

# Perturbation Analysis

(Backward error)

Since we cannot compute exactly we are concerned with effects of errors. (usually relative errors)

Consider comp.  $f(x)$  exact:  $y = f(x)$   
comp:  $\tilde{y}$

In general we want to know  $\|\tilde{y} - y\|$  or  $\frac{\|\tilde{y} - y\|}{\|y\|}$

(Forward error). It turns out that this approach

has problems (soon). Instead we assume we computed an exact answer to perturbed problem

(input) and assess the effect of that perturbation

by perturbation analysis (backward error = perturbation, it may not be unique)

computed exactly  $f(x+\varepsilon)$  (or  $f(x+\varepsilon) + \alpha f(x+\varepsilon)$ )

backward error is  $\varepsilon$ ; analyze  $|f(x+\varepsilon) - f(x)|$

Example: Compute  $u^T v$   $u, v \in \mathbb{R}^n$   
exact answer:  $\sum u_i v_i$   $\hookrightarrow u_i, v_i$ : machine numbers

floating pt.  $f_l(u^T v) = (u_1 v_1 (1+\varepsilon_1) + u_2 v_2 (1+\varepsilon_2)) / (1+\varepsilon_3) + \dots$

usually we just write  $(u_1 v_1 (1+\varepsilon) + u_2 v_2 (1+\varepsilon)) / (1+\varepsilon) + u_3 v_3 (1+\varepsilon) / (1+\varepsilon)$

and keep in mind that all " $\varepsilon$ " are different

but  $|\varepsilon| \leq \varepsilon_M$

$$f(u^T v) = u_1 v_1 (1+\epsilon)^n + u_2 v_2 (1+\epsilon)^{n-1} + u_3 v_3 (1+\epsilon)^{n-2} + \dots$$

(if  $u_i, v_i$  are not mach. numbers it adds factor  $(1+\epsilon)^2$  to each term)

$$u_i \rightarrow u_i (1+\epsilon) \quad u_2 \rightarrow u_2 (1+\epsilon) \quad \text{etc.}$$

$\rightarrow (\epsilon_1 + \epsilon_2 + \dots)$

$$f(u^T v) \approx u^T v + u_1 v_1 n \epsilon + u_2 v_2 n \epsilon + \dots$$

$$|f(u^T v) - u^T v| \leq n |u^T v| \epsilon$$

\* sign of each  $\epsilon$  may differ  $\rightarrow$  assume all errors accumulate (worst case)

\* assume  $\epsilon$  sufficiently small (relative to  $n$ ), that we can ignore  $\epsilon^2, \epsilon^3, \dots$  terms.

\* simplify by ~~ignoring~~ taking  $n \epsilon |u_3 v_3|$   
 i.e.  $(n-1) \epsilon |u_3 v_3|$  etc

$$\text{relative error} = \frac{n \epsilon |u^T v|}{|u^T v|} = \frac{|f(u^T v) - u^T v|}{|u^T v|}$$

rel. error can be huge if  $|u^T v|$  is very small compared with  $|u^T v|$  even if  $n \epsilon$  still small.

$$\|u\|_\infty = \max_i |u_i| = 1 \quad \|v\|_\infty = 1$$

but  $u^T v \approx 0$  (orthogonal)

So, there is no bound on (forward) <sup>relative</sup> errors

We may conclude that "simple" computation of dot product is therefore unreliable.

Backward error analysis:

$$fl(u^T v) \approx u^T v + \cancel{u_1 v_1 \epsilon} + \cancel{u_2 v_2 \epsilon} + \dots$$

$$\approx n \epsilon u_1 v_1 + n \epsilon u_2 v_2 + \dots$$

$$= u_1 v_1 (1 + n \epsilon) + u_2 v_2 (1 + n \epsilon) + u_3 v_3 (1 + n \epsilon) + \dots$$

$$\approx u_1 (1 + \frac{1}{2} n \epsilon) v_1 (1 + \frac{1}{2} n \epsilon) + u_2 (1 + \frac{1}{2} n \epsilon) v_2 (1 + \frac{1}{2} n \epsilon) + \dots$$

(assuming  $\frac{1}{4} n^2 \epsilon^2$  negligible)

$$= \tilde{u}_i^T \tilde{v}_i \text{ (exactly)}$$

$$\tilde{u}_i = u_i + \eta \quad \text{where } |\eta_i| \leq \frac{1}{2} n \epsilon_M$$

$$\text{so } \|\eta\|_{\infty} \leq \frac{1}{2} n \epsilon_M \quad \text{and } \|u\|_{\infty} = 1$$

So, floating pt. computation of dot product has small relative backward error.

