. Input is (A, \underline{b}) , output \underline{x}
- Given $\underline{x} \in X$, let $\delta \underline{x}$ be a perturbation of \underline{x} .
(Change in output $\delta \underline{f} = \underline{f}(\underline{x} + \delta \underline{x}) - \underline{f}(\underline{x})$)

Given $x \in X$, δx a perturbation of x
 Change in output $\delta f = f(x + \delta x) - f(x)$

The condition #, $K(x)$ (dep on x) satisfies

$$\frac{\|\delta f\|}{\|f\|} \leq K(x) \frac{\|\delta x\|}{\|x\|}$$

norm of $f(x)$ in Y .

Def: The relative condition # is

$$K(x) = \ln \sup_{\delta \rightarrow 0, \|\delta x\| < \delta} \left(\frac{\|\delta f\|}{\|f\|} \right) \left(\frac{\|\delta x\|}{\|x\|} \right)$$

Def The absolute condition # is (useful when $\|f\|=0$)

$$\hat{K}(x) = \ln \sup_{\delta \rightarrow 0, \|\delta x\| < \delta} \frac{\|\delta f\|}{\|\delta x\|}$$

If f is differentiable, $\hat{K}(x) = \|J\|$

where $J_{ij} = \frac{\partial f_i}{\partial x_j}$ and $\|\cdot\|$ is induced norm by norms on X and Y .

Fact: If f is differentiable, relative condition # is

$$K(x) = \|J\| / \left(\frac{\|f\|}{\|x\|} \right) = \|J\| \cdot \frac{\|x\|}{\|f\|}$$

norm of $f(x)$
in Y

Ex: 1) $f(x) = \sqrt{x} \Rightarrow f'(x) = \frac{1}{2} x^{-1/2} \Rightarrow K(x) = \frac{1}{2\sqrt{x}} \cdot \frac{x}{\sqrt{x}} = \frac{1}{2} \in$ well-cond

2) $f(x_1, x_2) = x_1 - x_2 \Rightarrow J = [1 \ -1] \rightarrow$ use $\|\cdot\|_2$ as induced norm

$$\Rightarrow \|J\|_2 = \sqrt{2} \quad \|x\|_2 = \sqrt{x_1^2 + x_2^2}$$

$\Rightarrow K(x) = \sqrt{2} \cdot \frac{\sqrt{x_1^2 + x_2^2}}{|x_1 - x_2|} \Rightarrow$ subtraction of nearly eq-4! #'s
is poorly-conditioned!

Lec 4 - MAT226A - 10/4

Bob's OH: 4:15-5:30 Weds.

Condition # for $f: X \rightarrow Y$: $\frac{\| \delta f \|}{\| f \|} \leq \kappa(x) \cdot \frac{\| \delta x \|}{\| x \|}$

Ex: Polynomial evaluation

$$p(x) = \sum_{k=0}^n a_k x^k \quad \leftarrow \text{look at condition # w.r.t. perturbations in the coefficients } a_k$$

$$\kappa(x) = \frac{\sum_{k=0}^n |a_k x^k|}{\left| \sum_{k=0}^n a_k x^k \right|} \quad \leftarrow \text{blows up near roots of } p(x)$$

e.g. $p(x) = (x-2)^9 = \sum_{k=0}^9 a_k x^k$

$$\text{At } x=2, \quad \kappa = \sum_{k=0}^9 |a_k x^k|$$

Notice that $\kappa \approx 10^{-11} \Rightarrow \kappa(2) \approx 262,144 = 3 \times 10^5$.
= same magnitude of "jiggies" in Bob's MATLAB solver example

Rule of thumb: rel. error $\leq (\text{cond \#})(\epsilon_{\text{machine}})$

\Rightarrow Abs. cond. # is approx. # of digits we lose.

Stability & Accuracy: $f: X \rightarrow Y$ problem

Hope \hat{f} approx f !

$\hat{f}: X \rightarrow Y$ algorithm

Relative accuracy of algorithm: $\frac{\| f(x) - \hat{f}(x) \|}{\| f(x) \|}$

Hope that $E_{f,\hat{f}} = \Theta(\epsilon_{\text{machine}})$

$\hat{E}_{f,\hat{f}}$

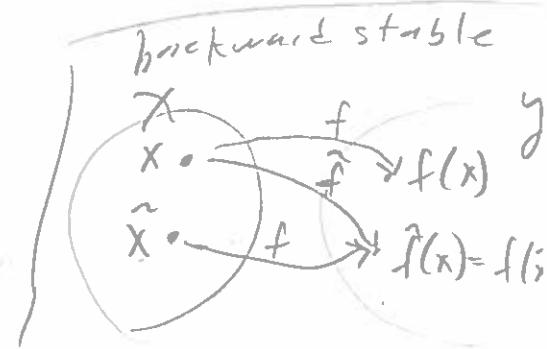
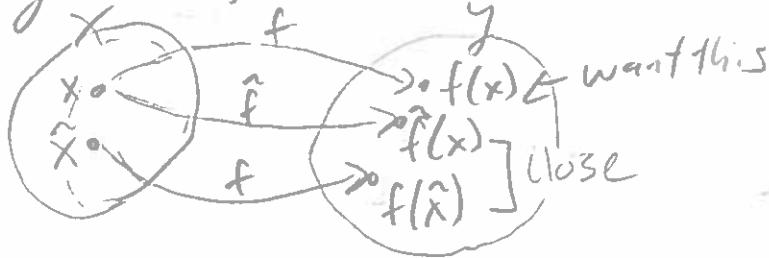
Recall: $\forall t: \mathcal{O}(\varphi(t))$ if $f \in S.t.$ $|f(t)| \leq C\varphi(t)$.

Unfortunately, $E_{f,f} = \mathcal{O}(E_{\text{mach}})$ in some assumed limit, usually $t \rightarrow 0$ or $t \rightarrow \infty$ is too much to ask.
Would like a statement about bounding effects of round-off error still!

Stability: An algorithm \hat{f} is stable if $\forall x \in X \exists \tilde{x} \in X$ s.t.

$$\frac{\|\hat{f}(x) - f(\tilde{x})\|}{\|f(\tilde{x})\|} = \mathcal{O}(E_{\text{mach}}) \quad \frac{\|x - \tilde{x}\|}{\|x\|} = \mathcal{O}(E_{\text{mach}})$$

A stable algorithm gives nearly the right answer to nearly the right question.



A stronger stability is backward stability: \hat{f} is backward stable if

$$\forall x \in X, \exists \tilde{x} \in X \text{ s.t. } \hat{f}(x) = f(\tilde{x}) \text{ and } \frac{\|x - \tilde{x}\|}{\|x\|} = \mathcal{O}(E_{\text{mach}})$$

↳ Backward stable alg gives the exact answer to nearly the right question.
(E.g. $\hat{f}(x) = \text{fl}(x)$ is backward stable since $\text{fl}(x) = \tilde{x}$ within E_{mach} !)

Thm: For a backward-stable alg, $\frac{\|f(x) - \hat{f}(x)\|}{\|f(x)\|} = \mathcal{O}(K(x) E_{\text{mach}})$

Pf:
$$\frac{\|f(x) - \hat{f}(x)\|}{\|f(x)\|} \stackrel{\text{backward stable}}{=} \frac{\|f(x) - f(\tilde{x})\|}{\|f(x)\|} \stackrel{\text{defn}}{\leq} K(x) \frac{\|\tilde{x} - x\|}{\|x\|} = \mathcal{O}(K(x) E_{\text{mach}})$$

Fact: Floating point arithmetic is backward stable!

Ex: $f(x,y) = x + y$, $\text{fl}(x) = x(1 + \varepsilon_1)$, $\text{fl}(y) = y(1 + \varepsilon_2)$ ($\varepsilon_1, \varepsilon_2 < \epsilon$)

$$\Rightarrow \hat{f}(x,y) = \text{fl}(x) \oplus \text{fl}(y) = (x(1 + \varepsilon_1) + y(1 + \varepsilon_2))(1 + \varepsilon_3)$$

$$= x(1 + \varepsilon_1)(1 + \varepsilon_3) + y(1 + \varepsilon_2)(1 + \varepsilon_3) = \tilde{x} + \tilde{y}$$

$$1 + \tilde{\varepsilon} = \sqrt{(1 + \varepsilon_1)(1 + \varepsilon_2)} \approx \sqrt{(1 + \varepsilon_1) / (1 + \varepsilon_2)}$$

• Find a backward-stable algorithm $\frac{\|\tilde{f}(x) - f(x)\|}{\|f(x)\|} \leq \epsilon_{\text{mach}}$
 \Rightarrow If $K = 10^r$, expect to lose r digits of acc.

• Last time, showed ⑦ is backward-stable: $\tilde{f}(x, y) = f(x) \oplus f(y) \leq \epsilon_{\text{mach}}$
 $\hookrightarrow \tilde{f}(x, y) = f(\tilde{x}, \tilde{y})$.

• Given $f(x) = x + 1$, $\tilde{f}(x) = f(x) \oplus 1$ is not backward-stable.

$$\begin{aligned} &= x(1 + \epsilon_1) \oplus 1 = (x(1 + \epsilon_1) + 1)(1 + \epsilon_2) \\ &= \underbrace{x(1 + \epsilon_1)(1 + \epsilon_2)}_{\tilde{x}} + \epsilon_2 + 1 \end{aligned}$$

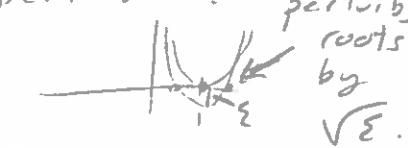
Then $\tilde{f}(x) = f(\tilde{x})$ and $\tilde{x} = x(1 + \epsilon_1)(1 + \epsilon_2) + \epsilon_2$

$$\text{but } \frac{|\tilde{x} - x|}{|x|} = \frac{|x + \epsilon_1 x + \epsilon_2 x + \epsilon_1 \epsilon_2 x + \epsilon_2 - x|}{|x|} = \epsilon_1 + \epsilon_2 + \epsilon_1 \epsilon_2 + \frac{\epsilon_2}{|x|} \neq O(\epsilon_{\text{mach}})$$

Not even at least stable! $\tilde{f}(x) = f(x) \oplus 1$

\Rightarrow Stable, not backward $= (\underbrace{x(1 + \epsilon_1)}_{\tilde{x}} + 1)(1 + \epsilon_2) = f(\tilde{x})(1 + \epsilon_2)$

An unstable algorithm: given square matrix A , find λ 's by finding roots of char poly.

Ex: $I = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \Rightarrow p(x) = x^2 - 2x + 1$ What if this is perturbed?
 \hookrightarrow shift by ϵ : \Rightarrow if $\epsilon = \epsilon_{\text{mach}} = O(10^{-10})$, roots shift by $O(10^{-8})$. 

Conditioning of a linear system $A\underline{x} = \underline{b}$ $A \in \mathbb{C}^{n \times n}$

Vector norm: $\underline{x} \in \mathbb{C}^n$, $\|\underline{x}\|_2 = \left(\sum_{j=1}^n |x_j|^2 \right)^{1/2} = \sqrt{\underline{x}^* \underline{x}} = \sqrt{\langle \underline{x}, \underline{x} \rangle}$, $\underline{x}^* \in \mathbb{C}^n$

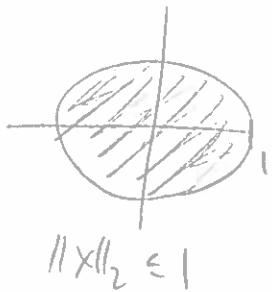
\underline{x}^* col vector $n \times 1$, \underline{x}^* is row vector $1 \times n$ of complex conjugate transpose.

- 1) norm, $\|\cdot\|$, is a fn. $\mathbb{C}^n \rightarrow \mathbb{R}$ s.t. i) $\|\underline{x}\| \geq 0$ & $\|\underline{x}\| = 0 \iff \underline{x} = 0$,
- 2) $\|\alpha \underline{x}\| = |\alpha| \|\underline{x}\| \quad \forall \alpha \in \mathbb{C}$, $\Rightarrow \|\underline{x} + \underline{y}\| \leq \|\underline{x}\| + \|\underline{y}\| \quad \forall \underline{x}, \underline{y} \in \mathbb{C}^n$

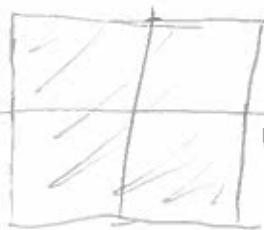
Norms we'll use in this class: $\|x\|_2 = \left(\sum_{j=1}^n |x_j|^2 \right)^{1/2}$

$$\|x\|_\infty = \max_{1 \leq j \leq n} |x_j| , \quad \|x\|_1 = \sum_{j=1}^n |x_j| , \quad \|x\|_p = \left(\sum_{j=1}^n |x_j|^p \right)^{1/p} \quad 1 \leq p < \infty$$

Unit ball in ℓ_2 :

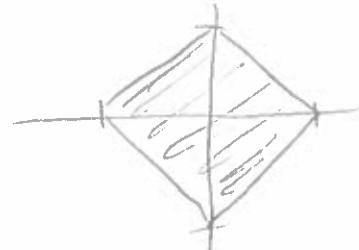


in ℓ_∞ :



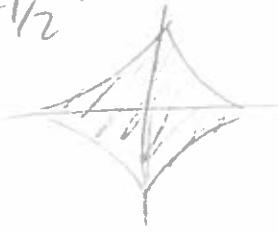
$$\|x\|_\infty \leq 1$$

in ℓ_1 :



$$\|x\|_1 \leq 1$$

in $\ell_{1/2}$:



in ℓ_0 :



Lec 6 - MAT226A - 10/9

Matrix norms : $A \in \mathbb{C}^{m \times n}$

• Could use a norm on \mathbb{C}^{mn} : Frobenius norm: $\|A\|_F = \left(\sum_{ij} |a_{ij}|^2 \right)^{1/2}$
 b) just a measure of the 'stuff' in matrix, no structure! $= \text{Tr}(A^*A)^{1/2}$

• An $m \times n$ matrix defines a linear function $\mathbb{C}^n \rightarrow \mathbb{C}^m$ by
 $f(\underline{x}) = A\underline{x}$

Induced matrix norms reflect the size of the function instead
 of size of entries
 Let's measure the size of the "stretch" \rightarrow look at max stretch

Def: Let $A \in \mathbb{C}^{m \times n}$, $\|\cdot\|_n$ a norm on \mathbb{C}^n , $\|\cdot\|_m$ norm on \mathbb{C}^m , then

$$\|A\|_{(m,n)} = \sup_{\substack{\underline{x} \in \mathbb{C}^n \\ \underline{x} \neq 0}} \frac{\|A\underline{x}\|_m}{\|\underline{x}\|_n} = \max_{\substack{\|\underline{x}\|_n=1}} \|A\underline{x}\|_m$$

↑ unit ball is compact, so $\sup = \max$

Generally, pick the same norms on \mathbb{C}^n and \mathbb{C}^m , & rotate w/ that.

c.7. $\|A\|_1 = \max_{1 \leq j \leq n} \|\underline{a}_{j1}\|_1$, $\|A\|_\infty = \max_{1 \leq i \leq m} \|\underline{a}_{i*}\|_\infty$

↑ \underline{a}_{j1} col of A ↑ \underline{a}_{i*} row of A

Max col. sum Max row sum

c.7. $\|A\|_2 = \sqrt{\text{largest mag. eigenvalue of } A^*A} = \sqrt{\rho(A^*A)} = \max_{\text{value}} \text{sing. } 0_i$

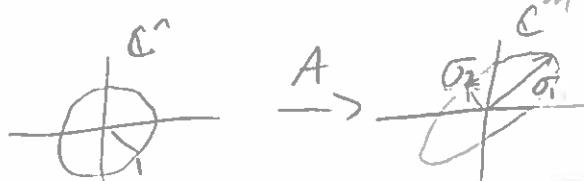


image of unit l_2 -ball is ellipse under A ,
 radii of the semi-axes are singular values

pf of $\|A\|_1 = \max_{1 \leq j \leq n} \|\underline{a}_{j1}\|_1$: $\|A\underline{x}\|_1 = \left\| \sum_{j=1}^n \underline{a}_{j1} x_j \right\|_1 \leq \sum_{j=1}^n \|\underline{a}_{j1}\|_1 |x_j|$
 choose \underline{x} s.t. $\|\underline{a}_{j1}\|_1 = \max_{1 \leq j \leq n} \|\underline{a}_{j1}\|_1$.

Then choose $\underline{x} = e_j$, so $\|A\underline{x}\|_1 = \left\| \sum_{j=1}^n \underline{a}_{j1} x_{jj} \right\|_1 \leq \max_{1 \leq j \leq n} \|\underline{a}_{j1}\|_1$
 $= \|\underline{a}_{11}\|_1 = \max_{1 \leq j \leq n} \|\underline{a}_{j1}\|_1$. □

2-norm. Recall: a matrix Q is called unitary (orthog. & real) if

$$Q^*Q = I_n \text{ i.e., } Q^* = Q^{-1}$$

each entry i,j of Q^*Q is an inner product $\langle q_i, q_j \rangle$

$$Q^*Q = I \Rightarrow q_i^* q_j = \delta_{ij} \Rightarrow \text{cols. are orthonormal}$$

$$\Rightarrow \|Qx\|_2 = \|\underline{x}\|_2 \rightarrow \|Q\|_2 = 1.$$

$$(\text{since } \|Qx\|_2^2 = \underline{x}^* Q^* Q \underline{x} = \underline{x}^* \underline{x} = \|\underline{x}\|_2^2)$$

Recall: If $A^* = A$, then A has real eigenvalues & orthog. eigenvectors
↑ Hermitian (symmetric if real)

& we can write $A = Q^* \Lambda Q$

$$\text{Let } \Lambda \text{ be diagonal. } \|\Lambda \underline{x}\|_2 = \left(\sum_{j=1}^n |\lambda_{jj} x_j|^2 \right)^{1/2} \leq \max_{1 \leq j \leq n} |\lambda_{jj}| \cdot \|\underline{x}\|_2$$

& get equality using $\underline{x} = e_j$ w/ j corrsp. to max λ .

$$\Rightarrow \|\Lambda\|_2 = \max_j |\lambda_{jj}|.$$

Claim: $\|A\|_2 = (\rho(A^*A))^{1/2}$ for $A \in \mathbb{C}^{m \times n}$

$$\begin{aligned} \|A\underline{x}\|_2 &= (\underline{x}^* \underbrace{A^* A \underline{x}}_{\text{Hermitian}})^{1/2} = (\underline{x}^* \underbrace{Q^* \Lambda Q \underline{x}}_{= A^* A})^{1/2} = (\underline{y}^* \Lambda \underline{y})^{1/2} \\ &= \|\Lambda \underline{y}\|_2 \leq \|\Lambda\|_2 \|\underline{y}\|_2 \end{aligned}$$

where $\underline{y} = Q\underline{x}$
note $\|\underline{y}\|_2 = \|\underline{x}\|_2$

Lec 7 - MAT 226A - 10/11

Claim: $\|A\|_2 = \sqrt{\rho(A^*A)} = \sigma_1$

Pf: $\|A_x\|_2 = (x^* A^* A x)^{1/2} = (\gamma^* \Lambda y)^{1/2}$

where $y = Qx$, $A^*A = Q^* \Lambda Q$

$$\Rightarrow \|y\|_2 = \|x\|_2 \text{ since } Q \text{ unitary}$$

$$\begin{aligned} \Rightarrow \|A_x\|_2 &= \|\Lambda^{1/2} y\|_2 \leq \|\Lambda^{1/2}\|_2 \cdot \|y\|_2 = \|x\|_2 \cdot \max_i \lambda_i^{1/2} \\ &= \|x\|_2 \cdot \sqrt{\rho(A^*A)} \\ &= \|x\|_2 \cdot \sigma_1 \end{aligned}$$

Choose x to be the first singular vector/eigenvector of A^*A corresponding to σ_1 , then $\|A_x\|_2 = \sigma_1 \|x\|_2$.

Hence $\|A\|_2 = \sigma_1 = \sqrt{\rho(A^*A)}$. □

From the def. of induced norm, $\|A_x\| \leq \|A\| \|x\|$, we can show that $\|AB\| \leq \|A\| \|B\|$.

Pf: $\|ABx\| \leq \|A\| \cdot \|Bx\| \leq \|A\| \|B\| \|x\|$

$$\Rightarrow \frac{\|ABx\|}{\|x\|} \leq \|A\| \|B\| \text{ & sup over } x \text{ & done.} \quad \square$$

Condition # of a matrix:

Fix $A \in \mathbb{C}^{m \times n}$, condition # of Ax from perturbations to x .

$$K(x) = \limsup_{\delta \rightarrow 0, \|x_\delta\| \leq \delta} \frac{\|A(x + \delta x) - Ax\|}{\|Ax\|} \cdot \frac{\|x\|}{\|\delta x\|}$$

$$= \frac{\|A\delta x\|}{\|Ax\|} \cdot \frac{\|x\|}{\|Ax\|} = \|A\| \cdot \frac{\|x\|}{\|Ax\|}$$

Note: Solving $Ax = b$ s.t. perturbations in b
 $\Rightarrow x = A^{-1}b \Rightarrow \|x\| \leq \|A^{-1}\| \cdot \|A\|$
 same thing!

Special case 1: Consider A square & invertible

$$\frac{\|x\|}{\|Ax\|} = \frac{\|A^{-1}Ax\|}{\|Ax\|} \leq \frac{\|A^{-1}\| \cdot \|Ax\|}{\|Ax\|} = \|A^{-1}\|.$$

$$\Rightarrow \|x\| \leq \|A\| \cdot \|A^{-1}\|.$$

Def: Condition # of the matrix : $\kappa(A) = \|A\| \cdot \|A^{-1}\|$.

