

# Rayleigh Quotient Based Numerical Methods For Eigenvalue Problems

Ren-Cang Li

University of Texas at Arlington

Gene Golub SIAM Summer School 2013  
10th Shanghai Summer School on Analysis and Numerics in  
Modern Sciences

July 22 – August 9, 2013

## Overview

- Motivating Examples
- Hermitian Eigenvalue Problem – Basics
- Steep Descent/Ascent Type Methods
- Conjugate Gradient Type Methods
- Extending Min-Max Principles: Indefinite  $B$
- Linear Response Eigenvalue Problem
- Quadratic Hyperbolic Eigenvalue Problem

# Motivating Examples

- Density Functional Theory – Kohn-Sham Equation
- Data Mining – Trace Ratio Maximization

More in Chapter 10 of

 Y. Saad. *Numerical Methods for Large Eigenvalue Problems*. SIAM, 2011.

**Kohn-Sham Equation** (Hohenberg and Kohn'64, Kohn and Sham'65):

$$\left[ -\frac{1}{2}\nabla^2 + \underbrace{\int \frac{n(\mathbf{r}')}{|\mathbf{r}-\mathbf{r}'|} d\mathbf{r}' + \frac{\delta E_{xc}(n(\mathbf{r}))}{\delta n(\mathbf{r})}}_{v_{KS}[n](\mathbf{r})} + v_{\text{ext}}(\mathbf{r}) \right] \phi_i(\mathbf{r}) = \lambda_i \phi_i(\mathbf{r}),$$

a remarkably successful theory to describe ground-state properties of condensed matter systems.

**A nonlinear eigenvalue problem:** Kohn-Sham (KS) operator depends on *electronic*

*density*  $n(\mathbf{r}) = \sum_{i=1}^{N_v} \phi_i(\mathbf{r})\phi_i^*(\mathbf{r})$  which depends on eigen-functions  $\phi_i(\mathbf{r})$ .

Usually solved by Self-Consistent-Field (SCF) iteration:

1) initial  $n_0(\mathbf{r}) = \sum_{i=1}^{N_v} \phi_i^{(0)}(\mathbf{r})\phi_i^{(0)*}(\mathbf{r})$ , and

2) repeat  $[-\nabla^2/2 + v_{KS}[n_j](\mathbf{r})] \phi_i^{(j+1)} = \lambda_i^{(j+1)} \phi_i^{(j+1)}(\mathbf{r})$ .

Each inner-iteration is an eigenvalue problem.

# Discretized Kohn-Sham Equation

Ways of discretizations: plane waves, finite differences, finite elements, localized orbitals, and wavelets.

**Discretized Kohn-Sham Equation:**  $H(X)X = SX\Lambda$ ,  $X^H SX = I_{N_v}$ .

$H(X)$  is symmetric, depends on  $X$ , eigenvalue matrix  $\Lambda$  is diagonal, and  $S \succ 0$ . Some discretizations:  $S = I$ .

Nonlinear eigenvalue problem, dependent on eigenvectors, as oppose to usually on the eigenvalues.

Usually solved by Self-Consistent-Field (SCF) iteration:

- 1) initial  $X_0$ , and
- 2) repeat  $H(X_j)X_{j+1} = SX_{j+1}\Lambda_j$  for  $j = 0, 1, \dots$  until convergence.

Each inner-iteration is a symmetric eigenvalue problem.

References more accessible to numerical analysts:



Yousef Saad, James R. Chelikowsky, Suzanne M. Shontz, *Numerical Methods for Electronic Structure Calculations of Materials*, SIAM Rev. 52:1 (2010), 3-54.



C. Yang, J. C. Meza, B. Lee, and L.-W. Wang. KSSOLV—a MATLAB toolbox for solving the Kohn-Sham equations. *ACM Trans. Math. Software*, 36(2):1–35, 2009.

# Trace Optimization

In Fisher linear discriminant analysis (LDA) for supervised learning, need to solve

$$\max_{V^T V = I_k} \frac{\text{trace}(V^T A V)}{\text{trace}(V^T B V)},$$

where  $A, B \in \mathbb{R}^{n \times n}$  symmetric,  $B$  positive semidefinite and  $\text{rank}(B) > n - k$ .

$\text{trace}(V^T A V)$  represents the in-between scatter, while  $\text{trace}(V^T B V)$  represents the within scatter. Maximizer  $V$  is used to project  $n$ -dimensional vectors (data) into  $k$ -dimensional vectors that best separates  $n$ -dimensional datasets into two or more datasets.

KKT condition for Maximizers:

$$\underbrace{\left[ A - \frac{\text{trace}(V^T A V)}{\text{trace}(V^T B V)} B \right]}_{=: E(V)} V = V [V^T E(V) V]$$

such that eigenvalues of  $V^T E(V) V$  are the  $k$  largest eigenvalues of  $E(V)$ .

Can be solved via SCF-like iteration; each inner iteration is a symmetric eigenvalue problem.

## References for trace ratio problem:

-  T. Ngo, M. Bellalij, and Y. Saad. The trace ratio optimization problem for dimensionality reduction. *SIAM J. Matrix Anal. Appl.*, 31(5):2950–2971, 2010.
-  L.-H. Zhang, L.-Z. Liao, and M. K. Ng. Fast algorithms for the generalized Foley-Sammon discriminant analysis. *SIAM J. Matrix Anal. Appl.*, 31(4):1584–1605, 2010.

More eigenvalues arising from Data mining can be found in chapter 2 of

-  S. Yu, L.-C. Tranchevent, B. De Moor, and Y. Moreau. *Kernel-based Data Fusion for Machine Learning: Methods and Applications in Bioinformatics and Text Mining*. Springer, Berlin, 2011.

# Basic Theory

- Hermitian  $Ax = \lambda x$
- Hermitian  $Ax = \lambda Bx$  ( $B \succ 0$ )
- Justifying Rayleigh-Ritz

# Hermitian $Ax = \lambda x$

Hermitian  $A = A^H \in \mathbb{C}^{n \times n}$ .

Eigenvalues  $\lambda_j$  and eigenvectors  $u_j \in \mathbb{C}^n$ .

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n, \quad u_j^H u_j = \delta_{ij}, \quad Au_j = \lambda_j u_j.$$

Rich, elegant, and well-developed theories in “every” aspect ...

## Popular References



R. Bhatia. *Matrix Analysis*. Springer, New York, 1996.



R. A. Horn, C. R. Johnson, *Matrix Analysis*, Cambridge University Press, Cambridge, 1985.



J. Demmel. *Applied Numerical Linear Algebra*. SIAM, Philadelphia, PA, 1997.



G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, Maryland, 3rd edition, 1996.



B. N. Parlett. *The Symmetric Eigenvalue Problem*. SIAM, Philadelphia, 1998.



Lloyd N. Trefethen and David Bau, III. *Numerical Linear Algebra*. SIAM, Philadelphia, 1997.



G. W. Stewart and Ji-Guang Sun. *Matrix Perturbation Theory*. Academic Press, Boston, 1990.

# Courant-Fischer Theorem

Hermitian  $A = A^H \in \mathbb{C}^{n \times n}$ . Eigenvalues:  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ .

Rayleigh quotient:  $\rho(x) = \frac{x^H A x}{x^H x}$ .

Courant (1920) and Fischer (1905)

$$\lambda_i = \min_{\dim \mathcal{X}=i} \max_{x \in \mathcal{X}} \rho(x), \quad \lambda_i = \max_{\text{codim } \mathcal{X}=i-1} \min_{x \in \mathcal{X}} \rho(x).$$

In particular,

$$\lambda_1 = \min_x \rho(x), \quad \lambda_n = \max_x \rho(x). \quad (1)$$

- Can be used to justify Rayleigh-Ritz approximations for computational purposes.
- (1) is the foundation for using optimization techniques: steepest descent/ascent, CG type, for computing  $\lambda_1$  and, with the help of deflation, other  $\lambda_j$ .

$A = A^H \in \mathbb{C}^{n \times n}$ . Eigenvalues:  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ .

## Trace Min/Trace Max

$$\sum_{j=1}^k \lambda_j = \min_{X^H X = I_k} \text{trace}(X^H A X),$$

$$\sum_{j=1}^k \lambda_{n-j+1} = \max_{X^H X = I_k} \text{trace}(X^H A X).$$

- Can be used to justify Rayleigh-Ritz approximations for computational purposes.
- **Rayleigh quotient matrix:**  $X^H A X$ , assuming  $X^H X = I_k$ .

# Cauchy Interlacing Theorem

Hermitian  $A = A^H \in \mathbb{C}^{n \times n}$ . Eigenvalues:  $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ .

$X \in \mathbb{C}^{n \times k}$ ,  $k \leq n$ ,  $X^H X = I_k$ . Eigenvalue of  $X^H A X$ :  $\mu_1 \leq \mu_2 \leq \dots \leq \mu_k$ .

Cauchy (1829)

$$\lambda_j \leq \mu_j \leq \lambda_{j+n-k} \quad \text{for } 1 \leq j \leq k.$$

Numerical implication: Pick  $X$  to “push” each  $\mu_j$  **down** to  $\lambda_j$  or **up** to  $\lambda_{j+n-k}$ .

# Hermitian $Ax = \lambda Bx$ ( $B \succ 0$ )

$A = A^H, B = B^H \in \mathbb{C}^{n \times n}$ , and  $B$  positive definite.

Equivalency:  $Ax = \lambda Bx \Leftrightarrow \underbrace{B^{-1/2}AB^{-1/2}}_{=: \hat{A}} \hat{x} = \lambda \hat{x}, \hat{x} = B^{1/2}x.$

so same eigenvalues, and eigenvectors related by  $\hat{x} = B^{1/2}x$ .

Eigenvalues  $\lambda_i$  and eigenvectors  $u_i \in \mathbb{C}^n$ .

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n, \quad u_i^H B u_j = \delta_{ij}, \quad A u_i = \lambda_i B u_i.$$

Rayleigh quotient:  $\rho(x) := \frac{x^H A x}{x^H B x} \equiv \frac{\hat{x}^H \hat{A} \hat{x}}{\hat{x}^H \hat{x}}.$

Verbatim translation of theoretical results for  $\hat{A}\hat{x} = \lambda\hat{x}$  to ones for  $Ax = \lambda Bx$ .

# Hermitian $Ax = \lambda Bx$ ( $B \succ 0$ )

Courant (1920) and Fischer (1905)

$$\lambda_i = \min_{\dim \mathcal{X}=i} \max_{x \in \mathcal{X}} \rho(x), \quad \lambda_i = \max_{\text{codim } \mathcal{X}=i-1} \min_{x \in \mathcal{X}} \rho(x).$$

In particular,  $\lambda_1 = \min_x \rho(x)$ , and  $\lambda_n = \max_x \rho(x)$ .

Trace Min/Trace Max

$$\sum_{j=1}^k \lambda_j = \min_{X^H B X = I_k} \text{trace}(X^H A X), \quad \sum_{j=1}^k \lambda_{n-j+1} = \max_{X^H B X = I_k} \text{trace}(X^H A X).$$

Cauchy (1829)

$X \in \mathbb{C}^{n \times k}$ ,  $k \leq n$ ,  $\text{rank}(X) = k$ . Eigenvalues of  $X^H A X - \lambda X^H B X$ :

$$\mu_1 \leq \mu_2 \leq \dots \leq \mu_k.$$

$$\lambda_j \leq \mu_j \leq \lambda_{j+n-k} \quad \text{for } 1 \leq j \leq k.$$

# Why Rayleigh-Ritz?

Two most important aspects in solving large scale eigenvalue problems:

- 1 building subspaces close to the desired eigenvectors (or invariant subspaces).  
E.g., Krylov subspaces.
- 2 seeking “*best possible*” approximations from the suitably built subspaces.

For 2nd aspect: given  $\mathcal{Y} \in \mathbb{C}^n$  and  $\dim \mathcal{Y} = m$ , find the “*best possible*” approximations to some of the eigenvalues of  $A - \lambda B$  using  $\mathcal{Y}$ .

Usually done by Rayleigh-Ritz Procedure. Let  $Y$  be  $\mathcal{Y}$ 's basis matrix.

## Rayleigh-Ritz Procedure

- 1 Solve the eigenvalue problem for  $Y^H A Y - \lambda Y^H B Y$ :  $Y^H A Y y_i = \mu_i Y^H B Y y_i$ ;
- 2 Approximate eigenvalues (Ritz values):  $\mu_i (\approx \lambda_i)$ ;  
approximate eigenvectors (Ritz vectors):  $Y y_i$ .

But in what sense and why are those approximations “*best possible*”?

**Courant-Fischer:**  $\lambda_i = \min_{\dim \mathcal{X}=i} \max_{x \in \mathcal{X}} \rho(x)$  suggests that best possible approximation to  $\lambda_i$  should be taken as

$$\mu_i = \min_{\mathcal{X} \subset \mathcal{Y}, \dim \mathcal{X}=i} \max_{x \in \mathcal{X}} \rho(x)$$

which is the  $i$ th eigenvalue of  $Y^H A Y - \lambda Y^H B Y$ .

**Trace min principle:**  $\sum_{j=1}^k \lambda_j = \min_{X^H B X = I_k} \text{trace}(X^H A X)$  suggests that best possible approximations to  $\lambda_i$  ( $1 \leq i \leq k$ ) should be gotten so that

$\text{trace}(X^H A X)$  is minimized, subject to  $\text{span}(X) \subset \mathcal{Y}$ ,  $X^H B X = I_k$ .

The optimal value is the sum of 1st  $k$  eigenvalues  $\mu_i$  of  $Y^H A Y - \lambda Y^H B Y$ . Consequently,  $\mu_i \approx \lambda_i$  are “best possible”.

# Steepest Descent Methods

- Standard Steepest Descent Method
- Extended Steepest Descent Method
- Convergence Analysis
- Preconditioning Techniques
- Deflation

# Problem: Hermitian $Ax = \lambda Bx$ ( $B \succ 0$ )

$A = A^H, B = B^H \in \mathbb{C}^{n \times n}$ , and  $B$  positive definite.

Eigenvalues  $\lambda_i$  and eigenvectors  $u_i \in \mathbb{C}^n$ .

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n, \quad u_i^H B u_j = \delta_{ij}, \quad A u_i = \lambda_i B u_i.$$

Rayleigh quotient:  $\rho(x) := \frac{x^H A x}{x^H B x}$ .

Interested in computing 1st eigenpair  $(\lambda_1, u_1)$ . Later: Other eigenpairs with the help of deflation.

Largest eigenpairs: through considering  $(-A) - \lambda B$  instead.

SD method: a general technique to solve  $\min f(x)$ .

Steepest descent direction: at given  $x_0$ , along which direction  $p$ ,  $f$  decreases fastest?

$$\min_p \left. \frac{d}{dt} f(x_0 + tp) \right|_{t=0} = \min_p p^T \nabla f(x_0) = -\|\nabla f(x_0)\|_2 \quad (2)$$

optimal  $p$  is in the opposite direction of the gradient  $\nabla f(x_0)$ .

Plain SD: Given  $x_0$ , for  $i = 0, 1, \dots$  until convergence

$$t_i = \arg \min_t f(x_i + t \nabla f(x_i)), \quad x_{i+1} = x_i + t_i \nabla f(x_i). \quad (3)$$

Major work: solve  $\min_t f(x_i + tp)$ , so-called *line search*.

**Food for thought.** Derivation in (2) not quite right for real-valued function  $f$  of complex vector  $x$ .  
In (3):  $t \in \mathbb{R}$  or  $t \in \mathbb{C}$  makes difference.  $t \in \mathbb{C}$  potentially much more complicated!

# Application to $\rho(x) = x^H Ax / x^H Bx$

Recall  $\lambda_1 = \min_x \rho(x)$ .

Gradient:  $\nabla \rho(x) = \frac{2}{x^H Bx} [Ax - \rho(x)Bx] =: \frac{2}{x^H Bx} r(x)$ . Note:  $x^H r(x) \equiv 0$ .

$\|q\|$  tiny, up to 1st order:

$$\begin{aligned}\rho(x+q) &= \frac{(x+q)^H A(x+q)}{(x+q)^H B(x+q)} = \frac{x^H Ax + q^H Ax + x^H Aq}{x^H Bx + q^H Bx + x^H Bq} \\ &= \frac{x^H Ax + q^H Ax + x^H Aq}{x^H Bx} \cdot \left[ 1 - \frac{q^H Bx + x^H Bq}{x^H Bx} \right] = \rho(x) + \frac{q^H r(x) + r(x)^H q}{x^H Bx}.\end{aligned}$$

Steepest descent direction:  $\nabla \rho(x)$  parallel to *residual*  $r(x) = A - \rho(x)Bx$ .

Plain SD: Given  $x_0$ , for  $i = 0, 1, \dots$  until convergence

$$t_i = \arg \inf_t \rho(x_i + t r(x_i)), \quad x_{i+1} = x_i + t_i r(x_i).$$

- Major work: solve  $\inf_t \rho(x_i + t r(x_i))$ , so-called *line search*.
- When to stop?

# Line Search $\inf_{t \in \mathbb{C}} \rho(x_i + t p)$

Can show  $\inf_{t \in \mathbb{C}} \rho(x + t p) = \min_{|\xi|^2 + |\zeta|^2 > 0} \rho(\xi x + \zeta p)$

which is smaller eigenvalue  $\mu$  of  $2 \times 2$  pencil  $X^H A X - \lambda X^H B X$ , where  $X = [x, p]$ .

Let  $v = \begin{bmatrix} \nu_1 \\ \nu_2 \end{bmatrix}$  be the eigenvector. Then  $\rho(Xv) = \mu$ , and  $Xv = \nu_1 x + \nu_2 p$ . So

$$\arg \inf_{t \in \mathbb{C}} \rho(x + t p) =: t_{\text{opt}} = \begin{cases} \nu_2 / \nu_1, & \text{if } \nu_1 \neq 0, \\ \infty, & \text{if } \nu_1 = 0. \end{cases}$$

Interpret  $t_{\text{opt}} = \infty$  in the sense  $\lim_{t \rightarrow \infty} \rho(x + t p) = \rho(p)$ .

$$\rho(y) = \inf_{t \in \mathbb{C}} \rho(x + t p), \quad y = \begin{cases} x + t_{\text{opt}} p & \text{if } t_{\text{opt}} \text{ is finite,} \\ p & \text{otherwise} \end{cases}$$

# A Theorem for Line Search

## Line Search

Suppose  $x, p$  are linearly independent. Then  $x^H r(y) = 0$  and  $p^H r(y) = 0$ .

## Proof

$p^H r(y) = 0$ : True if  $y = p$ , i.e.,  $t_{\text{opt}} = \infty$ . Otherwise

$$y = x + t_{\text{opt}}p, \quad \rho(y) = \min_{t \in \mathbb{C}} \rho(x + tp) = \min_{s \in \mathbb{C}} \rho(y + sp).$$

Optimal at  $s = 0$ .  $\rho(y + sp) = \rho(y) + \frac{2}{y^H B y} \Re(sp^H r(y)) + \mathcal{O}(s^2)$  implies  $p^H r(y) = 0$ .

$x^H r(y) = 0$ : True if  $y = x$ ,  $t_{\text{opt}} = 0$ . and thus  $x^H r(y) = x^H r(x) = 0$ . Otherwise

$$y = \arg \inf_{s \in \mathbb{C}} \rho(p + sx).$$

Therefore  $x^H r(y) = 0$ .

# Stopping Criteria

Common one: check if  $\|r(\mathbf{x})\|$  tiny enough. Reason: Easy to use and available.

$$\text{if } \frac{\|r(\mathbf{x})\|_2}{\|A\mathbf{x}\|_2 + |\rho(\mathbf{x})| \|B\mathbf{x}\|_2} \leq \text{rtol}.$$

Implication:  $(\rho(\mathbf{x}), \mathbf{x})$  is an exact eigenpair of  $(A + E) - \lambda B$  for some Hermitian matrix  $E$ .

Can prove that (suppose  $\|\mathbf{x}\|_2 = 1$ )

$$\min \|E\|_2 = \|r(\mathbf{x})\|_2, \quad \min \|E\|_F = \sqrt{2} \|r(\mathbf{x})\|_2.$$

More can be found in Chapter 5 of:



Zhaojun Bai, J. Demmel, J. Dongarra, A. Ruhe, and H. van der Vorst (editors).  
*Templates for the solution of Algebraic Eigenvalue Problems: A Practical Guide.*  
SIAM, Philadelphia, 2000.

## Steepest Descent method

Given an initial approximation  $\mathbf{x}_0$  to  $u_1$ , and a relative tolerance `rto1`, the algorithm attempts to compute an approximate eigenpair to  $(\lambda_1, u_1)$  with the prescribed `rto1`.

```
1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|_B$ ,  $\rho_0 = \mathbf{x}_0^H A \mathbf{x}_0$ ,  $\mathbf{r}_0 = A \mathbf{x}_0 - \rho_0 B \mathbf{x}_0$ ;  
2: for  $\ell = 0, 1, \dots$  do  
3:   if  $\|\mathbf{r}_\ell\|_2 / (\|A \mathbf{x}_\ell\|_2 + |\rho_\ell| \|B \mathbf{x}_\ell\|_2) \leq \text{rto1}$  then  
4:     BREAK;  
5:   else  
6:     compute the smaller eigenvalue  $\mu$  and corresponding eigenvector  $v$  of  
        $Z^H (A - \lambda B) Z$ , where  $Z = [\mathbf{x}_\ell, \mathbf{r}_\ell]$ ;  
7:      $\hat{\mathbf{x}} = Z v$ ,  $\mathbf{x}_{\ell+1} = \hat{\mathbf{x}} / \|\hat{\mathbf{x}}\|_B$ ;  
8:      $\rho_{\ell+1} = \mu$ ,  $\mathbf{r}_{\ell+1} = A \mathbf{x}_{\ell+1} - \rho_{\ell+1} B \mathbf{x}_{\ell+1}$ ;  
9:   end if  
10: end for  
11: return  $(\rho_\ell, \mathbf{x}_\ell)$  as an approximate eigenpair to  $(\lambda_1, u_1)$ .
```

Note: At Line 6,  $\text{rank}(Z) = 2$  always unless  $\mathbf{r}_\ell = 0$  because  $\mathbf{x}_\ell^H \mathbf{r}_\ell = 0$ .

Pros: Easy to implement; low memory requirement.

Cons: Possibly slow to converge, sometimes unbearably slow.

Well-known: SD slowly moves in zigzag towards an optimal point when the contours near the point are extremely flat.

Ways to rescue:

- Extended the search space: “line search” to “subspace search”
- Modify the search direction: move away from the steepest descent direction  $-\nabla\rho(x)$
- Combination

Seek to extend the search space *naturally*.

SD search space  $\text{span}\{x, r(x)\}$ . Note  $r(x) = Ax - \rho(x)Bx = [A - \rho(x)B]x$ .

$$\text{span}\{x, r(x)\} = \text{span}\{x, [A - \rho(x)B]x\} = \mathcal{K}_2([A - \rho(x)B], x)$$

the *2nd Krylov subspace* of  $A - \rho(x)B$  on  $x$ .

*Naturally* extend  $\mathcal{K}_2([A - \rho(x)B], x)$  to

$$\mathcal{K}_m([A - \rho(x)B], x) = \text{span}\{x, [A - \rho(x)B]x, \dots, [A - \rho(x)B]^{m-1}x\},$$

the *mth Krylov subspace* of  $A - \rho(x)B$  on  $x$ .

Call resulting method *extended steepest descent method* (ESD). It is in fact the inverse free Krylov subspace method of Golub and Ye (2002).

## Extended Steepest Descent method

Given an initial approximation  $\mathbf{x}_0$  to  $u_1$ , a relative tolerance `rtol`, and an integer  $m \geq 2$ , the algorithm attempts to compute an approximate eigenpair to  $(\lambda_1, u_1)$  with the prescribed `rtol`.

```
1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|_B$ ,  $\rho_0 = \mathbf{x}_0^H A \mathbf{x}_0$ ,  $\mathbf{r}_0 = A \mathbf{x}_0 - \rho_0 B \mathbf{x}_0$ ;  
2: for  $\ell = 0, 1, \dots$  do  
3:   if  $\|\mathbf{r}_\ell\|_2 / (\|A \mathbf{x}_\ell\|_2 + |\rho_\ell| \|B \mathbf{x}_\ell\|_2) \leq \text{rtol}$  then  
4:     BREAK;  
5:   else  
6:     compute a basis matrix  $Z \in \mathbb{C}^{n \times m}$  of Krylov subspace  $\mathcal{K}_m(A - \rho_\ell B, \mathbf{x}_\ell)$ ;  
7:     compute the smallest eigenvalue  $\mu$  and corresponding eigenvector  $v$  of  
        $Z^H(A - \lambda B)Z$ ;  
8:      $y = Zv$ ,  $\mathbf{x}_{\ell+1} = y / \|y\|_B$ ;  
9:      $\rho_{\ell+1} = \mu$ ,  $\mathbf{r}_{\ell+1} = A \mathbf{x}_{\ell+1} - \rho_{\ell+1} B \mathbf{x}_{\ell+1}$ ;  
10:   end if  
11: end for  
12: return  $(\rho_\ell, \mathbf{x}_\ell)$  as an approximate eigenpair to  $(\lambda_1, u_1)$ .
```

Note: If  $B = I$ , it is equivalent to *Restarted Lanczos*

# Basis of $\mathcal{K}_m(A - \rho B, \mathbf{x})$

$C = A - \rho B$  is Hermitian.

## Lanczos process

```
1:  $z_1 = \mathbf{x} / \|\mathbf{x}\|_2$ ,  $\beta_0 = 0$ ;  $z_0 = 0$ ;  
2: for  $j = 1, 2, \dots, k$  do  
3:    $z = Cz_j$ ,  $\alpha_j = z_j^H z$ ;  
4:    $z = z - \alpha_j z_j - \beta_{j-1} z_{j-1}$ ,  $\beta_j = \|z\|_2$ ;  
5:   if  $\beta_j = 0$  then  
6:     BREAK;  
7:   else  
8:      $z_{j+1} = z / \beta_j$ ;  
9:   end if  
10: end for
```

- Keep  $Az_j$  and  $Bz_j$  for projecting  $A$  and  $B$  later. Or, just  $z_j$  but solve  $Z^H CZ - \lambda Z^H BZ$  instead
- Implemented as is,  $Z = [z_1, \dots, z_m]$  may lose orthogonality — partial or full re-orthogonalization should be used. Pose little problem since usually  $m$  is modest.
- Possibly  $\dim \mathcal{K}_m(A - \rho B, \mathbf{x}) < m$ . Pose no problem —  $Z$  has fewer than  $m$  columns

## Convergence

For SD and ESD,  $\rho_\ell$  converges to some eigenvalue  $\hat{\lambda}$  of  $A - \lambda B$  and  $\|(A - \hat{\lambda}B)\mathbf{x}_\ell\|_2 \rightarrow 0$ .

