

Welcome to Math 221! Lecture 19, Fall 24

Symmetric Eigenproblem + SVD

$$\text{Real: } A = A^T = Q \Lambda Q^T, \quad Q Q^T = I$$

$$\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

$$\lambda_1 \geq \dots \geq \lambda_n$$

$$Q = [q_1, \dots, q_n]$$

$$\text{Def: Rayleigh Quotient } p(u, A) = \frac{u^T A u}{u^T u} \quad u \neq 0$$

$$\text{Props: } Au = \lambda u \Rightarrow p(u, A) = \lambda$$

$$u = Q Q^T v = Q b \Rightarrow p(u, A) = \frac{b^T \Lambda b}{b^T b}$$

= convex comb. of evals

Courant-Fischer Minimax Thm

R^j = j -dim subspace of \mathbb{R}^n

S^{n-j+1} = $n-j+1$ -dim subspace of \mathbb{R}^n

$$\max_{R^j} \min_{\substack{r \in R^j \\ r \neq 0}} p(r, A) = \lambda_j = \min_{S^{n-j+1}} \max_{\substack{s \in S^{n-j+1} \\ s \neq 0}} p(s, A)$$

Weyl's Thm: $A = A^T$ with $\lambda_1, \dots, \lambda_n$

and $E = E^T$

$A + E$ has evals $\mu_1 \geq \dots \geq \mu_n$ with

$$|\mu_i - \lambda_i| \leq \|E\|_2 \quad \forall i$$

Corollary for SVD: if A general,
 sing vals $\sigma_1 \geq \dots \geq \sigma_n$
 and $A+E$ has sing vals $\tau_1 \geq \dots \geq \tau_n$
 then $|\sigma_i - \tau_i| \leq \|E\|_2$

[follows from $\begin{bmatrix} 0 & A+E \\ A^T+E^T & 0 \end{bmatrix}$]

proof of Weyl:

$$\begin{aligned} \mu_j &= \min_{S^{n-j+1}} \max_{\substack{s \in S^{n-j+1} \\ s \neq 0}} \frac{s^T (A+E) s}{s^T s} \\ &= \min \max \frac{s^T A s}{s^T s} + \frac{s^T E s}{s^T s} \\ &\leq \min \max \frac{s^T A s}{s^T s} + \|E\|_2 \\ &= \lambda_j + \|E\|_2 \end{aligned}$$

for other inequality, swap λ_j and μ_j

Def: Inertia(A) =

(#neg_evals(A), #zero_evals(A), #pos_evals(A))

Sylvester's Thm: $A=A^T$, X nonsingular

$$\Rightarrow \text{Inertia}(A) = \text{Inertia}(X^T A X)$$

Fact: Suppose we factor $A = LDL^T$
(Gauss elim with symmetric or no pivoting)

$$\text{Inertia}(A) = \text{Inertia}(D) \\ = (\# D_{ii} < 0, \# D_{ii} = 0, \# D_{ii} > 0)$$

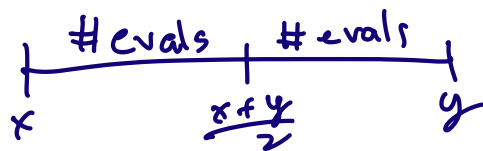
$$\text{Factor } A - xI = L' \cdot D' \cdot L'^T \\ \# D'_{ii} < 0 = \# \text{ evals of } A - xI < 0 \\ = \# \text{ evals of } A < x$$

Factor $A - yI$, $x < y$
get $\# \text{ evals of } A \leq y$

$$\# \text{ evals in } [x, y] = (\# \text{ evals } \leq y) - (\# \text{ evals } < x)$$

\Rightarrow count $\# \text{ evals}$ in any interval

\Rightarrow bisection works: $\text{Inertia}(A - \frac{x+y}{2}I)$



in practice, start by reducing $A = QTQ^T$
to tridiagonal T , then
computing $\text{inertia}(T - \alpha I)$ costs $O(n)$

Proof of Sylvester's Thm

Suppose $\# \text{ evals of } A < 0$ is m
but $\# \text{ evals of } X^T A X < 0$ is $m' < m$
seek contradiction

$N = m$ -dimensional subspace spanned by evecs of m negative evals of A

$$\Rightarrow 0 \neq x \in N \Rightarrow x^T A x < 0$$

$P = n-m'$ dimensional subspace spanned by $n-m'$ evecs of $n-m'$ nonnegative evals of $X^T A X$

$$\Rightarrow 0 \neq x \in P \Rightarrow x^T X^T A X x = \underbrace{(Xx)^T A (Xx)}_{XP} \geq 0$$

$$\dim(XP) + \dim(N)$$

$$= n-m' + m > n$$

\Rightarrow they intersect in nonzero z

$$\Rightarrow z^T A z \geq 0 \text{ and } z^T A z < 0$$

contradiction

Perturbation Theory for Evecs

$$\text{Thm: } A = Q \Lambda Q^T \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$$

$$A + E = Q' \Lambda' Q'^T \quad \Lambda' = \text{diag}(\lambda'_1, \dots, \lambda'_n)$$

$$\theta_i = \text{angle}(p_i, q_i)$$

$$\text{gap}(i, A) = \min_{j \neq i} |\lambda_i - \lambda_j|$$

$$|0.5 \sin(2\theta_i)| \leq \frac{\|E\|_2}{\text{gap}(i, A)}$$

$$\sim \theta_i$$

$$\text{if } \theta_i \ll 1$$

Why gap in denominator?