Now consider the condition # of solving $Ax = b$ s.t. perturbations in A .

$$(A + \delta A)(x + \delta x) = b$$

$$Ax + (\delta A)x + A(\delta x) + \text{h.o.t.} = b$$

$$\delta x = -A^{-1}(\delta A)x$$

$$\Rightarrow \|\delta x\| = \|A^{-1}(\delta A)x\| \leq \|A^{-1}\| \frac{\|\delta A\| \|x\|}{\|A\|}$$

take \limsup

$$\frac{\|\delta x\|}{\|x\|} \leq \|A^{-1}\| \cdot \frac{\|\delta A\|}{\|A\|}$$

$$\Rightarrow \kappa(A) \leq \|A^{-1}\| \cdot \|A\|$$

Equality is attained $\Rightarrow \kappa(A) = \|A\| \|A^{-1}\|$

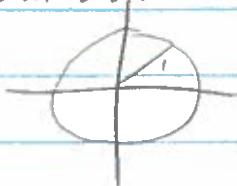
Condition # of solving $Ax = b$

s.t. perturb. in A .

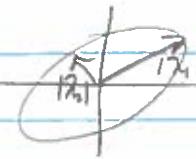
Special case 2: Supp. A Hermitian, i.e. $A^* = A$.

$$\|A\|_2 = \rho(A), \quad \|A^{-1}\| = \rho(A^{-1}) = \frac{1}{\min_j |\lambda_j|}$$

$$\Rightarrow \kappa(A) = \frac{\max_i |\lambda_i|}{\min_i |\lambda_i|} \leftarrow \begin{array}{l} \text{can be } \infty \text{ if} \\ \text{some } \lambda_i = 0 \text{ (i.e. A not-invertible)} \end{array}$$



\xrightarrow{A}



$\leftarrow \kappa(A) = \frac{\text{longest axis}}{\text{shortest axis}} = \text{aspect ratio}$

Then eigendecomp. of A is sum of projection onto axes

Poor conditioning \Rightarrow axes on diff. scales \Rightarrow addition of diff. scales \Rightarrow loss of precision

Lec 7 cont'd

10/11

MAT226A

For non-square matrix A , in 2-norm

$$\|A\|_2 = \frac{\sigma_1}{\sigma_m} = \frac{\text{largest singular value}}{\text{smallest singular value}}$$

Same idea: x is sum of projection onto right-singular vectors,
so if largest σ & smallest σ on diff scales + we lose precision
in our reconstruction of x from singular modes!

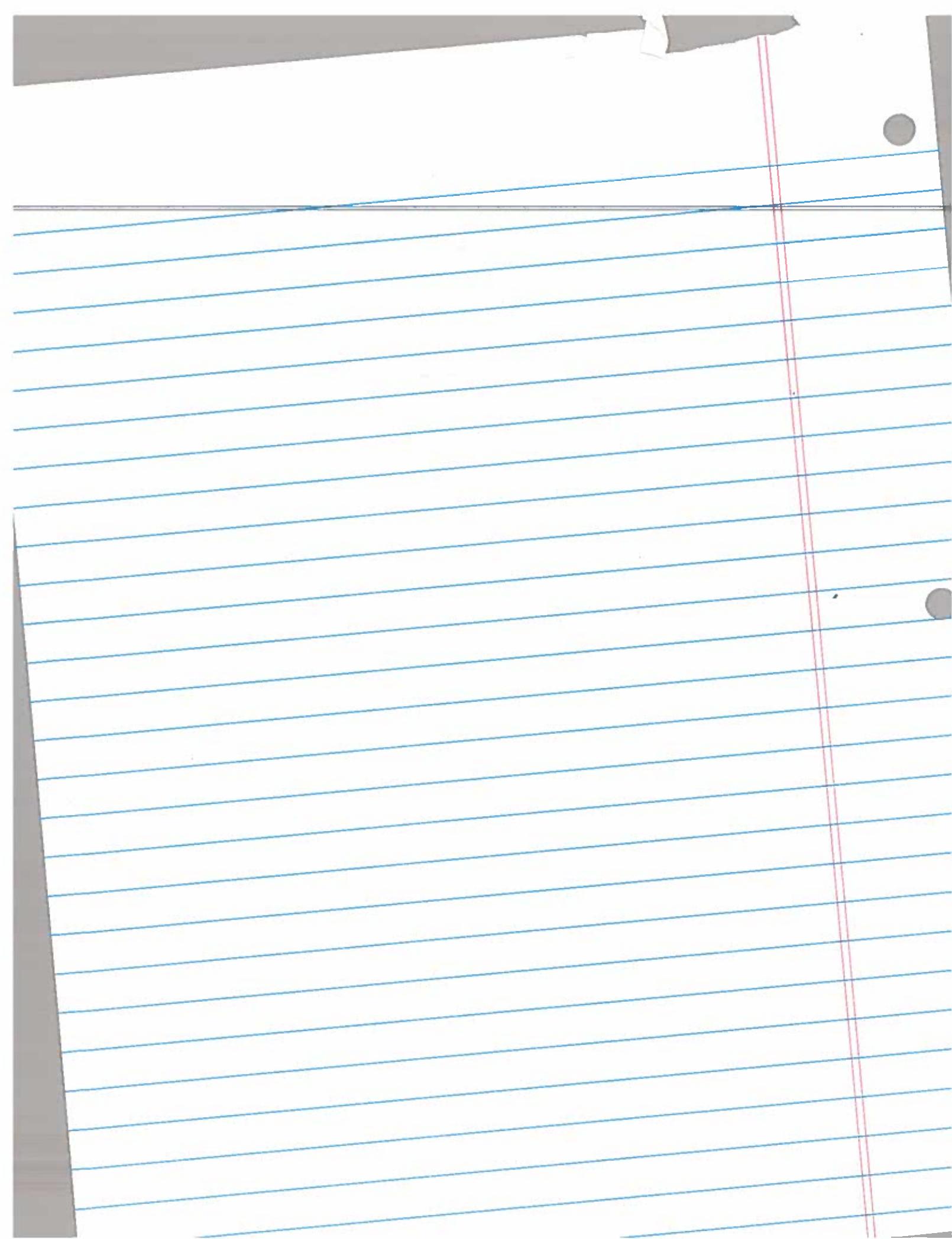
Consider solving $Ax = b$, $A \in \mathbb{C}^{n \times n}$; A invertible.
with no other information on special structure, use Gaussian elimination
(Gaussian elimination = row-reduction)

$$A = \begin{pmatrix} x & x & x \\ x & x & x \\ x & x & x \end{pmatrix} \xrightarrow{\text{mult by } L_1} \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & x & x \end{pmatrix} \xrightarrow{\text{mult by } L_2} \begin{pmatrix} x & x & x \\ 0 & x & x \\ 0 & 0 & x \end{pmatrix} = U$$

$\Rightarrow \underbrace{L_2 L_1}_\text{both are lower-triangular} A = U$ upper triangular \rightarrow generating triangular factorization

$$A = (L_2 L_1)^{-1} U = \underbrace{L_2^{-1} L_1^{-1}}_\text{still lower triangular} U = L U$$

upper triangular



MAT226A-Lec 8 - 10/13

Gaussian elimination is effectively triangular factorization

$$A = LU = \text{lower } \Delta \cdot \text{upper } \Delta$$

Ex:

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ -2 & 1 & 0 & 0 \\ -4 & 0 & 1 & 0 \\ -3 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 2 & 1 & 1 & 0 \\ 4 & 3 & 3 & 1 \\ 8 & 7 & 9 & 5 \\ 6 & 7 & 9 & 8 \end{pmatrix} = \begin{pmatrix} 2 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 3 & 5 & 5 \\ 0 & 4 & 6 & 8 \end{pmatrix}$$

L_1 Δ A $L_1 A$

To eliminate a_{21} , second row of L_1 : $-2R_1 + R_2$

To eliminate a_{31} , third row of L_1 : $-4R_1 + R_3$

To elim. a_{41} , fourth row of L_1 : $-3R_1 + R_4$

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & -3 & 1 & 0 \\ 0 & -4 & 0 & 1 \end{pmatrix} \begin{pmatrix} 2 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 3 & 5 & 5 \\ 0 & 4 & 6 & 8 \end{pmatrix} = \begin{pmatrix} 2 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 2 & 2 \\ 0 & 0 & 2 & 4 \end{pmatrix}$$

L_2 $L_1 A$ $L_2(L_1 A)$

To eliminate $(L_1 A)_{32}$, third row of L_2 : $-3R_2 + R_3$

To eliminate $(L_1 A)_{42}$, fourth row of L_2 : $-4R_2 + R_4$

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & -1 & 1 \end{pmatrix} \begin{pmatrix} 2 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 2 & 2 \\ 0 & 0 & 2 & 4 \end{pmatrix} = \begin{pmatrix} 2 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 2 & 2 \\ 0 & 0 & 0 & 2 \end{pmatrix}$$

L_3 $L_2 L_1 A$ U

$$\Rightarrow L_3 L_2 L_1 A = U \Rightarrow L = \underbrace{(L_3 L_2 L_1)^{-1}}_{\text{computation is trivial}} = L_1^{-1} L_2^{-1} L_3^{-1}$$

Compute inverses on the fly!

$$L_1^{-1} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ -2 & 1 & 0 & 0 \\ -4 & 0 & 1 & 0 \\ -3 & 0 & 0 & 1 \end{pmatrix}^{-1} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 4 & 0 & 1 & 0 \\ 3 & 0 & 0 & 1 \end{pmatrix} \quad \text{super trivial}$$

Product also trivial

$$L_1^{-1} L_2^{-1} L_3^{-1} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 2 & 1 & 0 & 0 \\ 4 & 3 & 1 & 0 \\ 3 & 4 & 1 & 1 \end{pmatrix} \quad \text{& just throwing in entries into } L \text{ as we go.}$$

In general, $L_k = \begin{pmatrix} 1 & & & 0 \\ 0 & \ddots & & \\ 0 & & \ddots & \\ 0 & & & 1 \end{pmatrix}$ w/ $l_{jk} = \frac{x_{jk}}{x_{kk}}$

ent. j
in mat.
 \downarrow
 k^{th}
step
of
elim

$$\underline{l}_k = \begin{pmatrix} 0 \\ \vdots \\ l_{k+1,k} \\ \vdots \\ l_{n,k} \end{pmatrix} \Rightarrow L_k = I - \underline{l}_k e_k^T$$

Show inverse in general: $(I - \underline{l}_k e_k^T)(I + \underline{l}_k e_k^T)$

$$\begin{aligned} &= I - \underline{l}_k \underline{l}_k^T + \underline{l}_k \underline{l}_k^T - \underline{l}_k (\underline{l}_k^T \underline{l}_k) \underline{l}_k^T \\ &= I. \quad \downarrow \text{if elem of } \underline{l}_k = 0. \end{aligned}$$

Hence $L_k^{-1} = (I + \underline{l}_k \underline{l}_k^T)$

General product: $L_k^{-1} L_{k+1}^{-1} = (I + \underline{l}_k \underline{l}_k^T)(I + \underline{l}_{k+1} \underline{l}_{k+1}^T)$

$$\begin{aligned} &= I + \underline{l}_k \underline{l}_k^T + \underline{l}_{k+1} \underline{l}_{k+1}^T + \underline{l}_k (\underline{l}_k^T \underline{l}_{k+1}) \underline{l}_{k+1}^T \\ &= I + \underline{l}_k \underline{l}_k^T + \underline{l}_{k+1} \underline{l}_{k+1}^T \quad \begin{matrix} \uparrow \\ \text{if elem of } \underline{l}_{k+1} \\ \text{is also } 0 \end{matrix} \end{aligned}$$

Lec 8 cont'd

Simple LU factorization algorithm (to be modified by Monday)

Given $A \in \mathbb{C}^{n \times n}$, compute $L+U$ for $A=LU$.

Init: $U=A, L=I$

Iterate: for $k=1$ to $n-1$ (loop over columns)

for $j=k+1$ to n

$l_{jk} = u_{jk}/u_{kk}$ & note we're in trouble here if $u_{kk}=0$

$u_{j,k:n} = u_{j,k:n} - l_{jk}u_{k,k:n}$ & now j^{th} row is lin. comb. of
end

end

Operation Count - count floating-point operations - flops



To compute l_{jk} 's in col k , do $n-k$ divisions
Then changing rows need $2(n-k)^2$ flops
(mult & addition)

$$\text{total flops} = \sum_{k=1}^{n-1} 2(n-k)^2 + (n-k) = \frac{2}{3}n^3 + \text{l.o.t.}$$

$$\text{Recall: } \sum_{k=1}^{n-1} 2k^2 + k \Rightarrow \sum_{k=1}^n k^2 = \frac{1}{3}n^3 + \text{l.o.t.}$$

$$\Rightarrow \text{total flops} = \mathcal{O}(n^3). \quad (\text{where } A \in \mathbb{C}^{n \times n})$$

The algorithm last lecture is unstable.

$\begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix}$ can't be LU-factorized w/ GE algorithm
but reordering rows yields a LU factorization

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}$$

P A L U

Perturb this example: $A = \begin{pmatrix} 10^{-20} & 1 \\ 1 & 1 \end{pmatrix} \leftarrow$ GE will work now

$$L = \begin{pmatrix} 1 & 0 \\ 10^{20} & 1 \end{pmatrix} \quad U = \begin{pmatrix} 10^{-20} & 1 \\ 0 & 1-10^{-20} \end{pmatrix} \quad \text{but actually } \tilde{U} = \begin{pmatrix} 10^{-20} & 1 \\ 0 & -10^{20} \end{pmatrix}$$

w'd get

$$\Rightarrow L \tilde{U} = \begin{pmatrix} 1 & 0 \\ 10^{20} & 1 \end{pmatrix} \begin{pmatrix} 10^{-20} & 1 \\ 0 & -10^{20} \end{pmatrix} = \begin{pmatrix} 10^{-20} & 1 \\ 1 & 0 \end{pmatrix} = \tilde{A}$$

massive error here!

Sup. $\underline{b} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$, solve $A\underline{x} = \underline{b}$

actual soln $\underline{x} \approx \begin{pmatrix} -1 \\ 1 \end{pmatrix}$, but numerical soln $\hat{\underline{x}} \approx \begin{pmatrix} 0 \\ 1 \end{pmatrix} \leftarrow$ 100% error in max-norm
why?
 $\hookrightarrow K(A) \approx 2.6$, but $K(L) \approx K(U) \approx 10^{40}$ awful!

Try w/ permuting rows!

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 10^{-20} & 1 \\ 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 10^{-20} & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 10^{20} & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 0 & 1-10^{-20} \end{pmatrix}$$

$K(L) \approx 1 \quad K(U) \approx 2.6$

No explosive growth in this case!

lec 9 cont'd

- 10/16

This idea of swapping rows to keep condition #s low is called Pivoting.

Gaussian elimination w/ Pivoting:

$$\begin{pmatrix} x & x & x & x \\ 0 & x_{kk} & x & x \\ 0 & x_{ik} & x & x \\ 0 & x & x & x \end{pmatrix} \xrightarrow{\text{usual}} \begin{pmatrix} x & x & x & x \\ 0 & x_{kk} & x & x \\ 0 & 0 & x & x \\ 0 & 0 & x & x \end{pmatrix}$$

$$\begin{pmatrix} x & x & x & x \\ 0 & 0 & x & x \\ 0 & x_{ik} & x & x \\ 0 & 0 & x & x \end{pmatrix}$$

↓ how about
* don't have to go down the line, can look at submatrix & decide which rows/cols to eliminate

$$\begin{pmatrix} x & x & x & x \\ 0 & \boxed{x} & x & x \\ 0 & x & x & x \\ 0 & x & x & x \end{pmatrix}$$

- Want to avoid dividing by numbers
- * Full-pivoting: pick the largest magnitude entry in submatrix for the "pivot"
 - * Partial-pivoting: pick biggest entry in column for the "pivot".

Note: Full-pivoting is more stable, but much more expensive!

To find pivot, compares n^2 entries n times $\Rightarrow \mathcal{O}(n^3)$ work

Partial pivoting compares n entries n times to find pivots $\Rightarrow \mathcal{O}(n^2)$ work.

Algorithm:

- Pick the largest entry in current column (diagonal & below)
- Swap rows to move the pivot to the diagonal.

$$\begin{pmatrix} x & x & x & x \\ 0 & \boxed{x_{kk}} & x_{k.} & x_{.k} \\ 0 & x & x & x \\ 0 & x_{ik} & x_{i.} & x_{.i} \end{pmatrix} \xrightarrow{\text{row swap}} \begin{pmatrix} x & x & x & x \\ 0 & x_{in} & x_{i.} & x_{.i} \\ 0 & x & x & x \\ 0 & x_{ik} & x_{i.} & x_{.i} \end{pmatrix} \leftarrow \begin{array}{l} \text{row swap represented by} \\ \text{mult. by a permutation matrix.} \\ \hookrightarrow \text{Identity matrix except} \\ \text{w/ permutation we want} \end{array}$$

Permutation matrix ex: Swap 2nd & 4th rows:

$$P = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$$

note P is symmetric & unitary
 $P = P^T = P^{-1}$. (only true for single swap)

Is this still LU factorization? Looks like we've permuted L so much we lost its structure.

$$\underbrace{L_n P_n \cdots L_1 P_1}_\text{can show } PA = LU \text{ where } P = P_n \cdots P_1 A$$

What is this then?

How do we compute these P 's & L 's?

$$\text{eg. } L_3 \underbrace{P_3 L_2 (P_3^{-1} P_2)}_{L'_2} P_2 L_1 P_1 A$$

$$= L_3 L'_2 \underbrace{P_3 P_2 L_1 (P_2^{-1} P_3^{-1} P_3 P_2)}_{L'_1} P_1 A$$

$$= \underbrace{L_3 L'_2 L'_1}_{L'} \underbrace{(P_3 P_2 P_1)}_P A = L' P A$$

P_3 permutes L_2 from 3rd row & below \rightarrow shifts entries, same structure

L' has same structure as L_2

$$P_3 L_2 P_3^{-1} = P_3 (I - \underline{e}_2 e_2^\top) P_3^\top = I - \underbrace{P_3 \underline{e}_2 e_2^\top P_3^\top}_{P_3 \text{ permutes } e_2 \text{ from 3rd row & below} \rightarrow \text{all zero!}} = I - \widehat{P}_3 \underline{e}_2 e_2^\top$$

Partial-Pivot GE algorithm Init $U=A$, $L=I$, $P=I$

for $k=1$ to $n-1$ (over cols)
 - select $\max_{i \geq k} |u_{ik}|$ & swap i & k rows in U, L, P

stable in practice

for $j=k+1$ to n (over rows)

$$l_{jk} = u_{jk}/u_{kk}$$

$$u_{j,k:n} = u_{j,k:n} - l_{jk} \cdot u_{k,k:n}$$

end

\rightarrow produces $PA=LU$.

GE w/ pivoting - small correction

- Step w/o pivoting pivot \rightarrow swap row in $i \& k$ in U, L, P
 \hookrightarrow cannot do this in $L \rightarrow$ don't want to mess w/ l's diag.

k^{th} step: $L:$

instead, just swap
i-th & k-th rows in L
for cols 1:k-1

Stability of GE w/ pivoting

Recall w/o pivoting, the entries of L and U can get arb. large
 \hookrightarrow unstable, solns can grow

W/ pivoting, we certainly avoid large entries in L , so
 what about U ?

Thm: Let A non-singular and $A = LU$ be computed w/o pivoting.
 If the algorithm completes (w/ no pivots), then the computed \tilde{L}, \tilde{U} satisfy
 $\tilde{U} = A + \delta A$, where $\frac{\|\delta A\|}{\|L\| \|U\|} = \mathcal{O}(\epsilon_{\text{machine}})$

* Looks like def of backwards stability, but don't have $\|A\|$ in denom!

We know $\|A\| \leq \|L\| \|U\|$, but this direction doesn't help...

With pivoting, L has 1's on diag. and $|l_{ij}| \leq 1$ below diag.

so $\|L\|_1 \leq n \Rightarrow \|L\| = \mathcal{O}(1)$

* norm of L is nice w/ pivoting (the entire point
 of pivoting was to keep the norm bounded)
 \Rightarrow mult. b/w L will be $\mathcal{O}(1)$ instead of $\mathcal{O}(n)$

What about $\|U\|$? Define $\rho = \frac{\max_{ij} |u_{ij}|}{\max_{ij} |a_{ij}|}$ as growth factor

Then $\|U\| \leq \rho \|A\|$, so $\|U\| = O(\rho \|A\|)$, hence

$$\tilde{L}\tilde{U} = \tilde{\rho}A + \delta A$$

$$\Rightarrow \frac{\|\delta A\|}{\|A\|} = O(\rho \epsilon_{\text{mach}})$$

In particular, if $\rho = O(1)$, then GE w/ partial pivoting is backward stable.

Ex (worst case scenario of A):

$$\begin{pmatrix} 1 & 0 & 0 & 1 \\ -1 & 1 & 0 & 1 \\ -1 & -1 & -1 & 1 \\ -1 & -1 & -1 & 1 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & -1 & 1 & 2 \\ 0 & -1 & -1 & 2 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 4 \\ 0 & 0 & -1 & 4 \end{pmatrix} \rightarrow \begin{pmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 2 \\ 0 & 0 & 1 & 4 \\ 0 & 0 & 0 & 8 \end{pmatrix}$$

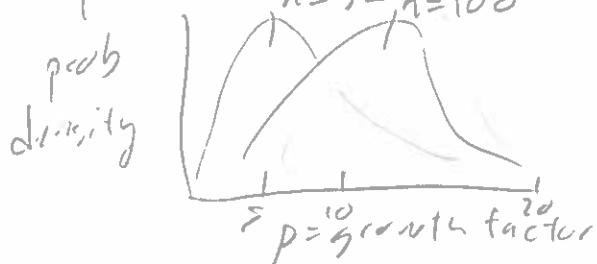
\Rightarrow for a matrix of this structure of size $n \times n$, $\rho \leq 2^{n-1}$

This still gives backwards stability! $\frac{\|\delta A\|}{\|A\|} = O(2^{n-1} \cdot \epsilon_{\text{mach}})$

Ex: $n=100, 2^{99}$ is theoretical bound $\approx 6 \times 10^{29} \Rightarrow$ gigantic & useless!

