

# The Hyperbolic Quadratic Eigenvalue Problem

Xin Liang  
Ren-Cang Li

Technical Report 2014-01

# The Hyperbolic Quadratic Eigenvalue Problem

Xin Liang\*      Ren-Cang Li†

January 8, 2014

## Abstract

The hyperbolic quadratic eigenvalue problem (HQEP) was shown to admit the Courant-Fischer type min-max principles in 1955 by Duffin and Cauchy type interlacing inequalities in 2010 by Veselić. It can be regarded as the closest analogue (among all kinds of quadratic eigenvalue problems) to the standard Hermitian eigenvalue problem (among all kinds of standard eigenvalue problems). In this paper, we conduct a systematic study on HQEP both theoretically and numerically. In the theoretic front, we generalize Wiedlandt-Lidskii type min-max principles and, as a special case, Ky-Fan type trace min/max principles and establish Weyl type and Mirsky type perturbation results when an HQEP is perturbed to another HQEP. In the numerical front, we justify the natural generalization of the Rayleigh-Ritz procedure with the existing and our new optimization principles and, as consequences of these principles, we extend various current optimization approaches – steepest descent/ascent and nonlinear conjugate gradient type methods for the Hermitian eigenvalue problem – to calculate few extreme quadratic eigenvalues (of both pos- and neg-type). A detailed convergent analysis is given on the steepest descent/ascent methods. The analysis reveals the intrinsic quantities that control convergence rates and consequently yields ways of constructing effective preconditioners. Numerical examples are presented to demonstrate the proposed theory and algorithms.

**Key words.** Hyperbolic quadratic eigenvalue problem, Rayleigh quotient, min-max principle, Cauchy interlacing inequality, eigenvalue perturbation, extended steepest descent/ascent method, locally optimal extended conjugate gradient method, preconditioning

**AMS subject classifications.** 15A18, 15A42, 65F08, 65F30, 65G99

---

\*School of Mathematical Sciences, Peking University, Beijing, 100871, P. R. China. E-mail: liangxinslm@pku.edu.cn. Supported in part by China Scholarship Council and National Natural Science Foundation of China NSFC-61075119. This work is primarily done while this author was a visiting student, from August 2011 to September 2013, at Department of Mathematics, University of Texas at Arlington, Arlington, TX 76019.

†Department of Mathematics, University of Texas at Arlington, P.O. Box 19408, Arlington, TX 76019. E-mail: rcli@uta.edu. Supported in part by NSF grants DMS-1115834 and DMS-1317330, and a Research Gift Grant from Intel Corporation.

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Hyperbolic quadratic matrix polynomial</b>	<b>5</b>
<b>3</b>	<b>Variational principles</b>	<b>14</b>
3.1	Courant-Fischer type min-max principles . . . . .	15
3.2	Wielandt-Lidskii type min-max principles . . . . .	15
3.3	Cauchy type interlacing inequalities . . . . .	19
<b>4</b>	<b>Perturbation analysis</b>	<b>21</b>
4.1	Setting the stage . . . . .	21
4.2	Asymptotical analysis . . . . .	23
4.3	Perturbation bounds in the spectral norm . . . . .	24
4.4	Perturbation bounds in unitarily invariant norms . . . . .	33
4.5	Perturbation bounds in the Frobenius norms . . . . .	38
<b>5</b>	<b>Best approximations from a subspace and Rayleigh-Ritz procedure</b>	<b>42</b>
<b>6</b>	<b>The steepest descent/ascent method</b>	<b>45</b>
6.1	Gradients . . . . .	45
6.2	The steepest descent/ascent method . . . . .	46
6.3	The extended steepest descent/ascent method . . . . .	48
6.4	Convergence analysis . . . . .	48
6.5	Proof of Theorem 6.2 . . . . .	53
<b>7</b>	<b>Preconditioned steepest descent/ascent method</b>	<b>63</b>
7.1	Preconditioning . . . . .	63
7.2	Preconditioned steepest descent/ascent method . . . . .	65
7.3	Convergence analysis . . . . .	66
<b>8</b>	<b>Block preconditioned steepest descent/ascent method</b>	<b>71</b>
<b>9</b>	<b>Conjugate gradient method</b>	<b>76</b>
9.1	Preconditioned conjugate gradient method . . . . .	76
9.2	Locally optimal block preconditioned extended conjugate gradient method . . . . .	77
<b>10</b>	<b>Numerical examples</b>	<b>79</b>
<b>11</b>	<b>Concluding remarks</b>	<b>81</b>
<b>A</b>	<b>Digression: positive semidefinite matrix pencil</b>	<b>83</b>

# 1 Introduction

It was argued in [26] that the hyperbolic quadratic eigenvalue problem (HQEP) is the closest analogue of the standard Hermitian eigenvalue problem when it comes to the quadratic eigenvalue problem (QEP)

$$(\lambda^2 A + \lambda B + C)x = 0. \quad (1.1)$$

In many ways, both problems share common properties: the eigenvalues are all real, and for HQEP there is a version of the min-max principles [12, 1955] that is very much like the Courant-Fischer min-max principles.

One source of QEPs (1.1) is dynamical systems with friction, where  $A$ ,  $C$  are associated with the kinetic-energy and potential-energy quadratic form, respectively, and  $B$  is associated with the Rayleigh dissipation function [16, 65]. When  $A$ ,  $B$ , and  $C$  are Hermitian, and  $A$  and  $B$  are positive definite and  $C$  positive semidefinite, we say the dynamical system is *overdamped* if

$$(x^H B x)^2 - 4(x^H A x)(x^H C x) > 0 \quad \text{for any nonzero vector } x.$$

Overdamped dynamical systems are common in elevator and car braking systems<sup>1</sup>. A HQEP is slightly more general than an overdamped QEP in that  $B$  and  $C$  are no longer required positive definite or positive semidefinite, respectively. However, a suitable shift in  $\lambda$  can turn a HQEP into an overdamped QEP [20].

If (1.1) is satisfied for a scalar  $\lambda$  and nonzero vector  $x$ , we call  $\lambda$  a *quadratic eigenvalue*,  $x$  an associated *quadratic eigenvector*, and  $(\lambda, x)$  a *quadratic eigenpair*.

In this paper, we will launch a systematic study of the HQEP both in theory and numerical computations that will further reinforce the belief that this class of QEP is the closest analogue to the standard Hermitian eigenvalue problem. In the theoretical front, we will

- review existing results of Courant-Fischer type min-max principles, Cauchy interlacing inequalities;
- establish Wielandt-Lidskii type min-max principles for the sums of selected quadratic eigenvalues and, as corollaries, trace min/max type principles;
- establish perturbation results in the spectral and Frobenius norm, as well as general unitarily invariant norms on how the quadratic eigenvalues will change if  $A$ ,  $B$ ,  $C$  are perturbed.

In the numerical front, we will

- justify a naturally extended Rayleigh-Ritz type procedure, with the existing and newly established min-max principles, why the procedure will produce the best approximations to quadratic eigenvalues/eigenvectors;
- propose extended steepest descent/ascent and CG type methods for computing extreme quadratic eigenpairs;

---

<sup>1</sup>W. Kahan, private communications, November 2013.

- establish convergence results, including the rate of convergence for the extended steepest descent/ascent methods, which shed light on preconditioning in what constitutes a good preconditioner and how to construct one.

In a separate paper, we will extend most of the development in this paper to the hyperbolic polynomial eigenvalue problem.

The rest of this paper is organized as follows. In section 2, we collect some properties for hyperbolic quadratic matrix polynomials and establish a few more about an HQEP. Wielandt-Lidskii type min-max principles, among others, are given in section 3. Eigen-perturbation analysis for HQEP is done in section 4. In section 5, we justify the use of the Rayleigh-Ritz procedure for extracting interested quadratic eigenvalues and their associated quadratic eigenvectors within a given subspace. The steepest descent/ascent method and its extended variation are studied in section 6, where a detailed convergence analysis is performed. Section 7 investigates the preconditioning techniques to speed up the extended steepest descent/ascent method and explain how an effective preconditioner should be constructed from two different perspectives. Section 8 introduces the block variations of the methods in the previous two sections. Various conjugate gradient methods – the plain, locally optimal, and extended subspace search versions combined with suitable preconditioners and blocking – are described in detail in section 9. Two numerical examples are presented in section 10 to demonstrate the effectiveness of the locally optimal block preconditioned conjugate gradient method in the previous section. Finally in section 11, we present our concluding remarks. In appendix section A, we review the Jordan canonical form of a positive semidefinite matrix pencil and establish a perturbation theory for a positive definite matrix pencil for use in section 4.

**Notation.** Throughout this paper,  $\mathbb{C}^{n \times m}$  is the set of all  $n \times m$  complex matrices,  $\mathbb{C}^n = \mathbb{C}^{n \times 1}$ , and  $\mathbb{C} = \mathbb{C}^1$ .  $\mathbb{R}$  is the set of all real numbers.  $I_n$  (or simply  $I$  if its dimension is clear from the context) is the  $n \times n$  identity matrix, and  $e_j$  is its  $j$ th column.  $X^H$  is the conjugate transpose of a vector or matrix. For  $X \in \mathbb{C}^{n \times m}$ ,  $\sigma_{\min}(X)$  is the smallest singular value of  $X$  ( $X$  has  $\min\{m, n\}$  singular values),  $\|X\|_2$  and  $\|X\|_F$  and  $\|X\|_{\text{ui}}$  are the spectral, Frobenius, and a general unitarily invariant norm of  $X$ , and  $\kappa_2(X) = \|X\|_2 \|X^{-1}\|_2$  is the condition number of  $X$ .

$A \succ 0$  ( $A \succeq 0$ ) means that  $A$  is Hermitian positive (semi-)definite, and  $A \prec 0$  ( $A \preceq 0$ ) if  $-A \succ 0$  ( $-A \succeq 0$ ).  $A^{1/2} \succeq 0$  is the unique square root of  $A \succeq 0$ .

The integer triplet  $(i_-(H), i_0(H), i_+(H))$  denotes the inertia of an Hermitian matrix  $H$ , meaning that  $H$  has  $i_-(H)$  negative,  $i_0(H)$  zero, and  $i_+(H)$  positive eigenvalues, respectively, and  $\lambda_{\min}(H)$  and  $\lambda_{\max}(H)$  are its smallest and largest eigenvalue.

Generic notation  $\text{eig}(\cdot)$  is the set of all eigenvalues, counting algebraic multiplicities, of a matrix or a matrix pencil, depending on its argument(s):  $\text{eig}(A)$  is for  $A$ , and  $\text{eig}(A, B)$  is for  $A - \lambda B$ . We use  $\text{polyeig}(A_0, A_1, \dots, A_k)$  as MATLAB's function `polyeig` for the set of all polynomial eigenvalues of  $\lambda^k A_k + \dots + \lambda A_1 + A_0$ . Note  $\text{polyeig}(A_0, A_1)$  is not the same of  $\text{eig}(A_0, A_1)$ .

## 2 Hyperbolic quadratic matrix polynomial

Given  $A, B, C \in \mathbb{C}^{n \times n}$ , define

$$\mathbf{Q}(\lambda) := \lambda^2 A + \lambda B + C, \quad (2.1)$$

a quadratic matrix polynomial of order  $n$ .

**Definition 2.1.**  $\mathbf{Q}(\lambda)$  is said *Hermitian* if  $A, B$ , and  $C$  are all Hermitian, *hyperbolic* if it is Hermitian,  $A \succ 0$ , and

$$(x^H B x)^2 - 4(x^H A x)(x^H C x) > 0, \quad \text{for all } 0 \neq x \in \mathbb{C}^n, \quad (2.2)$$

*overdamped* if it is hyperbolic as well as  $B \succ 0, C \succeq 0$ . For a hyperbolic  $\mathbf{Q}(\lambda)$ , define

$$\varsigma(x) := [(x^H B x)^2 - 4(x^H A x)(x^H C x)]^{1/2}, \quad \varsigma_0(x) := \frac{\varsigma(x)}{x^H x}. \quad (2.3)$$

The quadratic eigenvalue problem (QEP) for  $\mathbf{Q}(\cdot)$  is to find  $\lambda \in \mathbb{C}$  and  $0 \neq x \in \mathbb{C}^n$  such that

$$\mathbf{Q}(\lambda)x = 0.$$

When this equation is satisfied,  $\lambda$  is called a *quadratic eigenvalue* and  $x$  the associated *quadratic eigenvector*. Evidently all quadratic eigenvalues of  $\mathbf{Q}(\cdot)$  is the roots of  $\det \mathbf{Q}(\lambda) = 0$  which has  $2n$  (complex) roots, counting multiplicities.

The next theorem summarizes some of the relevant theoretical results on hyperbolic quadratic polynomials. They can be found in Guo and Lancaster [20] which is an excellent gateway to references of origins for these results. Item 3(c) can be found in [64, (0.7)].

**Theorem 2.1.** Let  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$  as in (2.1) be Hermitian with  $A \succ 0$ .

1.  $\mathbf{Q}(\lambda)$  is hyperbolic if and only if there exists  $\lambda_0 \in \mathbb{R}$  such that  $\mathbf{Q}(\lambda_0) \prec 0$ .
2. If  $\mathbf{Q}(\lambda)$  is hyperbolic, then its quadratic eigenvalues are all real.
3. Suppose  $\mathbf{Q}(\lambda)$  is hyperbolic. Denote its quadratic eigenvalues by  $\lambda_i^\pm$  and arrange them in the order of

$$\lambda_1^- \leq \cdots \leq \lambda_n^- < \lambda_1^+ \leq \cdots \leq \lambda_n^+. \quad (2.4)$$

Then

- (a)  $\mathbf{Q}(\lambda) \prec 0$  for all  $\lambda \in (\lambda_n^-, \lambda_1^+)$ ;
- (b)  $\mathbf{Q}(\lambda) \succ 0$  for all  $\lambda \in (-\infty, \lambda_1^-) \cup (\lambda_n^+, +\infty)$ ;
- (c) the inertia of  $\mathbf{Q}(\lambda)$  is  $(n - k, 0, k)$  for  $\lambda \in (\lambda_k^+, \lambda_{k+1}^+)$  or  $\lambda \in (\lambda_{n-k}^-, \lambda_{n+1-k}^-)$  for  $k = 1, \dots, n$ , concluding that  $\mathbf{Q}(\lambda)$  is indefinite for  $\lambda \in (\lambda_1^{\text{typ}}, \lambda_n^{\text{typ}})$ ;
- (d)  $\mathbf{Q}(\lambda)$  is overdamped if and only if  $\lambda_n^+ \leq 0$ .

An immediate consequence of Theorem 2.1 is a test to determine whether  $\mathbf{Q}(\lambda)$  is hyperbolic or not [20]: check if its quadratic eigenvalues are all real and, in the case they are all real, check if  $\mathbf{Q}(\lambda_0) \succ 0$ , where  $\lambda_0 = (\lambda_n^- + \lambda_1^+)/2$ .

A common technique of solving QEP (1.1), or more generally the polynomial eigenvalue problem, is *linearization* that converts a polynomial eigenvalue problem to an equivalent generalized (linear) eigenvalue problem of a matrix pencil [16, 25, 42].

Under the condition that  $A$  is nonsingular, QEP (1.1) is equivalent to the generalized eigenvalue problem of the following matrix pencil

$$\mathcal{L}_{\mathbf{Q}}(\lambda) := \begin{bmatrix} -C & 0 \\ 0 & A \end{bmatrix} - \lambda \begin{bmatrix} B & A \\ A & 0 \end{bmatrix} = \mathcal{A} - \lambda \mathcal{B}, \quad (2.5)$$

or

$$\mathcal{K}_{\mathbf{Q}}(\lambda) := \begin{bmatrix} 0 & -C \\ -C & -B \end{bmatrix} - \lambda \begin{bmatrix} -C & 0 \\ 0 & A \end{bmatrix} = \mathcal{A} - \lambda \mathcal{B} \quad (2.6)$$

in the sense that  $\text{polyeig}(C, B, A) = \text{eig}(\mathcal{A}, \mathcal{B})$  and associated eigenvectors of one can be recovered from those for the other. More can be said if  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$  is hyperbolic. Relevant results are summarized in the following lemma, where item 5 is essentially in [4] (see also [9], [26, Theorem 3.6], and [63, Theorem 5A]).

**Theorem 2.2.** *Let  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$  as in (2.1) and let  $\mathcal{L}_{\mathbf{Q}}(\lambda)$  be as in (2.5). Suppose  $A$  is nonsingular.*

1.  $\text{polyeig}(C, B, A) = \text{eig}(\mathcal{A}, \mathcal{B})$ .
2. If  $A \succ 0$  and  $B$  is Hermitian, then the inertia of  $\mathcal{B}$  is  $(n, 0, n)$ .
3. If  $(\mu, x)$  is an eigenpair of  $\mathbf{Q}(\lambda)$ , then  $(\mu, \begin{bmatrix} x \\ \mu x \end{bmatrix})$  is an eigenpair of  $\mathcal{L}_{\mathbf{Q}}(\lambda)$ .
4. If  $(\mu, \begin{bmatrix} x \\ y \end{bmatrix})$  is an eigenpair of  $\mathcal{L}_{\mathbf{Q}}(\lambda)$ , then  $(\mu, x)$  is an eigenpair of  $\mathbf{Q}(\lambda)$  and  $y = \mu x$ .
5. Suppose  $\mathbf{Q}(\lambda)$  is Hermitian.  $\mathbf{Q}(\lambda)$  is hyperbolic if and only if  $\mathcal{L}_{\mathbf{Q}}(\lambda)$  is a positive definite pencil.
6. Suppose  $\mathbf{Q}(\lambda)$  is hyperbolic, and adopt the notation in item 3 of Theorem 2.1. Then  $\mathcal{L}_{\mathbf{Q}}(\lambda) \succ 0$  for all  $\lambda \in (\lambda_n^-, \lambda_1^+)$ .

*Proof.* Since for any  $\lambda \in \mathbb{C}$ ,

$$\begin{bmatrix} I & 0 \\ -\lambda I & I \end{bmatrix}^T \begin{bmatrix} -\mathbf{Q}(\lambda) & 0 \\ 0 & A \end{bmatrix} \begin{bmatrix} I & 0 \\ -\lambda I & I \end{bmatrix} = \begin{bmatrix} -C - \lambda B & -\lambda A \\ -\lambda A & A \end{bmatrix} = \mathcal{L}_{\mathbf{Q}}(\lambda). \quad (2.7)$$

Thus  $(-1)^n \det \mathbf{Q}(\lambda) \cdot \det A \equiv \det \mathcal{L}_{\mathbf{Q}}(\lambda)$  and item 1 follows. For item 2,  $A \succ 0$  guarantees that there is a nonsingular matrix  $X \in \mathbb{C}^{n \times n}$  such that

$$X^H A X = I_n, \quad X^H B X = \text{diag}(\omega_1, \dots, \omega_n) =: \Omega,$$

where  $\omega_i \in \mathbb{R}$ . We have

$$\begin{bmatrix} X & \\ & X \end{bmatrix}^H \mathcal{B} \begin{bmatrix} X & \\ & X \end{bmatrix} = \begin{bmatrix} \Omega & I_n \\ I_n & 0 \end{bmatrix} \quad (2.8)$$

whose eigenvalues are the union of all the eigenvalues of

$$\begin{bmatrix} \omega_i & 1 \\ 1 & 0 \end{bmatrix} \quad \text{for } i = 1, 2, \dots, n.$$

But the two eigenvalues of each one of these  $2 \times 2$  matrices are

$$\frac{\omega_i - \sqrt{\omega_i^2 + 4}}{2} < 0, \quad \frac{\omega_i + \sqrt{\omega_i^2 + 4}}{2} > 0.$$

Therefore the last matrix in (2.8) has  $n$  positive and  $n$  negative eigenvalues, as expected. Items 3 and 4 can be verified in a straightforward way by using (2.7). Also by using (2.7), we see that  $\text{diag}(-\mathbf{Q}(\lambda), A)$  and  $\mathcal{L}_{\mathbf{Q}}(\lambda)$  are congruent for all  $\lambda \in \mathbb{R}$ , and hence items 5 and 6 follow from items 1 and 3(a) of Theorem 2.1, respectively.  $\square$

One consequence of Theorem 2.2 is that any hyperbolic  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$  gives rise to a positive definite matrix pencil  $\mathcal{L}_{\mathbf{Q}}(\lambda)$  as defined by (2.5) with  $\mathcal{B}$  having inertia  $(n, 0, n)$ . There is a converse to the statement, too.

**Theorem 2.3.** *Let  $L(\lambda) = \mathcal{A} - \lambda \mathcal{B}$  be a positive definite Hermitian pair of order  $2n$ . If the inertia of  $\mathcal{B}$  is  $(n, 0, n)$ , then there exists a hyperbolic  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$  and a nonsingular matrix  $U \in \mathbb{C}^{2n \times 2n}$  such that the following statements are true.*

1. *If  $(\mu, x)$  is a quadratic eigenpair of  $\mathbf{Q}(\lambda)$ , then  $(\mu, U \begin{bmatrix} x \\ \mu x \end{bmatrix})$  is an eigenpair of  $L(\lambda)$ .*
2. *If  $(\mu, \begin{bmatrix} \tilde{x} \\ \tilde{y} \end{bmatrix})$  is an eigenpair of  $L(\lambda)$  and we define  $\begin{bmatrix} x \\ y \end{bmatrix} = U^{-1} \begin{bmatrix} \tilde{x} \\ \tilde{y} \end{bmatrix}$ , where  $x \in \mathbb{C}^n$ , then  $(\mu, x)$  is a quadratic eigenpair of  $\mathbf{Q}(\lambda)$  and  $y = \mu x$ .*

*Proof.* Since  $L(\lambda)$  is positive definite and the inertia of  $\mathcal{B}$  is  $(n, 0, n)$ , by Theorem A.1 there exists a nonsingular matrix  $W$  such that  $W^H \mathcal{A} W = \text{diag}(\Lambda_+, -\Lambda_-)$  and  $W^H \mathcal{B} W = \text{diag}(I, -I)$ , where  $\Lambda_+ = \text{diag}(\lambda_1^+, \dots, \lambda_n^+)$ ,  $\Lambda_- = \text{diag}(\lambda_1^-, \dots, \lambda_n^-)$  and  $\lambda_i^\pm \in \mathbb{R}$  and  $\lambda_i^+ > \lambda_j^-$  for all  $i$  and  $j$ . Set

$$\begin{aligned} A &= I, \quad B = -(\Lambda_+ + \Lambda_-), \quad C = \Lambda_+ \Lambda_-, \\ S &= \begin{bmatrix} \Lambda_- & -I \\ \Lambda_+ & -I \end{bmatrix} \begin{bmatrix} (\Lambda_+ - \Lambda_-)^{-1/2} & 0 \\ 0 & (\Lambda_+ - \Lambda_-)^{-1/2} \end{bmatrix}, \end{aligned}$$

and  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$ . It can be verified that corresponding to this  $\mathbf{Q}(\lambda)$ ,  $\mathcal{L}_{\mathbf{Q}}(\lambda)$  of (2.5) satisfies  $\mathcal{L}_{\mathbf{Q}}(\lambda) = S^H W^H L(\lambda) W S$ . Since  $L(\lambda)$  is positive definite, there is a  $\lambda_0 \in \mathbb{R}$  such that  $L(\lambda_0) \succ 0$  which implies  $\mathcal{L}_{\mathbf{Q}}(\lambda_0) \succ 0$  and thus  $\mathbf{Q}(\lambda_0) \prec 0$  by (2.7). Consequently, this  $\mathbf{Q}(\lambda)$  is hyperbolic by item 1 of Theorem 2.1. Finally take  $U = W S$  for items 1 and 2.  $\square$



**Theorem 2.4.** Let  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$  be hyperbolic. Then for any  $X \in \mathbb{C}^{n \times m}$  satisfying  $X^H A X = I_m$ ,

$$(X^H B X)^2 - 4(X^H C X) \succ 0. \quad (2.9)$$

*Proof.* For any  $y \in \mathbb{C}^m$  with  $\|y\|_2 = 1$ , write  $x = Xy$ . We have

$$\begin{aligned} & y^H [(X^H B X)^2 - 4(X^H C X)] y \\ &= (X^H B X y)^H (X^H B X y) - 4(Xy)^H C (Xy) \\ &= \|y\|_2^2 \cdot \|X^H B X y\|_2^2 - 4(Xy)^H C (Xy) \cdot y^H (X^H A X) y \end{aligned} \quad (2.10)$$

$$\geq [y^H (X^H B X y)]^2 - 4(Xy)^H C (Xy) \cdot (Xy)^H A (Xy) \quad (2.11)$$

$$\begin{aligned} &= (x^H B x)^2 - 4x^H C x \cdot x^H A x \\ &> 0, \end{aligned} \quad (2.12)$$

where we have used  $\|y\|_2 = 1$  and  $X^H A X = I_m$  for (2.10), and used the Cauchy-Bunyakovsky-Schwarz inequality for (2.11). Therefore  $(X^H B X)^2 - 4(X^H C X) \succ 0$  by (2.12).  $\square$

**Theorem 2.5.** Let  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$  be a hyperbolic quadratic matrix polynomial of order  $n$ , and denote by  $\lambda_i^\pm$  its quadratic eigenvalues which are arranged as in (2.4). Set

$$\Lambda_+ = \text{diag}(\lambda_1^+, \dots, \lambda_n^+), \quad \Lambda_- = \text{diag}(\lambda_1^-, \dots, \lambda_n^-). \quad (2.13)$$

Then there exists nonsingular  $Z \in \mathbb{C}^{2n \times 2n}$  of the form

$$Z = \begin{bmatrix} U_+ & U_- \\ U_+ \Lambda_+ & U_- \Lambda_- \end{bmatrix}, \quad (2.14)$$

where  $U_+, U_- \in \mathbb{C}^{n \times n}$  are nonsingular and

$$\mathcal{U} := U_+^{-1} U_- \quad (2.15)$$

is unitary, such that

$$Z^H \mathcal{A} Z = Z^H \begin{bmatrix} -C & \\ & A \end{bmatrix} Z = \begin{bmatrix} \Lambda_+ & \\ & -\Lambda_- \end{bmatrix}, \quad (2.16a)$$

$$Z^H \mathcal{B} Z = Z^H \begin{bmatrix} B & A \\ A & \end{bmatrix} Z = \begin{bmatrix} I_n & \\ & -I_n \end{bmatrix}. \quad (2.16b)$$

Write

$$U_+ = [u_1^+, u_2^+, \dots, u_n^+], \quad U_- = [u_1^-, u_2^-, \dots, u_n^-].$$

As a consequence of (2.14) and (2.16), we have the following statements.

1.  $\mathbf{Q}(\lambda_i^+) u_i^+ = 0$ ,  $\mathbf{Q}(\lambda_i^-) u_i^- = 0$  for  $i = 1, 2, \dots, n$ . Thus there are  $n$  linearly independent quadratic eigenvectors associated with all  $\lambda_i^+$ , and the same can be said about quadratic eigenvectors associated with all  $\lambda_i^-$ .
2.  $\varsigma(u_i^\pm) = 1$  for  $i = 1, 2, \dots, n$ .

3.  $\mathbf{Q}(\lambda)$  admits

$$\mathbf{Q}(\lambda) = U_-^{-\text{H}}(\lambda I - \Lambda_-)U_-^{\text{H}}AU_+(\lambda I - \Lambda_+)U_+^{-1}, \quad (2.17\text{a})$$

$$\mathbf{Q}(\lambda) = U_+^{-\text{H}}(\lambda I - \Lambda_+)U_+^{\text{H}}AU_-(\lambda I - \Lambda_-)U_-^{-1}. \quad (2.17\text{b})$$

4.  $U_-^{\text{H}}AU_+ = (\Lambda_+\Upsilon - \Upsilon\Lambda_-)^{-1}$ . As a result,  $A, B, C$  and  $\mathbf{Q}(\lambda)$  can be expressed in terms of  $\Lambda_{\pm}$  and any two of  $U_+, U_-,$  and  $\Upsilon$ , assuming (2.15). In particular,

$$A = U_+^{-\text{H}}\Theta U_+^{-1}, \quad (2.18\text{a})$$

$$B = U_+^{-\text{H}}(I - \Theta\Lambda_+ - \Lambda_+\Theta)U_+^{-1}, \quad (2.18\text{b})$$

$$C = U_+^{-\text{H}}(\Lambda_+\Theta\Lambda_+ - \Lambda_+)U_+^{-1}, \quad (2.18\text{c})$$

$$\mathbf{Q}(\lambda) = U_+^{-\text{H}}\left[(\lambda I - \Lambda_+)\Theta(\lambda I - \Lambda_+) + (\lambda I - \Lambda_+)\right]U_+^{-1}, \quad (2.18\text{d})$$

where

$$\Theta = (\Lambda_+ - \Upsilon\Lambda_-\Upsilon^{\text{H}})^{-1}. \quad (2.18\text{e})$$

5. We have

$$\|U_+\|_2 = \|U_-\|_2 \leq \frac{\|A^{-1/2}\|_2}{\sqrt{\lambda_1^+ - \lambda_n^-}}, \quad (2.19\text{a})$$

$$\|U_+^{-1}\|_2 = \|U_-^{-1}\|_2 \leq \|A^{1/2}\|_2 \sqrt{\lambda_n^+ - \lambda_1^-}, \quad (2.19\text{b})$$

$$\kappa(U_+) = \kappa(U_-) \leq \sqrt{\kappa(A)} \sqrt{\frac{\lambda_n^+ - \lambda_1^-}{\lambda_1^+ - \lambda_n^-}}, \quad (2.19\text{c})$$

and

$$\|Z\|_2 \leq \Xi \|U_{\pm}\|_2, \quad \|Z^{-1}\|_2 \leq \frac{\Xi}{\lambda_1^+ - \lambda_n^-} \|U_{\pm}^{-1}\|_2, \quad (2.20)$$

where  $\xi_{\pm} = \max\{|\lambda_1^{\pm}|, |\lambda_n^{\pm}|\}$  and

$$\Xi = \frac{2 + \xi_+^2 + \xi_-^2 + \sqrt{[(\xi_+ - 1)^2 + (\xi_- + 1)^2][(\xi_+ + 1)^2 + (\xi_- - 1)^2]}}{2}.$$

The following converse to item 4 is also true: given diagonal matrices  $\Lambda_{\pm}$  as in (2.13) and two of  $U_+, U_-,$  and  $\Upsilon$ , where  $\Upsilon \in \mathbb{C}^{n \times n}$  as in (2.15) is unitary and  $U_+, U_- \in \mathbb{C}^{n \times n}$  are nonsingular, if  $\lambda_i^{\pm}$  can be arranged as in (2.4), then the quadratic matrix polynomial constructed by (2.18) is hyperbolic.

*Proof.* Since  $\mathbf{Q}(\lambda)$  is hyperbolic,  $\mathcal{L}_{\mathbf{Q}}(\lambda)$  in (2.5) is a positive definite pencil. By Theorem A.1, there exists a nonsingular  $Z \in \mathbb{C}^{2n \times 2n}$  to give (2.16). We have to show that  $Z$  must take the form (2.14).

Since each column of  $Z$  is an eigenvector of the pencil  $\mathcal{L}_{\mathbf{Q}}(\lambda)$ , by Theorem 2.2, we conclude that the  $i$ th column of  $Z$  can be expressed as  $\begin{bmatrix} u_i^+ \\ \lambda_i^+ u_i^+ \end{bmatrix}$  for  $1 \leq i \leq n$  and  $\begin{bmatrix} u_j^- \\ \lambda_j^- u_j^- \end{bmatrix}$

for  $1 \leq j = i - n \leq n$ , where  $u_i^+, u_j^-$  are the corresponding quadratic eigenvectors of  $\mathbf{Q}(\lambda)$  associated with  $\lambda_i^+$  and  $\lambda_j^-$ , respectively. Hence  $Z$  takes the form (2.14).

Blockwise, the equations in (2.16) yield

$$U_+^H C U_+ - \Lambda_+ U_+^H A U_+ \Lambda_+ = -\Lambda_+, \quad (2.21a)$$

$$U_-^H C U_- - \Lambda_- U_-^H A U_- \Lambda_- = \Lambda_-, \quad (2.21b)$$

$$U_+^H C U_- - \Lambda_+ U_+^H A U_- \Lambda_- = 0, \quad (2.21c)$$

$$U_+^H B U_+ + U_+^H A U_+ \Lambda_+ + \Lambda_+ U_+^H A U_+ = I, \quad (2.21d)$$

$$U_-^H B U_- + U_-^H A U_- \Lambda_- + \Lambda_- U_-^H A U_- = -I, \quad (2.21e)$$

$$U_+^H B U_- + U_+^H A U_- \Lambda_- + \Lambda_+ U_+^H A U_- = 0. \quad (2.21f)$$

We claim that  $U_+$  is nonsingular. Consider  $U_+ x = 0$  for some  $x \in \mathbb{C}^n$ . We will prove that  $x = 0$  and thus  $U_+$  is nonsingular. By (2.21d),

$$x^H x = x^H I x = x^H (U_+^H B U_+ + U_+^H A U_+ \Lambda_+ + \Lambda_+ U_+^H A U_+) x = 0$$

which implies  $x = 0$ , as was to be shown. Similarly,  $U_-$  is nonsingular.

Next, we define

$$\widehat{\Lambda}_+ := U_+ \Lambda_+ U_+^{-1}, \quad \widehat{\Lambda}_- := U_- \Lambda_- U_-^{-1}. \quad (2.22)$$

We deduce from (2.21c) and (2.21f) the expressions for  $C$  and  $B$  in (2.23a) below, and then use  $C = C^H$  and  $B = B^H$  to get (2.23b).

$$C = \widehat{\Lambda}_-^H A \widehat{\Lambda}_+, \quad B = -A \widehat{\Lambda}_+ - \widehat{\Lambda}_-^H A, \quad (2.23a)$$

$$C = \widehat{\Lambda}_+^H A \widehat{\Lambda}_-, \quad B = -A \widehat{\Lambda}_- - \widehat{\Lambda}_+^H A. \quad (2.23b)$$

Using the second equation in (2.23a), we deduce from (2.21d) and (2.21e) that

$$\begin{aligned} U_+^{-H} U_+^{-1} &= B + A \widehat{\Lambda}_+ + \widehat{\Lambda}_+^H A = (\widehat{\Lambda}_+ - \widehat{\Lambda}_-)^H A, \\ U_-^{-H} U_-^{-1} &= -B - A \widehat{\Lambda}_- - \widehat{\Lambda}_-^H A = A(\widehat{\Lambda}_+ - \widehat{\Lambda}_-). \end{aligned}$$

So  $U_+^{-H} U_+^{-1} = (U_-^{-H} U_-^{-1})^H = U_-^{-H} U_-^{-1}$ . Thus,

$$(U_+^{-1} U_-)^H U_+^{-1} U_- = U_-^H U_+^{-H} U_+^{-1} U_- = I,$$

which infers  $\mathcal{Y} := U_+^{-1} U_-$  is unitary.

Item 1 is straightforward. We now prove item 2 for  $u_i^+$  and the case for  $u_i^-$  can be handled in exactly the same way. Write  $a_i = (u_i^+)^H A u_i^+$ ,  $b_i = (u_i^+)^H B u_i^+$ , and  $c_i = (u_i^+)^H C u_i^+$ . By (2.21a) and (2.21d), we have

$$c_i - (\lambda_i^+)^2 a_i = -\lambda_i^+, \quad b_i + 2a_i \lambda_i^+ = 1$$

solving which for  $c_i$  and  $b_i$  to get

$$b_i^2 - 4a_i c_i = (1 - 2a_i \lambda_i^+)^2 - 4a_i [-\lambda_i^+ + (\lambda_i^+)^2 a_i] = 1.$$

For item 3, we have, by (2.23),

$$\mathbf{Q}(\lambda) = (\lambda I - \widehat{\Lambda}_-^H) A (\lambda I - \widehat{\Lambda}_+), \quad \mathbf{Q}(\lambda) = (\lambda I - \widehat{\Lambda}_+^H) A (\lambda I - \widehat{\Lambda}_-)$$

which, together with (2.22), yield (2.17). For item 4, write  $\Lambda_{-;\mathcal{Y}} = \mathcal{Y}\Lambda_-\mathcal{Y}^H$ , then  $\Lambda_+ - \Lambda_{-;\mathcal{Y}} \succ 0$  because for  $x \neq 0$ ,

$$x^H(\Lambda_+ - \Lambda_{-;\mathcal{Y}})x \geq \lambda_1^+ x^H x - \lambda_n^- x^H \mathcal{Y}^H \mathcal{Y} x = (\lambda_1^+ - \lambda_n^-) x^H x > 0$$

which also implies

$$0 \prec (\Lambda_+ - \Lambda_{-;\mathcal{Y}})^{-1} \preceq (\lambda_1^+ - \lambda_n^-)^{-1} I. \quad (2.24)$$

Substitute  $U_- = U_+\mathcal{Y}$  into (2.21c) to get  $U_+^H C U_+ - \Lambda_+ U_+^H A U_+ \Lambda_{-;\mathcal{Y}} = 0$  and thus by (2.21a), we have

$$\begin{aligned} 0 &= U_+^H C U_+ - \Lambda_+ U_+^H A U_+ \Lambda_+ + \Lambda_+ \\ &= \Lambda_+ U_+^H A U_+ \Lambda_{-;\mathcal{Y}} - \Lambda_+ U_+^H A U_+ \Lambda_+ + \Lambda_+ \\ &= \Lambda_+ \left[ I - U_+^H A U_+ (\Lambda_+ - \Lambda_{-;\mathcal{Y}}) \right]. \end{aligned} \quad (2.25)$$

Substitute  $U_+ = U_-\mathcal{Y}^H$  into (2.21c) to get  $U_-^H C U_- - \Lambda_{+;\mathcal{Y}} U_-^H A U_- \Lambda_- = 0$ , where  $\Lambda_{+;\mathcal{Y}} = \mathcal{Y}^H \Lambda_+ \mathcal{Y}$ . Thus by (2.21b), we have

$$\begin{aligned} 0 &= U_-^H C U_- - \Lambda_- U_-^H A U_- \Lambda_- - \Lambda_- \\ &= \Lambda_{+;\mathcal{Y}} U_-^H A U_- \Lambda_- - \Lambda_- U_-^H A U_- \Lambda_- - \Lambda_- \\ &= - \left[ I - (\Lambda_{+;\mathcal{Y}} - \Lambda_-) U_-^H A U_- \right] \Lambda_-. \end{aligned} \quad (2.26)$$

We note that at least one of  $\Lambda_+$  and  $\Lambda_-$  is nonsingular. If  $\Lambda_+$  is nonsingular, then (2.25) implies

$$U_+^H A U_+ (\Lambda_+ - \Lambda_{-;\mathcal{Y}}) = I \quad \Rightarrow \quad U_+^H A U_+ = (\Lambda_+ - \Lambda_{-;\mathcal{Y}})^{-1}. \quad (2.27)$$

If  $\Lambda_-$  is nonsingular, then (2.26) implies  $(\Lambda_{+;\mathcal{Y}} - \Lambda_-) U_-^H A U_- = I$  which, upon using  $U_- = U_+\mathcal{Y}$ , also implies the second equation in (2.27). Then  $U_+^H A U_+ = (\Lambda_+ \mathcal{Y} - \mathcal{Y} \Lambda_-)^{-1}$ . So,  $U_+^H A U_+ = \Theta$ ,  $U_+^H B U_+ = -\Theta \Lambda_+ - \Lambda_{-;\mathcal{Y}} \Theta$ , and  $U_+^H C U_+ = \Lambda_{-;\mathcal{Y}} \Theta \Lambda_+$ . Noticing

$$\Lambda_{-;\mathcal{Y}} \Theta = -(\Lambda_+ - \Lambda_{-;\mathcal{Y}}) \Theta + \Lambda_+ \Theta = -I + \Lambda_+ \Theta,$$

we have (2.18).

For item 5, the equalities in (2.19) is a consequence of  $U_- = U_+\mathcal{Y}$  and that  $\mathcal{Y}$  is unitary. We now prove (2.19) for  $U_+$ . Use  $(A^{1/2} U_+)^H (A^{1/2} U_+) = \Theta$  to get

$$\|U_+\|_2 \leq \|A^{-1/2}\|_2 \|A^{1/2} U_+\|_2 = \|A^{-1/2}\|_2 \sqrt{\|\Theta\|_2} \leq \frac{\|A^{-1/2}\|_2}{\sqrt{\lambda_1^+ - \lambda_n^-}},$$

and use  $(U_+^{-1} A^{-1/2})(U_+^{-1} A^{-1/2})^H = \Theta^{-1}$  to get

$$\|U_+^{-1}\|_2 \leq \|U_+^{-1} A^{-1/2}\|_2 \|A^{1/2}\|_2 = \sqrt{\|\Theta^{-1}\|_2} \|A^{1/2}\|_2 \leq \|A^{1/2}\|_2 \sqrt{\lambda_n^+ - \lambda_1^-}.$$

They give (2.19a) and (2.19b) for  $U_+$ . Combine (2.19a) and (2.19b) to get (2.19c). For the first inequality in (2.20), we have

$$\|Z\|_2 \leq \left\| \begin{bmatrix} \|U_+\|_2 & \|U_-\|_2 \\ \|U_+\|_2 \xi_+ & \|U_-\|_2 \xi_- \end{bmatrix} \right\|_2 = \|U_+\|_2 \left\| \begin{bmatrix} 1 & 1 \\ \xi_+ & \xi_- \end{bmatrix} \right\|_2 = \|U_+\|_2 \mathcal{E}.$$

For the second inequality, we notice by using  $U_- = U_+ \Upsilon$

$$Z = \begin{bmatrix} U_+ & 0 \\ 0 & U_+ \end{bmatrix} \begin{bmatrix} I & \Upsilon \\ \Lambda_+ & \Upsilon \Lambda_- \end{bmatrix} = \begin{bmatrix} U_+ & 0 \\ 0 & U_+ \end{bmatrix} \begin{bmatrix} I & 0 \\ \Lambda_+ & I \end{bmatrix} \begin{bmatrix} I & \Upsilon \\ 0 & S \end{bmatrix},$$

where  $S = \Upsilon \Lambda_- - \Lambda_+ \Upsilon = -\Theta^{-1} \Upsilon$ . This expression, after some calculations, leads to

$$\begin{aligned} Z^{-1} &= \begin{bmatrix} I & -\Upsilon S^{-1} \\ 0 & S^{-1} \end{bmatrix} \begin{bmatrix} I & 0 \\ -\Lambda_+ & I \end{bmatrix} \begin{bmatrix} U_+^{-1} & 0 \\ 0 & U_+^{-1} \end{bmatrix} \\ &= \begin{bmatrix} \Upsilon S^{-1} \Upsilon \Lambda_- \Upsilon^H & \Upsilon S^{-1} \\ -S^{-1} \Lambda_+ & S^{-1} \end{bmatrix} \begin{bmatrix} U_+^{-1} & 0 \\ 0 & U_+^{-1} \end{bmatrix}, \end{aligned}$$

and thus

$$\|Z^{-1}\|_2 \leq \|S^{-1}\|_2 \left\| \begin{bmatrix} \xi_- & 1 \\ \xi_+ & 1 \end{bmatrix} \right\|_2 \|U_+^{-1}\|_2 = \|U_+^{-1}\|_2 \|\Theta\|_2 \Xi$$

which implies the second inequality in (2.20).

We now prove the converse of item 4. First  $\Theta$  is Hermitian and  $\Theta \succ 0$  by (2.24). Obviously  $A, B, C$  in (2.18) is Hermitian and  $A \succ 0$ . Choose  $\lambda_0 = (\lambda_1^+ + \lambda_n^-)/2$ , then  $\Theta^{-1} \succ \Lambda_+ - \lambda_0 I \succ 0$  and  $\Theta \prec (\Lambda_+ - \lambda_0 I)^{-1}$ . Thus,

$$U_+^H \mathbf{Q}(\lambda_0) U_+ = (\Lambda_+ - \lambda_0 I) \Theta (\Lambda_+ - \lambda_0 I) - (\Lambda_+ - \lambda_0 I) \prec 0$$

which says  $\mathbf{Q}(\lambda_0) \prec 0$ . By item 1 of Theorem 2.1,  $\mathbf{Q}(\lambda)$  is hyperbolic.  $\square$

**Remark 2.1.** 1. Each of the decompositions in (2.17) doesn't reflect the symmetry property in  $\mathbf{Q}(\lambda)$  somewhat. However, using the fact that  $\Upsilon = U_+^{-1} U_-$  is unitary, we can turn them into

$$\mathbf{Q}(\lambda) = U_+^{-H} (\lambda I - \Upsilon \Lambda_- \Upsilon^H) (\Lambda_+ - \Upsilon \Lambda_- \Upsilon^H)^{-1} (\lambda I - \Lambda_+) U_+^{-1}, \quad (2.28a)$$

$$\mathbf{Q}(\lambda) = U_-^{-H} (\lambda I - \Upsilon^H \Lambda_+ \Upsilon) (\Upsilon \Lambda_+ \Upsilon^H - \Lambda_-)^{-1} (\lambda I - \Lambda_-) U_-^{-1}. \quad (2.28b)$$

These equations are essentially the decomposition in [43, Theorem 31.24] but with more detail.