What explains the huge relative forward error of  $u^T v \approx 0$  is the "sensitivity" or "conditioning" of the ~~compu~~ problem. If  $u^T v \approx 0$  the dot product is ill-conditioned. If  $u^T v$  not very small the problem is well-conditioned.

Backward error analysis let's us separate the accumulation of numerical error from

the sensitivity of the problem. So, we combine backw. error anal. with perturbation analysis to obtain bounds on computed answers. If an algorithm produces small (bounded) backward error the accuracy depends on the sensitivity.

ill-cond  $\rightarrow$  answer may be poor  
(but cannot be helped)

well-cond  $\rightarrow$  answer accurate.

# Perturbation Theory

The eig. vals./vec.s are (typically) exact for slightly perturbed matrix (small backward error)

We want to know how much eig. vals. and vec.s change under small perturbation

1) Eigenvalues are continuous function(s) of the matrix (coefficients)

Let  $A$  have eig. values  $d_1, \dots, d_n$  (counting mult.s) and  $\tilde{A} = A + E$  have eig. values  $\tilde{d}_1, \dots, \tilde{d}_n$ . Then there exists ordering  $j_1, \dots, j_n$  s.t.

$$|d_i - \tilde{d}_{j_i}| \leq 4(\|A\|_2 + \|\tilde{A}\|_2)^{1-\frac{1}{n}} \|E\|_2^{\frac{1}{n}}$$

\* The exponent  $\frac{1}{n}$  is in general pessimistic but necessary (sharp  $\rightarrow$  obtained in specific cases, esp. large Jordan block)

Gershgorin disks

$Ax = \lambda x$ , let  $x_i$  be largest abs comp. scale x s.t.  $x_i = 1$  ( $|x_j| \leq 1$   $j \neq i$ )

$$a_{ii} + \sum_{j \neq i} a_{ij} x_j = \lambda x_i = \lambda \Leftrightarrow$$

$$a_{ii} - \lambda = - \sum_{j \neq i} a_{ij} x_j$$

$$|a_{ii} - \lambda| \leq \sum_{j \neq i} |a_{ij}| |x_j| \leq \sum_{j \neq i} |a_{ij}|$$

So,  $\lambda$  lies inside disk  $G_i = \{z \mid |z - a_{ii}| \leq \sum_{j \neq i} |a_{ij}|\}$

Define disks  $G_1, \dots, G_n$  this way

$$\Lambda(A) \subset \bigcup_{i=1, \dots, n} G_i$$

Theo 3.1 p. 37

~~$$\text{hd}(A, \tilde{A}) \leq \frac{\|A\|_2 + \|\tilde{A}\|_2}{\|E\|_2^{\frac{1}{n}}}$$~~

"Proof" (From Stewart and Sun)  
 $\text{hd}(A, \tilde{A}) \leq \frac{\|A\|_2 + \|\tilde{A}\|_2}{\|E\|_2^{\frac{1}{n}}}$

$$\text{md}(A, \tilde{A}) \leq (2n-1) \cdot \text{hd}(A, \tilde{A})$$

ith row for  
 $Ax = \lambda x$

Theo 3.2

Moreover if  $G_{i_1} \cup \dots \cup G_{i_k} \cap B_{i_{k+1}} \dots G_{i_n} = \emptyset$

$$G_{i_1} \cup \dots \cup G_{i_k} \cap B_{i_{k+1}} \dots G_{i_n} = \emptyset$$

disjoint from complement for some ~~ordering~~ permutation  $i_1, \dots, i_n$  of  $\{1, \dots, n\}$ , then

$G_{i_1} \cup \dots \cup G_{i_k}$  contains exactly  $k$  eigenvalues

Proof by continuity  $A = D + B$   
 $\uparrow$   
 $\text{diag } A$

$$A(t) = D + tB$$

$$\Lambda(A(0)) = \{a_{11}, a_{22}, \dots, a_{nn}\}$$

By continuity eigenvalues stay in disks

$$G_i(t) = \left\{ z : |z - a_{ii}| \leq t \sum_{j \neq i} |a_{ij}| \right\}$$

If the column sum (off-diag) is much larger than the row sum we can improve the disk corresponding to that row by a diagonal similarity transformation

$$D^{-1} A D \text{ where } d_i \text{ is small}$$

(what limits how small we can make  $d_i$ ?)