In practice, growth factor is much smaller

prob. of getting huge ρ is small, pretty much zero.



See HW 2
#3c

- Consider special case examples w/ more info on A .

Case: Tridiagonal matrix $A \rightarrow$ only nonzero on diag, sub, & superdiag.

Assume no pivoting

When we do GE on tridiag, don't introduce any nonzero where there were zeros before

Structure becomes:

$$\begin{pmatrix} A \\ \vdots \end{pmatrix} = \begin{pmatrix} L \\ \vdots \\ U \end{pmatrix} \begin{pmatrix} D \\ \vdots \\ D \end{pmatrix}$$

superdiag is just superdiag of A .

(\sim) represents just 2 vectors $\Rightarrow \mathcal{O}(n)$ work.

Case: (Hermitian/Symm) Positive Definite Matrices

$A^* = A, x^* A x > 0, \forall x \neq 0 \Rightarrow$ all eigenvalues of A are real & positive

(Cholesky factorization) $A = R^* R$.

$$\text{Ex: } A = \begin{pmatrix} 1 & -w^* & -1 \\ w & K & 1 \\ -1 & 1 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ w & I & 0 \\ -1 & 1 & 1 \end{pmatrix} \begin{pmatrix} 1 & -w^* & -1 \\ 0 & K-ww^* & 1 \\ 1 & 1 & 1 \end{pmatrix}$$

$$= \begin{pmatrix} 1 & 0 \\ w & I \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & K-ww^* \end{pmatrix} \begin{pmatrix} 1 & w^* \\ 0 & I \end{pmatrix}$$

R_1^*

if this submatrix is symm posdef,
then we can just repeat this

to generate series R_1, R_2, R_3, \dots & get factorization $A = R^* R$

Lec 11 - 10/20 - MAT226A

Continue example, what if $a_{11} \neq 1$?

$$A = \begin{pmatrix} a_{11} & w^* \\ w & K \end{pmatrix}$$

Let $\alpha = \sqrt{a_{11}}$ (possible since $e_1^T A e_1 = a_{11} > 0$)

$$\text{Then } A = \begin{pmatrix} a_{11} & w^* \\ w & K \end{pmatrix} = \begin{pmatrix} \alpha & w^* \\ w & K \end{pmatrix} \underbrace{\begin{pmatrix} 1 & 0 \\ 0 & \frac{K-ww^*}{a_{11}} \end{pmatrix}}_{R_1^*} R_1$$

Repeat on matrix $\frac{K-ww^*}{a_{11}}$. get $A = R_1^* R_2^* \dots R_n^* R_n R_{n-1} \dots R_1$
easily show R_i 's combine naturally $\Rightarrow A = R^* R$.

This is the Cholesky factorization.

Exploits symmetry, so work $\sim \frac{1}{3}n^3$ (2x as fast as LU fact.)

* Cholesky doesn't need pivoting!

(last time, then: $\tilde{L}\tilde{U} = A + \delta A$, $\frac{\|\delta A\|}{\|L\| \|U\|} = O(\epsilon_{mach})$)

For chd, $\tilde{R}\tilde{R}^T = A + \delta A$, $\frac{\|\delta A\|}{\|\tilde{R}\| \|\tilde{R}^T\|} = O(\epsilon_{mach})$

$$\|\tilde{R}\|_2 = \|\tilde{R}\|_F = \sqrt{\rho(R^* R)} = \sqrt{\rho(A)} = \|A\|_2^{1/2} \text{ since } A \text{ is symmetric.}$$

$$\Rightarrow \frac{\|\delta A\|}{\|A\|_2} = O(\epsilon_{mach}) \quad \Rightarrow \text{Cholesky factorization is backwards stable!}$$

We talked about direct methods, but not iterative methods!

↳ ZGZB

Nonlinear Equations

$f: D \rightarrow \mathbb{R}^n$ solve $f(\underline{x}) = 0$
 $D \subset \mathbb{R}^n$

equivalent to $\begin{cases} f_1(x_1, \dots, x_n) = 0 \\ \vdots \\ f_n(x_1, \dots, x_n) = 0 \end{cases}$] solve n nonlinear equations for n unknowns

- For nonlinear eqns, cannot find a general algorithm to get the solution in a finite # of steps.
 - ↪ Know that there's no formula for the roots of polynomials of deg ≥ 5 .
- Use algorithms which give an approximate solution to a desired accuracy.

Ex: Want $\sqrt{3} = 1.?$ Could guess & check $(1, a)^2 \geq 3?$

More systematically:

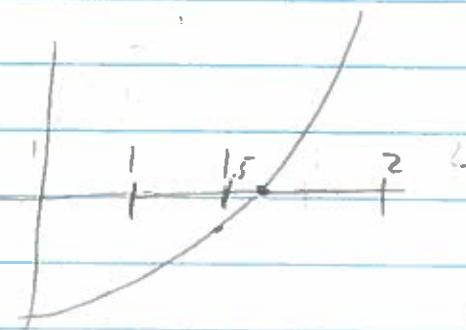
$$f(x) = x^2 - 3$$

have soln in $(1, 2)$

try $x = 1.5, f(1.5) < 0$

so soln in $(1.5, 2)$

If so, decrease a , else increase a



(desired accuracy)

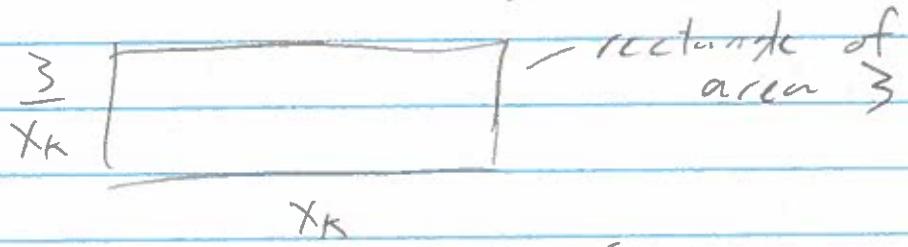
↪ Bisection Method

↪ great for scalar problems, guaranteed convergence

↪ not generalizable to higher dims. ↪ each step cuts ... in half

How fast does bisection method converge? Each step halves error.
 \Rightarrow 3-4 steps per digit of accuracy. \rightarrow linear convergence

Geometric soln to finding $\sqrt{3}$: Babylonian algorithm



Check if $\frac{3}{x_k}$ & x_k equal, if not, one is too high, the other low
 \therefore next guess $\frac{1}{2}(x_k + \frac{3}{x_k}) = x_{k+1}$

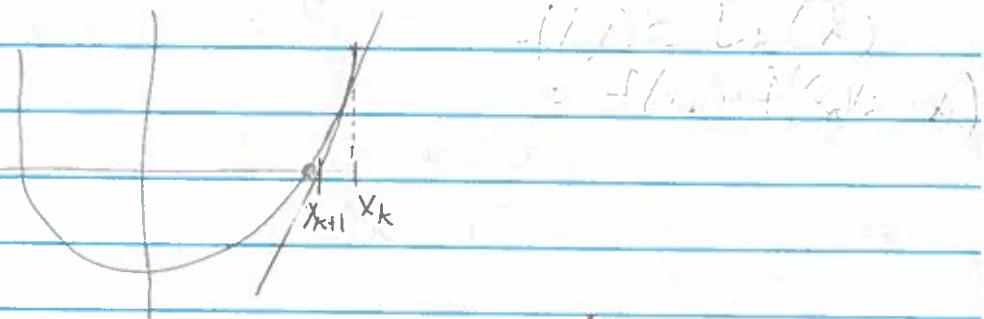
k	$f(x_k)$	$\frac{x_k - \sqrt{3}}{\sqrt{3}}$	
0	1	1.55×10^{-1}	
1	6.25×10^{-2}	1.04×10^{-2}	
2	3.19×10^{-4}	5.31×10^{-5}	
3	8.47×10^{-9}	1.41×10^{-9}	
4	4.44×10^{-16}	$< 10^{-16}$	

Super-linear convergence

Lec 12 - MAT226A - 10/23

Newton's Method

Def: Use linear approximation to find the root of $f(x)$



$$f(x) \approx L_K(x) = f(x_K) + f'(x_K)(x - x_K)$$

Solve $L_k(x_{k+1}) = 0$

$$\Rightarrow f(x_k) + f'(x_k)(x_{k+1} - x_k)$$

$$(b) \text{ Newton's Method: } x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}$$

$$\begin{aligned} \underline{\text{Ex:}} \quad f(x) &= x^2 - 3 \Rightarrow x_{k+1} = x_k - \frac{x_k^2 - 3}{2x_k} \\ f'(x) &= 2x \end{aligned}$$

Same as
geometric
algorithm!

$$x_{k+1} = \frac{1}{2}(x_k + 3/x_k)$$

How fast does $\epsilon_k = x_k - \sqrt{3} \rightarrow 0$?

$$x_{k+1} - \sqrt{3} = -\sqrt{3} + \frac{1}{2}x_k + \frac{3}{2x_k} = \frac{x_k^2 - 2\sqrt{3}x_k + 3}{2x_k}$$

$$x_{k+1} - \sqrt{3} = \frac{(x_k - \sqrt{3})^2}{2x_k}$$

$$l_{k+1} = l_k^2 / 2x_k$$

If x_k is "close" to $\sqrt{3}$, $|e_{k+1}| \leq C |e_k|^2$

Newton's Method for Systems of eqns.

$$\underline{F}(\underline{x}) : \mathbb{R}^n \rightarrow \mathbb{R}^n$$

Solve $\underline{F}(\underline{x}) = \underline{0}$. Guess an initial soln, linearize \underline{F} near guess.

$$\underline{F}(\underline{x}) \approx \underline{L}_k(\underline{x}) = \underline{F}(\underline{x}_k) + \underline{J}(\underline{x}_k)(\underline{x} - \underline{x}_k)$$

$$\text{Find } \underline{L}_k(\underline{x}_{k+1}) = \underline{0} \quad \rightarrow J_{ij} = \frac{\partial F_i}{\partial x_j}$$

$$\Rightarrow \underline{F}(\underline{x}_k) + \underline{J}(\underline{x}_k)(\underline{x}_{k+1} - \underline{x}_k) = \underline{0}$$

$$\text{Define } \delta_k = \underline{x}_{k+1} - \underline{x}_k \Rightarrow \underline{x}_{k+1} = \underline{x}_k + \delta_k$$

$$1. \text{ Solve } \underline{J}(\underline{x}_k) \delta_k = -\underline{F}(\underline{x}_k)$$

$$2. \text{ Update } \underline{x}_{k+1} = \underline{x}_k + \underline{\delta}_k$$

How accurate is the linear approximation?

Supp. f is C^2 . The approx $L(y) = f(x) + f'(x)(y-x)$ ignores quadratic terms. Taylor's thm says $\exists \gamma$ b/w x & y s.t.

$$f(y) = L(y) + R(y)$$

$$\text{where } R(y) = \frac{f''(\gamma)}{2}(y-x)^2,$$

$$\text{Thus } |f(y) - f(x) - f'(x)(y-x)| = \left| \frac{f''(\gamma)}{2}(y-x)^2 \right| \leq \frac{\gamma}{2}(y-x)^2$$

$$\text{where } \gamma = \max |f''|$$

For this result, don't need $f \in C^2$, only need f' to be Lipschitz cont's w/ Lipschitz constant γ .

Recall: f is Lipschitz cont's iff $\exists \gamma$ s.t. $\forall x, y \in D$

$$|f(x) - f(y)| \leq \gamma |x-y|$$

strictly w/ constant γ on D .

Lec 12 cont'd

10/23

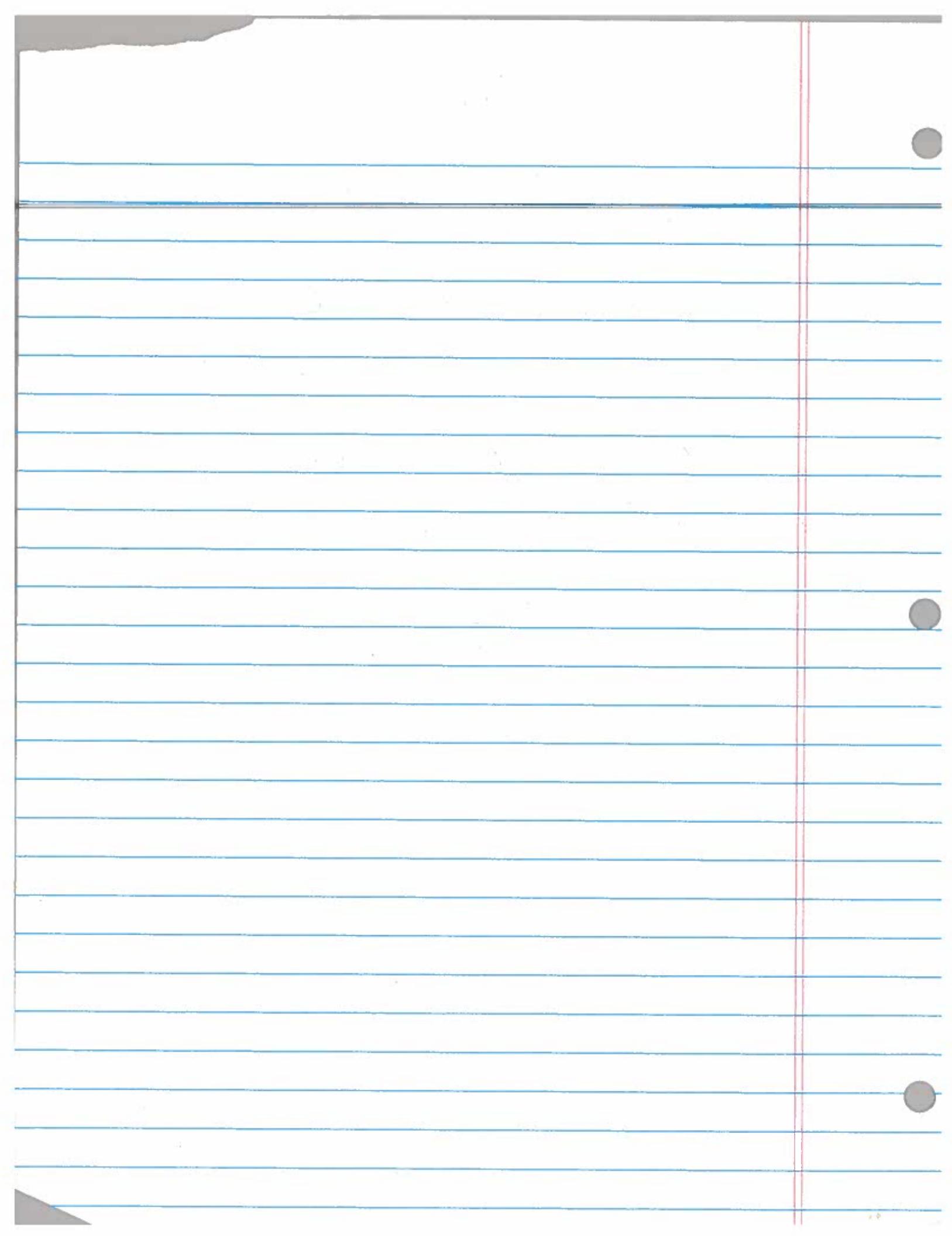
If g is C^1 , then g is Lipschitz cont's with constant $\gamma = \max |g'|$.

(Newton-Raphson) Convergence Thm: $f: D \rightarrow \mathbb{R}$ for an open interval D
 f' is Lipschitz on D w/ constant γ .
(or $|f''| \leq \gamma$ on D)

Suppose $|f'(x)| \geq \rho > 0$.

If $f(x^*) = 0$ for some $x^* \in D$, then there is some η s.t. if $|x_0 - x^*| < \eta$,
then $x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)} \rightarrow x^*$

$$\text{and } |x_{k+1} - x^*| \leq \frac{\gamma}{2\rho} |x_k - x^*|^2$$



Flexible beam in low-Reynolds # fluid

$$\text{discretize } \dot{y} \quad \frac{dy(x,t)}{dt} = -\kappa \frac{\partial^2 y}{\partial x^2} + f(t) ?$$

discrete space $x \rightarrow \text{approx } \frac{\partial^2 y}{\partial x^2} \rightarrow \text{get system of ODEs}$

$$\frac{dy}{dt} = Ay + f$$

$$\frac{y^{n+1} - y^n}{\Delta t} = Ay^n + f^n \quad \leftarrow \text{stiff ODE, requires tiny } \Delta t$$

Use Backward-Euler:

$$\frac{y^{n+1} - y^n}{\Delta t} = Ay^{n+1} + f^{n+1} \rightarrow \text{solve system for } y^{n+1} \text{ at each step}$$

large amplitude. $\frac{dx}{dt} = F(x) + g$ \leftarrow nonlinear function

Implicit method \rightarrow Backward-Euler: $\frac{x^{n+1} - x^n}{\Delta t} = F(x^{n+1}) + g^{n+1}$

stick all on one side & root find for x^{n+1}

\hookrightarrow use previous time step as initial guess in root find.

flow $\rightarrow \{ \{ \{ \{ u(x,t)$

$$Eu_{yy} - \Theta u - G = 0$$

$$u(0,t) = u(1,t) = 0$$

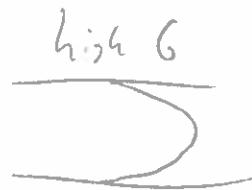
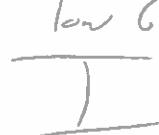
chemical fibers have reaction \rightarrow

$$\Theta_t = 1 - e^{-lu} \Theta$$

low G

high G

Alternate PDE & ODE solvers \rightarrow



Stable $|$ $x \in \text{expect}$ \leftarrow steady state $Eu_{yy} - e^{-lu} u - G = 0$

run time backward to get unstable branch
 \hookrightarrow doesn't work on saddle!

Instead, do root-find \rightarrow guess midpoint b/w stable states & other does to find unstable root! Find one, perturb G & use root as guess
 \hookrightarrow "continuation"

Convergence of Newton's Method Pf:

We have that $|f'(x)| \geq \rho > 0$ & $|f''(x)| < \delta$ or f'' Lipschitz continuous & const.

Statement: If we start close enough, quadratic convergence.

Let $\gamma \in (0, 1)$, $\tilde{\gamma}$ = radius of largest open interval around root x^* in D

Show by induction that: $|x_{n+1} - x^*| \leq \gamma |x_n - x^*| < \frac{\gamma}{\tilde{\gamma}}$ (convergence)

Let $\beta = \min\{\tilde{\gamma}, \gamma(2\rho)\}$.

Suppose $|x_0 - x^*| < \beta$.

$$\begin{aligned} \text{Then } x_1 = x_0 - \frac{f(x_0)}{f'(x_0)} &\Rightarrow x_1 - x^* = x_0 - x^* - \frac{f(x_0) - f(x^*)}{f'(x_0)} \\ &= \frac{|f'(x_0)(x_0 - x^*) - f(x_0) + f(x^*)|}{|f'(x_0)|} \end{aligned}$$

$$\text{Last time, } |f(x^*) - L_0(x^*)| \leq \frac{\delta}{2} (x_0 - x^*)^2$$

$$\Rightarrow |x_1 - x^*| \leq \frac{1}{|f'(x_0)|} \frac{\delta}{2} (x_0 - x^*)^2 \leq \frac{\delta}{2\rho} (x_0 - x^*)^2$$

$$\text{Have } |x_0 - x^*| < \beta \leq \gamma(2\rho) \Rightarrow \frac{\delta}{2\rho} |x_0 - x^*| < \gamma \quad \text{↗}$$

$$\Rightarrow |x_1 - x^*| < \gamma |x_0 - x^*|. \quad \# \text{base case}$$

Then assume for x_n , show x_{n+1} (same thing).

Hence convergence. □

If γ is small ($|f''| < \gamma$), γ gets big \rightarrow larger starting radius for convergence
 e.g. $\gamma = 0$, $|f''| = 0$ everywhere \Rightarrow linearization is great approx
 Newton's method on linear problem converges in 1 step
 $f(\underline{x}) = A\underline{x} - b$ indep. of starting value!

$$J = A \Rightarrow \underline{x}_* = \underline{x}_0 - A^{-1}(A\underline{x}_0 - b) = A^{-1}b \text{ (exact root)}$$

Imagine $f(\underline{x}) = A\underline{x} - b + \text{h.o.t.}$

$\Rightarrow \frac{\gamma}{\rho}$ measures linear term ratio quadratic.

This convergence thm. gives local convergence.

- Something like bisection always converges if you have an interval where f changes sign & is cont's.
- Would like a similar statement for Newton's method (cont'd)

Ex: Root of $\tan^{-1}(x) = 0$.

Newton's method: $x_{n+1} = x_n - (1+x_n^2)\tan^{-1}(x_n)$

K	x_k
0	1
1	-5.7×10^{-1}
2	1.2×10^{-1}
3	-1.1×10^{-3}
4	8×10^{-10}
5	$< \epsilon_{\text{machine}}$

K	x_k
0	2
1	-3.5
2	14
3	-280

What is close enough for an initial guess?

way overshoot!