2. [22, Lemma 6.1] and Problem `gen_hyper2` of [5] provide a different set of formulas for  $B$  and  $C$ :

$$B = U_+^{-H} [-\Theta(\Lambda_+^2 - \Upsilon \Lambda_-^2 \Upsilon^H) \Theta] U_+^{-1}, \quad (2.29a)$$

$$\begin{aligned} C &= U_+^{-H} [-\Theta(\Lambda_+^3 - \Upsilon \Lambda_-^3 \Upsilon^H) \Theta \\ &\quad + \Theta(\Lambda_+^2 - \Upsilon \Lambda_-^2 \Upsilon^H) \Theta (\Lambda_+^2 - \Upsilon \Lambda_-^2 \Upsilon^H) \Theta] U_+^{-1}. \end{aligned} \quad (2.29b)$$

[31, Corollary 6] provides yet another formula for  $C$ :

$$C = U_+^{-H} [-(\Lambda_+^{-1} - \Upsilon \Lambda_-^{-1} \Upsilon^H)^{-1}] U_+^{-1}. \quad (2.30)$$

Although both (2.29) and (2.30) look more complicated than ours for  $B$  and  $C$  in (2.18b) and (2.18c), they are actually the same in theory. In fact, we have

$$\Theta(\Lambda_+^2 - \Upsilon \Lambda_-^2 \Upsilon^H) \Theta = \Theta(\Lambda_+^2 - [\Lambda_+ - \Theta^{-1}]^2) \Theta$$

$$= \Lambda_+ \Theta + \Theta \Lambda_+ - I \quad (2.31)$$

which says (2.29a) is the same as (2.18b).

$$\begin{aligned} \Lambda_+^{-1} - \Upsilon \Lambda_-^{-1} \Upsilon^H &= \Lambda_+^{-1} - [\Lambda_+ - \Theta^{-1}]^{-1} \quad (\text{use (2.18e)}) \\ &= \Lambda_+^{-1} (-\Theta^{-1}) [\Lambda_+ - \Theta^{-1}]^{-1} \quad (\text{use } X^{-1} - Y^{-1} = X^{-1} [Y - X] Y^{-1}) \\ &= -(\Lambda_+ \Theta \Lambda_+ - \Lambda_+)^{-1}. \end{aligned}$$

So (2.30) is the same as (2.18c). Finally

$$\begin{aligned} \Theta(\Lambda_+^3 - \Upsilon \Lambda_-^3 \Upsilon^H) \Theta &= \Theta(\Lambda_+^3 - [\Lambda_+ - \Theta^{-1}]^3) \Theta \\ &= \Theta^{-1} + \Theta \Lambda_+^2 + \Lambda_+^2 \Theta + \Theta \Lambda_+ \Theta^{-1} \Lambda_+ \Theta - \Theta \Lambda_+ \Theta^{-1} \\ &\quad - \Theta^{-1} \Lambda_+ \Theta - \Lambda_+. \end{aligned}$$

Therefore use also (2.31) to get

$$\begin{aligned} &-\Theta(\Lambda_+^3 - \Upsilon \Lambda_-^3 \Upsilon^H) \Theta + \Theta(\Lambda_+^2 - \Upsilon \Lambda_-^2 \Upsilon^H) \Theta(\Lambda_+^2 - \Upsilon \Lambda_-^2 \Upsilon^H) \Theta \\ &= -(\Theta^{-1} + \Theta \Lambda_+^2 + \Lambda_+^2 \Theta + \Theta \Lambda_+ \Theta^{-1} \Lambda_+ \Theta - \Theta \Lambda_+ \Theta^{-1} - \Theta^{-1} \Lambda_+ \Theta - \Lambda_+) \\ &\quad + (\Theta \Lambda_+ + \Lambda_+ \Theta - I) \Theta^{-1} (\Theta \Lambda_+ + \Lambda_+ \Theta - I) \\ &= -(\Theta^{-1} + \Theta \Lambda_+^2 + \Lambda_+^2 \Theta + \Theta \Lambda_+ \Theta^{-1} \Lambda_+ \Theta - \Theta \Lambda_+ \Theta^{-1} - \Theta^{-1} \Lambda_+ \Theta - \Lambda_+) \\ &\quad + \Theta^{-1} - \Theta \Lambda_+ \Theta^{-1} - \Lambda_+ - \Lambda_+ + \Theta \Lambda_+^2 + \Lambda_+ \Theta \Lambda_+ \\ &\quad - \Theta^{-1} \Lambda_+ \Theta + \Theta \Lambda_+ \Theta^{-1} \Lambda_+ \Theta + \Lambda_+^2 \Theta \\ &= -\Lambda_+ + \Lambda_+ \Theta \Lambda_+ \end{aligned}$$

which proves that (2.29b) is the same as (2.18c).

3.  $\widehat{\Lambda}_\pm$  in (2.22) are two solutions of the matrix equation

$$AX^2 + BX + C = 0. \quad (2.32)$$

In fact,

$$A(U_+ \Lambda_+ U_+^{-1})^2 + B(U_+ \Lambda_+ U_+^{-1}) + C = (AU_+ \Lambda_+^2 + BU_+ \Lambda_+ + CU_+) U_+^{-1} = 0,$$

and similarly for  $A(U_- \Lambda_- U_-^{-1})^2 + B(U_- \Lambda_- U_-^{-1}) + C = 0$ . On the other hand, the ability of solving (2.32) factorizes  $\mathbf{Q}(\lambda)$  into the product of two linear matrix polynomials, based on which Guo and Lancaster [20] proposed their solvent approach for solving HQEP (1.1) of modest sizes.

### 3 Variational principles

Throughout this section,  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C \in \mathbb{C}^{n \times n}$  will be always assumed a hyperbolic quadratic matrix polynomial and the notations in Theorem 2.5 will be kept. The scalar  $\lambda_0$  is as in item 1 of Theorem 2.1 such that  $\mathbf{Q}(\lambda_0) \prec 0$ .

Consider the following equation in  $\lambda$

$$f(\lambda, x) := x^H \mathbf{Q}(\lambda) x = \lambda^2 (x^H A x) + \lambda (x^H B x) + (x^H C x) = 0, \quad (3.1)$$

given  $x \neq 0$ . Since  $\mathbf{Q}(\lambda)$  is hyperbolic, this equation always has two distinct real roots (as functions of  $x$ )

$$\rho_{\pm}(x) = \frac{-x^H B x \pm [(x^H B x)^2 - 4(x^H A x)(x^H C x)]^{1/2}}{2(x^H A x)}. \quad (3.2)$$

We shall call  $\rho_+(x)$  the *pos-type Rayleigh quotient* of  $\mathbf{Q}(\lambda)$  at  $x$ , and  $\rho_-(x)$  the *neg-type Rayleigh quotient* of  $\mathbf{Q}(\lambda)$  at  $x$ . It is easy to verify that for any  $x \neq 0$ ,  $\rho_{\pm}(x) \in \mathbb{R}$ , and  $\rho_{\pm}(\alpha x) = \rho_{\pm}(x)$  for any  $\alpha \in \mathbb{C}$ . By the elementary knowledge of scalar quadratic polynomials, we have

$$\rho_+(x) + \rho_-(x) = -\frac{x^H B x}{x^H A x}, \quad \rho_+(x) \cdot \rho_-(x) = \frac{x^H C x}{x^H A x}. \quad (3.3)$$

Both will be used later in this paper.

**Theorem 3.1.** *We have*

$$\rho_+(x) \in [\lambda_1^+, \lambda_n^+], \quad \rho_-(x) \in [\lambda_1^-, \lambda_n^-], \quad (3.4)$$

$$\varsigma(x) := [(x^H B x)^2 - 4(x^H A x)(x^H C x)]^{1/2} = \pm [2\rho_{\pm}(x) x^H A x + x^H B x], \quad (3.5)$$

$$\varsigma_0(x) = \frac{\varsigma(x)}{x^H x} \in [(\lambda_1^+ - \lambda_n^-) \lambda_{\min}(A), (\lambda_n^+ - \lambda_1^-) \lambda_{\max}(A)]. \quad (3.6)$$

Consequently,  $\lambda_i^+ = \rho_+(u_i^+)$  for the quadratic eigenpair  $(\lambda_i^+, u_i^+)$  and  $\rho_-(u_j^-) = \lambda_j^-$  for the quadratic eigenpair  $(\lambda_j^-, u_j^-)$ .

*Proof.* By item 3 of Theorem 2.1, for any fixed nonzero  $x$ ,  $f(\lambda, x) < 0$  for  $\lambda \in (\lambda_n^-, \lambda_1^+)$  and  $f(\lambda, x) > 0$  for  $\lambda \in (-\infty, \lambda_1^-) \cup (\lambda_n^+, +\infty)$ . Thus, the larger root of the scalar quadratic equation  $f(\lambda, x) = 0$  in  $\lambda$  must lie in  $[\lambda_1^+, \lambda_n^+]$  and the smaller one in  $[\lambda_1^-, \lambda_n^-]$ . That is (3.4). For (3.5), we have

$$\begin{aligned} 2\rho_{\pm}(x) x^H A x + x^H B x &= [-x^H B x \pm \sqrt{(x^H B x)^2 - 4(x^H A x)(x^H C x)}] + x^H B x \\ &= \pm \varsigma(x). \end{aligned}$$

Lastly, the inclusion (3.6) is a result of  $\varsigma(x) = [\rho_+(x) - \rho_-(x)] x^H A x$ .  $\square$

### 3.1 Courant-Fischer type min-max principles

Theorem 3.2 below is a restatement of [43, Theorem 32.10, Theorem 32.11 and Remark 32.13]. However, it is essentially due to Duffin [12, 1955] whose proof, although for overdamped  $\mathbf{Q}$ , works for the general hyperbolic case. Closely related ones for more general nonlinear eigenvalue problems (other than quadratic eigenvalue problems) can be found in [49, 50, 66, 67]. They can be considered as a generalization of the Courant-Fischer min-max principles (see [47, p.206], [56, p.201]).

**Theorem 3.2** ([12]). *We have for  $1 \leq i \leq n$*

$$\lambda_i^+ = \max_{\substack{\mathcal{X} \subseteq \mathbb{C}^n \\ \text{codim } \mathcal{X} = i-1}} \min_{\substack{x \in \mathcal{X} \\ x \neq 0}} \rho_+(x), \quad (3.7a)$$

$$\lambda_i^+ = \min_{\substack{\mathcal{X} \subseteq \mathbb{C}^n \\ \text{dim } \mathcal{X} = i}} \max_{\substack{x \in \mathcal{X} \\ x \neq 0}} \rho_+(x), \quad (3.7b)$$

$$\lambda_i^- = \max_{\substack{\mathcal{X} \subseteq \mathbb{C}^n \\ \text{codim } \mathcal{X} = i-1}} \min_{\substack{x \in \mathcal{X} \\ x \neq 0}} \rho_-(x), \quad (3.7c)$$

$$\lambda_i^- = \min_{\substack{\mathcal{X} \subseteq \mathbb{C}^n \\ \text{dim } \mathcal{X} = i}} \max_{\substack{x \in \mathcal{X} \\ x \neq 0}} \rho_-(x). \quad (3.7d)$$

*In particular,*

$$\lambda_1^+ = \min_{x \neq 0} \rho_+(x), \quad \lambda_n^+ = \max_{x \neq 0} \rho_+(x), \quad (3.8a)$$

$$\lambda_1^- = \min_{x \neq 0} \rho_-(x), \quad \lambda_n^- = \max_{x \neq 0} \rho_-(x). \quad (3.8b)$$

### 3.2 Wielandt-Lidskii type min-max principles

Theorems 3.3 and 3.4 which can be considered as generalizations of Amir-Moéz type min-max principles [1] and Theorem 3.5 which can be considered as generalizations of the Wielandt-Lidskii min-max principles ([39, 69] and also [6, p.67], [56, p.199]) and Ky-Fan trace min/max principles [15] are new. For the ease of stating them, let  $\lambda_{\pm} \in \mathbb{R}$  such that

$$\lambda_- \leq \lambda_1^- \leq \lambda_n^- \leq \lambda_0 \leq \lambda_1^+ \leq \lambda_n^+ \leq \lambda_+.$$

Such  $\lambda_{\pm}$  exist, e.g.,  $\lambda_- = \lambda_1^-$  or  $-\infty$  and  $\lambda_+ = \lambda_n^+$  or  $\infty$ . Set intervals

$$\mathcal{I}_+ = \begin{cases} [\lambda_0, \lambda_+], & \text{if } \lambda_+ < \infty, \\ [\lambda_0, \infty), & \text{otherwise,} \end{cases} \quad \mathcal{I}_- = \begin{cases} [\lambda_-, \lambda_0], & \text{if } \lambda_- > -\infty, \\ (-\infty, \lambda_0], & \text{otherwise.} \end{cases} \quad (3.9)$$

The following lemma is also essentially due to Duffin [12] whose proof, although for overdamped  $\mathbf{Q}$ , again works for the general hyperbolic case.

**Lemma 3.1.** *We have*

$$\lambda_i^+ \geq \rho_+(x) \text{ for any } x \in \text{span}\{u_1^+, u_2^+, \dots, u_i^+\}, \quad (3.10a)$$

$$\lambda_i^+ \leq \rho_+(x) \text{ for any } x \in \text{span}\{u_i^+, u_{i+1}^+, \dots, u_n^+\}. \quad (3.10b)$$



To generalize Amir-Moéz/Wielandt-Lidskii min-max principles, we introduce the following notations. For  $X \in \mathbb{C}^{n \times k}$  with  $\text{rank}(X) = k$ ,  $X^H \mathbf{Q}(\lambda) X$  is a  $k \times k$  hyperbolic quadratic matrix polynomial. Hence its quadratic eigenvalues are real. Denote them by  $\lambda_{i,X}^\pm$  arranged as

$$\lambda_{1,X}^- \leq \cdots \leq \lambda_{k,X}^- \leq \lambda_{1,X}^+ \leq \cdots \leq \lambda_{k,X}^+. \quad (3.11)$$

**Theorem 3.3.** *Let  $1 \leq i_1 < \cdots < i_k \leq n$ . For any*

$$\Phi : \underbrace{\mathcal{I}_+ \times \cdots \times \mathcal{I}_+}_k \rightarrow \mathbb{R}$$

that is non-decreasing in each of its arguments, we have<sup>2</sup>

$$\min_{\substack{\mathcal{X}_1 \subset \cdots \subset \mathcal{X}_k \\ \dim \mathcal{X}_j = i_j}} \sup_{\substack{x_j \in \mathcal{X}_j, j=1, \dots, k \\ X = [x_1, \dots, x_k] \\ \text{rank}(X) = k}} \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) = \Phi(\lambda_{i_1}^+, \dots, \lambda_{i_k}^+), \quad (3.12a)$$

$$\max_{\substack{\mathcal{X}_1 \supset \cdots \supset \mathcal{X}_k \\ \text{codim } \mathcal{X}_j = i_j - 1}} \inf_{\substack{x_j \in \mathcal{X}_j, j=1, \dots, k \\ X = [x_1, \dots, x_k] \\ \text{rank}(X) = k}} \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) = \Phi(\lambda_{i_1}^+, \dots, \lambda_{i_k}^+). \quad (3.12b)$$

If also  $\Phi$  is continuous, then “sup” in (3.12a) and “inf” in (3.12b) can be replaced by “max” and “min”, respectively. In particular, setting  $i_j = j$  in (3.12a) and setting  $i_j = j + n - k$  in (3.12b), respectively, give

$$\min_{\text{rank}(X) = k} \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) = \Phi(\lambda_1^+, \dots, \lambda_k^+), \quad (3.13a)$$

$$\max_{\text{rank}(X) = k} \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) = \Phi(\lambda_{n-k+1}^+, \dots, \lambda_n^+). \quad (3.13b)$$

*Proof.* For convenience, we define, for a matrix  $W = [w_1, \dots, w_p]$ ,

$$\mathcal{S}_{j,W} := \text{span}\{w_1, \dots, w_j\}, \mathcal{T}_{j,W} := \text{span}\{w_j, \dots, w_p\} \quad \text{for } j = 1, \dots, p.$$

In particular  $\mathcal{S}_W = \mathcal{S}_{p,W}$ ,  $\mathcal{T}_W = \mathcal{T}_{1,W}$ , and thus  $\mathcal{S}_W = \mathcal{T}_W$ .

First we prove (3.12b). Recall the quadratic eigenvectors  $u_j^+$  introduced in Theorem 2.5. Choose

$$\widehat{\mathcal{X}}_j = \text{span}\{u_{i_j}^+, \dots, u_n^+\} \quad \text{for } j = 1, 2, \dots, k. \quad (3.14)$$

Then  $\widehat{\mathcal{X}}_1 \supset \cdots \supset \widehat{\mathcal{X}}_k$  and  $\text{codim } \widehat{\mathcal{X}}_j = i_j - 1$ . By Lemma 3.1,  $\rho_+(x) \geq \lambda_{i_j}^+$  for any  $x \in \widehat{\mathcal{X}}_j$ . Therefore

$$\min_{\substack{x \in \widehat{\mathcal{X}}_j \\ x \neq 0}} \rho_+(x) = \lambda_{i_j}^+.$$

For any  $X = [x_1, \dots, x_k]$  with  $x_j \in \widehat{\mathcal{X}}_j$  for  $j = 1, \dots, k$  such that  $\text{rank}(X) = k$ , consider  $X^H \mathbf{Q}(\lambda) X$  which is a  $k \times k$  hyperbolic quadratic matrix polynomial. For  $j = 1, \dots, k$ , noticing  $\mathcal{T}_{j,X} \subset \widehat{\mathcal{X}}_j$ , we have by Theorem 3.2

$$\lambda_{j,X}^+ = \max_{\substack{\mathcal{X} \subset \mathcal{T}_{j,X} \\ \dim \mathcal{X} = k-j+1}} \min_{\substack{x \in \mathcal{X} \\ x \neq 0}} \rho_+(x) \geq \min_{\substack{x \in \mathcal{T}_{j,X} \\ x \neq 0}} \rho_+(x) \geq \min_{\substack{x \in \widehat{\mathcal{X}}_j \\ x \neq 0}} \rho_+(x) = \lambda_{i_j}^+.$$

<sup>2</sup>In (3.12a), it is not clear if the “sup” is attainable for any given  $\mathcal{X}_j$  satisfying the given assumptions, except for continuous  $\Phi$ . The same comment applies to the “inf” in (3.12b).

Since  $\Phi(\cdot)$  is non-decreasing in each of its arguments,

$$\Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) \geq \Phi(\lambda_{i_1}^+, \dots, \lambda_{i_k}^+)$$

which gives

$$\min_{\substack{x_j \in \widehat{\mathcal{X}}_j, j=1, \dots, k \\ X=[x_1, \dots, x_k] \\ \text{rank}(X)=k}} \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) \geq \Phi(\lambda_{i_1}^+, \dots, \lambda_{i_k}^+)$$

because  $x_j \in \widehat{\mathcal{X}}_j$  for  $1 \leq i \leq k$  are arbitrary, subject to  $\text{rank}(X) = k$ . Therefore

$$\sup_{\substack{\mathcal{X}_1 \supset \dots \supset \mathcal{X}_k \\ \text{codim } \mathcal{X}_j = i_j - 1}} \inf_{\substack{x_j \in \mathcal{X}_j, j=1, \dots, k \\ X=[x_1, \dots, x_k] \\ \text{rank}(X)=k}} \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) \geq \Phi(\lambda_{i_1}^+, \dots, \lambda_{i_k}^+). \quad (3.15)$$

On the other hand, let  $\mathcal{X}_j$  for  $j = 1, \dots, k$  be any subspaces that satisfy the assumptions:  $\mathcal{X}_1 \supset \dots \supset \mathcal{X}_k$  and  $\text{codim } \mathcal{X}_j = i_j - 1$ . Define  $\mathcal{Y}_j = \text{span}\{u_1^+, \dots, u_{i_j}^+\}$ . Then  $\mathcal{Y}_1 \subset \dots \subset \mathcal{Y}_k$  and  $\dim \mathcal{Y}_j = i_j$ . By [1, Corollary 2.2] (see also [37, Lemma 3.2]), there exists two  $A$ -orthonormal sets  $\{x_1, \dots, x_k\}$  and  $\{y_1, \dots, y_k\}$  with  $x_j \in \mathcal{X}_j$  for  $j = 1, \dots, k$  and  $y_j \in \mathcal{Y}_j$  for  $1 \leq j \leq k$  such that

$$\mathcal{T}_X := \text{span}\{x_1, \dots, x_k\} = \text{span}\{y_1, \dots, y_k\} =: \mathcal{S}_Y.$$

where  $X = [x_1, \dots, x_k]$  and  $Y = [y_1, \dots, y_k]$  satisfy  $X^H A X = Y^H A Y = I_k$ .  $Y^H \mathbf{Q}(\lambda) Y$  is a hyperbolic quadratic matrix polynomial whose pos-type quadratic eigenvalues are  $\lambda_{1,Y}^+ \leq \dots \leq \lambda_{k,Y}^+$ . Since  $\mathcal{S}_Y = \mathcal{T}_X$ ,  $\lambda_{j,Y}^+ = \lambda_{j,X}^+$  for  $j = 1, \dots, k$ . By Lemma 3.1,  $\rho_+(y) \leq \lambda_{i_j}^+$  for any  $y \in \mathcal{Y}_j$ . Therefore

$$\max_{\substack{y \in \mathcal{Y}_j \\ y \neq 0}} \rho_+(y) = \lambda_{i_j}^+.$$

By Theorem 3.2 and noticing  $\mathcal{S}_{j,Y} \subset \mathcal{Y}_j$ , we have, for  $j = 1, \dots, k$ ,

$$\lambda_{j,X}^+ = \lambda_{j,Y}^+ = \min_{\substack{\mathcal{Y} \subset \mathcal{S}_Y \\ \dim \mathcal{Y} = j}} \max_{\substack{y \in \mathcal{Y} \\ y \neq 0}} \rho_+(y) \leq \max_{\substack{y \in \mathcal{S}_{j,Y} \\ y \neq 0}} \rho_+(y) \leq \max_{\substack{y \in \mathcal{Y}_j \\ y \neq 0}} \rho_+(y) = \lambda_{i_j}^+.$$

Since  $\Phi(\cdot)$  is non-decreasing in each of its arguments,

$$\Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) \leq \Phi(\lambda_{i_1}^+, \dots, \lambda_{i_k}^+),$$

which gives

$$\inf_{\substack{x_j \in \mathcal{X}_j, j=1, \dots, k \\ X=[x_1, \dots, x_k] \\ \text{rank}(X)=k}} \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) \leq \Phi(\lambda_{i_1}^+, \dots, \lambda_{i_k}^+).$$

Since  $\mathcal{X}_j$  are arbitrary, we conclude

$$\sup_{\substack{\mathcal{X}_1 \supset \dots \supset \mathcal{X}_k \\ \text{codim } \mathcal{X}_j = i_j - 1}} \inf_{\substack{x_j \in \mathcal{X}_j, j=1, \dots, k \\ X=[x_1, \dots, x_k] \\ \text{rank}(X)=k}} \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) \leq \Phi(\lambda_{i_1}^+, \dots, \lambda_{i_k}^+). \quad (3.16)$$

Combine (3.15) and (3.16) to get

$$\sup_{\substack{\mathcal{X}_1 \supset \dots \supset \mathcal{X}_k \\ \text{codim } \mathcal{X}_j = i_j - 1}} \inf_{\substack{x_j \in \mathcal{X}_j, j=1, \dots, k \\ X = [x_1, \dots, x_k] \\ \text{rank}(X) = k}} \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) = \Phi(\lambda_{i_1}^+, \dots, \lambda_{i_k}^+). \quad (3.12b')$$

But the “sup” here is achievable by the selection in (3.14). Thus we have (3.12b).

Now we claim the “inf” can be replaced by “min” for a continuous  $\Phi$ . Let  $\mathcal{X}_j$  for  $j = 1, \dots, k$  be given and satisfy the assumptions:  $\mathcal{X}_1 \supset \dots \supset \mathcal{X}_k$  and  $\text{codim } \mathcal{X}_j = i_j - 1$ . There exist a sequence  $X^{(i)} \in \mathbb{C}^{n \times k}$  with  $\text{rank}(X^{(i)}) = k$  and its  $j$ th column in  $\mathcal{X}_j$  such that

$$\lim_{i \rightarrow \infty} \Phi(\lambda_{1,X^{(i)}}^+, \dots, \lambda_{k,X^{(i)}}^+) = \inf_{\substack{x_j \in \mathcal{X}_j, j=1, \dots, k \\ X = [x_1, \dots, x_k] \\ \text{rank}(X) = k}} \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+). \quad (3.17)$$

Without loss of generality, we may assume  $X^{(i)}$  has  $A$ -orthonormal columns, i.e.,

$$(X^{(i)})^H A X^{(i)} = I_k;$$

otherwise we can perform the Gram-Schmidt  $A$ -orthogonalization on the columns of  $X^{(i)}$  from the last column backwards, and the new  $X^{(i)}$  has the same property as the old  $X^{(i)}$ :  $\text{rank}(X^{(i)}) = k$  and its  $j$ th column in  $\mathcal{X}_j$ , and also  $\lambda_{j,X^{(i)}}^\pm$  remain the same. Since  $\{X^{(i)}\}$  is a bounded set in  $\mathbb{C}^{n \times k}$ , it has a convergent subsequence. Through renaming, we may assume that  $\{X^{(i)}\}$  itself is convergent, and let  $Y \in \mathbb{C}^{n \times k}$  be the limit. It is not hard to see that  $Y^H A Y = I_k$  which implies  $\text{rank}(Y) = k$  and that the  $j$ th column of  $Y$  is in  $\mathcal{X}_j$ . Since  $(X^{(i)})^H \mathbf{Q}(\lambda) X^{(i)}$  goes to  $Y^H \mathbf{Q}(\lambda) Y$ , by the continuity of quadratic eigenvalues with respect to the coefficient matrices we conclude

$$\lim_{i \rightarrow \infty} \lambda_{j,X^{(i)}}^\pm = \lambda_{j,Y}^\pm \quad \text{for } 1 \leq j \leq k.$$

Therefore the left-hand side of (3.17) is equal to  $\Phi(\lambda_{1,Y}^+, \dots, \lambda_{k,Y}^+)$ , and thus the “inf” in (3.17) is attainable.

For (3.12a), a proof similar to what we did above for (3.12b) works: choosing  $\widehat{\mathcal{X}}_j = \text{span}\{u_1^+, \dots, u_{i_j}^+\}$  will lead to that the left-hand side is no bigger than its right-hand side, and choosing  $\mathcal{Y}_j = \text{span}\{u_{i_j}^+, \dots, u_n^+\}$  will give the opposite.  $\square$

Similarly to Theorem 3.3, we have

**Theorem 3.4.** *Let  $1 \leq i_1 < \dots < i_k \leq n$ . For any*

$$\Psi : \underbrace{\mathcal{I}_- \times \dots \times \mathcal{I}_-}_k \rightarrow \mathbb{R}$$

*that is non-decreasing in each of its arguments, we have<sup>3</sup>*

$$\min_{\substack{\mathcal{X}_1 \subset \dots \subset \mathcal{X}_k \\ \dim \mathcal{X}_j = i_j}} \sup_{\substack{x_j \in \mathcal{X}_j, j=1, \dots, k \\ X = [x_1, \dots, x_k] \\ \text{rank}(X) = k}} \Psi(\lambda_{1,X}^-, \dots, \lambda_{k,X}^-) = \Psi(\lambda_{i_1}^-, \dots, \lambda_{i_k}^-), \quad (3.18a)$$

<sup>3</sup>In (3.18a), it is not clear if the “sup” is attainable for any given  $\mathcal{X}_j$  satisfying the given assumptions. The same comment applies to the “inf” in (3.18b).

$$\max_{\substack{\mathcal{X}_1 \supset \dots \supset \mathcal{X}_k \\ \text{codim } \mathcal{X}_j = i_j - 1}} \inf_{\substack{x_j \in \mathcal{X}_j, j=1, \dots, k \\ X = [x_1, \dots, x_k] \\ \text{rank}(X) = k}} \Psi(\lambda_{1,X}^-, \dots, \lambda_{k,X}^-) = \Psi(\lambda_{i_1}^-, \dots, \lambda_{i_k}^-). \quad (3.18b)$$

If also  $\Psi$  is continuous, then “sup” in (3.18a) and “inf” in (3.18b) can be replaced by “max” and “min”, respectively. In particular, setting  $i_j = j$  in (3.18a) and setting  $i_j = j + n - k$  in (3.18b), respectively, give

$$\min_{\text{rank}(X)=k} \Psi(\lambda_{1,X}^-, \dots, \lambda_{k,X}^-) = \Psi(\lambda_1^-, \dots, \lambda_k^-), \quad (3.19a)$$

$$\max_{\text{rank}(X)=k} \Psi(\lambda_{1,X}^-, \dots, \lambda_{k,X}^-) = \Psi(\lambda_{n-k+1}^-, \dots, \lambda_n^-). \quad (3.19b)$$

*Proof.* Consider the hyperbolic quadratic matrix polynomial  $\widehat{Q}(\lambda) = \lambda^2 A + \lambda(-B) + C$  whose quadratic eigenvalues are

$$\widehat{\lambda}_1^- \leq \dots \leq \widehat{\lambda}_n^- < \widehat{\lambda}_1^+ \leq \dots \leq \widehat{\lambda}_n^+,$$

where  $\widehat{\lambda}_i^- = -\lambda_{n-i+1}^+$  and  $\widehat{\lambda}_j^+ = -\lambda_{n-j+1}^-$ . Apply (3.12b) to  $\widehat{Q}(\lambda)$  with

$$\widehat{\Phi}(\xi_1, \dots, \xi_k) := -\Psi(-\xi_k, \dots, -\xi_1)$$

to get (3.18a), and apply (3.12a) to  $\widehat{Q}(\lambda)$  with the same  $\widehat{\Phi}$  to get (3.18b).  $\square$

Specializing Theorems 3.3 and 3.4 to the case where  $\widehat{\Phi}$  and  $\Psi$  are the sum of its arguments gives us Wielandt-Lidskii type min-max principles as summarized in the following theorem and Ky-Fan type trace min/max principles.

**Theorem 3.5.** *Let  $1 \leq i_1 < \dots < i_k \leq n$  and  $\text{typ} \in \{+, -\}$ . Then*

$$\min_{\substack{\mathcal{X}_1 \subset \dots \subset \mathcal{X}_k \\ \dim \mathcal{X}_j = i_j}} \max_{\substack{x_j \in \mathcal{X}_j \\ X = [x_1, \dots, x_k] \\ \text{rank}(X) = k}} \sum_{j=1}^k \lambda_{j,X}^{\text{typ}} = \sum_{j=1}^k \lambda_{i_j}^{\text{typ}}, \quad (3.20a)$$

$$\max_{\substack{\mathcal{X}_1 \supset \dots \supset \mathcal{X}_k \\ \text{codim } \mathcal{X}_j = i_j - 1}} \min_{\substack{x_j \in \mathcal{X}_j \\ X = [x_1, \dots, x_k] \\ \text{rank}(X) = k}} \sum_{j=1}^k \lambda_{j,X}^{\text{typ}} = \sum_{j=1}^k \lambda_{i_j}^{\text{typ}}. \quad (3.20b)$$

In particular, setting  $i_j = j$  in (3.20a) and setting  $i_j = j + n - k$  in (3.20b) give

$$\min_{\text{rank}(X)=k} \sum_{j=1}^k \lambda_{j,X}^{\text{typ}} = \sum_{j=1}^k \lambda_j^{\text{typ}}, \quad \max_{\text{rank}(X)=k} \sum_{j=1}^k \lambda_{j,X}^{\text{typ}} = \sum_{j=1}^k \lambda_{n-k+j}^{\text{typ}}. \quad (3.21)$$

### 3.3 Cauchy type interlacing inequalities

The Cauchy type interlacing inequalities in (3.22) were recently obtained by Veselić [64]. Here we present a simple proof, using our generalizations of Amir-Moéz type min-max principles in Theorems 3.3 and 3.4.

**Theorem 3.6** (Cauchy-type interlacing inequalities [64]). *Suppose  $X \in \mathbb{C}^{n \times k}$  with  $\text{rank}(X) = k$ . Denote the quadratic eigenvalues of  $X^H Q(\lambda) X$  by*

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

Then

$$\lambda_i^+ \leq \mu_i^+ \leq \lambda_{i+n-k}^+, \quad i = 1, \dots, k, \quad (3.22a)$$

$$\lambda_j^- \leq \mu_j^- \leq \lambda_{j+n-k}^-, \quad j = 1, \dots, k. \quad (3.22b)$$

*Proof.* Let

$$\Phi(\alpha_1, \dots, \alpha_k) = \text{the } i\text{th largest } \alpha_j.$$

Then this  $\Phi$  satisfies the condition of Theorem 3.3. Making use of (3.13a) and (3.13b) gives  $\mu_i^+ \geq \lambda_i^+$  and  $\mu_i^+ \leq \lambda_{i+n-k}^+$ , respectively. That is (3.22a). Similarly, we get (3.22b) by Theorem 3.4.  $\square$

**Remark 3.1.** The Cauchy type interlacing inequalities in Theorem 3.6 are sharper than those possibly derivable by linearization. Actually, through linearization and by item 1 of [38, Theorem 1.1] (which is, in fact, [30, Theorem 2.1]), we can only obtain

$$\begin{aligned} \lambda_i^+ &\leq \mu_i^+ \leq \lambda_{i+2n-2k}^+, & i = 1, \dots, k, \\ \lambda_{j-(n-k)}^- &\leq \mu_j^- \leq \lambda_{j+n-k}^-, & j = 1, \dots, k, \end{aligned}$$

where we set  $\lambda_i^+ = +\infty$  for  $i > n$  and  $\lambda_j^- = -\infty$  for  $j < 1$ .

## 4 Perturbation analysis

### 4.1 Setting the stage

Throughout this section, we suppose that Hermitian matrices  $A$ ,  $B$ , and  $C$  are perturbed to Hermitian matrices  $\tilde{A}$ ,  $\tilde{B}$ , and  $\tilde{C}$  and set

$$\Delta A = \tilde{A} - A, \quad \Delta B = \tilde{B} - B, \quad \Delta C = \tilde{C} - C. \quad (4.1)$$

This notational convention of placing a “ $\sim$ ” over a symbol for the corresponding perturbed quantity and a “ $\Delta$ ” before a symbol for the change in the quantity will be generalized to all quantities that depend on  $A$ ,  $B$ , and  $C$ . For example,  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$  is perturbed to  $\tilde{\mathbf{Q}}(\lambda) = \lambda^2 \tilde{A} + \lambda \tilde{B} + \tilde{C}$ , as a result, and

$$\Delta \rho_{\pm}(x) = \frac{-(x^{\text{H}} \tilde{B} x) \pm \left[ (x^{\text{H}} \tilde{B} x)^2 - 4(x^{\text{H}} \tilde{A} x)(x^{\text{H}} \tilde{C} x) \right]^{1/2}}{2(x^{\text{H}} \tilde{A} x)} - \frac{-(x^{\text{H}} B x) \pm \left[ (x^{\text{H}} B x)^2 - 4(x^{\text{H}} A x)(x^{\text{H}} C x) \right]^{1/2}}{2(x^{\text{H}} A x)}.$$

Besides  $A \succ 0$ , the other key condition for  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$  to be hyperbolic is

$$[\varsigma(x)]^2 = (x^{\text{H}} B x)^2 - 4(x^{\text{H}} A x)(x^{\text{H}} C x) > 0, \quad \text{for all } 0 \neq x \in \mathbb{C}^n. \quad (2.2)$$

We first establish a condition under which (2.2) is weakly<sup>4</sup> satisfied for all convex combination  $(1-t)\mathbf{Q}(\lambda) + t\tilde{\mathbf{Q}}(\lambda)$ . To this end, we define

$$\phi(x) := (x^{\text{H}} \Delta B x)^2 - 4(x^{\text{H}} \Delta A x)(x^{\text{H}} \Delta C x), \quad (4.2)$$

$$\psi(x) := (x^{\text{H}} B x)(x^{\text{H}} \Delta B x) - 2(x^{\text{H}} A x)(x^{\text{H}} \Delta C x) - 2(x^{\text{H}} C x)(x^{\text{H}} \Delta A x), \quad (4.3)$$

and define  $\tilde{\phi}(x)$  and  $\tilde{\psi}(x)$  in the same way, except by swapping the positions of  $A$ ,  $B$ ,  $C$  with those of  $\tilde{A}$ ,  $\tilde{B}$ , and  $\tilde{C}$ . It can be verified that

$$\tilde{\phi}(x) = \phi(x), \quad \tilde{\psi}(x) = -\psi(x) - \phi(x).$$

Also define

$$\begin{aligned} g(t) &:= (x^{\text{H}} [B + t \Delta B] x)^2 - 4(x^{\text{H}} [A + t \Delta A] x)(x^{\text{H}} [C + t \Delta C] x) \\ &= \varsigma(x)^2 + 2\psi(x)t + \phi(x)t^2. \end{aligned}$$

So  $g(0) = \varsigma(x)$  and  $g(1) = \zeta(x)$ . Correspondingly,

$$\begin{aligned} \tilde{g}(t) &:= (x^{\text{H}} [\tilde{B} - t \Delta B] x)^2 - 4(x^{\text{H}} [\tilde{A} - t \Delta A] x)(x^{\text{H}} [\tilde{C} - t \Delta C] x) \\ &= \tilde{\zeta}(x)^2 + 2\tilde{\psi}(x)t + \phi(x)t^2. \end{aligned}$$

Note that  $g(t) = \tilde{g}(1-t)$ .

By definition, if  $A \succ 0$ , then  $\mathbf{Q}(\lambda)$  is hyperbolic if and only if  $g(0) > 0$  for any nonzero  $x \in \mathbb{C}^n$ , and if  $\tilde{A} \succ 0$ , then  $\tilde{\mathbf{Q}}(\lambda)$  is hyperbolic if and only if  $\tilde{g}(1) > 0$  for any nonzero  $x \in \mathbb{C}^n$ .

<sup>4</sup>By *weakly*, we mean the strict positivity in (2.2) is given in to nonnegativity.

**Lemma 4.1.** *Suppose  $\min\{g(0), g(1)\} \geq 0$ . Then  $g(t) \geq 0$  for all  $0 \leq t \leq 1$  and nonzero  $x \in \mathbb{C}^n$  if and only if*

$$\min\{\phi(x), -\psi(x), -\tilde{\psi}(x), \psi(x)^2 - \phi(x)\varsigma(x)^2\} \leq 0 \text{ for all } x \neq 0. \quad (4.4)$$

*Proof.* The condition (4.4) is equivalent to that for any  $x$ , at least one of

$$\phi(x) \leq 0, \quad \psi(x) \geq 0, \quad \tilde{\psi}(x) = -\psi(x) - \phi(x) \geq 0, \quad \psi(x)^2 - \phi(x)\varsigma(x)^2 \leq 0$$

holds. Note that  $g(0) \geq 0$  and  $g(1) \geq 0$  by assumption. We first prove that (4.4) implies  $g(t) \geq 0$  for all  $0 \leq t \leq 1$  and for any nonzero  $x \in \mathbb{C}^n$ . To this end, we let  $0 \leq t \leq 1$  and  $0 \neq x \in \mathbb{C}^n$ .

1. If  $\phi(x) \leq 0$ , then  $g(t)$  is concave and thus  $g(t) \geq (1-t)g(0) + tg(1) \geq 0$ ;
2. If  $\psi(x) \geq 0$ , then

$$\begin{aligned} g(t) &= \varsigma(x)^2 + 2\psi(x)t + \phi(x)t^2 \\ &\geq \varsigma(x)^2 + 2\psi(x)t^2 + \phi(x)t^2 \\ &= (1-t^2)g(0) + t^2g(1) \\ &\geq 0; \end{aligned}$$

3. If  $\tilde{\psi}(x) \geq 0$ , then similarly  $\tilde{g}(t) \geq (1-t^2)\tilde{g}(0) + t^2\tilde{g}(1) \geq 0$ ;
4. Consider the case  $\psi(x)^2 - \phi(x)\varsigma(x)^2 \leq 0$ . Suppose<sup>5</sup>  $\phi(x) > 0$ . Then  $g(t)$  is a nontrivial quadratic function and has at most one zero in  $\mathbb{R}$ . Going through the cases either there is no zero or the zero is in  $(0, 1)$  or the zero is outside of  $(0, 1)$ , we can see  $g(t) \geq 0$  for all  $0 \leq t \leq 1$ .

Next for the necessity of (4.4), suppose there were an  $x \neq 0$  satisfying  $\phi(x) > 0$ ,  $\psi(x) < 0$ ,  $-\tilde{\psi}(x) = \psi(x) + \phi(x) > 0$ , and  $\psi(x)^2 - \phi(x)\varsigma(x)^2 > 0$ . Then

$$\min_t g(t) = -\frac{\psi(x)^2 - \phi(x)\varsigma(x)^2}{\phi(x)} < 0$$

and  $\min_t g(t)$  is attained at  $t_{\min} = -\frac{\psi(x)}{\phi(x)} \in (0, 1)$ , contradicting the assumption that  $g(t) \geq 0$  for  $0 \leq t \leq 1$ .  $\square$

Given a shift  $\lambda_0 \in \mathbb{R}$ , define

$$\mathbf{Q}_{\lambda_0}(\lambda) := \mathbf{Q}(\lambda + \lambda_0) = \lambda^2 A + \lambda(2\lambda_0 A + B) + \mathbf{Q}(\lambda_0) \quad (4.5)$$

$$= \lambda^2 A + \lambda B_{\lambda_0} + C_{\lambda_0}, \quad (4.6)$$

where

$$B_{\lambda_0} = 2\lambda_0 A + B, \quad C_{\lambda_0} = \mathbf{Q}(\lambda_0). \quad (4.7)$$

It can be verified that  $(\mu, x)$  is a quadratic eigenpair of  $\mathbf{Q}_{\lambda_0}(\lambda)$  if and only if  $(\mu + \lambda_0, x)$  is a quadratic eigenpair of  $\mathbf{Q}(\lambda)$ .

<sup>5</sup>The case  $\phi(x) \leq 0$  has already been dealt with.

**Lemma 4.2.** *Suppose that  $\mathbf{Q}(\lambda)$  is hyperbolic, and adopt the notations introduced in Theorem 2.5.*

1. *If  $\lambda_0 \in (\lambda_n^-, \lambda_1^+)$ , then  $\text{diag}(-C_{\lambda_0}, A) = \text{diag}(-\mathbf{Q}(\lambda_0), A) \succ 0$ .*
2. *If  $\lambda_0 \in [\lambda_n^+, +\infty)$ , then  $\mathbf{Q}_{\lambda_0}(\lambda)$  is overdamped, i.e.  $B_{\lambda_0} \succ 0$  and  $C_{\lambda_0} \succeq 0$ . Moreover,*

$$-(\lambda_n^- + \lambda_n^+ - 2\lambda_0)A \preceq B_{\lambda_0} \preceq -(\lambda_1^- + \lambda_1^+ - 2\lambda_0)A, \quad (4.8)$$

$$(\lambda_n^- - \lambda_0)(\lambda_n^+ - \lambda_0)A \preceq C_{\lambda_0} \preceq (\lambda_1^- - \lambda_0)(\lambda_1^+ - \lambda_0)A. \quad (4.9)$$

3. *If  $\|A^{-1/2}\Delta AA^{-1/2}\|_2 < 1$ , then  $\tilde{A} \succ 0$ .*

*Proof.* Item 1 is a consequence of Theorem 2.1 and (4.7). For (4.8) of item 2, we have for any  $x \neq 0$

$$\begin{aligned} x^H B_{\lambda_0} x &= 2\lambda_0 x^H A x + x^H B x \\ &= x^H A x \left( 2\lambda_0 + \frac{x^H B x}{x^H A x} \right) \\ &= x^H A x (2\lambda_0 - [\rho_+(x) + \rho_-(x)]) \end{aligned}$$

which, together with (3.4), yields (4.8). For (4.9), we have for any  $x \neq 0$

$$x^H C_{\lambda_0} x = x^H \mathbf{Q}(\lambda_0) x = x^H A x [\lambda_0 - \rho_+(x)] [\lambda_0 - \rho_-(x)]$$

which, together with (3.4), yields (4.9). For item 3, we notice the smallest eigenvalue of  $A^{-1/2}\tilde{A}A^{-1/2}$  satisfies

$$\lambda_{\min}(A^{-1/2}\tilde{A}A^{-1/2}) = 1 + \lambda_{\min}(A^{-1/2}\Delta AA^{-1/2}) \geq 1 - \|A^{-1/2}\Delta AA^{-1/2}\|_2 > 0$$

if  $\|A^{-1/2}\Delta AA^{-1/2}\|_2 < 1$ . □

## 4.2 Asymptotical analysis

It is a common technique to perform an asymptotical analysis in numerical analysis for at least three reasons:

1. it is mathematically sound, provided it is known that the interested quantities are continuous with respect to what is being perturbed;
2. it is relatively easy because it is a first order analysis, and
3. it is powerful in revealing the intrinsic sensitivity of the interested quantities.