## Proof.

- 1)  $\{\rho_\ell\}$  monotonically decreasing and  $\rho_\ell \geq \lambda_1 \Rightarrow \rho_\ell \rightarrow \hat{\lambda}$ .
- 2)  $\{\mathbf{x}_\ell\}$  bounded in  $\mathbb{C}^n \Rightarrow$  convergent  $\{\mathbf{x}_{n_\ell}\}, \mathbf{x}_{n_\ell} \rightarrow \hat{\mathbf{x}}$ .
- 3)  $\mathbf{x}_{n_\ell}^H (A - \rho_{n_\ell} B) \mathbf{x}_{n_\ell} = 0 \Rightarrow \hat{\mathbf{x}}^H \hat{\mathbf{r}} = \hat{\mathbf{x}}^H (A - \hat{\lambda} B) \hat{\mathbf{x}} = 0$ .
- 4) Claim  $\hat{\mathbf{r}} = 0$ . Otherwise  $\hat{\mathbf{r}} \neq 0$ ,  $\text{rank}([\hat{\mathbf{x}}, \hat{\mathbf{r}}]) = 2$ , and

$$\hat{A} - \hat{\lambda} \hat{B} := \begin{bmatrix} \hat{\mathbf{x}}^H \\ \hat{\mathbf{r}}^H \end{bmatrix} A[\hat{\mathbf{x}}, \hat{\mathbf{r}}] - \hat{\lambda} \begin{bmatrix} \hat{\mathbf{x}}^H \\ \hat{\mathbf{r}}^H \end{bmatrix} B[\hat{\mathbf{x}}, \hat{\mathbf{r}}] = \begin{bmatrix} 0 & \hat{\mathbf{r}}^H \hat{\mathbf{r}} \\ \hat{\mathbf{r}}^H \hat{\mathbf{r}} & \hat{\mathbf{r}}^H (A - \hat{\lambda} B) \hat{\mathbf{r}} \end{bmatrix} \text{ is indefinite.}$$

Smaller eigenvalue  $\mu$  of  $\hat{A} - \lambda \hat{B}$ :  $\mu < \hat{\lambda}$ . Let

$$\hat{A}_\ell = [\mathbf{x}_\ell, \mathbf{r}_\ell]^H A[\mathbf{x}_\ell, \mathbf{r}_\ell], \quad \hat{B}_\ell = [\mathbf{x}_\ell, \mathbf{r}_\ell]^H B[\mathbf{x}_\ell, \mathbf{r}_\ell],$$

$\mu_{\ell+1}$  smaller eigenvalue of  $\hat{A}_\ell - \lambda \hat{B}_\ell$ . Then

$$(i) \hat{A}_{n_\ell} \rightarrow \hat{A}, \quad \hat{B}_{n_\ell} \rightarrow \hat{B} \Rightarrow \mu_{n_\ell+1} \rightarrow \mu,$$

$$(ii) \rho_{n_\ell+1} \leq \mu_{n_\ell+1} \Rightarrow \hat{\lambda} = \lim_{i \rightarrow \infty} \rho_{n_\ell+1} \leq \lim_{i \rightarrow \infty} \mu_{n_\ell+1} = \mu,$$

a contradiction! So  $\hat{\mathbf{r}} = 0$ . □

# Digression: Chebyshev Polynomial

The  $j$ th Chebyshev polynomial of the first kind  $\mathcal{T}_j(t)$

$$\begin{aligned}\mathcal{T}_j(t) &= \cos(j \arccos t) && \text{for } |t| \leq 1, \\ &= \frac{1}{2} \left[ \left( t + \sqrt{t^2 - 1} \right)^j + \left( t + \sqrt{t^2 - 1} \right)^{-j} \right] && \text{for } t \geq 1.\end{aligned}$$

Or,  $\mathcal{T}_0(t) = 1$ ,  $\mathcal{T}_1(t) = t$ , and  $\mathcal{T}_j(t) = 2\mathcal{T}_{j-1}(t) - \mathcal{T}_{j-2}(t)$  for  $j \geq 2$ .

Numerous optimal properties among polynomials

- $\deg(p) \leq j$ ,  $|p(t)| \leq 1$  for  $t \in [-1, 1] \Rightarrow |p(t)| \leq |\mathcal{T}_j(t)|$  for  $t \notin [-1, 1]$ .

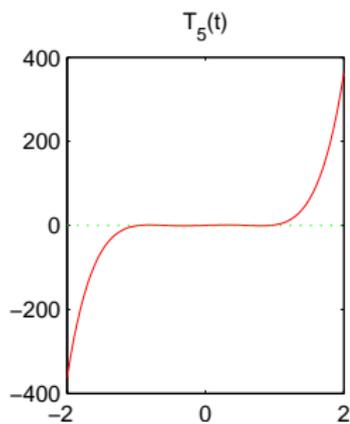
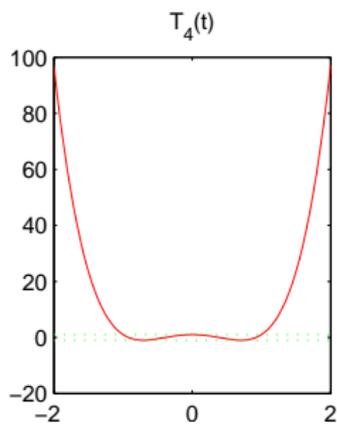
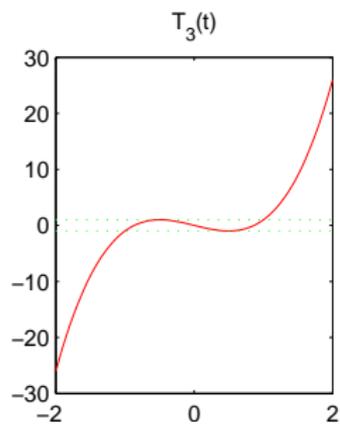
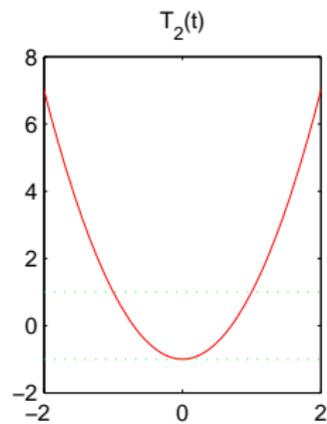
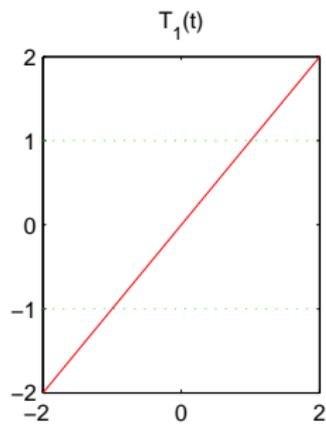
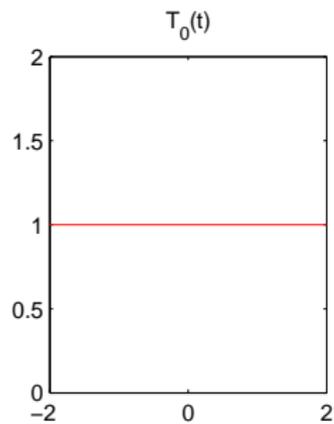
I.e.,  $|\mathcal{T}_j(t)| \leq 1$  for  $|t| \leq 1$  and  $|\mathcal{T}_j(t)|$  grows fastest. (Sample plots next slide.)

$$\left| \mathcal{T}_j \left( \frac{1+t}{1-t} \right) \right| = \left| \mathcal{T}_j \left( \frac{t+1}{t-1} \right) \right| = \frac{1}{2} \left[ \Delta_t^j + \Delta_t^{-j} \right] \quad \text{for } 1 \neq t > 0,$$

where  $\Delta_t := \frac{\sqrt{t} + 1}{|\sqrt{t} - 1|}$  for  $t > 0$ .

Frequently show up in numerical analysis and computations: Chebyshev acceleration in iterative methods, convergence of CG and Lanczos methods.

# Chebyshev Polynomial (sample plots)



# Chebyshev Polynomial: typical use

**Problem** Given  $[\alpha, \beta]$  and  $\gamma \notin [\alpha, \beta]$ , seek a polynomial  $p$  with  $\deg(p) \leq m$  such that  $p(\gamma) = 1$  and  $\max_{x \in [\alpha, \beta]} |p(x)|$  is minimized.

Define 1-1 mapping  $x \in [\alpha, \beta] \rightarrow t \equiv t(x) := \frac{2}{\beta - \alpha} \left( x - \frac{\alpha + \beta}{2} \right) \in [-1, 1]$ ,

$$t(\gamma) = -\frac{1 + \frac{\alpha - \gamma}{\beta - \gamma}}{1 - \frac{\alpha - \gamma}{\beta - \gamma}} \quad \text{for } \gamma < \alpha, \quad \text{and} \quad \frac{1 + \frac{\gamma - \beta}{\gamma - \alpha}}{1 - \frac{\gamma - \beta}{\gamma - \alpha}} \quad \text{for } \beta < \gamma.$$

Optimal  $p(x) = \frac{\mathcal{T}_m(t(x))}{\mathcal{T}_m(t(\gamma))}$ :

$$\rho(\gamma) = 1, \quad \max_{x \in [\alpha, \beta]} |p(x)| = \frac{1}{|\mathcal{T}_m(t(\gamma))|} = 2 \left[ \Delta_\eta^j + \Delta_\eta^{-j} \right]^{-1},$$
$$\Delta_\eta = \frac{1 + \sqrt{\eta}}{1 - \sqrt{\eta}}, \quad \eta = \frac{\alpha - \gamma}{\beta - \gamma} \quad \text{for } \gamma < \alpha, \quad \text{and} \quad \eta = \frac{\gamma - \beta}{\gamma - \alpha} \quad \text{for } \beta < \gamma.$$

## Theorem on Convergence Rate (Golub & Ye, 2002)

Suppose  $\lambda_1$  is simple, i.e.,  $\lambda_1 < \lambda_2$ , and  $\lambda_1 < \rho_\ell < \lambda_2$ . Let

$$\omega_1 < \omega_2 \leq \dots \leq \omega_n$$

be the eigenvalues of  $A - \rho_\ell B$  and  $v_1$  be an eigenvector corresponding to  $\omega_1$ . Then

$$\rho_{\ell+1} - \lambda_1 \leq (\rho_\ell - \lambda_1)\epsilon_m^2 + 2(\rho_\ell - \lambda_1)^{3/2}\epsilon_m \left( \frac{\|B\|_2}{\omega_2} \right)^{1/2} + \delta_\ell$$

where

$$0 \leq \delta_\ell := \rho_\ell - \lambda_1 + \omega_1 \frac{v_1^H v_1}{v_1^H B v_1} = \mathcal{O}(|\rho_\ell - \lambda_1|^2),$$

$$\epsilon_m := \min_{f \in \mathbb{P}_{m-1}, f(\omega_1)=1} \max_{j>1} |f(\omega_j)|.$$

$\epsilon_m := \min_{f \in \mathbb{P}_{m-1}, f(\omega_1)=1} \max_{j>1} |f(\omega_j)|$  usually unknown, except,

**1**  $m = 2$  for which the optimal  $f_{\text{opt}}$  is

$$f_{\text{opt}}(t) = \frac{t - (\omega_2 + \omega_n)/2}{\omega_1 - (\omega_2 + \omega_n)/2} \in \mathbb{P}_{m-1}, \quad f_{\text{opt}}(\omega_1) = 1, \quad \max_{j>1} |f_{\text{opt}}(\omega_j)| = |f_{\text{opt}}(\omega_2)| < 1.,$$

$$\epsilon_2 = \frac{1 - \eta}{1 + \eta}, \quad \eta = \frac{\omega_2 - \omega_1}{\omega_n - \omega_1}.$$

**2**  $\omega_2 = \dots = \omega_n$  for which  $f_{\text{opt}}(t) = (t - \omega_2)/(\omega_1 - \omega_2)$  and  $\epsilon_m = 0$  for all  $m \geq 2$ .

In general  $\epsilon_m$  can be bounded by using the Chebyshev polynomial

$$f(t) = \mathcal{T}_{m-1} \left( \frac{2t - (\omega_n + \omega_2)}{\omega_n - \omega_2} \right) / \mathcal{T}_{m-1} \left( \frac{1 + \eta}{1 - \eta} \right), \quad f(\omega_1) = 1,$$

$$\epsilon_m \leq \max_{\omega_2 \leq t \leq \omega_n} |f(t)| = \left[ \mathcal{T}_{m-1} \left( \frac{1 + \eta}{1 - \eta} \right) \right]^{-1} = 2 \left[ \Delta_\eta^{m-1} + \Delta_\eta^{-(m-1)} \right]^{-1}.$$

$$\rho_{\ell+1} - \lambda_1 \leq (\rho_\ell - \lambda_1)\epsilon_m^2 + 2(\rho_\ell - \lambda_1)^{3/2}\epsilon_m \left( \frac{\|B\|_2}{\omega_2} \right)^{1/2} + \mathcal{O}(|\rho_\ell - \lambda_1|^2).$$

Ignoring high order terms,  $\frac{\rho_{\ell+1} - \lambda_1}{\rho_\ell - \lambda_1} \lesssim \epsilon_m^2$ .

If  $\epsilon_m = 0$  (unlikely, however), then  $\rho_{\ell+1} - \lambda_1 = \mathcal{O}(|\rho_\ell - \lambda_1|^2)$ , quadratic convergence.

Locally,  $\rho_\ell \approx \lambda_1$ ,  $\text{eig}(A - \rho_\ell B) \approx \text{eig}(A - \lambda_1 B) = \{0 = \gamma_1 < \gamma_2 \leq \dots \leq \gamma_n\}$ . So

$$\epsilon_m \approx \min_{f \in \mathbb{P}_{m-1}, f(\gamma_1)=1} \max_{j>1} |f(\gamma_j)| \leq 2 \left[ \Delta_\eta^{m-1} + \Delta_\eta^{-(m-1)} \right]^{-1},$$

$$\Delta_\eta = \frac{1 + \sqrt{\eta}}{1 - \sqrt{\eta}}, \quad \eta = \frac{\gamma_2 - \gamma_1}{\gamma_n - \gamma_1}.$$

**Observation.**  $\epsilon_m$  depends on  $\text{eig}(A - \rho_\ell B)$ , not on  $\text{eig}(A, B)$ . This is the Key for preconditioning later: transforming  $A - \lambda B$  to preserve  $\text{eig}(A, B)$  but make  $\text{eig}(A - \rho_\ell B)$  more preferable.

Lemma (Golub & Ye, 2002)

Let  $(\omega_1, v_1)$  be the smallest eigenvalue of  $A - \rho_\ell B$ , i.e.. Then

$$-\omega_1 \frac{v_1^H v_1}{v_1^H B v_1} \leq \rho_\ell - \lambda_1 \leq -\omega_1 \frac{u_1^H u_1}{u_1^H B u_1}. \quad (4)$$

Asymptotically, if  $\lambda_1$  is a simple eigenvalue of  $A - \lambda B$ , then as  $\rho_\ell \rightarrow \lambda_1$ ,

$$\omega_1 \frac{v_1^H v_1}{v_1^H B v_1} = (\lambda_1 - \rho_\ell) + \mathcal{O}(|\lambda_1 - \rho_\ell|^2). \quad (5)$$

**Importance.** Relate  $\lambda_1 - \rho_\ell$  to  $\omega_1$ . Special case:  $B = I$ ,  $\omega_1 = \lambda_1 - \rho_\ell$ .

# Proof of Key Lemma

- 1)  $\rho_\ell \geq \lambda_1$  always. If  $\rho_\ell = \lambda_1$ , then  $\omega_1 = 0$ . No proof needed.
- 2) Suppose  $\rho_\ell > \lambda_1$ .  $A - \rho_\ell B$  is indefinite and hence  $\omega_1 < 0$ . We have

$$(A - \rho_\ell B)v_1 = \omega_1 v_1, \quad (A - \omega_1 I - \rho_\ell B)v_1 = 0, \quad A - \omega_1 I - \rho_\ell B \succeq 0.$$

Therefore  $(\rho_\ell, v_1)$  is the smallest eigenpair of  $(A - \omega_1 I) - \lambda B$ .

- 3) Note also  $(\lambda_1, u_1)$  is the smallest eigenpair of  $A - \lambda B$ .

$$\begin{aligned}\rho_\ell &= \frac{v_1^H (A - \omega_1 I) v_1}{v_1^H B v_1} = \frac{v_1^H A v_1}{v_1^H B v_1} + \frac{-\omega_1 v_1^H v_1}{v_1^H B v_1} \\ &\geq \min_x \frac{x^H A x}{x^H B x} + \frac{-\omega_1 v_1^H v_1}{v_1^H B v_1} = \lambda_1 + \frac{-\omega_1 v_1^H v_1}{v_1^H B v_1}, \\ \lambda_1 &= \frac{u_1^H A u_1}{u_1^H B u_1} = \frac{u_1^H (A - \rho_\ell B) u_1}{u_1^H u_1} \cdot \frac{u_1^H u_1}{u_1^H B u_1} + \rho_\ell \\ &\geq \min_x \frac{x^H (A - \rho_\ell B) x}{x^H x} \cdot \frac{u_1^H u_1}{u_1^H B u_1} + \rho_\ell = \omega_1 \cdot \frac{u_1^H u_1}{u_1^H B u_1} + \rho_\ell.\end{aligned}$$

Together yielding 
$$-\omega_1 \frac{v_1^H v_1}{v_1^H B v_1} \leq \rho_\ell - \lambda_1 \leq -\omega_1 \frac{u_1^H u_1}{u_1^H B u_1}.$$

# Proof of Key Lemma

- 4)  $\omega_1(t) = \lambda_{\min}(A - tB)$ , for  $t$  near  $\lambda_1$ . Then  $\omega_1(\lambda_1) = 0$  and  $\omega_1(\boldsymbol{\rho}_\ell) = \omega_1$ .
- 5)  $\omega_1(\lambda_1) = 0$  is a simple eigenvalue of  $A - \lambda_1 B$ . So  $\omega_1(t)$  is differentiable in a neighborhood of  $\lambda_1$ .
- 6) Expand  $\omega_1(t)$  at  $\boldsymbol{\rho}_\ell$ , sufficiently close to  $\lambda_1$ . Can prove  $\omega_1'(\boldsymbol{\rho}_\ell) = -\frac{v_1^H B v_1}{v_1^H v_1}$ . Hence

$$\begin{aligned}\omega_1(t) &= \omega_1(\boldsymbol{\rho}_\ell) + \sigma_1'(\boldsymbol{\rho}_\ell)(t - \boldsymbol{\rho}_\ell) + \mathcal{O}(|t - \boldsymbol{\rho}_\ell|^2) \\ &= \omega_1 - \frac{v_1^H B v_1}{v_1^H v_1}(t - \boldsymbol{\rho}_\ell) + \mathcal{O}(|t - \boldsymbol{\rho}_\ell|^2).\end{aligned}$$

Setting  $t = \lambda_1$ ,

$$0 = \omega_1(\lambda_1) = \omega_1 - \frac{v_1^H B v_1}{v_1^H v_1}(\lambda_1 - \boldsymbol{\rho}_\ell) + \mathcal{O}(|\lambda_1 - \boldsymbol{\rho}_\ell|^2),$$

from which  $\omega_1 \frac{v_1^H v_1}{v_1^H B v_1} = (\lambda_1 - \boldsymbol{\rho}_\ell) + \mathcal{O}(|\lambda_1 - \boldsymbol{\rho}_\ell|^2)$ . □

## Theorem on Convergence Rate (Golub & Ye, 2002)

Suppose  $\lambda_1$  is simple, i.e.,  $\lambda_1 < \lambda_2$ , and  $\lambda_1 < \rho_\ell < \lambda_2$ . Let

$$\omega_1 < \omega_2 \leq \dots \leq \omega_n$$

be the eigenvalues of  $A - \rho_\ell B$  and  $v_1$  be an eigenvector corresponding to  $\omega_1$ . Then

$$\rho_{\ell+1} - \lambda_1 \leq (\rho_\ell - \lambda_1)\epsilon_m^2 + 2(\rho_\ell - \lambda_1)^{3/2}\epsilon_m \left( \frac{\|B\|_2}{\omega_2} \right)^{1/2} + \delta_\ell$$

where

$$0 \leq \delta_\ell := \rho_\ell - \lambda_1 + \omega_1 \frac{v_1^H v_1}{v_1^H B v_1} = \mathcal{O}(|\rho_\ell - \lambda_1|^2),$$

$$\epsilon_m := \min_{f \in \mathbb{P}_{m-1}, f(\omega_1)=1} \max_{j>1} |f(\omega_j)|.$$

- 1)  $C = A - \rho_\ell B = V\Omega V^H$  (eigen-decomposition),  $V = [v_1, v_2, \dots, v_n]$  (orthogonal),  $\Omega = \text{diag}(\omega_1, \omega_2, \dots, \omega_n)$ .
- 2)  $\mathcal{K}_m \equiv \mathcal{K}_m(C, \mathbf{x}_\ell) = \{f(C)\mathbf{x}_\ell, f \in \mathbb{P}_{m-1}\}$ , and

$$\begin{aligned} \rho_{\ell+1} &= \min_{x \in \mathcal{K}_m} \frac{x^H A x}{x^H B x} = \rho_\ell + \min_{x \in \mathcal{K}_m} \frac{x^H (A - \rho_\ell B) x}{x^H B x} \\ &= \rho_\ell + \min_{f \in \mathbb{P}_{m-1}} \frac{\mathbf{x}_\ell^H f(C) C f(C) \mathbf{x}_\ell}{\mathbf{x}_\ell^H f(C) B f(C) \mathbf{x}_\ell}. \end{aligned} \quad (6)$$

- 3) Let  $f_{\text{opt}} \in \mathbb{P}_{m-1}$  be the minimizing polynomial that defines  $\epsilon_m$ . Then  $f_{\text{opt}}(\omega_1) = 1$  by the definition, and also  $\epsilon_m = \max_{j>1} |f_{\text{opt}}(\omega_j)| < 1$  because

$$f(t) = \frac{t - (\omega_2 + \omega_n)/2}{\omega_1 - (\omega_2 + \omega_n)/2} \in \mathbb{P}_{m-1}, \quad f(\omega_1) = 1, \quad \max_{j>1} |f(\omega_j)| = |f(\omega_2)| < 1.$$

- 4)  $\mathbf{x}_\ell^H (A - \rho_\ell B) \mathbf{x}_\ell = 0 \Rightarrow$  that  $v_1^H \mathbf{x}_\ell \neq 0$  and hence  $f_{\text{opt}}(C) \mathbf{x}_\ell \neq 0$  (why?). Thus

$$\begin{aligned} \rho_{\ell+1} &\leq \rho_\ell + \frac{\mathbf{x}_\ell^H f_{\text{opt}}(C) C f_{\text{opt}}(C) \mathbf{x}_\ell}{\mathbf{x}_\ell^H f_{\text{opt}}(C) B f_{\text{opt}}(C) \mathbf{x}_\ell} = \rho_\ell + \frac{\mathbf{x}_\ell^H V f_{\text{opt}}^2(\Omega) \Omega V^H \mathbf{x}_\ell}{\mathbf{x}_\ell^H V f_{\text{opt}}(\Omega) B_1 f_{\text{opt}}(\Omega) V^H \mathbf{x}_\ell} \\ &= \rho_\ell + \frac{y^H f_{\text{opt}}^2(\Omega) \Omega y}{y^H f_{\text{opt}}(\Omega) B_1 f_{\text{opt}}(\Omega) y}, \quad \text{where } B_1 = V^H B V, \quad y = V^H \mathbf{x}_\ell. \end{aligned}$$

5) Write  $y = [\xi_1, \xi_2, \dots, \xi_n]^T = \xi_1 \mathbf{e}_1 + \hat{y}$ ,  $\hat{y} = [0, \xi_2, \dots, \xi_n]^T$ . We have

$$\begin{aligned} y^H f_{\text{opt}}(\Omega) B_1 f_{\text{opt}}(\Omega) y &= (\xi_1 \mathbf{e}_1 + \hat{y})^H f_{\text{opt}}(\Omega) B_1 f_{\text{opt}}(\Omega) (\xi_1 \mathbf{e}_1 + \hat{y}) \\ &= \xi_1^2 f_{\text{opt}}(\omega_1)^2 \mathbf{e}_1^H B_1 \mathbf{e}_1 + 2\xi_1 f_{\text{opt}}(\omega_1) \mathbf{e}_1^H B_1 f_{\text{opt}}(\Omega) \hat{y} \\ &\quad + \hat{y}^H f_{\text{opt}}(\Omega) B_1 f_{\text{opt}}(\Omega) \hat{y} \\ &= \xi_1^2 \beta_1^2 + 2\xi_1 \beta_2 + \beta_3^2, \end{aligned}$$

where

$$\begin{aligned} \beta_1^2 &= \mathbf{e}_1^H B_1 \mathbf{e}_1 = \mathbf{v}_1^H B \mathbf{v}_1, \\ \beta_3^2 &= \hat{y}^H f_{\text{opt}}(\Omega) B_1 f_{\text{opt}}(\Omega) \hat{y} \leq \max_{j>1} f_{\text{opt}}(\omega_j)^2 \|B_1\|_2 \|\hat{y}\|_2^2 = \epsilon_m^2 \|B\|_2 \|\hat{y}\|_2^2, \\ |\beta_2| &= |\mathbf{e}_1^H B_1 f_{\text{opt}}(\Omega) \hat{y}| \leq \beta_1 \beta_3. \end{aligned}$$

Note  $\sum_j \omega_j \xi_j^2 = y^H \Omega y = \mathbf{x}_\ell^H (A - \rho_\ell B) \mathbf{x}_\ell = 0 \Rightarrow |\omega_1| \xi_1^2 = \sum_{j>1} \omega_j \xi_j^2 \geq \omega_2 \|\hat{y}\|_2^2$ .  
Hence

$$\beta_3 \leq \epsilon_m \|B\|_2^{1/2} \left( \frac{|\omega_1|}{\omega_2} \right)^{1/2} |\xi_1|.$$

6)  $y^H f_{\text{opt}}^2(\Omega)\Omega y = \xi_1^2 \omega_1 + \hat{y}^H f_{\text{opt}}(\Omega)^2 \Omega \hat{y}$ , and

$$y^H f_{\text{opt}}^2(\Omega)\Omega y = \sum_j \omega_j f_{\text{opt}}^2(\omega_j) \xi_\ell^2 \leq \sum_j \omega_j \xi_j^2 = y^H \Omega y = 0,$$

$$\hat{y}^H f_{\text{opt}}^2(\Omega)\Omega \hat{y} = \sum_{j>1} \omega_j f_{\text{opt}}^2(\omega_j) \xi_\ell^2 \leq \epsilon_m^2 \sum_{j>1} \omega_j \xi_j^2 = \epsilon_m^2 |\omega_1| \xi_1^2.$$

$$\begin{aligned} \frac{y^H f_{\text{opt}}^2(\Omega)\Omega y}{y^H f_{\text{opt}}(\Omega) B_1 f_{\text{opt}}(\Omega) y} &\leq \frac{\xi_1^2 \omega_1 + \hat{y}^H f_{\text{opt}}(\Omega)^2 \Omega \hat{y}}{\xi_1^2 \beta_1^2 + 2|\xi_1| \beta_1 \beta_3 + \beta_3^2} \\ &= \frac{\omega_1}{\beta_1^2} - \frac{\omega_1}{\beta_1^2} \cdot \frac{2|\xi_1| \beta_1 \beta_3 + \beta_3^2}{\xi_1^2 \beta_1^2 + 2|\xi_1| \beta_1 \beta_3 + \beta_3^2} + \frac{\hat{y}^H f_{\text{opt}}(\Omega)^2 \Omega \hat{y}}{\xi_1^2 \beta_1^2 + 2|\xi_1| \beta_1 \beta_3 + \beta_3^2} \\ &\leq \frac{\omega_1}{\beta_1^2} - \frac{\omega_1}{\beta_1^2} \cdot \frac{2|\xi_1| \beta_1 \beta_3}{\xi_1^2 \beta_1^2} + \frac{\hat{y}^H f_{\text{opt}}(\Omega)^2 \Omega \hat{y}}{\xi_1^2 \beta_1^2} \\ &\leq \frac{\omega_1}{\beta_1^2} + 2 \left( \frac{|\omega_1|}{\beta_1^2} \right)^{3/2} \epsilon_m \left( \frac{\|B\|_2}{\omega_2} \right)^{1/2} + \frac{|\omega_1|}{\beta_1^2} \epsilon_m^2. \end{aligned}$$

7) We have proved

$$\frac{y^H f_{\text{opt}}^2(\Omega) \Omega y}{y^H f_{\text{opt}}(\Omega) B_1 f_{\text{opt}}(\Omega) y} \leq \frac{\omega_1}{\beta_1^2} + 2 \left( \frac{|\omega_1|}{\beta_1^2} \right)^{3/2} \epsilon_m \left( \frac{\|B\|_2}{\omega_2} \right)^{1/2} + \frac{|\omega_1|}{\beta_1^2} \epsilon_m^2.$$

Note  $\frac{\omega_1}{\beta_1^2} = \omega_1 \frac{v_1^H v_1}{v_1^H B v_1} = (\lambda_1 - \rho_\ell) + \mathcal{O}(|\lambda_1 - \rho_\ell|^2)$  by the lemma. Therefore

$$\begin{aligned} \rho_{\ell+1} - \lambda_1 &\leq \rho_\ell - \lambda_1 + \frac{y^H f_{\text{opt}}^2(\Omega) \Omega y}{y^H f_{\text{opt}}(\Omega) B_1 f_{\text{opt}}(\Omega) y} \\ &\leq \underbrace{\rho_\ell - \lambda_1 + \frac{\omega_1}{\beta_1^2}}_{\delta_\ell} + 2 \left( \frac{|\omega_1|}{\beta_1^2} \right)^{3/2} \epsilon_m \left( \frac{\|B\|_2}{\omega_2} \right)^{1/2} + \frac{|\omega_1|}{\beta_1^2} \epsilon_m^2 \\ &\leq \delta_\ell + 2(\rho_\ell - \lambda_1)^{3/2} \epsilon_m \left( \frac{\|B\|_2}{\omega_2} \right)^{1/2} + (\rho_\ell - \lambda_1) \epsilon_m^2, \end{aligned}$$

as expected. □

# What is Preconditioning?

**Preconditioning.** Transforming a problem that is “easier” (e.g., taking less time) to solve iteratively.

Preconditioning **natural** for linear systems: transform  $Ax = b$  to  $KAx = Kb$  which is “easier” than before. Extreme case:  $KA = I$ , i.e.,  $K = A^{-1}$ , then  $x = Kb$ . But this is impractical!