if $A=I$, $\text{gap}=0$, worst case:

upper bound = ∞

$$\begin{bmatrix} 1.1 & 0 \\ 0 & .9 \end{bmatrix} \quad \text{gap} = .2$$

$$E = \begin{bmatrix} -.1 & \\ & +.1 \end{bmatrix} \Rightarrow A+E = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\frac{\|E\|_2}{\text{gap}} = \frac{.1}{.2} = \frac{1}{2}$$

$$|.5 \sin(2\theta)| \leq \frac{1}{2}$$

$$|\sin(2\theta)| \leq 1 \quad \text{no control on } \theta$$

proof of weaker result (full proof in text)

write evec of $A+E$ as $q_i + d$

where $d^T q_i = 0$



$$q_i' = \frac{q_i + d}{\|q_i + d\|_2} \quad ; \quad \text{goal: bound } \|d\|_2 = \tan \theta$$

$$(A+E)(q_i + d) = \lambda_i' (q_i + d)$$

$$A q_i + A d + \underbrace{E q_i + E d}_{\text{ignore second order term}} = \lambda_i' q_i + \lambda_i' d$$

$$\text{LHS} = (A + E - \lambda'_i I) q_i = (\lambda'_i I - A) d = \text{RHS}$$

$$d = \sum_{j \neq i} d_j q_j \quad \|d\|_2 = \sqrt{\sum_{j \neq i} d_j^2}$$

$$\begin{aligned} \text{LHS} &= (A + E - \lambda'_i I) q_i \\ &= (\lambda'_i I + E - \lambda'_i I) q_i \\ &= ((\lambda'_i - \lambda_i) I + E) q_i \end{aligned}$$

$$\begin{aligned} \|\text{LHS}\|_2 &\leq \|(\lambda'_i - \lambda_i) q_i\|_2 + \|E q_i\|_2 \quad \|q_i\|_2 = 1 \\ &\leq \underbrace{\|E\|_2}_{\text{Weyl}} + \|E\|_2 = 2\|E\|_2 \end{aligned}$$

$$\begin{aligned} \text{RHS} &= (\lambda'_i I - A) \sum_{j \neq i} d_j q_j \\ &= \sum_{j \neq i} (\lambda'_i d_j q_j - \lambda_j d_j q_j) \end{aligned}$$

$$\lambda'_i - \lambda_j = \underbrace{\lambda'_i - \lambda_i}_{\leq \|E\|_2, \text{Weyl}} + \underbrace{\lambda_i - \lambda_j}_{\geq \text{gap}}$$

$$\sum_{j \neq i} (\lambda'_i - \lambda_j) d_j q_j$$

$$\|\text{RHS}\|_2 \geq (\text{gap} - \|E\|) \|d\|_2$$

$$\frac{2\|E\|}{\text{gap}} \approx \frac{2\|E\|}{\text{gap} - \|E\|} \geq \|d\|_2 = \tan \theta$$

if $\|E\| < \text{gap}$

$$Q \neq \emptyset$$

More results on Rayleigh Quotient:
why good eval approximation:

Thm: Given $\|x\|_2 = 1$ and β

Then A has eval α :

$$|\alpha - \beta| \leq \|Ax - \beta x\|_2$$

Given only x , $\beta = p(x, A)$

minimizes $\|Ax - \beta x\|_2$

Given any unit vector x , \exists eval
within distance $\|Ax - p(x, A)x\|_2$ of $p(x, A)$
and $p(x, A)$ minimizes distance

Now let λ_i be eval of A closest
to $p(x, A)$ and $\text{gap} = \min_{j \neq i} |\lambda_j - p(x, A)|$

$$\text{then } |\lambda_i - p(x, A)| \leq \frac{\|Ax - p(x, A)x\|_2^2}{\text{gap}}$$