Let  $(\mu, x)$  is a simple quadratic eigenpair of HQEP (1.1) for  $\mathbf{Q}(\lambda)$ . Since HQEP (1.1) is equivalent to the eigenvalue problem for the regular matrix pencil  $\mathcal{L}_{\mathbf{Q}}(\lambda)$  in (2.5) and since the eigenvalues of a regular matrix pencil and the eigenvectors associated with simple eigenvalues are continuous with respect to the entries of the involved matrices [56],  $\tilde{\mathbf{Q}}(\lambda)$  has a simple quadratic eigenpair  $(\tilde{\mu}, \tilde{x}) = (\mu + \Delta\mu, x + \Delta x)$  such that  $\Delta\mu \rightarrow 0$  and  $\Delta x \rightarrow 0$  as  $\Delta A, \Delta B, \Delta C \rightarrow 0$ . Now suppose that  $\|\Delta A\|, \|\Delta B\|$ , and  $\|\Delta C\|$  are sufficiently tiny,



and so are  $\Delta\mu$  and  $\|\Delta x\|$ . Ignoring terms of order 2 or higher and noticing  $\mathbf{Q}(\mu)x = 0$ , we have from  $\tilde{\mathbf{Q}}(\mu + \Delta\mu)(x + \Delta x) = 0$

$$\Delta\mu[2\mu A + B]x + [\mu^2\Delta A + \mu\Delta B + \Delta C]x + [\mu^2A + \mu B + C]\Delta x \approx 0, \quad (4.10)$$

where the “ $\approx$ ” means the equation is true after ignoring terms of order 2 or higher. Premultiply (4.10) by  $x^H$  and use  $x^H\mathbf{Q}(\mu) = 0$  to get

$$\Delta\mu \approx -\frac{x^H[\mu^2\Delta A + \mu\Delta B + \Delta C]x}{x^H[2\mu A + B]x} \quad (4.11)$$

$$= -\frac{x^H[\mu^2\Delta A + \mu\Delta B + \Delta C]x}{\varsigma(x)} \quad (4.12)$$

$$= -\frac{\mu^2}{\pm\varsigma(x)} \cdot x^H\Delta Ax - \frac{\mu}{\pm\varsigma(x)} \cdot x^H\Delta Bx - \frac{1}{\pm\varsigma(x)} \cdot x^H\Delta Cx. \quad (4.13)$$

where the equality in (4.12) is due to (3.5). There is a clear interpretation of (4.13): the change  $\Delta\mu$  is proportional to  $\Delta A$ ,  $\Delta B$ ,  $\Delta C$  with multiplying factors  $|\mu^2/\varsigma(x)|$ ,  $|\mu/\varsigma(x)|$ , and  $1/|\varsigma(x)|$ , respectively. Our following strict bounds reflect this interpretation.

The expression (4.11) is not new and its derivation follows a rather standard technique (see, e.g., [62]). What is new here is the use of (3.5) to relate its denominator  $x^H[2\mu A + B]x$  to  $\varsigma(x)$ , a quantity that determines the hyperbolicity of  $\mathbf{Q}$ .

### 4.3 Perturbation bounds in the spectral norm

Throughout the rest of this section, we assume  $\mathbf{Q}(\lambda)$  and  $\tilde{\mathbf{Q}}(\lambda)$  are hyperbolic and

$$\|A^{-1/2}\Delta AA^{-1/2}\|_2 < 1 \quad (4.14)$$

which guarantees  $\tilde{A} \succ 0$ . We will adopt the notations introduced in Theorem 2.5. Our goal is to bound the norms of

$$\Delta\Lambda_+ = \text{diag}(\tilde{\lambda}_1^+ - \lambda_1^+, \dots, \tilde{\lambda}_n^+ - \lambda_n^+), \quad \Delta\Lambda_- = \text{diag}(\tilde{\lambda}_1^- - \lambda_1^-, \dots, \tilde{\lambda}_n^- - \lambda_n^-).$$

Bounds on norms of the change to  $\Lambda = \text{diag}(\Lambda_-, \Lambda_+)$  are easily derivable through

$$\begin{aligned} \|\Delta\Lambda\|_2 &= \max_{\pm} \|\Delta\Lambda_{\pm}\|_2, \quad \|\Delta\Lambda\|_F = \sqrt{\|\Delta\Lambda_+\|_F^2 + \|\Delta\Lambda_-\|_F^2}, \\ \|\Delta\Lambda\|_{\text{ui}} &\leq 2 \max_{\pm} \|\Delta\Lambda_{\pm}\|_{\text{ui}}. \end{aligned}$$

In this subsection, we will focus on the spectral norm, and leave the case for the Frobenius norms and more generally unitarily invariant norms to the next subsection. Our main results of this subsection are summarized in Theorem 4.1.

**Theorem 4.1.** *Let  $\text{typ} \in \{+, -\}$ , and*

$$\epsilon_a = \|A^{-1/2}\Delta AA^{-1/2}\|_2, \quad \epsilon_b = \frac{\|\Delta B\|_2}{\|B\|_2}, \quad \epsilon_c = \frac{\|\Delta C\|_2}{\|C\|_2}, \quad (4.15)$$

$$\lambda_{\max}^{\text{typ}} = \max\{|\lambda_1^{\text{typ}}|, |\lambda_n^{\text{typ}}|\}, \quad \tilde{\lambda}_{\max}^{\text{typ}} = \max\{|\tilde{\lambda}_1^{\text{typ}}|, |\tilde{\lambda}_n^{\text{typ}}|\}, \quad (4.16)$$

$$\chi_{\varsigma} = \min_{x \neq 0} \{\varsigma_0(x), \tilde{\varsigma}_0(x)\}, \quad \chi_{\lambda^{\text{typ}}} = \max\{\lambda_{\max}^{\text{typ}}, \tilde{\lambda}_{\max}^{\text{typ}}\}. \quad (4.17)$$

1. If  $\Delta A = \Delta B = 0$  and

$$\epsilon_c < \frac{\chi_\zeta^2}{4\|A\|_2\|C\|_2}, \quad (4.18)$$

then

$$\|\Delta A_{\text{typ}}\|_2 \leq \frac{1}{\chi_\zeta} \|\Delta C\|_2. \quad (4.19)$$

2. If  $\Delta B = \Delta C = 0$  and

$$\epsilon_a < \min \left\{ 1, \frac{\chi_\zeta^2}{4\|A\|_2\|C\|_2} \right\}, \quad (4.20)$$

then

$$\|\Delta A_{\text{typ}}\|_2 \leq \frac{\chi_{\lambda^{\text{typ}}}^2}{(1 - \epsilon_a)\chi_\zeta} \|\Delta A\|_2. \quad (4.21)$$

3. If  $\Delta A = \Delta C = 0$  and

$$\epsilon_b < \frac{\chi_\zeta^2}{\|B\|_2(\|B\|_2 + 2\sqrt{\|A\|_2\|C\|_2})}, \quad (4.22)$$

then

$$\|\Delta A_{\text{typ}}\|_2 \leq \frac{\chi_{\lambda^{\text{typ}}}}{\chi_\zeta} \|\Delta B\|_2 + \frac{\|C\|_2}{\chi_\zeta^3} \|\Delta B\|_2^2. \quad (4.23)$$

4. If  $\Delta A = \Delta C = 0$  and

$$\|\Delta B\|_2 < \frac{\chi_\zeta^2}{\|2\lambda_0 A + B\|_2 + 2\sqrt{\|A\|_2\|Q(\lambda_0)\|_2}}, \quad (4.24)$$

where  $\lambda_0 \in (-\infty, \min\{\lambda_1^-, \tilde{\lambda}_1^-\}] \cup [\max\{\lambda_n^+, \tilde{\lambda}_n^+\}, +\infty)$ , then

$$\|\Delta A_{\text{typ}}\|_2 \leq \frac{\chi_{\lambda^{\text{typ}}} + |\lambda_0|}{\chi_\zeta} \|\Delta B\|_2. \quad (4.25)$$

5. In general, without assuming two of  $\Delta A$ ,  $\Delta B$ , and  $\Delta C$  are zeros, if

$$\epsilon_a < \gamma \min \left\{ 1, \frac{\chi_\zeta^2}{4\|A\|_2\|C\|_2} \right\}, \quad (4.26a)$$

$$\epsilon_b < \gamma \frac{\chi_\zeta^2}{\|B\|_2(\|B\|_2 + 2\sqrt{\|A\|_2\|C\|_2})}, \quad (4.26b)$$

$$\epsilon_c < \gamma \frac{\chi_\zeta^2}{4\|A\|_2\|C\|_2}, \quad (4.26c)$$

where

$$\gamma = \frac{\chi_\zeta^2}{\|B\|_2^2 + \chi_\zeta^2 + \sqrt{(\|B\|_2^2 + \chi_\zeta^2)(\|B\|_2^2 + 2\chi_\zeta^2)}} < \sqrt{2} - 1, \quad (4.27)$$

then

$$\begin{aligned} \|\Delta A_{\text{typ}}\|_2 &\leq \frac{4}{(1 - \epsilon_a)\chi_\zeta^3} \|C\|_2 [\|A\|_2\|C\|_2(\epsilon_a + \epsilon_c)^2 + \|B\|_2^2(\epsilon_b + \epsilon_a)(\epsilon_b + \epsilon_c)] \\ &\quad + \frac{1}{(1 - \epsilon_a)\chi_\zeta} [(\chi_{\lambda^{\text{typ}}})^2 \|\Delta A\|_2 + \chi_{\lambda^{\text{typ}}} \|\Delta B\|_2 + \|\Delta C\|_2]. \end{aligned} \quad (4.28)$$

All bounds by this theorem are strict. They are consistent with the asymptotic expression (4.13) rather well after dropping terms of order 2 or higher in  $\epsilon_a$ ,  $\epsilon_b$ , and  $\epsilon_c$ . For example, (4.28) yields

$$\|\Delta A_{\text{typ}}\|_2 \lesssim \frac{1}{\chi_\varsigma} [(\chi_{\lambda_{\text{typ}}})^2 \|\Delta A\|_2 + \chi_{\lambda_{\text{typ}}} \|\Delta B\|_2 + \|\Delta C\|_2]. \quad (4.29)$$

The rest of this subsection is devoted for the proof of Theorem 4.1. Later in the next subsection, we will extend (4.19) to a general unitarily invariant norm.

Each of many expressions below is in its compact form for two. For example, (4.30) includes two displayed equations: one for  $\Delta\rho_+$  and one for  $\Delta\rho_-$  with all “ $\pm$ ” selected as either “+” or “-”, accordingly.

**Lemma 4.3.** *If (4.4) and (4.14) hold, then there exists  $0 \leq \xi \leq 1$  such that*

$$\Delta\rho_\pm(x) = \delta^\pm(x, \xi) := \pm \left[ \delta_3(x, \xi) - \frac{x^H A x}{x^H \tilde{A} x} \delta_2^\pm(x) \right] \quad (4.30)$$

for any  $x \neq 0$ , where

$$\delta_2^\pm(x) = \frac{\rho_\pm(x)^2 (x^H \Delta A x) + \rho_\pm(x) (x^H \Delta B x) + x^H \Delta C x}{\varsigma(x)}, \quad (4.31a)$$

$$\delta_3(x, \xi) = \frac{\varsigma(x)^2 \phi(x) - \psi(x)^2}{4(x^H \tilde{A} x) [\varsigma(x)^2 + 2\psi(x)\xi + \phi(x)\xi^2]^{3/2}}, \quad (4.31b)$$

$\phi(x)$  and  $\psi(x)$  are defined in (4.2) and (4.3). In addition,

$$\frac{1}{1 + \|A^{-1/2} \Delta A A^{-1/2}\|_2} \leq \frac{x^H A x}{x^H \tilde{A} x} \leq \frac{1}{1 - \|A^{-1/2} \Delta A A^{-1/2}\|_2}, \quad (4.32)$$

$$|\delta_2^\pm(x)| \leq \frac{\max\{|\lambda_1^\pm|^2, |\lambda_n^\pm|^2\} \|\Delta A\|_2 + \max\{|\lambda_1^\pm|, |\lambda_n^\pm|\} \|\Delta B\|_2 + \|\Delta C\|_2}{\min_{x \neq 0} \varsigma_0(x)}. \quad (4.33)$$

*Proof.* According to how the difference operator  $\Delta$  is defined at the beginning of subsection 4.1, we have

$$\pm \Delta\rho_\pm(x) = \frac{\Delta\varsigma(x) \mp x^H \Delta B x}{2(x^H A x)} + \frac{\tilde{\varsigma}(x) \mp x^H \tilde{B} x}{2} \Delta \left( \frac{1}{x^H A x} \right) =: \epsilon_1 + \epsilon_2. \quad (4.34)$$

The rest of this proof is to calculate  $\epsilon_1$  and  $\epsilon_2$ . By Lemma 4.1,

$$f(t; x) := [\varsigma(x)^2 + 2\psi(x)t + \phi(x)t^2]^{1/2} \quad (4.35)$$

is well-defined and differentiable for  $0 \leq t \leq 1$ . By the Taylor expansion, there exists  $0 \leq \xi \leq 1$  such that

$$\begin{aligned} \tilde{\varsigma}(x) &= f(1; x) = f(0; x) + f'(0; x) + \frac{1}{2} f''(\xi; x) \\ &= \varsigma(x) + \frac{\psi(x)}{\varsigma(x)} + \frac{\varsigma(x)^2 \phi(x) - \psi(x)^2}{2[f(\xi; x)]^3}. \end{aligned} \quad (4.36)$$

This  $\xi$  depends on  $x$ . Now we are ready to calculate  $\epsilon_1$  and  $\epsilon_2$ . We have

$$\begin{aligned}
\epsilon_1 &= \mp \frac{x^H \Delta B x}{2(x^H A x)} + \frac{1}{2(x^H A x)} \left( \frac{\psi(x)}{\varsigma(x)} + \frac{\varsigma(x)^2 \phi(x) - \psi(x)^2}{2[f(\xi; x)]^3} \right) \\
&= \mp \frac{x^H \Delta B x}{2(x^H A x)} + \frac{(x^H B x)(x^H \Delta B x)}{2(x^H A x)\varsigma(x)} - \frac{x^H \Delta C x}{\varsigma(x)} - \frac{x^H C x}{\varsigma(x)} \frac{x^H \Delta A x}{x^H A x} + \frac{\varsigma(x)^2 \phi(x) - \psi(x)^2}{4(x^H A x)[f(\xi; x)]^3} \\
&= -\frac{\pm \varsigma(x) - (x^H B x)}{2(x^H A x)} \frac{x^H \Delta B x}{\varsigma(x)} - \frac{x^H \Delta C x}{\varsigma(x)} - \frac{x^H C x}{\varsigma(x)} \frac{x^H \Delta A x}{x^H A x} + \frac{\varsigma(x)^2 \phi(x) - \psi(x)^2}{4(x^H A x)[f(\xi; x)]^3} \\
&= -\frac{\rho_{\pm}(x)(x^H \Delta B x)}{\varsigma(x)} - \frac{x^H \Delta C x}{\varsigma(x)} - \frac{x^H C x}{\varsigma(x)} \frac{x^H \Delta A x}{x^H A x} + \frac{x^H \tilde{A} x}{x^H A x} \frac{\varsigma(x)^2 \phi(x) - \psi(x)^2}{4(x^H \tilde{A} x)[f(\xi; x)]^3} \\
&= -\delta_2^{\pm}(x) + \frac{\rho_{\pm}(x)^2 (x^H \Delta A x)}{\varsigma(x)} - \frac{x^H C x}{\varsigma(x)} \frac{x^H \Delta A x}{x^H A x} + \frac{x^H \tilde{A} x}{x^H A x} \delta_3(x, \xi)
\end{aligned}$$

and

$$\epsilon_2 = -\frac{[\tilde{\varsigma}(x) \mp x^H \tilde{B} x](x^H \Delta A x)}{2(x^H \tilde{A} x)(x^H A x)} = \frac{\mp \tilde{\rho}_{\pm}(x)(x^H \Delta A x)}{x^H A x} = -[\pm \rho_{\pm}(x) \pm \Delta \rho_{\pm}(x)] \frac{x^H \Delta A x}{x^H A x}.$$

Noticing

$$\begin{aligned}
\frac{x^H C x}{\varsigma(x)} \pm \rho_{\pm}(x) &= \frac{x^H C x}{\varsigma(x)} \pm \frac{-x^H B x \pm \varsigma(x)}{2(x^H A x)} \\
&= \frac{2(x^H A x)(x^H C x) \mp x^H B x \varsigma(x) + \varsigma(x)^2}{2\varsigma(x)(x^H A x)} \\
&= \frac{(x^H B x)^2 - \varsigma(x)^2 \mp 2(x^H B x)\varsigma(x) + 2\varsigma(x)^2}{4\varsigma(x)(x^H A x)} \\
&= \frac{[x^H B x \mp \varsigma(x)]^2}{4\varsigma(x)(x^H A x)} = \frac{\rho_{\pm}(x)^2 (x^H A x)}{\varsigma(x)},
\end{aligned}$$

we have

$$\pm \Delta \rho_{\pm}(x) = \epsilon_1 + \epsilon_2 = -\delta_2^{\pm}(x) + \frac{x^H \tilde{A} x}{x^H A x} \delta_3(x, \xi) - [\pm \Delta \rho_{\pm}(x)] \frac{x^H \Delta A x}{x^H A x}$$

solving which for  $\pm \Delta \rho_{\pm}(x)$  leads to  $\Delta \rho_{\pm}(x) = \delta^{\pm}(x, \xi)$ .  $\square$

**Lemma 4.4.** *Suppose (4.4) and (4.14) hold. Let  $\delta_{\text{lb}}^{\pm}(x)$ ,  $\delta_{\text{ub}}^{\pm}(x)$ ,  $\tilde{\delta}_{\text{lb}}^{\pm}(x)$ , and  $\tilde{\delta}_{\text{ub}}^{\pm}(x)$  be functions satisfying*

$$\delta_{\text{lb}}^{\pm}(x) \leq \delta^{\pm}(x, \xi) \leq \delta_{\text{ub}}^{\pm}(x), \quad \tilde{\delta}_{\text{lb}}^{\pm}(x) \leq \tilde{\delta}^{\pm}(x, \xi) \leq \tilde{\delta}_{\text{ub}}^{\pm}(x) \quad (4.37)$$

for all  $x \in \mathbb{C}^n$ ,  $\xi \in [0, 1]$ , where  $\delta^{\pm}(x, \xi)$  is defined as in Lemma 4.3. Write

$$\begin{aligned}
\gamma_{\text{uu}}^{\pm} &= \max_{x \neq 0} \{ \delta_{\text{ub}}^{\pm}(x), \tilde{\delta}_{\text{ub}}^{\pm}(x) \}, & \gamma_{\text{ll}}^{\pm} &= \max_{x \neq 0} \{ -\delta_{\text{lb}}^{\pm}(x), -\tilde{\delta}_{\text{lb}}^{\pm}(x) \}, \\
\gamma_{\text{lu}}^{\pm} &= \max_{x \neq 0} \{ -\delta_{\text{lb}}^{\pm}(x), \delta_{\text{ub}}^{\pm}(x) \}, & \tilde{\gamma}_{\text{lu}}^{\pm} &= \max_{x \neq 0} \{ -\tilde{\delta}_{\text{lb}}^{\pm}(x), \tilde{\delta}_{\text{ub}}^{\pm}(x) \}.
\end{aligned}$$

Then

$$\|\Delta \Lambda_{\pm}\|_2 = \max_{1 \leq i \leq n} |\Delta \lambda_i^{\pm}| \leq \min \{ \gamma_{\text{uu}}^{\pm}, \gamma_{\text{ll}}^{\pm}, \gamma_{\text{lu}}^{\pm}, \tilde{\gamma}_{\text{lu}}^{\pm} \}. \quad (4.38)$$

*Proof.* We only consider the “+” case below; the “−” case is similar. In fact simply replacing “+” with “−” gives a proof for the “−” case.

By Lemma 4.3,

$$\delta_{\text{lb}}^+(x) \leq \Delta\rho_+(x) = \delta^+(x, \xi) \leq \delta_{\text{ub}}^+(x).$$

Let  $\mathcal{S}_i = \text{span}\{u_1^+, \dots, u_i^+\}$ ,  $\mathcal{T}_i = \text{span}\{u_i^+, \dots, u_n^+\}$  and similarly define  $\tilde{\mathcal{S}}_i$  and  $\tilde{\mathcal{T}}_i$ . By the Courant-Fischer type min-max principles in Theorem 3.2,

$$\begin{aligned} \lambda_i^+ &= \min_{\dim \mathcal{X}=i} \max_{0 \neq x \in \mathcal{X}} \rho_+(x) = \max_{0 \neq x \in \mathcal{S}_i} \rho_+(x) = \rho_+(u_i^+), \\ \tilde{\lambda}_i^+ &= \min_{\dim \mathcal{X}=i} \max_{0 \neq x \in \mathcal{X}} \tilde{\rho}_+(x) = \max_{0 \neq x \in \tilde{\mathcal{S}}_i} \tilde{\rho}_+(x) = \tilde{\rho}_+(\tilde{u}_i^+), \\ \lambda_i^+ &= \max_{\text{codim } \mathcal{X}=i-1} \min_{0 \neq x \in \mathcal{X}} \rho_+(x) = \min_{0 \neq x \in \mathcal{T}_i} \rho_+(x) = \rho_+(u_i^+), \\ \tilde{\lambda}_i^+ &= \max_{\text{codim } \mathcal{X}=i-1} \min_{0 \neq x \in \mathcal{X}} \tilde{\rho}_+(x) = \min_{0 \neq x \in \tilde{\mathcal{T}}_i} \tilde{\rho}_+(x) = \tilde{\rho}_+(\tilde{u}_i^+). \end{aligned}$$

Therefore,

$$\begin{aligned} \tilde{\lambda}_i^+ &= \min_{\dim \mathcal{X}=i} \max_{0 \neq x \in \mathcal{X}} \tilde{\rho}_+(x) \leq \max_{0 \neq x \in \mathcal{S}_i} \tilde{\rho}_+(x) \\ &\leq \max_{0 \neq x \in \mathcal{S}_i} [\rho_+(x) + \delta_{\text{ub}}^+(x)] \\ &\leq \max_{0 \neq x \in \mathcal{S}_i} \rho_+(x) + \max_{0 \neq x \in \mathcal{S}_i} \delta_{\text{ub}}^+(x) \\ &= \lambda_i^+ + \max_{0 \neq x \in \mathcal{S}_i} \delta_{\text{ub}}^+(x), \\ \tilde{\lambda}_i^+ &= \max_{\text{codim } \mathcal{X}=i-1} \min_{0 \neq x \in \mathcal{X}} \tilde{\rho}_+(x) \geq \min_{0 \neq x \in \mathcal{T}_i} \tilde{\rho}_+(x) \\ &\geq \min_{0 \neq x \in \mathcal{T}_i} [\rho_+(x) + \delta_{\text{lb}}^+(x)] \\ &\geq \min_{0 \neq x \in \mathcal{T}_i} \rho_+(x) + \min_{0 \neq x \in \mathcal{T}_i} \delta_{\text{lb}}^+(x) \\ &= \lambda_i^+ + \min_{0 \neq x \in \mathcal{T}_i} \delta_{\text{lb}}^+(x). \end{aligned}$$

They give (4.39a) below and (4.39b) as well, by switching the roles of  $\mathcal{Q}$  and  $\tilde{\mathcal{Q}}$ :

$$\min_{0 \neq x \in \mathcal{T}_i} \delta_{\text{lb}}^+(x) \leq \tilde{\lambda}_i^+ - \lambda_i^+ \leq \max_{0 \neq x \in \mathcal{S}_i} \delta_{\text{ub}}^+(x), \quad (4.39a)$$

$$\min_{0 \neq x \in \tilde{\mathcal{T}}_i} \tilde{\delta}_{\text{lb}}^+(x) \leq \lambda_i^+ - \tilde{\lambda}_i^+ \leq \max_{0 \neq x \in \tilde{\mathcal{S}}_i} \tilde{\delta}_{\text{ub}}^+(x). \quad (4.39b)$$

It follows from (4.39) that

$$\begin{aligned} |\Delta\lambda_i^+| &\leq \max \left\{ \max_{0 \neq x \in \mathcal{S}_i} \delta_{\text{ub}}^+(x), \max_{0 \neq x \in \tilde{\mathcal{S}}_i} \tilde{\delta}_{\text{ub}}^+(x) \right\} \\ &\leq \max_{x \neq 0} \{\delta_{\text{ub}}^+(x), \tilde{\delta}_{\text{ub}}^+(x)\} = \gamma_{\text{uu}}^+, \\ |\Delta\lambda_i^+| &\leq \max \left\{ -\min_{0 \neq x \in \mathcal{T}_i} \delta_{\text{lb}}^+(x), -\min_{0 \neq x \in \tilde{\mathcal{T}}_i} \tilde{\delta}_{\text{lb}}^+(x) \right\} \end{aligned}$$

$$\begin{aligned}
&\leq \max_{x \neq 0} \{-\delta_{\text{lb}}^+(x), -\tilde{\delta}_{\text{lb}}^+(x)\} = \gamma_{\text{ll}}^+, \\
|\Delta\lambda_i^+| &\leq \max \left\{ -\min_{0 \neq x \in \mathcal{T}_i} \delta_{\text{lb}}^+(x), \max_{0 \neq x \in \mathcal{S}_i} \delta_{\text{ub}}^+(x) \right\} \\
&\leq \max_{x \neq 0} \{-\delta_{\text{lb}}^+(x), \delta_{\text{ub}}^+(x)\} = \gamma_{\text{lu}}^+, \\
|\Delta\lambda_i^+| &\leq \max \left\{ -\min_{0 \neq x \in \tilde{\mathcal{T}}_i} \tilde{\delta}_{\text{lb}}^+(x), \max_{0 \neq x \in \tilde{\mathcal{S}}_i} \tilde{\delta}_{\text{ub}}^+(x) \right\} \\
&\leq \max_{x \neq 0} \{-\tilde{\delta}_{\text{lb}}^+(x), \tilde{\delta}_{\text{ub}}^+(x)\} = \tilde{\gamma}_{\text{lu}}^+.
\end{aligned}$$

This completes the proof of (4.38) for the “+” case.  $\square$

*Proof* of Theorem 4.1. We only prove the perturbation results for  $\Lambda_+$ . The case for  $\Lambda_-$  can be turned into one for  $\Lambda_+$  by considering the pos-type quadratic eigenvalues of  $\mathbf{Q}(-\lambda)$  and  $\tilde{\mathbf{Q}}(-\lambda)$ .

For any  $\alpha > 0, x \neq 0$ , we have

$$\epsilon_a < \alpha \quad \Rightarrow \quad |x^{\text{H}} \Delta A x| < \alpha x^{\text{H}} A x, \quad (4.40a)$$

$$\epsilon_a < \alpha \frac{\chi_\zeta^2}{4\|A\|_2\|C\|_2} \quad \Rightarrow \quad |x^{\text{H}} \Delta A x| < \alpha \frac{\zeta(x)^2}{4|x^{\text{H}} C x|}, \quad (4.40b)$$

$$\epsilon_c < \alpha \frac{\chi_\zeta^2}{4\|A\|_2\|C\|_2} \quad \Rightarrow \quad |x^{\text{H}} \Delta C x| < \alpha \frac{\zeta(x)^2}{4x^{\text{H}} A x}, \quad (4.40c)$$

$$\epsilon_b < \alpha \frac{\chi_\zeta^2}{\|B\|_2(\|B\|_2 + 2\sqrt{\|A\|_2\|C\|_2})} \quad \Rightarrow \quad |x^{\text{H}} \Delta B x| < \alpha |x^{\text{H}} B x|, \quad (4.40d)$$

where (4.40a) and (4.40b) hold because

$$\left| \frac{x^{\text{H}} \Delta A x}{x^{\text{H}} A x} \right| = \left| \frac{x^{\text{H}} A^{1/2} (A^{-1/2} \Delta A A^{-1/2}) A^{1/2} x}{x^{\text{H}} A^{1/2} A^{1/2} x} \right| \leq \|A^{-1/2} \Delta A A^{-1/2}\|_2 = \epsilon_a,$$

and (4.40d) holds because the left part tells

$$|x^{\text{H}} \Delta B x| < \alpha \frac{\zeta(x)^2}{|x^{\text{H}} B x| + \sqrt{4(x^{\text{H}} A x)|x^{\text{H}} C x|}} = \alpha \left( |x^{\text{H}} B x| - \sqrt{4(x^{\text{H}} A x)|x^{\text{H}} C x|} \right). \quad (4.41)$$

For item 1:  $\Delta A = \Delta B = 0$ ,  $\phi(x) = \tilde{\phi}(x) = 0$ ,  $\psi(x) = -2(x^{\text{H}} A x)(x^{\text{H}} \Delta C x)$  and (4.14) holds. Under the assumption (4.18), (4.40c) holds with  $\alpha = 1$ . Thus  $g(1) = \zeta(x)^2 + 2\psi(x) + \phi(x) > 0$ , or equivalently the perturbed quadratic polynomial is still hyperbolic. Note (4.4) holds for  $\phi(x) = 0$ . Thus  $\delta_3(x, \xi) \leq 0$  and  $\tilde{\delta}_3(x, \xi) \leq 0$ . We can take, in (4.37),

$$\delta_{\text{ub}}^+(x) = -\delta_2^+(x) = -\frac{x^{\text{H}} \Delta C x}{\zeta(x)}, \quad \tilde{\delta}_{\text{ub}}^+(x) = -\tilde{\delta}_2^+(x) = \frac{x^{\text{H}} \Delta C x}{\tilde{\zeta}(x)} \quad (4.42)$$

to give

$$|\delta_{\text{ub}}^+(x)| \leq \frac{\|\Delta C\|_2}{\min_{x \neq 0} \zeta_0(x)}, \quad |\tilde{\delta}_{\text{ub}}^+(x)| \leq \frac{\|\Delta C\|_2}{\min_{x \neq 0} \tilde{\zeta}_0(x)}.$$

Using (4.38), we have  $\|\Delta\Lambda_+\|_2 \leq \gamma_{\text{uu}}^+$  to get (4.19).

For item 2:  $\Delta B = \Delta C = 0$ ,  $\phi(x) = \tilde{\phi}(x) = 0$ ,  $\psi(x) = -2(x^{\text{H}}Cx)(x^{\text{H}}\Delta Ax)$ . Under the assumption (4.20), (4.14) holds; (4.40a) and (4.40b) hold with  $\alpha = 1$ . Thus  $g(1) = \varsigma(x)^2 + 2\psi(x) + \phi(x) > 0$ , or equivalently the perturbed quadratic polynomial is still hyperbolic. Note (4.4) holds for  $\phi(x) = 0$ . Thus  $\delta_3(x, \xi) \leq 0$  and  $\tilde{\delta}_3(x, \xi) \leq 0$ . We can take, in (4.37),

$$\begin{aligned}\delta_{\text{ub}}^+(x) &= -\frac{x^{\text{H}}Ax}{x^{\text{H}}\tilde{A}x} \delta_2^+(x) = -\frac{x^{\text{H}}Ax}{x^{\text{H}}\tilde{A}x} \frac{\rho_+(x)^2(x^{\text{H}}\Delta Ax)}{\varsigma(x)}, \\ \tilde{\delta}_{\text{ub}}^+(x) &= -\frac{x^{\text{H}}\tilde{A}x}{x^{\text{H}}Ax} \tilde{\delta}_2^+(x) = \frac{x^{\text{H}}\tilde{A}x}{x^{\text{H}}Ax} \frac{\tilde{\rho}_+(x)^2(x^{\text{H}}\Delta Ax)}{\tilde{\varsigma}(x)},\end{aligned}$$

along with (4.32), to give

$$|\delta_{\text{ub}}^+(x)| \leq \frac{1}{1 - \epsilon_a} \frac{(\lambda_{\text{max}}^+)^2 \|\Delta A\|_2}{\min_{x \neq 0} \varsigma_0(x)}, \quad |\tilde{\delta}_{\text{ub}}^+(x)| \leq (1 + \epsilon_a) \frac{(\tilde{\lambda}_{\text{max}}^+)^2 \|\Delta A\|_2}{\min_{x \neq 0} \tilde{\varsigma}_0(x)}.$$

Using (4.38), we have  $\|\Delta\Lambda_+\|_2 \leq \gamma_{\text{uu}}^+$  to get (4.21).

For item 3:  $\Delta A = \Delta C = 0$ ,  $\phi(x) = \tilde{\phi}(x) = (x^{\text{H}}Bx)(x^{\text{H}}\Delta Bx)$ ,  $\psi(x) = (x^{\text{H}}\Delta Bx)^2$  and (4.14) holds. Under the assumption (4.22), (4.40d) and (4.41) hold with  $\alpha = 1$ . (4.41) tells

$$\sqrt{4(x^{\text{H}}Ax)|x^{\text{H}}Cx|} < |x^{\text{H}}Bx| - |x^{\text{H}}\Delta Bx| \leq |x^{\text{H}}Bx + x^{\text{H}}\Delta Bx|.$$

Thus

$$\begin{aligned}g(1) &= \varsigma(x)^2 + 2\psi(x) + \phi(x) \\ &= (x^{\text{H}}\Delta Bx)^2 + 2(x^{\text{H}}\Delta Bx)(x^{\text{H}}Bx) + (x^{\text{H}}Bx)^2 - 4(x^{\text{H}}Ax)(x^{\text{H}}Cx) \\ &\geq \left[ x^{\text{H}}\Delta Bx + x^{\text{H}}Bx - \sqrt{4(x^{\text{H}}Ax)|x^{\text{H}}Cx|} \right] \left[ x^{\text{H}}\Delta Bx + x^{\text{H}}Bx + \sqrt{4(x^{\text{H}}Ax)|x^{\text{H}}Cx|} \right] \\ &> 0,\end{aligned}$$

or equivalently the perturbed quadratic polynomial is still hyperbolic. (4.40d) tells  $|\psi(x)| = |x^{\text{H}}Bx| > |x^{\text{H}}\Delta Bx| = \phi(x)$ . Thus (4.4) holds. Notice

$$\begin{aligned}\varsigma(x)^2\phi(x) - \psi(x)^2 &= \varsigma(x)^2(x^{\text{H}}\Delta Bx)^2 - [(x^{\text{H}}Bx)(x^{\text{H}}\Delta Bx)]^2 \\ &= -4(x^{\text{H}}Ax)(x^{\text{H}}Cx)(x^{\text{H}}\Delta Bx)^2\end{aligned}$$

to get

$$\delta_3(x, \xi) = -\frac{(x^{\text{H}}Cx)(x^{\text{H}}\Delta Bx)^2}{[f(\xi; x)]^3},$$

where  $f(\xi; x) = [\varsigma(x)^2 + 2\psi(x)\xi + \phi(x)\xi^2]^{1/2}$ . Since<sup>6</sup>

$$\min_{0 \leq \xi \leq 1} f(\xi; x) = \min\{f(0), f(1)\} = \min\{\varsigma(x), \tilde{\varsigma}(x)\}, \quad (4.43)$$

<sup>6</sup>For the quadratic function  $h(t) = a(t - c)^2 + b$  with  $a > 0$ , if  $|c| \geq 1$ , i.e.,  $c$ , the minimal point of  $h(t)$  for  $t \in \mathbb{R}$ , is not in the interval  $(0, 1)$ , then the minimal point of  $h(t)$  on  $[0, 1]$  must be either 0 or 1. For the case here,  $c = \psi(x)/\phi(x)$ .

we can take, in (4.37),

$$\begin{aligned}\delta_{\text{ub}}^+(x) &= -\delta_2^+(x) + \frac{|x^H C x| |x^H \Delta B x|^2}{\min\{\varsigma(x), \tilde{\varsigma}(x)\}^3} = -\frac{\rho_+(x)(x^H \Delta B x)}{\varsigma(x)} + \frac{|x^H C x| |x^H \Delta B x|^2}{\min\{\varsigma(x), \tilde{\varsigma}(x)\}^3}, \\ \tilde{\delta}_{\text{ub}}^+(x) &= -\tilde{\delta}_2^+(x) + \frac{|x^H \tilde{C} x| |x^H \Delta B x|^2}{\min\{\varsigma(x), \tilde{\varsigma}(x)\}^3} = \frac{\tilde{\rho}_+(x)(x^H \Delta B x)}{\tilde{\varsigma}(x)} + \frac{|x^H \tilde{C} x| |x^H \Delta B x|^2}{\min\{\varsigma(x), \tilde{\varsigma}(x)\}^3}\end{aligned}$$

to give

$$\begin{aligned}|\delta_{\text{ub}}^+(x)| &\leq \frac{\lambda_{\max}^+}{\min_{x \neq 0} \varsigma_0(x)} \|\Delta B\|_2 + \frac{\|C\|_2}{\chi_\varsigma^3} \|\Delta B\|_2^2, \\ |\tilde{\delta}_{\text{ub}}^+(x)| &\leq \frac{\tilde{\lambda}_{\max}^+}{\min_{x \neq 0} \tilde{\varsigma}_0(x)} \|\Delta B\|_2 + \frac{\|\tilde{C}\|_2}{\chi_{\tilde{\varsigma}}^3} \|\Delta B\|_2^2.\end{aligned}$$

Using (4.38), we have  $\|\Delta A_+\|_2 \leq \gamma_{\text{uu}}^+$  to get (4.23).

For item 4:  $\Delta A = \Delta C = 0$ , consider the shifted  $\mathbf{Q}_{\lambda_0}(\lambda)$ . By item 2 of Lemma 4.2,  $\mathbf{Q}_{\lambda_0}(\lambda)$  and  $\tilde{\mathbf{Q}}_{\lambda_0}(\lambda)$  are overdamped for  $\lambda_0 \in (-\infty, \min\{\lambda_1^-, \tilde{\lambda}_1^-\}] \cup [\max\{\lambda_n^+, \tilde{\lambda}_n^+\}, +\infty)$ . In particular,  $B_{\lambda_0} \succ 0, C_{\lambda_0} \succeq 0; \tilde{B}_{\lambda_0} \succ 0, \tilde{C}_{\lambda_0} \succeq 0$ . Note  $\varsigma_{\lambda_0}(x) \equiv \varsigma(x), \tilde{\varsigma}_{\lambda_0}(x) \equiv \tilde{\varsigma}(x)$ . Under the assumption (4.24), like<sup>7</sup> in item 3,  $|\psi_{\lambda_0}(x)| > \phi_{\lambda_0}(x)$ . Thus (4.4) for  $\mathbf{Q}_{\lambda_0}(\lambda)$  and  $\tilde{\mathbf{Q}}_{\lambda_0}(\lambda)$  holds. Just as in item 3 (note  $\Delta B_{\lambda_0} = \Delta B$  since  $\Delta A = 0$ ),

$$\varsigma_{\lambda_0}(x)^2 \phi_{\lambda_0}(x) - \psi_{\lambda_0}(x)^2 = -4(x^H A x)(x^H C_{\lambda_0} x)(x^H \Delta B x)^2 < 0$$

which infers  $\delta_{3;\lambda_0}(x, \xi) \leq 0$  and thus we can take, in (4.37),

$$\begin{aligned}\delta_{\text{ub};\lambda_0}^+(x) &= -\delta_{2;\lambda_0}^+(x) = -\frac{\rho_{+;\lambda_0}(x)(x^H \Delta B x)}{\varsigma(x)}, \\ \tilde{\delta}_{\text{ub};\lambda_0}^+(x) &= -\tilde{\delta}_{2;\lambda_0}^+(x) = -\frac{\tilde{\rho}_{+;\lambda_0}(x)(x^H \Delta B x)}{\tilde{\varsigma}(x)}\end{aligned}$$

to give

$$|\delta_{\text{ub};\lambda_0}^+(x)| \leq \frac{\lambda_{\max;\lambda_0}^+}{\min_{x \neq 0} \varsigma_0(x)} \|\Delta B\|_2, \quad |\tilde{\delta}_{\text{ub};\lambda_0}^+(x)| \leq \frac{\tilde{\lambda}_{\max;\lambda_0}^+}{\min_{x \neq 0} \tilde{\varsigma}_0(x)} \|\Delta B\|_2.$$

Using (4.38), we have  $\|\Delta A_{+;\lambda_0}\|_2 \leq \gamma_{\text{uu};\lambda_0}^+$  to get (4.25).

For item 5, under the assumption (4.26),  $\epsilon_a < \gamma < 1$  and (4.40) holds with  $\alpha = \gamma$ . Then (4.14) holds, and

$$\begin{aligned}|\psi(x)| &\leq |x^H B x| |x^H \Delta B x| + 2(x^H A x) |x^H \Delta C x| + 2|x^H C x| |x^H \Delta A x| \\ &< |x^H B x|^2 \gamma + \frac{\varsigma(x)^2}{2} \gamma + \frac{\varsigma(x)^2}{2} \gamma \\ &= [|x^H B x|^2 + \varsigma(x)^2] \gamma, \\ |\phi(x)| &\leq |x^H \Delta B x|^2 + 4|x^H \Delta A x| |x^H \Delta C x|\end{aligned}$$

<sup>7</sup>We will use the same symbols as those for  $\mathbf{Q}$  but with the subscript “ $\lambda_0$ ” to represent the corresponding quantities for  $\mathbf{Q}_{\lambda_0}$ .



$$\begin{aligned}
&< |x^H Bx|^2 \gamma^2 + |x^H \Delta Ax| \frac{\varsigma(x)^2 \gamma}{x^H Ax} \\
&< |x^H Bx|^2 \gamma^2 + \varsigma(x)^2 \gamma^2 \\
&= [|x^H Bx|^2 + \varsigma(x)^2] \gamma^2,
\end{aligned}$$

which infers

$$\begin{aligned}
g(1) &= \varsigma(x)^2 + 2\psi(x) + \phi(x) \\
&> \varsigma(x)^2(1 - 2\gamma - \gamma^2) - |x^H Bx|^2(2\gamma + \gamma^2) \\
&\geq (x^H x)^2 [\chi_\varsigma^2(1 - 2\gamma - \gamma^2) - \|B\|_2^2(2\gamma + \gamma^2)] \\
&= (x^H x)^2 [\chi_\varsigma^2 - (\|B\|_2^2 + \chi_\varsigma^2)(2\gamma + \gamma^2)] \\
&= 0,
\end{aligned}$$

or equivalently the perturbed quadratic polynomial is still hyperbolic. By the same reasoning we had for items 1, 2 and 3, (4.4) holds and at the same time, we have (4.43). Note that

$$\begin{aligned}
\varsigma(x)^2 \phi(x) - \psi(x)^2 &= -4[(x^H Ax)(x^H \Delta Cx) - (x^H Cx)(x^H \Delta Ax)]^2 \\
&\quad - 4[(x^H Ax)(x^H \Delta Bx) - (x^H Bx)(x^H \Delta Ax)] \times \\
&\quad \quad [(x^H Cx)(x^H \Delta Bx) - (x^H Bx)(x^H \Delta Cx)],
\end{aligned}$$

and similarly

$$\begin{aligned}
\tilde{\varsigma}(x)^2 \tilde{\phi}(x) - \tilde{\psi}(x)^2 &= -4[-(x^H \tilde{A}x)(x^H \Delta Cx) + (x^H \tilde{C}x)(x^H \Delta Ax)]^2 \\
&\quad - 4[-(x^H \tilde{A}x)(x^H \Delta Bx) + (x^H \tilde{B}x)(x^H \Delta Ax)] \times \\
&\quad \quad [-(x^H \tilde{C}x)(x^H \Delta Bx) + (x^H \tilde{B}x)(x^H \Delta Cx)] \\
&= \varsigma(x)^2 \phi(x) - \psi(x)^2.
\end{aligned}$$

Now take

$$\begin{aligned}
\delta_{\text{ub}}^+(x) &= -\frac{x^H Ax}{x^H \tilde{A}x} \delta_2^+(x) + \frac{|\varsigma(x)^2 \phi(x) - \psi(x)^2|}{(x^H \tilde{A}x) \min\{\varsigma(x), \tilde{\varsigma}(x)\}^3}, \\
\tilde{\delta}_{\text{ub}}^+(x) &= -\frac{x^H \tilde{A}x}{x^H Ax} \tilde{\delta}_2^+(x) + \frac{|\varsigma(x)^2 \phi(x) - \psi(x)^2|}{(x^H Ax) \min\{\varsigma(x), \tilde{\varsigma}(x)\}^3}
\end{aligned}$$

in (4.37). Note

$$\left| \frac{x^H \Delta Ax}{x^H Ax} \right| \leq \epsilon_a,$$

we have

$$|\varsigma(x)^2 \phi(x) - \psi(x)^2| \leq 4(x^H Ax)^2 \|C\|_2^2 [\epsilon_c + \epsilon_a]^2 + 4(x^H Ax) \|B\|_2^2 \|C\|_2 [\epsilon_b + \epsilon_a] [\epsilon_b + \epsilon_c].$$

Using (4.38), we have  $\|\Delta \Lambda_+\|_2 \leq \gamma_{\text{uu}}^+$  to get (4.28).  $\square$

#### 4.4 Perturbation bounds in unitarily invariant norms

Our main result of this subsection is Theorems 4.2 and 4.3. The proof of Theorem 4.2 is based on our new Wielandt-Lidskii min-max principles. Since it is rather long, we postpone it after stating both theorems.