A compromise: make  $KA \approx I$  as much as practical. Here  $KA \approx I$  is understood either  $\|KA - I\|$  is relatively small or  $KA - I$  is near a low rank matrix.

Preconditioning **not so natural** for eigenvalue problems: transform  $A - \lambda B$  to  $KA - \lambda KB$  or  $L^H A L - \lambda L^H B L$  which is “easier” than before.

- No straightforward explanation as to what  $K$  makes  $KA - \lambda KB$  “easier”
- No straightforward explanation as to what  $L$  makes  $L^H A L - \lambda L^H B L$  “easier”, except  $L$  being the eigenvector matrix that is unknown. No easy way to approximate the unknown eigenvector matrix either.

Will present two ways to understand eigen-problem preconditioning and construct preconditioners.

# Eigen-problem preconditioning, I

Ideal search direction  $p$ : starting at  $\mathbf{x}_\ell$ ,  $p$  points to the optimum, i.e., the optimum is on the line  $\{\mathbf{x}_\ell + t\mathbf{p} : t \in \mathbb{C}\}$ . How can it be done with unknown optimum?

Expand  $\mathbf{x}_\ell$  as a linear combination of  $u_\ell$

$$\mathbf{x}_\ell = \sum_{j=1}^n \alpha_j u_j =: \alpha_1 u_1 + \mathbf{v}, \quad \mathbf{v} = \sum_{j=2}^n \alpha_j u_j \perp_B u_1.$$

Then ideal  $p = \alpha u_1 + \beta \mathbf{v}$ ,  $\beta \neq 0$  such that  $\alpha_1 \beta - \alpha \neq 0$  (otherwise  $p = \beta \mathbf{x}_\ell$ ).

Ideal  $p$  has to be approximated to be practical. One such approximate  $p$  is

$$p = (A - \sigma B)^{-1} \mathbf{r}_\ell = (A - \sigma B)^{-1} [A - \rho_\ell B] \mathbf{x}_\ell,$$

where  $\rho_\ell \neq \sigma \approx \lambda_1$ , also reasonably we assume  $\sigma \neq \lambda_j$  for all  $j > 1$ . Why so?

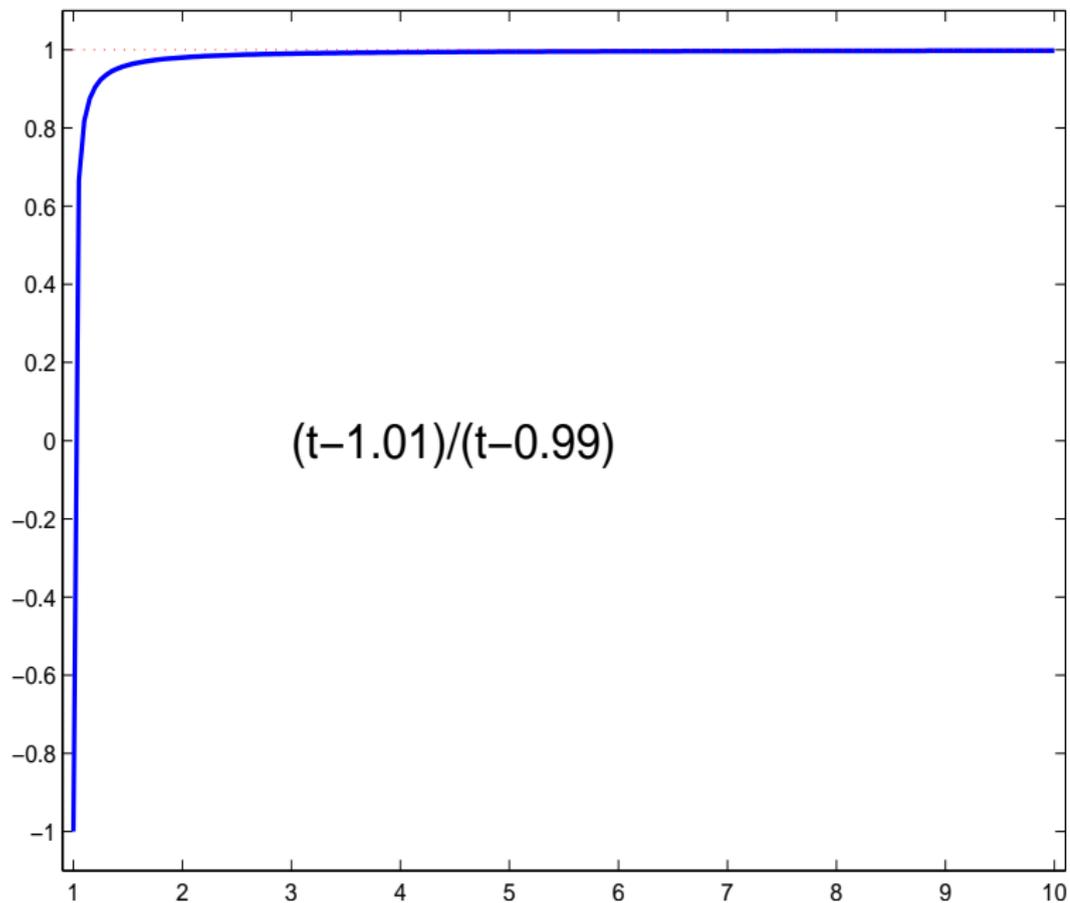
$$p = \sum_{j=1}^n \mu_j \alpha_j u_j, \quad \mu_j := \frac{\lambda_j - \rho_\ell}{\lambda_j - \sigma}.$$

Now if  $\lambda_1 \leq \rho_\ell < \lambda_2$  and if the gap  $\lambda_2 - \lambda_1$  is reasonably modest, then

$$\mu_j \approx 1 \quad \text{for } j > 1$$

to give a  $p \approx \alpha v_1 + \mathbf{v}$ , resulting in fast convergence.

$$\mu_j = (\lambda_j - \rho_\ell) / (\lambda_j - \sigma)$$



Preconditioner  $(A - \sigma B)^{-1}$

Let  $\mathbf{x}_\ell = \sum_{j=1}^n \alpha_j u_j$ , and suppose  $\alpha_1 \neq 0$ . If  $\sigma \neq \rho_\ell$  such that

either  $\mu_1 < \mu_j$  for  $2 \leq j \leq n$  or  $\mu_1 > \mu_j$  for  $2 \leq j \leq n$ ,

where  $\mu_j = \frac{\lambda_j - \rho_\ell}{\lambda_j - \sigma}$ , then

$$\begin{aligned}\tan \theta_B(u_1, \mathcal{K}_m) &\leq 2 \left[ \Delta_\eta^{m-1} + \Delta_\eta^{-(m-1)} \right]^{-1} \tan \theta_B(u_1, \mathbf{x}_\ell), \\ 0 \leq \rho_{\ell+1} - \lambda_1 &\leq 4 \left[ \Delta_\eta^{m-1} + \Delta_\eta^{-(m-1)} \right]^{-2} \tan \theta_B(u_1, \mathbf{x}_\ell),\end{aligned}$$

where  $\mathcal{K}_m := \mathcal{K}_m([A - \sigma B]^{-1}[A - \rho_\ell B], \mathbf{x}_\ell)$ , and

$$\eta = \begin{cases} \frac{\lambda_n - \sigma}{\lambda_n - \lambda_1} \cdot \frac{\lambda_2 - \lambda_1}{\lambda_2 - \sigma}, & \text{if } \mu_1 < \mu_j \text{ for } 2 \leq j \leq n, \\ \frac{\lambda_2 - \sigma}{\lambda_2 - \lambda_1} \cdot \frac{\lambda_n - \lambda_1}{\lambda_n - \sigma}, & \text{if } \mu_1 > \mu_j \text{ for } 2 \leq j \leq n. \end{cases}$$

$\eta \approx 1$  (fast convergence) if  $\sigma \approx \lambda_1$ . In fact  $\eta = 1$  (implying  $\Delta\eta = \infty$ ) if  $\sigma = \lambda_1$ .

But shift  $\sigma$  needs to make  $\mu_1$  either smallest or biggest among all  $\mu_j$ . Three interesting cases:

- $\sigma < \lambda_1 \leq \rho < \lambda_2$ ,  $\mu_1$  smallest
- $\lambda_1 < \sigma < \rho < \lambda_2$ ,  $\mu_1$  biggest
- $\lambda_1 < \rho < \sigma < \lambda_2$ ,  $\mu_1$  smallest.

$(A - \sigma B)^{-1}$  realized through linear system solving, but cost is high if solved accurately, thus only approximately, such as

- incomplete decompositions  $LDL^H$  of  $A - \sigma B$  with/without an iterative method
- CG, MINRES

(more from Sherry Li's lectures.)

Use Golub and Ye's Theorem (2002) as starting point:

$$\rho_{\ell+1} - \lambda_1 \leq (\rho_{\ell} - \lambda_1)\epsilon_m^2 + 2(\rho_{\ell} - \lambda_1)^{3/2}\epsilon_m \left( \frac{\|B\|_2}{\omega_2} \right)^{1/2} + \mathcal{O}(|\rho_{\ell} - \lambda_1|^2),$$
$$\epsilon_m := \min_{f \in \mathbb{P}_{m-1}, f(\omega_1)=1} \max_{j>1} |f(\omega_j)|,$$

where  $\omega_1 < \omega_2 \leq \dots \leq \omega_n$  are the eigenvalues of  $A - \rho_{\ell}B$ .

Idea: Transform  $A - \lambda B$  to  $L^{-1}(A - \lambda B)L^{-H}$  so that  $L^{-1}(A - \rho_{\ell}B)L^{-H}$  has “better” eigenvalue distribution, i.e., much smaller  $\epsilon_m$ .

**Ideal:**  $\omega_2 = \dots = \omega_n$ , then  $\epsilon_m = 0$  for  $m \geq 2$  and thus  $\rho_{\ell+1} - \lambda_1 = \mathcal{O}(|\rho_{\ell} - \lambda_1|^2)$ , quadratic convergence.

$A - \rho_{\ell}B = LDL^H \Rightarrow L^{-1}(A - \rho_{\ell}B)L^{-H} = D = \text{diag}(\pm 1)$ . Ideal but not practical:

- 1**  $A - \rho_{\ell}B = LDL^H$  may not exist at all. It exists if all leading principle submatrices are nonsingular.
- 2**  $A - \rho_{\ell}B = LDL^H$  may not be numerically stable to compute, especially when  $\rho_{\ell} \approx \lambda_1$ .
- 3**  $L$  significantly denser than  $A$  and  $B$  combined. Ensuing computations are too expensive.

**Compromise.**  $A - \rho_\ell B \approx LDL^H$  with a good chance that

one smallest isolated eigenvalue  $\omega_1$ , and the rest  $\omega_j$  ( $2 \leq j \leq n$ ) a few tight clusters.

Here  $A - \rho_\ell B \approx LDL^H$  includes not only the usual “approximately equal”, but also when  $(A - \rho_\ell B) - LDL^H$  approximately of a low rank.

$L$  varies from one iterative step to another; Can be expensive; Possible to use constant preconditioner, i.e., one  $L$  for all steps or change it every few steps.

Constant preconditioner: Use a shift  $\sigma \approx \lambda_1$ , and perform an incomplete  $LDL^H$  decomposition of  $A - \sigma B \approx LDL^H$ . Then

$$\widehat{C}_\ell = L^{-1}(A - \sigma B)L^{-H} + (\sigma - \rho_\ell)L^{-1}BL^{-H} \approx D$$

would have a better spectral distribution so long as  $(\sigma - \rho_\ell)L^{-1}BL^{-H}$  is small relative to  $\widehat{C}_\ell$ .

Insisted so far about applying ESD straightforwardly to the transformed problem  $L^{-1}(A - \lambda B)L^{-H}$ . But there is alternative, perhaps better, way.

# Eigen-problem preconditioning, II

$A - \lambda B$  to  $\widehat{A}_\ell - \lambda \widehat{B}_\ell := L_\ell^{-1}(A - \lambda B)L_\ell^{-H}$ . Typical step of ESD for  $\widehat{A}_\ell - \lambda \widehat{B}_\ell$ :

compute the smallest eigenvalue  $\mu$  and corresponding eigenvector  $v$  of  $\widehat{Z}^H(\widehat{A}_\ell - \lambda \widehat{B}_\ell)\widehat{Z}$ , where  $\widehat{Z} \in \mathbb{C}^{n \times m}$  is a basis matrix of Krylov subspace  $\mathcal{K}_m(\widehat{A}_\ell - \widehat{\rho}_\ell \widehat{B}_\ell, \widehat{\mathbf{x}}_\ell)$ .

Notice  $[\widehat{A}_\ell - \widehat{\rho}_\ell \widehat{B}_\ell]^j \widehat{\mathbf{x}}_\ell = L_\ell^H [(L_\ell L_\ell^H)^{-1}(A - \widehat{\rho}_\ell B)]^j (L_\ell^{-H} \widehat{\mathbf{x}}_\ell)$  to see

$$L_\ell^{-H} \cdot \mathcal{K}_m(\widehat{A}_\ell - \widehat{\rho}_\ell \widehat{B}_\ell, \widehat{\mathbf{x}}_\ell) = \mathcal{K}_m(K_\ell(A - \widehat{\rho}_\ell B), \mathbf{x}_\ell), \quad \mathbf{x}_\ell = L_\ell^{-H} \widehat{\mathbf{x}}_\ell, \quad K_\ell = (L_\ell L_\ell^H)^{-1}.$$

So  $Z = L_\ell^{-H} \widehat{Z}$  is a basis matrix of Krylov subspace  $\mathcal{K}_m(K_\ell(A - \widehat{\rho}_\ell B), \mathbf{x}_\ell)$ . Also

$$\widehat{Z}^H(\widehat{A}_\ell - \lambda \widehat{B}_\ell)\widehat{Z} = (L_\ell^{-H} \widehat{Z})^H(A - \lambda B)(L_\ell^{-H} \widehat{Z}) = Z^H(A - \lambda B)Z,$$

$$\widehat{\rho}_\ell = \frac{\widehat{\mathbf{x}}_\ell^H \widehat{A}_\ell \widehat{\mathbf{x}}_\ell}{\widehat{\mathbf{x}}_\ell^H \widehat{B}_\ell \widehat{\mathbf{x}}_\ell} = \frac{\mathbf{x}_\ell^H A \mathbf{x}_\ell}{\mathbf{x}_\ell^H B \mathbf{x}_\ell} = \rho_\ell.$$

The typical step can be reformulated equivalently to

compute the smallest eigenvalue  $\mu$  and corresponding eigenvector  $v$  of  $Z^H(A - \lambda B)Z$ , where  $Z \in \mathbb{C}^{n \times m}$  is a basis matrix of Krylov subspace  $\mathcal{K}_m(K_\ell(A - \rho_\ell B), \mathbf{x}_\ell)$ , where  $K_\ell = (L_\ell L_\ell^H)^{-1}$ .

## Extended Preconditioned Steepest Descent method

Given an initial approximation  $\mathbf{x}_0$  to  $u_1$ , a relative tolerance `rtol`, and an integer  $m \geq 2$ , the algorithm attempts to compute an approximate eigenpair to  $(\lambda_1, u_1)$  with the prescribed `rtol`.

---

```

1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|_B$ ,  $\rho_0 = \mathbf{x}_0^H A \mathbf{x}_0$ ,  $\mathbf{r}_0 = A \mathbf{x}_0 - \rho_0 B \mathbf{x}_0$ ;
2: for  $\ell = 0, 1, \dots$  do
3:   if  $\|\mathbf{r}_\ell\|_2 / (\|A \mathbf{x}_\ell\|_2 + |\rho_\ell| \|B \mathbf{x}_\ell\|_2) \leq \text{rtol}$  then
4:     BREAK;
5:   else
6:     construct a preconditioner  $K_\ell$ ;
7:     compute a basis matrix  $Z \in \mathbb{C}^{n \times m}$  of Krylov subspace
        $\mathcal{K}_m(K_\ell(A - \rho_\ell B), \mathbf{x}_\ell)$ ;
8:     compute the smallest eigenvalue  $\mu$  and corresponding eigenvector  $v$  of
        $Z^H(A - \lambda B)Z$ ;
9:      $y = Zv$ ,  $\mathbf{x}_{\ell+1} = y / \|y\|_B$ ;
10:     $\rho_{\ell+1} = \mu$ ,  $\mathbf{r}_{\ell+1} = A \mathbf{x}_{\ell+1} - \rho_{\ell+1} B \mathbf{x}_{\ell+1}$ ;
11:   end if
12: end for
13: return  $(\rho_\ell, \mathbf{x}_\ell)$  as an approximate eigenpair to  $(\lambda_1, u_1)$ .
```

---

It actually includes [Eigen-problem preconditioning, I & II](#).

# Convergence Rate

Theorem on Convergence Rate (Golub & Ye, 2002)

Suppose  $\lambda_1$  is simple, i.e.,  $\lambda_1 < \lambda_2$ , and  $\lambda_1 < \rho_\ell < \lambda_2$ , and preconditioner  $K_\ell \succ 0$ . Let

$$\omega_1 < \omega_2 \leq \dots \leq \omega_n$$

be the eigenvalues of  $K_\ell(A - \rho_\ell B)$  and  $v_1$  be an eigenvector corresponding to  $\omega_1$ . Then

$$\rho_{\ell+1} - \lambda_1 \leq (\rho_\ell - \lambda_1)\epsilon_m^2 + 2(\rho_\ell - \lambda_1)^{3/2}\epsilon_m \left( \frac{\|L_\ell^{-1}BL_\ell^{-H}\|_2}{\omega_2} \right)^{1/2} + \delta_\ell$$

where

$$0 \leq \delta_\ell := \rho_\ell - \lambda_1 + \omega_1 \frac{v_1^H K_\ell^{-1} v_1}{v_1^H B v_1} = \mathcal{O}(|\rho_\ell - \lambda_1|^2),$$

$$\epsilon_m := \min_{f \in \mathbb{P}_{m-1}, f(\omega_1)=1} \max_{j>1} |f(\omega_j)|.$$

$\epsilon_m := \min_{f \in \mathbb{P}_{m-1}, f(\omega_1)=1} \max_{j>1} |f(\omega_j)|$  usually unknown, except,

**1**  $m = 2$  for which the optimal  $f_{\text{opt}}$  is

$$f(t) = \frac{t - (\omega_2 + \omega_n)/2}{\omega_1 - (\omega_2 + \omega_n)/2} \in \mathbb{P}_{m-1}, \quad f(\omega_1) = 1, \quad \max_{j>1} |f(\omega_j)| = |f(\omega_2)| < 1.,$$

$$\epsilon_2 = \frac{1 - \eta}{1 + \eta}, \quad \eta = \frac{\omega_2 - \omega_1}{\omega_n - \omega_1}.$$

**2**  $\omega_2 = \dots = \omega_n$  for which  $f_{\text{opt}}(t) = (t - \omega_2)/(\omega_1 - \omega_2)$  and  $\epsilon_m = 0$  for all  $m \geq 2$ .

In general  $\epsilon_m$  can be bounded by using the Chebyshev polynomial

$$f(t) = \mathcal{T}_{m-1} \left( \frac{2t - (\omega_n + \omega_2)}{\omega_n - \omega_2} \right) / \mathcal{T}_{m-1} \left( \frac{1 + \eta}{1 - \eta} \right), \quad f(\omega_1) = 1,$$

$$\epsilon_m \leq \max_{\omega_2 \leq t \leq \omega_n} |f(t)| = \left[ \mathcal{T}_{m-1} \left( \frac{1 + \eta}{1 - \eta} \right) \right]^{-1} = 2 \left[ \Delta_\eta^{m-1} + \Delta_\eta^{-(m-1)} \right]^{-1}.$$

$$\rho_{\ell+1} - \lambda_1 \leq (\rho_\ell - \lambda_1)\epsilon_m^2 + 2(\rho_\ell - \lambda_1)^{3/2}\epsilon_m \left( \frac{\|L_\ell^{-1}BL_\ell^{-H}\|_2}{\omega_2} \right)^{1/2} + \mathcal{O}(|\rho_\ell - \lambda_1|^2).$$

Ignoring high order terms,  $\frac{\rho_{\ell+1} - \lambda_1}{\rho_\ell - \lambda_1} \approx \epsilon_m^2$ .

If  $\epsilon_m = 0$  (unlikely, however), then  $\rho_{\ell+1} - \lambda_1 = \mathcal{O}(|\rho_\ell - \lambda_1|^2)$ , quadratically convergence.

Locally,  $\rho_\ell \approx \lambda_1$ ,  $\text{eig}(K_\ell(A - \rho_\ell B)) \approx \text{eig}(K_\ell(A - \lambda_1 B)) = \{0 = \gamma_1 < \gamma_2 \leq \dots \leq \gamma_n\}$ , and

$$\epsilon_m \approx \min_{f \in \mathbb{P}_{m-1}, f(\gamma_1)=1} \max_{j>1} |f(\gamma_j)| \leq 2 \left[ \Delta_\eta^{m-1} + \Delta_\eta^{-(m-1)} \right]^{-1},$$

$$\Delta_\eta = \frac{1 + \sqrt{\eta}}{1 - \sqrt{\eta}}, \quad \eta = \frac{\gamma_2 - \gamma_1}{\gamma_n - \gamma_1}.$$

# Convergence Rate: Samokish Theorem

Theorem on Convergence Rate (Samokish, 1958)

$m = 2$ ,  $\text{eig}(K_\ell(A - \lambda_1 B)) = \{0 = \gamma_1 < \gamma_2 \leq \dots \leq \gamma_n\}$ . Suppose  $\lambda_1$  is simple, i.e.,  $\lambda_1 < \lambda_2$ , and  $\lambda_1 < \rho_\ell < \lambda_2$ ,

$$\delta = \sqrt{\|B^{1/2}K_\ell B^{1/2}\|_2 [\rho_\ell - \lambda_1]}, \quad \tau = \frac{2}{\gamma_2 + \gamma_n}.$$

If  $\tau(\sqrt{\gamma_n} + \delta)\delta < 1$ , then

$$\rho_{\ell+1} - \lambda_1 \leq \left[ \frac{\epsilon_2 + \tau\sqrt{\gamma_n}\delta}{1 - \tau(\sqrt{\gamma_n} + \delta)\delta} \right]^2 [\rho_\ell - \lambda_1], \quad (7)$$

where

$$\eta = \frac{\gamma_2 - \gamma_1}{\gamma_n - \gamma_1}, \quad \epsilon_2 = 2[\Delta_\eta + \Delta_\eta^{-1}]^{-1} = \frac{1 - \eta}{1 + \eta}.$$

Asymptotically  $\rho_{\ell+1} - \lambda_1 \lesssim \epsilon_2^2 (\rho_\ell - \lambda_1)$ , same as Golub & Ye (2002), but (7) is strict.