Theo 3.5

Let  $A = A^H$  (Hermitian) and let

$$d_1 \geq d_2 \geq \dots \geq d_n \text{ then}$$

$$\max_{\substack{W \\ \dim(W)=i}} \min_{\substack{w \in W \\ \|w\|_2=1}} w^H A w = d_i$$

$$\min_{\substack{W \\ \dim(W)=n-i}} \max_{\substack{w \in W \\ \|w\|_2=1}} w^H A w = d_i$$

Outline of Proof  
 from Horn &  
 Johnson I,  
 see also  
 Stewart & Sun  
 (Courant-Fischer/  
 Fischer's,  
 Rayleigh-Ritz)

A Hermitian  $\rightarrow A = U \Lambda U^H$ ,  $\|w\|_2 = 1$

$$w^H A w = w^H U \Lambda U^H w = \xi^H \Lambda \xi \text{ where}$$

$$\xi = U^H w. \text{ Also assume } d_1 \geq d_2 \geq \dots \geq d_n$$

$$\xi^H \Lambda \xi = \sum_i d_i |\xi_i|^2 \text{ and we have}$$

$$d_1 = \sum_i d_i |\xi_i|^2 \geq \sum_i d_i |\xi_i|^2 = w^H A w$$

note  $d_1 = u_1^H A u_1$  (so bound sharp)

$$\left( \sum |\xi_i|^2 = \xi^H \xi = w^H U^H U w = w^H w = 1 \right)$$

$$\text{Same way } d_n \leq \sum_i d_n |\xi_i|^2 \leq \sum_i d_i |\xi_i|^2 = w^H A w$$

(sharp since  $u_n^H A u_n = d_n$ )

Consider  $\max_{\substack{\|w\|_2=1 \\ w \perp u_1}} w^H A w$

$$\xi = U^H w = \begin{pmatrix} 0 \\ \xi_2 \\ \vdots \\ \xi_n \end{pmatrix} \quad \xi_1 = 0 \Leftrightarrow u_1^H w = 0$$

$$w^H A w = \sum_i d_i |\xi_i|^2 = \sum_{i=2}^n d_i |\xi_i|^2 \leq d_2 \sum_{i=2}^n |\xi_i|^2 = d_2$$

$$\text{Hence } \max_{\substack{\|w\|_2=1 \\ w \perp u_1}} w^H A w = d_2$$

and so on.

Now we'll show that these bounds correspond to the max/min & min/max characterizations.

let  $v$  be arbitrary vector and consider ( $v \neq 0$ )

$$v \in \mathbb{C}^n, v \neq 0$$

"max over space  $\Rightarrow$   
max over subspace"

$$(\Rightarrow \gamma_3 = \gamma_4 = \dots = 0)$$

$$1: \max_{\substack{\|w\|_2=1 \\ w \perp v}} w^H A w$$

$$2: \min_{\substack{v \in \mathbb{C}^n \\ v \neq 0}} \left( \max_{\substack{\|w\|_2=1 \\ w \perp v}} w^H A w \right)$$

$$1: \max_{\substack{w \perp v \\ \|w\|_2=1}} w^H A w = \max_{\substack{w^H u_1 u_1^H w = \dots \\ \|w\|_2=1}} w^H A w =$$

$$\max_{\substack{\gamma = u^H w \perp u^H v \\ \|\gamma\|_2=1}} \gamma^H \Lambda \gamma = \max_{\dots} \sum_{i=1}^n d_i |\gamma_i|^2$$

$$\geq \max_{\substack{\gamma \perp u^H v \\ |\gamma_1|^2 + |\gamma_2|^2 = 1}} d_1 |\gamma_1|^2 + d_2 |\gamma_2|^2 \geq d_2$$

$$\max_{\substack{\|w\|_2=1 \\ w \perp v}} w^H A w \geq d_2 \quad \text{and}$$

$$\max_{\substack{w \perp u_1 \\ \|w\|_2=1}} w^H A w = d_2$$

$$\text{So, } \min_{\substack{v \in \mathbb{C}^n \\ v \neq 0}} \max_{\substack{w \perp v \\ \|w\|_2=1}} w^H A w = d_2$$

"retroactively" all min/max (rather than inf/sup) okay since min/max

obtained (must be for min.../max... over sphere in finite dimensions)