Should try half-steps or line-search type alg to avoid overshooting

For $\arctan x \exists \bar{x} \approx 1.39\ldots$

s.t. Newton's method yields a period 2 cycle

$$x_k = \bar{x}, -\bar{x}, \bar{x}, \dots$$

Thus Newton's method will converge

for $\arctan x$ iff $x_0 \in (-\bar{x}, \bar{x})$

This is not great! Want an algorithm that doesn't rely on x_0 for convergence!

Let's make Newton's method globally convergent w/ a hybrid method to combine something like bisection method w/ Newton when close.

Simple idea for scalar problems: Require $|f(x_{n+1})| < |f(x_n)|$

$$\text{First: } x_+ = x_n - \frac{f(x_n)}{f'(x_n)}$$

while ($|f(x_+)| \geq |f(x_n)|$)

$$\text{end} \quad x_+ = \frac{x_n + x_+}{2}$$

$$\text{Next: } x_{n+1} = x_+$$

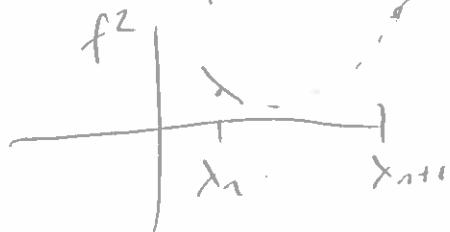
Does this work? Yes, but why?

$$\frac{d}{dx} f^2(x) = 2f(x)f'(x)$$

$$x_n - \frac{f(x_n)}{f'(x_n)} = x_n - \frac{2f(x_n)f'(x_n)}{2(f'(x_n))^2} = x_n - \beta \left(\frac{\frac{d}{dx} f^2}{\text{at } x_n} \right)$$

some direction
↓
co

Know that for β small enough, f^2 goes down in this iteration



Extends to higher dims:

$$\|\underline{F}\|^2 = \underline{F}^T \underline{F} = \sum_k F_k^2$$

$$(\nabla \|\underline{F}\|^2)_i = \frac{\partial}{\partial x_i} \sum_k F_k^2 = \sum_k 2F_k \cdot \frac{\partial F_k}{\partial x_i} = (2 \underline{J}^T \underline{F})_i$$

Newton step: $\underline{x}_{n+1} = \underline{x}_n - \underline{J}^{-1}(\underline{x}_n) \underline{F}(\underline{x}_n)$
 $= \underline{x}_n + \underline{p}_n$

Compute $\underline{p}_n \cdot \nabla \|\underline{F}\|^2 = 2(\underline{J}^T \underline{F})^T \underline{p} = 2 \underline{F}^T \underline{J} (-\underline{J}^{-1} \underline{F})$
Search direction uphill. $= -2 \underline{F}^T \underline{F} = -2 \|\underline{F}\|^2 < 0$

$\Rightarrow \underline{p}_n$ points in direction of decreasing $\|\underline{F}\|$!

Line Search/Backtracking algorithm:

Idea: Suppose $\underline{x}_+ = \underline{x}_n + \underline{p}_n$ w/ $\underline{p}_n = -\underline{J}^{-1}(\underline{x}_n) \underline{F}(\underline{x}_n)$
& supp. $\|\underline{F}(\underline{x}_+)\| \geq \|\underline{F}(\underline{x}_n)\|$

Search along line $\underline{x}_+(\lambda) = \underline{x}_n + \lambda \underline{p}_n \in \text{diag}(0,1)$

so that $\|\underline{F}(\underline{x}_+(\lambda))\| < \|\underline{F}(\underline{x}_n)\|$

Pseudocode:

$$\underline{p} = -\underline{J}^{-1}(\underline{x}_n) \underline{F}(\underline{x}_n)$$

$$\lambda = 1$$

$$\underline{x}_+ = \underline{x}_n + \lambda \underline{p}$$

while $(\|\underline{F}(\underline{x}_+)\| \geq \|\underline{F}(\underline{x}_n)\|)$ \leftarrow

$$\lambda = \lambda/2$$

$$\underline{x}_+ = \underline{x}_n + \lambda \underline{p}$$

$$\underline{x}_{n+1} = \underline{x}_+$$

make fancy by changing this to
 $\|\underline{F}(\underline{x}_n + \lambda \underline{p})\| \geq (1-\alpha\lambda) \|\underline{F}(\underline{x}_n)\|$
 $\alpha \in (0,1)$

e.g. $\alpha = 10^{-7}$ in practice

guarantees $\|\underline{F}\| \rightarrow 0$
not just decreasing

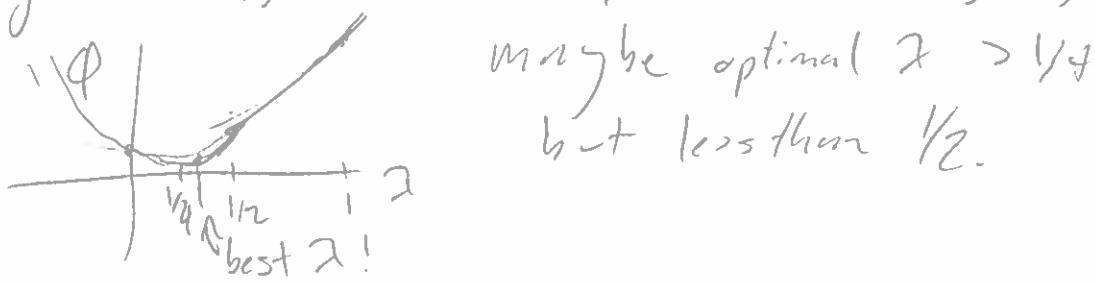
Lec 14 cont'd

In line search alg, use info from failed attempts!

Define $\phi(\lambda) = \|F(x_n + \lambda p)\|^2$

Sup. $\phi(0) < \phi(1)$

Try $\lambda = 1/2$, get $\phi(1/2) > \phi(0)$ (fails again)



Lec 15 - 10/30

(closing thoughts on Newton's method

- Terminating the iteration - want $F(x) = 0$
generically, stop when $\|F\|$ is small - absolutely/relatively?
(relative to what though? $\|F(x^*)\| = 0$ & can't divide!) $\rightarrow \|F(x_k)\| < \gamma_a$
- $\|F(x_k)\| < \gamma_r \|F(x_0)\|$. \leftarrow Use both in some combo
 - Sensitive to initial guess
 - bad guess \Rightarrow early stop
 - good guess \Rightarrow might be too hard a bound
- Also can stop if
 $\|x_{n+1} - x_n\| < \gamma \times \underbrace{\|x_n\|}_{\text{scale of problem, if dont know say } \|x_n\|}$

What if you don't have the Jacobian in Newton's method?

↳ can use finite differences to approx. Jacobian.

$$J_{ij} = \frac{F_i(x + h e_j) - F_i(x)}{h} \quad \rightarrow \text{Typically, } h = \sqrt{\epsilon_{\text{mach}} \cdot \max(|x_j|, 1)} \cdot \text{sign}(x_j)$$

Alternative: Broyden's method - each iteration update approx. Jacobian w/ current information

Jacobian-free Newton-Krylov.

- Approximately solve

$$J\delta x = -F$$

using Krylov method - e.g. gmres, conj. gradient (CG).

which requires only a routine to compute JW .

↳ why solve $J\delta x = -F$ tightly? Why not half-ass solve it?
Or tenth-ass solve it?

- Too tight a tolerance - spend a lot on linear iterations (inner) - small η

- Too loose a tolerance - spend a lot on nonlinear (outer) iterations - large (Newton)

* Solve $J\delta x = -F$ iteratively

$$\text{until } \|J\delta x + F\| \leq \gamma \|F\| \quad \rightarrow \gamma = 10^{-1} \text{ or } 10^{-2}?$$

* Good performance requires pre-conditioning!!! ~~Newton~~

$$\hookrightarrow BJ\delta x = -BF \text{ where } B \text{ invertible}$$

Ideally, $B \approx J^{-1} \Rightarrow$ system easily solved to start

Lec 15 - MAT226A - 10/30

Least-Squares - overdetermined linear systems

$$Ax = b \quad \text{w/ } A \in \mathbb{R}^{m \times n}, m > n$$

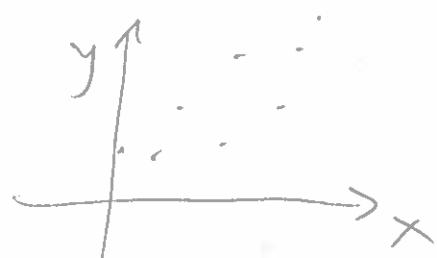
more equations than
unknowns!

$$x \in \mathbb{R}^{n \times 1}$$

$$b \in \mathbb{R}^{m \times 1}$$

overdetermined - generally, no soln.

Ex: Fit data w/ a polynomial



data $(x_i, y_i) \quad i=1, \dots, m$

Fit w/ a poly. of deg. $n-1$

$$p(x) = c_0 + c_1 x + c_2 x^2 + \dots + c_{n-1} x^{n-1}$$

would like:
 $p(x_1) = y_1$
 $p(x_2) = y_2$

$$\left. \begin{array}{l} p(x_1) = y_1 \\ p(x_2) = y_2 \\ \vdots \\ p(x_m) = y_m \end{array} \right\} \rightarrow \left. \begin{array}{l} c_0 + c_1 x_1 + c_2 x_1^2 + \dots + c_{n-1} x_1^{n-1} = y_1 \\ \vdots \\ c_0 + c_1 x_m + c_2 x_m^2 + \dots + c_{n-1} x_m^{n-1} = y_m \end{array} \right\}$$

$$\Rightarrow \left[\begin{array}{cccc|c} 1 & x_1 & x_1^2 & \dots & x_1^{n-1} \\ 1 & x_2 & x_2^2 & \dots & x_2^{n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_m & x_m^2 & \dots & x_m^{n-1} \end{array} \right] \left[\begin{array}{c} c_0 \\ \vdots \\ c_{n-1} \end{array} \right] = \left[\begin{array}{c} y_1 \\ \vdots \\ y_m \end{array} \right]$$

When $m > n$,
expect no solution!

$$Ax = b \quad \rightarrow r = b - Ax \quad \text{want residual as small as possible.}$$

$$\left(\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right)^{(n \times 1)} = \left(\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right)^{(m \times 1)}$$

Least squares soln

Find $x \in \mathbb{R}^n$ to minimize

$$\|r\|_2^2 = \|b - Ax\|_2^2$$

How to solve? Tg calculus:

$$\Sigma = \|b - Ax\|_2^2$$

Find x to minimize Σ .

$$\begin{aligned}\Sigma &= (b - Ax)^T(b - Ax) \\ &= b^T b - x^T A^T b - b^T A x + x^T A^T A x\end{aligned}$$

$$\begin{aligned}\frac{\partial \Sigma}{\partial x_j} &= -\frac{\partial}{\partial x_j} \sum_{i=1}^n x_i (a_i^T) b_i - \frac{\partial}{\partial x_j} \sum_{i=1}^n b_i a_i^T x_i \\ &\quad + \frac{\partial}{\partial x_j} \sum_{i=1}^n x_i^2 (a_i^T) a_i \\ &= -(a^T)_j b_j - b_j a_j + 2 x_j (a^T)_j a_j\end{aligned}$$

$$\nabla \Sigma = -2A^T b + 2A^T A x$$

Lec 16

$$\left\{ \begin{array}{l} \nabla \Sigma = 0 \\ A^T A x = A^T b \end{array} \right.$$

normal eqns \rightarrow if normal eqns have soln x , does it yield a min of Σ ?

$$\nabla \nabla \Sigma = 2A^T A$$

\times min Σ if $A^T A$ is pos-def.

$$\Rightarrow y^T A^T A y > 0 \text{ for } y \neq 0$$

$$\|A y\|_2^2 > 0 \text{ if } \text{null}(A) = \emptyset$$

and A full-rank

Lec 16 - MAT226A - 1/1

$\begin{pmatrix} A \\ \vdots \\ \vdots \end{pmatrix}(y) = \text{lin comb. of cols}$
 ↳ if cols of A are lin ind.
 $Ay \neq 0$ for $y \neq 0$
 $\Rightarrow A$ full rank $\Rightarrow A^T A$ is pos-def & hermitian

Suppose A full-rank, so $(A^T A)^{-1}$ exists

$$\text{hence } x = (A^T A)^{-1} A^T b$$

is the soln to the normal eqns: $Ds=0$
 $A^T A x = A^T b$

hence x is the unique minimizer of E .

$(A^T A)^{-1} A^T = A^+$ is the Moore-Penrose pseudoinverse of A .

$$\hookrightarrow \text{note } A^+ A = I_{n \times n} \quad (\text{since } \underset{n \times n}{A^T A} = \underset{n \times n}{I_{n \times n}}) \\ \Rightarrow \underset{n \times n}{(A^T A)^{-1}} \underset{n \times m}{A^+} = \underset{n \times n}{I_{n \times n}}$$

Solving normal eqns \rightarrow Assume A full rank

$$A^+ A x = A^+ b$$

pos-def. & hermitian \rightarrow use cholesky factorization to solve for $(A^T A)^{-1}$.

) fastest algorithm for least-squares.

Cost of solving normal eqns $\sim \frac{1}{3}n^3$ to solve $x = (A^T A)^{-1} A^+ b$

Cost of $A^+ b$ $\sim \underset{n \times n \times m}{2mn}$ number of elements of product
cost per element of product

Cost of $A^T A$ $\sim \underset{n \times n}{2mn^2}$ n of elements $\sim \underset{\text{cost of element}}{2mn^2}$ exploits symmetry \hookrightarrow most expensive step

Backslash in MATLAB for rectangular $A = m \times n, m \geq n$
 $x = A \backslash b$ gives Least-squares soln
 (but uses a different algorithm)

If the condition # of A is large, introduce error $A^* b$
 & condition # of $A^* A$ is worse than A roughly K_A^2 .

Other methods (of increasing expense & robustness) for LS:

"standard" method for Least-Squares (used in MATLAB)
 $x = A \backslash b$

↳ QR factorization

$$A = \underbrace{Q}_{\text{orthogonal}} \underbrace{R}_{\text{upper-triangular}}$$

(Q is invert. $K(Q) = 1 \rightarrow$ doesn't introduce error!)

Most robust way \rightarrow for singular problems / nearly-singular
 (very poorly conditioned A)

↳ SVD : $A = \underbrace{U}_{\text{orthogonal}} \underbrace{\Sigma}_{\text{diagonal}} \underbrace{V^*}_{\text{orthogonal}}$

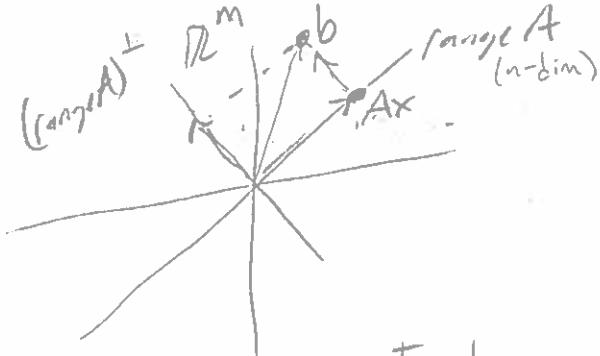
} most expensive

For square A , work for QR $\sim \frac{4}{3}n^3$.

SVD involves using QR repeatedly, so more expensive

Lec 16 cont'd - 11/1

Geometric view of Least Squares



A is $m \times n$, $m > n$

range of A is an n -dimensional subspace of \mathbb{R}^m

Generically, $b \in \mathbb{R}^m$ but $b \notin \text{range}(A)$.

Find x s.t. Ax is as close as possible to b .

looks like Ax should be the orthogonal projection of b onto $\text{range}(A)$.

Decompose $\underline{b} = \underline{b}_0 + \underline{b}_r$ s.t. $\langle \underline{b}_0, \underline{b}_r \rangle = 0$

$$(\text{range}(A))^\perp \cap \text{range}(A)$$

Then $\| \underline{b} - Ax \|_2^2 = \| \underline{b}_0 + \underline{b}_r - Ax \|_2^2 = \| \underline{b}_0 \|_2^2 + \| \underline{b}_r - Ax \|_2^2$

$\perp \text{range}(A) \in \text{range}(A)$ can make 0

for some x

Claim: For x the soln to the LS problem, the residual is orthogonal to the range of A ($r = b_0$).

$$r = b - Ax$$

$$\langle r, Ax \rangle = (Ax)^* (b - Ax) = x^* A^* (b - Ax) = x^* (A^* b - A^* Ax)$$

For x illie soln to the normal eqns, $x = (A^* A)^{-1} A^* b$

$$\Rightarrow \langle r, Ax \rangle = 0$$

□

Given \underline{b} , how to find its orthogonal projection onto $\text{range}(A)$?
You may be familiar w/ orthogonal projections onto single vectors

Lec 17 - MAT226A - 11/3

Given \underline{b} , how to find orthogonal projection onto range A ?

Recall: Given \underline{b} , projection onto vector \underline{q} :

$$pb = \underbrace{\frac{\underline{q}^* \underline{b}}{\underline{q}^* \underline{q}}}_{\text{Scalar}} \underline{q} = \underbrace{\frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}}}_{\text{distr. matrix}} \underline{b}$$

Claim: $P\underline{b} = \frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}} \underline{b}$, $\underline{b} - P\underline{b} \perp \underline{q}$.

$$\underline{q}^* (\underline{b} - P\underline{b}) = \underline{q}^* \left(\underline{b} - \frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}} \underline{b} \right) = \underline{q}^* \underline{b} - \left(\frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}} \right) \underline{q}^* \underline{b} = 0$$

Def: A square matrix P is a projector if $P^2 = P$.

A projector is orthogonal if $P^* (I - P) = 0$,

equivalently, if $P^* = P$.

Ex: $P = \frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}} \Rightarrow P^2 = \left(\frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}} \right)^2 = \frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}} \frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}} = \frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}} = P$.

$$P^* = \left(\frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}} \right)^* = \frac{\underline{q} \underline{q}^*}{\underline{q}^* \underline{q}} = P \Rightarrow \text{our example indeed was an orthogonal projector.}$$

Lec 17 Cont'd - 11/3

Project into n -dim'l space spanned by the orthonormal basis $\{q_1, \dots, q_n\}$

$$\hookrightarrow \text{recall } \langle z_i, z_j \rangle = \delta_{ij}$$

$$Pb = q_1(q_1^* b) + q_2(q_2^* b) + \dots + q_n(q_n^* b)$$

$$= \begin{pmatrix} q_1 & q_2 & \dots & q_n \end{pmatrix} \begin{pmatrix} q_1^* b \\ q_2^* b \\ \vdots \\ q_n^* b \end{pmatrix}$$

$$= \begin{pmatrix} q_1^T & q_2^T & \dots & q_n^T \end{pmatrix} \begin{pmatrix} 1 & q_1^* & 0 & \dots & 0 \\ 0 & 1 & q_2^* & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix} \underline{b}$$

$$= Q Q^* b$$

$$\Rightarrow P = Q Q^* \Rightarrow P^2 = \underbrace{Q Q^* Q Q^*}_{I \text{ by orthonormality of } \{q_1, \dots, q_n\}} = Q Q^* = P.$$

$$P^* = (Q Q^*)^* = Q (Q^*)^* = Q$$

P is an orthogonal projector.

If we get an orthonormal basis for range A , how does it help solve the LS problem?

This generates a reduced QR factorization

$$A = \underbrace{\tilde{Q} \tilde{R}}_{m \times n \text{ - orthonormal columns}} \text{ upper-triangular } n \times n$$

$m \times n$ $m \times n$ - orthonormal columns

$$Ax = b \Rightarrow \tilde{Q} \tilde{R} x = b \Rightarrow \underbrace{\tilde{Q}^* \tilde{Q} \tilde{R} x}_{I} = \tilde{Q}^* b$$

$\Rightarrow \tilde{R} x = \tilde{Q}^* b$ \leftarrow linear system on range A solve

Lec 18 - MAT226A - 11/6

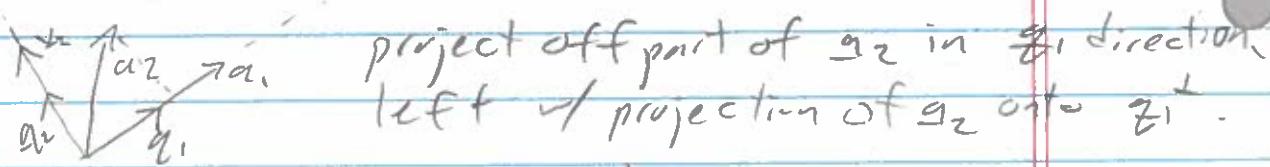
QR factorization - goal is to find orthonormal basis for range A

Given matrix A $m \times n$, $n \geq m$, find an orthonormal basis for range A.

' range A is spanned by the cols of A: $\{\underline{a}_1, \dots, \underline{a}_n\} \xrightarrow{\text{span}} \text{range A}$
 ↳ want to orthogonalize this basis

Gram-Schmidt Orthogonalization: Work sequentially.

Given $\{\underline{a}_1, \underline{a}_2\}$. Let $q_1 = \frac{\underline{a}_1}{\|\underline{a}_1\|_2} \leftarrow$ spans same space as $\{\underline{a}_1\}$


 project off part of \underline{a}_2 in \underline{q}_1 direction.
 left of projection of \underline{a}_2 onto \underline{q}_1^\perp .

$$v_2 = \underline{a}_2 - q_1(q_1^* \underline{a}_2)$$

$$q_2 = \frac{v_2}{\|v_2\|_2}$$

$$\begin{pmatrix} \underline{a}_1 & \underline{a}_2 \\ \underline{I} & \underline{I} \end{pmatrix} = \begin{pmatrix} \underline{q}_1 & \underline{q}_2 \\ \underline{I} & \underline{I} \end{pmatrix} \begin{pmatrix} r_{11} & r_{12} \\ 0 & r_{22} \end{pmatrix} \quad \text{w/} \quad \begin{aligned} r_{11} &= \|\underline{a}_1\|_2 \\ r_{12} &= q_1^* \underline{a}_2 \\ r_{22} &= \|\underline{a}_2 - r_{12} q_1\|_2 \end{aligned}$$

Easy to generalize - for n^{th} basis vector, subtract off $n-1$ projections

Classical Gram-Schmidt (unstable algorithm)

for $j=1:n$

$$v_j = \underline{a}_j$$

for $i=1:j-1$

$$\uparrow v_j = v_j - r_{ij} q_i$$

end loop

$$r_{ij} = \|v_j\|_2$$

$$q_i = v_j / r_{ij}$$

Lec 18 cont'd - 11/6

What did classical Gram-Schmidt produce?