- 1)  $t_{\text{opt}} = \arg \min_t \rho(\mathbf{x}_\ell + tK_\ell \mathbf{r}_\ell)$ ,  $y = \mathbf{x}_\ell + t_{\text{opt}}K_\ell \mathbf{r}_\ell$ . Thus  $\rho_{\ell+1} = \rho(y)$ .
- 2) Drop subscript  $\ell$  to  $\mathbf{x}$ ,  $\mathbf{r}$ , and  $K$ :  $\mathbf{r}_\ell = r(\mathbf{x})$ ,  $\rho_\ell = \rho(\mathbf{x})$ .
- 3)  $z = \mathbf{x} - \tau K r(\mathbf{x})$ . Then  $\lambda_1 \leq \rho(y) \leq \rho(z)$ , thus  $\rho(y) - \lambda_1 \leq \rho(z) - \lambda_1$ .
- 4) Suffices to show  $\rho(z) - \lambda_1 \leq \text{RHS of (7)}$ .
- 5)  $A - \lambda_1 B$  is symmetric positive semidefinite.  $\|\cdot\|_{A-\lambda_1 B}$  is a semi-norm.

$$\begin{aligned} \|w\|_{A-\lambda_1 B}^2 &= [\rho(w) - \lambda_1] \|w\|_B^2, \\ \|[I - \tau K(A - \lambda_1 B)]w\|_{A-\lambda_1 B} &\leq \epsilon_2 \|w\|_{A-\lambda_1 B}. \end{aligned}$$

- 6) Write  $z = [I - \tau K(A - \lambda_1 B)]\mathbf{x} + \tau[\rho(\mathbf{x}) - \lambda_1]K\mathbf{B}\mathbf{x}$ , and assume  $\|\mathbf{x}\|_B = 1$ .

$$\begin{aligned} \|z\|_{A-\lambda_1 B} &= \sqrt{\rho(z) - \lambda_1} \|z\|_B, \\ \|z\|_{A-\lambda_1 B} &\leq \|[I - \tau K(A - \lambda_1 B)]\mathbf{x}\|_{A-\lambda_1 B} + \tau[\rho(\mathbf{x}) - \lambda_1] \|K\mathbf{B}\mathbf{x}\|_{A-\lambda_1 B} \\ &\leq \epsilon_2 \|\mathbf{x}\|_{A-\lambda_1 B} + \tau[\rho(\mathbf{x}) - \lambda_1] \sqrt{\gamma_n} \|B\mathbf{x}\|_K \\ &\leq \epsilon_2 \sqrt{\rho(\mathbf{x}) - \lambda_1} + \tau[\rho(\mathbf{x}) - \lambda_1] \sqrt{\gamma_n \|B^{1/2}KB^{1/2}\|_2} \\ &= (\epsilon_2 + \tau\sqrt{\gamma_n} \delta) \sqrt{\rho(\mathbf{x}) - \lambda_1}. \end{aligned}$$

7)  $z = \mathbf{x} - \tau K r(\mathbf{x})$ , and  $\|\mathbf{x}\|_B = 1$ .

$$\begin{aligned} \|z\|_B &\geq \|\mathbf{x}\|_B - \tau \|Kr(\mathbf{x})\|_B = 1 - \tau \|Kr(\mathbf{x})\|_B, \\ \|Kr(\mathbf{x})\|_B &= \|K(A - \lambda_1 B)\mathbf{x} - [\rho(\mathbf{x}) - \lambda_1]KB\mathbf{x}\|_B \\ &\leq \|K(A - \lambda_1 B)\mathbf{x}\|_B + [\rho(\mathbf{x}) - \lambda_1] \|KB\mathbf{x}\|_B \\ &\leq \sqrt{\|K^{1/2}BK^{1/2}\|_2 \gamma_n} \|\mathbf{x}\|_{A-\lambda_1 B} + [\rho(\mathbf{x}) - \lambda_1] \|B^{1/2}KB^{1/2}\|_2 \|\mathbf{x}\|_B \\ &= \sqrt{\gamma_n} \delta + \delta^2. \end{aligned}$$

8) Finally use

$$\rho(z) - \lambda_1 = \frac{\|z\|_{A-\lambda_1 B}^2}{\|z\|_B^2} \leq \frac{\|z\|_{A-\lambda_1 B}^2}{[1 - \tau \|Kr(\mathbf{x})\|_B]^2}$$

to complete the proof.

So far computing  $(\lambda_1, u_1)$  by *single-vector* steepest descent type methods.

To compute any following eigenpairs, must incorporate deflation techniques.

Or use a multi-vector/block method (not discussed yet). Deflation is also a necessary tool to make a block method more efficient.

Assume acceptable approximations to  $(\lambda_j, u_j)$  for  $1 \leq j \leq k$  known.

Diagonal  $\mathbf{D} \in \mathbb{R}^{k \times k}$  holds known approximations of  $\lambda_j$ ,

$\mathbf{U} \in \mathbb{R}^{n \times k}$  holds known approximations of  $u_j$ .

Assume  $\mathbf{U}^H \mathbf{B} \mathbf{U} = \mathbf{I}_k$ .

Deflation: avoid computing  $(\lambda_j, u_j)$  for  $1 \leq j \leq k$ , and seek approximation to  $(\lambda_{k+1}, u_{k+1})$ . Will discuss two deflation techniques.

# Through orthogonalizing against $U$

When the basis matrix  $Z$  is computed, make sure that  $Z$  is  $B$ -orthogonal to  $U$ . E.g., build a basis matrix  $Z$  for  $\mathcal{K}_m(K(A - \rho B), x)$  such that  $U^H \perp_B Z = 0$ . Suppose  $x \perp_B U = 0$  already.

## Arnoldi-like process

```
1:  $Z_{(:,1)} = x / \|x\|_B$ ,  $\rho = x^H A x / \|x\|_B^2$ ;  
2: for  $i = 2$  to  $m$  do  
3:    $q = K(AZ_{(:,i-1)} - \rho BZ_{(:,i-1)})$ ;  
4:    $q = q - U(U^H(Bq))$ ;  
5:   for  $j = 1$  to  $i - 1$  do  
6:      $t = Z_{(:,i-1)}^H Bq$ ,  $q = q - Z_{(:,i-1)} t$ ;  
7:   end for  
8:    $t = \|q\|_B$ ;  
9:   if  $t > 0$  then  
10:     $Z_{(:,i)} = q / t$ ;  
11:   else  
12:    BREAK;  
13:   end if  
14: end for
```

Note: Keep  $AZ$  for later use.

# Through shifting $\lambda_i$ away

## Lemma

Let  $U_1 = U_{(:,1:k)} = [u_1, \dots, u_k]$ .  $(A + \zeta BU_1 U_1^H B) - \lambda B$  and  $A - \lambda B$  share same eigenvectors  $u_i$ , but the eigenvalues of  $(A + \zeta BU_1 U_1^H B) - \lambda B$  are

$$\lambda_i + \zeta \text{ for } 1 \leq i \leq k \text{ and } \lambda_i \text{ for } k+1 \leq i \leq n.$$

Modify  $A - \lambda B$  in form, but not explicitly, to  $(A + \zeta BUU^H B) - \lambda B$ , where  $\zeta$  should be selected such that  $\zeta + \lambda_1 \geq \lambda_{k+2}$ .

But  $\lambda_{k+2}$  is unknown. What we can do in practice to pick  $\zeta$  a sufficiently large number.

# Block Steepest Descent Method

- Block Steepest Descent Method
- Block Extended Steepest Descent Method
- Block Preconditioned Extended Steepest Descent Method

## Single-vector SD and variations:

- compute  $(\lambda_1, u_1)$ , and, with deflations, other  $(\lambda_i, u_i)$ , one pair at a time. Most computations are of matrix-vector type.
- Slow convergence if  $(\gamma_2 - \gamma_1)/(\gamma_n - \gamma_1)$  tiny; usually happens when  $\lambda_2$  very close to  $\lambda_1$ . ( $\gamma_i$  are eigenvalues of  $K(A - \lambda_1 B)$ .)
- Often in practice, there are needs to compute the first few eigenpairs, not just the first one.

## Block versions:

- Can simultaneously compute the first  $k$  eigenpairs  $(\lambda_j, u_j)$ ;
- Run more efficiently on modern computer architecture: more computations in matrix-matrix multiplication type;
- Better rates of convergence; can save overall cost by using a block size that is slightly bigger than the number of asked eigenpairs.

In summary, the benefits of using a block variation are similar to those of using the simultaneous subspace iteration vs. the power method.

Start with  $X_0 \in \mathbb{C}^{n \times n_b}$ ,  $\text{rank}(X_0) = n_b \geq k$ , instead of just one vector  $\mathbf{x}_0 \in \mathbb{C}^n$ .

May assume  $j$ th column of  $X_0$  approximates  $u_j$ ; otherwise  $\mathcal{R}(X_0)$  approximates  $\text{span}\{u_1, \dots, u_{n_b}\}$ . In the latter, preprocessing  $X_0$ :

- 1 compute eigen-decomposition  $(X_0^H A X_0) W = (X_0^H B X_0) W \Omega$ , where  $\Omega = \text{diag}(\rho_{0;1}, \rho_{0;2}, \dots, \rho_{0;n_b})$ ;
- 2 Reset  $X_0 := X_0 W$ .

Can always assume  $j$ th column of  $X_0$  approximates  $u_j$ .

Typical  $\ell$ th iterative step: already have

$$X_\ell = [x_{\ell;1}, x_{\ell;2}, \dots, x_{\ell;n_b}] \in \mathbb{C}^{n \times n_b}, \quad j\text{th column } x_{\ell;j} \text{ approximates } u_j,$$

$$\Omega_\ell = \text{diag}(\rho_{\ell;1}, \rho_{\ell;2}, \dots, \rho_{\ell;n_b}), \quad \rho_{\ell;j} = \rho(x_{\ell;j}) \approx \lambda_j.$$

To compute new approximations as follows.

- 1 Compute a basis matrix  $Z$  of  $\mathcal{R}([X_\ell, R_\ell])$  by, e.g., MGS in the  $B$ -inner product, keeping in mind that  $X_\ell$  is  $B$ -orthonormal already;
- 2 Find the first  $n_b$  eigenpairs of  $Z^H A Z - \lambda Z^H B Z$  to get  $(Z^H A Z)W = (Z^H B Z)W \Omega_{\ell+1}$ ,  $\Omega_{\ell+1} = \text{diag}(\rho_{\ell+1;1}, \rho_{\ell+1;2}, \dots, \rho_{\ell+1;n_b})$ ;
- 3 Set  $X_{\ell+1} = ZW$ .

Block SD (previous slide) is the stronger version of **Simultaneous Rayleigh Quotient Minimization Method** of Longsine and McCormick (1980).

Note that  $r(x_{\ell;j}) = (A - \rho_{\ell;j}B)x_{\ell;j}$  and thus

$$\mathcal{R}([X_{\ell}, R_{\ell}]) = \sum_{j=1}^{n_b} \mathcal{R}([x_{\ell;j}, (A - \rho_{\ell;j}B)x_{\ell;j}]) = \sum_{j=1}^{n_b} \mathcal{K}_2(A - \rho_{\ell;j}B, x_{\ell;j}).$$

**Naturally**, as before, to expand search space,  $\mathcal{R}([X_{\ell}, R_{\ell}])$  through extending each  $\mathcal{K}_2(Ax_{\ell;j} - \rho_{\ell;j}B, x_{\ell;j})$  to a high order one. The new extended search subspace now is

$$\sum_{j=1}^{n_b} \mathcal{K}_m(A - \rho_{\ell;j}B, x_{\ell;j}) = \text{span}\{X_{\ell}, \mathcal{R}_{\ell}(X_{\ell}), \dots, \mathcal{R}_{\ell}^{m-1}(X_{\ell})\} =: \mathcal{K}_m(\mathcal{R}_{\ell}, X_{\ell}),$$

where the linear operator  $\mathcal{R}_{\ell} : X \in \mathbb{C}^{n \times n_b} \rightarrow \mathcal{R}_{\ell}(X) = AX - BX\Omega_{\ell} \in \mathbb{C}^{n \times n_b}$ .

$\mathcal{R}_{\ell}^i(\cdot) = \mathcal{R}_{\ell}^{i-1}(\mathcal{R}_{\ell}(\cdot))$ , e.g.,  $\mathcal{R}_{\ell}^2(X) = \mathcal{R}_{\ell}(\mathcal{R}_{\ell}(X))$ .

Block Extended SD: make  $Z$  basis matrix of  $\mathcal{K}_m(\mathcal{R}_{\ell}, X_{\ell})$ .

In light of extensive discussions on preconditioning, **natural** to modify the search subspace to

$$\sum_{j=1}^{n_b} \mathcal{K}_m(K_{\ell;j}(A - \rho_{\ell;j}B), x_{\ell;j}),$$

where  $K_{\ell;j}$  is the preconditioner intended to move  $(\rho_{\ell;j}, x_{\ell;j})$  towards  $(\lambda_j, u_j)$  faster for each  $j$ .

Two ways to construct  $K_{\ell;j}$ :

- $K_{\ell;j} \approx (A - \tilde{\rho}_{\ell;j}B)^{-1}$  for some  $\tilde{\rho}_{\ell;j} \neq \rho_{\ell;j}$ , ideally closer to  $\lambda_j$  than to any other eigenvalue of  $A - \lambda B$ .

Since the eigenvalues of  $A - \lambda B$  are unknown, practically make  $\tilde{\rho}_{\ell;j}$  closer but not equal to  $\rho_{\ell;j}$  than to any other  $\rho_{\ell;k}$ .

- Perform incomplete  $LDL^H$  factorization:  $A - \rho_{\ell;j}B \approx L_{\ell;j}D_{\ell;j}L_{\ell;j}^H$ , where “ $\approx$ ” includes not only the usual “approximately equal”, but also the case when  $(A - \rho_{\ell;j}B) - L_{\ell;j}D_{\ell;j}L_{\ell;j}^H$  is approximately a low rank matrix, and  $D_{\ell;j} = \text{diag}(\pm 1)$ .

Finally,  $K_{i;j} = L_{\ell;j}L_{\ell;j}^H$ .

## Block Preconditioned Extended Steepest Descent method

Given an initial approximation  $X_0 \in \mathbb{C}^{n \times n_b}$  with  $\text{rank}(X_0) = n_b$ , and an integer  $m \geq 2$ , the algorithm attempts to compute approximate eigenpair to  $(\lambda_j, u_j)$  for  $1 \leq j \leq n_b$ .

- 1: compute the eigen-decomposition:  $(X_0^H A X_0) W = (X_0^H B X_0) W \Omega_0$ ,  
where  $W^H (X_0^H B X_0) W = I$ ,  $\Omega_0 = \text{diag}(\rho_{0;1}, \rho_{0;2}, \dots, \rho_{0;n_b})$ ;
- 2:  $X_0 = X_0 W$ ;
- 3: **for**  $\ell = 0, 1, \dots$  **do**
- 4:     test convergence and lock up the converged (detail to come later);
- 5:     construct preconditioners  $K_{\ell;j}$  for  $1 \leq j \leq n_b$ ;
- 6:     compute a basis matrix  $Z \in \mathbb{C}^{n \times m n_b}$  of  $\sum_{j=1}^{n_b} \mathcal{K}_m(K_{\ell;j}(A - \rho_{\ell;j} B), x_{\ell;j})$ ;
- 7:     compute the  $n_b$  smallest eigenvalues and corresponding eigenvectors of  $Z^H(A - \lambda B)Z$  to get  $(Z^H A Z) W = (Z^H B Z) W \Omega_\ell$ , where  $W^H (Z^H B Z) W = I$ ,  
 $\Omega_{\ell+1} = \text{diag}(\rho_{\ell+1;1}, \rho_{\ell+1;2}, \dots, \rho_{\ell+1;n_b})$ ;
- 8:      $X_{\ell+1} = ZW$ ;
- 9: **end for**
- 10: **return** approximate eigenpairs to  $(\lambda_j, u_j)$  for  $1 \leq j \leq n_b$ .

Different preconditioner  $K_{\ell;j}$  for each different approximate eigenpair  $(\rho_{\ell;j}, x_{\ell;j})$  good for convergence rates, but may not reduce overall time:

- expensive to construct all preconditioners
- cannot compute  $Z$  mostly by matrix-matrix multiplications (more later)

Use  $K_{\ell;j} \equiv K_{\ell}$ , one preconditioner for all speeding up the convergence of  $(\rho_{\ell;1}, x_{\ell;1})$ . At the same time other  $(\rho_{\ell;j}, x_{\ell;j})$  are making progress, too, but at a slower speed.

Usually  $(\rho_{\ell;1}, x_{\ell;1})$  converges first and quickly.

Once  $(\rho_{\ell;1}, x_{\ell;1})$  (or the first few in the case of a tight cluster) is determined to be sufficiently accurate, the converged eigenpair is **locked** up and **deflated**.

A new preconditioner is computed to aim at the next non-converged eigenpair, and the process continues.

Need to compute basis matrix  $Z \in \mathbb{C}^{n \times mn_b}$  of  $\sum_{j=1}^{n_b} \mathcal{K}_m(K_{\ell;j}(A - \rho_{\ell;j}B), x_{\ell;j})$ .

$Z$  can be gotten by packing the basis matrices of all  $\mathcal{K}_m(K_{\ell;j}(A - \rho_{\ell;j}B), x_{\ell;j})$  for  $1 \leq j \leq n_b$  together. Two drawbacks:

- Such a  $Z$  could be ill-conditioned, i.e., columns of  $Z$  may not be sufficiently numerically linearly independent; Possible cure: re-orthogonalize packed  $Z$  – too costly.
- Building basis for each  $\mathcal{K}_m(K_{\ell;j}(A - \rho_{\ell;j}B), x_{\ell;j})$  uses mostly BLAS-2 operations. Have to be this way if  $K_{\ell;j}$  are different.

Different situation if  $K_{\ell;j} \equiv K$  (drop iteration step index  $\ell$ ). Then

$$\begin{aligned} \sum_{j=1}^{n_b} \mathcal{K}_m(K(A - \rho_{\ell;j}B), x_{\ell;j}) &= \mathcal{K}_m(K\mathcal{R}, X) \\ &\equiv \text{span}\{X, K\mathcal{R}(X), \dots, [K\mathcal{R}]^{m-1}(X)\}, \end{aligned}$$

where  $\mathcal{R}(X) = AX - BX\Omega$ ,  $[K\mathcal{R}]^i(\cdot) = [K\mathcal{R}]^i(K\mathcal{R}(\cdot))$ , e.g.,  $[K\mathcal{R}]^2(X) = K\mathcal{R}_\ell(K\mathcal{R}(X))$ .

$Z = [Z_1, Z_2, \dots, Z_m]$  can be computed by the following block Arnoldi-like process in the  $B$ -inner product.

## Arnoldi-like process for $Z$

```
1:  $Z_1 = X$  (recall  $X^H B X = I_{n_b}$  already);  
2: for  $i = 2$  to  $m$  do  
3:    $Y = K(AZ_{i-1} - B\Omega Z_{i-1})$ ;  
4:   for  $j = 1$  to  $i - 1$  do  
5:      $T = Z_j^H B Y$ ;  $Y = Y - Z_j T$ ;  
6:   end for  
7:    $Z_i T = Y$  (MGS in the  $B$ -inner product);  
8: end for
```

Note: At Line 7,  $Y$  may not be numerically of full column rank – not a problem.

Anytime if a column is deemed linearly dependent on previous columns, that column should be deleted, along with corresponding  $\rho_j$  from  $\Omega$ .

At completion of MGS,  $Z_{j+1}$  will have fewer columns than  $Y$  and the size of  $\Omega$  is shrunk accordingly.

$(\rho_{\ell;j}, x_{\ell;j})$  is considered acceptable if  $\frac{\|r_{\ell;j}\|_2}{\|Ax_{\ell;j}\|_2 + |\rho_{\ell;j}| \|Bx_{\ell;j}\|_2} \leq \text{rtol}$ .

Usually  $\lambda_j$  are converged to in order, i.e., the smallest eigenvalues emerge first.

**Lock** all acceptable approximate eigenpairs in  $k_{\text{cvgd}} \times k_{\text{cvgd}}$  diagonal matrix  $\mathbf{D}$  for eigenvalues and  $n \times k_{\text{cvgd}}$  tall matrix  $\mathbf{U}$  for eigenvectors.

Every time a converged eigenpair is detected, delete the converged  $\rho_{\ell;j}$  and  $x_{\ell;j}$  from  $\Omega_\ell$  and  $X_\ell$ , respectively, and expand  $\mathbf{D}$  and  $\mathbf{U}$  to lock up the pair, accordingly.

At the same time, either reduce  $n_b$  by 1 or append a (random)  $B$ -orthogonal column to  $X$  to maintain  $n_b$  unchanged.

**Deflate** to avoid recomputing converged eigenpairs:

- 1 At Line 7 in the Arnoldi-like process, each column of  $Z_{j+1}$  is  $B$ -orthogonalized against  $\mathbf{U}$ .
- 2 Modify  $A - \lambda B$  in form, but not explicitly, to  $(A + \zeta B \mathbf{U} \mathbf{U}^H B) - \lambda B$ , where  $\zeta$  should be selected such that  $\zeta + \lambda_1 \geq \lambda_{k_{\text{cvgd}} + n_b + 1}$ . Here we pre-assume the  $k_{\text{cvgd}}$  converged eigenpairs are indeed those for  $(\lambda_j, u_j)$  for  $1 \leq j \leq k_{\text{cvgd}}$ . This is usually so, but with no guarantee, of course.

# Conjugate Gradient Methods

- Digression: CG for Linear System  $Ax = b$
- Conjugate Gradient Method
- Preconditioned Conjugate Gradient Method
- Locally Optimal Conjugate Gradient Method
- Locally Optimal Extended Conjugate Gradient Method
- Locally Optimal Block Preconditioned Extended Conjugate Gradient Method

# CG for Linear System $Ax = b$

$A$  is  $n \times n$ , symmetric, and positive definite. Let

$$\phi(x) = \frac{1}{2}x^T Ax - x^T b,$$

quadratic in  $x$ , convex, a unique local and global minimum at  $x = A^{-1}b$ ,  
 $\nabla\phi(x) = r(x) \equiv Ax - b$ .

CG Algorithm (Hestenes and Stiefel, 1950s):

- 1 Given  $x_0$ , compute  $r_0 = Ax_0 - b$ , and set  $p_0 = -r_0$ ;
- 2 For  $i = 0, 1, \dots$ , do

$$\begin{aligned}\alpha_i &= \arg \min_{\alpha} \phi(x_i + \alpha p_i), & x_{i+1} &= x_i + \alpha_i p_i, \\ r_{i+1} &= r_i + \alpha_i A p_i, & p_{i+1} &= -r_{i+1} + \beta_i p_i.\end{aligned}$$

$\beta_i$  chosen so that  $p_{i+1}^T A p_i = 0$ ; equivalent expressions:

$$\beta_i = \frac{p_i^T A r_{i+1}}{p_i^T A p_i} = \frac{r_{i+1}^T r_{i+1}}{r_i^T r_i} = \frac{r_{i+1}^T (r_{i+1} - r_i)}{r_i^T r_i}.$$

*Verbatim* translations of Hestenes' and Stiefel's CG to solve

$$\min_x \phi(x), \quad \phi(x) \text{ not necessarily quadratic,}$$

replacing all  $r(x_i)$  by  $\nabla\phi(x_i)$ .

Nonlinear CG Algorithm (Fletcher and Reeves, 1964):

- 1 Given  $x_0$ , compute  $\nabla\phi_0 = \nabla\phi(x_0)$ , and set  $p_0 = -\nabla\phi_0$ ;
- 2 For  $i = 0, 1, \dots$ , do

$$\alpha_i = \arg \min_{\alpha} \phi(x_i + \alpha p_i), \quad x_{i+1} = x_i + \alpha_i p_i,$$
$$\text{evaluate } \nabla\phi_{i+1} = \nabla\phi(x_{i+1}), \quad p_{i+1} = -\nabla\phi_{i+1} + \beta_i p_i.$$

Several choices for  $\beta_i$ :

$$\beta_i = \frac{\nabla\phi_{i+1}^T \nabla\phi_{i+1}}{\nabla\phi_i^T \nabla\phi_i}, \quad \beta_i = \frac{\nabla\phi_{i+1}^T (\nabla\phi_{i+1} - \nabla\phi_i)}{\nabla\phi_i^T \nabla\phi_i}.$$

**Linear CG:** choices of  $\beta_i$  make

- search directions  $p_i$  conjugate, i.e.,  $p_i^T A p_j = 0$  for  $i \neq j$ .
- CG method terminates in at most  $n$  steps.

**Nonlinear CG:** many nice properties no longer hold for any choice of  $\beta_i$ .

Observe

$$\begin{aligned}x_{i+2} &= x_{i+1} + \alpha_{i+1}(-\nabla\phi_{i+1} + \beta_i p_i) \\ &\in \text{span}\{x_{i+1}, \nabla\phi_{i+1}, p_i\} = \text{span}\{x_{i+1}, \nabla\phi_{i+1}, x_i\}.\end{aligned}$$

Since many nice properties in linear CG are lost anyway in the nonlinear case, why not pick  $\beta_i$ , implicitly, such that (Takahashi, 1965)

$$x_{i+2} = \arg \min_{y \in \text{span}\{x_{i+1}, \nabla\phi_{i+1}, x_i\}} \phi(y).$$

This gives **locally optimal CG**. But search over  $y \in \text{span}\{x_{i+1}, \nabla\phi_{i+1}, x_i\}$  harder than before.

Minimize  $\rho(x)$  to compute  $(\lambda_1, u_1)$ :

$$\rho(x) = \frac{x^H Ax}{x^H Bx}, \quad \nabla \rho(x) = \frac{2}{x^H Bx} r(x), \quad r(x) := Ax - \rho(x) Bx.$$

Line-search  $\rho(y) = \inf_{t \in \mathbb{C}} \rho(x + tp)$

- 1: compute the smaller eigenvalue  $\mu$  of  $X^H Ax - \lambda X^H Bx$ , where  $X = [x, p]$ , and eigenvector  $v = [\nu_1, \nu_2]^T$ ;
- 2:  $\arg \inf_{t \in \mathbb{C}} \rho(x + tp) =: t_{\text{opt}} = \begin{cases} \nu_2/\nu_1, & \text{if } \nu_1 \neq 0, \\ \infty, & \text{if } \nu_1 = 0; \end{cases}$
- 3:  $y = \begin{cases} x + t_{\text{opt}} p & \text{if } t_{\text{opt}} \text{ is finite,} \\ p & \text{otherwise.} \end{cases}$

CG for  $Ax = \lambda Bx$ : in nonlinear CG simply replace  $\nabla \phi(x)$  by  $r(x) := Ax - \rho(x) Bx$ .

CG for  $Ax = \lambda Bx$ 

Given an initial approximation  $\mathbf{x}_0$  to  $u_1$ , and a relative tolerance  $\text{rtol}$ , the algorithm attempts to compute an approximate eigenpair to  $(\lambda_1, u_1)$  with the prescribed  $\text{rtol}$ .

- 1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|_B$ ,  $\boldsymbol{\rho}_0 = \mathbf{x}_0^H A \mathbf{x}_0$ ,  $\mathbf{r}_0 = A \mathbf{x}_0 - \boldsymbol{\rho}_0 B \mathbf{x}_0$ ,  $\boldsymbol{\rho}_0 = -\mathbf{r}_0$ ;
- 2: **for**  $\ell = 0, 1, \dots$  **do**
- 3:   **if**  $\|\mathbf{r}_\ell\|_2 / (\|A \mathbf{x}_\ell\|_2 + |\boldsymbol{\rho}_\ell| \|B \mathbf{x}_\ell\|_2) \leq \text{rtol}$  **then**
- 4:     **BREAK**;
- 5:   **else**
- 6:     compute  $\alpha_\ell = t_{\text{opt}} := \inf_{t \in \mathbb{C}} \rho(\mathbf{x}_\ell + t \boldsymbol{\rho}_\ell)$ , and then
 
$$y = \begin{cases} \mathbf{x}_\ell + \alpha_\ell \boldsymbol{\rho}_\ell & \text{if } \alpha_\ell \text{ is finite,} \\ \boldsymbol{\rho}_\ell & \text{otherwise.} \end{cases}$$
- 7:      $\mathbf{x}_{\ell+1} = y / \|y\|_B$ ;
- 8:     set  $\boldsymbol{\rho}_{\ell+1} = \mathbf{x}_{\ell+1}^H A \mathbf{x}_{\ell+1}$ ,  $\mathbf{r}_{\ell+1} = A \mathbf{x}_{\ell+1} - \boldsymbol{\rho}_{\ell+1} B \mathbf{x}_{\ell+1}$ ,  $\boldsymbol{\rho}_{\ell+1} = -\mathbf{r}_{\ell+1} + \beta_\ell \boldsymbol{\rho}_\ell$ ,  
 where  $\beta_\ell = \frac{\mathbf{r}_{\ell+1}^H \mathbf{r}_{\ell+1}}{\mathbf{r}_\ell^H \mathbf{r}_\ell}$  or  $\frac{\mathbf{r}_{\ell+1}^H (\mathbf{r}_{\ell+1} - \mathbf{r}_\ell)}{\mathbf{r}_\ell^H \mathbf{r}_\ell}$
- 9:   **end if**
- 10: **end for**
- 11: **return**  $(\boldsymbol{\rho}_\ell, \mathbf{x}_\ell)$  as an approximate eigenpair to  $(\lambda_1, u_1)$ .