**Theorem 4.2.** *Suppose  $\Delta A = \Delta B = 0$  and (4.18) holds, and let*

$$\gamma = (\lambda_1^+ - \lambda_n^-)\lambda_{\min}(A), \quad \tilde{\gamma} = (\tilde{\lambda}_1^+ - \tilde{\lambda}_n^-)\lambda_{\min}(A). \quad (4.44)$$

Then

$$\|\Delta A_{\pm}\|_{\text{ui}} \leq c \cdot \frac{\|\Delta C\|_{\text{ui}}}{\min\{\gamma, \tilde{\gamma}\}}, \quad (4.45)$$

where the constant  $c = 1$  if  $\Delta C$  is semidefinite and  $c = 2$  in general.

The inequality (4.45) can be considered as an extension of (4.19), but a little bit less satisfying in that it does not become (4.19) after specializing the unitarily invariant norm to the spectral norm in two aspects:  $c$  is not always 1 and

$$\min_{x \neq 0} \varsigma_0(x) \geq \gamma$$

which can be a strict inequality. Thus it makes us wonder if the stronger version of (4.45) upon setting  $c = 1$  always and replacing  $\min\{\gamma, \tilde{\gamma}\}$  by  $\chi_{\varsigma}$  holds. But how to settle this question eludes us for now.

Recall Theorem 2.5. The next theorem is a straightforward application of Theorem A.2, where  $\|Z\|_2$  and  $\|\tilde{Z}\|_2$  can be bounded using item 5 of Theorem 2.5.

**Theorem 4.3.** *Let  $\mathcal{A} - \lambda\mathcal{B} = \mathcal{L}_{\mathcal{Q}}(\lambda)$  and  $\tilde{\mathcal{A}} - \lambda\tilde{\mathcal{B}} = \mathcal{L}_{\tilde{\mathcal{Q}}}(\lambda)$ , admitting the eigen-decomposition in (2.16). Then*

$$\|\tilde{A} - A\|_{\text{ui}} \leq \|Z\|_2 \|\tilde{Z}\|_2 \left( \|\tilde{\mathcal{A}} - \mathcal{A}\|_{\text{ui}} + \xi \|\tilde{\mathcal{B}} - \mathcal{B}\|_{\text{ui}} \right), \quad (4.46)$$

where  $\xi = \max\{|\lambda_{\max}^+|, |\lambda_{\max}^-|, |\tilde{\lambda}_{\max}^+|, |\tilde{\lambda}_{\max}^-|\}$ , and  $\lambda_{\max}^{\pm}$  and  $\tilde{\lambda}_{\max}^{\pm}$  are defined by (4.16).

The rest of this subsection is devoted to the proof of Theorem 4.2.

**Lemma 4.5.** *Suppose  $\Delta A = \Delta B = 0$  and (4.18) holds. Let  $\varepsilon_1 \leq \varepsilon_2 \leq \dots \leq \varepsilon_n$  be the eigenvalues of  $\Delta C$ , and  $\gamma$  and  $\tilde{\gamma}$  be given by (4.44).*

1. *Given  $X \in \mathbb{C}^{n \times k}$  with  $\text{rank}(X) = k$ , denote the quadratic eigenvalues of  $X^{\text{H}}\mathcal{Q}(\lambda)X$  by*

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

and the quadratic eigenvalues of  $X^{\text{H}}\tilde{\mathcal{Q}}(\lambda)X$  by  $\tilde{\lambda}_{j,X}^{\pm}$  arranged in the same way. Then

$$-\sum_{i=1}^k \frac{\max\{0, -\varepsilon_1\} + \varepsilon_{n-1+i}}{\tilde{\gamma}} \leq \sum_{i=1}^k \Delta\lambda_{i,X}^+ \leq -\sum_{i=1}^k \frac{\min\{0, -\varepsilon_n\} + \varepsilon_i}{\gamma}, \quad (4.47a)$$

$$\sum_{i=1}^k \frac{\min\{0, -\varepsilon_n\} + \varepsilon_i}{\gamma} \leq \sum_{i=1}^k \Delta\lambda_{i,X}^- \leq \sum_{i=1}^k \frac{\max\{0, -\varepsilon_1\} + \varepsilon_{n-1+i}}{\tilde{\gamma}}. \quad (4.47b)$$

2. For any  $1 \leq i_1 < \dots < i_k \leq n$ ,

$$-\sum_{i=1}^k \frac{\max\{0, -\varepsilon_1\} + \varepsilon_{n+1-i}}{\tilde{\gamma}} \leq \sum_{i=1}^k \Delta\lambda_{i_k}^+ \leq -\sum_{i=1}^k \frac{\min\{0, -\varepsilon_n\} + \varepsilon_i}{\gamma}, \quad (4.48a)$$

$$\sum_{i=1}^k \frac{\min\{0, -\varepsilon_n\} + \varepsilon_i}{\gamma} \leq \sum_{i=1}^k \Delta\lambda_{i_k}^- \leq \sum_{i=1}^k \frac{\max\{0, -\varepsilon_1\} + \varepsilon_{n+1-i}}{\tilde{\gamma}}, \quad (4.48b)$$

*Proof.* The assumption (4.18) guarantees that  $\tilde{\mathbf{Q}}(\lambda)$  is still hyperbolic. Without loss of generality, we may assume that  $X$  has orthonormal columns; otherwise, we consider  $V^H\mathbf{Q}(\lambda)V$  instead, where  $V$  is from a QR decomposition  $X = VR$  of  $X$ ,  $V^H V = I_k$  and  $R \in \mathbb{C}^{k \times k}$ . Evidently  $X^H\mathbf{Q}(\lambda)X$  and  $V^H\mathbf{Q}(\lambda)V$  have the same quadratic eigenvalues.

Recall the linearization (2.5) for  $\mathbf{Q}(\lambda)$ . We linearize

$$\mathbf{Q}_X(\lambda) := X^H\mathbf{Q}(\lambda)X \equiv A_X\lambda^2 + B_X\lambda + C_X$$

in the same way to get

$$\mathcal{A}_X - \lambda\tilde{\mathcal{B}}_X \equiv \begin{bmatrix} -C_X & 0 \\ 0 & A_X \end{bmatrix} - \lambda \begin{bmatrix} B_X & A_X \\ A_X & 0 \end{bmatrix} = \mathcal{L}_{\mathbf{Q}_X}(\lambda).$$

Next we apply Theorem 2.5 to  $\mathbf{Q}_X(\lambda)$  to obtain various associated eigen-decompositions and denote the corresponding quantities by the same symbols as those for  $\mathbf{Q}(\lambda)$  but with the subscript  $X$  to indicate them for  $\mathbf{Q}_X(\lambda)$ . In particular, we will have

$$U_X = [u_{1,X}^+, \dots, u_{k,X}^+], \quad A_{+,X} = \text{diag}(\lambda_{1,X}^+, \lambda_{2,X}^+, \dots, \lambda_{k,X}^+),$$

where  $u_{i,X}^+$  are quadratic eigenvectors of  $\mathbf{Q}_X(\lambda)$ ,  $\varsigma_X(u_{i,X}^+) = 1$ , and

$$S_X = \begin{bmatrix} U_X \\ U_X A_{+,X} \end{bmatrix}, \quad S_X^H \tilde{\mathcal{B}}_X S_X = I_k.$$

Also  $S_X^H \tilde{\tilde{\mathcal{B}}}_X S_X = I_k$  since  $\tilde{\tilde{\mathcal{B}}}_X = \tilde{\mathcal{B}}_X$ . Note that  $U_X \in \mathbb{C}^{k \times k}$  is nonsingular. By Theorems 2.2 and [37, Corollary 2.1],

$$\inf_{Z^H \tilde{\mathcal{B}}_X Z = I_k} \text{trace}(Z^H \mathcal{A}_X Z) = \sum_{i=1}^k \lambda_{i,X}^+ = \text{trace}(S_X^H \mathcal{A}_X S_X).$$

Let  $\varepsilon_{1,X} \leq \dots \leq \varepsilon_{k,X}$  be the eigenvalues of  $\Delta C_X = X^H \Delta C X$ . Since  $X$  has orthonormal columns, we have  $\varepsilon_i \leq \varepsilon_{i,X} \leq \varepsilon_{n-k+i}$  by the Cauchy interlacing theorem, and thus

$$\sum_{i=1}^k \varepsilon_i \leq \sum_{i=1}^k \varepsilon_{i,X} \leq \sum_{i=1}^k \varepsilon_{n+1-i}.$$

For the sake of presentation, we will drop the superscript “+” to  $u_{i,X}^+$  in the rest of this proof. We have

$$\sum_{i=1}^k \tilde{\lambda}_{i,X}^+ = \inf_{Z^H \tilde{\tilde{\mathcal{B}}}_X Z = I_k} \text{trace}(Z^H \tilde{\mathcal{A}}_X Z)$$

$$\begin{aligned}
&\leq \text{trace}(S_X^H \widetilde{\mathcal{A}}_X S_X) \quad (\text{since } S_X^H \widetilde{\mathcal{B}}_X S_X = I_k) \\
&= \text{trace}(S_X^H \mathcal{A}_X S_X) + \text{trace}(S_X^H \Delta \mathcal{A}_X S_X) \\
&= \sum_{i=1}^k \lambda_{i,X}^+ - \text{trace}(U_X^H \Delta C_X U_X). \tag{4.49}
\end{aligned}$$

Let  $\mu = \min\{0, -\varepsilon_n\} \leq 0$ . For any scalar  $\tau_0 \in (0, 1)$ , set  $\tau^2 = \tau_0^2 \gamma = \tau_0^2 (\lambda_1^+ - \lambda_n^-) \lambda_{\min}(A)$ , and

$$\begin{aligned}
E_X &= -\mu U_X^H U_X, & D_X &= U_X^H (U_X^{-H} U_X^{-1} - \tau^2 I) U_X, \\
\mathcal{C}_X &= \begin{bmatrix} \tau^{-2}(\Delta C_X + \mu I) & 0 \\ 0 & E_X \end{bmatrix} \in \mathbb{C}^{2k \times 2k}, & \mathcal{D}_X &= \begin{bmatrix} I & 0 \\ 0 & D_X \end{bmatrix} \in \mathbb{C}^{2k \times 2k}.
\end{aligned}$$

Note that by (2.18a), (2.18e), and (2.24),

$$U_X^H A_X U_X \preceq (\lambda_{1,X}^+ - \lambda_{k,X}^-)^{-1} I \preceq (\lambda_1^+ - \lambda_n^-)^{-1} I$$

which infers

$$U_X^{-H} U_X^{-1} \succeq (\lambda_1^+ - \lambda_n^-) A_X \succeq (\lambda_1^+ - \lambda_n^-) \lambda_{\min}(A_X) I \succeq (\lambda_1^+ - \lambda_n^-) \lambda_{\min}(A) I = \gamma I \succ \tau^2 I.$$

Thus,  $D_X \succ 0$ , and so  $\mathcal{D}_X \succ 0$ . Hence the matrix pencil  $\mathcal{C}_X - \lambda \mathcal{D}_X$  has  $2k$  finite eigenvalues  $\nu_i$  ( $i = 1, \dots, 2k$ ). By the choice of  $\mu$ ,  $\Delta C_X + \mu I \preceq 0$  and  $E_X \succeq 0$ . Therefore these  $\nu_i$  can be ordered as

$$\nu_1 \leq \dots \leq \nu_k \leq 0 \leq \nu_{k+1} \leq \dots \leq \nu_{2k},$$

where  $\nu_i$  for  $i = 1, \dots, k$  are the eigenvalues of  $\tau^{-2}(\Delta C_X + \mu I)$  and  $\nu_i$  for  $i = k+1, \dots, 2k$  are the generalized eigenvalues of  $E_X - \lambda D_X$ . By the Courant-Fischer min-max principle, we have for  $i = 1, \dots, k$

$$\begin{aligned}
\nu_i &= \min_{\dim \mathcal{X}=i} \max_{0 \neq x \in \mathcal{X}} \frac{x^H (\Delta C_X + \mu I) x}{\tau^2 x^H x} \\
&= \frac{1}{\tau^2} \left[ \mu + \min_{\dim \mathcal{X}=i} \max_{0 \neq x \in \mathcal{X}} \frac{x^H \Delta C_X x}{x^H x} \right] \\
&= \frac{1}{\tau^2} [\mu + \varepsilon_{i,X}] \\
&\geq \frac{1}{\tau^2} [\mu + \varepsilon_i] \\
&= \frac{1}{\tau_0^2 \gamma} [\mu + \varepsilon_i].
\end{aligned}$$

By the arbitrary choice of  $\tau_0 \in (0, 1)$ ,

$$\nu_i \geq \frac{\mu + \varepsilon_i}{\gamma}.$$

For the matrix  $T_X := \begin{bmatrix} \tau U_X \\ I \end{bmatrix}$ , we have

$$T_X^H \mathcal{D}_X T_X = \tau^2 U_X^H U_X + D_X = I,$$

$$T_X^H \mathcal{C}_X T_X = \tau^2 \tau^{-2} U_X^H (\Delta C_X + \mu I) U_X + E_X = U_X^H \Delta C_X U_X.$$

Therefore

$$\begin{aligned} \text{trace}(U_X^H \Delta C_X U_X) &= \text{trace}(T_X^H \mathcal{C}_X T_X) \\ &\geq \min_{Z^H \mathcal{D}_X Z = I} \text{trace}(Z^H \mathcal{C}_X Z) \\ &= \sum_{i=1}^k \nu_i. \end{aligned}$$

Thus, (4.49) becomes

$$\sum_{i=1}^k \Delta \lambda_{i,X}^+ \leq - \sum_{i=1}^k \nu_i \leq - \sum_{i=1}^k \frac{\mu + \varepsilon_i}{\gamma} = - \sum_{i=1}^k \frac{\min\{0, -\varepsilon_n\} + \varepsilon_i}{\gamma}. \quad (4.50)$$

Think of  $\mathbf{Q}$  as obtained from perturbing  $\tilde{\mathbf{Q}}$  and apply (4.50) to get

$$- \sum_{i=1}^k \Delta \lambda_{i,X}^+ \leq - \sum_{i=1}^k \frac{\min\{0, -(-\varepsilon_1)\} + (-\varepsilon_{n-1+i})}{\tilde{\gamma}} \quad (4.51)$$

which, combined with (4.50), leads to (4.47a). Apply (4.47a) to  $\mathbf{Q}(-\lambda)$  and  $\tilde{\mathbf{Q}}(-\lambda)$  to get (4.47b).

Now we prove (4.48). With all ‘‘sup’’ being taken over  $\mathcal{X}_1 \subset \dots \subset \mathcal{X}_k$  and  $\text{codim } \mathcal{X}_j = i_j - 1$ , and all ‘‘inf’’ over  $x_j \in \mathcal{X}_j$ ,  $X = [x_1, \dots, x_k]$ , and  $\text{rank}(X) = k$ , we have by Theorem 3.3

$$\begin{aligned} \sum_{j=1}^k \tilde{\lambda}_{i_k}^+ &= \sup \inf \sum_{j=1}^k \tilde{\lambda}_{k,X}^+ \\ &\leq \sup \inf \left[ \sum_{j=1}^k \lambda_{k,X}^+ - \sum_{i=1}^k \frac{\min\{0, -\varepsilon_n\} + \varepsilon_i}{\gamma} \right] \quad (\text{by (4.50)}) \\ &= \sup \inf \sum_{j=1}^k \lambda_{k,X}^+ - \sum_{i=1}^k \frac{\min\{0, -\varepsilon_n\} + \varepsilon_i}{\gamma} \\ &\leq \sum_{j=1}^k \lambda_{i_k}^+ - \sum_{i=1}^k \frac{\min\{0, -\varepsilon_n\} + \varepsilon_i}{\gamma}. \end{aligned} \quad (4.52)$$

Similarly,

$$\sum_{j=1}^k \lambda_{i_k}^+ \leq \sum_{j=1}^k \tilde{\lambda}_{i_k}^+ - \sum_{i=1}^k \frac{\min\{0, -(-\varepsilon_1)\} + (-\varepsilon_{n-1+i})}{\tilde{\gamma}}. \quad (4.53)$$

The inequalities in (4.48a) is a consequence of (4.52) and (4.53). Apply (4.48a) to  $\mathbf{Q}(-\lambda)$  and  $\tilde{\mathbf{Q}}(-\lambda)$  to get (4.48b).  $\square$

**Lemma 4.6.** *Suppose  $\Delta A = \Delta B = 0$  and (4.18) holds. We have for  $1 \leq j \leq n$*

$$\tilde{\lambda}_j^+ \leq \lambda_j^+ \text{ and } \tilde{\lambda}_j^- \geq \lambda_j^- \text{ if } \Delta C \succeq 0, \quad (4.54a)$$

$$\tilde{\lambda}_j^+ \geq \lambda_j^+ \text{ and } \tilde{\lambda}_j^- \leq \lambda_j^- \text{ if } \Delta C \preceq 0. \quad (4.54b)$$

Consequently  $\tilde{\gamma} \leq \gamma$  if  $\Delta C \succeq 0$ , and  $\tilde{\gamma} \geq \gamma$  if  $\Delta C \preceq 0$ .

*Proof.* The assumption (4.18) guarantees that  $\tilde{\mathbf{Q}}(\lambda)$  is still hyperbolic. By (3.2), we see

$$\begin{aligned} \tilde{\rho}_+(x) &\leq \rho_+(x) \text{ and } \tilde{\rho}_-(x) \geq \rho_-(x) \text{ if } \Delta C \succeq 0, \\ \tilde{\rho}_+(x) &\geq \rho_+(x) \text{ and } \tilde{\rho}_-(x) \leq \rho_-(x) \text{ if } \Delta C \preceq 0. \end{aligned}$$

Now use Theorem 3.2 to get (4.54).  $\square$

*Proof* of Theorem 4.2. The assumption (4.18) guarantees that  $\tilde{\mathbf{Q}}(\lambda)$  is still hyperbolic.

As in Lemma 4.5, let  $\varepsilon_1 \leq \varepsilon_2 \leq \dots \leq \varepsilon_n$  be the eigenvalues of  $\Delta C$ .

Consider first the case  $\Delta C \succeq 0$ . Then  $0 \leq \varepsilon_1$ . Also  $\Delta \lambda_i^+ \leq 0$  for all  $i$  by Lemma 4.6. Therefore the leftmost inequality in (4.48a) gives

$$\sum_{i=1}^k |\Delta \lambda_{i_k}^+| \leq \sum_{i=1}^k \frac{\varepsilon_{n+1-i}}{\tilde{\gamma}}$$

for any  $1 \leq i_1 < \dots < i_k \leq n$ . As a result of [56, Theorem II.3.6 and Theorem II.3.17], we have

$$\|\Delta \Lambda_+\|_{\text{ui}} \leq \frac{\|\Delta C\|_{\text{ui}}}{\tilde{\gamma}}. \quad (4.55)$$

Similarly, use the rightmost inequality in (4.48b) to get

$$\|\Delta \Lambda_-\|_{\text{ui}} \leq \frac{\|\Delta C\|_{\text{ui}}}{\tilde{\gamma}}. \quad (4.56)$$

Now we turn to the case  $\Delta C \preceq 0$ . Then  $\varepsilon_n \leq 0$ . Also  $\Delta \lambda_i^+ \geq 0$  for all  $i$  by Lemma 4.6. Therefore the rightmost inequality in (4.48a) gives

$$\sum_{i=1}^k |\Delta \lambda_{i_k}^+| \leq \sum_{i=1}^k \frac{|\varepsilon_i|}{\gamma}$$

for any  $1 \leq i_1 < \dots < i_k \leq n$ . Again as a result of [56, Theorem II.3.6 and Theorem II.3.17], we have

$$\|\Delta \Lambda_+\|_{\text{ui}} \leq \frac{\|\Delta C\|_{\text{ui}}}{\gamma}. \quad (4.57)$$

Similarly, use the leftmost inequality in (4.48b) to get

$$\|\Delta \Lambda_-\|_{\text{ui}} \leq \frac{\|\Delta C\|_{\text{ui}}}{\gamma}. \quad (4.58)$$

The inequalities (4.55) – (4.56) together give (4.45) for the case when  $\Delta C$  is semidefinite.

For the general case when  $\Delta C$  is indefinite, we can decompose  $\Delta C = \Delta C_+ - \Delta C_-$ , where  $\Delta C_{\pm} \succeq 0$  and

$$\text{eig}(\Delta C_+) = \{\max\{0, \varepsilon_i\}, 1 \leq i \leq n\}, \quad \text{eig}(\Delta C_-) = \{\max\{0, -\varepsilon_i\}, 1 \leq i \leq n\}.$$

In particular,  $\|\Delta C_{\pm}\|_{\text{ui}} \leq \|\Delta C\|_{\text{ui}}$ . Let  $\widehat{C} = C - \Delta C_-$  and  $\widehat{Q}(\lambda) = \lambda^2 A + \lambda B + \widehat{C}$ . We claim  $\widehat{Q}(\lambda)$  is hyperbolic. This is because  $\widetilde{C} = C + \Delta C_+ - \Delta C_- \succeq C - \Delta C_- = \widehat{C}$  and thus for any  $x \neq 0$

$$0 < (x^H B x)^2 - 4(x^H A x)(x^H \widetilde{C} x) \leq (x^H B x)^2 - 4(x^H A x)(x^H \widehat{C} x),$$

where the first inequality holds because  $\widetilde{Q}(\lambda)$  is hyperbolic. Apply what we just proved to  $\widehat{Q}$  and  $\widetilde{Q}$  to get

$$\|\widehat{\Lambda}_{\pm} - \Lambda_{\pm}\|_{\text{ui}} \leq \frac{\|\Delta C_-\|_{\text{ui}}}{\gamma} \leq \frac{\|\Delta C\|_{\text{ui}}}{\gamma}, \quad (4.59)$$

where  $\widehat{\Lambda}_{\pm}$  are similarly defined for  $\widehat{Q}$  to  $\Lambda_{\pm}$  for  $\widetilde{Q}$ . Notice  $\widetilde{C} = \widehat{C} + \Delta C_+$  and apply what we just proved to  $\widetilde{Q}$  and  $\widehat{Q}$  to get

$$\|\widetilde{\Lambda}_{\pm} - \widehat{\Lambda}_{\pm}\|_{\text{ui}} \leq \frac{\|\Delta C_+\|_{\text{ui}}}{\widetilde{\gamma}} \leq \frac{\|\Delta C\|_{\text{ui}}}{\widetilde{\gamma}}. \quad (4.60)$$

Finally

$$\|\widetilde{\Lambda}_{\pm} - \Lambda_{\pm}\|_{\text{ui}} \leq \|\widetilde{\Lambda}_{\pm} - \widehat{\Lambda}_{\pm}\|_{\text{ui}} + \|\widehat{\Lambda}_{\pm} - \Lambda_{\pm}\|_{\text{ui}} \leq 2 \cdot \frac{\|\Delta C\|_{\text{ui}}}{\min\{\gamma, \widetilde{\gamma}\}},$$

as was to be shown.  $\square$

#### 4.5 Perturbation bounds in the Frobenius norms

**Theorem 4.4.** *Suppose (4.26) holds and  $\lambda_0 \in (\lambda_n^-, \lambda_1^+) \cap (\widetilde{\lambda}_n^-, \widetilde{\lambda}_1^+)$ . Then*

$$\|\Delta A\|_{\text{F}}^2 \leq 2 \left( \chi_1^2 \zeta_1^2 \chi_2^4 \|\Delta A\|_{\text{F}}^2 + \chi_2^2 \|\Delta B_{\lambda_0}\|_{\text{F}}^2 + \chi_3^2 \zeta_2^2 \|\Delta C_{\lambda_0}\|_{\text{F}}^2 \right), \quad (4.61)$$

where

$$\begin{aligned} \chi_1 &= \sqrt{\|C_{\lambda_0}\|_2 + \|\widetilde{C}_{\lambda_0}\|_2 + \left( \|\widetilde{A}^{-1/2} \widetilde{B}_{\lambda_0}\|_2 + \|A^{-1/2} B_{\lambda_0}\|_2 \right)^2}, \\ \chi_2 &= \sqrt{\|A^{-1}\|_2 \|\widetilde{A}^{-1}\|_2}, \\ \chi_3 &= \sqrt{\|A^{-1}\|_2 + \|\widetilde{A}^{-1}\|_2}, \\ \zeta_1 &= \frac{1}{\|A\|_2^{-1/2} + \|\widetilde{A}\|_2^{-1/2}}, \\ \zeta_2 &= \frac{1}{\|C_{\lambda_0}^{-1}\|_2^{-1/2} + \|\widetilde{C}_{\lambda_0}^{-1}\|_2^{-1/2}}. \end{aligned}$$

In particular, if  $\Delta A = 0$ , then the scalar 2 in (4.61) can be replaced by 1 to give

$$\|\Delta A\|_{\text{F}}^2 \leq \|A^{-1}\|_2^2 \|\Delta B\|_{\text{F}}^2 + 2 \|A^{-1}\|_2 \zeta_2^2 \|\lambda_0 \Delta B + \Delta C\|_{\text{F}}^2. \quad (4.62)$$

*Proof.* The assumptions in (4.26) guarantee that  $\tilde{\mathbf{Q}}(\lambda)$  is still hyperbolic. The assumption  $\lambda_0 \in (\lambda_n^-, \lambda_1^+) \cap (\tilde{\lambda}_n^-, \tilde{\lambda}_1^+)$  ensures  $\tilde{C}_{\lambda_0} \succeq 0$  and  $C_{\lambda_0} \succeq 0$ . Without loss of generality, we may assume that both  $\mathbf{Q}(\lambda)$  and  $\tilde{\mathbf{Q}}(\lambda)$  have already been shifted, or equivalently  $\lambda_0 = 0$ . This allows to drop the potential subscript “ $\lambda_0$ ” to  $B_{\lambda_0}$ ,  $C_{\lambda_0}$ , etc.

Note that  $\mathcal{K}_{\mathbf{Q}}(\lambda) = \mathcal{A} - \lambda\mathcal{B}$  as in (2.6), where  $\mathcal{B} = \text{diag}(-C, A) \succ 0$ . Thus, the hyperbolic eigenvalue problem  $\mathbf{Q}(\lambda)$  is equivalent to the Hermitian eigenvalue problem  $\mathbf{K}$  where

$$\mathbf{K} = \mathcal{B}^{-1/2} \mathcal{A} \mathcal{B}^{-1/2} = \begin{bmatrix} 0 & [-C]^{1/2} A^{-1/2} \\ A^{-1/2} [-C]^{1/2} & -A^{-1/2} B A^{-1/2} \end{bmatrix}.$$

By the Hoffman-Wielandt theorem [27, 56],

$$\|\Delta A\|_{\mathbb{F}}^2 \leq \|\Delta \mathbf{K}\|_{\mathbb{F}}^2 = 2 \left\| \Delta \left( [-C]^{1/2} A^{-1/2} \right) \right\|_{\mathbb{F}}^2 + \left\| \Delta \left( A^{-1/2} B A^{-1/2} \right) \right\|_{\mathbb{F}}^2. \quad (4.63)$$

The rest of the proof is just to bound the two terms in the right-hand side of (4.63). To this end, we note that

$$\begin{aligned} \left\| \Delta \left( [-C]^{1/2} A^{-1/2} \right) \right\|_{\mathbb{F}} &= \left\| [-\tilde{C}]^{1/2} \Delta(A^{-1/2}) + \Delta([-C]^{1/2}) A^{-1/2} \right\|_{\mathbb{F}} \\ &\leq \|[-\tilde{C}]^{1/2}\|_2 \|\Delta(A^{-1/2})\|_{\mathbb{F}} + \|A^{-1/2}\|_2 \left\| \Delta \left( [-C]^{1/2} \right) \right\|_{\mathbb{F}}, \end{aligned} \quad (4.64)$$

and similarly

$$\begin{aligned} \left\| \Delta \left( A^{-1/2} B A^{-1/2} \right) \right\|_{\mathbb{F}} &= \left\| \Delta([-C]^{1/2}) \tilde{A}^{-1/2} + [-C]^{1/2} \Delta(A^{-1/2}) \right\|_{\mathbb{F}} \\ &\leq \|\tilde{A}^{-1/2}\|_2 \left\| \Delta \left( [-C]^{1/2} \right) \right\|_{\mathbb{F}} + \|[-C]^{1/2}\|_2 \left\| \Delta(A^{-1/2}) \right\|_{\mathbb{F}}. \end{aligned} \quad (4.65)$$

Also,

$$\begin{aligned} \left\| \Delta \left( A^{-1/2} B A^{-1/2} \right) \right\|_{\mathbb{F}} &= \left\| \tilde{A}^{-1/2} \tilde{B} \Delta(A^{-1/2}) + \tilde{A}^{-1/2} \Delta B A^{-1/2} + \Delta(A^{-1/2}) B A^{-1/2} \right\|_{\mathbb{F}} \\ &\leq \left( \|\tilde{A}^{-1/2} \tilde{B}\|_2 + \|A^{-1/2} B\|_2 \right) \|\Delta(A^{-1/2})\|_{\mathbb{F}} \\ &\quad + \|\tilde{A}^{-1/2}\|_2 \|A^{-1/2}\|_2 \|\Delta B\|_{\mathbb{F}}. \end{aligned} \quad (4.66)$$

Combine<sup>8</sup> (4.63) – (4.66) to get

$$\|\Delta A\|_{\mathbb{F}}^2 \leq 2 \left[ \chi_1^2 \|\Delta(A^{-1/2})\|_{\mathbb{F}}^2 + \chi_2^2 \|\Delta B\|_{\mathbb{F}}^2 + \chi_3^2 \left\| \Delta \left( [-C]^{1/2} \right) \right\|_{\mathbb{F}}^2 \right]. \quad (4.67)$$

By [53],

$$\begin{aligned} \|\Delta(A^{-1/2})\|_{\mathbb{F}} &\leq \zeta_1 \|\Delta(A^{-1})\|_{\mathbb{F}} \\ &\leq \zeta_1 \|\tilde{A}^{-1}\|_2 \|\Delta A\|_{\mathbb{F}} \|\tilde{A}^{-1}\|_2, \end{aligned} \quad (4.68)$$

$$\left\| \Delta \left( [-C]^{1/2} \right) \right\|_{\mathbb{F}} \leq \zeta_2 \|\Delta C\|_{\mathbb{F}}, \quad (4.69)$$

where the inequality sign in (4.68) is due to  $\Delta(A^{-1}) = -\tilde{A}^{-1} \Delta A A^{-1}$ . Now substitute (4.68) and (4.69) into (4.67) to yield the desired inequality.  $\square$

<sup>8</sup>Actually we only use this:  $(a+b)^2 \leq 2(a^2+b^2)$  which results in the scalar 2 in (4.67).



Theorem 4.4 gives a perturbation result for all quadratic eigenvalues of  $\mathbf{Q}(\lambda)$ . However, using a different approach, we can obtain results in the Frobenius for only pos- or neg-type quadratic eigenvalues of  $\mathbf{Q}(\lambda)$ .

Following [20], we know the matrix equation

$$AX^2 + BX + C = 0$$

has two special solutions. One has all pos-type quadratic eigenvalues of  $\mathbf{Q}(\lambda)$  as its eigenvalues while the other has all neg-type quadratic eigenvalues of  $\mathbf{Q}(\lambda)$  as its eigenvalues. We call the first special solution the *pos-type* solution and the second special solution the *neg-type* solution.

Consider  $\mathbf{Q}_{\lambda_0}(\lambda)$  and set

$$B_A = A^{-1/2}B_{\lambda_0}A^{-1/2}, \quad C_A = A^{-1/2}C_{\lambda_0}A^{-1/2}. \quad (4.70)$$

Because  $A^{-1/2}\mathbf{Q}_{\lambda_0}(\lambda)A^{-1/2} = \lambda^2I + \lambda B_A + C_A$  is hyperbolic, the following equation

$$X^2 + B_AX + C_A = 0, \quad (4.71)$$

has the pos- and neg-type solutions. Denote them by  $R_{\pm}$ , respectively, in the rest of this section. Both  $R_{\pm}$  can be expressed explicitly by the quantities defined in Theorem 2.5. In fact,

$$R_{\pm} := A^{1/2}U_{\pm}(A_{\pm} - \lambda_0I)U_{\pm}^{-1}A^{-1/2}. \quad (4.72)$$

**Lemma 4.7.** *Suppose (4.26) holds and  $\lambda_0 \geq \max\{\lambda_n^+, \tilde{\lambda}_n^+\}$ . Let  $\text{typ} \in \{+, -\}$ . If*

$$\eta := 2\lambda_0 - \tilde{\lambda}_n^+ - \tilde{\lambda}_n^- - \|\tilde{R}_{\text{typ}}\|_2 - \|R_{\text{typ}}\|_2 > 0, \quad (4.73)$$

then

$$\|\Delta R_{\text{typ}}\|_F \leq \frac{\chi_4 \zeta_1 \chi_2^2}{\eta} \|\Delta A\|_F + \frac{\chi_2}{\eta} (\|R_{\text{typ}}\|_2 \|\Delta B_{\lambda_0}\|_F + \|\Delta C_{\lambda_0}\|_F), \quad (4.74)$$

where

$$\chi_4 = \|R_{\text{typ}}\|_2 (\|\tilde{A}^{-1/2}\tilde{B}_{\lambda_0}\|_2 + \|A^{-1/2}B_{\lambda_0}\|_2) + \|\tilde{A}^{-1/2}\tilde{C}_{\lambda_0}\|_2 + \|A^{-1/2}C_{\lambda_0}\|_2,$$

and  $\chi_2, \zeta_1$  are as in Theorem 4.4.

*Proof.* The assumptions in (4.26) guarantee that  $\tilde{\mathbf{Q}}(\lambda)$  is still hyperbolic. By (4.8) and (4.9),

$$\text{eig}(B_A) \in [2\lambda_0 - \lambda_n^- - \lambda_n^+, 2\lambda_0 - \lambda_1^- - \lambda_1^+], \quad (4.75)$$

$$\text{eig}(C_A) \in [(\lambda_0 - \lambda_n^-)(\lambda_0 - \lambda_n^+), (\lambda_0 - \lambda_1^-)(\lambda_0 - \lambda_1^+)]. \quad (4.76)$$

Subtract  $\tilde{R}_{\text{typ}}^2 + \tilde{B}_A\tilde{R}_{\text{typ}} + \tilde{C}_A = 0$  from  $R_{\text{typ}}^2 + B_AR_{\text{typ}} + C_A = 0$  to get

$$(\tilde{R}_{\text{typ}} + \tilde{B}_A)\Delta R_{\text{typ}} + (\Delta R_{\text{typ}})R_{\text{typ}} = -(\Delta B_A)R_{\text{typ}} - \Delta C_A,$$

or equivalently

$$\left[ I \otimes (\tilde{R}_{\text{typ}} + \tilde{B}_A) + R_{\text{typ}}^T \otimes I \right] \text{vec}(\Delta R_{\text{typ}}) = -\text{vec}((\Delta B_A)R_{\text{typ}} - \Delta C_A), \quad (4.77)$$

where  $\text{vec}(\cdot)$  turns a matrix to a vector by appending the columns of the matrix one after another with the first column followed by the second column and so on. The equation (4.77) yields

$$\begin{aligned}\|\Delta R_{\text{typ}}\|_{\text{F}} &\leq \left\| \left[ I \otimes (\tilde{R}_{\text{typ}} + \tilde{B}_A) + R_{\text{typ}}^{\text{T}} \otimes I \right]^{-1} \right\|_2 \|(\Delta B_A)R_{\text{typ}} - \Delta C_A\|_{\text{F}} \\ &\leq \left\| \left[ I \otimes (\tilde{R}_{\text{typ}} + \tilde{B}_A) + R_{\text{typ}}^{\text{T}} \otimes I \right]^{-1} \right\|_2 (\|R_{\text{typ}}\|_2 \|\Delta B_A\|_{\text{F}} + \|\Delta C_A\|_{\text{F}}). \quad (4.78)\end{aligned}$$

Choose a  $\tau \leq \tilde{\lambda}_1^+ + \tilde{\lambda}_1^- - 2\lambda_0 \leq \tilde{\lambda}_n^+ + \tilde{\lambda}_n^- - 2\lambda_0 = -\eta - \|\tilde{R}_{\text{typ}}\|_2 - \|R_{\text{typ}}\|_2 < 0$ . Then

$$\begin{aligned}\|I \otimes \tilde{R}_{\text{typ}} + I \otimes (\tilde{B}_A + \tau I) + R_{\text{typ}}^{\text{T}} \otimes I\|_2 &\leq \|\tilde{R}_{\text{typ}}\|_2 + \|R_{\text{typ}}^{\text{T}}\|_2 + \|\tilde{B}_A + \tau I\|_2 \\ &\leq \|\tilde{R}_{\text{typ}}\|_2 + \|R_{\text{typ}}\|_2 + \tilde{\lambda}_n^+ + \tilde{\lambda}_n^- - 2\lambda_0 - \tau \\ &< -\eta - \tau < -\tau = |\tau|\end{aligned}$$

from which we infer

$$\begin{aligned}&\left\| \left( I \otimes (\tilde{R}_{\text{typ}} + \tilde{B}_A) + R_{\text{typ}}^{\text{T}} \otimes I \right)^{-1} \right\|_2 \\ &= \left\| \tau^{-1} \left( \frac{I \otimes \tilde{R}_{\text{typ}} + I \otimes (\tilde{B}_A + \tau I) + R_{\text{typ}}^{\text{T}} \otimes I}{\tau} - I \otimes I \right)^{-1} \right\|_2 \\ &\leq \frac{|\tau|^{-1}}{1 - |\tau|^{-1} \|I \otimes \tilde{R}_{\text{typ}} + I \otimes (\tilde{B}_A + \tau I) + R_{\text{typ}}^{\text{T}} \otimes I\|_2} \\ &= \frac{1}{-\tau - \|I \otimes \tilde{R}_{\text{typ}} + I \otimes (\tilde{B}_A + \tau I) + R_{\text{typ}}^{\text{T}} \otimes I\|_2} \\ &\leq \frac{1}{-\tau - (-\eta - \tau)} = \frac{1}{\eta}. \quad (4.79)\end{aligned}$$

Like (4.66), (4.68) and (4.69), we can obtain the estimates of  $\|\Delta B_A\|_{\text{F}}$  and  $\|\Delta C_A\|_{\text{F}}$ . Then (4.74) follows.  $\square$

$R_{\text{typ}}$  is diagonalizable by (4.72). By [7, Theorem 3.1], we have

$$\|\Delta \Lambda_{\text{typ}}\|_{\text{F}} \leq \kappa \|\Delta R_{\text{typ}}\|_{\text{F}}, \quad (4.80)$$

where

$$\kappa = \sqrt{\kappa_2(A^{1/2}U_{\text{typ}})\kappa_2(\tilde{A}^{1/2}\tilde{U}_{\text{typ}})}. \quad (4.81)$$

**Theorem 4.5.** *Suppose  $\lambda_0 \geq \max\{\lambda_n^+, \tilde{\lambda}_n^+\}$ . If (4.73) holds, then*

$$\|\Delta \Lambda_{\text{typ}}\|_{\text{F}} \leq \frac{\kappa \chi_4 \zeta_1 \chi_2^2}{\eta} \|\Delta A\|_{\text{F}} + \frac{\kappa \chi_2}{\eta} (\|R_{\text{typ}}\|_2 \|\Delta B_{\lambda_0}\|_{\text{F}} + \|\Delta C_{\lambda_0}\|_{\text{F}}) \quad (4.82)$$

where  $\kappa, \eta, \chi_4, \chi_2, \zeta_1$  are as in Lemma 4.7 and (4.81).

## 5 Best approximations from a subspace and Rayleigh-Ritz procedure

Two most important aspects in solving a large scale eigenvalue problem are

1. building subspaces to which the desired eigenvectors (or invariant subspaces) are close, and
2. seeking “*best possible*” approximations from the suitably built subspaces.

In this section, we shall address the second aspect for our current problem at hand, i.e., seeking “*best possible*” approximations to a few quadratic eigenvalues of  $\mathbf{Q}(\lambda)$  and their associated quadratic eigenvectors from a given subspace of  $\mathbb{C}^n$ . We leave the first aspect to the later sections when we present our computational algorithms.

The concept of “*best possible*” comes with a quantitative measure as to what constitutes “*best possible*”. There may not be such a measure in general. In [47, section 11.4], Parlett uses three different ways to justify the use of the Rayleigh-Ritz procedure for the symmetric eigenvalue problem. For the HQEP here, each of the minimization principles in section 3 provides a quantitative measure.

Let  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C \in \mathbb{C}^{n \times n}$  be a hyperbolic quadratic matrix polynomial, and let  $\mathcal{Y} \subset \mathbb{C}^n$  be a subspace of dimension  $m$ . We are seeking “*best possible*” approximations to a few quadratic eigenvalues of  $\mathbf{Q}(\lambda)$  using  $\mathcal{Y}$ . Let  $Y \in \mathbb{C}^{n \times m}$  be a basis matrix of  $\mathcal{Y}$ .

According to (3.7a) which says (upon substituting  $i = n - j + 1$ )

$$\lambda_{n-j+1}^+ = \max_{\substack{\mathcal{X} \subseteq \mathbb{C}^n \\ \dim \mathcal{X} = j}} \min_{\substack{x \in \mathcal{X} \\ x \neq 0}} \rho_+(x), \quad (3.7a')$$

it is natural to approximate  $\lambda_{n-j+1}^+$ , given  $\mathcal{Y} \subset \mathbb{C}^n$ , by

$$\mu_{m-j+1}^+ := \max_{\substack{\mathcal{X} \subseteq \mathcal{Y} \\ \dim \mathcal{X} = j}} \min_{\substack{x \in \mathcal{X} \\ x \neq 0}} \rho_+(x), \quad (5.1)$$

via replacing  $\mathcal{X} \subseteq \mathbb{C}^n$  in (3.7a') by  $\mathcal{X} \subseteq \mathcal{Y}$ . Any  $x \in \mathcal{X} \subseteq \mathcal{Y}$  can be written as  $x = Yy$  for some  $y \in \mathbb{C}^m$ , and thus

$$\rho_+(x) = \rho_+(Yy) = \frac{-(y^H Y^H B Y y) + [(y^H Y^H B Y y)^2 - 4(y^H Y^H A Y y)(y^H Y^H C Y y)]^{1/2}}{2(y^H Y^H A Y y)}.$$

Combined with (3.7a') and this expression for  $\rho_+(x)$ , (5.1) implies that  $\mu_1^+, \dots, \mu_m^+$  are the  $m$  pos-type quadratic eigenvalues of  $Y^H \mathbf{Q}(\lambda) Y$ . What this means is that  $\mu_j^+$  for  $1 \leq j \leq m$  provide the best approximations to the  $m$  largest  $\lambda_j^+$ , given  $\mathcal{Y}$ , in the sense of (3.7a). Of course, some approximations  $\mu_j^+ \approx \lambda_{n-m+j}^+$  are more accurate than others.

Similarly, given  $\mathcal{Y}$ ,  $\mu_j^+$  for  $1 \leq j \leq m$  provide the best approximations to the  $m$  smallest  $\lambda_j^+$  in the sense of (3.7b).

Let  $\mu_1^-, \dots, \mu_m^-$  are the  $m$  neg-type quadratic eigenvalues of  $Y^H \mathbf{Q}(\lambda) Y$ . The same argument shows, given  $\mathcal{Y}$ ,  $\mu_j^-$  for  $1 \leq j \leq m$  provide the best approximations to the  $m$  largest  $\lambda_j^-$  in the sense of (3.7c), and the best approximations to the  $m$  smallest  $\lambda_j^-$  in the sense of (3.7d).

---

**Algorithm 5.1** Rayleigh-Ritz procedure

---

Given  $Y \in \mathbb{C}^{n \times m}$  which is a basis matrix of  $\mathcal{Y} \subset \mathbb{C}^n$ , this algorithm returns approximations to  $k$  extreme quadratic eigenpairs (of pos- or neg-type) of  $\mathbf{Q}(\lambda)$ .

---

1: solve the QEP for  $Y^H \mathbf{Q}(\lambda) Y$  to get its quadratic eigenvalues  $\mu_j^\pm$  and associated quadratic eigenvectors  $y_j^\pm$ .