↳ Gives the reduced QR factorization

$$(A) = (Q)(\tilde{R})(U)$$

$m \times n$ $m \times n$ $n \times n$
 orthonormal
cols. upper-
triangular

Could adapt this to generate full QR-factorization

$$(A) = \begin{pmatrix} Q \\ \bar{Q} \end{pmatrix} \begin{pmatrix} \text{\diagdown\diagup} \\ \bar{0} \end{pmatrix} \quad \begin{matrix} n \text{ rows} \\ = R \end{matrix}$$

$n \times n$ $n \times n$ $m \times n$

orthogonal
cols

extra vectors from last $m-n$ cols of \tilde{Q} are a basis for $(\text{range } A)^\perp = \text{null}(A^*)$.

Modified Gram-Schmidt (stable) - project off q_i

$$V = A$$

for i = 1:n

$$r_{ii} = \|\mathbf{v}_i\|_2$$

$$q_i = \Sigma_i / r_{ii}$$

-for $j = i+1 : n$

$$c_{ij} = q_i^* v_j$$

$$v_j = v_j - r_{ij} q_i$$

- end; loop

-end i loop

Modified GS

1st step:

$$\begin{pmatrix} I & & T \\ V_1 & V_2 & \dots & V_n \\ I & I & \dots & I \end{pmatrix} \begin{pmatrix} I & -\frac{r_{12}}{r_{11}} & -\frac{r_{1n}}{r_{11}} \\ 1 & 1 & \dots \\ \vdots & \vdots & \ddots \end{pmatrix} = \begin{pmatrix} I & T^{(2)} \\ q_1 & V_2^{(2)} \\ I & I \\ \vdots & \vdots \end{pmatrix}$$

Next step:

$$R_2 = \begin{pmatrix} 1 & 0 & \dots & -\frac{r_{2n}}{r_{22}} \\ & I & & \dots \\ & & I & \\ & & & \ddots \end{pmatrix}$$

$$\Rightarrow A \underbrace{R_1 R_2 \dots R_n}_{R^{-1}} = \hat{Q}$$

How expensive is GS?

Work of inner loop: $r_{ij} = q_i^* v_j$ ← inner prod of vectors length
 ↳ 4m flops inside $- \quad = 2m$ flops
 $v_j = v_j - r_{ij} q_i = 2m$ flops

$$\text{Total work: } \sum_{i=1}^n \sum_{j=i+1}^m 4m \approx 2mn^2$$

3rd method for QR: Householder method

$$Q_1 \cdots Q_n R = A$$

GS = triangular orthogonalization - use repeated triangular RS to get Q

HH = orthogonal triangularization - use repeated orth. Q's to get R

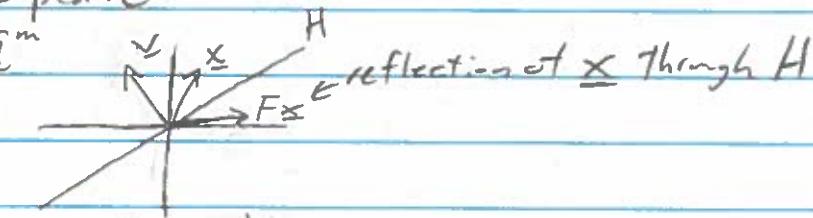
★ HH will be faster & more accurate since all Q's have $\lambda = 1$
 ↳ Superior for LS when A has poor conditioning.

Lec 19 - MAT226A - 11/8

Gram-Schmidt: $AR_1R_2 \cdots R_n = Q$
 Householder: $Q_n \cdots Q_2Q_1 A = R$

The Q_k 's involve Householder reflections through a hyperplane.
 Define hyperplane through the origin by the normal vector \underline{v} .

In \mathbb{C}^m , the space orthogonal to \underline{v} is a $m-1$ dim'l space,
 i.e. a hyperplane



$$\underline{x} = \underbrace{P_v \underline{x}}_{\perp \text{ projection onto } \text{span}\{\underline{v}\}} + \underbrace{(\mathbf{I} - P_v) \underline{x}}_{\perp \text{ projection onto } H}$$

$$\Rightarrow F\underline{x} = -P_v \underline{x} + (\mathbf{I} - P_v) \underline{x} = (\mathbf{I} - 2P_v) \underline{x}$$

$$\Rightarrow F = \mathbf{I} - 2P_v$$

Recall, P_v is an orthogonal projector, so $P_v^2 = P_v$, $P_v^* = P_v$.

Claim: F is orthogonal.

$$\begin{aligned} \text{Pf: } F^*F &= (\mathbf{I} - 2P_v)^* (\mathbf{I} - 2P_v) \\ &= (\mathbf{I} - 2\underbrace{P_v^*}_{=P_v}) (\mathbf{I} - 2P_v) \\ &= \mathbf{I} - 4P_v + 4\underbrace{P_v^2}_{=P_v} = \mathbf{I} - 4P_v \\ &= \mathbf{I} \end{aligned}$$

□

How do we use F to triangulate A ?

Goal: $\begin{pmatrix} * & * & * \\ * & * & * \\ * & * & * \\ * & * & * \end{pmatrix} \xrightarrow{Q_1} \begin{pmatrix} * & * & * \\ 0 & * & * \\ 0 & * & * \\ 0 & * & * \end{pmatrix} \xrightarrow{Q_2} \begin{pmatrix} * & * & * \\ 0 & * & * \\ 0 & 0 & * \\ 0 & 0 & * \end{pmatrix} \xrightarrow{Q_3} \begin{pmatrix} * & * & * \\ 0 & * & * \\ 0 & 0 & * \\ 0 & 0 & 0 \end{pmatrix}$

A $Q_1 A$ $Q_2 Q_1 A$ $Q_3 Q_2 Q_1 A$ $= R$

Q_2 only looks at this

In designing Q_k , only care about col i submatrix & setting subdiag entries to 0.

↳ Figure out Q_1 , then forget Q_2 consider submatrix & so on.

$$Q_k = \begin{bmatrix} I \\ F \end{bmatrix} \quad \text{where } I \text{ is } k-1 \times k-1 \text{ and } F \text{ is } m(k-1) \times m-(k-1) \text{ reflector}$$

Why? $\begin{bmatrix} I & 0 \\ 0 & F \end{bmatrix} \begin{bmatrix} T & B \\ 0 & C \end{bmatrix} = \begin{bmatrix} T & B \\ 0 & FC \end{bmatrix}$ preserves upper-triangular structure!

Goal: $Fy = F \begin{pmatrix} x \\ \vdots \\ x \end{pmatrix} = \begin{pmatrix} \pm \|x\|_2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$ since F orthogonal

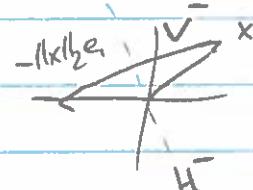
vector v connecting x & $\|x\|_2 e_1$ is \perp to H
 $\Rightarrow P_v = v v^* / v^* v$

$$\Rightarrow v = \|x\|_2 e_1 - u$$

$$F = I - \frac{v v^*}{v^* v}$$

Could use $v = z \|x\|_2 e_1 - u$ where $|z|=1$.

(choice of $z = \pm 1$ matters for stability).



Lec 19 cont'd - 11/8

In Householder reflection, reflection H^+ , H^- matters for stability!
 $\hookrightarrow v^+$ $\hookrightarrow v^-$

Suppose $x_1 > 0$:

$$v^+ = \|x\|_2 e_1 - x \rightarrow \text{if } x_1 \text{ is close to } \|x\|_2 \rightarrow \text{subtracting}$$

$$v^- = -\|x\|_2 e_1 - x \quad \begin{matrix} \text{2 nearly equal numbers!} \\ \text{= big loss of precision.} \end{matrix}$$

\downarrow
adding 2 nearly equal #s
= no big deal

So choose

$$v = \text{sign}(x_1) \|x\|_2 e_1 + x$$

$$\text{w/ } \text{sign}(x_1) = \begin{cases} +1, & x_1 \geq 0 \\ -1, & x_1 < 0 \end{cases}$$

Householder Algorithm

for $k=1:n$ (loop over cols)

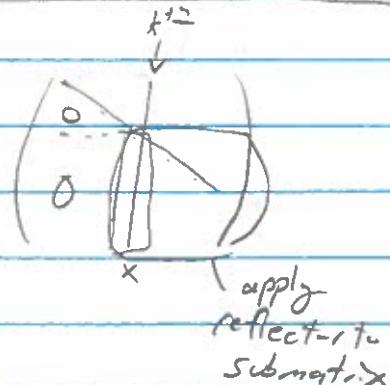
$$x = A(k:m, k)$$

$$v_k = \text{sign}(x_1) \|x\|_2 + x$$

$$v_k = v_k / \|v_k\|_2$$

$$A(k:m, k:n) = \underbrace{A(k:m, k:n)}_{IA} - 2 \underbrace{v_k (v_k^* A(k:m, k:n))}_{-2P_v A}$$

end (kloop)



Work $\sim 2mn^2 - \frac{2}{3}n^3$

If $m=n$, work $\sim \frac{4}{3}n^3$ ($2 \times$ cost of LU)

This method generates R .

$$QR \underline{x} = b$$

$$R \underline{x} = Q^* b$$

inverted this for linear solve $\dots \Rightarrow$

Solve linear system $Ax = b$ w/ QR:

$$QRx = b$$

$$Rx = Q^*b$$

only need this for linear solve

now Householder gives

$$\underbrace{Q_n \cdots Q_2 Q_1}_Q A = R$$

Add:

$\cancel{Q} \cancel{Q^*}$ we had the application of Q^* in alg.

for $k=1=n$

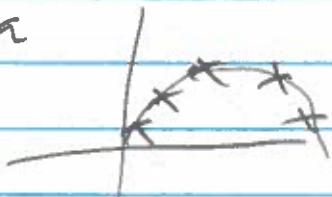
$$b(k=m) = b(k=m) - 2v_k(\cancel{v}_k^* b(k=m))$$

end (k loop)

Lec 20 - MAT226A - 11/13

Least-Squares & Orthogonal Polynomials

Intro - LS polygon fit to data



Now, given cont's fn. f on some interval,
find a polynomial to approximate f .

Let $f \in C([a, b])$

The "data" to fit is a cont's fn.

Use a degree n poly. $P_n(x) = \sum_{k=0}^n c_k x^k$

$r(x) = f(x) - P_n(x)$ \rightarrow residual is a cont's fn.

Same approach - use 2-norm to measure residual

& minimize

$$\mathcal{E} = \|r(x)\|_2^2 = \int_a^b (r(x))^2 dx$$

$$= \int_a^b (f(x) - P_n(x))^2 dx$$

$$= \int_a^b \left(f(x) - \sum_{k=0}^n c_k x^k \right)^2 dx$$

Find c_k 's where $\nabla_{c_k} \mathcal{E} = 0$

Find c_k 's where $\nabla_{c_k} \mathcal{E} = 0$, i.e.

$$\frac{\partial \mathcal{E}}{\partial c_j} = 0 \quad \text{for } j=0, \dots, n$$

$$\mathcal{E} = \int_a^b \left(f(x) - \sum_{k=0}^n c_k x^k \right)^2 dx$$

$$\frac{\partial \mathcal{E}}{\partial c_j} = \int_a^b 2 \left(f(x) - \sum_{k=0}^n c_k x^k \right) \cdot (-x^j) dx = 0$$

$$\int_a^b x^j \sum_{k=0}^n c_k x^k dx = \int_a^b f(x) x^j dx$$

$$\sum_{k=0}^n c_k \int_a^b x^j x^k dx = \int_a^b f(x) x^j dx$$

$$\sum_{k=0}^n c_k \langle x^j, x^k \rangle = \langle f(x), x^j \rangle \rightarrow \text{all of these eqns}$$

Do this for $j=0, \dots, n$ to get

$$A \underline{c} = \underline{b} \quad \rightarrow \text{normal eqns}$$

$$\text{w/ } b_j = \langle f(x), x^j \rangle$$

$$\text{& } A_{jk} = \langle x^j, x^k \rangle$$

If $a=0, b=1$, i.e. $f \in C([0,1])$

$A_{jk} = \int_0^1 x^{j+k} dx$ is the Hilbert matrix

\Rightarrow Very poorly conditioned!

n	$K(A)$
5	1.5×10^7
10	5.2×10^{14}

✗ obviously
should have
chosen
orthog

Lec 20 cont'd - 11/13

Back to discrete case $\underline{Ax} = \underline{b}$

where A is $n \times n$ w/ $n \geq n$

$$A^* A \underline{x} = A^* \underline{b}$$

What if cols of A were orthogonal?

$\hookrightarrow A^* A$ is diagonal & $\lambda(A^*) = 1$ so $A^* \underline{b}$ is nice!

$$\hookrightarrow A = Q \mathbb{I}$$

$$\Rightarrow A^* A \underline{x} = \underline{x} = A^* \underline{b} \text{ done!}$$

~~Start w/ orthonormal cols. LS soln is trivial!~~

Don't use basis $\{1, x, x^2, \dots, x^n\}$
for $T_n = \text{space of polys of deg. } \leq n$.

Use an orthogonal basis!

\hookrightarrow When we wrote out problem, we weren't given matrix A , we chose a basis & derived it.

So choose orthogonal basis & re-derive A .

$$\hookrightarrow \{\phi_0(x), \phi_1(x), \dots, \phi_n(x)\}$$

$$P_n = \sum_{k=0}^n c_k \phi_k(x)$$

$$\mathcal{E} = \int_a^b \left(f(x) - \sum_{k=0}^n c_k \phi_k(x) \right)^2 dx$$

$$\frac{\partial \mathcal{E}}{\partial c_j} = \int_a^b 2 \left(f(x) - \sum_{k=0}^n c_k \phi_k(x) \right) (-\phi_j(x)) dx = 0$$

$$\Rightarrow \sum_{k=0}^n c_k \langle \phi_j(x), \phi_k(x) \rangle = \langle f(x), \phi_j(x) \rangle$$

$$\text{A: orthogonal} \Rightarrow c_j \langle \phi_j(x), \phi_j(x) \rangle = \langle f(x), \phi_j(x) \rangle$$

Hence system of normal eqns becomes

$$A \underline{c} = \underline{b}$$

$$\text{w/ } b_j = \langle f(x), \phi_j(x) \rangle$$

$$A_{jk} = \langle \phi_j(x), \phi_k(x) \rangle \quad j, k \rightarrow, h = I.$$

$\{\phi_j\}$ orthogonal $\Rightarrow A$ diagonal.

$$\Rightarrow c_j = \langle f, \phi_j \rangle / \langle \phi_j, \phi_j \rangle$$

This is the size of projection of f onto ϕ_j !

Where do we get an orthogonal basis?

\hookrightarrow Gram-Schmidt works great analytically.

\hookrightarrow Use GS on $\{1, x, x^2, \dots, x^n\}$ to generate orthogonal basis $\{\phi_0, \phi_1, \dots, \phi_n\}$

\hookrightarrow obtain the Legendre polynomials if we chose inner product & interval $[-1, 1]$.

Legendre polys: $\{1, x, \frac{3}{2}x^2 - \frac{1}{2}, \frac{5}{2}x^3 - \frac{3}{2}x, P_4, \dots, P_n\}$

$$\text{where } P_k = \left(\frac{2k-1}{k}\right)x P_{k-1} - \left(\frac{k-1}{k}\right)P_{k-2}$$

(If we did GS on $\{1, x, x^2\}$, we'd get $P_2 = x^2 - \frac{1}{3}$ instead)

\Rightarrow monic Legendre polys (same up to constant multiple.)

Lec 21 - MAT226A - 11/15

Legendre Polynomials - orthogonal in $\langle f, g \rangle = \int_{-1}^1 f(x)g(x)dx$
 w/ P_n a n^{th} degree poly.

Say $\{ \phi_j(x) \}_j$ orthogonal on $(-1, 1)$ \rightarrow generate on (a, b) ?

map $(a, b) \rightarrow (-1, 1)$

$$s = \frac{2(x-a) - (b-a)}{b-a}$$

The polynomials $\{ \phi_j(s(x)) \}_j$ are orthogonal on (a, b) .

$$\hookrightarrow \int_a^b \phi_i\left(\frac{2(x-a)-(b-a)}{b-a}\right) \phi_j\left(\frac{2(x-a)-(b-a)}{b-a}\right) dx$$

Introduce s as above, $ds = \frac{2}{b-a} dx$

$$\Rightarrow \int_{-1}^1 \phi_i(s) \phi_j(s) \frac{b-a}{2} ds = 0 \quad i \neq j.$$

In HW3, solve discrete least-squares problem

by sampling a function & fitting poly. to data (w/ monomial basis)

\hookrightarrow get matrix $A = 50 \times 12$, w/ $\kappa(A) \approx 10^8$.

\hookrightarrow Could have used Legendre polys. as basis

but these are not orthogonal wrt. discrete inner product
 (but they are nearly orthogonal so $\kappa(A_L) \approx 5$).

\hookrightarrow If switched to 20^{th} -degree monomial basis, have more

basis vectors linearly dependent \Rightarrow worse conditioning
 $\kappa(A) \approx 10^5$. Or w/ 20^{th} -deg. Legendre polys., $\kappa(A_L) \approx 10$.

Other types of orthogonal polynomials

diff. domains

= diff. inner product
 $\langle f, g \rangle = \int_a^b f(x)g(x)w(x)dx$

Type

Basis

Inner Product

$w(x) \geq 0$

Legendre

$$P_n = \left(\frac{2k+1}{k}\right)xP_{k-1} - \left(\frac{k-1}{k}\right)P_{k-2}, \int_{-1}^1 f(x)g(x)dx \quad (w(x)=1)$$

Tchebyshov

$$T_k(x) = \cos(k\cos^{-1}(x)) \int_{-1}^1 f(x)g(x) \frac{1}{\sqrt{1-x^2}} dx$$

Trigonometric $\rightarrow \{1, \cos kx, \sin kx\}$

$$\int_{-\pi}^{\pi} f(x)g(x)dx$$

$$\rightarrow \{e^{ikx}\}$$

$$\int_{-\pi}^{\pi} f(x)\bar{g}(x)dx$$

* Tchebyshov polys: orthog. in $\int_{-1}^1 f(x)g(x) \frac{1}{\sqrt{1-x^2}} dx$

$$T_0(x) = 1, T_1(x) = x, T_2(x) = 2x^2 - 1, T_3(x) = 4x^3 - 3x$$

$T_k(x) = \cos(k\cos^{-1}(x)) \rightarrow$ closely related to Fourier series
 $\approx k^{\text{th}}$ degree poly.

* Approx of periodic fns \rightarrow f cont's w/ $f(x+p) = f(x)$

(can restrict to a single period to approx. f on $[-\frac{P}{2}, \frac{P}{2}]$)

Want $P_n(x) = \sum_{k=0}^n c_k \phi_k(x)$ \leftarrow use periodic $\phi_k(x)$'s & orthogonal

e.g. $\{1, \cos x, \sin x, \cos 2x, \sin 2x, \dots, \cos nx, \sin nx\}$ - periodic & orthogonal

use $P_n(x) = a_0 + \sum_{k=1}^n a_k \cos kx + b_k \sin kx$
 $\approx n^{\text{th}}$ degree trigonometric poly.

Approx f w/ P_n by minimizing

$$\|f(x) - P_n(x)\|_2^2 = \int_{-\pi}^{\pi} (f(x) - P_n(x))^2 dx$$

Lec 21 cont'd - 11/15

Approximating periodic f w/ orthogonal basis of trigonometric polynomials, get the

normal eqns : $a_0 = \frac{\langle f, 1 \rangle}{\langle 1, 1 \rangle} = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) dx$

$$a_k = \frac{\langle f, \cos kx \rangle}{\langle \cos kx, \cos kx \rangle} = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cos(kx) dx$$

$$b_k = \frac{\langle f, \sin kx \rangle}{\langle \sin kx, \sin kx \rangle} = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \sin(kx) dx$$

This $p_n(x)$ is a truncated Fourier series of $f(x)$.

$$\text{so } \|f(x) - p_n(x)\|_2^2 \rightarrow 0 \text{ as } n \rightarrow \infty \quad \text{for } f \in L^2.$$

(uniform convergence! \Rightarrow pointwise convergence).

Could use complex rep'n

$$\cos(kx) = \frac{e^{ikx} + e^{-ikx}}{2}, \quad \sin(kx) = \frac{e^{ikx} - e^{-ikx}}{2i}$$

$$\Rightarrow e^{ikx} = \cos(kx) + i \sin(kx)$$

to use the basis $\{e^{ikx}\}_{k=-n}^n$ for trig. polys instead!

$\{e^{ikx}\}$ are orthogonal in the inner product

$$\langle f, g \rangle = \int_{-\pi}^{\pi} f(x) \bar{g}(x) dx$$

Normal eqns for least-squares problem

$$p_n(x) = \sum_{k=-n}^n c_k e^{ikx} \quad \text{where } c_k \text{ complex \& if } p_n \text{ is real, then } \overline{c_k} = c_k.$$

$$c_k = \langle f, e^{ikx} \rangle = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-ikx} dx \leftarrow \text{Fourier transform!}$$

Lec 22 - MAT226A - IV/17

Recall: Chebyshev polys are orthogonal in

$$\langle f, g \rangle = \int_{-1}^1 \frac{f(x)g(x)}{\sqrt{1-x^2}} dx$$

$$T_k(x) = \cos(k \cos^{-1}(x)) \quad \text{e will show this is a poly!}$$

Analog to FS on non-periodic fns.