# A Convergence Theorem

## Convergence Theorem (Yang, 1993)

With  $\beta_\ell = \frac{\mathbf{r}_{\ell+1}^H \mathbf{r}_{\ell+1}}{\mathbf{r}_\ell^H \mathbf{r}_\ell}$ ,  $\boldsymbol{\rho}_\ell$  converges to some eigenvalue  $\hat{\lambda}$  of  $A - \lambda B$  and there is a convergent subsequence  $\{\mathbf{x}_{\ell_i}\}$  of  $\{\mathbf{x}_\ell\}$  such that

$$\|(A - \hat{\lambda}B)\mathbf{x}_{\ell_i}\|_2 \rightarrow 0 \quad \text{as } i \rightarrow \infty,$$

i.e.,  $\mathbf{x}_{\ell_i}$  converges in direction to a corresponding eigenvector.

If  $\hat{\lambda} = \lambda_1$ , then  $\|(A - \hat{\lambda}B)\mathbf{x}_\ell\|_2 \rightarrow 0$  as  $\ell \rightarrow \infty$ , i.e.,  $\mathbf{x}_\ell$  converges in direction to a corresponding eigenvector. (*seem new*)

- First part due to (Yang, 1993); second part seems [new](#).
- Only for  $\beta_\ell = \frac{\mathbf{r}_{\ell+1}^H \mathbf{r}_{\ell+1}}{\mathbf{r}_\ell^H \mathbf{r}_\ell}$ , however.
- Proof much more complicated than for SD.
- Rate of convergence: mostly heuristic, none rigorous proven.

# Preconditioned CG for $Ax = \lambda Bx$

As in SD, Preconditioned CG = *vanilla* CG on  $L^{-H}AL^{-1} - \lambda L^{-H}BL^{-1}$ .

Let  $\tilde{A} - \lambda\tilde{B} := L^{-H}AL^{-1} - \lambda L^{-H}BL^{-1}$ . Adopt notation convention for  $\tilde{A} - \lambda\tilde{B}$ : same symbols but with *tildes*. E.g.,  $\tilde{x} = Lx$ ,

$$\tilde{\rho}(\tilde{x}) = \frac{\tilde{x}^H L^{-H} A L^{-1} \tilde{x}}{\tilde{x}^H L^{-H} B L^{-1} \tilde{x}} \equiv \rho(x), \quad \tilde{r}(\tilde{x}) = L^{-H} A L^{-1} \tilde{x} - \tilde{\rho}(\tilde{x}) L^{-H} B L^{-1} \tilde{x} \equiv L^{-H} r(x).$$

Key CG step:

$$\begin{aligned} \tilde{\alpha}_\ell &= \arg \min_{\tilde{\alpha}} \tilde{\rho}(\tilde{\mathbf{x}}_\ell + \tilde{\alpha} \tilde{\mathbf{p}}_\ell), & \tilde{\mathbf{x}}_{\ell+1} &= \tilde{\mathbf{x}}_\ell + \tilde{\alpha}_\ell \tilde{\mathbf{p}}_\ell, \\ \tilde{\mathbf{r}}_{\ell+1} &= L^{-H} A L^{-1} \tilde{\mathbf{x}}_{\ell+1} - \tilde{\rho}(\tilde{\mathbf{x}}_{\ell+1}) L^{-H} B L^{-1} \tilde{\mathbf{x}}_{\ell+1}, & \tilde{\mathbf{p}}_{\ell+1} &= -\tilde{\mathbf{r}}_{\ell+1} + \tilde{\beta}_\ell \tilde{\mathbf{p}}_\ell. \end{aligned}$$

Perform substitutions  $\tilde{\mathbf{x}}_\ell = L\mathbf{x}_\ell$  and  $\tilde{\mathbf{r}}_\ell = L^{-H}\mathbf{r}_\ell$ :

$$\begin{aligned} \tilde{\alpha}_\ell &= \arg \min_{\tilde{\alpha}} \rho(\mathbf{x}_\ell + \underbrace{\tilde{\alpha} L^{-1} \tilde{\mathbf{p}}_\ell}_{=: \mathbf{p}_\ell}), & \mathbf{x}_{\ell+1} &= \mathbf{x}_\ell + \tilde{\alpha}_\ell L^{-1} \tilde{\mathbf{p}}_\ell, \\ \mathbf{r}_{\ell+1} &= A\mathbf{x}_{\ell+1} - \rho(\mathbf{x}_{\ell+1}) B\mathbf{x}_{\ell+1}, & \underbrace{L^{-1} \tilde{\mathbf{p}}_{\ell+1}}_{=: \mathbf{p}_{\ell+1}} &= - \underbrace{(L^H L)^{-1}}_{=: K} \mathbf{r}_{\ell+1} + \tilde{\beta}_\ell \underbrace{L^{-1} \tilde{\mathbf{p}}_\ell}_{=: \mathbf{p}_\ell}. \end{aligned}$$

# Preconditioned CG for $Ax = \lambda Bx$

## Preconditioned CG for $Ax = \lambda Bx$

Given an initial approximation  $\mathbf{x}_0$  to  $u_1$ , a (positive definite) preconditioner  $K$ , and a relative tolerance  $\text{rtol}$ , the algorithm attempts to compute an approximate pair to  $(\lambda_1, u_1)$  with the prescribed  $\text{rtol}$ .

```
1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|_B$ ,  $\boldsymbol{\rho}_0 = \mathbf{x}_0^H A \mathbf{x}_0$ ,  $\mathbf{r}_0 = A \mathbf{x}_0 - \boldsymbol{\rho}_0 B \mathbf{x}_0$ ,  $\mathbf{p}_0 = -K \mathbf{r}_0$ ;  
2: for  $\ell = 0, 1, \dots$  do  
3:   if  $\|\mathbf{r}_\ell\|_2 / (\|A \mathbf{x}_\ell\|_2 + |\boldsymbol{\rho}_\ell| \|B \mathbf{x}_\ell\|_2) \leq \text{rtol}$  then  
4:     BREAK;  
5:   else  
6:     compute  $\alpha_\ell = t_{\text{opt}} := \inf_{t \in \mathbb{C}} \rho(\mathbf{x}_\ell + t \mathbf{p}_\ell)$ , and then  
       
$$y = \begin{cases} \mathbf{x}_\ell + \alpha_\ell \mathbf{p}_\ell & \text{if } \alpha_\ell \text{ is finite,} \\ \mathbf{p}_\ell & \text{otherwise.} \end{cases}$$
  
7:      $\mathbf{x}_{\ell+1} = y / \|y\|_B$ ;  
8:     set  $\boldsymbol{\rho}_{\ell+1} = \mathbf{x}_{\ell+1}^H A \mathbf{x}_{\ell+1}$ ,  $\mathbf{r}_{\ell+1} = A \mathbf{x}_{\ell+1} - \boldsymbol{\rho}_{\ell+1} B \mathbf{x}_{\ell+1}$ ,  
       
$$\mathbf{p}_{\ell+1} = -K \mathbf{r}_{\ell+1} + \beta_\ell \mathbf{p}_\ell$$
, where  $\beta_\ell = \frac{\mathbf{r}_{\ell+1}^H K \mathbf{r}_{\ell+1}}{\mathbf{r}_\ell^H K \mathbf{r}_\ell}$  or  $\frac{\mathbf{r}_{\ell+1}^H K (\mathbf{r}_{\ell+1} - \mathbf{r}_\ell)}{\mathbf{r}_\ell^H K \mathbf{r}_\ell}$ .  
9:   end if  
10: end for  
11: return  $(\boldsymbol{\rho}_\ell, \mathbf{x}_\ell)$  as an approximate eigenpair to  $(\lambda_1, u_1)$ .
```

Earlier discussions on selecting a good preconditioner for PSD should apply:

- $A - \sigma B = LDL^H$ ,  $D = \text{diag}(\pm 1)$ ,  $K = (L^H L)^{-1}$ .

Various heuristics on the convergence rates of the preconditioned CG, but none is rigorously proved. Even less can be said about the theoretical analysis of block (or subspace) versions of the preconditioned CG method (to come soon).

But since preconditioned CG is CG for  $L^{-H}AL^{-1} - \lambda L^{-H}BL^{-1}$ , previous convergence theorem for CG remains valid.

# Locally Optimal CG for $Ax = \lambda Bx$

In writing down CG for  $Ax = \lambda Bx$ , we did

- gradient-to-residual replacement: replacing the gradient by the eigen-residual  $r(x) = Ax - \rho(x)Bx$  which differs by a scalar factor  $2/x^H Bx$  from the gradient  $\nabla \rho(x) = \frac{2}{x^H Bx} [Ax - \rho(x) Bx]$ ;
- also normalizing  $x_\ell$ . No theory around as to why we should normalize  $x_\ell$ , beside that they are some eigenvector approximations.

We made a couple of “arbitrary choices”. Their effects on the rate of convergence are not clear.

Locally optimal CG eliminates the “arbitrariness” altogether: compute  $x_{\ell+1}$  from the subspace  $\text{span}\{x_{\ell-1}, x_\ell, r_\ell\}$  by

$$\min_{x \in \text{span}\{x_{\ell-1}, x_\ell, r_\ell\}} \rho(x),$$

which is solvable through the Rayleigh-Ritz procedure.

## Locally Optimal CG for $Ax = \lambda Bx$

Given an initial approximation  $\mathbf{x}_0$  to  $u_1$ , and a relative tolerance  $\text{rtol}$ , the algorithm attempts to compute an approximate eigenpair to  $(\lambda_1, u_1)$  with the prescribed  $\text{rtol}$ .

```
1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|_B$ ,  $\rho_0 = \mathbf{x}_0^H A \mathbf{x}_0$ ,  $\mathbf{r}_0 = A \mathbf{x}_0 - \rho_0 B \mathbf{x}_0$ ,  $\mathbf{x}_{-1} = 0$ ;  
2: for  $\ell = 0, 1, \dots$  do  
3:   if  $\|\mathbf{r}_\ell\|_2 / (\|A \mathbf{x}_\ell\|_2 + |\rho_\ell| \|B \mathbf{x}_\ell\|_2) \leq \text{rtol}$  then  
4:     BREAK;  
5:   else  
6:     compute a basis matrix  $Z \in \mathbb{C}^{n \times k}$  ( $k = 2$  or  $3$ ) of the subspace  
        $\text{span}\{\mathbf{x}_\ell, \mathbf{x}_{\ell-1}, \mathbf{r}_\ell\}$ ;  
7:     compute the smallest eigenvalue  $\mu$  and corresponding eigenvector  $v$  of  
        $Z^H (A - \lambda B) Z$ ;  
8:      $y = Zv$ ,  $\mathbf{x}_{\ell+1} = y / \|y\|_B$ ;  
9:      $\rho_{\ell+1} = \mu$ ,  $\mathbf{r}_{\ell+1} = A \mathbf{x}_{\ell+1} - \rho_{\ell+1} B \mathbf{x}_{\ell+1}$ ;  
10:   end if  
11: end for  
12: return  $(\rho_\ell, \mathbf{x}_\ell)$  as an approximate eigenpair to  $(\lambda_1, u_1)$ .
```

$\mathbf{x}_\ell$  moves closer and closer to  $\mathbf{u}_1$ ;  $\mathbf{x}_\ell, \mathbf{x}_{\ell-1}$  increasingly move towards being linearly dependent.

Line 6:  $Z$  contaminated more and more by rounding errors. How to mitigate that?

To replace  $\mathbf{x}_{\ell-1}$  by some  $\mathbf{y}_\ell := \xi_{\ell,1}\mathbf{x}_\ell - \xi_{\ell,2}\mathbf{x}_{\ell-1}$  such that

$$\text{span}\{\mathbf{x}_\ell, \mathbf{x}_{\ell-1}, \mathbf{r}_\ell\} = \text{span}\{\mathbf{x}_\ell, \mathbf{y}_\ell, \mathbf{r}_\ell\}.$$

Then same  $(\mu, \nu)$  at Line 7. But need to generate  $\mathbf{y}_{\ell+1}$ , given  $\mathbf{x}_\ell, \mathbf{y}_\ell, \mathbf{r}_\ell$ .

$Z = [z_1, z_2, z_3]$  is  $B$ -orthonormal (by MGS), and  $z_1 = \mathbf{x}_\ell$ . Then  $y = Zv = \nu_1 z_1 + \nu_2 z_2 + \nu_3 z_3 = \nu_1 \mathbf{x}_\ell + \nu_2 z_2 + \nu_3 z_3$ .

Set  $\mathbf{y}_{\ell+1} := y - \nu_1 \mathbf{x}_\ell = \|y\|_B \mathbf{x}_{\ell+1} - \nu_1 \mathbf{x}_\ell =: \xi_{\ell+1,1} \mathbf{x}_{\ell+1} - \xi_{\ell+1,2} \mathbf{x}_\ell$ .

Modify Lines 1, 6, and 8 as follows while keeping others the same.

- 
- 1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|_B, \boldsymbol{\rho}_0 = \mathbf{x}_0^H \mathbf{A} \mathbf{x}_0, \mathbf{r}_0 = \mathbf{A} \mathbf{x}_0 - \boldsymbol{\rho}_0 \mathbf{B} \mathbf{x}_0, \mathbf{y}_0 = 0$ ;
  - 6: compute a basis matrix  $Z \in \mathbb{C}^{n \times k}$  ( $k = 2$  or  $3$ ) of the subspace  $\text{span}\{\mathbf{x}_\ell, \mathbf{y}_\ell, \mathbf{r}_\ell\}$ ;
  - 8:  $y = Zv, \mathbf{x}_{\ell+1} = y / \|y\|_B, \mathbf{y}_{\ell+1} = Z\hat{v}$ , where  $\hat{v}$  is  $v$  with its 1st entry zeroed;
-

## Convergence Theorem for LOCG)

$\rho_\ell$  converges to some eigenvalue  $\hat{\lambda}$  of  $A - \lambda B$  and  $\|(A - \hat{\lambda}B)\mathbf{x}_\ell\|_2 \rightarrow 0$  as  $\ell \rightarrow \infty$ , i.e.,  $\mathbf{x}_\ell$  converges in direction to a corresponding eigenvector.

- Same convergence theorem for SD;
- For CG, PCG, only  $\|(A - \hat{\lambda}B)\mathbf{x}_{\ell_i}\|_2 \rightarrow 0$ ;
- Inclusion of the residual  $\mathbf{r}_\ell$  makes the difference.

Three ideas for improving SD **naturally** apply here:

- 1 Incorporate a preconditioner  $K$ : simply modify  $\mathbf{r}_\ell$  to  $K\mathbf{r}_\ell$ ;
- 2 Extend search space from currently

$$\text{span}\{\mathbf{x}_{\ell-1}\} + \mathcal{K}_2(A - \rho_\ell B, \mathbf{x}_\ell) \quad \text{to} \quad \text{span}\{\mathbf{x}_{\ell-1}\} + \mathcal{K}_m(A - \rho_\ell B, \mathbf{x}_\ell);$$

- 3 Use block  $X_0 \in \mathbb{C}^{n \times n_b}$ .

The ideas can be applied in any combination ( $2^3 = 8$  of them): E.g.,

- Locally Optimal Preconditioned CG (LOPCG):  $m = 2$ ,  $n_b = 1$ ,  $K \neq I$ ;
- Locally Optimal Block Preconditioned CG (**LOBPCG**):  $m = 2$ ,  $n_b > 1$ ,  $K \neq I$ ;
- Locally Optimal Extended CG (LOECG):  $m > 2$ ,  $n_b = 1$ ,  $K = I$ ;
- Locally Optimal Preconditioned Extended CG (LOPECG):  $m > 2$ ,  $n_b = 1$ ,  $K \neq I$ ;
- Locally Block Optimal Preconditioned Extended CG (LOBPECG):  $m > 2$ ,  $n_b > 1$ ,  $K \neq I$ .

## Extended Locally Block Optimal Preconditioned CG

Given an initial approximation  $X_0 \in \mathbb{C}^{n \times n_b}$  with  $\text{rank}(X_0) = n_b$ , and an integer  $m \geq 2$ , the algorithm attempts to compute approximate eigenpairs to  $(\lambda_j, u_j)$  for  $1 \leq j \leq n_b$ .

- 1: compute the eigen-decomposition:  $(X_0^H A X_0)W = (X_0^H B X_0)W \Omega_0$ , where  $W^H (X_0^H B X_0) W = I$ ,  $\Omega_0 = \text{diag}(\rho_{0;1}, \rho_{0;2}, \dots, \rho_{0;n_b})$ ;
- 2:  $X_0 = X_0 W$ , and  $X_{-1} = 0$ ;
- 3: **for**  $\ell = 0, 1, \dots$  **do**
- 4:   test convergence and lock up the converged (detail as in EBPSD);
- 5:   construct preconditioners  $K_{\ell;j}$  for  $1 \leq j \leq n_b$ ;
- 6:   compute a basis matrix  $Z \in \mathbb{C}^{n \times (m+1)n_b}$  of the subspace 
$$\sum_{j=1}^{n_b} \mathcal{K}_m(K_{\ell;j}(A - \rho_{\ell;j}B), x_{\ell;j}) + \mathcal{R}(X_{\ell-1});$$
- 7:   compute the  $n_b$  smallest eigenvalues and corresponding eigenvectors of  $Z^H(A - \lambda B)Z$  to get  $(Z^H A Z)W = (Z^H B Z)W \Omega_\ell$ , where  $W^H (Z^H B Z) W = I$ ,  $\Omega_{\ell+1} = \text{diag}(\rho_{\ell+1;1}, \rho_{\ell+1;2}, \dots, \rho_{\ell+1;n_b})$ ;
- 8:    $X_{\ell+1} = ZW$ ;
- 9: **end for**
- 10: **return** approximate eigenpairs to  $(\lambda_j, u_j)$  for  $1 \leq j \leq n_b$ .

Three important implementation issues earlier for XBPSD essentially apply here, but more need to be said about  $Z$  at Line 6 here.

$X_{\ell-1}$  can be replaced by something else, using the idea earlier for LOCG. Specifically, Lines 2, 6, and 8 should be modified to

---

2:  $X_0 = X_0 W$ , and  $Y_0 = 0$ ;

6: compute a basis matrix  $Z \in \mathbb{C}^{n \times (m+1)n_b}$  of the subspace

$$\sum_{j=1}^{n_b} \mathcal{K}_m(K_{\ell;j}(A - \rho_{\ell;j}B), x_{\ell;j}) + \mathcal{R}(Y_\ell) \text{ such that } \mathcal{R}(Z_{(:,1:n_b)}) = \mathcal{R}(X_\ell);$$

8:  $X_{\ell+1} = ZW$ ,  $Y_{\ell+1} = Z\widehat{W}$ , where  $\widehat{W}$  is  $W$  with its  $n_b$  rows zeroed;

---

For  $K_{\ell;j} \equiv K_\ell$ ,  $Z$  is basis matrix of (dropping the subscript  $\ell$ )

$$\mathcal{K}_m(K\mathcal{R}, X) + \mathcal{R}(Y) = \text{span}\{X, K\mathcal{R}(X), \dots, [K\mathcal{R}]^{m-1}(X)\} + \mathcal{R}(Y).$$

- 1 compute a basis matrix  $[Z_1, Z_2, \dots, Z_m]$  for  $\mathcal{K}_m(K\mathcal{R}, X)$  by the Block Arnoldi-like process in the  $B$ -inner product. In particular,  $Z_1 = X$ .
- 2  $B$ -orthogonalize  $Y$  against  $[Z_1, Z_2, \dots, Z_m]$  to get  $Z_{m+1}$  satisfying  $Z_{m+1}^H B Z_{m+1} = I$ .
- 3  $Z = [Z_1, Z_2, \dots, Z_{m+1}]$ .

Precise rates of convergence for various CG methods are scarce and not well understood, especially so for methods of block version. The existing research on the convergence of various SD and CG-type methods, although fragmental and incomplete, should be helpful and provide heuristic insights. Some of the references are



L. Bergamaschi, G. Gambolati, and G. Pini. Asymptotic convergence of conjugate gradient methods for the partial symmetric eigenproblem. *Numer. Linear Algebra Appl.*, 4(2):69–84, 1997.



J. H. Bramble, J. E. Pasciak, and A. V. Knyazev. A subspace preconditioning algorithm for eigenvector/eigenvalue computation. *Adv. in Comput. Math.*, 6:159–189, 1996.



Andrew V. Knyazev and Klaus Neymeyr. A geometric theory for preconditioned inverse iteration III: A short and sharp convergence estimate for generalized eigenvalue problems. *Linear Algebra Appl.*, 358(1-3):95–114, 2003.



D. E. Longsine and S. F. McCormick. Simultaneous Rayleigh-quotient minimization methods for  $Ax = \lambda Bx$ . *Linear Algebra Appl.*, 34:195–234, 1980.



S. Oliveira. On the convergence rate of a preconditioned subspace eigensolver. *Computing*, 63:219–231, 1999.



E. E. Ovtchinnikov. Jacobi correction equation, line search, and conjugate gradients in hermitian eigenvalue computation I: Computing an extreme eigenvalue. *SIAM J. Numer. Anal.*, 46(5):2567–2592, 2008.



E. E. Ovtchinnikov. Jacobi correction equation, line search, and conjugate gradients in hermitian eigenvalue computation II: Computing several extreme eigenvalues. *SIAM J. Numer. Anal.*, 46(5):2593–2619, 2008.



H. Yang. Conjugate gradient methods for the Rayleigh quotient minimization of generalized eigenvalue problems. *Computing*, 51:79–94, 1993.

# Extending Min-Max Principles: Indefinite $B$

- Early Extensions
- Positive Semi-definite Pencil

Min-max principles, Cauchy interlace inequalities are Foundations for optimization approaches to solve few extreme eigenpairs of  $A - \lambda B$  with  $B \succ 0$ .

How far can these theoretical results be extended?

Early extensions (before 1982) of Courant-Fischer min-max principles:

- **$Ax = \lambda x$**  with  $A \succeq 0$  in an indefinite inner product



R. S. Phillips. A minimax characterization for the eigenvalues of a positive symmetric operator in a space with an indefinite metric. *J. Fac. Sci. Univ. Tokyo Sect. I*, 17:51–59, 1970.



B. Textorius. Minimaxprinzip zur Bestimmung der Eigenwerte  $J$ -nichtnegativer Operatoren. *Math. Scand.*, 35:105–114, 1974.

It turns out to be a special case of  $A - \lambda B$  with indefinite and nonsingular  $B$ .

- **Hyperbolic**  $Q(\lambda) = A\lambda^2 + B\lambda + C$ :



R. Duffin. A minimax theory for overdamped networks. *Indiana Univ. Math. J.*, 4:221–233, 1955.

- **More general nonlinear eigenvalue problems:**



H. Voss and B. Werner. A minimax principle for nonlinear eigenvalue problems with applications to nonoverdamped systems. *Math. Meth. Appl. Sci.*, 4:415–424, 1982.

and references therein.

Will focus on  $A - \lambda B$  with indefinite  $B$  and hyperbolic  $Q(\lambda)$  ...

## Difficulties:

- $B$  is indefinite —  $Ax = \lambda Bx$  not equivalent to standard Hermitian eigenvalue problem;
- $Ax = \lambda Bx$  may have complex eigenvalues — no min-max for complex eigenvalues;

*The case we consider will turn out to have only real eigenvalues.*

- $A - \lambda B$  may be a singular pencil:  $\det(A - \lambda B) \equiv 0$  for  $\lambda \in \mathbb{C}$ .  
Need a definition for eigenvalues.

- 1 Extending **Courant-Fischer** for regular Hermitian pencil:
  - Nonsingular  $B$ : Lancaster & Ye (1989), Ye's thesis (1989), Najman & Ye (1991), Binding & Ye (1995)
  - Singular  $B$ : Najman & Ye (1993), Binding & Najman & Ye (1999)Only certain real semi-simple eigenvalues admit a Courant-Fischer type characterization.
- 2 Extending **trace min** for positive semi-definite Hermitian pencil (i.e.,  $A - \lambda_0 B \succeq 0$  for some  $\lambda_0 \in \mathbb{R}$ ):
  - Nonsingular  $B$ : Kovač-Striko & Veselić (1995)
  - (Possibly) singular pencil  $A - \lambda B$ : Liang & Li & Bai (2012)
- 3 Extending **Wielandt's min-max** for positive semi-definite Hermitian pencil:
  - Nonsingular  $B$ : Nakić and Veselić (2003) (actually for regular Hermitian  $A - \lambda B$ , but beware of inaccurate/incorrect statements/equations there)
  - (Possibly) singular pencil  $A - \lambda B$ : Liang & Li (2012)
- 4 Extending **trace min** for linear response eigenvalue problem: Bai & Li (2011).

# Positive Semi-definite Pencil

$$A = A^H, B = B^H \in \mathbb{C}^{n \times n}.$$

- 1 **Positive semi-definite pencil:**  $A - \lambda_0 B \succeq 0$  for some  $\lambda_0 \in \mathbb{R}$ ;  
 $A - \lambda B$  will be assumed so hereafter.
- 2 **Finite eigenvalue**  $\mu (\neq \infty)$ :  $\text{rank}(A - \mu B) < \max_{\lambda \in \mathbb{C}} \text{rank}(A - \lambda B)$ ; This allows singular pencil  $A - \lambda B$ .
- 3 **Eigenvector**  $x$ :  $Ax = \mu Bx$  and  $x \notin \mathcal{N}(A) \cap \mathcal{N}(B)$  and .
- 4  $B$ 's **Inertia**  $(n_+, n_0, n_-)$ :  $n_+$  positive,  $n_0$  zero, and  $n_-$  negative eigenvalues, respectively.
- 5 Can prove: *positive semi-definite pencil*  $A - \lambda B$  has only

$$r := \text{rank}(B) = n_+ + n_-$$

*finite eigenvalues all of which are real:*

$$\lambda_{n_-}^- \leq \dots \leq \lambda_1^- \leq \lambda_0 \leq \lambda_1^+ \leq \dots \leq \lambda_{n_+}^+.$$

- 6 In more detail ...