2: **return**

- $(\mu_i^\pm, Y y_i^\pm)$  for  $1 \leq i \leq k$  as approximations to  $(\lambda_i^\pm, u_i^\pm)$  for  $1 \leq i \leq k$ , or
- $(\mu_i^\pm, Y y_i^\pm)$  for  $m - k + 1 \leq i \leq m$  as approximations to  $(\lambda_i^\pm, u_i^\pm)$  for  $n - k + 1 \leq i \leq n$ ,

depending on what kind of extreme quadratic eigenpairs are desired.

---

In summary, we have justified that the quadratic eigenvalues of  $Y^H \mathbf{Q}(\lambda) Y$  yield the best approximations to some of the largest or smallest pos- or neg-type quadratic eigenvalues of  $\mathbf{Q}(\lambda)$  in certain respective senses. This statement could sound confusing: how could the same set of values be the best approximations to some of both the largest and smallest eigenvalues at the same time? But we point out this is not what the statement is saying. The key to understand the subtlety is not to forget that they provide the best approximations under the mentioned senses, and being the best approximations (under a particular sense) does not necessarily imply that the approximates are good, just that they are the best (under that particular sense). In practice,  $\mathcal{Y}$  is built to approximate either the largest or smallest eigenvalues well, but unlikely both.

Theorems 3.3, 3.4, and 3.5, generalizing Amir-Moéz's min-max principles and the Wielandt-Lidskii min-max principles, can also be used to justify that the quadratic eigenvalues of  $Y^H \mathbf{Q}(\lambda) Y$  are candidates for best approximating the largest or smallest pos- or neg-type quadratic eigenvalues of  $\mathbf{Q}(\lambda)$ , too. For example, according to (3.13a) with any pre-chosen  $\Phi$ , we should seek best approximations to  $\lambda_i^+$  for  $1 \leq i \leq k$  by

$$\text{minimizing } \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) \text{ subject to } \mathcal{R}(X) \subseteq \mathcal{Y} \text{ and } \text{rank}(X) = k. \quad (3.13a')$$

Noticing that any  $X \in \mathbb{C}^{n \times k}$  satisfying  $\mathcal{R}(X) \subseteq \mathcal{Y}$  and  $\text{rank}(X) = k$  can be written as  $X = Y \hat{X}$  for some  $\hat{X} \in \mathbb{C}^{m \times k}$  with  $\text{rank}(\hat{X}) = k$ , we see that  $\lambda_{j,X}^+$  are pos-type quadratic eigenvalues of  $[Y \hat{X}]^H \mathbf{Q}(\lambda) [Y \hat{X}] = \hat{X}^H Y^H \mathbf{Q}(\lambda) Y \hat{X}$ . Varying  $X$  subject to  $\mathcal{R}(X) \subseteq \mathcal{Y}$  and  $\text{rank}(X) = k$  is transferred to varying  $\hat{X} \in \mathbb{C}^{m \times k}$  subject to  $\text{rank}(\hat{X}) = k$ . Consequently,

$$\min_X \Phi(\lambda_{1,X}^+, \dots, \lambda_{k,X}^+) = \min_{\hat{X}} \Phi(\mu_{1,\hat{X}}^+, \dots, \mu_{k,\hat{X}}^+), \quad (5.2)$$

where  $\mu_{j,\hat{X}}^+$  are pos-type quadratic eigenvalues of  $\hat{X}^H Y^H \mathbf{Q}(\lambda) Y \hat{X}$ . Apply Theorem 3.3 to see the right-hand side of (5.2) is  $\Phi(\mu_1^+, \dots, \mu_k^+)$ , indicating  $\mu_j^+$  for  $1 \leq j \leq k$  provide the best approximations to the  $k$  smallest  $\lambda_j^+$ , as expected.

The same statement can be made about  $\mu_j^+$  as approximations to the largest  $\lambda_j^+$ ,  $\mu_j^-$  as approximations to the smallest  $\lambda_j^-$  or as approximations to the largest  $\lambda_j^-$ , using other min-max principles in Theorems 3.3, 3.4, and 3.5.

In summary, our discussion so far lead to a Rayleigh-Ritz type procedure detailed in Algorithm 5.1 to compute the best approximations to the desired quadratic eigenpairs of  $Q(\lambda)$ , given a pre-built subspace  $\mathcal{Y}$ .

## 6 The steepest descent/ascent method

A common approach to solve a quadratic eigenvalue problem in general, as well as any polynomial eigenvalue problem, is through *linearization* which converts the problem into a linear generalized eigenvalue problem of a matrix pencil [25, 42, 41]. The latter can be either solved by some iterative methods for a large scale problem or by the QZ algorithm [2, 44] for a problem of small to modest size ( $n$  up to around a few thousands for example). This approach is usually adopted for QEP without much structure to exploit. For HQEP, however, it is a different story – there is much to exploit. Most recent development includes the solvent approach [10, 21, 24, 61] for certain kinds of QEPs among which is HQEP [20]. Numerical evidence indicates that this solvent approach is rather efficient for QEP of small to modest sizes.

In this paper, we focus on optimization approaches based on various min-max principles previously established and the new ones established here. They are iterative methods and intended for solving large scale HQEP.

The equations in (3.8):

$$\lambda_1^+ = \min_{x \neq 0} \rho_+(x), \quad \lambda_n^+ = \max_{x \neq 0} \rho_+(x), \quad (3.8a)$$

$$\lambda_1^- = \min_{x \neq 0} \rho_-(x), \quad \lambda_n^- = \max_{x \neq 0} \rho_-(x). \quad (3.8b)$$

naturally suggest using some optimization techniques, including the steepest descent/ascent or CG-type method, to compute the first or last quadratic eigenpair  $(\lambda_j^\pm, u_j^\pm)$  as in the case of the standard Hermitian eigenvalue problem [3, 14]. Block variations can also be devised to simultaneously compute the first or last few quadratic eigenpairs  $(\lambda_j^\pm, u_j^\pm)$  again as in the case of the standard Hermitian eigenvalue problem [3, 40].

### 6.1 Gradients

To apply any of optimization techniques, we need to compute the gradients of  $\rho_\pm(x)$ . To this end, we use  $\rho(x)$  for either  $\rho_+(x)$  or  $\rho_-(x)$ . As  $x$  is perturbed to  $x + p$ , where  $p$  is assumed small in magnitude,  $\rho(x + p)$  is changed to  $\rho(x + p) = \rho(x) + \eta$ , where the magnitude  $\eta$  is comparable to  $\|p\|$ . We have by (3.1)

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

which gives, upon noticing  $f(\rho(x), x) = 0$ , that

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

and thus

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

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

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

or equivalently

$$\nabla \rho_{\pm}(x) = \mp \frac{2\mathbf{Q}(\rho_{\pm}(x))x}{\varsigma(x)}, \quad (6.1)$$

where we have used (3.5).

It is important to notice that the gradient  $\nabla \rho_{\pm}(x)$  is parallel to the residual vector

$$r_{\pm}(x) := [\rho_{\pm}(x)^2 A + \rho_{\pm}(x)B + C]x = \mathbf{Q}(\rho_{\pm}(x))x \quad (6.2)$$

whose normalized norm is commonly used to determine if the approximate eigenpair  $(\rho_{\pm}(x), x)$  meets a pre-set tolerance `rtol`:

$$\frac{\|r_{\pm}(x)\|}{|\rho_{\pm}(x)|^2 \|Ax\| + |\rho_{\pm}(x)| \|Bx\| + \|Cx\|} < \text{rtol}. \quad (6.3)$$

If (6.3) holds for  $(\rho_+(x), x)$ , then it is accepted as a converged pos-type quadratic eigenpairs, and similarly for  $(\rho_-(x), x)$ . Here which vector norm  $\|\cdot\|$  to use is usually inconsequential, but for the sake of convenience. More conservatively,  $\|Ax\|$  in the denominator should be replaced by  $\|A\| \|x\|$ , and likewise for  $\|Bx\|$  and  $\|Cx\|$  there. For large sparse matrices, the use of  $\|Ax\|$ ,  $\|Bx\|$ , and  $\|Cx\|$  is more economical because of their availability.

Beside being easily implementable, the use of (6.3) can also be rationalized by the existing backward error analysis of approximate eigenpairs for polynomial eigenvalue problems [25, 36, 62].

## 6.2 The steepest descent/ascent method

Now the steepest descent/ascent method for computing one of  $\lambda_{\ell}^{\pm}$  for  $\ell \in \{1, n\}$  can be readily given. For this purpose, we fix two parameters “typ” and  $\ell$  with varying ranges as

$$\text{typ} \in \{+, -\}, \quad \ell \in \{1, n\} \quad (6.4)$$

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

$$t_{\text{opt}} = \underset{t \in \mathbb{C}}{\text{argopt}} \rho_{\text{typ}}(x + tp), \quad (6.5)$$

where  $x$  is the current approximation to  $u_{\ell}^{\text{typ}}$  (thus no reason to let  $x = 0$ ),  $p$  is the search direction, and

$$\text{argopt} = \begin{cases} \text{argmin}, & \text{for } \ell = 1, \\ \text{argmax}, & \text{for } \ell = n. \end{cases} \quad (6.6)$$

The next approximate quadratic eigenvector is

$$y = \begin{cases} x + t_{\text{opt}} p, & \text{if } t_{\text{opt}} \text{ is finite,} \\ p, & \text{otherwise.} \end{cases} \quad (6.7)$$

But the line-search problem (6.5) doesn’t seem to be solvable straightforwardly by simple calculus as for the standard symmetric eigenvalue problem (see, e.g., [3, 14, 40, 70]), given the (complicated) expressions for  $\rho_{\text{typ}}$  in (3.2). Fortunately, the theory we developed in

---

**Algorithm 6.1** Steepest descent/ascent method
 

---

Given an initial approximation  $\mathbf{x}_0$  to  $u_\ell^{\text{typ}}$ , and a relative tolerance  $\text{rtol}$ , the algorithm computes an approximate pair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$  with the prescribed  $\text{rtol}$ .

---

```

1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|$ ,  $\boldsymbol{\rho}_0 = \rho_{\text{typ}}(\mathbf{x}_0)$ ,  $\mathbf{r}_0 = r_{\text{typ}}(\mathbf{x}_0)$ ;
2: for  $i = 0, 1, \dots$  do
3:   if  $\|\mathbf{r}_i\| / (|\boldsymbol{\rho}_i|^2 \|A\mathbf{x}_i\| + |\boldsymbol{\rho}_i| \|B\mathbf{x}_i\| + \|C\mathbf{x}_i\|) \leq \text{rtol}$  then
4:     BREAK;
5:   else
6:     solve QEP for  $Y_i^H \mathbf{Q}(\lambda) Y_i$ , where  $Y_i = [\mathbf{x}_i, \mathbf{r}_i]$  to get its quadratic eigenvalues  $\mu_j^\pm$ 
       as in (6.8) and corresponding quadratic eigenvectors  $y_j^\pm$ ;
7:     select the next approximate quadratic eigenpair  $(\mu, y) = (\mu_j^{\text{typ}}, Y_i y_j^{\text{typ}})$  according
       to the table (6.9);
8:      $\mathbf{x}_{i+1} = y / \|y\|$ ,  $\boldsymbol{\rho}_{i+1} = \mu$ ,  $\mathbf{r}_{i+1} = r_{\text{typ}}(\mathbf{x}_{i+1})$ ;
9:   end if
10: end for
11: return  $(\boldsymbol{\rho}_i, \mathbf{x}_i)$  as an approximate quadratic eigenpair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$ .

```

---

section 5 points us another way to look at it and thus solve it with ease. In fact, the problem is equivalent to find the best possible approximation within the subspace  $\mathcal{Y} = \mathcal{R}([x, p])$ . Suppose  $x$  and  $p$  are linearly independent<sup>9</sup> and let  $Y = [x, p]$ . Solve the 2-by-2 HQEP for  $Y^H \mathbf{Q}(\lambda) Y$  to get its quadratic eigenvalues

$$\mu_1^- \leq \mu_2^- < \mu_1^+ \leq \mu_2^+ \quad (6.8)$$

and corresponding quadratic eigenvectors  $y_j^\pm \in \mathbb{C}^2$ . We then have the following table for selecting the next approximate quadratic eigenpair, according to the parameter pair  $(\text{typ}, \ell)$ .

$(\text{typ}, \ell)$	current approx.	next approx.
(+, 1)	$(\rho_+(x), x)$	$(\mu_1^+, Y y_1^+)$
(+, $n$ )	$(\rho_+(x), x)$	$(\mu_2^+, Y y_2^+)$
(-, 1)	$(\rho_-(x), x)$	$(\mu_1^-, Y y_1^-)$
(-, $n$ )	$(\rho_-(x), x)$	$(\mu_2^-, Y y_2^-)$

(6.9)

In light of this alternative way to solve (6.5), the resulting steepest descent/ascent method is summarized in Algorithm 6.1.

**Lemma 6.1.** For (6.5) – (6.7),  $p^H r_{\text{typ}}(y) = 0$ .

*Proof.* If  $x$  and  $p$  are linearly dependent (the trivial case  $p = 0$  included), then  $p = \alpha x$  and  $y = \beta x$  for some scalars  $\alpha$  and  $\beta$ . Thus  $\rho_{\text{typ}}(y) = \rho_{\text{typ}}(x)$ ,  $r_{\text{typ}}(y) = \beta r_{\text{typ}}(x)$ , and  $p^H r_{\text{typ}}(y) = \alpha \beta x^H r_{\text{typ}}(x) = 0$  by the definition of  $\rho_{\text{typ}}(x)$ .

Suppose  $x$  and  $p$  are linearly independent. If  $|t_{\text{opt}}| = \infty$ , then  $y = p$ . Thus  $p^H r_{\text{typ}}(y) = y^H r_{\text{typ}}(y) = 0$ . Consider the case that  $t_{\text{opt}}$  is finite. Let  $t = t_{\text{opt}} + s$ . For tiny  $s$ , we have

$$\rho(y + sp) = \rho(y) - \frac{2\text{RE} \left( s[\rho(y)^2 Ay + \rho(y)By + Cy]^H p \right)}{2\rho(y) y^H Ay + y^H By} + \mathcal{O}(s^2),$$

---

<sup>9</sup>Otherwise, no improvement is expected by optimizing  $\rho_{\text{typ}}(x + tp)$  because then  $\rho_{\text{typ}}(x + tp) \equiv \rho_{\text{typ}}(x)$  for all scalar  $t$ .



where we drop the subscript “typ” to  $\rho_{\text{typ}}(\cdot)$  for convenience. Since  $\min_s \rho(y + sp)$  over  $s \in \mathbb{C}$  is attained at  $s = 0$ , it must hold that  $[\rho(y)^2 Ay + \rho(y)By + Cy]^H p = 0$ , as was to be shown.  $\square$

### 6.3 The extended steepest descent/ascent method

In Algorithm 6.1, the search space is spanned by

$$\mathbf{x}_i, \mathbf{r}_i = \mathbf{Q}(\boldsymbol{\rho}_i)\mathbf{x}_i.$$

Thus it is the second order Krylov subspace  $\mathcal{K}_2(\mathbf{Q}(\boldsymbol{\rho}_i), \mathbf{x}_i)$  of  $\mathbf{Q}(\boldsymbol{\rho}_i)$  on  $\mathbf{x}_i$ . Inspired by the inverse free Krylov subspace method [18] which seeks to improve the steepest descent method for the Hermitian generalized eigenvalue problem by extending the search space to a Krylov subspace, we may improve Algorithm 6.1 in the same way, i.e., using a high order Krylov subspace

$$\mathcal{K}_m(\mathbf{Q}(\boldsymbol{\rho}_i), \mathbf{x}_i) = \text{span}\{\mathbf{x}_i, \mathbf{Q}(\boldsymbol{\rho}_i)\mathbf{x}_i, \dots, [\mathbf{Q}(\boldsymbol{\rho}_i)]^{m-1}\mathbf{x}_i\} \quad (6.10)$$

as the search space. Let  $Y_i$  be a basis matrix of this Krylov subspace. We then solve<sup>10</sup> the  $m$ -by- $m$  HQEP for  $Y_i^H \mathbf{Q}(\lambda) Y_i$  to get its quadratic eigenvalues

$$\mu_1^- \leq \dots \leq \mu_m^- < \mu_1^+ \leq \dots \leq \mu_m^+ \quad (6.11)$$

and corresponding quadratic eigenvectors  $y_j^\pm$ . We then have the following table for selecting the next approximate quadratic eigenpair, according to the parameter pair  $(\text{typ}, \ell)$ .

$(\text{typ}, \ell)$	current approx.	next approx.
$(+, 1)$	$(\rho_+(\mathbf{x}_i), \mathbf{x}_i)$	$(\mu_1^+, Y_i y_1^+)$
$(+, n)$	$(\rho_+(\mathbf{x}_i), \mathbf{x}_i)$	$(\mu_m^+, Y_i y_m^+)$
$(-, 1)$	$(\rho_-(\mathbf{x}_i), \mathbf{x}_i)$	$(\mu_1^-, Y_i y_1^-)$
$(-, n)$	$(\rho_-(\mathbf{x}_i), \mathbf{x}_i)$	$(\mu_m^-, Y_i y_m^-)$

(6.12)

We summarize the resulting method, called the *Extended Steepest Descent/Ascent method*, into Algorithm 6.2.

When  $m = 2$ , Algorithm 6.2 reduces to the steepest descent/ascent method given in Algorithm 6.1.

### 6.4 Convergence analysis

While our convergent results are stated for all four possible  $(\text{typ}, \ell) \in \{(\pm, 1), (\pm, n)\}$ , our proofs will be presented mostly for one  $(\text{typ}, \ell)$

$$(\text{typ}, \ell) = (+, 1), \quad \text{and thus } \text{argopt} = \text{argmin} \text{ in (6.6)} \quad (6.13)$$

to save space. Proofs for other  $(\text{typ}, \ell)$  can be obtained with minor changes accordingly. For convenience, in our proofs we will drop the pos-type sign “+” in  $r_+(\cdot)$ ,  $\rho_+(\cdot)$ , and

<sup>10</sup>Often  $Y_i \in \mathbb{C}^{n \times m}$ , but there is a possibility that  $\dim \mathcal{K}_m(\mathbf{Q}(\boldsymbol{\rho}_i), \mathbf{x}_i) < m$ . When this occurs,  $Y_i$  will have fewer columns than  $m$ , and the rest of the development is still valid with minor changes. This is rare, especially in actual computations. For simplicity of presentation, we will assume that  $Y_i$  has  $m$  columns.

---

**Algorithm 6.2** Extended steepest descent/ascent method
 

---

Given an initial approximation  $\mathbf{x}_0$  to  $u_\ell^{\text{typ}}$ , and a relative tolerance  $\text{rtol}$ , and the search space dimension  $m$ , the algorithm computes an approximate pair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$  with the prescribed  $\text{rtol}$ .

---

```

1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|$ ,  $\boldsymbol{\rho}_0 = \rho_{\text{typ}}(\mathbf{x}_0)$ ,  $\mathbf{r}_0 = r_{\text{typ}}(\mathbf{x}_0)$ ;
2: for  $i = 0, 1, \dots$  do
3:   if  $\|\mathbf{r}_i\| / (|\boldsymbol{\rho}_i|^2 \|A\mathbf{x}_i\| + |\boldsymbol{\rho}_i| \|B\mathbf{x}_i\| + \|C\mathbf{x}_i\|) \leq \text{rtol}$  then
4:     BREAK;
5:   else
6:     compute a basis matrix  $Y_i$  for the Krylov subspace  $\mathcal{K}_m(\mathbf{Q}(\boldsymbol{\rho}_i), \mathbf{x}_i)$  in (6.10);
7:     solve QEP for  $Y_i^H \mathbf{Q}(\lambda) Y_i$  to get its quadratic eigenvalues  $\mu_j^\pm$  as in (6.11) and
       corresponding quadratic 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 in (6.12);
9:      $\mathbf{x}_{i+1} = y / \|y\|$ ,  $\boldsymbol{\rho}_{i+1} = \mu$ ,  $\mathbf{r}_{i+1} = r_{\text{typ}}(\mathbf{x}_{i+1})$ ;
10:  end if
11: end for
12: return  $(\boldsymbol{\rho}_i, \mathbf{x}_i)$  as an approximate quadratic eigenpair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$ .

```

---

$u_j^\pm$  with an understanding that they are all for the pos-type, even though occasionally, the sign is still written out at critical places.

By Theorem 2.5,  $\mathbf{Q}(\lambda)$  has  $n$  linearly independent pos-type quadratic eigenvectors  $u_j^+$  for  $1 \leq j \leq n$  and  $n$  linearly independent neg-type quadratic eigenvectors  $u_j^-$  for  $1 \leq j \leq n$ . Define for each (pos/neg-type) quadratic eigenvalue  $\mu$  its corresponding quadratic eigenspace

$$\mathcal{U}_\mu = \{x \in \mathbb{C}^n \mid \mathbf{Q}(\mu)x = 0\} = \bigoplus_{\lambda_i^{\text{typ}} = \mu} \text{span}\{u_i^{\text{typ}}\}.$$

We'll use the angle  $\theta(\mathbf{x}_i, \mathcal{U}_\mu)$  from  $\mathbf{x}_i$  to an eigenspace  $\mathcal{U}_\mu$ :

$$\cos \theta(\mathbf{x}_i, \mathcal{U}_\mu) := \min_{0 \neq u \in \mathcal{U}_\mu} \frac{|u^H \mathbf{x}_i|}{\|\mathbf{x}_i\|_2 \|u\|_2}$$

to measure the convergence of  $\mathbf{x}_i$  towards  $\mathcal{U}_\mu$ . Note  $0 \leq \theta(\mathbf{x}_i, \mathcal{U}_\mu) \leq \pi/2$ .

For the sake of our convergence analysis, it is convenient for us to execute Algorithms 6.1 and 6.2 without their Lines 3 and 4 so that  $\mathbf{x}_i$ ,  $\mathbf{r}_i$ , and  $\boldsymbol{\rho}_i$  are defined for all  $i \geq 0$ . But without the two lines, we need to be clear about the case when  $\mathbf{r}_i = 0$  for some  $i$ . When it occurs,  $\mathcal{K}_m(\mathbf{Q}(\boldsymbol{\rho}_i), \mathbf{x}_i) = \text{span}\{\mathbf{x}_i\}$  for any  $m \geq 2$ . For Algorithm 6.2, all subsequent  $\mathbf{x}_j$ ,  $\boldsymbol{\rho}_j$ , and  $\mathbf{r}_j$  for  $j > i$  are well-defined. In fact, we will have

$$\boldsymbol{\rho}_i = \boldsymbol{\rho}_{i+1} = \dots, \quad \mathbf{x}_i = \mathbf{x}_{i+1} = \dots, \quad \mathbf{r}_i = \mathbf{r}_{i+1} = \dots = 0. \quad (6.14)$$

But for Algorithm 6.1, all we have to do is to modify its Line 6 to “ $Y_i = \mathbf{x}_i$  if  $\mathbf{r}_i = 0$ ” and then  $\mathbf{x}_j$ ,  $\boldsymbol{\rho}_j$ , and  $\mathbf{r}_j$  for  $j > i$  are again well-defined and they again satisfy (6.14).

**Theorem 6.1.** *Let the sequences  $\{\boldsymbol{\rho}_i\}$ ,  $\{\mathbf{r}_i\}$ ,  $\{\mathbf{x}_i\}$  be produced by Algorithm 6.1/6.2.*

1. Only one of the following two mutually exclusive situations can occur:
  - (a) For some  $i$ , (6.14) holds, and  $(\boldsymbol{\rho}_i, \mathbf{x}_i)$  is a quadratic eigenpair of  $\mathbf{Q}(\lambda)$ .
  - (b)  $\boldsymbol{\rho}_i$  is strictly monotonically decreasing for  $(\text{typ}, \ell) \in \{(\pm, 1)\}$  or strictly monotonically increasing for  $(\text{typ}, \ell) \in \{(\pm, n)\}$ ,  $\mathbf{r}_i \neq 0$  for all  $i$ , and no two  $\mathbf{x}_i$  are linearly dependent.
2.  $\mathbf{x}_i^H \mathbf{r}_i = 0$ ,  $\mathbf{r}_i^H \mathbf{r}_{i+1} = 0$ ,  $\mathbf{x}_i^H \mathbf{r}_{i+1} = 0$  for Algorithm 6.1;
3.  $\mathbf{x}_i^H \mathbf{r}_i = 0$ ,  $Y_i^H \mathbf{r}_{i+1} = 0$  for Algorithm 6.2;
4. In the case of 1(b),
  - (a)  $\boldsymbol{\rho}_i \rightarrow \hat{\rho} \in [\lambda_1^{\text{typ}}, \lambda_n^{\text{typ}}]$  as  $i \rightarrow \infty$ ,
  - (b)  $\mathbf{r}_i \neq 0$  for all  $i$  but  $\mathbf{r}_i \rightarrow 0$  as  $i \rightarrow \infty$ ,
  - (c)  $\hat{\rho}$  is a quadratic eigenvalue of  $\mathbf{Q}(\lambda)$ , and any limit point  $\hat{\mathbf{x}}$  of  $\{\mathbf{x}_i\}$  is a corresponding quadratic eigenvector, i.e.,  $\mathbf{Q}(\hat{\rho})\hat{\mathbf{x}} = 0$ ,
  - (d)  $\theta(\mathbf{x}_i, \mathcal{U}_{\hat{\rho}}) \rightarrow 0$  as  $i \rightarrow \infty$ .

*Proof.* As we remarked at the beginning of this subsection, we will prove the claims only for  $(\text{typ}, \ell) = (+, 1)$ .

There are only two possibilities: either  $\mathbf{r}_i = 0$  for some  $i$  or  $\mathbf{r}_i \neq 0$  for all  $i$ . If  $\mathbf{r}_i = 0$  for some  $i$ , then  $\boldsymbol{\rho}_i = \boldsymbol{\rho}_{i+1}$  and  $\mathbf{x}_i = \mathbf{x}_{i+1}$  because  $\rho(\mathbf{x}_i + t\mathbf{r}_i) \equiv \rho(\mathbf{x}_i)$ . Consequently  $\mathbf{r}_{i+1} = 0$ , and the equations in (6.14) hold. Consider now  $\mathbf{r}_i \neq 0$  for all  $i$ . Note that  $\mathbf{r}_i \neq 0$  implies  $\nabla \boldsymbol{\rho}_i \neq 0$ , and so  $\rho(\mathbf{x}_i - s\nabla \boldsymbol{\rho}_i) < \rho(\mathbf{x}_i)$  for some  $s$  with sufficiently tiny  $|s|$ . This in turn implies  $\rho(\mathbf{x}_i + t\mathbf{r}_i) < \rho(\mathbf{x}_i)$  for some  $t$  with sufficiently tiny  $|t|$  and thus

$$\boldsymbol{\rho}_{i+1} = \inf_t \rho(\mathbf{x}_i + t\mathbf{r}_i) < \rho(\mathbf{x}_i).$$

Therefore  $\boldsymbol{\rho}_i$  is strictly monotonically decreasing. No two  $\mathbf{x}_i$  are linear dependent because linear dependent  $\mathbf{x}_i$  and  $\mathbf{x}_j$  produce  $\boldsymbol{\rho}_i = \boldsymbol{\rho}_j$ . This proves item 1.

For item 2,  $\mathbf{x}_i^H \mathbf{r}_i = \mathbf{x}_i^H \mathbf{Q}(\boldsymbol{\rho}_i) \mathbf{x}_i = 0$ . Since  $\rho(\mathbf{x}_{i+1}) = \min_t \rho(\mathbf{x}_i + t\mathbf{r}_i)$ , by Lemma 6.1,  $\mathbf{r}_i^H \mathbf{r}_{i+1} = 0$ . We now prove  $\mathbf{x}_i^H \mathbf{r}_{i+1} = 0$ . If  $\mathbf{r}_i = 0$ , then all  $\mathbf{r}_j = 0$  for  $j > i$  – no proof is necessary. Consider  $\mathbf{r}_i \neq 0$ . Then  $\boldsymbol{\rho}_{i+1} < \boldsymbol{\rho}_i$ . Note  $\mathbf{x}_{i+1}$  is a linear combination of  $\mathbf{x}_i$  and  $\mathbf{r}_i$ ; so we write  $\mathbf{x}_{i+1} = \alpha_i \mathbf{x}_i + \beta_i \mathbf{r}_i$  for some scalar  $\alpha_i$  and  $\beta_i$ . We know  $\beta_i \neq 0$ ; otherwise  $\mathbf{x}_{i+1} = \alpha_i \mathbf{x}_i$  to yield  $\boldsymbol{\rho}_{i+1} = \boldsymbol{\rho}_i$  which contradicts  $\boldsymbol{\rho}_{i+1} < \boldsymbol{\rho}_i$ . Therefore

$$\boldsymbol{\rho}_{i+1} = \rho(\mathbf{r}_i + (\alpha_i/\beta_i)\mathbf{x}_i) = \inf_t \rho(\mathbf{r}_i + t\mathbf{x}_i).$$

Apply Lemma 6.1 with  $x = \mathbf{r}_i$  and  $p = \mathbf{x}_i$  to get  $\mathbf{x}_i^H \mathbf{r}_{i+1} = 0$ .

For item 3, again  $\mathbf{x}_i^H \mathbf{r}_i = \mathbf{x}_i^H \mathbf{Q}(\boldsymbol{\rho}_i) \mathbf{x}_i = 0$ . Let  $\mathbf{x}_{i+1} = Y_i y$ . Then for each column  $z$  of  $Y_i$ , we have

$$\boldsymbol{\rho}_{i+1} = \rho(Y_i y) = \inf_t \rho(Y_i y + tz).$$

Apply Lemma 6.1 with  $x = Y_i y$  and  $p = z$  to get  $z^H \mathbf{r}_{i+1} = 0$ . Since  $z$  is any column of  $Y_i$ , we conclude  $Y_i^H \mathbf{r}_{i+1} = 0$ .

Now for item 4(a), since  $\boldsymbol{\rho}_i$  is strictly monotonically decreasing and bounded from below since  $\boldsymbol{\rho}_i \geq \lambda_1^+$ , it is convergent and  $\boldsymbol{\rho}_i \rightarrow \hat{\rho} \in [\lambda_1^+, \lambda_n^+]$  because  $\boldsymbol{\rho}_i = \rho(\mathbf{x}_i) \in [\lambda_1^+, \lambda_n^+]$  for all  $i$  by Theorem 3.1.

For item 4(b), we have  $\|\mathbf{r}_i\| = \|(A\rho_i^2 + B\rho_i + C)\mathbf{x}_i\| \leq \|A\|(\lambda_n^+)^2 + \|B\|\lambda_n^+ + \|C\|$  since  $\|\mathbf{x}_i\| = 1$ ; so both  $\{\mathbf{r}_i\}$  and  $\{\mathbf{x}_i\}$  are bounded sequences. It suffices to show that any limit point of  $\{\mathbf{r}_i\}$  is the zero vector. Assume, to the contrary,  $\{\mathbf{r}_i\}$  has a nonzero limit point  $\hat{r}$ , i.e.,  $\mathbf{r}_{i_j} \rightarrow \hat{r}$ , where  $\{\mathbf{r}_{i_j}\}$  is a subsequence of  $\{\mathbf{r}_i\}$ . Since  $\{\mathbf{x}_{i_j}\}$  is bounded, it has a convergent subsequence. Without loss of generality, we may assume  $\mathbf{x}_{i_j}$  itself is convergent and  $\mathbf{x}_{i_j} \rightarrow \hat{x}$  as  $j \rightarrow \infty$ . We have  $\hat{r}^H \hat{x} = 0$  and  $\|\hat{x}\| = 1$  because  $\mathbf{r}_{i_j}^H \mathbf{x}_{i_j} = 0$  and  $\|\mathbf{x}_{i_j}\| = 1$ . Now consider the quadratic eigenvalue problem for

$$\mathbf{Q}_{i_j}(\lambda) := Y_{i_j}^H \mathbf{Q}(\lambda) Y_{i_j} = \begin{bmatrix} \mathbf{x}_{i_j}^H \mathbf{Q}(\lambda) \mathbf{x}_{i_j} & \mathbf{x}_{i_j}^H \mathbf{Q}(\lambda) \mathbf{r}_{i_j} \\ \mathbf{r}_{i_j}^H \mathbf{Q}(\lambda) \mathbf{x}_{i_j} & \mathbf{r}_{i_j}^H \mathbf{Q}(\lambda) \mathbf{r}_{i_j} \end{bmatrix}, \quad (6.15)$$

where  $Y_{i_j} = [\mathbf{x}_{i_j}, \mathbf{r}_{i_j}]$ . Since  $\mathbf{r}_{i_j}^H \mathbf{x}_{i_j} = 0$ ,  $\text{rank}(Y_{i_j}) = 2$ , and thus  $\mathbf{Q}_{i_j}(\lambda)$  is hyperbolic. Denote by  $\mu_{j;k}^\pm$  its quadratic eigenvalues. It can be seen that

$$\lambda_1^- \leq \mu_{j;1}^- \leq \mu_{j;2}^- \leq \lambda_n^- < \lambda_1^+ \leq \mu_{j;1}^+ \leq \mu_{j;2}^+ \leq \lambda_n^+. \quad (6.16)$$

Then<sup>11</sup>  $\lambda_1^+ \leq \boldsymbol{\rho}_{i_j+1} \leq \mu_{j;1}^+$ . Let

$$\hat{\mathbf{Q}}(\lambda) = \lim_{j \rightarrow \infty} \mathbf{Q}_{i_j}(\lambda)$$

whose quadratic eigenvalues are denoted by  $\hat{\mu}_i^\pm$ . By the continuity of the quadratic eigenvalues with respect to the entries of coefficient matrices of a quadratic polynomial with a nonsingular leading coefficient matrix, we know  $\mu_{j;i}^\pm \rightarrow \hat{\mu}_i^\pm$  as  $j \rightarrow \infty$ , and thus

$$\lambda_1^- \leq \hat{\mu}_1^- \leq \hat{\mu}_2^- \leq \lambda_n^- < \lambda_1^+ \leq \hat{\mu}_1^+ \leq \hat{\mu}_2^+ \leq \lambda_n^+. \quad (6.17)$$

Notice by (6.16) and (6.17)

$$\lambda_1^+ \leq \boldsymbol{\rho}_{i_j+1} \leq \mu_{j;1}^+ \Rightarrow \hat{\mu}_2^- < \lambda_1^+ \leq \hat{\rho} \leq \hat{\mu}_1^+. \quad (6.18)$$

On the other hand, by (6.16), we have

$$\hat{\mathbf{Q}}(\hat{\rho}) = \lim_{j \rightarrow \infty} \mathbf{Q}_{i_j}(\boldsymbol{\rho}_{i_j}) = \lim_{j \rightarrow \infty} \begin{bmatrix} 0 & \mathbf{r}_{i_j}^H \mathbf{r}_{i_j} \\ \mathbf{r}_{i_j}^H \mathbf{r}_{i_j} & \mathbf{r}_{i_j}^H \mathbf{Q}(\boldsymbol{\rho}_{i_j}) \mathbf{r}_{i_j} \end{bmatrix} = \begin{bmatrix} 0 & \hat{r}^H \hat{r} \\ \hat{r}^H \hat{r} & \hat{r}^H \mathbf{Q}(\hat{\rho}) \hat{r} \end{bmatrix}$$

which is indefinite because  $\hat{r}^H \hat{r} > 0$ . But by (6.18) and Theorem 2.1,  $\hat{\mathbf{Q}}(\hat{\rho}) \preceq 0$ , a contradiction. So  $\hat{r} = 0$ , as was to be shown.

For item 4(c), since  $\|\mathbf{x}_i\| = 1$ ,  $\{\mathbf{x}_i\}$  has at least one limit point. Let  $\hat{x}$  be any limit point of  $\mathbf{x}_i$ , i.e.,  $\mathbf{x}_{i_j} \rightarrow \hat{x}$ . Take limit at the both sides of  $\mathbf{Q}(\boldsymbol{\rho}_{i_j})\mathbf{x}_{i_j} = \mathbf{r}_{i_j}$  to get  $\mathbf{Q}(\hat{\rho})\hat{x} = 0$ , i.e.,  $(\hat{\rho}, \hat{x})$  is a quadratic eigenpair.

For item 4(d), write  $\theta_i = \theta(\mathbf{x}_i, \mathcal{U}_{\hat{\rho}})$  for convenience and write<sup>12</sup>  $\mathbf{x}_i = \hat{u}_i \cos \theta_i + \hat{v}_i \sin \theta_i$ , where  $\hat{u}_i \in \mathcal{U}_{\hat{\rho}}$ ,  $\hat{v}_i \in \mathcal{U}_{\hat{\rho}}^\perp$  (the orthogonal complement of  $\mathcal{U}_{\hat{\rho}}$ ), and  $\|\hat{u}_i\|_2 = \|\hat{v}_i\|_2 = 1$ . Then

$$\mathbf{r}_i = \mathbf{Q}(\boldsymbol{\rho}_i)\mathbf{x}_i = (\boldsymbol{\rho}_i - \hat{\rho}) [(\boldsymbol{\rho}_i + \hat{\rho})A + B] \hat{u}_i \cos \theta_i + \mathbf{Q}(\boldsymbol{\rho}_i)\hat{v}_i \sin \theta_i. \quad (6.19)$$

<sup>11</sup>For Algorithm 6.1,  $\boldsymbol{\rho}_{i_j+1} = \mu_{j;1}^+$ .

<sup>12</sup>Without loss of generality, we may assume  $\|\cdot\|_2$  is used in the algorithms.

We claim that  $\mathbf{Q}(\boldsymbol{\rho}_i)\hat{v}_i \sin \theta_i \rightarrow 0$ . To see this, we notice

$$\|(\boldsymbol{\rho}_i + \hat{\rho})A + B\|_2 \leq 2 \max\{|\lambda_1^+|, |\lambda_n^+|\} \|A\|_2 + \|B\|_2,$$

$\mathbf{r}_i \rightarrow 0$ , and  $\boldsymbol{\rho}_i - \hat{\rho} \rightarrow 0$ . Thus  $\mathbf{Q}(\boldsymbol{\rho}_i)\hat{v}_i \sin \theta_i \rightarrow 0$  by (6.19). The null space of  $\mathbf{Q}(\hat{\rho})$  is  $\mathcal{U}_{\hat{\rho}}$ . Since  $\mathbf{Q}(\hat{\rho})$  is Hermitian,

$$\|\mathbf{Q}(\hat{\rho})v\|_2 \geq \gamma \|v\|_2 \quad \text{for any } v \in \mathcal{U}_{\hat{\rho}}^\perp,$$

where  $\gamma = \min |\xi|$  taken over all nonzero  $\xi \in \text{eig}(\mathbf{Q}(\hat{\rho}))$ . Therefore  $\|\mathbf{Q}(\hat{\rho})\hat{v}_i\|_2 \geq \gamma$ . Because  $\boldsymbol{\rho}_i \rightarrow \hat{\rho}$ , for sufficiently large  $i$  we have  $\|\mathbf{Q}(\boldsymbol{\rho}_i)\hat{v}_i\|_2 \geq \gamma/2$  and thus

$$\|\mathbf{Q}(\boldsymbol{\rho}_i)\hat{v}_i \sin \theta_i\|_2 \geq (\gamma/2) \sin \theta_i,$$

implying  $\sin \theta_i \rightarrow 0$  which leads to  $\theta_i \rightarrow 0$  because  $0 \leq \theta_i \leq \pi/2$ .  $\square$

Theorem 6.1 ensures us the global convergence of Algorithm 6.1/6.2, but gives no indication as how fast the convergence may be. For that, we turn to our next theorem – Theorem 6.2 – which provides an asymptotic rate of the sequences  $\{\boldsymbol{\rho}_i\}$  generated by the algorithms. Both theorems are reminiscent of [18, Theorem 3.2] and [18, Theorem 3.4], respectively. But Theorem 6.2 about the rate of convergence is much more difficult to prove than [18, Theorem 3.4]. Because of that, we will devote the entire subsection 6.5 for its proof.

We introduce a few new notations: for any  $x \neq 0$ ,

$$a(x) = \frac{x^H A x}{x^H x}, \quad b(x) = \frac{x^H B x}{x^H x}, \quad c(x) = \frac{x^H C x}{x^H x}. \quad (6.20)$$

Also recall  $\mathbf{Q}_{\lambda_0}(\lambda) := \mathbf{Q}(\lambda + \lambda_0)$  in (4.5) for a given shift  $\lambda_0$ . Accordingly,

$$b_0(x) = \frac{x^H B_{\lambda_0} x}{x^H x} = \frac{x^H (2\lambda_0 A + B) x}{x^H x}, \quad c_0(x) = \frac{x^H C_{\lambda_0} x}{x^H x} = \frac{x^H \mathbf{Q}(\lambda_0) x}{x^H x}. \quad (6.21)$$

**Theorem 6.2.** *Suppose  $\lambda_1^{\text{typ}} \leq \boldsymbol{\rho}_0 < \lambda_2^{\text{typ}}$  if  $\ell = 1$  or  $\lambda_{n-1}^{\text{typ}} < \boldsymbol{\rho}_0 \leq \lambda_n^{\text{typ}}$  if  $\ell = n$ , and let the sequences  $\{\boldsymbol{\rho}_i\}, \{\mathbf{r}_i\}, \{\mathbf{x}_i\}$  be produced by Algorithm 6.2. Given a shift  $\lambda_0 \geq \lambda_n^+$ , define  $B_{\lambda_0}, C_{\lambda_0}$  by (4.5).*

1. *As  $i \rightarrow \infty$ ,  $\boldsymbol{\rho}_i$  monotonically converges to  $\hat{\rho} = \lambda_\ell^{\text{typ}}$ , and  $\mathbf{x}_i$  converges to  $u_\ell^{\text{typ}}$  in direction, i.e.,  $\theta(\mathbf{x}_i, u_\ell^{\text{typ}}) \rightarrow 0$ .*
2. *The eigenvalues<sup>13</sup>  $\omega_j$  of  $\mathbf{Q}(\boldsymbol{\rho}_i)$  can be ordered as*

$$\omega_1 > 0 > \omega_2 \geq \cdots \geq \omega_n \quad \text{if } (\text{typ}, \ell) \in \{(+, 1), (-, n)\}, \quad \text{or}, \quad (6.22a)$$

$$\omega_1 < 0 < \omega_2 \leq \cdots \leq \omega_n \quad \text{if } (\text{typ}, \ell) \in \{(+, n), (-, 1)\}. \quad (6.22b)$$

*Denote by  $v_1$  the eigenvector of  $\mathbf{Q}(\boldsymbol{\rho}_i)$  associated with its eigenvalue  $\omega_1$ . If  $\boldsymbol{\rho}_i$  is sufficiently close to  $\lambda_\ell^{\text{typ}}$ , then*

$$|\boldsymbol{\rho}_{i+1} - \lambda_\ell^{\text{typ}}| \leq \varepsilon_m^2 |\boldsymbol{\rho}_i - \lambda_\ell^{\text{typ}}| + (1 - \varepsilon_m^2) \varepsilon_m \eta(v_1) |\boldsymbol{\rho}_i - \lambda_\ell^{\text{typ}}|^{3/2} + O(|\boldsymbol{\rho}_i - \lambda_\ell^{\text{typ}}|^2), \quad (6.23)$$

<sup>13</sup>Their dependency upon  $i$  is suppressed for clarity.

where

$$\varepsilon_m = \min_{g \in \mathbb{P}_{m-1}, g(\omega_1) \neq 0} \max_{i \neq 1} \frac{|g(\omega_i)|}{|g(\omega_1)|}, \quad (6.24)$$

$$\tau_A = \frac{1}{|\omega_2|} \frac{\|A\|_2}{a(v_1)}, \quad \tau_B = \frac{1}{|\omega_2|} \frac{\|B_{\lambda_0}\|_2}{b_0(v_1)}, \quad \tau_C = \frac{1}{|\omega_2|} \frac{\|C_{\lambda_0}\|_2}{c_0(v_1)}, \quad (6.25)$$

$$\eta(v_1) = 3\tau_A^{1/2} + 2 \frac{(b_0(v_1))^2 \tau_B^{1/2} + 2a(v_1)c_0(v_1)(\tau_A^{1/2} + \tau_C^{1/2})}{\varsigma_0(v_1)^2}, \quad (6.26)$$

and  $\mathbb{P}_{m-1}$ , the set of polynomials of degree no higher than  $m-1$ .