Recall: Deg. n trig approx to f which is 2π -periodic is a truncated Fourier

$$P_n(x) = \sum_{k=-n}^n c_k e^{ikx} \quad \text{where } c_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} f(x) e^{-ikx} dx$$

If $f(x)$ is real, then $\bar{c}_{-k} = c_k$.

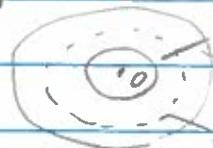
~~FS~~ Fourier series are related to Laurent series

$$f(z) = \sum_{k=-\infty}^{\infty} c_k z^k \quad \text{If have this repn., f is analytic on annulus}$$

evaluating Laurent series of f on

circle of radius 1 $\rightarrow z = e^{i\theta}$

$$\Rightarrow f(z) = f(e^{i\theta}) = \sum_{k=-\infty}^{\infty} c_k e^{ik\theta} = \text{Fourier series}$$



disc of rad. 1.

The c_k in Laurent series: $c_k = \frac{1}{2\pi i} \oint_{\gamma} \frac{f(z)}{z^{k+1}} dz$

for $z = e^{i\theta}$, this yields

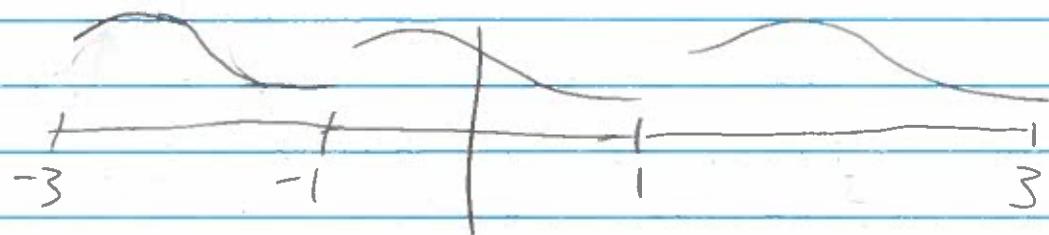
same c_k as FS.

Back to approximating f on $[-1, 1]$

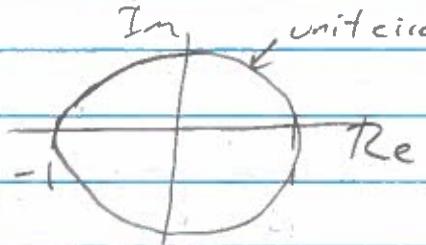


Can think of f as a periodic (but discontinuous) function

Lec 22 cont'd



FS converges in $L^2 \Rightarrow$ not pointwise, very slow
 ↳ get Gibbs phenomena at discontinuities



In unit circle in \mathbb{C} * Extend x to \mathbb{C}
 Evaluate f on the unit circle to use Fourier series.

(change of vars: $f(x) = F(z)$)

$$x = \frac{1}{2}(z + \frac{1}{z})$$

Note if $z = e^{i\theta}$, $x = \frac{1}{2}(e^{i\theta} + e^{-i\theta}) = \cos \theta$

F has the symmetry $F(z) = F(1/z) = F(\bar{z})$
 (on the unit circle)

We can compute the Fourier series of $F(z)$

$$c_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} F(e^{i\theta}) e^{-ik\theta} d\theta = \frac{1}{2\pi} \int_{-\pi}^{\pi} F(e^{i\theta}) e^{ik\theta} d\theta$$

exploiting symmetry.

$$\Rightarrow c_k = \frac{1}{2\pi} \int_{\pi}^{\pi} F(\bar{z}) \left(\frac{e^{ik\theta} + e^{-ik\theta}}{2} \right) d\theta$$

change back to x in integral

$$x = \cos \theta \Rightarrow \theta = \cos^{-1}(x) \quad d\theta = \frac{dx}{\sqrt{1-x^2}}$$

$$\frac{e^{ik\theta} + e^{-ik\theta}}{2} = \cos(k\theta) = \cos(k\cos^{-1}(x))$$

$$\Rightarrow c_k = \frac{2}{\pi} \int_{-1}^1 f(x) \frac{\cos(k\cos^{-1}(x))}{\sqrt{1-x^2}} dx, \quad c_0 = \frac{1}{\pi} \int_{-1}^1 f(x) dx$$

$$C_k = \frac{2}{\pi} \int_{-1}^1 f(x) \frac{\cos(k \cos^{-1}(x))}{\sqrt{1-x^2}} dx, \quad C_0 = \frac{1}{\pi} \int_{-1}^1 \frac{f(x)}{\sqrt{1-x^2}} dx$$

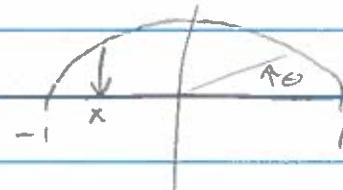
Thus the C_k 's are the orthogonal projection of f onto the Chebyshev polynomials in the appropriate inner product. $T_k(x) = \cos(k \cos^{-1}(x))$ is the k^{th} Chebyshev poly.

Thm: If f is Lipschitz conts on $[-1,1]$, it has a unique Chebyshev series representation

$$f(x) = \sum_{k=-\infty}^{\infty} C_k T_k(x) \leftarrow \text{series converges uniformly (pointwise)}$$

$$\hookrightarrow C_k = \frac{\langle f, T_k \rangle_{ch}}{\langle T_k, T_k \rangle_{ch}}, \text{ where } \langle f, g \rangle_{ch} = \int_{-1}^1 \frac{f(x)g(x)}{\sqrt{1-x^2}} dx$$

Note: $T_k(x) = \operatorname{Re}(e^{ik\theta})$
 $= \cos(k\theta)$
 $= \cos(k \cos^{-1}(x))$



Recursion Relation for Chebyshev

$$T_0(x) = \cos(0 \cdot \cos^{-1}(x)) = 1$$

$$T_1(x) = \cos(\cos^{-1}(x)) = x$$

$$\theta = \cos^{-1}(x)$$

$$T_{k+1}(x) = \cos((k+1)\theta) = \cos(k\theta)\cos(\theta) - \sin(k\theta)\sin(\theta)$$

$$T_{k-1}(x) = \cos((k-1)\theta) = \cos(k\theta)\cos(\theta) + \sin(k\theta)\sin(\theta)$$

$$\Rightarrow T_{k+1}(x) + T_{k-1}(x) = 2x \cdot T_k(x)$$

$$\boxed{T_{k+1}(x) = 2xT_k(x) - T_{k-1}(x)}$$

$$T_0 = 1, T_1 = x$$

Clearly, these are polys.

Lec 23 - MAT226A - 11/20

• Chebyshev polys: $T_n(x) = \cos(n \cos^{-1}(x))$

Recursive form: $T_{n+1}(x) = 2xT_n(x) - T_{n-1}(x)$
 $T_0(x) = 1, T_1(x) = x$

Properties of Chebyshev polys

• Roots & Extreme points:

$$T_n(x) = \cos(n \cos^{-1}(x)) = \cos(n\theta)$$

$$\cos \theta = x, x \text{ on } [-1, 1], \theta \text{ on } [0, \pi]$$

• Max/min of T_n : $n\theta = j\pi$

$$\Leftrightarrow 1, -1 \quad \theta = \frac{j\pi}{n}, j=0, \dots, n$$

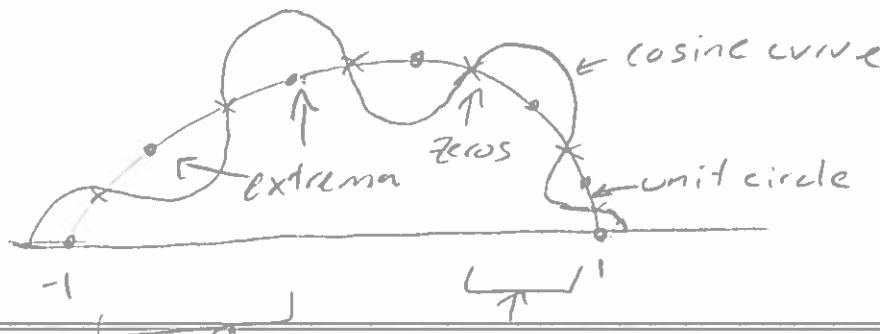
↳ n fixed, degree of poly.

• \Rightarrow Extrema $x_j = \cos\left(\frac{j\pi}{n}\right), j=0, \dots, n$
 $T'(x_j) = (-1)^j$

• Zeros of $T_n(x)$: $n\theta = (j - \frac{1}{2})\pi \rightarrow \theta \in [0, \pi] \text{ for } j=1, \dots, n$
 $\Rightarrow x_j = \cos\left((j - \frac{1}{2})\frac{\pi}{n}\right) \quad (\text{not } j=0!)$

$\Rightarrow n$ roots of $T_n(x)$ in $[-1, 1]$, $n+1$ extrema

↳ roots & extreme points useful for \Rightarrow Chebyshev interpolation points



$T_n(x)$ equioscillates on $[-1, 1]$, i.e. the oscillations vary based on the magnitude between the points.

↳ (chebyshev points clustered near the endpoints)

What's the leading coefficient of $\tilde{T}_n(x)$?

$$\tilde{T}_0(x) = 1, \quad \tilde{T}_1(x) = x$$

$$\tilde{T}_2(x) = 2x(x) - 1 = 2x^2 - 1$$

$$\tilde{T}_3(x) = 2x(2x^2 - 1) - x = 4x^3 - 3x$$

$$\tilde{T}_4(x) = 2x(4x^3 - 3x) - (2x^2 - 1) = 8x^4 - \dots$$

⇒ leading coeff. of \tilde{T}_n is 2^{n-1} .

Monic polys have leading coeff 1, so monic Chebyshev are:

$$\tilde{T}_0(x) = 1, \quad \tilde{T}_n(x) = \frac{1}{2^{n-1}} \tilde{T}_n(x) \quad n=1, 2, \dots$$

↳ have same roots & locations of extrema, but w/ max/min values of $(-1)^j / 2^{n-1}$
 \Rightarrow smaller & shrinking b/c of our rescaling.

$$\text{For } n \geq 1, \quad \frac{1}{2^{n-1}} = \max_{x \in [-1, 1]} |\tilde{T}_n(x)| \leq \max_{x \in [-1, 1]} |\hat{P}_n(x)|$$

A degree n monic polys. $\hat{P}_n(x)$.

⇒ Chebyshev polys have the smallest extreme values of all monic polys \Rightarrow smallest oscillations

Lec 23 cont'd - 11/20

Thm Let $\tilde{\Pi}_n$ be the space of monic, deg. n polys.

then $\min_{\tilde{P} \in \tilde{\Pi}_n} \left(\max_{x \in [-1, 1]} |\tilde{P}(x)| \right) = \frac{1}{2^{n-1}}$

Pf: (By contradiction). Suppose $P_n(x) \in \tilde{\Pi}_n$ and

$$\max_{x \in [-1, 1]} |P_n(x)| \leq \frac{1}{2^{n-1}} = \max_{x \in [-1, 1]} |\tilde{T}_n(x)|$$

Let $Q(x) = \tilde{T}_n(x) - P_n(x)$, note Q is at most of
deg. $n-1$
(since \tilde{T}_n, P_n monic)

Let x_j be the j^{th} extreme-point of $\tilde{T}_n(x)$ ($n+1$ such x_j)

$$Q(x_j) = \tilde{T}_n(x_j) - P_n(x_j) = \frac{(-1)^j}{2^{n-1}} - P_n(x_j)$$

Then $Q(x_j) \geq 0$ for j even

& $Q(x_j) \leq 0$ for j odd

$\Rightarrow Q$ must have roots in each interval (x_j, x_{j+1})

for $j = 0, \dots, n-1 \Rightarrow n$ roots

but $\deg(Q) \leq n-1 \Rightarrow Q = 0$!

Hence $P_n(x) = \tilde{T}_n(x)$



Polynomial Interpolation: Given 2 points $(x_0, f(x_0)), (x_1, f(x_1))$.

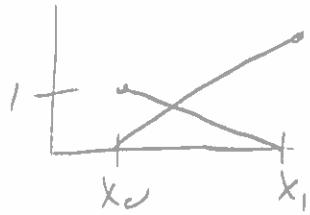
~~if $x_0 \neq x_1$, } unique linear fn.~~

$$\text{Newton form: } p(x) = f(x_0) + \frac{f(x_1) - f(x_0)}{x_1 - x_0} (x - x_0)$$

This is a rearrangement of point-slope form. Labeled Newton form
(could also write in Lagrange form)

$$\text{Lagrange Form: } p(x) = \left(\frac{x - x_1}{x_0 - x_1} \right) f(x_0) + \left(\frac{x - x_0}{x_1 - x_0} \right) f(x_1)$$

which is a "weighted average" or sum of these 2 fns:



Given x_0, \dots, x_n distinct data points & values $f(x_0), \dots, f(x_n)$
Unique poly. of deg $\leq n$, $p(x)$, s.t. $p(x_j) = f(x_j)$ $\forall j$.

How to find $p(x)$?

$$\begin{bmatrix} 1 & x_0 & x_0^2 & \cdots & x_0^n \\ 1 & x_1 & x_1^2 & \cdots & x_1^n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_n & x_n^2 & \cdots & x_n^n \end{bmatrix} \begin{bmatrix} c_0 \\ c_1 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_n \end{bmatrix} \quad \leftarrow \text{Solve for } \vec{c}$$

Dont just solve for \vec{c} generically, instead
we find the determinant $\prod_{0 \leq i < j \leq n} (x_i - x_j)$

For $n+1$ points,

$$p(x) = \sum_{k=0}^n l_k(x) f_k, \text{ where } l_k(x_j) = \delta_{jk}$$

$$\& l_k(x) = \prod_{\substack{j=0 \\ j \neq k}}^n \left[\frac{(x-x_j)}{(x_k-x_j)} \right]$$

Note $l_k(x)$ are n^{th} deg. polys.

This gives a unique poly!

Pf: Assume $\exists g$ s.t. $g(x_j) = f(x_j)$ for $j=0, \dots, n$, $\deg(g) \leq n$

Have $R(x) = p(x) - g(x)$ of $\deg(R) \leq n$

Then $R(x_j) = 0 \quad \forall j \Rightarrow n+1$ roots $\Rightarrow \deg(R) \leq n$

unless $R(x) = 0 \Rightarrow g(x) = p(x)$. \square

What about lower deg. polys for higher # of data points?

Suppose $f \in C^{n+1}[a,b]$, $x_0, x_1, \dots, x_n \in [a,b]$

For each $x \in [a,b]$, $\exists z \in (a,b)$ s.t.

$$f(x) = p(x) + \underbrace{\frac{f^{(n+1)}(z)}{(n+1)!} (x-x_0) \cdots (x-x_n)}$$

Remainder - don't know how big
so not as accurate

Lagrange form: $p(x) = \sum_{k=0}^n l_k(x) f_k$ \leftarrow numerically too much work

Newton form: $p(x) = c_0 + c_1(x-x_0) + c_2(x-x_0)(x-x_1)$
 $+ \dots + c_n \prod_{j=0}^{n-1} (x-x_j)$

Lec 24 cont'd

(3)

Fastest ways to find c_i 's in Newton Form:

Divided difference (chooses a basis that triangularizes the matrix)

$$c_0 = f[x_0]$$

$$c_1 = f[x_0, x_1] = \frac{f(x_1) - f(x_0)}{x_1 - x_0}$$

$$c_2 = f[x_0, x_1, x_2] = \frac{f[x_1, x_2] - f[x_0, x_1]}{x_2 - x_0}$$

:

:

Other polyforms:

Barycentric form: $\ell_k(x) = \prod_{\substack{j=0 \\ j \neq k}}^{n-1} \frac{(x - x_j)}{(x_k - x_j)}$

$$= \underbrace{\left(\prod_{j=0}^n (x - x_j) \right)}_{\phi(x)} \left(\frac{w_k}{x - x_k} \right)$$

where $w_k = \prod_{\substack{j=0 \\ j \neq k}}^{n-1} (x_k - x_j)$

$$\text{Then } \ell_k(x) = \phi(x) \frac{w_k}{x - x_k} \Rightarrow p(x) = \begin{cases} \phi(x) \sum_{k=0}^n \frac{w_k}{x - x_k} f_k, & x \neq x_k \\ f_k, & x = x_k \end{cases}$$

Note: w_k only depends on point distribution, can precompute!

Note: If $f_k = 1 \forall k$, then $p(x) = 1$ (do the algebra)

Lec 25 - MAT226A - 11/27

Global Polynomial Interpolation

▷ Performance of polynomial interpolation depends on point spacing

Compare interpolation to given fn. using equidistant points

$$x_k = -1 + \frac{2k}{n}, k=0, 1, \dots, n$$

or Chebyshev points \rightarrow generally better for high deg. polys.

$$x_k = \cos\left(\frac{k\pi}{n}\right), k=0, 1, \dots, n$$

Four examples on handouts:

1) $y = |x|$ C^0 , not C^1

2) $y = |x|^3$ C^2

3) $y = e^{-x}$ C^∞

4) $y = \frac{1}{25x^2+1}$ C^∞ \rightarrow has poles in complex plane, so convergence is slow for Ch. & non-existent for equidist. See Bernstein loops & contour integrals.

Thm. For smooth functions, interpolation error of n^{th} deg. approx is:

$$\underbrace{\frac{(n+1)!}{(n+1)!} \prod_{k=0}^n (x-x_k)}_{\text{depends on } h, \text{ } n+1 \text{ deg. monic polynomial}}$$

▷ If the x_k 's are the zeros of Chebyshev polys, this term is minimized

Conditioning of the interpolation problem

↳ Sensitivity of the interpolant to the data

Consider the x_k 's fixed (interpolant $f \in \mathcal{V}_n$) & wiggle the data:

$$p(x) = \sum_{k=0}^n f_k l_k(x)$$

$$\|p(x)\|_\infty = \max_{x \in [-1,1]} \left| \sum_{k=0}^n f_k l_k(x) \right|$$

$$\leq \max_{x \in [-1,1]} \sum_{k=0}^n |f_k| |l_k(x)|$$

$$\leq \|f\|_\infty \cdot \max_{x \in [-1,1]} \sum_{k=0}^n |l_k(x)|$$

depends on data depends on point distribution

$$G(\tilde{f})_k = f_k$$

Define Lebesgue constant $\Lambda = \max_{x \in [-1,1]} \sum_{k=0}^n |l_k(x)|$.

Can show $\Lambda = \sup_f \frac{\|p\|_\infty}{\|f\|_\infty}$ $\leftarrow \infty$ -norm of interpolation operator
that sends const. $f_n \rightarrow$ poly. f_n .

Λ depends on the points $x_k \rightarrow \# \& \text{distribution}$.

Can also show that Λ is the condition # of the interpolation problem

Λ_n = Lebesgue constant for n^{th} deg. poly interpolant

Equidistant points: $\Lambda_n \sim \frac{2^{n+1}}{\ln(n)}$ → awful, exponential growth

(Chebyshev points: $\Lambda_n \sim \frac{2}{\pi} \ln(n)$ → grows, but small for relevant n .

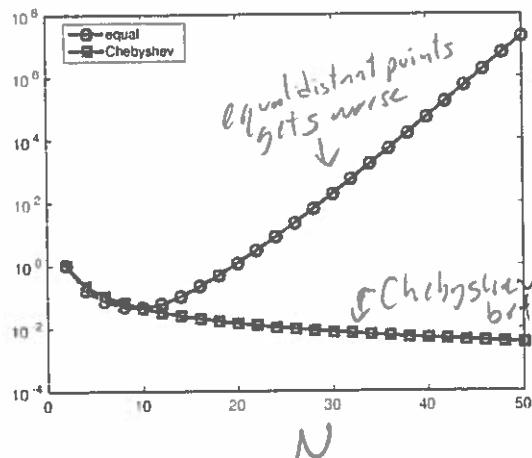
Lec25 - 11/27

$$N = \# \text{ points}$$

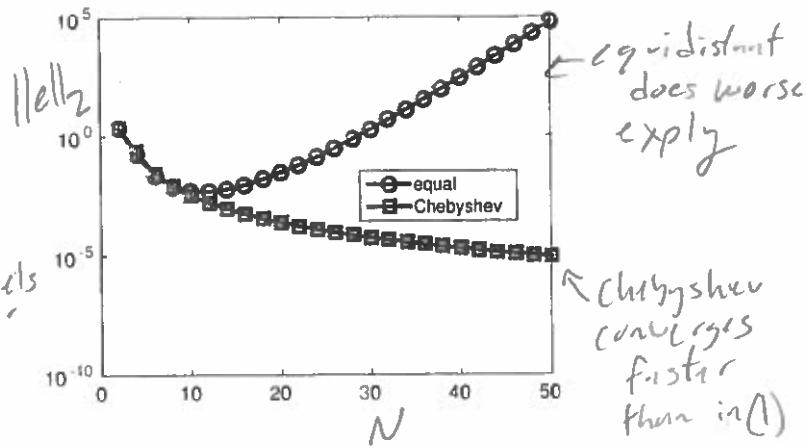
$$\|y\|_2 = \left(\sum_{n=1}^{2000} |y(x_n) - \hat{y}(x_n)|^2 \right)^{1/2} \approx \|y - \hat{y}\|_2$$

(evaluated on 2000 points numerically)

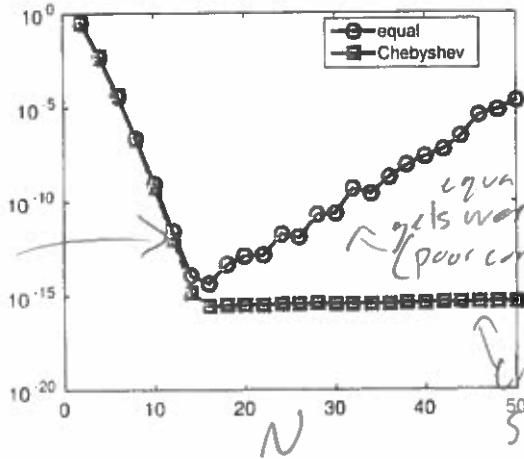
1) $y = |x|$



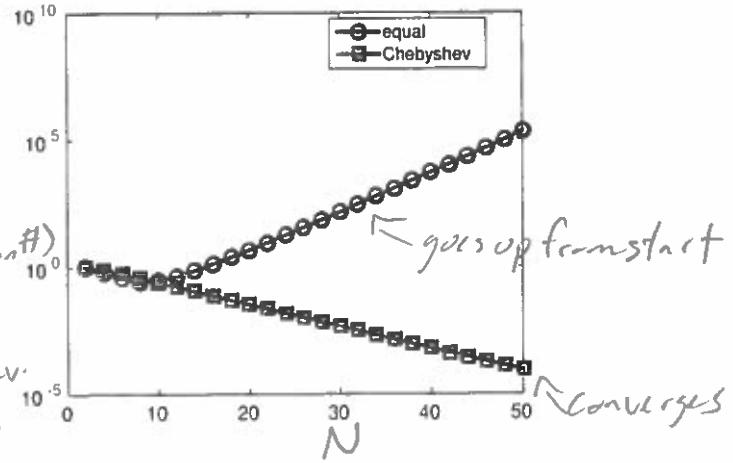
2) $y = |x|^3$



3) $y = \exp(-x)$



4) $y = \frac{1}{25x^2 + 1}$



Here, both converge to machine precision. After machine, the equispaced points interpolation begins to grow in error, due to its poor conditioning.