Note that  $\dim(W) = n \Rightarrow W = \mathbb{C}^n$ ,

$\dim(W) = n-1 \Rightarrow W \perp v (=0)$  for some  $v \in \mathbb{C}^n$

$\dim(W) = n-2 \Rightarrow W \perp v_1, v_2 \in \mathbb{C}^n$ , etc

Based on these ideas we can prove a many other properties of eigenvalues of  $A$  and of representations of  $A$  over subspace

theo 3.6:  $A$  Hermitian, ~~with~~ with eig. vals

$$d_1 \geq d_2 \geq \dots \geq d_n$$

$U^{n \times m}$  orthonormal and

$U^H A U$  has eig. vals  $\mu_1 \geq \mu_2 \geq \dots \geq \mu_m$

$$d_1 \geq \mu_1 \geq d_{n-m+1}$$

$$d_2 \geq \mu_2 \geq d_{n-m+2}$$

$$d_i \geq \mu_i \geq d_{n-m+i} \quad i=1..m$$

When  $m = n-1 \rightarrow$

$$d_1 \geq \mu_1 \geq d_2 \geq \mu_2 \geq \dots \geq d_{n-1} \geq \mu_{n-1} \geq d_n$$

$$d_i \in \Delta(A): d_1 \geq d_2 \geq \dots$$

$$d_i \in \Delta(\tilde{A}): \tilde{d}_1 \geq \tilde{d}_2 \geq \dots$$

$$\varepsilon_i \in \Delta(E): \varepsilon_1 \geq \varepsilon_2 \geq \dots$$

Consider  $A, \tilde{A} = A + E$  Hermitian

$$x_i^H \tilde{A} x_i = x_i^H A x_i + x_i^H E x_i \quad (\exists A x_i = d_i x_i)$$

$$\tilde{d}_i = d_i + x_i^H E x_i \quad \varepsilon_i \geq x_i^H E x_i \geq \varepsilon_n$$

$$\tilde{d}_i \leq d_i + \varepsilon_1 \quad \text{and} \quad \tilde{d}_i \geq d_i + \varepsilon_n$$

More generally (theo 3.8)

$$d_i + \varepsilon_n \leq \tilde{d}_i \leq d_i + \varepsilon_1 \quad i=1, \dots, n$$

$$\text{So, } |\tilde{d}_i - d_i| \leq \|E\|_2 \quad (= \varepsilon_i)$$

Also (theo 3.10: Hoffman-Wielandt)

$$\left( \sum_i (\tilde{d}_i - d_i)^2 \right)^{1/2} \leq \|E\|_F$$

(simple eig. pairs)

$$\begin{pmatrix} y^H \\ y^H \end{pmatrix} A(x \ X) = \begin{pmatrix} d & 0 \\ 0 & M \end{pmatrix}$$

$$y^H x = 1$$

$$y^H X = I$$

$$\begin{cases} Ax = dx \\ \text{neglect "small" } x \\ \text{"small" } \end{cases}$$

Next we consider eigenpairs:

$$A \rightarrow (d, x) \quad \tilde{A} = A + E \rightarrow (\tilde{d}, \tilde{x})$$

$$\tilde{d} = d + \varphi$$

$$\tilde{x} = x + Xp \quad (X \text{ complementary right eigenspace})$$

$$(A+E)(x+Xp) = (d+\varphi)(x+Xp) \Leftrightarrow$$

$$Ax + Ex + AXp + EXp = dx + \varphi x + dXp + \varphi Xp$$

$$AXp + Ex \approx dXp + \varphi x$$

$$\text{I: } y^H AXp + y^H Ex = d y^H Xp + \varphi y^H x$$

$$\varphi y^H Ex = \varphi$$

$$\text{II } y^H AXp + y^H Ex = d y^H Xp + \varphi y^H x$$

$$Mp + y^H Ex = dp \Leftrightarrow (M - dI)p = y^H Ex \Leftrightarrow$$

$$p = (M - dI)^{-1} y^H Ex$$

$$(\tilde{d}, \tilde{x}) \approx (d + y^H Ex, x + X(M - dI)^{-1} y^H Ex)$$

Theo 3.11

let  $(d, x)$  be simple eigen pair of  $A$ ,  
 $(x, X)$  nonsingular and  $(y, Y)^H = (x, X)^{-1}$

$\tilde{A} = A + E$ . let

$$\begin{pmatrix} y^H \\ y^H \end{pmatrix} E \begin{pmatrix} x \\ X \end{pmatrix} = \begin{pmatrix} \varphi_{11} & F_{12}^H \\ F_{21} & F_{22} \end{pmatrix}$$