3. Denote<sup>14</sup> by  $\gamma$  and  $\Gamma$  the smallest and largest positive eigenvalue of

$$\begin{cases} -\mathbf{Q}(\lambda_\ell^{\text{typ}}) & \text{for } (\text{typ}, \ell) \in \{(+, 1), (-, n)\}, \\ \mathbf{Q}(\lambda_\ell^{\text{typ}}) & \text{for } (\text{typ}, \ell) \in \{(+, n), (-, 1)\}. \end{cases}$$

If  $\rho_i$  is sufficiently close to  $\lambda_\ell^{\text{typ}}$ , then

$$|\rho_{i+1} - \lambda_\ell^{\text{typ}}| \leq \varepsilon^2 |\rho_i - \lambda_\ell^{\text{typ}}| + (1 - \varepsilon^2) \varepsilon \eta |\rho_i - \lambda_\ell^{\text{typ}}|^{3/2} + O(|\rho_i - \lambda_\ell^{\text{typ}}|^2), \quad (6.27)$$

where

$$\varepsilon = 2 \left[ \left( \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^{m-1} + \left( \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^{-(m-1)} \right]^{-1}, \quad \kappa = \frac{\Gamma}{\gamma}, \quad (6.28)$$

$$\eta = \sqrt{\frac{1}{|\gamma|}} \left[ 3 \sqrt{\frac{\|A\|_2}{a(u)}} + 2 \frac{b_0(u)^2}{\varsigma_0(u)^2} \sqrt{\frac{\|B_{\lambda_0}\|_2}{b_0(u)}} + 4 \frac{a(u)c_0(u)}{\varsigma_0(u)^2} \left( \sqrt{\frac{\|A\|_2}{a(u)}} + \sqrt{\frac{\|C_{\lambda_0}\|_2}{c_0(u)}} \right) \right] \quad (6.29)$$

$$\leq \sqrt{\frac{1}{|\gamma|}} \left[ 3 \sqrt{\frac{\|A\|_2}{a(u)}} + 2 \frac{\|B_{\lambda_0}\|_2^2 + 4\|A\|_2\|C_{\lambda_0}\|_2}{b(u)^2 - 4a(u)c(u)} \right], \quad (6.30)$$

and  $u = u_\ell^{\text{typ}}$  for short.

## 6.5 Proof of Theorem 6.2

We recall (3.5) to see

$$\begin{aligned} \varsigma(x) &:= [(x^H B x)^2 - 4(x^H A x)(x^H C x)]^{1/2} \\ &= \pm x^H [2\rho_\pm(x)A + B]x \\ &= \pm x^H \mathbf{Q}'(\rho_\pm(x))x, \end{aligned} \quad (6.31)$$

and  $\varsigma_0(x) = \varsigma(x)/\|x\|_2^2$ . For a perturbation  $E \in \mathbb{C}^{n \times n}$  which is assumed Hermitian, we define

$$\mathbf{Q}_E(\lambda) := \mathbf{Q}(\lambda) + E = \lambda^2 A + \lambda B + C + E. \quad (6.32)$$

<sup>14</sup> $\mathbf{Q}(\lambda_\ell^{\text{typ}})$  is singular and, by Theorem 2.1, negative semidefinite if  $(\text{typ}, \ell) \in \{(+, 1), (-, n)\}$  or positive semidefinite if  $(\text{typ}, \ell) \in \{(+, n), (-, 1)\}$ .

When  $\mathbf{Q}_E(\lambda)$  is also hyperbolic, the pos- and neg-type Rayleigh quotients, denoted by  $\rho_{E;\pm}$ , can be defined for  $\mathbf{Q}_E(\lambda)$ . Accordingly, we will define  $\varsigma_E$  and  $\varsigma_{E;0}$ , too. Specifically,

$$\rho_{E;\pm}(x) = \frac{-(x^H B x) \pm \{(x^H B x)^2 - 4(x^H A x)(x^H [C + E] x)\}^{1/2}}{2(x^H A x)}, \quad (6.33)$$

and

$$\begin{aligned} \varsigma_E(x) &:= \{(x^H B x)^2 - 4(x^H A x)(x^H [C + E] x)\}^{1/2} \\ &= \pm x^H [2\rho_{E;\pm}(x) A + B] x, \end{aligned} \quad (6.34a)$$

$$\varsigma_{E;0}(x) := \frac{\varsigma_E(x)}{\|x\|_2^2}. \quad (6.34b)$$

**Lemma 6.2.** *Suppose  $\mathbf{Q}_E(\lambda)$  in (6.32) is also hyperbolic.*

1. Let  $(\lambda_1^+, u_1^+)$  and  $(\mu_1^+, v_1^+)$  be the smallest quadratic eigenpair<sup>15</sup> with the pos-type of  $\mathbf{Q}(\lambda)$  and  $\mathbf{Q}_E(\lambda)$ , respectively. Then

$$\frac{\lambda_{\min}(E)}{\varsigma_0(u_1^+)} \leq \lambda_1^+ - \mu_1^+ \leq \frac{\lambda_{\max}(E)}{\varsigma_{E;0}(v_1^+)}. \quad (6.35)$$

2. Let  $(\lambda_n^+, u_n^+)$  and  $(\mu_n^+, v_n^+)$  be the largest quadratic eigenpair with the pos-type of  $\mathbf{Q}(\lambda)$  and  $\mathbf{Q}_E(\lambda)$ , respectively. Then

$$\frac{\lambda_{\min}(E)}{\varsigma_0(v_n^+)} \leq \lambda_n^+ - \mu_n^+ \leq \frac{\lambda_{\max}(E)}{\varsigma_{E;0}(u_n^+)}. \quad (6.36)$$

3. Let  $(\lambda_1^-, u_1^-)$  and  $(\mu_1^-, v_1^-)$  be the smallest quadratic eigenpair with the neg-type of  $\mathbf{Q}(\lambda)$  and  $\mathbf{Q}_E(\lambda)$ , respectively. Then

$$\frac{\lambda_{\min}(E)}{\varsigma_0(v_1^-)} \leq \mu_1^- - \lambda_1^- \leq \frac{\lambda_{\max}(E)}{\varsigma_{E;0}(u_1^-)}. \quad (6.37)$$

4. Let  $(\lambda_n^-, u_n^-)$  and  $(\mu_n^-, v_n^-)$  be the largest quadratic eigenpair with the neg-type of  $\mathbf{Q}(\lambda)$  and  $\mathbf{Q}_E(\lambda)$ , respectively. Then

$$\frac{\lambda_{\min}(E)}{\varsigma_0(u_n^-)} \leq \mu_n^- - \lambda_n^- \leq \frac{\lambda_{\max}(E)}{\varsigma_{E;0}(v_n^-)}. \quad (6.38)$$

*Proof.* As in the proof of Lemma 4.4, we have

$$\mu_1^+ = \min_x \rho_{E;+}(x) \leq \rho_{E;+}(u_1^+) \leq \rho_+(u_1^+) + \delta_{\text{ub}}^+(u_1^+) = \lambda_1^+ + \delta_{\text{ub}}^+(u_1^+)$$

which gives

$$\mu_1^+ - \lambda_1^+ \leq \delta_{\text{ub}}^+(u_1^+), \quad \lambda_1^+ - \mu_1^+ \leq \tilde{\delta}_{\text{ub}}^+(v_1^+), \quad (6.39)$$

---

<sup>15</sup>By the smallest (largest) pos/neg-type quadratic eigenpair, we mean the quadratic eigenvalue in question is the smallest (largest) of that given type. The same naming is used for the usual linear eigenpair, too.

where the second inequality is actually obtained from the first one there by switching the roles of  $\mathbf{Q}(\lambda)$  and  $\mathbf{Q}_E(\lambda)$ . Now use (4.42) in the proof of Theorem 4.1 for  $\Delta A = \Delta B = 0$  and  $\Delta C = E$  to get item 1.

Similarly, we have

$$\lambda_n^+ = \max_x \rho_+(x) \geq \rho_+(v_n^+) \geq \rho_{E;+}(v_n^+) - \delta_{\text{ub}}^+(v_n^+) = \mu_n^+ - \delta_{\text{ub}}^+(v_n^+)$$

which gives

$$\mu_n^+ - \lambda_n^+ \leq \delta_{\text{ub}}^+(v_n^+), \quad \lambda_n^+ - \mu_n^+ \leq \tilde{\delta}_{\text{ub}}^+(u_n^+), \quad (6.40)$$

where the second inequality is actually obtained from switching the roles of  $\mathbf{Q}(\lambda)$  and  $\mathbf{Q}_E(\lambda)$ . Now use (4.42) in the proof of Theorem 4.1 for  $\Delta A = \Delta B = 0$  and  $\Delta C = E$  to get item 2.

Items 3 and 4 are corollaries of items 2 and 1 applied to  $\mathbf{Q}(-\lambda)$  and  $\mathbf{Q}_E(-\lambda)$ .  $\square$

**Lemma 6.3.**  $\mathbf{Q}_E(\lambda)$  with  $E = -\sigma I$  is hyperbolic if

$$\sigma > -\frac{(\lambda_1^+ - \lambda_n^-)^2 \lambda_{\min}(A)}{4}. \quad (6.41)$$

*Proof.* For any vector  $x \neq 0$ , we have

$$\begin{aligned} (x^H B x)^2 - 4(x^H A x)(x^H [C - \sigma I] x) &= (x^H B x)^2 - 4(x^H A x)(x^H C x) + 4\sigma(x^H A x)(x^H x) \\ &= [\rho_+(x) - \rho_-(x)]^2 (x^H A x)^2 + 4\sigma(x^H A x)(x^H x) \\ &\geq (x^H A x)(x^H x) \left[ (\lambda_1^+ - \lambda_n^-)^2 \frac{x^H A x}{x^H x} + 4\sigma \right] \\ &\geq (x^H A x)(x^H x) [(\lambda_1^+ - \lambda_n^-)^2 \lambda_{\min}(A) + 4\sigma] \\ &> 0, \end{aligned}$$

where the last inequality holds because of (6.41).  $\square$

So  $\varsigma_E$  and  $\varsigma_{E;0}$  are well-defined for any  $E = -\sigma I$  satisfying (6.41). To emphasize such special  $E = -\sigma I$ , we introduce notations

$$\varsigma_\sigma(x) := \varsigma_E(v), \quad \varsigma_{\sigma;0}(v) := \varsigma_{E;0}(v) \quad \text{for } E = -\sigma I. \quad (6.42)$$

For  $\boldsymbol{\rho} \in (\lambda_1^{\text{typ}}, \lambda_n^{\text{typ}})$ , it follows from Theorem 2.1 that the largest eigenvalue, denoted by  $\omega_1$ , of  $\mathbf{Q}(\boldsymbol{\rho})$  is nonnegative, and thus this  $\sigma = \omega_1$  automatically satisfies (6.41). But the smallest eigenvalue, denoted also by  $\omega_1$ , of  $\mathbf{Q}(\boldsymbol{\rho})$  is non-positive and (6.41) may fail for  $\sigma = \omega_1$  unless  $|\omega_1|$  is sufficiently tiny.

**Lemma 6.4.** Given  $\lambda_1^{\text{typ}} \leq \boldsymbol{\rho} \leq \lambda_n^{\text{typ}}$ , let  $(\omega_1, v_1)$  be the largest eigenpair  $\mathbf{Q}(\boldsymbol{\rho})$  if  $(\text{typ}, \ell) \in \{(+, 1), (-, n)\}$  or the smallest eigenpair  $\mathbf{Q}(\boldsymbol{\rho})$  if  $(\text{typ}, \ell) \in \{(+, n), (-, 1)\}$ . If (6.41) holds with  $\sigma = \omega_1$ , then for the four different  $(\text{typ}, \ell)$

$$\frac{\varsigma_0(u_1^+)}{\varsigma_{\omega_1;0}(v_1)} (\boldsymbol{\rho} - \lambda_1^+) \leq \frac{\omega_1}{\varsigma_{\omega_1;0}(v_1)} \leq \boldsymbol{\rho} - \lambda_1^+ \quad \text{for } (\text{typ}, \ell) = (+, 1), \quad (6.43a)$$

$$\frac{\varsigma_{\omega_1;0}(u_n^+)}{\varsigma_0(v_1)} (\lambda_n^+ - \boldsymbol{\rho}) \leq \frac{-\omega_1}{\varsigma_0(v_1)} \leq \lambda_n^+ - \boldsymbol{\rho} \quad \text{for } (\text{typ}, \ell) = (+, n), \quad (6.43b)$$



$$\frac{\varsigma_{\omega_1;0}(u_1^-)}{\varsigma_0(v_1)}(\boldsymbol{\rho} - \lambda_1^-) \leq \frac{-\omega_1}{\varsigma_0(v_1)} \leq \boldsymbol{\rho} - \lambda_1^- \quad \text{for } (\text{typ}, \ell) = (-, 1), \quad (6.43c)$$

$$\frac{\varsigma_0(u_n^-)}{\varsigma_{\omega_1;0}(v_1)}(\lambda_n^- - \boldsymbol{\rho}) \leq \frac{\omega_1}{\varsigma_{\omega_1;0}(v_1)} \leq \lambda_n^- - \boldsymbol{\rho} \quad \text{for } (\text{typ}, \ell) = (-, n). \quad (6.43d)$$

Moreover, for  $\boldsymbol{\rho}$  sufficiently close to  $\lambda_\ell^{\text{typ}}$ ,

$$\frac{\omega_1}{\varsigma_{\omega_1;0}(v_1)} = \boldsymbol{\rho} - \lambda_1^+ + O([\boldsymbol{\rho} - \lambda_1^+]^2) \quad \text{for } (\text{typ}, \ell) = (+, 1), \quad (6.44a)$$

$$\frac{-\omega_1}{\varsigma_0(v_1)} = \lambda_n^+ - \boldsymbol{\rho} + O([\lambda_n^+ - \boldsymbol{\rho}]^2) \quad \text{for } (\text{typ}, \ell) = (+, n), \quad (6.44b)$$

$$\frac{-\omega_1}{\varsigma_0(v_1)} = \boldsymbol{\rho} - \lambda_1^- + O([\boldsymbol{\rho} - \lambda_1^-]^2) \quad \text{for } (\text{typ}, \ell) = (-, 1), \quad (6.44c)$$

$$\frac{\omega_1}{\varsigma_{\omega_1;0}(v_1)} = \lambda_n^- - \boldsymbol{\rho} + O([\lambda_n^- - \boldsymbol{\rho}]^2) \quad \text{for } (\text{typ}, \ell) = (-, n). \quad (6.44d)$$

*Proof.* Consider the case  $(\text{typ}, \ell) = (+, 1)$ . We have  $\omega_1 \geq 0$  and  $[\mathbf{Q}(\boldsymbol{\rho}) - \omega_1 I] v_1 = 0$ . Since  $\omega_1$  is the largest eigenvalue of  $\mathbf{Q}(\boldsymbol{\rho})$ ,  $\mathbf{Q}(\boldsymbol{\rho}) - \omega_1 I \preceq 0$ . Thus,  $(\boldsymbol{\rho}, v_1)$  is the smallest pos-type quadratic eigenpair of  $\mathbf{Q}_E(\lambda)$  with  $E = -\omega_1 I$ . By Lemma 6.2,

$$\frac{\omega_1}{\varsigma_{E;0}(v_1)} \leq \boldsymbol{\rho} - \lambda_1^+ \leq \frac{\omega_1}{\varsigma_0(u_1)}$$

which gives (6.43a). To prove (6.44a), we denote by  $\alpha(t)$  the largest eigenvalue of  $\mathbf{Q}(t)$  near  $t = \lambda_1^+$ . Then  $\alpha(\lambda_1^+) = 0$  and  $\alpha(\boldsymbol{\rho}) = \omega_1$ . Note that

$$\mathbf{Q}(\boldsymbol{\rho})v_1 = \omega_1 v_1 \quad \Rightarrow \quad v_1^H \mathbf{Q}(\boldsymbol{\rho})v_1 = \omega_1 v_1^H v_1 \quad \Rightarrow \quad v_1^H [\mathbf{Q}(\boldsymbol{\rho}) - \omega_1 I] v_1 = 0,$$

i.e.,  $\boldsymbol{\rho}$  is a Rayleigh quotient of  $\mathbf{Q}_E(\lambda)$  with  $E = -\omega_1 I$ . Therefore

$$\alpha'(\boldsymbol{\rho}) = \frac{v_1^H \mathbf{Q}'(\boldsymbol{\rho})v_1}{v_1^H v_1} = \frac{v_1^H \mathbf{Q}'_E(\boldsymbol{\rho})v_1}{v_1^H v_1} = \varsigma_{\omega_1;0}(v_1),$$

where the first equality is due to [56, p.183], and the third equality due to (6.31). Finally  $\alpha(\lambda_1^+) = \alpha(\boldsymbol{\rho}) + \varsigma_{\omega_1;0}(v_1)(\lambda_1^+ - \boldsymbol{\rho}) + O(|\lambda_1^+ - \boldsymbol{\rho}|^2)$  which leads to (6.44a).  $\square$

**Remark 6.1.** There is a different proof of Lemma 6.4, without using Lemma 6.2. For the case  $(\text{typ}, \ell) = (+, 1)$ ,  $(\boldsymbol{\rho}, v_1)$  is the smallest pos-type quadratic eigenpair of  $\mathbf{Q}_E(\lambda) =$

$\lambda^2 A + \lambda B + C - \omega_1 I$ . By direct calculations<sup>16</sup>,

$$\omega_1 = \omega_1 - \frac{u_1^H \mathbf{Q}(\boldsymbol{\rho}) u_1}{u_1^H u_1} + \varsigma_0(u_1)(\boldsymbol{\rho} - \lambda_1^+) + \frac{u_1^H A u_1}{u_1^H u_1} (\boldsymbol{\rho} - \lambda_1^+)^2, \quad (6.45a)$$

$$\omega_1 = \frac{v_1^H \mathbf{Q}(\lambda_1^+) v_1}{v_1^H v_1} + \varsigma_{\omega_1;0}(v_1)(\boldsymbol{\rho} - \lambda_1^+) - \frac{v_1^H A v_1}{v_1^H v_1} (\boldsymbol{\rho} - \lambda_1^+)^2. \quad (6.45b)$$

Along with  $\mathbf{Q}(\boldsymbol{\rho}) - \omega_1 I \preceq 0$ ,  $\mathbf{Q}(\lambda_1^+) \preceq 0$ , they yield

$$\frac{\omega_1}{\varsigma_{\omega_1;0}(v_1)} \leq \boldsymbol{\rho} - \lambda_1^+ \leq \frac{\omega_1}{\varsigma_0(u_1)}$$

and then

$$\frac{\varsigma_0(u_1)}{\varsigma_{\omega_1;0}(v_1)} (\boldsymbol{\rho} - \lambda_1^+) \leq \frac{\omega_1}{\varsigma_{\omega_1;0}(v_1)} \leq \boldsymbol{\rho} - \lambda_1^+$$

which is (6.43a).

While Lemmas 6.5 and 6.6 are stated for any  $g \in \mathbb{P}_{m-1}$  with the specified conditions satisfied, in their eventual application, it will be taken to be the one that minimizes  $\varepsilon_g$ .

**Lemma 6.5.** *Given  $x \in \mathbb{C}^n$ , assign  $\boldsymbol{\rho}_\pm = \rho_\pm(x)$  and  $\boldsymbol{\rho}_{g;\pm} = \rho_\pm(g(\mathbf{Q}(\boldsymbol{\rho}_+))x)$  for any  $g \in \mathbb{P}_{m-1}$ . Suppose  $\lambda_1^{\text{typ}} \leq \boldsymbol{\rho}_{\text{typ}} < \lambda_2^{\text{typ}}$  if  $\ell = 1$  or  $\lambda_{n-1}^{\text{typ}} < \boldsymbol{\rho}_{\text{typ}} \leq \lambda_n^{\text{typ}}$  if  $\ell = n$ , and let the eigenvalues of  $\mathbf{Q}(\boldsymbol{\rho}_{\text{typ}})$  be  $\omega_j$  for  $1 \leq j \leq n$  which can be arranged as*

$$\begin{aligned} \omega_1 > 0 > \omega_2 \geq \cdots \geq \omega_n & \quad \text{if } (\text{typ}, \ell) \in \{(+, 1), (-, n)\}, \quad \text{or,} \\ \omega_1 < 0 < \omega_2 \leq \cdots \leq \omega_n & \quad \text{if } (\text{typ}, \ell) \in \{(+, n), (-, 1)\}. \end{aligned}$$

Denote by  $v_1$  the eigenvector of  $\mathbf{Q}(\boldsymbol{\rho}_{\text{typ}})$  associated with its eigenvalue  $\omega_1$ . Then for a  $g \in \mathbb{P}_{m-1}$  such that  $g(\omega_1) \neq 0$  and

$$\varepsilon_g := \max_{i \neq 1} \frac{|g(\omega_i)|}{|g(\omega_1)|} < 1, \quad (6.46)$$

we have

$$|\boldsymbol{\rho}_{g;\text{typ}} - \lambda_\ell^{\text{typ}}| \leq |\boldsymbol{\rho}_{\text{typ}} - \lambda_\ell^{\text{typ}}| - \frac{|\omega_1|}{|\boldsymbol{\rho}_{\text{typ}} - \boldsymbol{\rho}_{g;\text{typ}}| a(v_1)} + \frac{|\omega_1|}{|\boldsymbol{\rho}_{\text{typ}} - \boldsymbol{\rho}_{g;\text{typ}'}| a(v_1)} h(\varepsilon_g, \omega_1), \quad (6.47)$$

<sup>16</sup>In fact,

$$\begin{aligned} u_1^H A u_1 (\boldsymbol{\rho} - \lambda_1^+)^2 + \varsigma(u_1)(\boldsymbol{\rho} - \lambda_1^+) &= u_1^H A u_1 [\boldsymbol{\rho}^2 - 2\boldsymbol{\rho}\lambda_1^+ + (\lambda_1^+)^2] + (2\lambda_1^+ u_1^H A u_1 + u_1^H B u_1)(\boldsymbol{\rho} - \lambda_1^+) \\ &= \boldsymbol{\rho}^2 u_1^H A u_1 + \boldsymbol{\rho} u_1^H B u_1 - (\lambda_1^+)^2 u_1^H A u_1 - \lambda_1^+ u_1^H B u_1 \\ &= u_1^H \mathbf{Q}(\boldsymbol{\rho}) u_1 - u_1^H \mathbf{Q}(\lambda_1^+) u_1 \\ &= u_1^H \mathbf{Q}(\boldsymbol{\rho}) u_1, \\ v_1^H A v_1 (\boldsymbol{\rho} - \lambda_1^+)^2 - \varsigma_{\omega_1}(v_1)(\boldsymbol{\rho} - \lambda_1^+) &= v_1^H A v_1 [\boldsymbol{\rho}^2 - 2\boldsymbol{\rho}\lambda_1^+ + (\lambda_1^+)^2] - (2\boldsymbol{\rho} v_1^H A v_1 + v_1^H B v_1)(\boldsymbol{\rho} - \lambda_1^+) \\ &= (\lambda_1^+)^2 v_1^H A v_1 + \lambda_1^+ v_1^H B v_1 - \boldsymbol{\rho}^2 v_1^H A v_1 - \boldsymbol{\rho} v_1^H B v_1 \\ &= v_1^H \mathbf{Q}(\lambda_1^+) v_1 - v_1^H \mathbf{Q}(\boldsymbol{\rho}) v_1 \\ &= v_1^H \mathbf{Q}(\lambda_1^+) v_1 - \omega_1 v_1^H v_1. \end{aligned}$$

They lead to the equations in (6.45) right away.

where  $\text{typ}'$  is the opposite type of  $\text{typ}$ , and

$$h(\varepsilon_g, \omega_1) = 1 - \frac{1 - \varepsilon_g^2}{\left(1 + \varepsilon_g |\omega_1|^{1/2} \tau_A^{1/2}\right)^2}, \quad \tau_A = \frac{1}{|\omega_2|} \frac{\|A\|_2}{a(v_1)}. \quad (6.48)$$

*Proof.* Consider the case  $(\text{typ}, \ell) = (+, 1)$ , and write  $\boldsymbol{\rho} = \boldsymbol{\rho}_+$ . Without loss of generality, we may assume  $\|v_1\|_2 = 1$ . Let the eigenvalue decomposition of  $\mathbf{Q}(\boldsymbol{\rho})$  be

$$\mathbf{Q}(\boldsymbol{\rho}) = V \Sigma V^H, \quad V = [v_1, \dots, v_n], \quad \Sigma = \text{diag}(\omega_1, \dots, \omega_n),$$

where  $\omega_1 > 0 > \omega_2 \geq \dots \geq \omega_n$  and  $V^H V = I_n$ . Set

$$\hat{x} = V^H x = \begin{bmatrix} \xi_1 \\ \xi_2 \\ \vdots \\ \xi_n \end{bmatrix}, \quad \hat{x}_2 = \hat{x} - \xi_1 e_1 = \begin{bmatrix} 0 \\ \xi_2 \\ \vdots \\ \xi_n \end{bmatrix}.$$

Then

$$0 = x^H \mathbf{Q}(\boldsymbol{\rho}) x = \hat{x}^H \Sigma \hat{x} = \omega_1 |\xi_1|^2 + \sum_{i \neq 1} \omega_i |\xi_i|^2. \quad (6.49)$$

Note that for any vector  $z$ ,  $z^H \mathbf{Q}(\lambda) z = z^H A z [\lambda - \rho_+(z)][\lambda - \rho_-(z)]$ . Substitute  $\lambda = \boldsymbol{\rho}$  and  $z = g(\mathbf{Q}(\boldsymbol{\rho}))x$  to get

$$\begin{aligned} \boldsymbol{\rho}_g - \lambda_1^+ &= \boldsymbol{\rho} - \lambda_1^+ - \frac{1}{\boldsymbol{\rho} - \boldsymbol{\rho}_{g;-}} \cdot \frac{x^H g(\mathbf{Q}(\boldsymbol{\rho}))^H \mathbf{Q}(\boldsymbol{\rho}) g(\mathbf{Q}(\boldsymbol{\rho})) x}{x^H g(\mathbf{Q}(\boldsymbol{\rho}))^H A g(\mathbf{Q}(\boldsymbol{\rho})) x} \\ &= \boldsymbol{\rho} - \lambda_1^+ - \frac{1}{\boldsymbol{\rho} - \boldsymbol{\rho}_{g;-}} \cdot \frac{\hat{x}^H g(\Sigma)^H \Sigma g(\Sigma) \hat{x}}{\hat{x}^H g(\Sigma)^H \widehat{A} g(\Sigma) \hat{x}}, \end{aligned} \quad (6.50)$$

where  $\widehat{A} = V^H A V$  and  $\boldsymbol{\rho}_g = \boldsymbol{\rho}_{g;+}$ . We need to estimate the right-hand side of (6.50). We have

$$\begin{aligned} \hat{x}^H g(\Sigma)^H \Sigma g(\Sigma) \hat{x} &= \omega_1 |g(\omega_1)|^2 |\xi_1|^2 + \sum_{i \neq 1} \omega_i |g(\omega_i)|^2 |\xi_i|^2 \\ &\geq \omega_1 |g(\omega_1)|^2 |\xi_1|^2 + \varepsilon_g^2 |g(\omega_1)|^2 \sum_{i \neq 1} \omega_i |\xi_i|^2 \\ &= \omega_1 |g(\omega_1)|^2 |\xi_1|^2 - \varepsilon_g^2 |g(\omega_1)|^2 \omega_1 |\xi_1|^2 \quad (\text{by (6.49)}) \\ &= (1 - \varepsilon_g^2) \omega_1 |g(\omega_1)|^2 |\xi_1|^2, \end{aligned} \quad (6.51)$$

$$\begin{aligned} \hat{x}^H g(\Sigma)^H \widehat{A} g(\Sigma) \hat{x} &= \|g(\Sigma) \hat{x}\|_{\widehat{A}}^2 \\ &= \|g(\omega_1) \xi_1 e_1 + g(\Sigma) \hat{x}_2\|_{\widehat{A}}^2 \\ &\leq \left[ |g(\omega_1)| |\xi_1| \|e_1\|_{\widehat{A}} + \|g(\Sigma) \hat{x}_2\|_{\widehat{A}} \right]^2 \\ &\leq \left[ |g(\omega_1)| |\xi_1| \|e_1\|_{\widehat{A}} + \varepsilon_g |g(\omega_1)| \|\hat{x}_2\|_{\widehat{A}} \right]^2 \\ &\leq \left[ |g(\omega_1)| |\xi_1| \|e_1\|_{\widehat{A}} + \varepsilon_g |g(\omega_1)| \left( \|A\|_2 \frac{\omega_1}{-\omega_2} |\xi_1|^2 \right)^{1/2} \right]^2 \end{aligned} \quad (6.52)$$

$$= |g(\omega_1)|^2 |\xi_1|^2 v_1^H A v_1 \left[ 1 + \varepsilon_g \left( \frac{\omega_1 \|A\|_2}{-\omega_2 v_1^H A v_1} \right)^{1/2} \right]^2, \quad (6.53)$$

where the inequality sign at (6.52) holds because

$$\|\hat{x}_2\|_{\hat{A}}^2 \leq \|\hat{A}\|_2 \|\hat{x}_2\|_2^2 = \|V^H A V\|_2 \sum_{i \neq 1} |\xi_i|^2 \leq \|A\|_2 \frac{\sum_{i \neq 1} \omega_i |\xi_i|^2}{\omega_2} = \|A\|_2 \frac{\omega_1}{-\omega_2} |\xi_1|^2$$

by (6.49). Thus, from (6.50), (6.51), and (6.53),

$$\boldsymbol{\rho}_g - \lambda_1^+ \leq \boldsymbol{\rho} - \lambda_1^+ - \frac{\omega_1}{(\boldsymbol{\rho} - \boldsymbol{\rho}_{g;-}) v_1^H A v_1} \frac{1 - \varepsilon_g^2}{\left[ 1 + \varepsilon_g \left( \frac{\omega_1 \|A\|_2}{-\omega_2 v_1^H A v_1} \right)^{1/2} \right]^2} \quad (6.54)$$

which gives (6.47) for the case  $(\text{typ}, \ell) = (+, 1)$ .  $\square$

**Lemma 6.6.** *Under the conditions of Lemma 6.5, we have*

$$|\boldsymbol{\rho}_{g;\text{typ}} - \lambda_\ell^{\text{typ}}| \leq \frac{|\omega_1|}{\varsigma_0(v_1)} \varepsilon_g^2 + \frac{1 - \varepsilon_g^2}{\varsigma_0(v_1)} \left( 3\tau_A^{1/2} + 2\chi_1 \right) \varepsilon_g |\omega_1|^{3/2} + O(\omega_1^2), \quad (6.55)$$

provided

$$\varepsilon_g |\omega_1|^{1/2} \max\{\tau_A^{1/2}, \zeta \chi_1\} < 1, \quad 4a(v_1) |\omega_1| < \varsigma_0(v_1)^2, \quad (6.56)$$

where  $\tau_A$ ,  $\tau_B$ , and  $\tau_C$  are defined in (6.25), and

$$\chi_1 = \frac{b_0(v_1)^2 \tau_B^{1/2} + 2a(v_1) c_0(v_1) (\tau_A^{1/2} + \tau_C^{1/2})}{\varsigma_0(v_1)^2}, \quad (6.57)$$

$$\zeta = 4 + 6\varepsilon_g \omega_1^{1/2} \tau_B^{1/2} + 4\varepsilon_g^2 \omega_1 \tau_B + \varepsilon_g^3 \omega_1^{3/2} \tau_B^{3/2}, \quad (6.58)$$

and the shift  $\lambda_0 \geq \lambda_n^+$  in defining  $b_0(\cdot)$  and  $c_0(\cdot)$  in (6.21). Alternatively,

$$|\boldsymbol{\rho}_{g;\text{typ}} - \lambda_\ell^{\text{typ}}| \leq \varepsilon_g^2 |\boldsymbol{\rho}_{\text{typ}} - \lambda_\ell^{\text{typ}}| + (1 - \varepsilon_g^2) (3\tau_A^{1/2} + 2\chi_1) \varepsilon_g |\boldsymbol{\rho}_{\text{typ}} - \lambda_\ell^{\text{typ}}|^{3/2} + O(|\boldsymbol{\rho}_{\text{typ}} - \lambda_\ell^{\text{typ}}|^2), \quad (6.59)$$

provided

$$|\boldsymbol{\rho}_{\text{typ}} - \lambda_\ell^{\text{typ}}| < \max \left\{ \frac{\varsigma_0(v_1)}{4a(v_1)}, \frac{1}{\varsigma_0(v_1) \varepsilon_g^2 \max\{\tau_A, \zeta^2 \chi_1^2\}} \right\}. \quad (6.60)$$

*Proof.* Consider the case  $(\text{typ}, \ell) = (+, 1)$ , and write  $\boldsymbol{\rho} = \boldsymbol{\rho}_+$ . Without loss of generality, we may assume  $\|v_1\|_2 = 1$ . Write  $x_g = g(\mathbf{Q}(\boldsymbol{\rho}))x$ , and

$$\begin{aligned} t_M &= \omega_1^{1/2} \tau_M^{1/2} \quad \text{for } M = A, B, C, \\ \mathbf{a} &= a(v_1), \quad \mathbf{b} = b(v_1), \quad \mathbf{c} = c(v_1), \\ \mathbf{b}_0 &= b_0(v_1), \quad \mathbf{c}_0 = c_0(v_1). \end{aligned}$$

By Lemma 6.5,  $\boldsymbol{\rho}_g \leq \boldsymbol{\rho}$  (see (6.54)) and

$$\boldsymbol{\rho}_g - \lambda_1^+ \leq \delta_0 + \delta_1 + \delta_2 + \delta_3, \quad (6.61)$$

where

$$\begin{aligned}
0 \leq \delta_0 &= \boldsymbol{\rho} - \lambda_1^+ - \frac{\omega_1}{\varsigma_{\omega_1;0}(v_1)} = O(|\boldsymbol{\rho} - \lambda_1^+|^2) = O(\omega_1^2), \\
\delta_1 &= \frac{\omega_1}{\varsigma_{\omega_1;0}(v_1)} - \frac{\omega_1}{(\boldsymbol{\rho}_g - \boldsymbol{\rho}_{g;-})\mathbf{a}}, \\
\delta_2 &= \frac{\omega_1}{(\boldsymbol{\rho}_g - \boldsymbol{\rho}_{g;-})\mathbf{a}} - \frac{\omega_1}{(\boldsymbol{\rho} - \boldsymbol{\rho}_{g;-})\mathbf{a}}, \\
\delta_3 &= \frac{\omega_1}{(\boldsymbol{\rho} - \boldsymbol{\rho}_{g;-})\mathbf{a}} h(\varepsilon_g, \omega_1).
\end{aligned} \tag{6.62}$$

The rest of the proof is mainly to estimate  $\delta_1$ ,  $\delta_2$ , and  $\delta_3$ .

For  $\delta_2$ , we have

$$0 \leq \delta_2 = \frac{\omega_1}{\mathbf{a}} \frac{\boldsymbol{\rho} - \boldsymbol{\rho}_g}{(\boldsymbol{\rho}_g - \boldsymbol{\rho}_{g;-})(\boldsymbol{\rho} - \boldsymbol{\rho}_{g;-})} \leq \frac{\omega_1}{\mathbf{a}} \frac{\boldsymbol{\rho} - \lambda_1^+}{(\boldsymbol{\rho}_g - \boldsymbol{\rho}_{g;-})(\boldsymbol{\rho} - \boldsymbol{\rho}_{g;-})} = O(\omega_1^2), \tag{6.63}$$

where we have used (6.44a).

Consider  $\delta_1$ . If  $4a\omega_1 < \mathbf{b}^2 - 4ac$  which holds for sufficiently tiny  $\omega_1$ , then

$$\frac{1}{\varsigma_{\omega_1}(v_1)} = \frac{1}{\sqrt{\mathbf{b}^2 - 4a(\mathbf{c} - \omega_1)}} = \frac{1}{\sqrt{\mathbf{b}^2 - 4ac}} \left[ 1 - \frac{2a}{\mathbf{b}^2 - 4ac} \omega_1 + O(\omega_1^2) \right]. \tag{6.64}$$

By item 2 of Lemma 4.2, any shift  $\lambda_0 \geq \lambda_n^+$  makes  $\mathbf{Q}_{\lambda_0}(\lambda)$  overdamped, i.e.,  $B_{\lambda_0} \succ 0$  and  $C_{\lambda_0} \succeq 0$ . It can be verified that

$$\mathbf{b}_0^2 - 4ac_0 = \mathbf{b}^2 - 4ac = [\varsigma(v_1)]^2.$$

We get, similarly to (6.53),

$$\begin{aligned}
\mathbf{a} |g(\omega_1)|^2 |\xi_1|^2 (1 - 2\varepsilon_g t_A) &\leq x_g^H A x_g \leq \mathbf{a} |g(\omega_1)|^2 |\xi_1|^2 (1 + \varepsilon_g t_A)^2, \\
\mathbf{b}_0 |g(\omega_1)|^2 |\xi_1|^2 (1 - 2\varepsilon_g t_B) &\leq x_g^H B_{\lambda_0} x_g \leq \mathbf{b}_0 |g(\omega_1)|^2 |\xi_1|^2 (1 + \varepsilon_g t_B)^2, \\
\mathbf{c}_0 |g(\omega_1)|^2 |\xi_1|^2 (1 - 2\varepsilon_g t_C) &\leq x_g^H C_{\lambda_0} x_g \leq \mathbf{c}_0 |g(\omega_1)|^2 |\xi_1|^2 (1 + \varepsilon_g t_C)^2.
\end{aligned}$$

Note that  $\boldsymbol{\rho}_g - \lambda_0$  (recalling  $\boldsymbol{\rho}_g$  is the shorthand for  $\boldsymbol{\rho}_{g;+}$ ) and  $\boldsymbol{\rho}_{g;-} - \lambda_0$  are two distinct roots of  $x_g^H A x_g \lambda^2 + x_g^H B_{\lambda_0} x_g \lambda + x_g^H C_{\lambda_0} x_g = 0$  in  $\lambda$ . So

$$\begin{aligned}
\frac{1}{(\boldsymbol{\rho}_g - \boldsymbol{\rho}_{g;-})\mathbf{a}} &= \frac{x_g^H A x_g}{\mathbf{a} \sqrt{(x_g^H B_{\lambda_0} x_g)^2 - 4(x_g^H A x_g)(x_g^H C_{\lambda_0} x_g)}} \\
&\geq \frac{1 - 2\varepsilon_g t_A}{\sqrt{\mathbf{b}_0^2 (1 + \varepsilon_g t_B)^4 - 4ac_0 (1 - 2\varepsilon_g t_A)(1 - 2\varepsilon_g t_C)}} \\
&= \frac{1 - 2\varepsilon_g t_A}{\sqrt{\mathbf{b}_0^2 - 4ac_0 + 4\varepsilon_g (\mathbf{b}_0^2 t_B + 2ac_0 t_A + 2ac_0 t_C) + 2\varepsilon_g^2 (3\mathbf{b}_0^2 t_B^2 - 8ac_0 t_A t_C) + 4\varepsilon_g^3 \mathbf{b}_0^2 t_B^3 + \varepsilon_g^4 \mathbf{b}_0^2 t_B^4}} \\
&= \frac{1 - 2\varepsilon_g t_A}{\sqrt{(\mathbf{b}_0^2 - 4ac_0)(1 + 4\varepsilon_g \chi_1 \omega_1^{1/2} + 2\varepsilon_g^2 \chi_2 \omega_1) + 4\varepsilon_g^3 \mathbf{b}_0^2 t_B^3 + \varepsilon_g^4 \mathbf{b}_0^2 t_B^4}} \\
&= \frac{1}{\sqrt{\mathbf{b}_0^2 - 4ac_0}} (1 - 2\varepsilon_g \omega_1^{1/2} \tau_A^{1/2}) \left[ 1 - 2\varepsilon_g \chi_1 \omega_1^{1/2} + \varepsilon_g^2 (6\chi_1^2 - \chi_2) \omega_1 + \dots \right] \tag{6.65}
\end{aligned}$$

$$= \frac{1}{\sqrt{\mathbf{b}^2 - 4\mathbf{ac}}} \left[ 1 - 2\varepsilon_g(\tau_A^{1/2} + \chi_1)\omega_1^{1/2} + \varepsilon_g^2(6\chi_1^2 - \chi_2 + 4\tau_A^{1/2}\chi_1)\omega_1 + O(\omega_1^{3/2}) \right], \quad (6.66)$$

where

$$\chi_1 = \frac{\mathbf{b}_0^2\tau_B^{1/2} + 2\mathbf{ac}_0(\tau_A^{1/2} + \tau_C^{1/2})}{\mathbf{b}^2 - 4\mathbf{ac}}, \quad \chi_2 = \frac{3\mathbf{b}_0^2\tau_B - 8\mathbf{ac}_0\tau_A^{1/2}\tau_C^{1/2}}{\mathbf{b}^2 - 4\mathbf{ac}}.$$

In obtaining (6.65), we need<sup>17</sup>  $\zeta\varepsilon_g\chi_1\omega_1^{1/2} < 1$ , where  $\zeta = 4 + 6\varepsilon_g t_B + 4\varepsilon_g^2 t_B^2 + \varepsilon_g^3 t_B^3$ . Using (6.66), we have for  $\delta_1$

$$\begin{aligned} \delta_1 &= \frac{\omega_1}{\varsigma_{\omega_1;0}(v_1)} - \frac{\omega_1}{(\boldsymbol{\rho}_g - \boldsymbol{\rho}_{g;-})\mathbf{a}} \\ &= \frac{\omega_1}{\sqrt{\mathbf{b}^2 - 4\mathbf{ac}}} \left[ 1 - \frac{2\mathbf{a}}{\mathbf{b}^2 - 4\mathbf{ac}}\omega_1 + O(\omega_1^2) \right] \\ &\quad - \frac{\omega_1}{\sqrt{\mathbf{b}^2 - 4\mathbf{ac}}} \left[ 1 - 2\varepsilon_g(\tau_A^{1/2} + \chi_1)\omega_1^{1/2} + \varepsilon_g^2(6\chi_1^2 - \chi_2 + 4\tau_A^{1/2}\chi_1)\omega_1 + O(\omega_1^{3/2}) \right] \\ &= \frac{2\varepsilon_g(\tau_A^{1/2} + \chi_1)\omega_1^{3/2}}{\sqrt{\mathbf{b}^2 - 4\mathbf{ac}}} + O(\omega_1^2). \end{aligned} \quad (6.67)$$

Now we turn to  $\delta_3$ . If  $\varepsilon_g t_A < 1$ , then

$$\begin{aligned} h(\varepsilon_g, \omega_1) &= 1 - (1 - \varepsilon_g^2)(1 + \varepsilon_g t_A)^{-2} \\ &= 1 - (1 - \varepsilon_g^2)(1 - \varepsilon_g t_A + 2\varepsilon_g^2 t_A^2 - 3\varepsilon_g^3 t_A^3 + \dots) \\ &= \varepsilon_g^2 + (1 - \varepsilon_g^2)(\varepsilon_g t - 2\varepsilon_g^2 t_A^2 + \dots) \\ &= \varepsilon_g^2 + \varepsilon_g(1 - \varepsilon_g^2)t_A + O(t_A^2) \\ &= \varepsilon_g^2 + \varepsilon_g(1 - \varepsilon_g^2)\omega_1^{1/2}\tau_A^{1/2} + O(\omega_1), \\ h(\varepsilon_g, \omega_1) &= 1 - (1 - \varepsilon_g^2)(1 + t_A \varepsilon_g)^{-2} \\ &\geq 1 - (1 - \varepsilon_g^2) \\ &= \varepsilon_g^2 \geq 0. \end{aligned}$$

Therefore

$$\begin{aligned} \delta_3 &= \frac{\omega_1}{(\boldsymbol{\rho} - \boldsymbol{\rho}_{g;-})\mathbf{a}} h(\varepsilon_g, \omega_1) \\ &= \frac{\omega_1 \varepsilon_g^2 + \varepsilon_g(1 - \varepsilon_g^2)\omega_1^{3/2}\tau_A^{1/2}}{(\boldsymbol{\rho} - \boldsymbol{\rho}_{g;-})\mathbf{a}} + O(\omega_1^2). \end{aligned} \quad (6.68)$$

<sup>17</sup>For the expansion in (6.65), it is needed that

$$4\varepsilon_g\chi_1\omega_1^{1/2} + 2\varepsilon_g^2\chi_2\omega_1 + \frac{4\varepsilon_g^3\mathbf{b}_0^2 t_B^3}{\mathbf{b}^2 - 4\mathbf{ac}} + \frac{\varepsilon_g^4\mathbf{b}_0^2 t_B^4}{\mathbf{b}^2 - 4\mathbf{ac}} < 1.$$

However,

$$\frac{2\varepsilon_g^2\chi_2\omega_1 + \frac{4\varepsilon_g^3\mathbf{b}_0^2 t_B^3}{\mathbf{b}^2 - 4\mathbf{ac}} + \frac{\varepsilon_g^4\mathbf{b}_0^2 t_B^4}{\mathbf{b}^2 - 4\mathbf{ac}}}{4\varepsilon_g\chi_1\omega_1^{1/2}} \leq \frac{2\varepsilon_g^2 3\mathbf{b}_0^2 t_B^2 + 4\varepsilon_g^3\mathbf{b}_0^2 t_B^3 + \varepsilon_g^4\mathbf{b}_0^2 t_B^4}{4\varepsilon_g\mathbf{b}_0^2 t_B} = \frac{\varepsilon_g t_B}{4} (6 + 4\varepsilon_g t_B + \varepsilon_g^2 t_B^2).$$

We have finished estimating  $\delta_i$  for  $i = 0, 1, 2, 3$ . Now, combine (6.61), (6.62), (6.63), (6.67), and (6.68) to get

$$\begin{aligned}\rho_g - \lambda_1^+ &\leq \frac{2\varepsilon_g(\tau_A^{1/2} + \chi_1)\omega_1^{3/2}}{\sqrt{\mathbf{b}^2 - 4\mathbf{a}\mathbf{c}}} + \frac{\omega_1\varepsilon_g^2 + \varepsilon_g(1 - \varepsilon_g^2)\omega_1^{3/2}\tau_A^{1/2}}{(\rho - \rho_{g;-})\mathbf{a}} + O(\omega_1^2) \\ &= \frac{\varepsilon_g^2}{(\rho - \rho_{g;-})\mathbf{a}}\omega_1 + \left( \frac{2(\tau_A^{1/2} + \chi_1)}{\sqrt{\mathbf{b}^2 - 4\mathbf{a}\mathbf{c}}} + \frac{(1 - \varepsilon_g^2)\tau_A^{1/2}}{(\rho - \rho_{g;-})\mathbf{a}} \right) \varepsilon_g\omega_1^{3/2} + O(\omega_1^2),\end{aligned}$$

which, along with

$$\frac{1}{(\rho - \rho_{g;-})\mathbf{a}} = \frac{1}{(\rho_g - \rho_{g;-})\mathbf{a}} - \frac{\delta_2}{\omega_1} = \frac{1}{\sqrt{\mathbf{b}^2 - 4\mathbf{a}\mathbf{c}}} \left[ 1 - 2\varepsilon_g(\tau_A^{1/2} + \chi_1)\omega_1^{1/2} \right] + O(\omega_1),$$

yield (6.55). Use (6.64) to see

$$\frac{1}{\varsigma_0(v_1)} = \frac{1}{\varsigma_{\omega_1;0}(v_1)} \left[ 1 + \frac{2\mathbf{a}}{\mathbf{b}^2 - 4\mathbf{a}\mathbf{c}}\omega_1 + O(\omega_1^2) \right]$$

substituting which and (6.44a) into (6.55) to get (6.59).  $\square$

We are now ready to prove Theorem 6.2.