On the other hand, Chebyshev spaced points interpol. stays at machine, suggesting condition # is small.

Chebyshev is analog of Fourier series, so smoothness of f_T leads to speed of convergence of approx.

'Not true for equispaced internal'

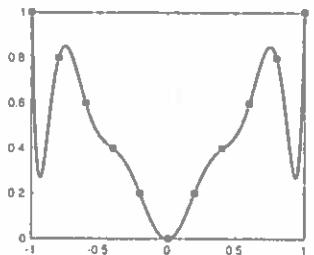
Comparing performance of poly. interpolation using different point spacings

Lec 25

11/27

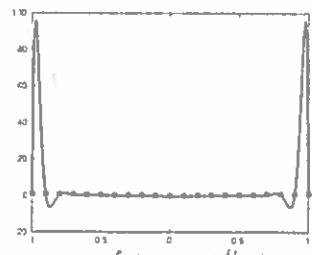
$$y = |x|$$

$N = 11$ Equal



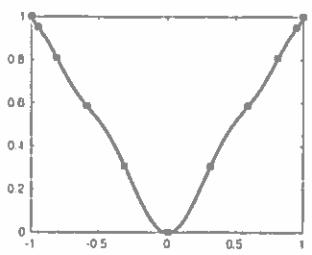
wiggles blown
points - bad!

$N = 21$ Equal



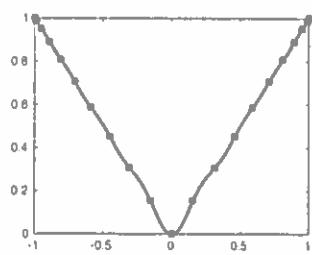
awful! oscillations
Tend to mag higher
near bdry $y = |x|^3$

$N = 11$ Chebyshev



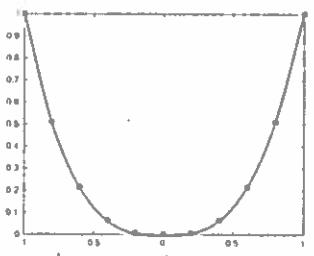
looks okay

$N = 21$ Chebyshev



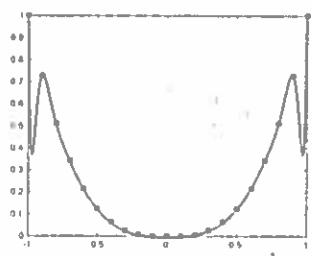
respectable

$N = 11$ Equal



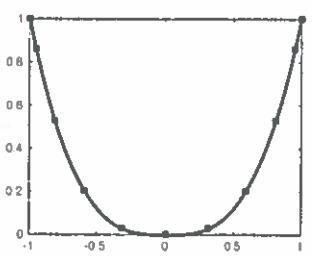
looks okay

$N = 21$ Equal



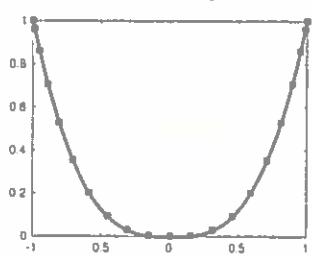
wiggles appear!
but $y = \exp(-x)$

$N = 11$ Chebyshev

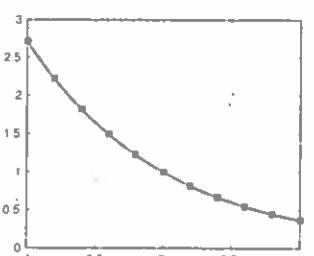


both look pretty good
2 the eye

$N = 21$ Chebyshev

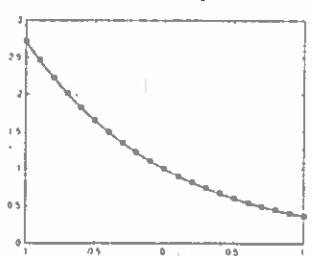


$N = 11$ Equal



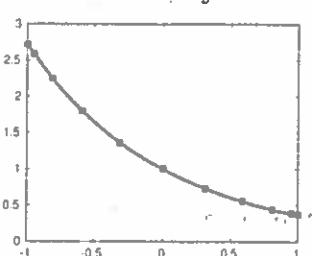
looks great!

$N = 21$ Equal



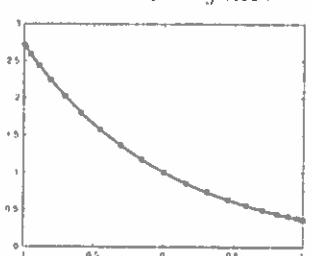
" " "

$N = 11$ Chebyshev

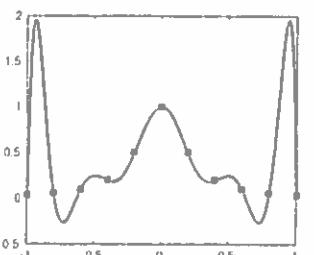


$y = \frac{1}{25x^2 + 1} \rightarrow$ has poles in complex plane

$N = 21$ Chebyshev

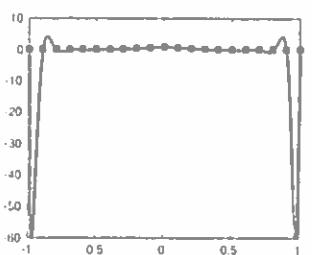


$N = 11$ Equal



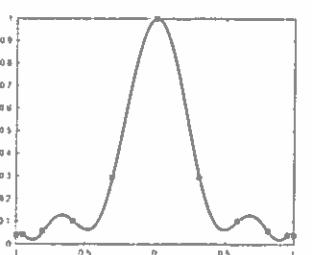
awful, wild
oscillations

$N = 21$ Equal



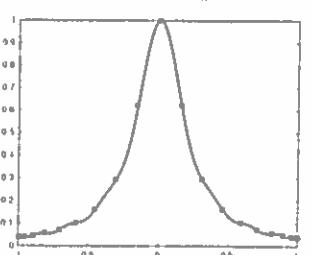
even larger
errors near
boundaries

$N = 11$ Chebyshev



Still not perfect,
wiggles exist

$N = 21$ Chebyshev



looks great!

Piecewise Polynomial Approximation

Previously, used single polynomial to approx fn/data on the entire domain

Piecewise linear interpolant

spline of order 2 → connect the dots spline

generically, given $a = x_0 < x_1 < \dots < x_n = b$ locations & data y_0, y_1, \dots, y_n

↳ underlying function $y_j = f(x_j)$ potentially.

$$S(x) = S_j(x) \text{ on } x \in [x_j, x_{j+1}]$$

$$S_j(x) = a_j + b_j(x - x_j)$$

Require S to interpolate data: $S(x_j) = y_j$
 & S is cont's

We have n intervals, $2n$ coefficients to determine S

Interpolated data → $n+1$ conditions to satisfy → $2n$ conditions

Continuity at interior nodes → $n-1$ conditions

Enough information + constraints to uniquely determine S

$$\text{Linear spline} \Rightarrow S_j(x) = y_j + \left(\frac{y_{j+1} - y_j}{x_{j+1} - x_j} \right)(x - x_j)$$

(order 2)

Example of spline: Given distinct $a = x_0 < x_1 < \dots < x_n = b$, a spline $S(x)$ of order $k \geq 1$ is

- a) poly. of deg k on each interval $[x_j, x_{j+1}]$
- b) for $k=1$, S is PW constant, for $k \geq 2$ S & its $k-2$ deris are cont's.

* PW linear spline minimizes length

History: a spline is a drafting tool - a thin, flexible wood/mica strip.

Pin some points called knots, place wood & then trace.



From elasticity, the curve will minimize strain energy

$$\Sigma = \int_a^b (K(x))^2 dx$$

curvature

$$\text{If } y' \text{ is small, } K \approx y'' \Rightarrow \Sigma = \int_a^b (K(x))^2 dx \approx \int_a^b (y''(x))^2 dx$$

Find $y(x)$ to minimize this functional

s.t. the interpolation constraints

↳ gives the natural interpolating cubic spline.

Cubic Spline Interpolants (order 4)

Given $a = x_0 < x_1 < \dots < x_n = b$ & data $y_j = f(x_j)$

$S(x)$ is a cubic poly. on each $[x_j, x_{j+1}]$

$$S_j(x) = a_j + b_j(x-x_j) + c_j(x-x_j)^2 + d_j(x-x_j)^3$$

Require S interpolates the data $S(x_j) = y_j \rightarrow n+1$ constraints

& $S \in C^2 \Rightarrow S, S', S''$ cont's at interior points $\rightarrow 3(n-1)$ constraints

$$S_j(x_{j+1}) = S_{j+1}(x_{j+1}), \quad S'_j(x_{j+1}) = S'_{j+1}(x_{j+1}), \quad S''_j(x_{j+1}) = S''_{j+1}(x_{j+1})$$

$\Rightarrow 4n-2$ constraints

Have n intervals & 4 coeffs on each \rightarrow need $4n$ constraints \rightarrow need 2 more constr.

Need boundary conditions! Natural $\rightarrow S''(a) = S''(b) = 0$ (energy min.)
clamped / ... molo etc. $S'(1), S''(1), S'(L), S''(L), \dots, S'(n), S''(n), S'(n+1), S''(n+1)$

Lec 26 cont'd 11/28

Cubic spline interpolants

$$S_j(x) = a_j + b_j(x - x_j) + c_j(x - x_j)^2 + d_j(x - x_j)^3$$

Require $S_j(x_j) = y_j$, $S_j(x_{j+1}) = S_{j+1}(x_{j+1})$, $S'_j(x_{j+1}) = f'_{j+1}(x_{j+1})$
 $\& S''_j(x_{j+1}) = f''_{j+1}(x_{j+1})$

Need 2 more constraints - choices:

Natural BC's : $S''(a) = S''(b) = 0$

Clamped/complete : $S'(a) = f'(a)$, $S'(b) = f'(b)$

Periodic BC's : $S'(a) = S'(b)$, $S''(a) = S''(b)$

Not-a-knot condns : $S'''_0(x_i) = S'''_i(x_i)$ $\rightarrow S_0 \& S_i$ same

$S'''_{n-2}(x_{n-1}) = S'''_{n-1}(x_{n-1})$ $\rightarrow S_{n-2} \& S_{n-1}$ same

Lec 27 - 12/1

Accuracy for clamped (complete) interpolating cubic splines.

Let $f \in C^4[a, b]$, let S be the unique clamped cubic spline interpolant to f through the nodes

$a = x_0 < x_1 < \dots < x_n = b$, then

$$|f(x) - S(x)| \leq \frac{5}{384} \max_{x \in [a, b]} |f^{(4)}(x)| (\max_j (x_{j+1} - x_j))^4$$

If $x_{j+1} - x_j = h \forall j$, then $\|f(x) - S(x)\|_\infty \leq Ch^4$

\hookrightarrow not as sensitive to point spacing compared w/ poly interpolation
guaranteed to converge as we add more data points & $h \rightarrow 0$.

Lec 27 contd

System of eqns. to find cubic interpolating spline

$$S_j(x) = a_j + b_j(x - x_j) + c_j(x - x_j)^2 + d_j(x - x_j)^3$$

$$S'_j(x) = b_j + 2c_j(x - x_j) + 3d_j(x - x_j)^2$$

$$S''_j(x) = 2c_j + 6d_j(x - x_j)$$

constraints: interpolation: $S_j(x_j) = y_j = a_j \quad \star \quad \text{for } j=0, \dots, n$

continuity of S : $S_j(x_{j+1}) = S_{j+1}(x_{j+1}) = a_{j+1}$

$$\Rightarrow S_j(x_{j+1}) = a_j + b_j \underbrace{(x_{j+1} - x_j)}_{h_j} + c_j h_j^2 + d_j h_j^3 = a_{j+1} \quad (1)$$

continuity of S' : $S'_j(x_{j+1}) = S'_{j+1}(x_{j+1}) = b_{j+1}$

$$\Rightarrow S'_j(x_{j+1}) = b_j + 2c_j h_j + 3d_j h_j^2 = b_{j+1} \quad (2)$$

continuity of S'' : $S''_j(x_{j+1}) = S''_{j+1}(x_{j+1}) = 2c_{j+1}$

$$\Rightarrow S''_j(x_{j+1}) = 2c_j + 6d_j h_j = 2c_{j+1} \quad (3)$$

$4n \times 4n$ system \rightarrow get down to $n \times n$

Solve (3): $d_j = \frac{c_{j+1} - c_j}{3h_j} \quad \star \rightarrow \text{plug into (1) \& (2) + eliminate } d_j$

$$(1'): a_j + b_j h_j + c_j h_j^2 + \frac{1}{3}(c_{j+1} - c_j) h_j^2 = a_{j+1}$$

$$(2'): b_j + 2c_j h_j + (c_{j+1} - c_j) h_j = b_{j+1}$$

Use (1') to solve for b_j in terms of c_j 's

Lec 27 cont'd

12/1

Solve for b_j from (1'):

$$\cancel{b_j = \frac{1}{h_j} (a_{j+1} - a_j) - \frac{h_j}{3} (2c_j + c_{j+1})}$$

Plug into (2'):

$$\text{System: } h_{j-1}c_{j-1} + 2(h_{j-1} + h_j)c_j + h_jc_{j+1} = \frac{3}{h_j}(a_{j+1} - a_j) - \frac{3}{h_{j-1}}(a_j - a_{j-1})$$

$j=1, \dots, n-1 \rightarrow$ need c_0 & c_n , define this

\hookrightarrow need BC's

Natural BC's: $s''(a) = s''(b) = 0$

$$s''(a) = s''(x_0) = c_0 \implies c_0 = 0$$

$$s''(b) = s''(x_n) = 2c_{n-1} + 6d_{n-1}h_{n-1} = 0$$

Recall: $d_j = \frac{c_{j+1} - c_j}{3h_j}$ is not valid for $j=n-1$ since no c_n

\Rightarrow use BC for $d_{n-1} = \frac{-c_{n-1}}{3h_{n-1}}$ \Rightarrow Define $c_n = 0$, these expressions sum to $j=n-1$.

\Rightarrow We can define $c_0 = 0$ & $d_j = \frac{c_{j+1} - c_j}{3h_j}, j=0, \dots, n-1$

Structure of linear system for c 's

$$A \subseteq = R$$

$$A = \begin{pmatrix} 2(h_0 + h_1) & h_1 & & & & \\ h_1 & 2(h_1 + h_2) & h_2 & & & \\ & \ddots & \ddots & \ddots & & \\ & & & & h_{n-2} & \\ & & & & & 2(h_{n-2} + h_{n-1}) \end{pmatrix} \begin{pmatrix} \frac{3}{h_1}(a_2 - a_1) - \frac{3}{h_0}(a_1 - a_0) \\ \vdots \\ \vdots \\ \vdots \\ \frac{3}{h_{n-1}}(a_n - a_{n-1}) \\ -\frac{3}{h_{n-2}}(a_{n-1} - a_{n-2}) \end{pmatrix}$$

note A is strictly diagonally dominant \Rightarrow invertible
 \rightarrow cause ... slim ... in condition

Lec 27 cont'd

Given f + data + given f'' at all nodes

↳ can use Piecewise Hermite poly. interpolant

↳ only C' , but 4th order accurate & way easier
to compute than cubic spline \rightarrow no linear system
solve.

$f(a)$ $f(b)$ $f'(a)$ $f'(b)$ \rightarrow 4 constraints can pin 4 coeffs to determine
a cubic Hermite polynomial

Cubic Hermite poly is easy to evaluate w/ this data.

Poly. int. error

Approx $f(x)$ w/ poly. interpolation

poly interpolant of deg n , $n+1$ pts

$$p_n(x_j) = f(x_j) \text{ for } j=0, \dots, n$$

Thm: $|f(x) - p_n(x)| = \frac{f^{(n+1)}(\xi_x)}{(n+1)!} \cdot \prod_{j=0}^n (x-x_j)$

For $x = x_j$, $j=0, \dots, n$, $|f(x) - p_n(x)| = 0$.

~~$$F(t) = f(t) - p_n(t) - \frac{[f(t) - p_n(t)]}{w(x)} w(t)$$~~

for $x \neq x_j$.

For $t = x_j$, $f(x_j) = p_n(x_j) = w(x_j) = 0 \Rightarrow F(x) = 0$

& for $t = x_i$, $F(x) = 0$

\Rightarrow $n+2$ zeros in $[a, b]$

By MVT, $\exists 0$ of F' between zeros of F .

~~\Rightarrow~~ F' has $n+1$ zeros

repeat until $F^{(n+1)}$ has 1 zero, ξ_x

Then $0 = F^{(n+1)}(\xi_x) = f^{(n+1)}(\xi_x) - 0 - \frac{[f(x) - p_n(x)]}{w(x)} (n+1)!$

$$\Rightarrow f(x) - p_n(x) = \frac{f^{(n+1)}(\xi_x)}{(n+1)!} \cdot w(x)$$

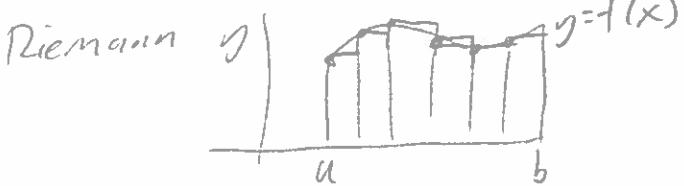
$$(A^\top A)^\top = A^\top A$$



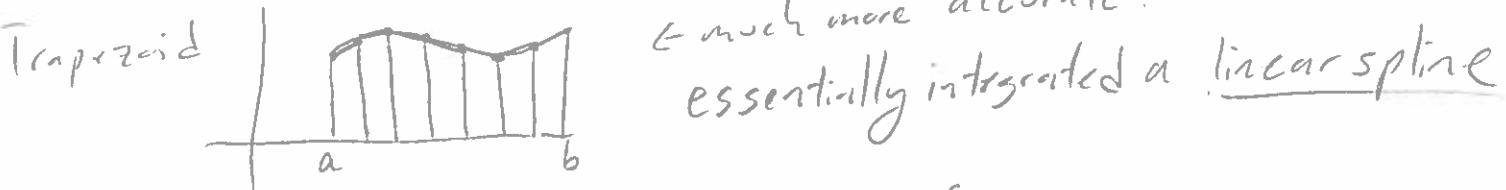
Lec 28 - MAT226A - 12/4

Numerical Integration → given function $f(x)$
 (not just data like interpolation)

$$\int_a^b f(x) dx \leftarrow \text{approximate numerically}$$



could use Riemann sum,
 or do better w/ trapezoids
 instead of rectangles



Rethink trapezoid rule as approximating f w/ a piecewise linear function which interpolates f at the endpoints + knots. Then integrating the interpolant
 \Rightarrow Allows us to extend to higher deg. approximations + estimate error

$$\int_a^b f(x) dx = \sum_k \underbrace{\int_{x_k}^{x_{k+1}} f(x) dx}_{\text{numerical quadrature}} + \underbrace{\epsilon}_{\text{error}}$$

- ↳ we will use PW polynomial, not necessarily a spline
- ↳ don't worry about matching derivatives at knots

We will cover:

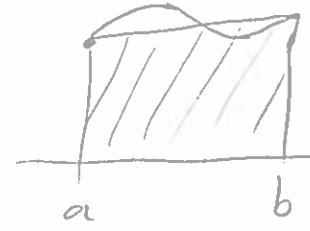
- 1) Newton-Cotes quadrature - use equidistant points.
 - ↳ e.g. trapezoidal rule, Simpson's rule, midpoint rule
 - approx fr. w/ interpolating poly on each interval
 - integrate
- 2) Gaussian quadrature - use unequally spaced points, pick quadrature points to minimize the error
 - related to orthogonal polynomials
 - use high deg. poly.
- 3) Low order Newton-Cotes & estimate errors → refine intervals w/ more points as needed
 (Adaptive quadrature)