# Canonical Form of Positive Semi-definite Pencil

1 There exists a nonsingular  $W \in \mathbb{C}^{n \times n}$  such that

$$W^H A W = \begin{matrix} n_1 & & & \\ r-n_1 & & & \\ n-r & & & \end{matrix} \begin{bmatrix} \Lambda_1 & & \\ & \Lambda_0 & \\ & & \Lambda_\infty \end{bmatrix}, \quad W^H B W = \begin{matrix} n_1 & & & \\ r-n_1 & & & \\ n-r & & & \end{matrix} \begin{bmatrix} \Omega_1 & & & \\ & \Omega_0 & & \\ & & & 0 \end{bmatrix},$$

- $\Lambda_1 = \text{diag}(s_1 \alpha_1, \dots, s_\ell \alpha_\ell)$ ,  $\Omega_1 = \text{diag}(s_1, \dots, s_\ell)$ ,  $s_i = \pm 1$ , and  $\Lambda_1 - \lambda_0 \Omega_1 \succ 0$ ;
- $\Lambda_0 = \text{diag}(\Lambda_{0,1}, \dots, \Lambda_{0,m+m_0})$ ,  $\Omega_0 = \text{diag}(\Omega_{0,1}, \dots, \Omega_{0,m+m_0})$ ,

$$\Lambda_{0,i} = t_i \lambda_0, \quad \Omega_{0,i} = t_i = \pm 1, \quad \text{for } 1 \leq i \leq m,$$
$$\Lambda_{0,i} = \begin{bmatrix} 0 & \lambda_0 \\ \lambda_0 & 1 \end{bmatrix}, \quad \Omega_{0,i} = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \quad \text{for } m+1 \leq i \leq m+m_0.$$

There are no such pair  $(\Lambda_0, \Omega_0)$  if  $A - \lambda_0 B \succ 0$ .

- $\Lambda_\infty = \text{diag}(\alpha_{r+1}, \dots, \alpha_n) \succeq 0$  with  $\alpha_i \in \{1, 0\}$  for  $r+1 \leq i \leq n$ .

2  $A - \lambda B$  has  $n_+ + n_-$  finite eigenvalues all of which are real. Denote these finite eigenvalues by  $\lambda_i^\pm$  and arrange them

$$\lambda_{n_-}^- \leq \dots \leq \lambda_1^- \leq \lambda_1^+ \leq \dots \leq \lambda_{n_+}^+.$$

3  $\{\gamma \in \mathbb{R} \mid A - \gamma B \succeq 0\} = [\lambda_{n_-}^-, \lambda_1^+]$ . Moreover, if  $A - \lambda B$  is regular, then  $A - \lambda B$  is a positive definite pencil if and only if  $\lambda_{n_-}^- < \lambda_1^+$ , in which case  $\{\gamma \in \mathbb{R} \mid A - \gamma B \succ 0\} = (\lambda_{n_-}^-, \lambda_1^+)$ .

(Proof in Liang, Li, & Bai, *Linear Algebra and its Applications*, 438 (2013), 3085-3106)

## Courant-Fischer type min-max principle

$$\lambda_i^+ = \sup_{\text{codim } \mathcal{X}=i-1} \inf_{\substack{x \in \mathcal{X} \\ x^H B x > 0}} \frac{x^H A x}{x^H B x}, \quad \lambda_i^+ = \inf_{\text{dim } \mathcal{X}=i} \sup_{\substack{x \in \mathcal{X} \\ x^H B x > 0}} \frac{x^H A x}{x^H B x} \quad \text{for } 1 \leq i \leq n_+,$$

$$\lambda_i^- = \inf_{\text{codim } \mathcal{X}=i-1} \sup_{\substack{x \in \mathcal{X} \\ x^H B x < 0}} \frac{x^H A x}{x^H B x}, \quad \lambda_i^- = \sup_{\text{dim } \mathcal{X}=i} \inf_{\substack{x \in \mathcal{X} \\ x^H B x < 0}} \frac{x^H A x}{x^H B x} \quad \text{for } 1 \leq i \leq n_-.$$

In particular, 
$$\lambda_1^+ = \inf_{x^H B x > 0} \frac{x^H A x}{x^H B x}, \quad \lambda_1^- = \sup_{x^H B x < 0} \frac{x^H A x}{x^H B x}.$$

- Lancaster & Ye (1989), Ye's thesis (1989) for **diagonalizable**  $A - \lambda B$  and  $B$  **nonsingular**. (Actually studied  $A - \lambda B$  not necessarily positive semi-definite, but then only some of the eigenvalues can be characterized.)
- Najman & Ye (1993), Binding & Najman & Ye (1999) for **regular**  $A - \lambda B$ . (Actually studied  $A - \lambda B$  not necessarily positive semi-definite, but then only some of the real eigenvalues can be characterized.)
- Liang & Li (2012) for allowing **singular pencil**  $A - \lambda B$ .

$$k_+ \leq n_+, \quad k_- \leq n_-, \quad k := k_+ + k_- \geq 1, \quad J_k = \begin{bmatrix} I_{k_+} & \\ & -I_{k_-} \end{bmatrix}.$$

## Trace minimization principle

$$\inf_{\substack{X_+ = [x_1, \dots, x_{k_+}] \\ X_- = [y_1, \dots, y_{k_-}] \\ X = [X_+, X_-], \\ \text{subject to (8)}}$$
$$\text{trace}(X^H A X) = \sum_{i=1}^{k_+} \lambda_i^+ - \sum_{i=1}^{k_-} \lambda_i^-.$$

$$\text{either } X^H B X = J_k, \text{ or } X_+^H B X_+ = I_{k_+} \text{ and } X_-^H B X_- = -I_{k_-}. \quad (8)$$

**A converse:**  $\inf_{X^H B X = J_k} \text{trace}(X^H A X) > -\infty \Rightarrow A - \lambda B$  positive semi-definite.

- Kovač-Striko & Veselić (1995) for  $B$  nonsingular, subject to  $X^H B X = J_k$ .
- Liang & Li & Bai (2012) for allowing singular pencil  $A - \lambda B$ .

Unfortunately no **Trace Max** in general.

Eigenvalues of  $A - \lambda B$ :  $\lambda_{n_-}^- \leq \dots \leq \lambda_1^- \leq \lambda_1^+ \leq \dots \leq \lambda_{n_+}^+$ .

$$k_+ \leq n_+, \quad k_- \leq n_-, \quad k := k_+ + k_- \geq 1, \quad J_k = \begin{bmatrix} I_{k_+} & \\ & -I_{k_-} \end{bmatrix};$$

$X \in \mathbb{C}^{k \times k}$ ,  $X^H B X = J_k$ , or the inertia of  $X^H B X$  is  $(k_+, 0, k_-)$ ;

Eigenvalues of  $X^H(A - \lambda B)X$ :  $\mu_{k_-}^- \leq \dots \leq \mu_1^- \leq \mu_1^+ \leq \dots \leq \mu_{k_+}^+$ .

## Cauchy-type interlacing inequality

$$\begin{aligned} \lambda_i^+ &\leq \mu_i^+ \leq \lambda_{i+n-k}^+, & \text{for } 1 \leq i \leq k_+, \\ \lambda_{j+n-k}^- &\leq \mu_j^- \leq \lambda_j^-, & \text{for } 1 \leq j \leq k_-, \end{aligned}$$

where undefined  $\lambda_i^+ = \infty$  for  $i > n_+$  and undefined  $\lambda_j^- = -\infty$  for  $j > n_-$ .

- Kovač-Striko & Veselić (1995) for  $B$  nonsingular.
- Liang & Li & Bai (2012) for allowing singular pencil  $A - \lambda B$ .

These results potentially lead to optimization approaches to compute 1st few  $\lambda_i^\pm$  (these are interior eigenvalues!). See Bai & Li (2011, 2012, 2013 for linear response eigenvalue problem), Kressner, Pandur, & Shao (2013).

# Linear Response Eigenvalue Problem

- Background
- Basic Theory
- Minimization Principles
- 4D SD and 4D CG type Methods

DFT, strictly a ground-state theory, cannot be applied to study the excitations of systems that are involved in **Optical Absorption Spectra** (OAS) calculations.

Runge and Gross (1984) generalized DFT to **Time-Dependent** Density Functional Theory (TD-DFT):

$$i \frac{\partial}{\partial t} \phi_i(\mathbf{r}, t) = \left[ -\frac{1}{2} \nabla^2 + \underbrace{\int \frac{n(\mathbf{r}', t)}{|\mathbf{r} - \mathbf{r}'|} d\mathbf{r}' + \frac{\delta A_{xc}(n(\mathbf{r}, t))}{\delta n(\mathbf{r}, t)} + v_{\text{ext}}(\mathbf{r}, t)}_{v_{\text{KS}}(\mathbf{r}, t)} \right] \phi_i(\mathbf{r}, t).$$

Now KS operator depends on time  $t$ :

$$\text{electronic density: } n(\mathbf{r}, t) = \sum_{i=1}^{N_v} \phi_i(\mathbf{r}, t) \phi_i^*(\mathbf{r}, t).$$



G. Onida, L. Reining and A. Rubio, *Electronic excitations: density-functional versus many-body Green's function approaches*, Rev. Mod. Phys. 74, 2002, (59 pages).

$$\text{DFT: } H_{\text{KS}}^{\text{GS}} \phi_i(\mathbf{r}) \equiv \left[ -\frac{1}{2} \nabla^2 + v_{\text{KS}}(\mathbf{r}) \right] \phi_i(\mathbf{r}) = \lambda_i \phi_i(\mathbf{r}),$$
$$v_{\text{KS}}(\mathbf{r}) = v_{\text{H}}(\mathbf{r}) + v_{\text{xc}}(\mathbf{r}) + v_{\text{ext}}(\mathbf{r}).$$

Perturb  $v_{\text{ext}}(\mathbf{r})$  slightly to  $v_{\text{ext}}(\mathbf{r}, t) = v_{\text{ext}}(\mathbf{r}) + \dot{v}_{\text{ext}}(\mathbf{r}, t)$ , which in turn induce perturbations to  $v_{\text{Hxc}}(\mathbf{r}) \equiv v_{\text{H}}(\mathbf{r}) + v_{\text{xc}}(\mathbf{r})$ :

$$v_{\text{Hxc}}(\mathbf{r}, t) = v_{\text{Hxc}}(\mathbf{r}) + \dot{v}_{\text{Hxc}}(\mathbf{r}, t).$$

$$\text{TD-DFT: } i \frac{\partial}{\partial t} \phi_i(\mathbf{r}, t) = H_{\text{KS}}(t) \phi_i(\mathbf{r}, t) \equiv \left[ -\frac{1}{2} \nabla^2 + v_{\text{KS}}(\mathbf{r}, t) \right] \phi_i(\mathbf{r}, t),$$
$$H_{\text{KS}}(t) = -\frac{1}{2} \nabla^2 + v_{\text{H}}(\mathbf{r}, t) + v_{\text{xc}}(\mathbf{r}, t) + v_{\text{ext}}(\mathbf{r}, t)$$
$$= H_{\text{KS}}^{\text{GS}} + \dot{v}_{\text{Hxc}}(\mathbf{r}, t) + \dot{v}_{\text{ext}}(\mathbf{r}, t).$$

Seek information on first order change in  $n(\mathbf{r}, t) = n(\mathbf{r}) + \dot{n}(\mathbf{r}, t)$ :

$$\begin{aligned}\phi_i(\mathbf{r}, t) &= \phi_i(\mathbf{r}) + \dot{\phi}_i(\mathbf{r}, t), \\ n(\mathbf{r}, t) &\equiv n(\mathbf{r}) + \dot{n}(\mathbf{r}, t) \\ &= n(\mathbf{r}) + \sum_{i=1}^{N_v} \left[ \dot{\phi}_i^*(\mathbf{r}, t) \phi_i(\mathbf{r}) + \phi_i^*(\mathbf{r}) \dot{\phi}_i(\mathbf{r}, t) \right].\end{aligned}$$

Better to explain using the single-particle density matrix which reads

$$\begin{aligned}\rho(\mathbf{r}, t) &= \sum_{i=1}^{N_v} |\phi_i(\mathbf{r}, t)\rangle \langle \phi_i(\mathbf{r}, t)| = \rho(\mathbf{r}) + \dot{\rho}(\mathbf{r}, t), \\ \rho(\mathbf{r}) &= \sum_{i=1}^{N_v} |\phi_i(\mathbf{r})\rangle \langle \phi_i(\mathbf{r})|, \\ \dot{\rho}(\mathbf{r}, t) &= \sum_{i=1}^{N_v} \left( |\dot{\phi}_i(\mathbf{r}, t)\rangle \langle \phi_i(\mathbf{r})| + |\phi_i(\mathbf{r})\rangle \langle \dot{\phi}_i(\mathbf{r}, t)| \right).\end{aligned}$$

(For Dirac **Bra-ket** notation, google bra-ket.)

Differentiate  $\rho(\mathbf{r}, t)$  with respect to  $t$  to get

$$\begin{aligned}i\frac{\partial}{\partial t}\rho(\mathbf{r}, t) &= \sum_{i=1}^{N_v} [H_{\text{KS}}(t) |\phi_i(\mathbf{r}, t)\rangle \langle \phi_i(\mathbf{r}, t)| - |\phi_i(\mathbf{r}, t)\rangle \langle \phi_i(\mathbf{r}, t)| H_{\text{KS}}(t)] \\ &= [H_{\text{KS}}(t), \rho(\mathbf{r}, t)].\end{aligned}$$

Substitute  $H_{\text{KS}}(t) = H_{\text{KS}}^{\text{GS}} + \hat{v}_{\text{Hxc}}(\mathbf{r}, t) + \hat{v}_{\text{ext}}(\mathbf{r}, t)$  and  $\rho(\mathbf{r}, t) = \rho(\mathbf{r}) + \dot{\rho}(\mathbf{r}, t)$  to get

$$\begin{aligned}i\frac{\partial}{\partial t}\dot{\rho}(\mathbf{r}, t) &= [H_{\text{KS}}^{\text{GS}}, \dot{\rho}(\mathbf{r}, t)] + [\hat{v}_{\text{Hxc}}(\mathbf{r}, t), \rho(\mathbf{r})] + [\hat{v}_{\text{ext}}(\mathbf{r}, t), \rho(\mathbf{r})] \\ &= \mathcal{L}\dot{\rho}(\mathbf{r}, t) + [\hat{v}_{\text{ext}}(\mathbf{r}, t), \rho(\mathbf{r})],\end{aligned}$$

where  $\mathcal{L}$  is the *Liouvillian super-operator*:

$$\mathcal{L}\dot{\rho}(\mathbf{r}, t) := [H_{\text{KS}}^{\text{GS}}, \dot{\rho}(\mathbf{r}, t)] + [\hat{v}_{\text{Hxc}}(\mathbf{r}, t), \rho(\mathbf{r})].$$

$$\begin{aligned}i \frac{\partial}{\partial t} \hat{\rho}(\mathbf{r}, t) &= \left[ H_{KS}^{GS}, \hat{\rho}(\mathbf{r}, t) \right] + [\hat{v}_{Hxc}(\mathbf{r}, t), \rho(\mathbf{r})] + [\hat{v}_{ext}(\mathbf{r}, t), \rho(\mathbf{r})] \\ &= \mathcal{L} \hat{\rho}(\mathbf{r}, t) + [\hat{v}_{ext}(\mathbf{r}, t), \rho(\mathbf{r})],\end{aligned}$$

Apply the Fourier transformation to get

$$\begin{aligned}\omega \hat{\rho}(\mathbf{r}, \omega) &= \left[ H_{KS}^{GS}, \hat{\rho}(\mathbf{r}, \omega) \right] + [\hat{v}_{Hxc}(\mathbf{r}, \omega), \rho(\mathbf{r})] + [\hat{v}_{ext}(\mathbf{r}, \omega), \rho(\mathbf{r})], \\ &= \mathcal{L} \hat{\rho}(\mathbf{r}, \omega) + [\hat{v}_{ext}(\mathbf{r}, \omega), \rho(\mathbf{r})],\end{aligned}$$

where  $\mathcal{L} \hat{\rho}(\mathbf{r}, \omega) = \left[ H_{KS}^{GS}, \hat{\rho}(\mathbf{r}, \omega) \right] + [\hat{v}_{Hxc}(\mathbf{r}, \omega), \rho(\mathbf{r})]$ . Therefore

$$(\omega - \mathcal{L}) \hat{\rho}(\mathbf{r}, \omega) = [\hat{v}_{ext}(\mathbf{r}, \omega), \rho(\mathbf{r})].$$

Set  $\hat{v}_{ext}(\mathbf{r}, \omega) = 0 \Rightarrow$  an eigenvalue problem; the smallest positive eigenvalues and associated eigenvectors give excitation states.

# Matrix representation of $\mathcal{L}$

Can show:

$$\tilde{\rho}(\mathbf{r}, t) = \begin{array}{c} \mathcal{R}(\phi_1(r)) \\ \vdots \\ \mathcal{R}(\phi_{N_V}(r)) \\ \mathcal{N}(\rho(\mathbf{r})) \end{array} \begin{pmatrix} \mathcal{R}(\phi_1(r)) & \cdots & \mathcal{R}(\phi_{N_V}(r)) & \mathcal{N}(\rho(\mathbf{r})) \\ 0 & \cdots & 0 & \langle \tilde{\phi}_1(\mathbf{r}, t) | \rho^\perp(\mathbf{r}) \rangle \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \cdots & 0 & \langle \tilde{\phi}_{N_V}(\mathbf{r}, t) | \rho^\perp(\mathbf{r}) \rangle \\ \rho^\perp(\mathbf{r}) | \tilde{\phi}_1(\mathbf{r}, t) \rangle & \cdots & \rho^\perp(\mathbf{r}) | \tilde{\phi}_{N_V}(\mathbf{r}, t) \rangle & 0 \end{pmatrix}.$$

Hence basis functions of the “vector space” of all possible  $\tilde{\rho}(\mathbf{r}, t)$

$$x_i(\mathbf{r}, t) = \rho^\perp(\mathbf{r}) | \tilde{\phi}_i(\mathbf{r}, t) \rangle, \quad y_i(\mathbf{r}, t) = \langle \tilde{\phi}_i(\mathbf{r}, t) | \rho^\perp(\mathbf{r}).$$

In the frequency space:

$$\tilde{\rho}(\mathbf{r}, \omega) = \begin{array}{c} \mathcal{R}(\phi_1(r)) \\ \vdots \\ \mathcal{R}(\phi_{N_V}(r)) \\ \mathcal{N}(\rho(\mathbf{r})) \end{array} \begin{pmatrix} \mathcal{R}(\phi_1(r)) & \cdots & \mathcal{R}(\phi_{N_V}(r)) & \mathcal{N}(\rho(\mathbf{r})) \\ 0 & \cdots & 0 & y_1(\mathbf{r}, \omega) \\ \vdots & \ddots & \vdots & \vdots \\ 0 & \cdots & 0 & y_{N_V}(\mathbf{r}, \omega) \\ x_1(\mathbf{r}, \omega) & \cdots & x_{N_V}(\mathbf{r}, \omega) & 0 \end{pmatrix}.$$

$$\mathcal{L} \begin{pmatrix} x_1(\mathbf{r}, \omega) \\ \vdots \\ x_{N_v}(\mathbf{r}, \omega) \\ y_1(\mathbf{r}, \omega) \\ \vdots \\ y_{N_v}(\mathbf{r}, \omega) \end{pmatrix} = \begin{pmatrix} \mathcal{D} + \mathcal{K} & \mathcal{K} \\ -\mathcal{K} & -\mathcal{D} - \mathcal{K} \end{pmatrix} \begin{pmatrix} x_1(\mathbf{r}, \omega) \\ \vdots \\ x_{N_v}(\mathbf{r}, \omega) \\ y_1(\mathbf{r}, \omega) \\ \vdots \\ y_{N_v}(\mathbf{r}, \omega) \end{pmatrix},$$

$$\mathcal{D} = \text{diag} \left( \rho^\perp(\mathbf{r}) H_{\text{KS}}^{\text{GS}} \rho^\perp(\mathbf{r}) - \epsilon_1, \dots, \rho^\perp(\mathbf{r}) H_{\text{KS}}^{\text{GS}} \rho^\perp(\mathbf{r}) - \epsilon_{N_v} \right),$$

$$\mathcal{K} \begin{pmatrix} z_1(\mathbf{r}) \\ \vdots \\ z_{N_v}(\mathbf{r}) \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^{N_v} \rho^\perp(\mathbf{r}) \int \kappa(\mathbf{r}, \mathbf{r}') \phi_i(\mathbf{r}') z_i(\mathbf{r}') d\mathbf{r}' |\phi_1(\mathbf{r})\rangle \\ \vdots \\ \sum_{i=1}^{N_v} \rho^\perp(\mathbf{r}) \int \kappa(\mathbf{r}, \mathbf{r}') \phi_i(\mathbf{r}') z_i(\mathbf{r}') d\mathbf{r}' |\phi_{N_v}(\mathbf{r})\rangle \end{pmatrix}.$$

# Linear Response Eigenvalue Problem

First several smallest positive eigenvalues and corresponding eigenvectors of

$$\mathcal{H} \begin{bmatrix} u \\ v \end{bmatrix} \equiv \begin{bmatrix} A & B \\ -B & -A \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \lambda \begin{bmatrix} u \\ v \end{bmatrix},$$
$$A^T = A, B^T = B \in \mathbb{R}^{n \times n}, \quad \begin{bmatrix} A & B \\ B & A \end{bmatrix} \succ 0.$$

Equivalently,  $H z = \lambda z$ :

$$J = \frac{1}{\sqrt{2}} \begin{bmatrix} I_n & I_n \\ I_n & -I_n \end{bmatrix}, \quad J^T J = J^2 = I_{2n},$$
$$J^T \begin{bmatrix} A & B \\ -B & -A \end{bmatrix} J = \begin{bmatrix} 0 & A - B \\ A + B & 0 \end{bmatrix} =: \begin{bmatrix} 0 & K \\ M & 0 \end{bmatrix} =: H.$$
$$K = A - B \succ 0, \quad M = A + B \succ 0.$$

$H$  non-symmetric, but rich structure to take advantage of.

## For Linear Response Eigenvalue Problem in general



G. Onida, L. Reining, and A. Rubio. Electronic excitations: density-functional versus many-body Green's function approaches. *Rev. Mod. Phys.*, 74(2):601–659, 2002.



Dario Rocca. *Time-Dependent Density Functional Perturbation Theory: New algorithms with Applications to Molecular Spectra*. PhD thesis, The International School for Advanced Studies, Trieste, Italy, 2007.



D. J. Thouless. Vibrational states of nuclei in the random phase approximation. *Nuclear Physics*, 22(1):78–95, 1961.

## Material in what follows on Linear Response Eigenvalue Problem largely taken from



Zhaojun Bai and Ren-Cang Li. Minimization principle for linear response eigenvalue problem, I: Theory. *SIAM J. Matrix Anal. Appl.*, 33(4):1075–1100, 2012.



Zhaojun Bai and Ren-Cang Li. Minimization principle for linear response eigenvalue problem, II: Computation. *SIAM J. Matrix Anal. Appl.*, 34(2):392–416, 2013.



D. Rocca, Z. Bai, R.-C. Li, and G. Galli. A block variational procedure for the iterative diagonalization of non-Hermitian random-phase approximation matrices. *J. Chem. Phys.*, 136:034111, 2012.

# Linear Response Eigenvalue Problem

First several smallest positive eigenvalues, eigenvectors of

$$\mathcal{H} \begin{bmatrix} u \\ v \end{bmatrix} \equiv \begin{bmatrix} A & B \\ -B & -A \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \lambda \begin{bmatrix} u \\ v \end{bmatrix},$$

$A^T = A, B^T = B \in \mathbb{R}^{n \times n}$ ,  $\begin{bmatrix} A & B \\ B & A \end{bmatrix}$  positive definite.

$$J = \frac{1}{\sqrt{2}} \begin{bmatrix} I_n & I_n \\ I_n & -I_n \end{bmatrix}, \quad J^T J = J^2 = I_{2n},$$

$$J^T \begin{bmatrix} A & B \\ -B & -A \end{bmatrix} J = \begin{bmatrix} 0 & A - B \\ A + B & 0 \end{bmatrix} =: \begin{bmatrix} 0 & K \\ M & 0 \end{bmatrix} =: H.$$

$K = A - B, M = A + B \in \mathbb{R}^{n \times n}$  definite because

$$J^T \begin{bmatrix} A & B \\ B & A \end{bmatrix} J = \begin{bmatrix} A + B & 0 \\ 0 & A - B \end{bmatrix} \equiv \begin{bmatrix} M & \\ & K \end{bmatrix}.$$

Eigenvalue problem for  $\mathcal{H}$  – original LR:

$$\mathcal{H} \begin{bmatrix} u \\ v \end{bmatrix} \equiv \begin{bmatrix} A & B \\ -B & -A \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \lambda \begin{bmatrix} u \\ v \end{bmatrix},$$

Eigenvalue problem for  $H$  – transformed LR:

$$Hz \equiv \begin{bmatrix} 0 & K \\ M & 0 \end{bmatrix} \begin{bmatrix} y \\ x \end{bmatrix} = \lambda \begin{bmatrix} y \\ x \end{bmatrix} \equiv \lambda z,$$

Eigenvalue Problems for  $\mathcal{H}$  and  $H$  equivalent:

- Same eigenvalues, and
- Eigenvectors related by

$$\begin{bmatrix} y \\ x \end{bmatrix} = J^T \begin{bmatrix} u \\ v \end{bmatrix}, \quad \begin{bmatrix} u \\ v \end{bmatrix} = J \begin{bmatrix} y \\ x \end{bmatrix}$$

**Problem.**  $H = \begin{bmatrix} 0 & K \\ M & 0 \end{bmatrix}$ ,  $0 \prec K^T = K, 0 \prec M^T = M \in \mathbb{R}^{n \times n}$ .

$KM$  and  $MK$  have positive eigenvalues:

$$0 < \lambda_1^2 \leq \lambda_2^2 \leq \dots \leq \lambda_n^2,$$

where  $0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ , because

$$\text{eig}(KM) = \text{eig}(K^{1/2}K^{1/2}M) = \text{eig}(K^{1/2}MK^{1/2}).$$

$H^2 = \begin{bmatrix} KM & 0 \\ 0 & MK \end{bmatrix} \Rightarrow H$  has eigenvalues

$$-\lambda_n \leq \dots \leq -\lambda_1 < +\lambda_1 \leq \dots \leq +\lambda_n.$$

$\exists X = Y^{-T} \in \mathbb{R}^{n \times n}$  such that

$$K = Y\Lambda^2Y^T, \quad M = XX^T, \quad \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n).$$

Critically important later.

Proof.

- 1) Cholesky decomposition:  $M^{-1} = R^T R$ .
- 2) Eigendecomposition:  $R^{-T} K R^{-1} = Q\Lambda^2 Q^T$ ,  $Q^T Q = I_n$ .
- 3) Finally  $Y = R^T Q$ ,  $X = Y^{-T}$ . □

$H$  is diagonalizable with Eigendecomposition:

$$H \begin{bmatrix} Y\Lambda & Y\Lambda \\ X & -X \end{bmatrix} = \begin{bmatrix} Y\Lambda & Y\Lambda \\ X & -X \end{bmatrix} \begin{bmatrix} \Lambda & \\ & -\Lambda \end{bmatrix}.$$

# Thouless' Minimization Principle

Eigenvalue problem for  $\mathcal{H}$ : special case of *Hamiltonian eigenvalue problem*.

Eigenvalues appear in  $\pm\lambda$  pairs:

$$-\lambda_n \leq \dots \leq -\lambda_1 < +\lambda_1 \leq \dots \leq +\lambda_n.$$

Thouless' Minimization Principle (1961):

$$\lambda_1 = \min_{u,v} \varrho(u, v), \quad \varrho(u, v) = \frac{\begin{bmatrix} u \\ v \end{bmatrix}^T \begin{bmatrix} A & B \\ B & A \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}}{|u^T u - v^T v|}.$$

Many of today's minimization approaches for computing  $\lambda_1$  are results of this principle.