*Proof* of Theorem 6.2. Item 1 is a direct consequence of item 4 of Theorem 6.1.

Item 2 is a consequence of Lemma 6.6 upon letting  $g$  be the minimizer that gives the minimal  $\varepsilon_m$  and using  $|\rho_{i+1} - \lambda_\ell^{\text{typ}}| \leq |\rho_g - \lambda_\ell^{\text{typ}}|$ .

For item 3, again let  $g$  be the minimizer that gives the minimal  $\varepsilon_m$ . As  $i \rightarrow \infty$  in item 2, we have  $\omega_1 \rightarrow 0$ ,  $\omega_2 \rightarrow \gamma$ , and  $v_1 \rightarrow u_\ell^{\text{typ}}$  in direction, and thus

$$\lim_{i \rightarrow \infty} \eta(v_1) = \lim_{i \rightarrow \infty} 3\tau_A^{1/2} + 2 \frac{(b_0(v_1))^2 \tau_{B\lambda_0}^{1/2} + 2a(v_1)c_0(v_1)(\tau_A^{1/2} + \tau_{C\lambda_0}^{1/2})}{\varsigma_0(v_1)^2} = \eta$$

as given by (6.29). Now let

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

where  $\mathcal{T}_{m-1}(t)$  is the  $(m-1)$ st Chebyshev polynomial of the first kind. Then [34, section 2]

$$\varepsilon_m \leq \varepsilon_{\hat{g}} \leq \max_{\omega_2 \leq t \leq \omega_n} |\hat{g}(t)| = 2 \left[ \left( \frac{1 + \sqrt{\hat{\kappa}}}{1 - \sqrt{\hat{\kappa}}} \right)^{m-1} + \left( \frac{1 + \sqrt{\hat{\kappa}}}{1 - \sqrt{\hat{\kappa}}} \right)^{-(m-1)} \right]^{-1}$$

which goes to  $\varepsilon$  as  $i \rightarrow \infty$  because  $\hat{\kappa} \rightarrow \kappa$ .  $\square$

## 7 Preconditioned steepest descent/ascent method

### 7.1 Preconditioning

We will explain the idea of preconditioning for computing  $(\lambda_1^+, u_1^+)$  only, via two different points of view. The same argument can be made for other extreme pos- and neg-quadratic eigenpairs.

It is well-known that when the contours of the objective function near its optimum are extremely elongated, at each step of the conventional steepest descent/ascent method, following the search direction which is the opposite of the gradient gets closer to the optimum on the line for a very short while and then starts to get away because the direction doesn't point "towards the optimum", resulting in a long zigzag path of a large number of steps. The ideal search direction  $p$  is therefore the one such that with its starting point at  $\mathbf{x}$ ,  $p$  points to the optimum, i.e., the optimum is on the line  $\{\mathbf{x} + tp : t \in \mathbb{C}\}$ . Specifically, expand  $\mathbf{x}$  as a linear combination of  $u_j^+$

$$\mathbf{x} = \sum_{j=1}^n \alpha_j u_j^+ =: \alpha_1 u_1^+ + \mathbf{v}, \quad \mathbf{v} = \sum_{j=2}^n \alpha_j u_j^+. \quad (7.1)$$

Then the ideal search direction is

$$p = \alpha u_1^+ + \beta \mathbf{v}$$

for some scalar  $\alpha$  and  $\beta \neq 0$  such that  $\alpha_1 \beta - \alpha \neq 0$  (otherwise  $p = \beta \mathbf{x}$ ). Of course, this is impractical because we don't know  $u_1^+$  and  $\mathbf{v}$ . But we can construct one that is close to it. One such  $p$  is

$$p = [\mathbf{Q}(\sigma)]^{-1} r_+(\mathbf{x}) = [\mathbf{Q}(\sigma)]^{-1} \mathbf{Q}(\boldsymbol{\rho}_+) \mathbf{x},$$

where  $\boldsymbol{\rho}_+ = \rho_+(\mathbf{x})$  and<sup>18</sup>  $\sigma$  is some shift near  $\lambda_1^+$  but not equal to  $\boldsymbol{\rho}_+$ . Let us analyze this  $p$ . By (2.17a), we have

$$[\mathbf{Q}(\sigma)]^{-1} \mathbf{Q}(\boldsymbol{\rho}_+) = U_+(\sigma I - \Lambda_+)^{-1} (U_-^H A U_+)^{-1} (\sigma I - \Lambda_-)^{-1} (\boldsymbol{\rho}_+ I - \Lambda_-) U_-^H A U_+ (\boldsymbol{\rho}_+ I - \Lambda_+) U_+^{-1}.$$

Suppose now that both  $\sigma$  and  $\boldsymbol{\rho}_+$  are near  $\lambda_1^+$ . Then

$$(\sigma I - \Lambda_-)^{-1} (\boldsymbol{\rho}_+ I - \Lambda_-) = I + (\boldsymbol{\rho}_+ - \sigma)(\sigma I - \Lambda_-)^{-1} \approx I.$$

Therefore  $[\mathbf{Q}(\sigma)]^{-1} \mathbf{Q}(\boldsymbol{\rho}_+) \approx U_+(\sigma I - \Lambda_+)^{-1} (\boldsymbol{\rho}_+ I - \Lambda_+) U_+^{-1}$ , and thus

$$p = [\mathbf{Q}(\sigma)]^{-1} \mathbf{Q}(\boldsymbol{\rho}_+) \mathbf{x} \approx \sum_{j=1}^n \mu_j \alpha_j u_j^+, \quad \mu_j := \frac{\lambda_j^+ - \boldsymbol{\rho}_+}{\lambda_j^+ - \sigma}. \quad (7.2)$$

Now if  $\lambda_1^+ \leq \boldsymbol{\rho}_+ < \lambda_2^+$  and if the gap  $\lambda_2^+ - \lambda_1^+$  is reasonably modest, then

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

to give a  $p \approx \alpha u_1^+ + \mathbf{v}$ , resulting in fast convergence. This rough but intuitive analysis suggests that  $K = [\mathbf{Q}(\sigma)]^{-1}$  with a suitably chosen shift  $\sigma$  can be used to serve as a

<sup>18</sup>We reasonably assume also  $\sigma \neq \lambda_j^+$  for all  $j$ , too.



good preconditioner to improve the steepest descent/ascent method – Algorithm 6.1 by simply modifying  $Y_i = [\mathbf{x}_i, \mathbf{r}_i]$  at Line 6 there to  $Y_i = [\mathbf{x}_i, K\mathbf{r}_i]$ . We caution the reader that implementing  $K\mathbf{r}_i$  is amount to solving a linear system. This is usually done approximately by, e.g., some iterative methods such as the linear conjugate gradient method, MINRES [11, 17, 19].

The second view point is similar to the one proposed by Golub and Ye [18] for the generalized linear eigenvalue problem. Theorem 6.2 reveals that the rates of convergence for Algorithms 6.1 and 6.2 depend on the distribution of the eigenvalues  $\omega_j$  of  $\mathbf{Q}(\boldsymbol{\rho}_i)$ , not the quadratic eigenvalues of  $\mathbf{Q}(\lambda)$ . In particular, if all  $\omega_2 = \dots = \omega_n$ , then  $\epsilon_m = 0$  for  $m \geq 2$  and thus

$$\boldsymbol{\rho}_{i+1} - \lambda_1^+ = \mathcal{O}(|\boldsymbol{\rho}_i - \lambda_1^+|^2),$$

suggesting quadratic convergence. Such an extreme case, though highly welcome, is unlikely to happen in practice, but it gives us an idea that if somehow we could transform an eigenvalue problem towards such an extreme case, the transformed problem would be easier to solve. Specifically we should seek equivalent transformations that change the eigenvalues of  $\mathbf{Q}(\boldsymbol{\rho}_i)$  as much as possible to,

$$\boxed{\text{one isolated eigenvalue } \omega_1, \text{ and the rest } \omega_j \ (2 \leq j \leq n) \text{ tightly clustered,}} \quad (7.3)$$

but leave the quadratic eigenvalues of  $\mathbf{Q}(\lambda)$  unchanged.

We would like to equivalently transform the QEP for  $\mathbf{Q}(\lambda)$  to for  $L^{-1}\mathbf{Q}(\lambda)L^{-H}$  by some nonsingular  $L$  (whose inverse or any linear system with  $L$  is easy to solve) so that the eigenvalues of  $L^{-1}\mathbf{Q}(\boldsymbol{\rho}_i)L^{-H}$  distribute more or less like (7.3). Then apply one step of Algorithm 6.1 or 6.2 to the pencil  $L^{-1}\mathbf{Q}(\lambda)L^{-H}$  to find the next approximation  $\boldsymbol{\rho}_{i+1}$ . The process repeats, i.e., find a new  $L$  to transform the problem and apply one step of Algorithm 6.1 or 6.2 to the transformed problem.

Such an  $L$  may be constructed using the  $LDL^H$  decomposition of  $\mathbf{Q}(\boldsymbol{\rho}_i)$  [17, p.139] if the decomposition exists:  $\mathbf{Q}(\boldsymbol{\rho}_i) = LDL^H$ , where  $L$  is lower triangular and  $D = \text{diag}(\pm 1)$ . Then  $L^{-1}\mathbf{Q}(\boldsymbol{\rho}_i)L^{-H} = D$  has the ideal eigenvalue distribution that gives  $\epsilon_m = 0$  for any  $m \geq 2$ . Unfortunately, this simple solution is impractical in practice for the following reasons:

1. The decomposition may not exist at all. In theory, the decomposition exists if all the leading principle submatrices of  $\mathbf{Q}(\boldsymbol{\rho}_i)$  are nonsingular.
2. If the decomposition does exist, it may not be numerically stable to compute, especially when  $\boldsymbol{\rho}_i$  comes closer and closer to  $\lambda_1^+$ .
3. The sparsity in  $\mathbf{Q}(\boldsymbol{\rho}_i)$  is most likely destroyed, leaving  $L$  significantly denser than  $\mathbf{Q}(\boldsymbol{\rho}_i)$ . This makes all ensuing computations much more expensive.

A more practical solution is, however, through an incomplete  $LDL^H$  factorization (see [51, Chapter 10]), to get

$$\mathbf{Q}(\boldsymbol{\rho}_i) \approx LDL^H,$$

where “ $\approx$ ” includes not only the usual “approximately equal”, but also the case when  $\mathbf{Q}(\boldsymbol{\rho}_i) - LDL^H$  is approximately a low rank matrix, and  $D = \text{diag}(\pm 1)$ . Such an  $L$  changes from one step of the algorithm to another. In practice, often we may use one

fixed preconditioner for all or a number of consecutive iterative steps. Using a constant preconditioner is certainly not optimal: it likely doesn't give the best rate of convergence per step and thus increases the number of total iterative steps but it may reduce overall cost because it saves work in preconditioner constructions and thus reduces cost per step. The basic idea of using a step-independent preconditioner is to find a  $\sigma$  that is close to  $\lambda_1^+$ , and perform an incomplete  $LDL^H$  decomposition:

$$\mathbf{Q}(\sigma) \approx LDL^H$$

and transform  $\mathbf{Q}(\lambda)$  accordingly before applying Algorithm 6.1 or 6.2. Now the rate of convergence is determined by the eigenvalues of

$$L^{-1}\mathbf{Q}(\rho_i)L^{-H} = L^{-1}\mathbf{Q}(\sigma)L^{-H} + (\rho_i - \sigma)L^{-1}\mathbf{Q}'(\sigma)L^{-H} + O(|\rho_i - \sigma|^2)$$

which would have a better spectral distribution so long as the last two terms is small relative to  $L^{-1}\mathbf{Q}(\rho_i)L^{-H}$ . When  $\lambda_n^- < \sigma < \lambda_1^+$ ,  $-\mathbf{Q}(\sigma) \succ 0$  and the incomplete  $LDL^H$  factorization becomes incomplete Cholesky factorization.

## 7.2 Preconditioned steepest descent/ascent method

We have insisted so far about applying Algorithm 6.1 or 6.2 straightforwardly to the transformed problem. There is another way, perhaps, better: only symbolically applying Algorithm 6.1 or 6.2 to the transformed problem as a derivation tool for a preconditioned method that always projects the original pencil  $\mathbf{Q}(\lambda)$  directly every step. The only difference is now the projecting subspaces are preconditioned. Again we will explain it for the case of computing the first pos-type quadratic eigenpair  $(\lambda_1^+, u_1^+)$ .

Suppose  $\mathbf{Q}(\lambda)$  is transformed to  $\widehat{\mathbf{Q}}(\lambda) := L^{-1}\mathbf{Q}(\lambda)L^{-H}$ . Consider a typical step of Algorithm 6.2 applied to  $\widehat{\mathbf{Q}}(\lambda)$ . For the purpose of distinguishing notational symbols, we will put hats on all those for  $\widehat{\mathbf{Q}}(\lambda)$ . The typical step of Algorithm 6.2 on  $\widehat{\mathbf{Q}}$  is

compute the smallest pos-type quadratic eigenvalue $\mu$ and corresponding quadratic eigenvector $\hat{v}$ of $\hat{Z}^H\widehat{\mathbf{Q}}(\lambda)\hat{Z}$ , where $\hat{Z} \in \mathbb{C}^{n \times m}$ is a basis matrix of Krylov subspace $\mathcal{K}_m(\widehat{\mathbf{Q}}(\hat{\rho}), \hat{\mathbf{x}})$ .	(7.4)
--	-------

Notice  $[\widehat{\mathbf{Q}}(\hat{\rho})]^j \hat{\mathbf{x}} = L^H [(LL^H)^{-1}\mathbf{Q}(\hat{\rho})]^j (L^{-H}\hat{\mathbf{x}})$  to see

$$L^{-H} \cdot \mathcal{K}_m(\widehat{\mathbf{Q}}(\hat{\rho}), \hat{\mathbf{x}}) = \mathcal{K}_m(K\mathbf{Q}(\hat{\rho}), \mathbf{x}),$$

where  $\mathbf{x} = L^{-H}\hat{\mathbf{x}}$  and  $K = (LL^H)^{-1}$ . So  $Z = L^{-H}\hat{Z}$  is a basis matrix of Krylov subspace  $\mathcal{K}_m(K\mathbf{Q}(\hat{\rho}), \mathbf{x})$ . Since also

$$\begin{aligned} \hat{Z}^H\widehat{\mathbf{Q}}(\lambda)\hat{Z} &= (L^{-H}\hat{Z})^H\mathbf{Q}(\lambda)(L^{-H}\hat{Z}), \\ \hat{\rho} &= \hat{\rho}_+(\hat{\mathbf{x}}) = \rho_+(\mathbf{x}) = \rho, \end{aligned}$$

the typical step (7.4) can be reformulated equivalently to

compute the smallest pos-type quadratic eigenvalue $\mu$ and corresponding quadratic eigenvector $v$ of $Z^H\mathbf{Q}(\lambda)Z$ , where $Z \in \mathbb{C}^{n \times m}$ is a basis matrix of Krylov subspace $\mathcal{K}_m(K\mathbf{Q}(\rho), \mathbf{x})$ , where $K = (LL^H)^{-1}$ .	(7.5)
---	-------

---

**Algorithm 7.1** Preconditioned extended steepest descent/ascent method
 

---

Given an initial approximation  $\mathbf{x}_0$  to  $u_\ell^{\text{typ}}$ , and a relative tolerance  $\text{rtol}$ , and the search space dimension  $m$ , the algorithm computes an approximate pair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$  with the prescribed  $\text{rtol}$ .

---

```

1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|$ ,  $\boldsymbol{\rho}_0 = \rho_{\text{typ}}(\mathbf{x}_0)$ ,  $\mathbf{r}_0 = r_{\text{typ}}(\mathbf{x}_0)$ ;
2: for  $i = 0, 1, \dots$  do
3:   if  $\|\mathbf{r}_i\| / (|\boldsymbol{\rho}_i|^2 \|A\mathbf{x}_i\| + |\boldsymbol{\rho}_i| \|B\mathbf{x}_i\| + \|C\mathbf{x}_i\|) \leq \text{rtol}$  then
4:     BREAK;
5:   else
6:     construct a preconditioner  $K_i$ ;
7:     compute a basis matrix  $Y_i$  for the Krylov subspace  $\mathcal{K}_m(K_i \mathbf{Q}(\boldsymbol{\rho}_i), \mathbf{x}_i)$ ;
8:     solve HQEP for  $Y_i^H \mathbf{Q}(\lambda) Y_i$  to get its quadratic eigenvalues  $\mu_j^\pm$  as in (6.11) and
       quadratic eigenvectors  $y_j^\pm$ ;
9:     select the next approximate quadratic eigenpair  $(\mu, y) = (\mu_j^{\text{typ}}, Y y_j^{\text{typ}})$  according
       to the table in (6.12);
10:     $\mathbf{x}_{i+1} = y / \|y\|$ ,  $\boldsymbol{\rho}_{i+1} = \mu$ ,  $\mathbf{r}_{i+1} = r_{\text{typ}}(\mathbf{x}_{i+1})$ ;
11:   end if
12: end for
13: return  $(\boldsymbol{\rho}_i, \mathbf{x}_i)$  as an approximate quadratic eigenpair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$ .

```

---

We are now ready to state a version of the *preconditioned extended steepest descent/ascent method*. To make it be inclusive, in Algorithm 7.1 we use  $K_i$  to denote the preconditioner at the  $i$ th iterative step. Once again, they may all be the same or vary from one iterative step to another. Although the derivation of this algorithm was for the preconditioners obtained from the second view point above, its final form includes the preconditioners from the first view point.

### 7.3 Convergence analysis

If  $K_i \succ 0$ , the  $i$ th iterative step of Algorithm 7.1 is just one step of the extended steepest descent/ascent method applied to  $K_i^{1/2} \mathbf{Q}(\lambda) K_i^{1/2}$ . Therefore Theorem 6.2 implies the following theorem for Algorithm 7.1.

**Theorem 7.1.** *Suppose  $\lambda_1^{\text{typ}} \leq \boldsymbol{\rho}_0 < \lambda_2^{\text{typ}}$  if  $\ell = 1$  or  $\lambda_{n-1}^{\text{typ}} < \boldsymbol{\rho}_0 \leq \lambda_n^{\text{typ}}$  if  $\ell = n$ , and let the sequences  $\{\boldsymbol{\rho}_i\}, \{\mathbf{r}_i\}, \{\mathbf{x}_i\}$  be produced by Algorithm 7.1. Suppose  $K_i \succ 0$ .*

1. *As  $i \rightarrow \infty$ ,  $\boldsymbol{\rho}_i$  monotonically converges to  $\hat{\rho} = \lambda_\ell^{\text{typ}}$ , and  $\mathbf{x}_i$  converges to  $u_\ell^{\text{typ}}$  in direction, i.e.,  $\theta(\mathbf{x}_i, u_\ell^{\text{typ}}) \rightarrow 0$ .*
2. *The eigenvalues<sup>19</sup>  $\omega_j$  of  $K_i \mathbf{Q}(\boldsymbol{\rho}_i)$  can be ordered as*

$$\omega_1 > 0 > \omega_2 \geq \dots \geq \omega_n \quad \text{if } (\text{typ}, \ell) \in \{(+, 1), (-, n)\}, \quad \text{or}, \quad (7.6a)$$

$$\omega_1 < 0 < \omega_2 \leq \dots \leq \omega_n \quad \text{if } (\text{typ}, \ell) \in \{(+, n), (-, 1)\}. \quad (7.6b)$$

---

<sup>19</sup>Their dependency upon  $i$  is suppressed for clarity.

If  $\rho_i$  is sufficiently close to  $\lambda_\ell^{\text{typ}}$ , then

$$|\rho_{i+1} - \lambda_\ell^{\text{typ}}| \leq \varepsilon_m^2 |\rho_i - \lambda_\ell^{\text{typ}}| + O\left(\varepsilon_m |\rho_i - \lambda_\ell^{\text{typ}}|^{3/2} + |\rho_i - \lambda_\ell^{\text{typ}}|^2\right), \quad (7.7)$$

where  $\varepsilon_m$  is defined as in (6.24).

3. Denote<sup>20</sup> by  $\gamma$  and  $\Gamma$  the smallest and largest positive eigenvalue of

$$\begin{cases} -K_i \mathbf{Q}(\lambda_\ell^{\text{typ}}) & \text{for } (\text{typ}, \ell) \in \{(+, 1), (-, n)\}, \\ K_i \mathbf{Q}(\lambda_\ell^{\text{typ}}) & \text{for } (\text{typ}, \ell) \in \{(+, n), (-, 1)\}. \end{cases}$$

If  $\rho_i$  is sufficiently close to  $\lambda_\ell^{\text{typ}}$ , then

$$|\rho_{i+1} - \lambda_\ell^{\text{typ}}| \leq \varepsilon^2 |\rho_i - \lambda_\ell^{\text{typ}}| + O\left(\varepsilon |\rho_i - \lambda_\ell^{\text{typ}}|^{3/2} + |\rho_i - \lambda_\ell^{\text{typ}}|^2\right), \quad (7.8)$$

where  $\varepsilon$  is defined as in (6.28).

There is a convergence rate estimate, essentially due to Samokish [52, 1958], for the preconditioned steepest descent/ascent method in the case of the standard Hermitian eigenvalue problem. The reader is referred to [29, 46] for detail. Theorem 7.2 below is an extension of Samokish's result for our case.

**Theorem 7.2.** *Suppose  $K \succ 0$ . Let  $\ell \in \{1, n\}$  and  $\text{typ}, \text{typ}' \in \{+, -\}$  such that  $\text{typ}$  and  $\text{typ}'$  are opposite, and denote by  $\gamma$  and  $\Gamma$  the smallest and largest positive eigenvalue of*

$$\begin{cases} -K \mathbf{Q}(\lambda_\ell^{\text{typ}}) & \text{for } (\text{typ}, \ell) \in \{(+, 1), (-, n)\}, \\ K \mathbf{Q}(\lambda_\ell^{\text{typ}}) & \text{for } (\text{typ}, \ell) \in \{(+, n), (-, 1)\}, \end{cases}$$

and

$$\tau = \frac{2}{\gamma + \Gamma}, \quad \kappa = \frac{\Gamma}{\gamma}, \quad \varepsilon = \frac{\kappa - 1}{\kappa + 1}.$$

Let  $\text{argopt}$  be as given in (6.6), and

$$\begin{aligned} t_{\text{opt}} &= \underset{t \in \mathbb{C}}{\text{argopt}} \rho_{\text{typ}}(x + tKr_{\text{typ}}(x)), \quad y = x + t_{\text{opt}}Kr_{\text{typ}}(x), \\ z &= \begin{cases} x + \tau Kr_{\pm}(x) & \text{for } (\text{typ}, \ell) \in \{(+, 1), (-, n)\}, \\ x - \tau Kr_{\pm}(x) & \text{for } (\text{typ}, \ell) \in \{(+, n), (-, 1)\}. \end{cases} \end{aligned}$$

We have

$$\begin{aligned} |\rho_{\text{typ}}(y) - \lambda_\ell^{\text{typ}}| &\leq |\rho_{\text{typ}}(z) - \lambda_\ell^{\text{typ}}| \\ &\leq \frac{1}{|\lambda_\ell^{\text{typ}} - \rho_{\text{typ}'}(z)|} \left[ \frac{\varepsilon \sqrt{|\lambda_\ell^{\text{typ}} - \rho_{\text{typ}'}(x)|} + \tau \sqrt{\Gamma} \delta_1}{1 - \tau (\sqrt{\Gamma} \delta_2 + \delta_3^2)} \right]^2 |\rho_{\text{typ}}(x) - \lambda_\ell^{\text{typ}}|, \end{aligned} \quad (7.9)$$

<sup>20</sup>It is worth emphasizing that  $K_i \mathbf{Q}(\lambda_\ell^{\text{typ}})$  is singular and, by Theorem 2.1,  $K_i^{-1/2} \mathbf{Q}(\lambda_\ell^{\text{typ}}) K_i^{1/2}$  is negative semidefinite if  $(\text{typ}, \ell) \in \{(+, 1), (-, n)\}$  and positive semidefinite if  $(\text{typ}, \ell) \in \{(+, n), (-, 1)\}$ .

provided  $\tau \left( \sqrt{\Gamma} \delta_2 + \delta_3^2 \right) < 1$ , where

$$\begin{aligned}\delta_1 &= \sqrt{|\rho_{\text{typ}}(x) - \lambda_\ell^{\text{typ}}| \|K^{1/2}\{A[\rho_{\text{typ}}(x) + \lambda_\ell^{\text{typ}}] + B\}A^{-1/2}\|_2}, \\ \delta_2 &= \sqrt{\|K^{1/2}AK^{1/2}\|_2 |\rho_{\text{typ}}(x) - \lambda_\ell^{\text{typ}}| \cdot |\lambda_\ell^{\text{typ}} - \rho_{\text{typ}}(x)|}, \\ \delta_3 &= \sqrt{\|A^{1/2}K\{A[\rho_{\text{typ}}(x) + \lambda_\ell^{\text{typ}}] + B\}A^{-1/2}\|_2 |\rho_{\text{typ}}(x) - \lambda_\ell^{\text{typ}}|}.\end{aligned}$$

*Proof.* We will prove the case:  $(\text{typ}, \ell) = (+, 1)$  only. The other cases can be handled in the same way.

Note  $z = x + \tau K r_+(x) = x + \tau K \mathbf{Q}(\rho_+(x))x$ . We have  $\lambda_1^+ \leq \rho_+(y) \leq \rho_+(z)$  and thus  $\rho_+(y) - \lambda_1^+ \leq \rho_+(z) - \lambda_1^+$ . So it remains to show that  $\rho_+(z) - \lambda_1^+$  is no bigger than the right-hand side of (7.9).

Let  $M = -\mathbf{Q}(\lambda_1^+) \succeq 0$ . For any vector  $w$ , we have

$$\begin{aligned}\|w\|_M^2 &= -w^H \mathbf{Q}(\lambda_1^+) w \\ &= [\rho_+(w) - \lambda_1^+][\lambda_1^+ - \rho_-(w)] \|w\|_A^2,\end{aligned}\tag{7.10}$$

$$\begin{aligned}\|[I + \tau K \mathbf{Q}(\lambda_1^+)]w\|_M &= \|[I - \tau K M]w\|_M \\ &\leq \varepsilon \|w\|_M.\end{aligned}\tag{7.11}$$

Write

$$\begin{aligned}z &= [I + \tau K \mathbf{Q}(\lambda_1^+)]x - \tau K [\mathbf{Q}(\lambda_1^+) - \mathbf{Q}(\rho_+(x))]x \\ &= [I + \tau K \mathbf{Q}(\lambda_1^+)]x + \tau [\rho_+(x) - \lambda_1^+] K [A(\rho_+(x) + \lambda_1^+) + B]x.\end{aligned}$$

Without loss of generality, we may assume  $\|x\|_A = 1$ . We have

$$\begin{aligned}\|z\|_M &= \sqrt{[\rho_+(z) - \lambda_1^+][\lambda_1^+ - \rho_-(z)]} \|z\|_A, \quad \text{by (7.10)} \\ \|z\|_M &\leq \|[I + \tau K \mathbf{Q}(\lambda_1^+)]x\|_M + \tau [\rho_+(x) - \lambda_1^+] \|K[A(\rho_+(x) + \lambda_1^+) + B]x\|_M \\ &\leq \varepsilon \|x\|_M + \tau [\rho_+(x) - \lambda_1^+] \sqrt{\Gamma} \| [A(\rho_+(x) + \lambda_1^+) + B]x \|_K \\ &\leq \varepsilon \sqrt{[\rho_+(x) - \lambda_1^+][\lambda_1^+ - \rho_-(x)]} \\ &\quad + \tau [\rho_+(x) - \lambda_1^+] \sqrt{\Gamma} \|K^{1/2}[A(\rho_+(x) + \lambda_1^+) + B]A^{-1/2}\|_2 \\ &= \left[ \varepsilon \sqrt{\lambda_1^+ - \rho_-(x)} + \tau \sqrt{\Gamma} \delta_1 \right] \sqrt{\rho_+(x) - \lambda_1^+},\end{aligned}\tag{7.12}$$

$$\begin{aligned}\|z\|_A &\geq \|x\|_A - \tau \|K r_+(x)\|_A \\ &= 1 - \tau \|K r_+(x)\|_A, \\ \|K r_+(x)\|_A &= \|K \mathbf{Q}(\lambda_1^+)x - K [\mathbf{Q}(\lambda_1^+) - \mathbf{Q}(\rho_+(x))]x\|_A \\ &\leq \|K \mathbf{Q}(\lambda_1^+)x\|_A + [\rho_+(x) - \lambda_1^+] \|K [A(\rho_+(x) + \lambda_1^+) + B]x\|_A \\ &\leq \sqrt{\|K^{1/2}AK^{1/2}\|_2 \Gamma} \|x\|_M \\ &\quad + [\rho_+(x) - \lambda_1^+] \|A^{1/2}K [A(\rho_+(x) + \lambda_1^+) + B]A^{-1/2}\|_2 \|x\|_A \\ &= \sqrt{\Gamma} \delta_2 + \delta_3^2.\end{aligned}\tag{7.13}$$

Finally use

$$\rho_+(z) - \lambda_1^+ = \frac{\|z\|_M^2}{[\lambda_1^+ - \rho_-(z)]\|z\|_A^2} \leq \frac{\|z\|_M^2}{[\lambda_1^+ - \rho_-(z)] \cdot [1 - \tau\|Kr_+(x)\|_A]^2}$$

and (7.12) and (7.13) to complete the proof.  $\square$

Similarly, we have the following result for Algorithm 7.1.

**Theorem 7.3.** *Suppose  $K \succ 0$ . Let  $\ell \in \{1, n\}$  and  $\text{typ}, \text{typ}' \in \{+, -\}$  such that  $\text{typ}$  and  $\text{typ}'$  are opposite, and let  $\gamma$  and  $\Gamma$  be the ones in Theorem 7.2, and*

$$\varepsilon = 2 \left[ \left( \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^{m-1} + \left( \frac{\sqrt{\kappa} + 1}{\sqrt{\kappa} - 1} \right)^{-(m-1)} \right]^{-1}, \quad \kappa = \frac{\Gamma}{\gamma}.$$

Let  $\text{argopt}$  be as given in (6.6), and

$$\begin{aligned} g_{\text{opt}} &= \underset{g \in \mathbb{P}_{m-1}, g(0)=1}{\text{argopt}} \rho_{\text{typ}}(g(K\mathbf{Q}(\rho_{\text{typ}}(x))x)), \\ y &= g_{\text{opt}}(K\mathbf{Q}(\rho_{\text{typ}}(x))x), \\ z &= \hat{g}(K\mathbf{Q}(\rho_{\text{typ}}(x))x), \end{aligned}$$

where

$$\begin{aligned} \hat{g}(t) &= \mathcal{T}_{m-1} \left( \frac{2t - (\Gamma + \gamma)}{\Gamma - \gamma} \right) / \mathcal{T}_{m-1} \left( -\frac{1 + \kappa}{1 - \kappa} \right) \\ &= 1 + c_1 t + \dots + c_{m-1} t^{m-1} \end{aligned}$$

since  $\hat{g}(0) = 1$ , and  $\mathcal{T}_{m-1}(t)$  is the  $(m-1)$ st Chebyshev polynomial of the first kind. We have

$$\begin{aligned} |\rho_{\text{typ}}(y) - \lambda_\ell^{\text{typ}}| &\leq |\rho_{\text{typ}}(z) - \lambda_\ell^{\text{typ}}| \\ &\leq \frac{1}{|\lambda_\ell^{\text{typ}} - \rho_{\text{typ}'}(z_{\text{typ}'})|} \left[ \frac{\varepsilon \sqrt{|\lambda_\ell^{\text{typ}} - \rho_{\text{typ}'}(x)|} + \eta \sqrt{|\rho_{\text{typ}}(x) - \lambda_\ell^{\text{typ}}|}}{1 - \eta |\rho_{\text{typ}}(x) - \lambda_\ell^{\text{typ}}|} \right]^2 \\ &\quad \times |\rho_{\text{typ}}(x) - \lambda_\ell^{\text{typ}}|, \end{aligned} \tag{7.14}$$

provided

$$\begin{aligned} \eta &:= \sum_{i=1}^{m-1} |c_i| \cdot \|K[A\rho_+(x) + \lambda_1^+] + B\|_2 \sum_{j=0}^{i-1} \|K\mathbf{Q}(\lambda_1^+)\|_2^{i-j-1} \|K\mathbf{Q}(\rho_+(x))\|^j \\ &< \frac{1}{|\rho_{\text{typ}}(x) - \lambda_\ell^{\text{typ}}|}. \end{aligned}$$

*Proof.* We will prove the case:  $(\text{typ}, \ell) = (+, 1)$  only. The other cases can be handled in the same way.

We have  $\lambda_1^+ \leq \rho_+(y) \leq \rho_+(z)$  and thus  $\rho_+(y) - \lambda_1^+ \leq \rho_+(z) - \lambda_1^+$ . So it suffices to show that  $\rho_+(z) - \lambda_1^+$  is no bigger than the right-hand side of (7.14).

Let  $M = -\mathbf{Q}(\lambda_1^+) \succeq 0$ . For any vector  $w$ , we have

$$\begin{aligned} \|w\|_M^2 &= -w^H \mathbf{Q}(\lambda_1^+) w \\ &= [\rho_+(w) - \lambda_1^+][\lambda_1^+ - \rho_-(w)] \|w\|_A^2, \end{aligned} \quad (7.15)$$

$$\|\hat{g}(-K\mathbf{Q}(\lambda_1^+))w\|_M \leq \max_{\gamma \leq t \leq T} |\hat{g}(\sigma)| \|w\|_M = \varepsilon \|w\|_M. \quad (7.16)$$

Write

$$z = \hat{g}(-K\mathbf{Q}(\lambda_1^+))x - \sum_{i=1}^{m-1} (-1)^i c_i \{ [K\mathbf{Q}(\lambda_1^+)]^i - [K\mathbf{Q}(\rho_+(x))]^i \} x.$$

Note that

$$\begin{aligned} [K\mathbf{Q}(\lambda_1^+)]^i - [K\mathbf{Q}(\rho_+(x))]^i &= \sum_{j=0}^{i-1} \left\{ [K\mathbf{Q}(\lambda_1^+)]^{i-j} [K\mathbf{Q}(\rho_+(x))]^j \right. \\ &\quad \left. - [K\mathbf{Q}(\lambda_1^+)]^{i-j-1} [K\mathbf{Q}(\rho_+(x))]^{j+1} \right\} \\ &= \sum_{j=0}^{i-1} [K\mathbf{Q}(\lambda_1^+)]^{i-j-1} [K\mathbf{Q}(\lambda_1^+) - K\mathbf{Q}(\rho_+(x))] [K\mathbf{Q}(\rho_+(x))]^j. \end{aligned}$$

Therefore

$$\begin{aligned} \|[K\mathbf{Q}(\lambda_1^+)]^i - [K\mathbf{Q}(\rho_+(x))]^i\|_2 &\leq \xi_i \|K\mathbf{Q}(\lambda_1^+) - K\mathbf{Q}(\rho_+(x))\|_2 \\ &\leq \xi_i (\rho_+(x) - \lambda_1^+) \|K[A\rho_+(x) + \lambda_1^+] + B\|_2, \end{aligned}$$

where  $\xi_i = \sum_{j=0}^{i-1} \|K\mathbf{Q}(\lambda_1^+)\|_2^{i-j-1} \|K\mathbf{Q}(\rho_+(x))\|_2^j$ . Without loss of generality, we may assume  $\|x\|_A = 1$ . We have

$$\begin{aligned} \|z\|_M &= \sqrt{[\rho_+(z) - \lambda_1^+][\lambda_1^+ - \rho_-(z)]} \|z\|_A, \quad \text{by (7.15)} \\ \|z\|_M &\leq \varepsilon \|x\|_M + \eta [\rho_+(x) - \lambda_1^+] \\ &= \left( \varepsilon \sqrt{\lambda_1^+ - \rho_-(x)} + \eta \sqrt{\rho_+(x) - \lambda_1^+} \right) \sqrt{\rho_+(x) - \lambda_1^+}, \end{aligned} \quad (7.17)$$

$$\begin{aligned} \|z\|_A &\geq \|x\|_A - \sum_{i=1}^{m-1} c_i \|[K\mathbf{Q}(\lambda_1^+)]^i - [K\mathbf{Q}(\rho_+(x))]^i\|_2 \|x\|_A \\ &\geq 1 - \eta (\rho_+(x) - \lambda_1^+), \end{aligned} \quad (7.18)$$

where  $\eta = \sum_{i=1}^{m-1} |c_i| \xi_i \|K[A\rho_+(x) + \lambda_1^+] + B\|_2$ . Finally use

$$\rho_+(z) - \lambda_1^+ = \frac{\|z\|_M^2}{[\lambda_1^+ - \rho_-(z)] \|z\|_A^2}$$

and (7.17) and (7.18) to complete the proof.  $\square$

## 8 Block preconditioned steepest descent/ascent method

The convergence of any of the previous steepest descent/ascent methods can be very slow if  $\lambda_1^\pm \approx \lambda_2^\pm$  or  $\lambda_{n-1}^\pm \approx \lambda_n^\pm$ . This is reflected by  $\omega_1 \approx \omega_2$  in Theorem 6.2 and 7.1. Often in practice, there are needs to compute the first few extreme quadratic eigenpairs, not just the first one. For that purpose, block variations of the methods become particularly attractive for at least the following reasons:

1. they can simultaneously compute the first  $k$  extreme quadratic eigenpairs  $(\lambda_j^\pm, u_j^\pm)$ ;
2. they run more efficiently on modern computer architecture because more computations can be organized into the matrix-matrix multiplication type;
3. they have better rates of convergence to the desired eigenpairs and 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 [55].

In what follows, we will explain a block steepest descent/ascent method for computing the first few  $(\lambda_j^\pm, u_j^\pm)$ . The same reasoning applies to other extreme quadratic eigenpairs.

Any block variation starts with a given  $X_0 \in \mathbb{C}^{n \times n_b}$  with  $\text{rank}(X_0) = n_b$ , instead of just one vector  $\mathbf{x}_0 \in \mathbb{C}^n$  previously for the single-vector steepest descent type methods. Here either the  $j$ th column of  $X_0$  is already an approximation to  $u_j^\pm$  or the subspace  $\mathcal{R}(X_0)$  is a good approximation to the subspace spanned by  $u_j^\pm$  for  $1 \leq j \leq n_b$  or the canonical angles from  $\mathcal{R}([u_1^\pm, \dots, u_k^\pm])$  to  $\mathcal{R}(X_0)$  are nontrivial, where  $k \leq n_b$  is the number of desired eigenpairs. In the latter two cases, a preprocessing is needed to turn the case into the first case:

1. solve the HQEP  $X_0^H \mathbf{Q}(\lambda) X_0$  to get its pos-type quadratic eigenpairs  $(\boldsymbol{\rho}_{0;j}^\pm, y_j^\pm)$ ;
2. reset  $X_0 := X_0[y_1^\pm, \dots, y_{n_b}^\pm]$ .

So we will assume henceforth the  $j$ th column of the given  $X_0$  is an approximation to  $u_j^\pm$ . Now consider generalizing the steepest descent method to a block version. Its typical  $i$ th iterative step may well look like the following. Suppose we have already computed

$$X_i = [x_{i;1}, x_{i;2}, \dots, x_{i;n_b}] \in \mathbb{C}^{n \times n_b}$$

whose  $j$ th column  $x_{i;j}$  approximates  $u_j^\pm$  and

$$\Omega_i = \text{diag}(\boldsymbol{\rho}_{i;1}^\pm, \boldsymbol{\rho}_{i;2}^\pm, \dots, \boldsymbol{\rho}_{i;n_b}^\pm)$$

whose  $j$ th diagonal entry  $\boldsymbol{\rho}_{i;j}^\pm = \rho_+(x_{i;j})$  approximates  $\lambda_j^\pm$ . Define the residual matrix

$$R_i = [r_+(x_{i;1}), r_+(x_{i;2}), \dots, r_+(x_{i;n_b})] = AX_i \Omega_i^2 + BX_i \Omega_i + CX_i.$$

The next set of approximations are computed as follows:

1. compute a basis matrix  $Z$  of  $\mathcal{R}([X_i, R_i])$  by, e.g., MGS;



2. solve the QEP  $Z^H \mathbf{Q}(\lambda) Z$  to get its pos-type quadratic eigenpairs  $(\boldsymbol{\rho}_{i+1;j}^+, y_j^+)$ , and let  $\Omega_{i+1} = \text{diag}(\boldsymbol{\rho}_{i+1;1}^+, \boldsymbol{\rho}_{i+1;2}^+, \dots, \boldsymbol{\rho}_{i+1;n_b}^+)$ ;
3. set  $X_{i+1} = Z[y_1^+, \dots, y_{n_b}^+]$ .

In the same way as we explained before, this block steepest descent method can be improved in two directions – extending the search space is one and incorporating preconditioners is the other.

Note that  $r_+(x_{i;j}) = \mathbf{Q}(\boldsymbol{\rho}_{i;j}^+)x_{i;j}$  and thus

$$\begin{aligned} \mathcal{R}([X_i, R_i]) &= \sum_{j=1}^{n_b} \mathcal{R}([x_{i;j}, \mathbf{Q}(\boldsymbol{\rho}_{i;j}^+)x_{i;j}]) \\ &= \sum_{j=1}^{n_b} \mathcal{K}_2(\mathbf{Q}(\boldsymbol{\rho}_{i;j}^+), x_{i;j}). \end{aligned}$$

So it is natural to extend the search space,  $\mathcal{R}([X_\ell, R_\ell])$  through extending each Krylov subspace  $\mathcal{K}_2(\mathbf{Q}(\boldsymbol{\rho}_{i;j}^+), x_{i;j})$  to a high order one, and of course different Krylov subspaces can be extended to different orders. For simplicity, we will extend each to the  $m$ th order. The new extended search subspace now is

$$\sum_{j=1}^{n_b} \mathcal{K}_m(\mathbf{Q}(\boldsymbol{\rho}_{i;j}^+), x_{i;j}). \quad (8.1)$$

Define the linear operator

$$\mathcal{R}_i : X \in \mathbb{C}^{n \times n_b} \rightarrow \mathcal{R}_i(X) = AX\Omega_i^2 + BX\Omega_i + CX \in \mathbb{C}^{n \times n_b}.$$

Then the subspace in (8.1) can be compactly written as

$$\mathcal{K}_m(\mathcal{R}_i, X_i) = \text{span}\{X_i, \mathcal{R}_i(X_i), \dots, \mathcal{R}_i^{m-1}(X_i)\}, \quad (8.2)$$

where  $\mathcal{R}_i^j(\cdot)$  is understood as successively applying the operator  $\mathcal{R}_i$   $j$  times, e.g.,  $\mathcal{R}_i^2(X) = \mathcal{R}_i(\mathcal{R}_i(X))$ .