$\|\cdot\|$  is a consistent norm and

$$\text{sep}(d, M) = \|(dI - M)^{-1}\|^{-1}$$

$$\forall \rho \quad 4\|F_{21}\|\|F_{12}^H\| < (\text{sep}(d, M) - |\varphi_{11}| - \|F_{22}\|)^2$$

then there exists  $\varphi$  and  $p$  such that

$(\tilde{d}, \tilde{x}) = (d + \varphi, x + Xp)$  is an eig. pair  
of  $\tilde{A}$  with  $\tilde{d} \notin \Lambda(M)$

Moreover,

$$\|p\| < \frac{2\|F_{21}\|}{\text{sep}(d, M) - |\varphi_{11}| - \|F_{22}\|}$$

$$\|p - (dI - M)^{-1} F_{21}\| < \frac{2\|p\|^2 \|F_{12}^H\|}{\text{sep}(d, M) - |\varphi_{11}| - \|F_{22}\|}$$

and

$$\|\varphi - y^H E x\| \leq \|p\| \|F_{12}\|$$

Note:

$$|\varphi| = |\tilde{d} - d| = O(\|E\|) \text{ and}$$

$$\|\varphi - y^H E x\| = O(\|E\|^2)$$

Proof see

Stewart & Sun

(discuss dep. on  
time)

So, simple eigenvalue <sup>is</sup> ~~are~~ much better behaved than ~~a~~ defect eig. val.

$$O(\|E\|) \quad \text{vs} \quad O(\|E\|^{1/k})$$

Moreover

$$\tilde{\lambda} - y^H \tilde{A} x = \tilde{\lambda} - d - y^H E x$$

$$\text{So, } |\tilde{\lambda} - y^H \tilde{A} x| = O(\|E\|^2)$$

$y^H \tilde{A} x / y^H x$  Rayleigh quotient

$$\hookrightarrow = y^H A x \quad (\text{since } y^H x = 1)$$

Let  $E = \varepsilon e_i e_j^H$  then

$$\tilde{\lambda} - d = \varepsilon \bar{y}_i x_j + O(\varepsilon^2)$$

Hence  $d(a_{ij})$  differentiable and

$$\frac{\partial d}{\partial a_{ij}} = \bar{y}_i x_j$$

eigenvector  $\tilde{x}$ :

$$\|\tilde{x} - x\| = O(\|E\|) \quad \text{and}$$

$$\|\tilde{x} - [x + X(dI - M)^{-1} f_{21}]\| = O(\|E\|^2)$$

Conditioning/

Sensitivity

(for small pert.)

$$\tilde{\lambda} = d + y^H E x + (\varphi - y^H E x)$$

$$|\tilde{\lambda} - d| \leq \|x\| \|y\| \|E\| + O(\|E\|^2)$$

$$\leq \sec \Delta(x, y) \|E\| + O(\|E\|^2)$$

We call  $\sec \Delta(x, y)$  condition number

$$\text{Note } \cos \Delta(x, y) = \left| \left( \frac{x}{\|x\|} \right)^H \left( \frac{y}{\|y\|} \right) \right| =$$

$$\frac{1}{\|x\| \cdot \|y\|} \Rightarrow \|x\| \cdot \|y\| = \sec \Delta(x, y)$$

Also common to take  $\|x\|=1, \|y\|=1$   
and  
 $\text{cond}(d) = \frac{1}{|y^H x|}$

$x, y$  not unique  $\rightarrow$  assume  $S$   
s.t.  $yS^{-H}$  orthonormal

replace  $x$  by  $xS, \Pi$  by  $S^{-1}MS$

Now assume  $y$  orthonormal.

(Lemma 3.12)

$$\text{Then } \sin \Delta(x, \tilde{x}) = \frac{\|y^H \tilde{x}\|_2}{\|\tilde{x}\|_2}$$

Let  $\|x\|_2 = \|\tilde{x}\|_2 = 1 \rightarrow (x, y)$  unitary

$$\left\| \begin{pmatrix} x^H \\ y^H \end{pmatrix} \tilde{x} \right\|_2^2 = |x^H \tilde{x}|^2 + \|y^H \tilde{x}\|_2^2 = 1$$

$$|x^H \tilde{x}|^2 = \cos^2 \Delta(x, \tilde{x}) \Rightarrow$$

$$\|y^H \tilde{x}\|_2^2 = \sin^2 \Delta(x, \tilde{x})$$

Assume Theo 3.11: Theo 3.13  
wrt  $\|\cdot\|_2$

$$\sin(\alpha, \tilde{\alpha}) \leq \frac{\|E\|_2}{\text{sep}(\tilde{d}, M)}$$

Proof:

$$(A+E)\tilde{x} = \tilde{d}\tilde{x}$$

$$y^H A \tilde{x} - \lambda y^H \tilde{x} = (M - \tilde{d}I) y^H \tilde{x} = -y^H E \tilde{x}$$