Newton-Cotes Quadrature

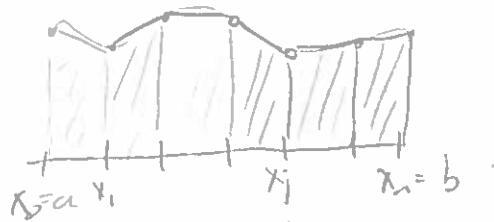
Trapezoidal Rule

$$\int_a^b f(x) dx$$

$$= \frac{(b-a)}{2} (f(a) + f(b)) + \Sigma$$



Composite trapezoidal rule: divide $[a, b]$ into equal-length subintervals & do trapezoidal rule on each subinterval



$$\rightarrow x_j = a + jh \quad \text{for } j=0, \dots, n$$

$$\text{w/ } h = \frac{b-a}{n}$$

$$\Rightarrow \int_a^b f(x) dx \approx \frac{h}{2} (f(x_0) + f(x_1)) + \frac{h}{2} (f(x_1) + f(x_2)) + \dots + \frac{h}{2} (f(x_{n-1}) + f(x_n)) \\ = h \left(\frac{f(x_0)}{2} + \sum_{j=1}^{n-1} f(x_j) + \frac{f(x_n)}{2} \right)$$

Error analysis \rightarrow look at one interval: $\int_a^b f(x) dx$

Approx f w/ linear fn. interpolating f at the endpoints

interpolate: $f(x) = \underbrace{\frac{(x-b)}{a-b} f(a) + \frac{(x-a)}{b-a} f(b)}_{\text{Lagrange interpolant}} + \underbrace{\frac{f''(3(x))}{2}(x-a)(x-b)}_{\text{error term}}$

integrate:

$$\int_a^b f(x) dx = f(x_0)w_0 + f(x_1)w_1 + \Sigma$$

$x_0=a, x_1=b$

$$w_0 = \int_a^b \frac{(x-b)}{(a-b)} dx = \frac{b-a}{2}$$

$$w_1 = \int_a^b \frac{x-a}{b-a} dx = \frac{b-a}{2}$$

$$\Sigma = \int_a^b \frac{f''(3(x))}{2}(x-a)(x-b) dx = \int_a^b \frac{f[a, b, x]}{2} (x-a)(x-b) dx$$

divided difference
easier to integrate

Lec 28 cont'd - 12/4

For higher deg. poly approx. in Newton-Cotes quadrature:

$$f(x) = P_n(x) + R(x)$$

$$\text{w/ } P_n(x) = \sum_{k=0}^n L_k(x) f(x_k) \quad \text{w/ } L_k(x) = \frac{\prod_{j \neq k} (x - x_j)}{\prod_{j \neq k} (x_k - x_j)}$$

$$\therefore R(x) = \prod_{k=0}^{n+1} \frac{f^{(n+1)}(x_k)}{(n+1)!} (x - x_k)$$

$$\Rightarrow \int_a^b f(x) dx = \sum_{k=0}^n f(x_k) w_k + \int_a^b R(x) dx$$

$$\text{w/ } w_k = \int_a^b L_k(x) dx \quad \leftarrow \text{only depends on points, not function!}$$

(closed) Newton-Cotes formulas - use n subintervals on $[a, b]$
including endpoints: $x_j = a + jh$
 $h = \frac{b-a}{n}$
to evaluate function & interpolate

$n=1 \Rightarrow$ Trapezoidal rule

$$\int_a^b f(x) dx = \frac{h}{2} (f(a) + f(b))$$

$n=2 \Rightarrow$ Simpson's rule

$$\begin{aligned} \int_a^b f(x) dx &= \frac{h}{3} (f(x_0) + 4f(x_1) + f(x_2)) \\ &= \frac{b-a}{6} (f(x_0) + 4f(x_1) + f(x_2)) \end{aligned} \quad h = \frac{b-a}{2}$$

Open Newton-Cotes : exclude endpoints from interpolant

$n=0$: Midpoint rule

$$\int_a^b f(x) dx = (b-a) f\left(\frac{a+b}{2}\right)$$



$n=1$: next one uses 2 interior points & a line through them



Standard Adaptive Simpson's Rule \rightarrow MATLAB's quad command

Accuracy

Degree of precision of a quadrature formula is the largest integer m s.t. the quadrature is exact for all polynomials of degree $\leq m$ but not for some poly. of deg. $m+1$.

- Expect trapezoidal rule to have deg. of precision 1 (it does)
- Expect Simpson's rule to have deg. of precision 2, but in fact, Simpson's rule has degree of precision 3!

Will see that composite trapezoidal rule is $\mathcal{O}(n^2)$ when $n = \# \text{ points}$
& Simpson's rule is $\mathcal{O}(n^4)$!

- Trap. rule can integrate linear polys exactly, but not quadratic
- Simpson's rule can integrate cubic polys. exactly, but not quartic

Error Analysis of Numerical Quadrature

→ Use weighted Mean Value Theorem for integrals

Thm: Suppose f is continuous on $[a, b]$ and g is integrable on $[a, b]$ & doesn't change sign

There exists $\bar{z} \in (a, b)$ s.t. $\int_a^b f(x)g(x)dx = f(\bar{z}) \int_a^b g(x)dx$

Trapezoidal Rule:

$$\int_a^b f(x)dx = \frac{b-a}{2} (f(a) + f(b)) + \int_a^b \frac{f''(\bar{z}(x))}{2} (x-a)(x-b)dx$$

error term from accuracy
of interpolating polynomial

Use weighted MVT on error term. $\exists \bar{z} \in (a, b)$ s.t.

$$\mathcal{E} = \frac{f''(\bar{z})}{2} \int_a^b (x-a)(x-b)dx = -\frac{f''(\bar{z})}{12} (b-a)^3 = -\frac{f''(\bar{z})}{12} h^3$$

Composite Trapezoidal Rule:



↓ extend to n subintervals,
 $\bar{z}_k \in [x_{k-1}, x_k]$

$$\int_a^b f(x)dx = h \left(\frac{1}{2}f(a) + \sum_{k=1}^{n-1} f(\bar{z}_k) + \frac{1}{2}f(b) \right) - \frac{h^3}{12} \sum_{k=0}^{n-1} f''(\bar{z}_k)$$

Note: $\min_{x \in [a, b]} f''(x) \leq \underbrace{\frac{1}{n} \sum_{k=0}^{n-1} f''(\bar{z}_k)}_{\text{average}} \leq \max_{x \in [a, b]} f''(x)$

Since f'' continuous, by intermediate value theorem, $\exists m \in (a, b)$ s.t.

$$\frac{1}{n} \sum_{k=0}^{n-1} f''(\bar{z}_k) = f''(m)$$

$$\mathcal{E} = -\frac{h^3}{12} \sum_{k=0}^{n-1} f''(\bar{z}_k) = \frac{(b-a)}{12} \frac{h^2}{n} \sum_{k=0}^{n-1} f''(\bar{z}_k) = \frac{b-a}{12} h^2 f'''(m)$$

$\Rightarrow \mathcal{E} = \mathcal{O}(h^2)$ so comp. trap. rule is second-order accurate.

If we double the points, get error reduction factor 4)

Simpson's Rule:

$$x_0 = a \quad x_i = \frac{a+b}{2} \quad b = x_n$$

$$\int_a^b f(x) dx = \frac{b-a}{6} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right) + \int_a^b \frac{f''(g(x))}{6} \cdot (x-a)(x-x_i)(x-b) dx$$

changes

$$(x-a)(x-x_i)(x-b) \rightarrow$$



apply MVT naively on each half-interval
get f'' with opposite signs
→ looks like 4th derivative

Idea: Integrate error term by parts to get 4th deriv & g(x) poly that doesn't change sign

$$\text{After: } \xi = -\frac{1}{90} \left(\frac{b-a}{2} \right)^5 f^{(4)}(m) = -\frac{(b-a)^5}{2880} f^{(4)}(m)$$

→ cubic polys have $f^{(4)}=0 \Rightarrow \xi=0$
extra degree of precision from symmetry

Composite Simpson's rule: same analysis as comp. trap rule.
add up n error terms to drop a power of h out of the error term

$$\xi = -\frac{(b-a)}{180} h^4 f^{(4)}(m) \quad \text{where } h \text{ is the point spacing}$$

\downarrow divide into intervals,
 \Rightarrow apply Simpson's rule
to triples of points

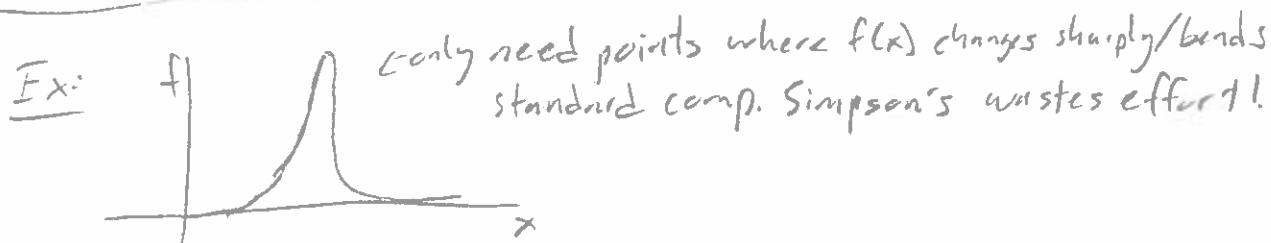
$$\xi = O(h^4)$$

\Rightarrow Comp. Simpson's rule is 4th-order accurate

So if we double the points \Rightarrow error goes down by factor of 16.

Adaptive Quadrature

Adaptive Simpson's : Idea: estimate the error using two diff. approx & use more points locally as needed



$$S(a, b) = \frac{b-a}{6} \left(f(a) + 4f\left(\frac{a+b}{2}\right) + f(b) \right), \quad \Sigma_1 = -\frac{1}{90} \left(\frac{b-a}{2} \right)^5 f^{(4)}\left(\frac{a+b}{2}\right)$$

$$\int_a^b f(x) dx = S(a, b) + \Sigma_1,$$

← apply 2 interval Simpson's rule

$$\begin{array}{ccccccc} a & & \bullet & & \bullet & & b \\ & & | & & | & & \\ & & c = \frac{a+b}{2} & & & & \end{array} \quad \int_a^b f(x) dx = S(a, c) + S(c, b) + \Sigma_2$$

$$\text{where } \Sigma_2 = -\frac{1}{90} \left(\frac{b-a}{2} \right)^5 \frac{2}{2^5} f^{(4)}\left(\frac{a+b}{2}\right)$$

$$\text{Assume } f^{(4)}\left(\frac{a+b}{2}\right) = f^{(4)}\left(\frac{a+b}{2}\right), \text{ then } \Sigma_2 = \frac{1}{16} \Sigma_1$$

$$\Rightarrow \int_a^b f(x) dx = S(a, b) + 16 \Sigma_2 \quad \text{Subtract these equations}$$

$$\int_a^b f(x) dx = S(a, c) + S(c, b) + \Sigma_2$$

$$S(a, c) + S(c, b) - S(a, b) - 15 \Sigma_2 = 0$$

$$\Rightarrow \Sigma_2 = \frac{S(a, c) + S(c, b) - S(a, b)}{15}$$

Idea:

$$\begin{array}{ccccccc} a & & \bullet & & \bullet & & b \\ & & | & & | & & \\ & & c & & & & \end{array}$$

← see if error Σ_2 is small enough

$$\begin{array}{ccccccc} a & & \bullet & & \bullet & & b \\ & & | & & | & & \\ & & c & & c' & & b \end{array}$$

← do same check by applying 2 Simpson's to each interval, refine until error in each part small enough

Adaptive Simpson's rule

$$[a \quad \bullet \quad c \quad \bullet \quad b] \rightarrow \int_a^b f(x) dx$$

do 2 Simpson's

rule approxs $Q_1 = S(a,b)$

$Q_2 = S(a,c) + S(c,b)$

$$\text{& obtained } \varepsilon_2 = \frac{Q_2 - Q_1}{15}$$

$$\text{If } |\varepsilon_2| < \epsilon, \quad Q = Q_2 + \frac{1}{15}(Q_2 - Q_1)$$

o/w, repeat on $[a,c]$ & $[c,b]$ separately

$$\int_a^b f(x) dx = \underbrace{\int_a^c f(x) dx}_{\substack{\text{apply same} \\ \text{procedure w/ tol } \epsilon/2}} + \underbrace{\int_c^b f(x) dx}_{\substack{\text{apply same} \\ \text{procedure w/ tol } \epsilon/2}}$$

$$[a \quad \bullet \quad c \quad \bullet \quad b] \rightarrow \begin{cases} \text{use 5 pts} \\ \text{for 2 Simpson rules} \end{cases}$$

$$[a \quad \bullet \quad c \quad \bullet \quad b] \rightarrow \begin{cases} \text{5 pts on each interval,} \\ \text{some function evals} \rightarrow \text{inefficiency can arise} \\ \text{if not coded well,} \\ \text{the same!} \end{cases}$$

$$[a \quad \bullet \quad c \quad \bullet \quad b]$$

Gaussian Quadrature - unlike Newton-Cotes, which uses equally spaced points, GQ uses different spacing points.

GQ uses different spacing points.

- Recall for high-order polynomial interpolation, using equidistant points led to wiggles/error in interpolant! Use others, e.g. Chebyshev points

- Choose point locations to get the greatest degree of precision

$$\text{Ex: 2 points: } \int_a^b f(x) dx \approx f(x_1)w_1 + f(x_2)w_2 \leftarrow \begin{array}{l} \text{4 degrees of freedom} \\ \rightarrow \text{can get deg. precision} \end{array}$$

$$\text{Require } \int_a^b 1 dx = w_1 + w_2, \quad \int_a^b x dx = x_1 w_1 + x_2 w_2, \quad \int_a^b x^2 dx = x_1^2 w_1 + x_2^2 w_2$$

$$\int_a^b x^3 dx = x_1^3 w_1 + x_2^3 w_2$$

Lec 30 cont'd

Solve system: $\int_a^b 1 dx = w_1 + w_2$

$$\int_a^b x dx = x_1 w_1 + x_2 w_2$$

$$\int_a^b x^2 dx = x_1^2 w_1 + x_2^2 w_2$$

$$\int_a^b x^3 dx = x_1^3 w_1 + x_2^3 w_2$$

} 2 point
Gaussian
quadrature

On $[-1, 1]$, we get $w_1 = w_2 = 1$, $x_{1,2} = \pm 1/\sqrt{3}$

Other GQ's tabulated for $[-1, 1]$ → look them up!

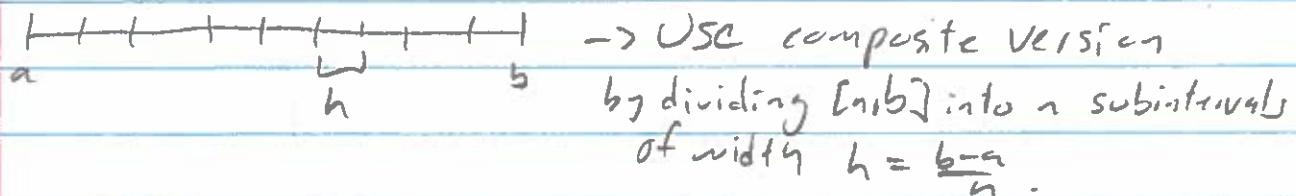
On $[a, b]$, use transformation

$$x = \left(\frac{b-a}{2}\right)z + \left(\frac{b+a}{2}\right)$$

to transform $[-1, 1]$ to $[a, b]$

$$\int_a^b f(x) dx = \frac{b-a}{2} \int_{-1}^1 f\left(\left(\frac{b-a}{2}\right)z + \left(\frac{b+a}{2}\right)\right) dz$$

!pt. GQ error: $E_2 = \frac{(b-a)^5}{4320} f^{(4)}(z)$ → looks like a factor of 2 more accurate than Simpson's (who cares?)



Obtain 2 pt. comp. GQ error

$$E = \frac{(b-a)}{4320} h^4 f^{(4)}(z) \leftarrow \text{factor 2 better than comp. Simpson's rule}$$

The diff. in accuracy b/w GQ & Newton-Cotes is more striking as we increase the # of points

→ GQ standard is ~7 pts → How do we get the points? Table, duh but where do they come from?

The points $\{x_k\}$ for n-point Gaussian quadrature $\sum_{k=1}^n f(x_k) w_k$ on $[-1, 1]$ are the n roots of the n^{th} degree Legendre polynomial.

$$\text{The weights } w_k = \int_{-1}^1 \prod_{j \neq k} \frac{(x-x_j)}{(x_k-x_j)} dx$$

The n -point Gaussian quadrature has degree of precision $2n-1$.

$$\text{On the interval } [a, b], \quad \epsilon_n = \frac{(b-a)^{2n+1} (n!)^4 f^{(2n)}(\bar{x})}{(2n+1) [(2n)!]^3} \rightarrow \text{lots of derivatives}$$

For composite n -pt GQ w/ interval length h . $\epsilon = \mathcal{O}(h^{2n})$

GQ is awesome! Needs smooth f since it requires so many derivatives
 ↴ so not perfect for non-smooth $f \rightarrow$ non-smooth, f better w/ Simpson's
 ↴ vectorizes easier than adaptive Simpson's.

Why is n -pt. GQ deg. precision n ?

pf: Let $p \in \Pi_{2n+1}$, then $p(x) = q(x) \phi_n(x) + r(x)$
 \uparrow n^{th} deg. Legendre poly

& q, r are $n-1$ deg. polys

$$\Rightarrow \int p(x) dx = \int \underbrace{q(x) \phi_n(x)}_0 + r(x) dx = \int r(x) dx$$

$$\Rightarrow \sum_{k=1}^n p(x_k) w_k = \sum_{k=1}^n \left(q(x_k) \phi_n(x_k) + r(x_k) \right) w_k = \sum_{k=1}^n r(x_k) w_k \\ = \int r(x) dx = \int p(x) dx.$$

Final Exam
Math 226A
Fall 2017

1. Suppose you run the following MATLAB commands to form three related random matrices and random vector.

```
N=4000; % size of the problem
r = rand(N,1); % random Nx1 vector
A = rand(N); % form a random NxN square matrix
B = A'*A;
C = A + A';
```

You then use the built-in linear solver to solve the three systems $Ax = r$, $Bx = r$, and $Cx = r$, and you record the time to solve these systems. The results are

```
>> tic; x=A\r; toc
Elapsed time is 1.051728 seconds.

>> tic; x=B\r; toc
Elapsed time is 0.498099 seconds.

>> tic; x=C\r; toc
Elapsed time is 0.989374 seconds.
```

e explain why $A^T A$ is posdef & symm. ✓

Comment on the results, and explain why one system is solved faster than the others.

2. The 2-point open Newton-Coats approximation to the integral

$$\int_0^1 f(x) dx \approx f(x_0)w_0 + f(x_1)w_1$$

uses the function values at $x_0 = 1/3$ and $x_1 = 2/3$.

- (a) Derive formulas for the weights w_0 and w_1 .
- (b) What is the degree of precision of the 2-point open Newton-Coats approximation to the integral?

3. A square matrix Q is orthogonal if $Q^T = Q^{-1}$. Show that the condition number (in 2-norm) of an orthogonal matrix is 1.

4. In class we examined the three different matrix factorizations: LU, Cholesky, and QR. All of these factorizations can be used to solve the linear system $A\vec{x} = \vec{b}$. For each factorization comment on when it would be advantageous to use one over the other and explain why.
5. Suppose that f is a C^3 function on \mathbb{R} , $f(x^*) = 0$, and $f'(x^*) \neq 0$.
 - Will Newton's method applied to f converge to x^* for an initial guess sufficiently close to x^* ?
 - Define $g(x) = (f(x))^2$. Then $g(x^*) = 0$. Discuss the convergence of Newton's method applied to g instead of f . Does it converge at all? If so, is the convergence faster or slower than the convergence for f ?
6. In class we observed the error in using a polynomial interpolant to approximate a function sometimes grew exponentially as the number of points increased. Is the growth of error due to stability or conditioning? Explain why exponential growth of the error is expected.
7. Suppose A is a full rank m by n matrix with $m > n$. The least squares solution to $A\vec{x} = \vec{b}$ can be found by either solving the normal equations $A^*A\vec{x} = \vec{b}$ or by using the QR factorization. Discuss the advantages and disadvantages of each approach. Specifically, comment on when you might use one method over the other.
8. The IEEE standard for quad precision floating-point numbers uses 128 bits: 1 for the sign bit (s), 15 for the exponent (e), and 112 for the mantissa (or fractional part) (q). Floating point numbers are represented as

$$x = (-1)^s(1 + q) \times 2^{e-16383}.$$
 - Recall that one definition of ϵ_{mach} is the relative difference between two floating point numbers so that for $x \in \mathbb{R}$ (in an appropriate range)

$$|x - fl(x)| \leq \epsilon_{\text{mach}}|x|,$$
 where $fl(x)$ is the floating point representation of x . What is ϵ_{mach} for quad precision?
 - Approximately how many decimal digits of accuracy are achieved by quad precision?
Note that $\log_{10}(2) \approx 0.3$. 10
33 or 34?
 - Suppose we use a backward stable algorithm applied to a problem with a condition number of $\kappa = 10^p$ to perform a calculation in quad precision. How many digits of accuracy do you expect in the computed solution?
9. Suppose you approximate $f(x) = \cosh(x)$ on $[-1, 1]$ by constructing a polynomial interpolant using the function values at the $n+1$ points $x_k = \cos(k\pi/n)$ for $k = 0, \dots, n$. Derive a bound on the maximum error in using the interpolant to approximate the function. Hint: these are the extreme points of the Chebyshev polynomials, and in homework 4, you examined a polynomial with these roots.