(result later: cubic convergence)
in shifted QR

proof: Part 1: $\|x\|_2 = 1$

$$1 = \|x\|_2 = \|(A - \beta I)^{-1} (A - \beta I)x\|_2$$

$$\leq \|(A - \beta I)^{-1}\|_2 \cdot \|(A - \beta I)x\|_2$$

$$A = QLQ^T$$

$$= \| (\Lambda - \beta I)^{-1} \|_2 \cdot \| A x - B x \|$$

$$= \frac{1}{\min_i |\lambda_i - \beta|} \cdot \| A x - B x \|$$

$$\min_i |\lambda_i - \beta| \leq \| A x - B x \|_2$$

Part 2: show that $\beta = \rho(x, A)$

minimizes $\| A x - B x \|_2$

$$A x - B x = \underbrace{A x - \rho(x, A) x}_y + \underbrace{\rho(x, A) x - B x}_z$$

$$\text{if } z^T y = 0 \Rightarrow$$

$$\| A x - B x \|_2^2 = \| y \|_2^2 + \| z \|_2^2$$

$$\geq \| y \|_2^2 = \| A x - \rho(x, A) x \|_2^2$$

$$\begin{aligned} z^T y &= (\rho(x, A) - \beta) x^T (A x - \rho(x, A) x) \\ &= (\rho(x, A) - \beta) (x^T A x - \rho(x, A) x^T x) \\ &= 0 \end{aligned}$$

by def of $\rho(x, A)$

Part 3: consider special case of

2×2 diagonal (captures all ideas of full proof):

$$A = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \quad x = \begin{bmatrix} c \\ s \end{bmatrix} \quad c^2 + s^2 = 1$$

$$\rho(x, A) = c^2 d_1 + s^2 d_2$$

assume $c > s \Rightarrow \rho(x, A)$ closer to d_1

can show
$$\frac{\|Ax - \rho(x, A)x\|_2^2}{\text{gap}} = |\lambda_1 - \rho(x, A)|$$

vs \leq in
general case

Algorithms: Design Space

- (1) "Usual Accuracy": backward stable
exact evals, evecs of $A+E$, $\|E\|_2 = O(\epsilon) \cdot \|A\|_2$
- (1.1) all evals (with or without evecs)
 - (1.2) just evals in $[x, y]$ (w or w/o evecs)
 - (1.3) just evals d_i, d_{i+1}, \dots, d_j
eg d_1, \dots, d_{10} : 10 largest evals
(w or w/o evecs)
- (1.2) and (1.3) can be much cheaper than
(1.1) when only few evals/evecs desired
- (2) "High accuracy": get tiny evals (and their evecs)
with more leading digits

Ex: if A well conditioned

i.e. all evals \sim same magnitude
then usual accuracy \Rightarrow error bound
is $O(\epsilon) \|A\| \Rightarrow$ all computed evals
have correct leading digits

What structures let us compute
tiny sing values or evals accurately?

$$B = D \cdot A, \quad D = \text{diag}(d_1, \dots, d_n)$$

some $d_i \gg$ other d_j

\Rightarrow some $\sigma_i \gg$ other σ_j

\Rightarrow usual accuracy guarantees no
correct digits for small σ_j

] perturbation theory and algs that
say small normwise perturbation to A
keeps all σ_i to high relative accuracy
if A well conditioned

Symmetric case: DAD same idea

(1.1) Given T : find all evals,
w or w/o evcs

many algs, all cost $O(n^2)$ for
evals alone, If evcs desired
costs range from $O(n^2)$ to $O(n^3)$,
varying numerical stability

(1.1.1) Oldest is QR iteration
with shift as in Chap 4

Thm (Wilkinson) with right shift,
tridiagonal QR is globally convergent,
usually cubic (# digits triples
at every iteration)

Cost: $O(n^2)$ for evals alone,
but $O(n^3)$ for evcs too

Lapack: `ssyev` (`sssteqr`)

Variant for SVD (Golub + Kahan,
1965)

LAPACK uses a variant (for
bidiagonal SVD) guaranteeing
high relative accuracy in all
sing values (D., Kahan)

(1.1.2) Improve cost of $O(n^3)$ for evecs to $O(n^2)$, but not guarantee orthogonal evecs if evals close:

$\frac{\|Ax - \lambda x\|}{\|A\|}$ will be $O(\epsilon)$, but $x_i^T x_j$ may be large

(1) compute evals: $O(n^2)$... LAPACK: sstebz

(2) compute each evec using 'inverse iteration' (Chap 4)

$$x_{i+1} = (T - \lambda I)^{-1} x_i$$

LAPACK: sstein

fast convergence if λ accurate
if λ, λ' very close, no guarantee of orthogonality; worst case if two evals round to same floating point number

Long term Goal: $O(n^2)$ and orthogonality: solved by MRRR (Parlett) LAPACK: ssyevr (sstemr)

(1.1.3) Divide + conquer:

(Cuppen, Gu + Eisenstat)

faster than QR, not as

fast as MRRR or inverse
iteration, guarantees orthogonal
evecs: cost $O(n^2)$, $2 \leq g \leq 3$

same idea is used for
updating evals evecs of $A + xx^T$
LAPACK: `ssyevd` (`ssstedc`)