~~$$\tilde{x} = (M - \tilde{d}I)^{-1} (-y^H E \tilde{x})$$~~

$$y^H \tilde{x} = (dI - M)^{-1} y^H E \tilde{x}$$

$$\frac{\|y^H \tilde{x}\|_2}{\|\tilde{x}\|_2} = \|(dI - M)^{-1}\|_2 \|y^H E \tilde{x}\|_2 / \|\tilde{x}\|_2$$

$$\frac{\|y^H \tilde{x}\|_2}{\|\tilde{x}\|_2} \leq \frac{\|E\|_2 \|\tilde{x}\|_2}{\|(dI - M)^{-1}\|^{-1} \|\tilde{x}\|_2} = \frac{\|E\|_2}{\text{sep}(\tilde{d}, M)}$$

$\uparrow$   
 $\sin(x, \tilde{x})$

For ~~small~~  $E \rightarrow 0$ ,  $\tilde{d} \rightarrow d$ . So,

$$\sin(x, \tilde{x}) \lesssim \frac{\|E\|_2}{\text{sep}(d, M)}$$

So,  $\text{sep}^{-1}(d, M)$  is a condition number

for eig. vector  $x$

Properties of  $\text{sep}$

$\text{sep}$  ?

$\text{sep}(d, M)$  is lower bound on separation of  $d$  from spectrum of  $M$ .

Theo 3.14

In any consistent norm,

$$\text{sep}(d, M) \leq \min_{\mu \in \Lambda(M)} |d - \mu|$$

$$\text{sep}(d, M)^{-1} = \|(dI - M)^{-1}\| \geq \rho((dI - M)^{-1}) =$$

$$\max_{\mu \in \Lambda(M)} |d - \mu|^{-1} \Rightarrow$$

$$\text{sep}(d, M) \leq \min_{\mu \in \Lambda(M)} |d - \mu|$$

However,  $\text{sep}(d, M)$  can be much smaller than  $\min |d - \mu|$  if an eigenvalue of  $M$  is ill-conditioned

Hermitian case

Since eigenvalues of Hermitian matrix are perfectly well-conditioned

(cond. nr = 1)

we have

Theo 3.16

Let  $d$  be real,  $M$  Hermitian. Then in 2-norm

$$\text{sep}(d, M) = \min_{\mu \in \Delta(M)} |d - \mu|$$

$\nabla$  If  $A$  Hermitian, we can take  $Y = X$  (See above), so that  $M$  Herm.

$$\sin \Delta(x, \tilde{x}) \lesssim \frac{\|E\|_2}{\min_{\mu \in \Delta(M)} |d - \mu|}$$

So, it's the separation of eigenvalues that counts.

# Perturbation Theory for Eig. Spaces (Chap. 2. Sec. 2)

Canonical angles  
p. 250

Let  $X$  and  $Y$  be subspaces of same dim,  $P$  orthonormal basis for  $X$ ,  $Y$  for  $Y$

Then  $\theta_i = \arccos y_i$  are the canonical angles between  $X$  and  $Y$ , where

$y_i$  are singular values of  $Y^H X$ .

We write  $\theta_i(X, Y)$  for  $i$ -th can. angle in descending order.

(also  $\Theta(X, Y)$ )

$$\Theta(X, Y) = \text{diag}(\theta_1, \dots, \theta_p)$$

If  $Y_\perp$  orthonormal basis for  $Y^\perp$

(theo 2.2)

then  $\sigma_i(Y_\perp^H X)$  are  $\sin(\theta_i)$

Also

$$\begin{pmatrix} Y^H \\ Y_\perp^H \end{pmatrix} X = \begin{pmatrix} C \\ S \end{pmatrix} \quad \text{where} \quad C^H C + S^H S = I$$

$$\sin \Delta(x, Y) = \min_{y \in Y} \|x - y\|_2$$

$$\rightarrow \cancel{\|x - y\|_2} \quad \|x - Y Y^H x\|_2$$

Theo 2.4

Let  $X$  and  $Y$  be orthonormal,  $X^H Y = 0$ ,

$$Z = X + Y Q$$

$\sigma_i(Z)$  nonzero singular vals descending  
 $\sigma_i(Q)$  " " " descending

$\theta_i(X, Z)$  nonzero canon. angles between  $R(X)$  and  $R(Z)$  in descending order