# Thouless' Minimization Principle

Use  $\begin{bmatrix} y \\ x \end{bmatrix} = J^T \begin{bmatrix} u \\ v \end{bmatrix}$ ,  $\begin{bmatrix} u \\ v \end{bmatrix} = J \begin{bmatrix} y \\ x \end{bmatrix}$  to get

Thouless' Minimization Principle (in different form)

$$\lambda_1 = \min_{x,y} \rho(x,y), \quad \rho(x,y) \stackrel{\text{def}}{=} \frac{x^T Kx + y^T My}{2|x^T y|}.$$

Will call both  $\varrho(u,v)$  and  $\rho(x,y)$  *Thouless' Functional*.

Proof.

Recall  $K = Y\Lambda^2 Y^T$ ,  $M = XX^T$ , and  $X = Y^{-T}$ . We have

$$\begin{aligned} \min_{x,y} \frac{x^T Kx + y^T My}{2|x^T y|} &= \min_{x,y} \frac{x^T Y\Lambda^2 Y^T x + y^T Y^{-T} Y^{-1} y}{2|x^T Y Y^{-1} y|} \\ &= \min_{\tilde{x}, \tilde{y}} \frac{\tilde{x}^T \Lambda^2 \tilde{x} + \tilde{y}^T \tilde{y}}{2|\tilde{x}^T \tilde{y}|} \\ &\geq \min_{\tilde{x}, \tilde{y}} \frac{2 \sum_i \lambda_i |\tilde{x}_{(i)} \tilde{y}_{(i)}|}{2|\sum_i \tilde{x}_{(i)} \tilde{y}_{(i)}|} \geq \lambda_1. \end{aligned}$$

Careful analysis  $\Rightarrow$  equality signs realizable, and optimal argument pair produces eigenvector. □

# Previous Work - summary

Four decades' researches by computational (quantum) physicists and chemists and numerical analysts.

Following three eigenvalue problems are equivalent:

$$Hz \equiv \begin{bmatrix} 0 & K \\ M & 0 \end{bmatrix} \begin{bmatrix} y \\ x \end{bmatrix} = \lambda \begin{bmatrix} y \\ x \end{bmatrix} \equiv \lambda z, \quad (\text{Eig-H})$$

$$KM y = \lambda^2 y, \quad (\text{Eig-KM})$$

$$MK x = \lambda^2 x. \quad (\text{Eig-MK})$$

By computational (quantum) physicists and chemists:

- Chi (1970): solve (Eig-KM) through Symmetric Eigenvalue Problem (SEP)  $RKR^T$ , where  $M = R^T R$  (Cholesky decomposition).
- Davidson-type algorithms (1980s & 1990s), Lanczos-like algorithms (1990s & 2000s)
- CG-like algorithms (more recently, based on Thouless' principle)

By numerical analysts:

- Wilkinson (1960s) discussed (Eig-KM) and (Eig-MK). Implemented as LAPACK `xSYGVD`
- GR algorithm for product eigenvalue problems, generalizing well-known QR algorithm (Watkins, Kressner)
- Krylov-Schur, Jacobi-Davidson, Hamiltonian Krylov-Schur-type, symplectic Lanczos, ...

**Trend.** Huge size –  $n$  in the order  $10^6$  or larger; pose tremendous challenge.

Despite four decades' researches, it is still challenging to robustly and efficiently compute first several positive eigenvalues and eigenvectors.

To come:

- New theory for  $H$  that parallels Symmetric Eigenvalue Problem (SEP)
- New algorithms capable of computing first several positive eigenvalues and eigenvectors *simultaneously*.

# Deflating Subspaces

$\mathcal{U}, \mathcal{V} \subseteq \mathbb{R}^n$ , subspaces. Call  $\{\mathcal{U}, \mathcal{V}\}$  a *pair of deflating subspaces* of  $\{K, M\}$  if

$$K\mathcal{U} \subseteq \mathcal{V} \quad \text{and} \quad M\mathcal{V} \subseteq \mathcal{U}.$$

Let  $U \in \mathbb{R}^{n \times k}$ ,  $V \in \mathbb{R}^{n \times k}$ , basis matrices for  $\mathcal{U}$  and  $\mathcal{V}$ , resp.

$\exists K_R, M_R \in \mathbb{R}^{k \times k}$  such that

$$KU = VK_R, \quad MV = UM_R.$$

In fact, for left generalized inverses  $U^\dagger, V^\dagger$  of  $U, V$ , resp.,

For example,  $U^\dagger U = I$  for  $K_R = V^\dagger KU, \quad M_R = U^\dagger MV.$

$$U^\dagger = (U^T U)^{-1} U^T, \quad \text{but we prefer}$$

$$U^\dagger = (V^T U)^{-1} V^T \quad \text{if } (V^T U)^{-1} \text{ exists.}$$

# Basics: Deflating Subspaces

$$KU = VK_R, MV = UM_R \Rightarrow \begin{bmatrix} 0 & K \\ M & 0 \end{bmatrix} \begin{bmatrix} V \\ U \end{bmatrix} = \begin{bmatrix} V \\ U \end{bmatrix} \begin{bmatrix} 0 & K_R \\ M_R & 0 \end{bmatrix}$$

$H_R = \begin{bmatrix} 0 & K_R \\ M_R & 0 \end{bmatrix}$  is a **restriction** of  $H$  onto  $\mathcal{V} \oplus \mathcal{U}$ .

$H_R$  same block structure as  $H$ ; but lose symmetry in  $K, M$ .

Suppose  $W \stackrel{\text{def}}{=} U^T V$  nonsingular. Factorize  $W = W_1^T W_2$ , where  $W_i$  nonsingular. Define

$$H_{SR} = \begin{bmatrix} 0 & W_1^{-T} U^T K U W_1^{-1} \\ W_2^{-T} V^T M V W_2^{-1} & 0 \end{bmatrix},$$

another **restriction** of  $H$  onto  $\mathcal{V} \oplus \mathcal{U}$ , too:

$$H \begin{bmatrix} V W_2^{-1} \\ U W_1^{-1} \end{bmatrix} = \begin{bmatrix} V W_2^{-1} \\ U W_1^{-1} \end{bmatrix} H_{SR}.$$

$H_{SR}$  same block structure as  $H$  and retain symmetry in  $K, M$ . Major role to come.

## Trace Minimization Principle

$$\sum_{i=1}^k \lambda_i = \frac{1}{2} \min_{U^T V = I_k} \text{trace}(U^T K U + V^T M V).$$

If  $\lambda_k < \lambda_{k+1}$ , optimal  $\{\text{span}(U), \text{span}(V)\}$  gives deflating subspaces of  $\{K, M\}$  corresponding to  $\pm\lambda_i$ ,  $1 \leq i \leq k$ .

Quite similar to the Trace Minimization Principle for Symmetric Eigenvalue Problem (SEP) discussed earlier.

A lengthy proof can be found in



Zhaojun Bai and Ren-Cang Li. Minimization principle for linear response eigenvalue problem, I: Theory. *SIAM J. Matrix Anal. Appl.*, 33(4):1075–1100, 2012.

## Cauchy-like Interlacing Inequalities

$U, V \in \mathbb{R}^{n \times k}$ ;  $W = U^T V$  nonsingular;  $W = W_1^T W_2$ ,  $W_i \in \mathbb{R}^{k \times k}$  nonsingular;  
 $\mathcal{U} = \text{span}(U)$ ,  $\mathcal{V} = \text{span}(V)$ ; Eigenvalues of

$$H_{\text{SR}} = \begin{bmatrix} 0 & W_1^{-T} U^T K U W_1^{-1} \\ W_2^{-T} V^T M V W_2^{-1} & 0 \end{bmatrix}.$$

are  $\pm \mu_i$  ( $1 \leq i \leq k$ ), where  $0 \leq \mu_1 \leq \dots \leq \mu_k$ . Then for  $1 \leq i \leq k$

$$\lambda_i \leq \mu_i \leq \min \left\{ \lambda_{i+2(n-k)}, \frac{\sqrt{\min\{\kappa(K), \kappa(M)\}}}{\cos \angle(\mathcal{U}, \mathcal{V})} \lambda_{i+n-k} \right\},$$

where  $\lambda_j = \infty$  for  $j > n$ . If either  $\mathcal{U} = M\mathcal{V}$  or  $\mathcal{V} = K\mathcal{U}$ , then  
 $\lambda_i \leq \mu_i \leq \lambda_{i+n-k}$ .

Quite similar to the Cauchy Interlacing Inequalities for Symmetric Eigenvalue Problem (SEP) discussed earlier.

A lengthy proof can be found in Bai and Li (2012) mentioned in the previous slide.

Seeking “*best possible*” approximations from the suitably built subspaces.

Given  $\{\mathcal{U}, \mathcal{V}\}$ , a pair of subspaces,  $\dim(\mathcal{U}) = \dim(\mathcal{V}) = n_b$ .

Minimization principles motivate us to seek

- the best approximation to  $\lambda_1$  in the sense of

$$\min_{x \in \mathcal{U}, y \in \mathcal{V}} \rho(x, y)$$

and its associated approximate eigenvector;

- the best approximations to  $\lambda_j$  ( $1 \leq j \leq k$ ) in the sense of

$$\frac{1}{2} \min_{\substack{\text{span}(\hat{U}) \subseteq \mathcal{U}, \text{span}(\hat{V}) \subseteq \mathcal{V} \\ \hat{U}^T \hat{V} = I_k}} \text{trace}(\hat{U}^T K \hat{U} + \hat{V}^T M \hat{V})$$

and their associated approximate eigenvectors. Necessarily  $k \leq n_b$ .

# Best Approximation: $\lambda_1$

$U, V \in \mathbb{R}^{n \times n_b}$ , basis matrices of  $\mathcal{U}$  and  $\mathcal{V}$ . Assume  $W = U^T V$  nonsingular.  
Factorize  $W = W_1^T W_2$ ,  $W_i \in \mathbb{R}^{n_b \times n_b}$  nonsingular.

$x \in \mathcal{U}, y \in \mathcal{V} \Leftrightarrow x = U\hat{u}, y = V\hat{v}$  for  $\hat{u}, \hat{v} \in \mathbb{R}^{n_b}$ .

$$\begin{aligned}\rho(x, y) &= \frac{x^T K x + y^T M y}{2|x^T y|} = \frac{\hat{u}^T U^T K U \hat{u} + \hat{v}^T V^T M V \hat{v}}{2|\hat{u}^T W \hat{v}|} \\ &= \frac{\hat{x}^T W_1^{-T} U^T K U W_1^{-1} \hat{x} + \hat{y}^T W_2^{-T} V^T M V W_2^{-1} \hat{y}}{2|\hat{x}^T \hat{y}|},\end{aligned}$$

where  $\hat{x} = W_1 \hat{u}$  and  $\hat{y} = W_2 \hat{v}$ .

$$\min_{x \in \mathcal{U}, y \in \mathcal{V}} \rho(x, y) = \min_{\hat{x}, \hat{y}} \frac{\hat{x}^T W_1^{-T} U^T K U W_1^{-1} \hat{x} + \hat{y}^T W_2^{-T} V^T M V W_2^{-1} \hat{y}}{2|\hat{x}^T \hat{y}|}$$

which is the smallest positive eigenvalue of

$$H_{\text{SR}} = \begin{bmatrix} 0 & W_1^{-T} U^T K U W_1^{-1} \\ W_2^{-T} V^T M V W_2^{-1} & 0 \end{bmatrix}.$$

# Best Approximation: $\lambda_j$ ( $1 \leq j \leq k$ )

$$\left\{ \begin{array}{l} \text{span}(\hat{U}) \subseteq \mathcal{U}, \text{span}(\hat{V}) \subseteq \mathcal{V}, \\ \hat{U}^T \hat{V} = I_k \end{array} \right\} \Leftrightarrow \left\{ \begin{array}{l} \hat{U} = UW_1^{-1} \hat{X}, \hat{V} = VW_2^{-1} \hat{Y}, \\ \hat{X}, \hat{Y} \in \mathbb{R}^{n_b \times k}, \hat{X}^T \hat{Y} = I_k \end{array} \right\}$$

$$\hat{U}^T K \hat{U} + \hat{V}^T M \hat{V} = \hat{X}^T W_1^{-T} U^T K U W_1^{-1} \hat{X} + \hat{Y}^T W_2^{-T} V^T M V W_2^{-1} \hat{Y}.$$

$$\begin{aligned} & \min_{\substack{\text{span}(\hat{U}) \subseteq \mathcal{U}, \text{span}(\hat{V}) \subseteq \mathcal{V} \\ \hat{U}^T \hat{V} = I_k}} \text{trace}(\hat{U}^T K \hat{U} + \hat{V}^T M \hat{V}) \\ &= \min_{\hat{X}^T \hat{Y} = I_k} \text{trace}(\hat{X}^T W_1^{-T} U^T K U W_1^{-1} \hat{X} + \hat{Y}^T W_2^{-T} V^T M V W_2^{-1} \hat{Y}). \end{aligned}$$

which is the sum of 1st  $k$  smallest positive eigenvalue of

$$H_{\text{SR}} = \begin{bmatrix} 0 & W_1^{-T} U^T K U W_1^{-1} \\ W_2^{-T} V^T M V W_2^{-1} & 0 \end{bmatrix}.$$

# Best Approximation: Eigenvectors

Positive eigenvalues of  $H_{\text{SR}}$ :  $0 \leq \rho_1 \leq \dots \leq \rho_{n_b}$ . Associated eigenvectors  $\hat{\mathbf{z}}_j$ .

$$H_{\text{SR}}\hat{\mathbf{z}}_j = \rho_j\hat{\mathbf{z}}_j, \quad \hat{\mathbf{z}}_j = \begin{bmatrix} \hat{y}_j \\ \hat{x}_j \end{bmatrix}.$$

Then  $\rho(UW_1^{-1}\hat{x}_j, VW_2^{-1}\hat{y}_j) = \rho_j$  for  $j = 1, \dots, n_b$ .

Naturally, take  $\lambda_j \approx \rho_j$ , and corresponding approximate eigenvectors of  $H$ :

$$\tilde{\mathbf{z}}_j \equiv \begin{bmatrix} \tilde{y}_j \\ \tilde{x}_j \end{bmatrix} = \begin{bmatrix} VW_2^{-1}\hat{y}_j \\ UW_1^{-1}\hat{x}_j \end{bmatrix} \quad \text{for } j = 1, \dots, n_b.$$

What if  $U^T V$  is **singular**? Still can do, just more complicated



Zhaojun Bai and Ren-Cang Li. Minimization principle for linear response eigenvalue problem, II: Computation. *SIAM J. Matrix Anal. Appl.*, 34(2):392–416, 2013.

Perturb  $x, y$  to  $\hat{x} = x + p$ ,  $\hat{y} = y + q$ ,  $p$  and  $q$  tiny. Assume  $x^T y \neq 0$ .

Up to the first order in  $p$  and  $q$ ,

$$\begin{aligned} \rho(\hat{x}, \hat{y}) &= \frac{(x + p)^T K(x + p) + (y + q)^T M(y + q)}{2|(x + p)^T (y + q)|} \\ &= \frac{x^T Kx + 2p^T Kx + y^T My + 2q^T My}{2|x^T y + p^T y + q^T x|} \\ &= \frac{x^T Kx + 2p^T Kx + y^T My + 2q^T My}{2|x^T y|} \left[ 1 - \frac{p^T y + q^T x}{x^T y} \right] \\ &= \rho(x, y) + \frac{1}{x^T y} p^T [Kx - \rho(x, y) y] + \frac{1}{x^T y} q^T [My - \rho(x, y) x]. \end{aligned}$$

Partial gradients:  $\nabla_x \rho = \frac{1}{x^T y} [Kx - \rho(x, y) y]$ ,  $\nabla_y \rho = \frac{1}{x^T y} [My - \rho(x, y) x]$ .

Closely related to residual:

$$Hz - \rho(x, y)z \equiv \begin{bmatrix} 0 & K \\ M & 0 \end{bmatrix} \begin{bmatrix} y \\ x \end{bmatrix} - \rho(x, y) \begin{bmatrix} y \\ x \end{bmatrix} = x^T y \begin{bmatrix} \nabla_x \rho \\ \nabla_y \rho \end{bmatrix}.$$

Interested in solving  $\min_{x,y} \rho(x,y) = \min_{x,y} \frac{x^T Kx + y^T My}{2|x^T y|}$  to compute  $\lambda_1$ .

Standard line search: Given current position  $\begin{bmatrix} y \\ x \end{bmatrix}$ , search direction  $\begin{bmatrix} q \\ p \end{bmatrix}$ , seek to minimize  $\rho$  along line

$$\left\{ \begin{bmatrix} y \\ x \end{bmatrix} + t \begin{bmatrix} q \\ p \end{bmatrix} : t \in \mathbb{R} \right\}$$

Doable via Calculus. But not flexible enough to have subspace extensions.

We will do differently.

Minimize  $\rho$  within the *4-dimensional subspace*

$$\left\{ \begin{bmatrix} \beta y + tq \\ \alpha x + sp \end{bmatrix} \text{ for all scalars } \alpha, \beta, s, \text{ and } t \right\}$$

to get

$$\min_{\alpha, \beta, s, t} \rho(\alpha x + sp, \beta y + tq) = \min_{u \in \text{span}(U), v \in \text{span}(V)} \rho(u, v),$$

where  $U = [x, p]$  and  $V = [y, q]$ . Returned to *Best Approximation*.

Naturally take

$$\begin{bmatrix} q \\ p \end{bmatrix} = \begin{bmatrix} \nabla_y \rho \\ \nabla_x \rho \end{bmatrix},$$

as in the *standard steepest descent* (SD) algorithm.

- Lead to plain 4-D SD algorithm for  $H$
- Can design block versions for computing several eigenpairs
- Can incorporate pre-conditioners

All can be viewed as variants of locally optimal 4-D CG algorithms which we will discuss.

Notation:  $\ell$  iteration index;  $j$  eigenpair index.

**Standard:** search next approximations within

$$\text{span} \left\{ \begin{bmatrix} y_j^{(\ell)} \\ x_j^{(\ell)} \end{bmatrix}, \begin{bmatrix} y_j^{(\ell-1)} \\ x_j^{(\ell-1)} \end{bmatrix}, \begin{bmatrix} q_j \\ p_j \end{bmatrix}, \quad j = 1 : k \right\},$$

where

$$\begin{bmatrix} q_j \\ p_j \end{bmatrix} = \Phi \begin{bmatrix} \nabla_x \rho \\ \nabla_y \rho \end{bmatrix} \Big|_{(x,y)=(x_j^{(\ell)}, y_j^{(\ell)})},$$

and  $\Phi$  is a preconditioner to be discussed later.

**We do differently:** search next approximations within

$$\text{span} \left\{ \begin{bmatrix} y_j^{(\ell)} \\ 0 \end{bmatrix}, \begin{bmatrix} y_j^{(\ell-1)} \\ 0 \end{bmatrix}, \begin{bmatrix} q_j \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ x_j^{(\ell)} \end{bmatrix}, \begin{bmatrix} 0 \\ x_j^{(\ell-1)} \end{bmatrix}, \begin{bmatrix} 0 \\ p_j \end{bmatrix} \quad j = 1 : k \right\}.$$

Breaking vectors into two this way is a common technique today in developing structure-preserving alg.:



Kevin J. Kerns and Andrew T. Yang. Preservation of passivity during RLC network reduction via split congruence transformations. *IEEE Trans. Computer-Aided Design of Integrated Circuits and Systems*, 17(7):582–591, July 1998.



R. W. Freund. SPRIM: Structure-preserving reduced-order interconnect macromodeling. In *Proc. Int. Conf. Computer Aided Design*, pages 80–87, Nov. 2004.



Ren-Cang Li and Zhaojun Bai. Structure-preserving model reductions using a Krylov subspace projection formulation. *Commun. Math. Sciences*, 3(2):179–199, 2005.

## Locally Optimal Block Preconditioned 4-D CG

Given an initial approximation  $Z_0 = [X_0^T, Y_0^T]^T \in \mathbb{C}^{2n \times n_b}$  with  $\text{rank}(X_0) = \text{rank}(Y_0) = n_b$ , the algorithm attempts to compute approximate eigenpair to  $(\lambda_j, z_j)$  for  $1 \leq j \leq n_b$ .

- 1: If  $j$ th column of  $Z_0$  isn't an approximation to  $z_j$  already, compute initial approximation with  $\{U, V\} := \{\mathcal{R}(X_0), \mathcal{R}(Y_0)\}$  to give a new  $Z_0$ ;
- 2: **for**  $\ell = 0, 1, \dots$  **do**
- 3:   test convergence and lock up the converged (to discuss later)
- 4:   construct a preconditioner  $\Phi_\ell$ ;
- 5:    $\begin{bmatrix} Q_\ell \\ P_\ell \end{bmatrix} \leftarrow \Phi_\ell \begin{bmatrix} KX_\ell - Y_\ell \Omega_\ell \\ MY_\ell - X_\ell \Omega_\ell \end{bmatrix}$ , where  $\Omega_\ell = \text{diag}(\rho_{\ell,j})$ ;
- 6:    $U = (X_\ell, X_{\ell-1}, P_\ell)$ ,  $V = (Y_\ell, Y_{\ell-1}, Q_\ell)$  (drop  $X_{\ell-1}$  and  $Y_{\ell-1}$  for  $\ell = 0$ );
- 7:   orthogonalize the columns of  $U$  and  $V$  and decompose  $W = U^T V = W_1^T W_2$ ;
- 8:   construct  $H_{\text{SR}}$  (assume  $W$  is nonsingular);
- 9:   compute  $n_b$  smallest positive eigenvalues  $\rho_{\ell+1,j}$  of  $H_{\text{SR}}$ , and associated eigenvectors  $\hat{z}_j$ ;
- 10:    $Z_{\ell+1} := \begin{bmatrix} Y_{\ell+1} \\ X_{\ell+1} \end{bmatrix} = \begin{bmatrix} VW_2^{-1}[\hat{y}_1, \dots, \hat{y}_k] \\ UW_1^{-1}[\hat{x}_1, \dots, \hat{x}_k] \end{bmatrix}$  (normalize each column).
- 11: **end for**
- 12: **return** approximate eigenpairs to  $(\lambda_j, z_j)$  for  $1 \leq j \leq n_b$ .

For convenience, drop iteration index  $\ell$ .

To compute eigenvalues close to  $\mu$ :  $\Phi = (H - \mu I_{2n})^{-1}$ , and

$$\begin{bmatrix} Q \\ P \end{bmatrix} = \Phi R, \quad R = HZ - Z\Omega = \begin{bmatrix} KX - Y\Omega \\ MY - X\Omega \end{bmatrix},$$

one step of the inverse power iteration on the residual.

Interested in the smallest positive eigenvalues, naturally  $\mu = 0$ :

$$\Phi R = \begin{bmatrix} 0 & M^{-1} \\ K^{-1} & 0 \end{bmatrix} R = \begin{bmatrix} M^{-1}[MY - X \text{diag}(\rho_j)] \\ K^{-1}[KX - Y \text{diag}(\rho_j)] \end{bmatrix}.$$

Both  $P$  and  $Q$  computable (column-by-column) by (linear) CG.

In general for  $\mu \neq 0$ , multiplying by  $\Phi$  involves solving linear system:  $(H - \mu I_{2n})z = b$ . Next slides consider this for  $0 < \mu < \lambda_1$ .

# Generic Pre-conditioner: $(H - \mu I_{2n})z = b$

Can verify

$$\begin{bmatrix} I & 0 \\ M & \mu I \end{bmatrix} (H - \mu I) = \begin{bmatrix} I & 0 \\ M & \mu I \end{bmatrix} \begin{bmatrix} -\mu I & K \\ M & -\mu I \end{bmatrix} = \begin{bmatrix} -\mu I & K \\ 0 & MK - \mu^2 I \end{bmatrix}$$

to get

$$(H - \mu I_{2n})^{-1} = \begin{bmatrix} -\mu I & K \\ 0 & MK - \mu^2 I \end{bmatrix}^{-1} \begin{bmatrix} I & 0 \\ M & \mu I \end{bmatrix}$$

Write  $b = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$  and  $z = \begin{bmatrix} y \\ x \end{bmatrix}$ .  $(H - \mu I_{2n})z = b$  can be solved as

solve  $(MK - \mu^2 I)x = Mb_1 - \mu b_2$  for  $x$ , and then  $y = \frac{1}{-\mu}(b_1 - Kx)$ .

Remain to solve  $(MK - \mu^2 I)x = c$  efficiently. Write  $A = MK - \mu^2 I$ , and symbolically

$$A = M^{1/2} \underbrace{(M^{1/2} K M^{1/2} - \mu^2 I)}_{=: \hat{A}} M^{-1/2} = M^{1/2} \hat{A} M^{-1/2}$$

$Ax = c$  is equivalent to  $M^{1/2} \hat{A} M^{-1/2} x = c \Leftrightarrow \underbrace{\hat{A} M^{-1/2} x}_{=: \hat{x}} = \underbrace{M^{-1/2} c}_{=: \hat{c}}$ .

$\hat{A}$  is SPD because  $0 < \mu < \lambda_1$ . Can apply linear CG to  $\hat{A}\hat{x} = \hat{c}$  **symbolically** first and then translate to  $Ax = c$ .

$$Ax \equiv (MK - \mu^2 I)x = c$$

Transform to  $\widehat{A}\widehat{x} = \widehat{c}$ ,  $\widehat{A} = M^{-1/2}AM^{1/2}$ ,  $\widehat{x} = M^{-1/2}x$ ,  $\widehat{c} = M^{-1/2}c$ .

Linear CG to  $\widehat{A}\widehat{x} = \widehat{c}$ :  $\widehat{r}_0 = \widehat{A}\widehat{x}_0 - \widehat{c}$ ,  $\widehat{p}_0 = -\widehat{r}_0$ , and for  $i \geq 0$

$$\widehat{x}_{i+1} = \widehat{x}_i + \alpha_i \widehat{p}_i, \quad \widehat{r}_{i+1} = \widehat{r}_i + \alpha_i \widehat{A}\widehat{p}_i, \quad \widehat{p}_{i+1} = -\widehat{r}_{i+1} + \beta_i \widehat{p}_i$$

$$\text{where } \alpha_i = -\frac{\widehat{p}_i^T \widehat{r}_i}{\widehat{p}_i^T \widehat{A}\widehat{p}_i} = \frac{\widehat{r}_i^T \widehat{r}_i}{\widehat{p}_i^T \widehat{A}\widehat{p}_i}, \quad \beta_i = \frac{\widehat{p}_i^T \widehat{A}\widehat{r}_{i+1}}{\widehat{p}_i^T \widehat{A}\widehat{p}_i} = \frac{\widehat{r}_{i+1}^T \widehat{r}_{i+1}}{\widehat{r}_i^T \widehat{r}_i}.$$

To convert them back to  $x$ -space (so-to-speak): Note

$\widehat{r} := \widehat{A}\widehat{x} - \widehat{c} = M^{-1/2}(Ax - c) =: M^{-1/2}r$ . So  $M^{1/2}\widehat{p}_0 = -r_0$ , and for  $i \geq 0$ ,

$$x_{i+1} = x_i + \alpha_i M^{1/2}\widehat{p}_i, \quad r_{i+1} = r_i + \alpha_i AM^{1/2}\widehat{p}_i, \quad M^{1/2}\widehat{p}_{i+1} = -r_{i+1} + \beta_i M^{1/2}\widehat{p}_i.$$

Two possible choices for  $p$ -vectors (drop subscripts):

$$p = M^{1/2}\widehat{p} \quad (\text{natural}),$$

$$p = M^{-1/2}\widehat{p} \quad (\text{not-so-natural}).$$

Difference in new formulas for  $\alpha_i$  and  $\beta_i$ .