Note  $K = U \Sigma V^H$

$$K^H K = V \Sigma^H \Sigma V^H$$

↑  
eig. val.s  
↑  
eig. vec.s

$$\text{Then } \sigma_i = \sec \varphi_i$$

$$\tau_i = \tan \varphi_i$$

$$Z^H Z = (X^H + Q^H Y^H)(X + YQ)$$

$$= X^H X + \underbrace{X^H Y Q}_{=0} + \underbrace{Q^H Y^H X}_{=0} + Q^H Q$$

$$= I + Q^H Q$$

$$\sigma_i^2 = 1 + \tau_i^2$$

To find canon. angles  $\varphi_i(X, Z)$  we

need orthonormal basis for  $R(Z)$

$$\text{Take } \hat{Z} = Z (Z^H Z)^{-1/2} \rightarrow \hat{Z}^H \hat{Z} = I$$

$$X^H \hat{Z} = X^H (X + YQ) (I + Q^H Q)^{-1/2}$$

$$= (I + Q^H Q)^{-1/2}$$

$$\cos \varphi_i = \frac{1}{\sigma_i} \Rightarrow \sigma_i = \frac{1}{\cos \varphi_i} = \sec \varphi_i$$

~~or~~

$$\sigma_i^2 = \frac{1}{\cos^2 \varphi_i} = \frac{\cos^2 \varphi_i + \sin^2 \varphi_i}{\cos^2 \varphi_i}$$

$$= 1 + \tan^2 \varphi_i = 1 + \tau_i^2$$

Corollary 2.5

Let  $X$  be orthonorm. basis for simple eig. space  $X$  of  $A$  and  $Y$  be basis for corresponding left eig. space s.t.

$$Y^H X = I$$

$$\text{Then } \sigma_i(Y) = \sec \varphi_i(X, Y)$$

$$\text{and } \|Y\|_2 = \sigma_{\max}(Y) = \sec \varphi_1(X, Y)$$

## Residual Analysis

Approx. eigenbasis  $X$  (basis for appr. eig. space):

$$R = AX - XL$$

Best  $L$  ?

Theorem 2.6

$(X \ X_{\perp})$  unitary  $R = AX - XL$

Then  $\|R\|$  minimized for any unit. invar. norm if

$$L = X^H A X \rightarrow \text{Rayleigh quotient}$$

More generally  $X^L$  left inverse of  $X$   
then  $X^L A X$  is Rayleigh quot. of  $A$

## Backward Error

Theo 2.8

$X$  orthonormal and  $R = AX - XL$

$$\text{Then } E = -R X^H \rightarrow (A + E)X = XL$$

$$\|E\|_2 = \|R\|_2, \quad \|E\|_F = \|R\|_F$$

If  $A$  Hermitian and  $L = X^H A X$

$$\text{then } E = -(R X^H + X^H R)$$

$$(A + E)X = XL$$

$$\|E\|_2 = \|R\|_2, \quad \|E\|_F = \sqrt{2} \|R\|_F$$

Proof: by construction + taking norms

(From previous typically take  
 $L = X^H A X$

II

$$R = AX - XL, \quad X^H X = I$$

We can proceed in 2 ways

$$(A+E)X = XL \quad \text{and consider}$$

~~perturbation~~ eigen pair of perturbed matrix

I

Consider block deflation / similarity transform and changes to  $X$  and  $L$

I

$(X \ X_{\perp})$  unitary

$$\begin{pmatrix} X^H \\ X_{\perp}^H \end{pmatrix} A \begin{pmatrix} X & X_{\perp} \end{pmatrix} = \begin{pmatrix} L & H \\ G & M \end{pmatrix}$$

$$A \begin{pmatrix} X & X_{\perp} \end{pmatrix} = \begin{pmatrix} XL + X_{\perp} G & X_{\perp} M + X H \end{pmatrix}$$

$$AX - XL = X_{\perp} G = R$$

$$\Rightarrow \|G\| = \|R\| \quad \text{for u.i.n.}$$

$$X_{\perp}^H (AX - XL) = G = X_{\perp}^H R$$

Sim. Transform:

$$\begin{pmatrix} I & 0 \\ -P & I \end{pmatrix} \begin{pmatrix} L & H \\ G & M \end{pmatrix} \begin{pmatrix} I & 0 \\ P & I \end{pmatrix} =$$

$$\begin{pmatrix} L+HP & H \\ G+MP-PL-PHP & M-PH \end{pmatrix} = \begin{pmatrix} * & * \\ 0 & * \end{pmatrix}$$

$$P \text{ s.t. } PL - MP = G - PHP$$

~~Define~~ Define

$$\text{sep}(L, M) = \min_{\|Q\|=1} \|QL - MQ\|$$

unitarily  
invariant  
norm

If  $4\|G\|\|H\| < \text{sep}^2(L, M)$  then

unique  $P$  exists s.t.  $PL - MP = G - PHP$   
and

$$\|P\| < \frac{2\|G\|}{\text{sep}(L, M)}$$

Furthermore

$$\Lambda(L + \overset{HP}{P}) \cap \Lambda(M - PH) = \emptyset$$

Proof hinges on convergence of  
fixed point iteration (solving for  $P$ )

$G_0 = G \rightarrow$  Solve  $P_0 : P_0 L - MP_0 = G_0$   
if  $\|G\|$  suff. small then  $\|P_0\|$  small

and  $\|P_0\|^2$  very small  $\rightarrow$

$G_1 = G_0 - P_0 H P_0$  Solve  $P_1 L - MP_1 = G_1$

Overall Sim. Transform for  $A$ :

$$\begin{pmatrix} I & 0 \\ -P & I \end{pmatrix} (X \ X_\perp)^H A (X \ X_\perp) \begin{pmatrix} I & 0 \\ P & I \end{pmatrix} = \begin{pmatrix} * & * \\ 0 & * \end{pmatrix}$$

$A \hat{X} = \hat{X} \hat{L}$  where

$$\hat{X} = X + X_\perp P, \quad \hat{L} = L + HP$$

From theo 2.4 (chap 4) :  $\tan \nu_i(X, \hat{X}) = \sigma_i(P)$

$$\rightarrow \tan \nu_i = \sigma_i(P) < \frac{2\|G\|}{\text{sep}(L, M)}$$

(details in  
Steward & Sun)

note earlier alg.  
for Sylvester eq.

normalize  $\|A\| \leq 1$   
(by scaling)

then  $\|H\| \leq 1$

Theo 2.12

$$R = AX - XL$$

$$S = X^H A - L X^H$$

$$\text{If } 4\|R\|\|S\| < \text{sep}^2(L, M)$$

$(\hat{L}, \hat{X})$  exists (for A)

$$A\hat{X} = \hat{X}\hat{L}$$

$$\|L - \hat{L}\| < \frac{2\|R\|\|S\|}{\text{sep}(L, M)}$$

$$\|\tan \Theta(X, \hat{X})\| < \frac{2\|R\|}{\text{sep}(L, M)}$$

sep is complicated but has nice properties. Note for  $\|\cdot\|_2$  sep is smallest singular value of Sylvester operator (which can be written as a matrix)

Theo 2.11 lists many properties.

Here we need that sep itself is well-conditioned

$$\frac{1}{\|E\| + \|F\|} |\text{sep}(L+E, M+F) - \text{sep}(L, M)| \leq$$

(sep defined for consistent norm)

$$X^H X = I \quad \text{II}$$

$$\tilde{A}X = XL + R \rightarrow \tilde{A}X_1 - X_1L = R$$

$$(\tilde{A} - RX_1^H)X_1 = XL + R - R = X_1L \text{ (exact)}$$

$$E = RX_1^H \rightarrow \tilde{A} = A + E \text{ and } AX_1 = X_1L$$

$$\text{Further let } A = X_1LY_1^H + X_2MY_2^H$$

$$\begin{pmatrix} Y_1^H \\ Y_2^H \end{pmatrix} (A + E) \begin{pmatrix} X_1 & X_2 \end{pmatrix} = \begin{pmatrix} L + F_{11} & F_{12} \\ F_{21} & M + F_{22} \end{pmatrix}$$

(For consistent family of norms)

$$\delta = \text{sep}(L, M) - \|F_{11}\| - \|F_{22}\|$$

$$\text{If } 4\|F_{21}\|\|F_{12}\| < \delta^2 \text{ then}$$

$$P \text{ exists s.t. } \|P\| < \frac{2\|F_{21}\|}{\delta}$$

$$(\tilde{L}, \tilde{X}_1) = (L + F_{11} + F_{12}P, X_1 + X_2P)$$

$$(\tilde{M}, \tilde{Y}_2) = (M + F_{22} - PF_{12}, Y_2 - Y_1PH)$$

are simple, complementary, right and left eigenpairs of  $\tilde{A}$

$$\begin{aligned} &\leftarrow \\ \text{sep}(L + F_{11}, M + F_{22}) &\leq \delta \end{aligned}$$