## CG(I) for $Ax \equiv (MK - \mu^2 I)x = c$

Take  $p = M^{1/2}\hat{p}$  (natural). Already  $\hat{r} = M^{-1/2}r$ ,  $\hat{A} = M^{-1/2}AM^{1/2}$ . So

$$\hat{p}^T \hat{r} = p^T M^{-1} r, \quad \hat{p}^T \hat{A} \hat{p} = p^T M^{-1} A p, \quad \hat{r}^T \hat{r} = r^T M^{-1} r, \quad \hat{p}^T \hat{A} \hat{r} = p^T M^{-1} A r.$$

Therefore

$$\alpha_i = -\frac{p_i^T M^{-1} r_i}{p_i^T M^{-1} A p_i} = \frac{r_i^T M^{-1} r_i}{p_i^T M^{-1} A p_i}, \quad \beta_i = \frac{p_i^T M^{-1} A r_{i+1}}{p_i^T M^{-1} A p_i} = \frac{r_{i+1}^T M^{-1} r_{i+1}}{r_i^T M^{-1} r_i}.$$

## CG(I) for $Ax \equiv (MK - \mu^2 I)x = c$

Given an initial approximation  $x_0$ , a relative tolerance `rto1`, the algorithm solves  $Ax \equiv (MK - \mu^2 I)x = c$ .

- 
- 1:  $r_0 = Ax_0 - c$ ,  $p_0 = -r_0$ ;
  - 2: **for**  $i = 0, 1, \dots$  **do**
  - 3:      $q_i = M^{-1}p_i$  by (linear) CG;
  - 4:      $\alpha_i = -(q_i^T r_i)/(q_i^T A p_i)$ ,  $x_{i+1} = x_i + \alpha_i p_i$ ,  $r_{i+1} = r_i + \alpha_i A p_i$ ;
  - 5:     **if**  $\|r_{i+1}\|_1/\|c\|_1 \leq \text{rto1}$ , **BREAK**;
  - 6:      $\beta_i = (q_i^T A r_{i+1})/(q_i^T A p_i)$ ,  $p_{i+1} = -r_{i+1} + \beta_i p_i$ ;
  - 7: **end for**
  - 8: **return** last  $x_i$  as an approximate solution.
-

## CG(II) for $Ax \equiv (MK - \mu^2 I)x = c$

Take  $p = M^{-1/2} \hat{p}$  (not-so-natural). Already  $\hat{r} = M^{-1/2} r$ ,  $\hat{A} = M^{-1/2} A M^{1/2}$ .  
So

$$\hat{p}^T \hat{r} = p^T r, \quad \hat{p}^T \hat{A} \hat{p} = p^T A M p, \quad \hat{r}^T \hat{r} = r^T M^{-1} r, \quad \hat{p}^T \hat{A} \hat{r} = p^T A r.$$

Therefore

$$\alpha_i = -\frac{p_i^T r_i}{p_i^T A M p_i} = \frac{r_i^T M^{-1} r_i}{p_i^T A M p_i}, \quad \beta_i = \frac{p_i^T A r_{i+1}}{p_i^T A M p_i} = \frac{r_{i+1}^T M^{-1} r_{i+1}}{r_i^T M^{-1} r_i}. \quad (9)$$

## CG(II) for $Ax \equiv (MK - \mu^2 I)x = c$

Given an initial approximation  $x_0$ , a relative tolerance `rto1`, the algorithm solves  $Ax \equiv (MK - \mu^2 I)x = c$ .

- 1:  $r_0 = Ax_0 - c$ ,  $q_0 = M^{-1} r_0$  (by linear CG),  $p_0 = -q_0$ ;
- 2: **for**  $i = 0, 1, \dots$  **do**
- 3:   compute  $\alpha_i$  by (9),  $x_{i+1} = x_i + \alpha_i p_i$ ,  $r_{i+1} = r_i + \alpha_i A M p_i$ ;
- 4:   if  $\|r_{i+1}\|_1 / \|c\|_1 \leq \text{rto1}$ , **BREAK**;
- 5:    $q_{i+1} = M^{-1} r_{i+1}$  by (linear) CG;
- 6:   compute  $\beta_i$  by (9),  $p_{i+1} = -q_{i+1} + \beta_i p_i$ ;
- 7: **end for**
- 8: **return** last  $x_i$  as an approximate solution.

Both CG require solving  $Mq = p$  in the inner iteration.

Another alternative is to rewrite symbolically

$$A = K^{-1/2} \underbrace{(K^{1/2}MK^{1/2} - \mu^2 I)}_{=: \hat{A}} K^{1/2} = K^{-1/2} \hat{A} K^{1/2}$$

$$Ax = c \text{ is equivalent to } K^{-1/2} \hat{A} K^{1/2} x = c \Leftrightarrow \underbrace{\hat{A} K^{1/2} x}_{=: \hat{x}} = \underbrace{K^{1/2} c}_{=: \hat{c}}.$$

$\hat{A}$  is SPD because  $0 < \mu < \lambda_1$ . Can apply linear CG to  $\hat{A}\hat{x} = \hat{c}$  **symbolically** first and then translate to  $Ax = c$ .

Detail is omitted.

# Convergence, Deflation (lock)

$(\rho_{\ell;j}, z_{\ell;j})$  is considered acceptable if  $\frac{\|Hz_{\ell;j} - \rho_{\ell;j}z_{\ell;j}\|_2}{\|Hz_{\ell;j}\|_2 + |\rho_{\ell;j}| \|z_{\ell;j}\|_2} \leq \text{rtol}$ .

Usually  $\lambda_j$  are converged to in order, i.e., the smallest eigenvalues emerge first.

**Lock** all acceptable approximate eigenpairs in  $k_{\text{cvgd}} \times k_{\text{cvgd}}$  diagonal matrix  $\mathbf{D}$  for eigenvalues and  $2n \times k_{\text{cvgd}}$  tall matrix  $\mathbf{Z}$  for eigenvectors.

Every time a converged eigenpair is detected, delete the converged  $\rho_{\ell;j}$  and  $z_{\ell;j}$  from  $\Omega_\ell$  and  $Z_\ell$ , respectively, and expand  $\mathbf{D}$  and  $\mathbf{Z}$  to lock up the pair, accordingly.

At the same time, either reduce  $n_b$  by 1 or append a new column to  $Z$  to maintain  $n_b$  unchanged. The latter can be done by computing more than  $n_b$  eigenpairs at Line 9.

**Deflate** to avoid recomputing converged eigenpairs: Write  $\mathbf{Z} = \begin{bmatrix} \mathbf{Y} \\ \mathbf{X} \end{bmatrix}$  and suppose

$$\mathbf{X}^T \mathbf{Y} = I_{k_{\text{cvgd}}}.$$

- Modify  $K$  and  $M$  in form, but not explicitly, to  $K + \zeta \mathbf{Y} \mathbf{Y}^T$  and  $M + \zeta \mathbf{X} \mathbf{X}^T$ , where  $\zeta$  should be selected such that  $\zeta + \lambda_1 \geq \lambda_{k_{\text{cvgd}} + n_b + 1}$ . Here we pre-assume the  $k_{\text{cvgd}}$  converged eigenpairs are indeed those for  $(\lambda_j, z_j)$  for  $1 \leq j \leq k_{\text{cvgd}}$ . This is usually so, but with no guarantee, of course.

# Hyperbolic Quadratic Eigenvalue Problem

- Basics
- Rayleigh Quotients
- Min-Max Principles
- SD and CG type Method

# Hyperbolic $Q(\lambda) = A\lambda^2 + B\lambda + C$ : Basics

$0 \prec A = A^H \in \mathbb{C}^{n \times n}$ , and  $B = B^H, C = C^H \in \mathbb{C}^{n \times n}$ .

$Q(\lambda)$  is *hyperbolic* if

$$(x^H Bx)^2 - 4(x^H Ax)(x^H Cx) > 0 \quad \text{for } 0 \neq x \in \mathbb{C}^n.$$

This type  $Q$  arises, e.g., from dynamical systems that are overly damped.

*Quadratic Eigenvalue Problem* (QEP):

find  $\lambda \in \mathbb{C}, 0 \neq x \in \mathbb{C}^n$  such that  $Q(\lambda)x = 0$ .

$\lambda$ : *quadratic eigenvalue*;  $x$ : *quadratic eigenvector*.

All quadratic eigenvalues of hyperbolic  $Q(\lambda)$  are real:

$$\lambda_n^- \leq \dots \leq \lambda_1^- < \lambda_1^+ \leq \dots \leq \lambda_n^+.$$

For more basic properties of Hyperbolic QEP, see



C.-H. Guo and P. Lancaster. Algorithms for hyperbolic quadratic eigenvalue problems. *Math. Comp.*, 74:1777–1791, 2005.



Nicholas J. Higham, Françoise Tisseur, and Paul M. Van Dooren. Detecting a definite Hermitian pair and a hyperbolic or elliptic quadratic eigenvalue problem, and associated nearness problems. *Linear Algebra Appl.*, 351-352:455–474, 2002.



A.S. Markus. *Introduction to the spectral theory of polynomial operator pencils*. Translations of mathematical monographs, vol. 71. AMS, Providence, RI, 1988.

Given  $x \neq 0$ , consider

$$f(\lambda, x) := x^H Q(\lambda)x = \lambda^2(x^H Ax) + \lambda(x^H Bx) + (x^H Cx) = 0.$$

Always has two distinct real roots (as functions of  $x$ )

$$\rho_{\pm}(x) = \frac{-(x^H Bx) \pm [(x^H Bx)^2 - 4(x^H Ax)(x^H Cx)]^{1/2}}{2(x^H Ax)}.$$

Can show

$$\rho_+(x) \in [\lambda_1^+, \lambda_n^+], \quad \rho_-(x) \in [\lambda_n^-, \lambda_1^-].$$

Reasonable to define  $\rho_{\pm}(x)$  as the *Rayleigh quotients* for the problem.

## Courant-Fischer type min-max principle (Duffin, 1955)

$$\lambda_i^+ = \max_{\text{codim } \mathcal{X}=i-1} \min_{x \in \mathcal{X}} \rho_+(x), \quad \lambda_i^+ = \min_{\text{dim } \mathcal{X}=i} \max_{x \in \mathcal{X}} \rho_+(x),$$
$$\lambda_i^- = \min_{\text{codim } \mathcal{X}=i-1} \max_{x \in \mathcal{X}} \rho_-(x), \quad \lambda_i^- = \max_{\text{dim } \mathcal{X}=i} \min_{x \in \mathcal{X}} \rho_-(x).$$

In particular,

$$\lambda_1^+ = \min_x \rho_+(x), \quad \lambda_n^+ = \max_x \rho_+(x),$$
$$\lambda_n^- = \min_x \rho_-(x), \quad \lambda_1^- = \max_x \rho_-(x).$$

- **Duffin (1955)** (though stated for hyperbolic  $Q$  with  $B \succ 0$ ,  $C \succ 0$ , Duffin's proof works for the general hyperbolic case.)
- **Also Markus (1988)** (mostly about hyperbolic matrix polynomial of any degree),  
**Voss (1982)** (about certain nonlinear  $Q$ ).

$Q(\lambda) = A\lambda^2 + B\lambda + C$  hyperbolic. Its quadratic eigenvalues:

$$\lambda_n^- \leq \cdots \leq \lambda_1^- < \lambda_1^+ \leq \cdots \leq \lambda_n^+.$$

$k \leq n$ ,  $X \in \mathbb{C}^{n \times k}$ ,  $\text{rank}(X) = k$ .  $X^H Q(\lambda) X$  also hyperbolic.

Quadratic eigenvalues of  $X^H Q(\lambda) X$ :

$$\lambda_{k,X}^- \leq \cdots \leq \lambda_{1,X}^- \leq \lambda_{1,X}^+ \leq \cdots \leq \lambda_{k,X}^+.$$

### Trace Min/Max type principle

$$\begin{aligned} \inf_{\text{rank}(X)=k} \sum_{j=1}^k \lambda_{j,X}^+ &= \sum_{j=1}^k \lambda_j^+, & \sup_{\text{rank}(X)=k} \sum_{j=1}^k \lambda_{j,X}^+ &= \sum_{j=1}^k \lambda_{n-k+j}^+, \\ \sup_{\text{rank}(X)=k} \sum_{j=1}^k \lambda_{j,X}^- &= \sum_{j=1}^k \lambda_j^-, & \inf_{\text{rank}(X)=k} \sum_{j=1}^k \lambda_{j,X}^- &= \sum_{j=1}^k \lambda_{n-k+j}^-. \end{aligned}$$

Corollary of a more general Wielandt type max/max principle (Liang & Li, 2013).

$Q(\lambda) = A\lambda^2 + B\lambda + C$  hyperbolic. Its quadratic eigenvalues:

$$\lambda_n^- \leq \dots \leq \lambda_1^- < \lambda_1^+ \leq \dots \leq \lambda_n^+.$$

$k \leq n$ ;  $X \in \mathbb{C}^{k \times k}$ ,  $\text{rank}(X) = k$ ;

Quadratic eigenvalues of  $X^H Q(\lambda) X$ :

$$\mu_k^- \leq \dots \leq \mu_1^- < \mu_1^+ \leq \dots \leq \mu_k^+.$$

## Cauchy-type interlacing inequality

$$\begin{aligned} \lambda_i^+ &\leq \mu_i^+ \leq \lambda_{i+n-k}^+, & i = 1, \dots, k, \\ \lambda_{j+n-k}^- &\leq \mu_j^- \leq \lambda_j^-, & j = 1, \dots, k. \end{aligned}$$

- Veselić (2010).
- Also derivable from **Wielandt type min-max principles** (not presented here) (Liang & Li, 2013).

# Rayleigh-Ritz Procedure for Hyperbolic $Q$

Recall two most important aspects in solving large scale eigenvalue problems: **building good subspaces** and **seeking “best possible” approximations**.

Given  $\mathcal{Y} \in \mathbb{C}^n$  and  $\dim \mathcal{Y} = m$ , find the “best possible” approximations to some of  $Q(\cdot)$ 's quadratic eigenpairs “using  $\mathcal{Y}$ ”.

Can be done by a new “Rayleigh-Ritz” procedure. Let  $Y$  be  $\mathcal{Y}$ 's basis matrix.

## Rayleigh-Ritz Procedure

- 1 Solve the QEP for  $Y^H Q(\lambda) Y$ :  $Y^H Q(\mu_i^\pm) Y y_i^\pm = 0$ , where

$$\mu_m^- \leq \cdots \leq \mu_1^- < \mu_1^+ \leq \cdots \leq \mu_m^+.$$

- 2 Approximate quadratic eigenvalues:  $\mu_i^\pm \approx \lambda_i^\pm$ , approximate quadratic eigenvectors:  $Y y_i^\pm$ .

But in what sense and why are those  $\mu_i^\pm$  and  $Y y_i^\pm$  “best possible”?

**Trace Min/Max principle:**  $\inf_{\text{rank}(X)=k} \sum_{j=1}^k \lambda_{j,X}^+ = \sum_{j=1}^k \lambda_j^+$  suggests that best possible approximations to  $\lambda_i^+$  ( $1 \leq i \leq k$ ) should be gotten so that

$$\sum_{j=1}^k \lambda_{j,X}^+ \text{ is minimized, subject to } \text{span}(X) \subset \mathcal{Y}, \text{rank}(X) = k.$$

The optimal value is  $\sum_{j=1}^k \mu_j^+$ .

Consequently, first few  $\mu_i^+ \approx \lambda_i^+$  are “best possible”. **Surprise:** “interior” eigenvalues are usually hard to compute but this is not the case here.

Similarly to argue for last few  $\mu_j^+ \approx \lambda_{j+n-k}^+$  are “best possible”.

Similarly to argue for first few  $\mu_i^- \approx \lambda_i^-$  are “best possible”. **Surprise:** “interior” eigenvalues are usually hard to compute but this is not the case here.

Similarly to argue for last few  $\mu_j^- \approx \lambda_{j+n-k}^-$  are “best possible”.

# Gradients of Rayleigh quotients $\rho_{\pm}(x)$

Use  $\rho(x)$  for either  $\rho_+(x)$  or  $\rho_-(x)$ , and perturb  $x$  to  $x + p$ ,  $\|p\|$  tiny.

$\rho(x)$  is changed to  $\rho(x + p) = \rho(x) + \eta + O(\|p\|^2)$ . Then

$$[\rho(x) + \eta]^2 (x + p)^H A (x + p) + [\rho(x) + \eta] (x + p)^H B (x + p) + (x + p)^H C (x + p) = 0$$

which gives, upon noticing  $x^H Q(\rho(x))x = 0$ , that

$$\begin{aligned} [2\rho(x) x^H A x + x^H B x] \eta + \rho^H [\rho(x)^2 A x + \rho(x) B x + C x] \\ + [\rho(x)^2 A x + \rho(x) B x + C x]^H p + O(\|p\|^2) = 0 \end{aligned}$$

and thus

$$\eta = - \frac{\rho^H [\rho(x)^2 A x + \rho(x) B x + C x] + [\rho(x)^2 A x + \rho(x) B x + C x]^H p}{2\rho(x) x^H A x + x^H B x}.$$

Therefore the gradient of  $\rho(x)$  at  $x$  is

$$\nabla \rho(x) = - \frac{2[\rho(x)^2 A + \rho(x) B + C]x}{2\rho(x) x^H A x + x^H B x}.$$

Important to notice that  $\nabla \rho(x)$  is parallel to the residual vector

$$r_{\pm}(x) := [\rho_{\pm}(x)^2 A + \rho_{\pm}(x) B + C]x = Q(\rho_{\pm}(x))x.$$

Steepest descent/ascent method for computing one of  $\lambda_1^\pm$  can be readily given.

Fix two parameters “typ” and  $\ell$  with varying ranges as

$$\text{typ} \in \{+, -\}, \quad \ell \in \{1, n\}$$

to mean that we are to compute the quadratic eigenpair  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$ .

A key step of the method is the following line-search problem

$$t_{\text{opt}} = \underset{t \in \mathbb{C}}{\text{argopt}} \rho_{\text{typ}}(x + t p), \quad \text{argopt} = \begin{cases} \text{arg min,} & \text{for } (\text{typ}, \ell) \in \{(-, n), (+, 1)\}, \\ \text{arg max,} & \text{for } (\text{typ}, \ell) \in \{(-, 1), (+, n)\}. \end{cases}$$

where  $x$  is the current approximation to  $u_\ell^{\text{typ}}$ ,  $p$  is the search direction.

Not easy to do: Rayleigh quotient  $\rho_{\text{typ}}$  too complicated, unlike for (linear) eigenvalue problems.

Better way to solve by using min-max principle.

**Line Search** is equivalent to find the best possible approximation within the subspace  $\text{span}([x, p])$ .

Suppose  $x$  and  $p$  are linearly independent and let  $Y = [x, p]$ .

Solve the 2-by-2 hyperbolic QEP for  $Y^H Q(\lambda) Y$  to get its quadratic eigenvalues

$$\mu_2^- \leq \mu_1^- < \mu_1^+ \leq \mu_2^+$$

and corresponding quadratic eigenvector  $y_j^\pm$ .

Table for selecting the next approximate quadratic eigenpair:

(typ, $\ell$ )	current approx.	next approx.
(+, 1)	$(\rho_+(\mathbf{x}), \mathbf{x})$	$(\mu_1^+, Yy_1^+)$
(+, $n$ )	$(\rho_+(\mathbf{x}), \mathbf{x})$	$(\mu_2^+, Yy_2^+)$
(-, 1)	$(\rho_-(\mathbf{x}), \mathbf{x})$	$(\mu_1^-, Yy_1^-)$
(-, $n$ )	$(\rho_-(\mathbf{x}), \mathbf{x})$	$(\mu_2^-, Yy_2^-)$

# Steepest Descent/Ascent method

Basically it is **Line Search** along gradient direction.

## Steepest Descent/Ascent method

Given an initial approximation  $\mathbf{x}_0$  to  $u_\ell^{\text{typ}}$ , and a relative tolerance `rtol`, the algorithm attempts to compute an approximate pair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$  with the prescribed `rtol`.

```
1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|$ ,  $\rho_0 = \rho_{\text{typ}}(\mathbf{x}_0)$ ,  $\mathbf{r}_0 = r_{\text{typ}}(\mathbf{x}_0)$ ;  
2: for  $i = 0, 1, \dots$  do  
3:   if  $\|\mathbf{r}_i\| / (|\rho_i|^2 \|A\mathbf{x}_i\| + |\rho_i| \|B\mathbf{x}_i\| + \|C\mathbf{x}_i\|) \leq \text{rtol}$  then  
4:     BREAK;  
5:   else  
6:     solve QEP for  $Y_i^H Q(\lambda) Y_i$ , where  $Y_i = [\mathbf{x}_i, \mathbf{r}_i]$ ;  
7:     select the next approximate quadratic eigenpair  $(\mu, y) = (\mu_j^{\text{typ}}, Y_i y_j^{\text{typ}})$   
       according to the table;  
8:      $\mathbf{x}_{i+1} = y / \|y\|$ ,  $\rho_{i+1} = \mu$ ,  $\mathbf{r}_{i+1} = r_{\text{typ}}(\mathbf{x}_{i+1})$ ;  
9:   end if  
10: end for  
11: return  $(\rho_i, \mathbf{x}_i)$  as an approximate eigenpair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$ .
```

# Extended Steepest Descent/Ascent method

In **Steepest Descent/Ascent method**, the search space is spanned by

$$\mathbf{x}_i, \mathbf{r}_i = Q(\boldsymbol{\rho}_i)\mathbf{x}_i.$$

It is the second order Krylov subspace  $\mathcal{K}_2(Q(\boldsymbol{\rho}_i), \mathbf{x}_i)$  of  $Q(\boldsymbol{\rho}_i)$  on  $\mathbf{x}_i$ .

One way to improve the method is to use a higher order Krylov subspace

$$\mathcal{K}_m(Q(\boldsymbol{\rho}_i), \mathbf{x}_i) = \text{span}\{\mathbf{x}_i, Q(\boldsymbol{\rho}_i)\mathbf{x}_i, \dots, [Q(\boldsymbol{\rho}_i)]^{m-1}\mathbf{x}_i\}.$$

Let  $Y_i$  be a basis matrix of  $\mathcal{K}_m(Q(\boldsymbol{\rho}_i), \mathbf{x}_i)$ . Solve  $m$ -by- $m$  hyperbolic QEP for  $Y_i^H Q(\lambda) Y_i$  to get its quadratic eigenvalues

$$\mu_m^- \leq \dots \leq \mu_1^- < \mu_1^+ \leq \dots \leq \mu_m^+$$

and corresponding quadratic eigenvectors  $y_j^\pm$ .

Table for selecting the next approximate quadratic eigenpair:

(typ, $\ell$ )	current approx.	next approx.
(+, 1)	$(\rho_+(\mathbf{x}), \mathbf{x})$	$(\mu_1^+, Y_i y_1^+)$
(+, $n$ )	$(\rho_+(\mathbf{x}), \mathbf{x})$	$(\mu_m^+, Y_i y_m^+)$
(-, 1)	$(\rho_-(\mathbf{x}), \mathbf{x})$	$(\mu_1^-, Y_i y_1^-)$
(-, $n$ )	$(\rho_-(\mathbf{x}), \mathbf{x})$	$(\mu_m^-, Y_i y_m^-)$

## Extended Steepest Descent/Ascent method

Given an initial approximation  $\mathbf{x}_0$  to  $u_\ell^{\text{typ}}$ , and a relative tolerance `rto1`, and the search space dimension  $m$ , the algorithm attempts to compute an approximate pair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$  with the prescribed `rto1`.

---

```

1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|$ ,  $\rho_0 = \rho_{\text{typ}}(\mathbf{x}_0)$ ,  $\mathbf{r}_0 = r_{\text{typ}}(\mathbf{x}_0)$ ;
2: for  $i = 0, 1, \dots$  do
3:   if  $\|\mathbf{r}_i\| / (|\rho_i|^2 \|A\mathbf{x}_i\| + |\rho_i| \|B\mathbf{x}_i\| + \|C\mathbf{x}_i\|) \leq \text{rto1}$  then
4:     BREAK;
5:   else
6:     compute a basis matrix  $Y_i$  for  $\mathcal{K}_m(Q(\rho_i), \mathbf{x}_i)$ ;
7:     solve QEP for  $Y_i^H Q(\lambda) Y_i$  to get its quadratic eigenvalues  $\mu_j^\pm$  and
       eigenvectors  $y_j^\pm$ ;
8:     select the next approximate quadratic eigenpair  $(\mu, y) = (\mu_j^{\text{typ}}, Y y_j^{\text{typ}})$ 
       according to the table;
9:      $\mathbf{x}_{i+1} = y / \|y\|$ ,  $\rho_{i+1} = \mu$ ,  $\mathbf{r}_{i+1} = r_{\text{typ}}(\mathbf{x}_{i+1})$ ;
10:  end if
11: end for
12: return  $(\rho_i, \mathbf{x}_i)$  as an approximate eigenpair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$ .

```

---

## Rate of Convergence (Liang & Li, 2013)

$$|\boldsymbol{\rho}_{i+1} - \lambda_\ell^{\text{typ}}| \leq \varepsilon_m^2 |\boldsymbol{\rho}_i - \lambda_\ell^{\text{typ}}| + O(\varepsilon_m |\boldsymbol{\rho}_i - \lambda_\ell^{\text{typ}}|^{3/2} + |\boldsymbol{\rho}_i - \lambda_\ell^{\text{typ}}|^2),$$

where

$$\varepsilon_m = \min_{g \in \mathbb{P}_{m-1}, g(\sigma_1) \neq 0} \max_{i \neq 1} \frac{|g(\sigma_i)|}{|g(\sigma_1)|},$$

and  $\sigma_j$  for  $1 \leq j \leq n$  are eigenvalues of  $Q(\boldsymbol{\rho}_i)$  arranged as in

$$\begin{aligned} \sigma_1 > 0 > \sigma_2 \geq \cdots \geq \sigma_n & \quad \text{if } (\text{typ}, \ell) \in \{(+, 1), (-, 1)\}, \quad \text{or,} \\ \sigma_1 < 0 < \sigma_2 \leq \cdots \leq \sigma_n & \quad \text{if } (\text{typ}, \ell) \in \{(+, n), (-, n)\}. \end{aligned}$$

- While the result is similar to the one for  $A - \lambda B$  ( $B \succ 0$ ), it is *much much* more complicated to prove.
- Important: rate depends on eigenvalue distribution of  $Q(\boldsymbol{\rho}_i)$ . Shed light to preconditioning:

$$Q(\boldsymbol{\rho}_i) \approx L_i D_i L_i^H, \quad D_i = \text{diag}(\pm 1),$$

and use **Extended Steepest Descent/Ascent method** on  $L^{-1}Q(\lambda)L^{-H}$ .

Should reformulate for implementation sake. Detail omitted.

Straightforward applications of ideas presented for  $A - \lambda B$  earlier.  
Left as exercises ...