As to incorporate suitable preconditioners, in light of our extensive discussions in subsection 7.1, the search subspace should be modified to

$$\sum_{j=1}^{n_b} \mathcal{K}_m(K_{i;j}\mathbf{Q}(\boldsymbol{\rho}_{i;j}^+), x_{i;j}), \quad (8.3)$$

where  $K_{i;j}$  are the preconditioners, one for each approximate eigenpair  $(\boldsymbol{\rho}_{i;j}^+, x_{i;j})$  for  $1 \leq j \leq n_b$  in the  $i$ th iterative step. As before,  $K_{i;j}$  can be constructed in one of the following two ways:

- $K_{i;j}$  is an approximate inverse of  $\mathbf{Q}(\tilde{\boldsymbol{\rho}}_{i;j}^+)$  for some  $\tilde{\boldsymbol{\rho}}_{i;j}^+$  different from  $\boldsymbol{\rho}_{i;j}^+$ , ideally closer to  $\lambda_j^+$  than to any other quadratic eigenvalue of  $\mathbf{Q}(\lambda)$ . But this requirement on  $\tilde{\boldsymbol{\rho}}_{i;j}^+$  is impractical because the quadratic eigenvalue  $\lambda_j^+$  of  $\mathbf{Q}(\lambda)$  is unknown. A compromise would be to make  $\tilde{\boldsymbol{\rho}}_{i;j}^+$  closer but not equal to  $\boldsymbol{\rho}_{i;j}^+$  than to any other  $\boldsymbol{\rho}_{i;j}^+$ .

---

**Algorithm 8.1** Block preconditioned extended steepest descent/ascent 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 computes approximate quadratic eigenpairs to  $(\lambda_j^{\text{typ}}, u_j^{\text{typ}})$  for  $j \in \mathbb{J}$ , where  $\mathbb{J} = \{1 \leq j \leq n_b\}$  for computing the few smallest quadratic eigenpairs of the given type or  $\{n - n_b + 1 \leq j \leq n\}$  for computing the few largest quadratic eigenpairs of the given type.

---

- 1: solve the HQEP  $X_0^H \mathbf{Q}(\lambda) X_0$  to get its quadratic eigenpairs  $(\rho_{0;j}^{\text{typ}}, y_j^{\text{typ}})$ ;
- 2:  $X_0 = X_0[y_1^{\text{typ}}, \dots, y_{n_b}^{\text{typ}}]$ ,  $\hat{\mathbb{J}} = \{1 \leq j \leq n_b\}$ ;
- 3: **for**  $i = 0, 1, \dots$  **do**
- 4:   construct preconditioners  $K_{i;j}$  for  $j \in \hat{\mathbb{J}}$ ;
- 5:   compute a basis matrix  $Z$  of the subspace

$$\sum_{j \in \hat{\mathbb{J}}} \mathcal{K}_m(K_{i;j} \mathbf{Q}(\rho_{i;j}^{\text{typ}}), x_{i;j}), \quad (8.4)$$

and let  $n_Z$  be its dimension and  $\hat{\mathbb{J}} = \{1 \leq j \leq n_b\}$  for computing the few smallest quadratic eigenpairs of the given type or  $\{n_Z - n_b + 1 \leq j \leq n_Z\}$  for computing the few largest quadratic eigenpairs of the given type;

- 6:   compute the  $n_b$  quadratic eigenpairs of  $Z^H \mathbf{Q}(\lambda) Z$ :  $(\rho_{i+1;j}^{\text{typ}}, y_j^{\text{typ}})$  for  $j \in \hat{\mathbb{J}}$  and let  $\Omega_{i+1} = \text{diag}(\dots, \rho_{i+1;j}^{\text{typ}}, \dots)$  whose diagonal entries are those for  $j \in \hat{\mathbb{J}}$ ;
  - 7:    $X_{i+1} = ZW$ , where  $W = [\dots, y_j^{\text{typ}}, \dots]$  whose columns are those for  $j \in \hat{\mathbb{J}}$ ;
  - 8: **end for**
  - 9: **return** approximate quadratic eigenpairs to  $(\lambda_j^{\text{typ}}, u_j^{\text{typ}})$  for  $j \in \mathbb{J}$ .
- 

- Perform an incomplete  $LDL^H$  factorization (see [51, Chapter 10])  $\mathbf{Q}(\tilde{\rho}_{i;j}^+) \approx L_{i;j} D_{i;j} L_{i;j}^H$ , where “ $\approx$ ” includes not only the usual “approximately equal”, but also the case when  $\mathbf{Q}(\tilde{\rho}_{i;j}^+) - L_{i;j} D_{i;j} L_{i;j}^H$  is approximately a low rank matrix, and  $D_{i;j} = \text{diag}(\pm 1)$ . Finally set  $K_{i;j} = L_{i;j} L_{i;j}^H$ .

Algorithm 8.1 is the general framework of a Block Preconditioned Extended Steepest Descent method (BPeSD) which embeds four methods into one:

1. Block Steepest Descent method:  $m = 2$  and all preconditioners  $K_{i;j} = I$ ;
2. Block Preconditioned Steepest Descent method:  $m = 2$  and nontrivial  $K_{i;j}$ ;
3. Block Extended Steepest Descent method:  $m > 2$  and all preconditioners  $K_{i;j} = I$ ;
4. Block Preconditioned Extended Steepest Descent method:  $m > 2$  and nontrivial  $K_{i;j}$ .

There are three important implementation issues to worry about in turning this general framework into a piece of working code.

1. In (8.3), a different preconditioner is used for each and every approximate eigenpair  $(\rho_{i;j}^+, x_{i;j})$  for  $1 \leq j \leq n_b$ . While, conceivably, doing so will speed up convergence for each approximate eigenpair because each preconditioner can be constructed to make that

approximate eigenpair converge faster, but the cost in constructing these preconditioners may likely be too heavy to bear. A more practical approach would be to use one preconditioner  $K_i$  for all  $K_{i;j}$  aiming at speeding up the convergence of  $(\rho_{i;1}^+, x_{i;1})$  (or the first few approximate quadratic eigenpairs for tightly clustered quadratic eigenvalues). Once it (or the first few in the case of a tightly cluster) is determined to be sufficiently accurate, the converged eigenpairs are locked up and deflated and a new preconditioner is computed to aim at the next non-converged eigenpairs, and the process continues.

**2.** Consider implementing Line 5, i.e., generating a basis matrix for the subspace (8.4). In the most general case,  $Z$  can be gotten by packing the basis matrices of all

$$\mathcal{K}_m(K_{i;j}\mathbf{Q}(\rho_{i;j}^+), x_{\ell;j}) \quad \text{for } 1 \leq j \leq n_b$$

together. There could be two problems with this: 1) such  $Z$  could be ill-conditioned, i.e., the columns of  $Z$  may not be sufficiently numerically linearly independent, and 2) the arithmetic operations in building a basis for each  $\mathcal{K}_m(K_{i;j}\mathbf{Q}(\rho_{i;j}^+), x_{i;j})$  are mostly matrix-vector multiplications, straying from one of the purposes: performing most arithmetic operations through matrix-matrix multiplications in order to achieve high performance on modern computers. To address these two problems, we may do a tradeoff by using  $K_{i;j} \equiv K_i$  for all  $j$ . This may likely degrade the effectiveness of the preconditioner per step in terms of rates of convergence for all approximate eigenpairs  $(\rho_{i;j}^+, x_{i;j})$  but may achieve overall gain in using less time because then the code will run much faster in matrix-matrix operations, not to mention the saving in constructing just one preconditioner  $K_i$  instead of  $n_b$  different preconditioners  $K_{i;j}$ . To simplify our discussion below, we will drop the subscript  $i$  for readability. Since  $K_{i;j} \equiv K$  for all  $j$ , (8.4) is the same as

$$\mathcal{K}_m(K\mathcal{R}, X) = \text{span}\{X, K\mathcal{R}(X), \dots, [K\mathcal{R}]^{m-1}(X)\}, \quad (8.5)$$

where  $[K\mathcal{R}]^j(\cdot)$  is understood as successively applying the operator  $K\mathcal{R}$   $j$  times, e.g.,  $[K\mathcal{R}]^2(X) = K\mathcal{R}_\ell(K\mathcal{R}(X))$ . A basis matrix

$$Z = [Z_1, Z_2, \dots, Z_m]$$

can be computed by the following block Arnoldi-like process.

---

```

1:  $Z_1 T = X$  (MGS);
2: for  $i = 2$  to  $m$  do
3:    $Y = K(AZ_{i-1}\Omega^2 + BZ_{i-1}\Omega + CZ_{i-1})$ ;
4:   for  $j = 1$  to  $i - 1$  do
5:      $T = Z_j^H Y$ ;  $Y = Y - Z_j T$ ;
6:   end for
7:    $Z_i T = Y$  (MGS);
8: end for

```

---

There is a possibility that at Line 7  $Y$  is numerically not of full column rank. If it happens, it poses no difficulty at all. In running MGS on  $Y$ 's columns, anytime if a column is deemed linearly dependent on previous columns, that column should be deleted, along with the corresponding  $\rho_j^+$  from  $\Omega$  to shrink its size by 1 as well. At the completion of MGS,  $Z_i$  will have fewer columns than  $Y$  and the size of  $\Omega$  is shrunk accordingly. Finally, at the

end, the columns of  $Z$  are orthonormal, i.e.,  $Z^H Z = I$  (of apt size) which may fail to an unacceptably level due to roundoff; so some form of re-orthogonalization should be incorporated.

---

**Algorithm 9.1** Preconditioned conjugate gradient method

---

Given an initial approximation  $\mathbf{x}_0$  to  $u_\ell^{\text{typ}}$ , a (positive definite) preconditioner  $K$ , and a relative tolerance `rtol`, the algorithm computes an approximate pair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$  with the prescribed `rtol`.

---

```
1:  $\mathbf{x}_0 = \mathbf{x}_0 / \|\mathbf{x}_0\|_2$ ,  $\rho_0 = \rho_{\text{typ}}(\mathbf{x}_0)$ ,  $\mathbf{r}_0 = r_{\text{typ}}(\mathbf{x}_0)$ ,  $\mathbf{p}_0 = -K\mathbf{r}_0$ ;  
2: for  $i = 0, 1, \dots$  do  
3:   if  $\|\mathbf{r}_i\|_2 / (|\rho_i|^2 \|\mathbf{A}\mathbf{x}_i\| + |\rho_i| \|\mathbf{B}\mathbf{x}_i\| + \|\mathbf{C}\mathbf{x}_i\|) \leq \text{rtol}$  then  
4:     BREAK;  
5:   else  
6:     solve the HQEP for  $Y_i^H \mathbf{Q}(\lambda) Y_i$ , where  $Y_i = [\mathbf{x}_i, \mathbf{p}_i]$  to get its quadratic eigenvalues  
7:      $\mu_j^\pm$  as in (6.8) and quadratic eigenvectors  $y_j^\pm$ ;  
8:     select the next approximate quadratic eigenpair  $(\mu, Y_i v)$  according to the table  
9:     (6.9);  
10:    compute  $\alpha_i = t_{\text{opt}}$  as in (9.2) and then  $y$  as in (6.7) with  $x = \mathbf{x}_i$  and  $p = \mathbf{p}_i$ ;  
11:     $\mathbf{x}_{i+1} = y / \|y\|_2$ ;  
12:    set  $\rho_{i+1} = \rho_{\text{typ}}(\mathbf{x}_{i+1})$ ,  $\mathbf{r}_{i+1} = r_{\text{typ}}(\mathbf{x}_{i+1})$ ,  $\mathbf{p}_{i+1} = -K\mathbf{r}_{i+1} + \beta_i \mathbf{p}_i$ , where  $\beta_i$  is  
13:    commonly given by either one of
```

$$\text{either } \beta_i = \frac{\mathbf{r}_{i+1}^H K \mathbf{r}_{i+1}}{\mathbf{r}_i^H K \mathbf{r}_i} \quad \text{or} \quad \beta_i = \frac{\mathbf{r}_{i+1}^H K (\mathbf{r}_{i+1} - \mathbf{r}_i)}{\mathbf{r}_i^H K \mathbf{r}_i}; \quad (9.1)$$

```
14:  end if  
15: end for  
16: return  $(\rho_i, \mathbf{x}_i)$  as an approximate eigenpair to  $(\lambda_\ell^{\text{typ}}, u_\ell^{\text{typ}})$ .
```

---

## 9 Conjugate gradient method

Again because of the equations in (3.8), the nonlinear CG type method [45, 59] and its variations are natural candidates for computing the first or last quadratic eigenpair  $(\lambda_j^\pm, u_j^\pm)$ , and their block variations can also be devised to simultaneously compute the first or last few quadratic eigenpairs  $(\lambda_j^\pm, u_j^\pm)$ . Since much of the machinery including gradients and preconditioning has already been built up, what remain are more or less simple adaptations of CG type methods [35] for the generalized Hermitian eigenvalue problem to the current case.

### 9.1 Preconditioned conjugate gradient method

The single-vector CG type methods heavily rely on the line-search problem (6.5) – (6.7) which was solved by projecting the original  $n \times n$  HQEP for  $\mathbf{Q}(\lambda)$  to a  $2 \times 2$  HQEP  $Y^H \mathbf{Q}(\lambda) Y$  without actually computing the optimal parameter  $t_{\text{opt}}$  and thus the next approximation  $y$  as in (6.7) for the steepest descent/ascent method and its variations. The outcome of it is that the computed next approximation is a (complex) scalar multiply of  $y$  in (6.7). This is good enough for the steepest descent/ascent method but not for the CG method for which  $y$  in (6.7) needs to be computed. We now show how this  $y$  can be recovered from the approximation given in the table (6.9). Let  $(\mu, Yv)$  is selected according to the

---

**Algorithm 9.2** Locally optimal block preconditioned extended conjugate gradient 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 computes approximate eigenpairs to  $(\lambda_j^{\text{typ}}, u_j^{\text{typ}})$  for  $j \in \mathbb{J}$ , where  $\mathbb{J} = \{1 \leq j \leq n_b\}$  for computing the few smallest quadratic eigenpairs of the given type or  $\{n - n_b + 1 \leq j \leq n\}$  for computing the few largest quadratic eigenpairs of the given type.

---

- 1: solve the HQEP  $X_0^H \mathbf{Q}(\lambda) X_0$  to get its quadratic eigenpairs  $(\rho_{0;j}^{\text{typ}}, y_j^{\text{typ}})$ ;
- 2:  $X_0 = X_0[y_1^{\text{typ}}, \dots, y_{n_b}^{\text{typ}}]$ ,  $X_{-1} = 0$ ,  $\hat{\mathbb{J}} = \{1 \leq j \leq n_b\}$ ;
- 3: **for**  $i = 0, 1, \dots$  **do**
- 4:   construct preconditioners  $K_{i;j}$  for  $j \in \hat{\mathbb{J}}$ ;
- 5:   compute a basis matrix  $Z$  of the subspace

$$\sum_{j \in \hat{\mathbb{J}}} \mathcal{K}_m(K_{i;j} \mathbf{Q}(\rho_{i;j}), x_{i;j}) + \mathcal{R}(X_{i-1}), \quad (9.3)$$

and let  $n_Z$  be its dimension and  $\hat{\mathbb{J}} = \{1 \leq j \leq n_b\}$  for computing the few smallest quadratic eigenpairs of the given type or  $\{n_Z - n_b + 1 \leq j \leq n_Z\}$  for computing the few largest quadratic eigenpairs of the given type;

- 6:   compute the  $n_b$  quadratic eigenpairs of  $Z^H \mathbf{Q}(\lambda) Z$ :  $(\rho_{i+1;j}^{\text{typ}}, y_j^{\text{typ}})$  for  $j \in \hat{\mathbb{J}}$  and let  $\Omega_{i+1} = \text{diag}(\dots, \rho_{i+1;j}^{\text{typ}}, \dots)$  whose diagonal entries are those for  $j \in \hat{\mathbb{J}}$ ;
  - 7:    $X_{i+1} = ZW$ , where  $W = [\dots, y_j^{\text{typ}}, \dots]$  whose columns are those for  $j \in \hat{\mathbb{J}}$ ;
  - 8: **end for**
  - 9: **return** approximate quadratic eigenpairs to  $(\lambda_j^{\text{typ}}, u_j^{\text{typ}})$  for  $j \in \mathbb{J}$ .
- 

table, and write  $v = [\nu_1, \nu_2]^T$  and  $\hat{y} = Yv = \nu_1 x + \nu_2 p$ . Thus

$$t_{\text{opt}} = \nu_2 / \nu_1 \text{ if } \nu_1 \neq 0, \text{ and } \infty \text{ otherwise.} \quad (9.2)$$

With this, set  $y$  as in (6.7).

Our discussions on selecting a good preconditioner in subsection 7.1 should be followed. Algorithm 9.1 presents the framework for the single-vector preconditioned conjugate gradient method for  $\mathbf{Q}(\lambda)$ .

## 9.2 Locally optimal block preconditioned extended conjugate gradient method

When it comes to eigenvalue computations by CG type methods, CG's locally optimal variations [48, 60] combined with preconditioning and blocking are more preferable than the usual single-vector CG method as in Algorithm 9.1 [3, 28, 35]. In Algorithm 9.2, we present a framework of the so-called *Locally Optimal Block Preconditioned Extended Conjugate Gradient Method* (LOBPeCG) whose different implementation choice gives rise to a list of CG-type methods which we will elaborate.

The three important implementation issues we discussed for Algorithm 8.1 (Block Preconditioned Extended Steepest Descent method) after its introduction essentially apply here, except some changes are needed in the computation of  $Z$  at Line 5 here.

First  $X_{i-1}$  can be replaced by something else. Specifically, we modify Lines 2, 6, and 8 of Algorithm 9.2 to

- 
- 2:  $X_0 = X_0W$ , and  $Y_0 = 0$ ,  $\hat{\mathbb{J}} = \{1 \leq j \leq n_b\}$ ;  
5: compute a basis matrix  $Z$  of the subspace

$$\sum_{j \in \hat{\mathbb{J}}} \mathcal{K}_m(K_{i;j} \mathbf{Q}(\boldsymbol{\rho}_{i;j}), x_{i;j}) + \mathcal{R}(Y_i), \quad (9.4)$$

such that  $\mathcal{R}(Z_{(:,1:n_b)}) = \mathcal{R}(X_i)$ . Let  $n_Z$  be its dimension and  $\hat{\mathbb{J}} = \{1 \leq j \leq n_b\}$  for computing the few smallest quadratic eigenpairs of the given type or  $\{n_Z - n_b + 1 \leq j \leq n_Z\}$  for computing the few largest quadratic eigenpairs of the given type;

- 7:  $X_{i+1} = ZW$ , where  $W = [\dots, y_j^{\text{typ}}, \dots]$  whose columns are those for  $j \in \hat{\mathbb{J}}$ ,  
 $Y_{i+1} = Z_{(:,n_b+1:(m+1)n_b)} W_{(n_b+1:(m+1)n_b,:)}$ ;
- 

Next we will compute a basis matrix for the subspace (9.3) or (9.4). For better performance (by using more matrix-matrix multiplications), we will assume  $K_{i;j} \equiv K_i$  for all  $j$  for simplification. Dropping the subscript  $i$  for readability, we see (9.4) is the same as

$$\mathcal{K}_m(K\mathcal{R}, X) + \mathcal{R}(Y) = \text{span}\{X, K\mathcal{R}(X), \dots, [K\mathcal{R}]^{m-1}(X)\} + \mathcal{R}(Y). \quad (9.5)$$

We will first 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 outlined at the end of section 8. In particular,  $\mathcal{R}(Z_1) = \mathcal{R}(X)$ . Then orthogonalize  $Y$  against  $[Z_1, Z_2, \dots, Z_m]$  to get  $Z_{m+1}$  satisfying  $Z_{m+1}^H Z_{m+1} = I$ . Finally take  $Z = [Z_1, Z_2, \dots, Z_{m+1}]$ .

So far, we have not mentioned any convergence properties of these CG type methods.

## 10 Numerical examples

In this section, we will present a couple of examples to demonstrate the numerical behavior of Algorithm 9.2 which often performs much better than the steepest descent/ascent type methods. In presenting numerical results, we will use the normalized residuals

$$\frac{\|\mathbf{Q}(\mu_i)x_i\|_2}{(\|A\|_1\mu_i^2 + \|B\|_1|\mu_i| + \|C\|_1)\|x_i\|_2}$$

to show the convergent progress for approximations  $(\mu_i, x_i)$  to a particular quadratic eigen-pair *vs.* the iteration index, where using the matrix  $\ell_1$  operator norms  $\|A\|_1$ ,  $\|B\|_1$ , and  $\|C\|_1$  is more for computational convenience than anything else as any other norm would serve the same purpose just as well.

**Example 10.1.** This is the problem **Wiresaw1** in the collection [5]. It is actually a gyroscopic QEP arising in the vibration analysis of a wiresaw [68], but leads to an HQEP. Here

$$A = \frac{1}{2}I_n, \quad C = \frac{(\nu^2 - 1)\pi^2}{2} \text{diag}(1^2, 2^2, \dots, n^2),$$

$$B = \iota(b_{ij}) \quad \text{with} \quad b_{ij} = \begin{cases} \nu \frac{4ij}{i^2 - j^2}, & \text{if } i + j \text{ is odd,} \\ 0, & \text{otherwise,} \end{cases}$$

where  $\iota = \sqrt{-1}$  is the imaginary unit,  $\nu$  is a real nonnegative parameter corresponding to the speed of the wire. For  $0 < \nu < 1$ ,  $\mathbf{Q}(0) = C$  is negative definite, and thus  $\mathbf{Q}(\lambda) = \lambda^2 A + \lambda B + C$  is hyperbolic by Theorem 2.1. Moreover

$$\lambda_i^- < 0 < \lambda_j^+ \quad \text{for all } i, j.$$

Therefore it is rather natural to use  $K = -C^{-1}$  as a preconditioner when it comes to compute the few smallest  $\lambda_j^+$  or largest  $\lambda_i^-$ , or for testing purpose some approximations to  $C^{-1}$  such as those corresponding to the linear conjugate gradient methods.

We ran Algorithm 9.2 with  $n_b = 10$ ,  $m = 2$  and random  $X_0 = \text{randn}(n, n_b)$  on this example for  $n = 1,000$  and  $\nu = 0.8$  without or with preconditioners

$$K \approx \begin{cases} [\mathbf{Q}(\pm 6.0 \cdot 10^3)]^{-1}, & \text{for largest } \lambda_j^+ \text{ or smallest } \lambda_j^-, \\ -[\mathbf{Q}(0)]^{-1} = -C^{-1}, & \text{for smallest } \lambda_j^+ \text{ or largest } \lambda_j^-, \end{cases} \quad (10.1)$$

implemented through the linear conjugate gradient method with stopping criteria of normalized residuals for the involved linear systems being no bigger than  $10^{-1}$  or reaching the maximum number CG steps which is 10. We have already explained the use of  $-C^{-1}$  or its approximations as possible preconditioners. After running Algorithm 9.2 without any preconditioner, we noticed that all  $\lambda_j^\pm$  lie in  $(-6.0 \cdot 10^3, 6.0 \cdot 10^3)$  which leads to the use of  $[\mathbf{Q}(\pm 6.0 \cdot 10^3)]^{-1}$  in (10.1).

Figure 10.1 plots the residual history for computing the largest or smallest few  $\lambda_i^\pm$ , where the left column is for without any preconditioner while the right column is for with the preconditioners as given in (10.1). We notice without using any preconditioner



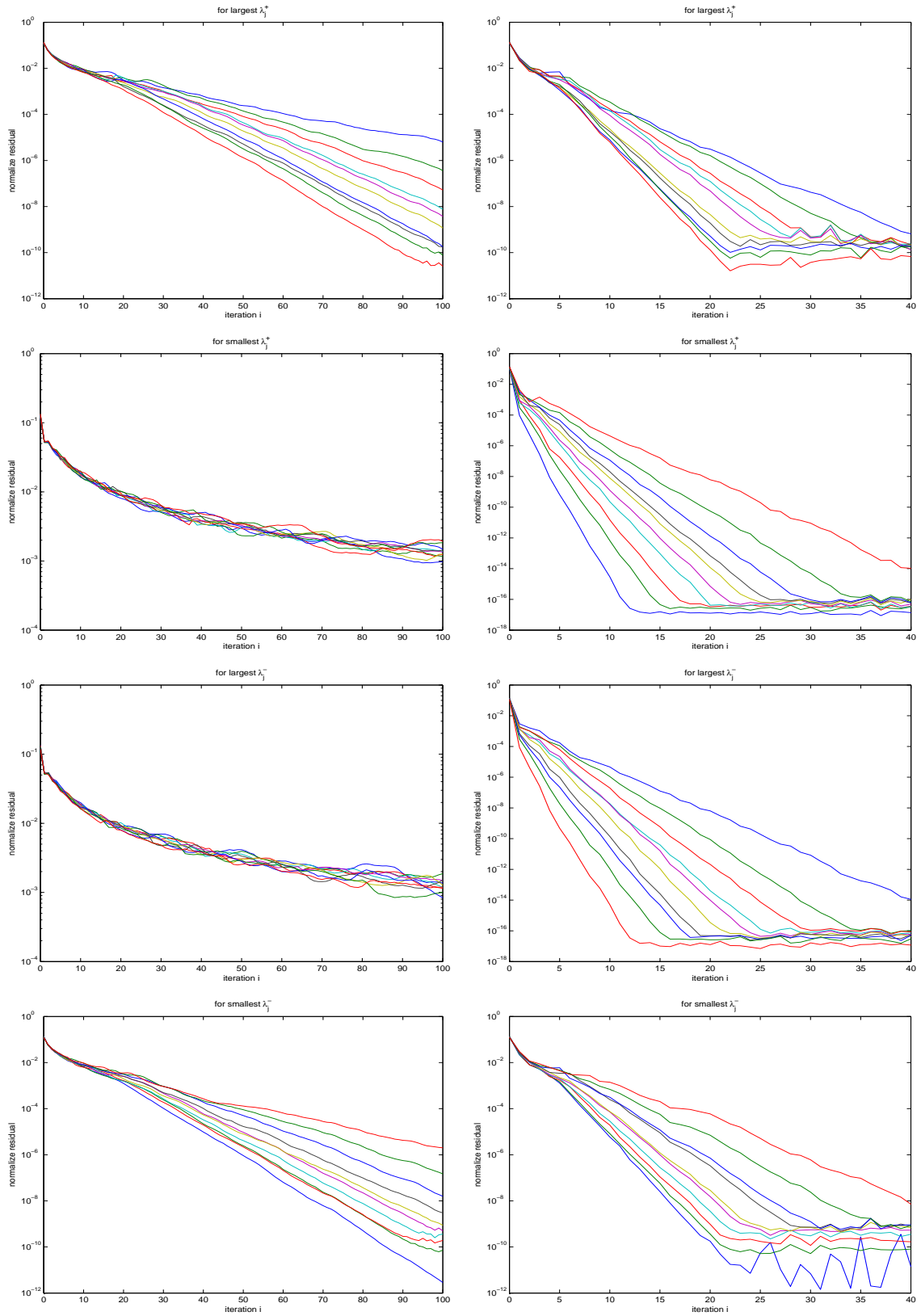


Figure 10.1: Residual history for running Algorithm 9.2 on Example 10.1

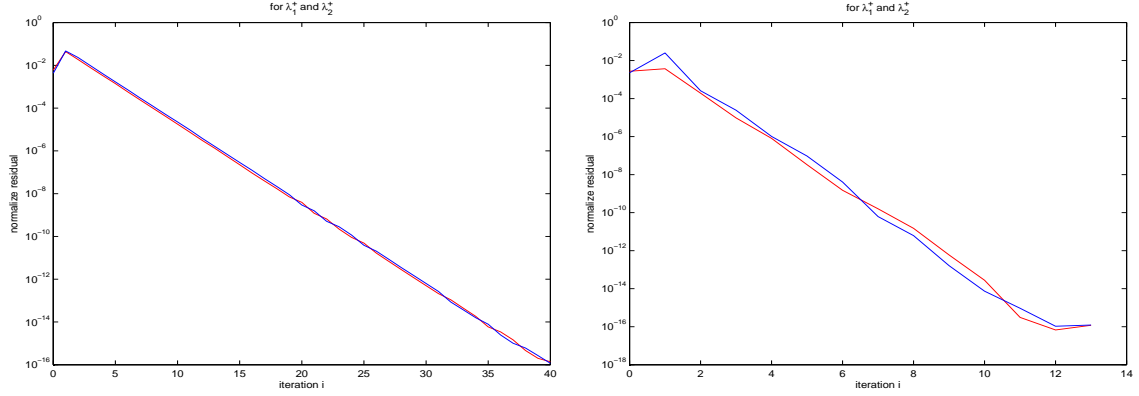


Figure 10.2: Residual history for running Algorithm 9.2 on Example 10.2 for computing  $\lambda_1^+$  and  $\lambda_2^+$

Algorithm 9.2 performed poorly for computing smallest  $\lambda_j^+$  or largest  $\lambda_j^-$  but reasonably well for largest  $\lambda_j^+$  or smallest  $\lambda_j^-$ . The effectiveness of the preconditioners as in (10.1) is rather evident by comparing the plots in the two columns.

**Example 10.2.** This is [20, Example 5], where  $A = I_n$ ,

$$B = \xi \begin{bmatrix} 20 & -10 & & & & & & \\ -10 & 30 & -10 & & & & & \\ & & \ddots & \ddots & \ddots & & & \\ & & & -10 & 30 & -10 & & \\ & & & & -10 & 20 & & \end{bmatrix}, C = \begin{bmatrix} 15 & -5 & & & & & & \\ -5 & 15 & -5 & & & & & \\ & & \ddots & \ddots & \ddots & & & \\ & & & -5 & 15 & -5 & & \\ & & & & -5 & 15 & & \end{bmatrix},$$

and  $\xi$  is a parameter. We take  $n = 1000$  and  $\xi = 1.1$ . This is a pathological example in the sense that most quadratic eigenvalues are close to one another – share about 3 significant decimal digits with their neighbors, except  $\lambda_1^+$  and  $\lambda_2^+$  which has a gap from the rest. When running Algorithm 9.2 with  $m = 2$  and various different  $n_b$ , we noticed the algorithm really had hard time computing all extreme  $\lambda_j^\pm$  even with some preconditioner  $K = \pm[Q(\mu)]^{-1}$  with  $\mu \in (\lambda_n^-, \lambda_1^+)$  or  $\mu > \lambda_n^+$  or  $\mu < \lambda_1^-$  purposely selected, except for  $\lambda_1^+$  and  $\lambda_2^+$  which are rather easy to compute actually. Figure 10.2 plots the residual history for computing  $\lambda_1^+$  and  $\lambda_2^+$ , where the left plot is for without any preconditioner while the right plot is for with a preconditioner  $K \approx [Q(-8.0)]^{-1}$  implemented through the linear conjugate gradient method with the same stopping criteria as in the previous example.

## 11 Concluding remarks

We have performed a systematic study of the hyperbolic quadratic eigenvalue problem  $Q(\lambda) = \lambda^2 A + \lambda B + C$ . Such a problem usually arises from dynamical systems with heavy friction. Such a system appears, for example, in an elevator or car braking system. It shares many characteristics with the standard Hermitian eigenvalue problem in the category of usual standard linear eigenvalue problems, and had attracted quite some attention in the past. Most of the results were collected in [16, 43, 65].

Our contributions in this paper lie in two fronts. Theoretically, we have established Amir-Moéz/Wielandt-Lidskii type min-max principles for the sums of selected quadratic eigenvalues and, as corollaries, trace min/max type principles, and also perturbation results in the spectral and Frobenius norm, as well as general unitarily invariant norms on how the quadratic eigenvalues will change if  $A$ ,  $B$ ,  $C$  are perturbed. Numerically, we have justified a naturally extended Rayleigh-Ritz type procedure, with the existing and newly established min-max principles, why the procedure will produce the best approximations to quadratic eigenvalues/eigenvectors, proposed steepest descent/ascent and CG type methods for computing extreme quadratic eigenpairs, and established convergence results, including the rate of convergence for the steepest descent/ascent methods, which shed light on preconditioning in what constitutes a good preconditioner and how to construct one.

Block steepest descent/ascent type methods often perform much better in practice than their single-vector counterparts, as they should be. But their exact rates of convergence are hard to establish. Experience shows that their corresponding locally optimal CG type methods perform even better, but then again we do not know the exact rates of convergence locally optimal CG type methods, either. It is recommended that locally optimal CG type methods should be preferred to their corresponding steepest descent/ascent type methods.

Despite many successes we have so far in this paper in extending, as many as we can, the important results (both theoretically and numerically) for the standard Hermitian eigenvalue problem, there are more to be done. We list a few here.

- We established perturbation bounds for quadratic eigenvalues, but didn't do so for quadratic eigenvectors/eigenspaces. The latter is worth investigating, too. We expect that  $\min_x \zeta_0(x)$  will play a role.
- Higham, Mackey, and Tisseur [23] expanded hyperbolic quadratic matrix polynomials to include the case when  $A$  is positive semidefinite. Conceivably, many results in this paper may be extensible to quadratic definite matrix polynomials in the sense of [23], but care must be taken to deal with infinite quadratic eigenvalues.
- Many results in this paper should be extensible to hyperbolic matrix polynomials of degrees higher than 2 [43]. We are working on it and results will be detailed in a separate paper.

## A Digression: positive semidefinite matrix pencil

Let  $\mathbf{A} - \lambda\mathbf{B}$  be a matrix pencil of order  $n$ , i.e.,  $\mathbf{A}, \mathbf{B} \in \mathbb{C}^{n \times n}$ .

**Definition A.1** ([38]).  $\mathbf{A} - \lambda\mathbf{B}$  is said *Hermitian* if both  $\mathbf{A}, \mathbf{B}$  are Hermitian, *positive (semi)definite* if it is Hermitian and there exists  $\lambda_0 \in \mathbb{R}$  such that  $\mathbf{A} - \lambda_0\mathbf{B} \succ 0$  ( $\mathbf{A} - \lambda_0\mathbf{B} \succeq 0$ ).

The concept of positive semidefinite pencil is closely related to that of the so-called *definite pencil* in the past literature [54, 57, 58]. The latter only requires that some linear combination (with possibly complex coefficients) is positive definite and thus is necessarily a regular pencil, i.e.,  $\det(\mathbf{A} - \lambda\mathbf{B}) \not\equiv 0$ . Definition A.1 uses more restrictive linear combinations, and also a positive semidefinite pencil of this definition may possibly be singular, i.e.,  $\det(\mathbf{A} - \lambda\mathbf{B}) \equiv 0$ .

To include, possibly, the case in which  $\mathbf{A} - \lambda\mathbf{B}$  is a singular pencil, we say  $\mu \neq \infty$  is a *finite eigenvalue* of  $\mathbf{A} - \lambda\mathbf{B}$  if

$$\text{rank}(\mathbf{A} - \mu\mathbf{B}) < \max_{\lambda \in \mathbb{C}} \text{rank}(\mathbf{A} - \lambda\mathbf{B}), \quad (\text{A.1})$$

and  $x \in \mathbb{C}^n$  is a corresponding *eigenvector* if  $0 \neq x \notin \mathcal{N}(\mathbf{A}) \cap \mathcal{N}(\mathbf{B})$  satisfies

$$\mathbf{A}x = \mu\mathbf{B}x, \quad (\text{A.2})$$

or equivalently,  $0 \neq x \in \mathcal{N}(\mathbf{A} - \mu\mathbf{B}) \setminus (\mathcal{N}(\mathbf{A}) \cap \mathcal{N}(\mathbf{B}))$ , where  $\mathcal{N}(\cdot)$  is the null space of a matrix.

In the rest of this subsection,  $\mathbf{A} - \lambda\mathbf{B}$  is assumed to be a positive semidefinite pencil. Let the inertia of  $\mathbf{B}$  be  $(i_-(\mathbf{B}), i_0(\mathbf{B}), i_+(\mathbf{B}))$ , meaning that  $\mathbf{B}$  has  $i_-(\mathbf{B})$  negative,  $i_0(\mathbf{B})$  zero, and  $i_+(\mathbf{B})$  positive eigenvalues, respectively, and set

$$n_- := i_-(\mathbf{B}), \quad n_+ := i_+(\mathbf{B}), \quad r := \text{rank}(\mathbf{B}) = n_+ + n_-.$$

Given  $0 \leq k_+ \leq n_+$  and  $0 \leq k_- \leq n_-$ , set

$$J_k = \begin{bmatrix} I_{k_+} & \\ & -I_{k_-} \end{bmatrix}.$$

We proved the following theorem in [38, Lemma 3.8], but later found out that it had been obtained in [13, Theorem 4.1] for the regular pencil case. This theorem play a major role in this paper.

**Theorem A.1** ([13, 38]). *Let  $\mathbf{A} - \lambda\mathbf{B}$  be a positive semidefinite Hermitian pencil of order  $n$ , and suppose that  $\lambda_0 \in \mathbb{R}$  such that  $\mathbf{A} - \lambda_0\mathbf{B} \succeq 0$ .*

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

$$W^H \mathbf{A} W = \begin{matrix} & \begin{matrix} n_1 & r-n_1 & n-r \end{matrix} \\ \begin{matrix} n_1 \\ r-n_1 \\ n-r \end{matrix} & \begin{bmatrix} A_1 & & \\ & A_0 & \\ & & A_\infty \end{bmatrix} \end{matrix}, \quad W^H \mathbf{B} W = \begin{matrix} & \begin{matrix} n_1 & r-n_1 & n-r \end{matrix} \\ \begin{matrix} n_1 \\ r-n_1 \\ n-r \end{matrix} & \begin{bmatrix} \Omega_1 & & \\ & \Omega_0 & \\ & & 0 \end{bmatrix} \end{matrix}, \quad (\text{A.3})$$

where

(a)  $A_1 = \text{diag}(s_1\alpha_1, \dots, s_{n_1}\alpha_{n_1})$ ,  $\Omega_1 = \text{diag}(s_1, \dots, s_{n_1})$ ,  $s_i = \pm 1$ , and  $A_1 - \lambda_0\Omega_1 \succ 0$ ;

(b)  $A_0 = \text{diag}(A_{0,1}, \dots, A_{0,m+m_0})$  and  $\Omega_0 = \text{diag}(\Omega_{0,1}, \dots, \Omega_{0,m+m_0})$  with

$$A_{0,i} = t_i\lambda_0, \quad \Omega_{0,i} = t_i = \pm 1, \quad \text{for } 1 \leq i \leq m,$$

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

There is no such pair  $(A_0, \Omega_0)$  if  $A - \lambda_0 B \succ 0$ . Evidently  $m + 2m_0 = r - n_1$ .

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

The representations in (A.3) are uniquely determined by  $A - \lambda B$ , up to a simultaneous permutation of the corresponding  $1 \times 1$  and  $2 \times 2$  diagonal block pairs  $(s_i\alpha_i, s_i)$  for  $1 \leq i \leq n_1$ ,  $(A_{0,i}, \Omega_{0,i})$  for  $1 \leq i \leq m + m_0$ , and  $(\alpha_i, 0)$  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 as<sup>21</sup>

$$\lambda_1^- \leq \dots \leq \lambda_{n_-}^- \leq \lambda_1^+ \leq \dots \leq \lambda_{n_+}^+. \quad (\text{A.4})$$

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

$$\{\gamma \in \mathbb{R} \mid A - \gamma B \succ 0\} = (\lambda_{n_-}^-, \lambda_1^+).$$

The next perturbation theorem for positive definite pencils seem to be new. It resembles various perturbation bounds in [8, 32, 33, 54, 57]. For the definition and properties of such unitarily invariant norms, the reader is referred to [6, 56] for details. In this article, for convenience, any  $\|\cdot\|_{\text{ui}}$  we use is generic to matrix sizes in the sense that it applies to matrices of all sizes. Examples include the matrix spectral norm  $\|\cdot\|_2$  and the Frobenius norm  $\|\cdot\|_{\text{F}}$ . Two important properties of unitarily invariant norms are

$$\|X\|_2 \leq \|X\|_{\text{ui}}, \quad \|XYZ\|_{\text{ui}} \leq \|X\|_2 \cdot \|Y\|_{\text{ui}} \cdot \|Z\|_2 \quad (\text{A.5})$$

for any matrices  $X$ ,  $Y$ , and  $Z$  of compatible sizes.

**Theorem A.2.** Let  $A - \lambda B$  and  $\tilde{A} - \tilde{\lambda B}$  be two positive definite Hermitian pencils of order  $n$ , admitting the following eigen-decompositions<sup>22</sup>:

$$W^H A W = J \Lambda, \quad W^H B W = J, \quad (\text{A.6a})$$

$$\tilde{W}^H \tilde{A} \tilde{W} = \tilde{J} \tilde{\Lambda}, \quad \tilde{W}^H \tilde{B} \tilde{W} = \tilde{J}, \quad (\text{A.6b})$$

where  $\Lambda$  is diagonal with diagonal entries consisting eigenvalues of  $A - \lambda B$  in ascending order,  $J = \text{diag}(-I_{i_-(B)}, I_{i_+(B)})$ , and similarly for  $\tilde{\Lambda}$  and  $\tilde{J}$ . Then for any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ ,

$$\|\tilde{A} - A\|_{\text{ui}} \leq \|W\|_2 \|\tilde{W}\|_2 \left( \|\tilde{A} - A\|_{\text{ui}} + \xi \|\tilde{B} - B\|_{\text{ui}} \right), \quad (\text{A.7})$$

where  $\xi = \max\{\|A\|_2, \|\tilde{A}\|_2\}$ .

<sup>21</sup>This ordering is different from the one we used in [38, 37] for the neg-type eigenvalues, in order to be consistent with what we will be using later for hyperbolic matrix polynomials. See Theorem 2.1.

<sup>22</sup>Such decompositions are guaranteed by Theorem A.1

*Proof.* We have

$$\begin{aligned} \mathbf{A}W W^H \mathbf{B} - \mathbf{B}W W^H \mathbf{A} &= 0, \\ \tilde{\mathbf{A}}W W^H \mathbf{B} - \tilde{\mathbf{B}}W W^H \mathbf{A} &= \tilde{\mathbf{A}}W W^H \mathbf{B} - \tilde{\mathbf{B}}W W^H \mathbf{A} - (\mathbf{A}W W^H \mathbf{B} - \mathbf{B}W W^H \mathbf{A}) \\ &= (\tilde{\mathbf{A}} - \mathbf{A})W W^H \mathbf{B} - (\tilde{\mathbf{B}} - \mathbf{B})W W^H \mathbf{A}. \end{aligned} \quad (\text{A.8})$$

Pre- and post-multiply (A.8) by  $\tilde{J}\tilde{W}^H$  and  $WJ$ , and plug the decompositions in (A.6) into (A.8) to get

$$\tilde{\Lambda}\tilde{W}^{-1}W - \tilde{W}^{-1}W\Lambda = \tilde{J}\tilde{W}^H(\tilde{\mathbf{A}} - \mathbf{A})W - \tilde{J}\tilde{W}^H(\tilde{\mathbf{B}} - \mathbf{B})W\Lambda. \quad (\text{A.9})$$

Switching the roles of  $\mathbf{A} - \lambda\mathbf{B}$  and  $\tilde{\mathbf{A}} - \lambda\tilde{\mathbf{B}}$ , we conclude from (A.9) that

$$\Lambda W^{-1}\tilde{W} - W^{-1}\tilde{W}\tilde{\Lambda} = JW^H(\mathbf{A} - \tilde{\mathbf{A}})\tilde{W} - JW^H(\mathbf{B} - \tilde{\mathbf{B}})\tilde{W}\tilde{\Lambda}. \quad (\text{A.10})$$

It follows from (A.9) and (A.10) that

$$\|\tilde{\Lambda}\tilde{W}^{-1}W - \tilde{W}^{-1}W\Lambda\|_{\text{ui}} \leq \|W\|_2 \|\tilde{W}\|_2 \left( \|\tilde{\mathbf{A}} - \mathbf{A}\|_{\text{ui}} + \xi \|\tilde{\mathbf{B}} - \mathbf{B}\|_{\text{ui}} \right), \quad (\text{A.11a})$$

$$\|\Lambda W^{-1}\tilde{W} - W^{-1}\tilde{W}\tilde{\Lambda}\|_{\text{ui}} \leq \|W\|_2 \|\tilde{W}\|_2 \left( \|\tilde{\mathbf{A}} - \mathbf{A}\|_{\text{ui}} + \xi \|\tilde{\mathbf{B}} - \mathbf{B}\|_{\text{ui}} \right). \quad (\text{A.11b})$$

Let  $\tilde{W}^{-1}W = U\Sigma V^H$  be the SVD of  $\tilde{W}^{-1}W$  and set  $\mathbf{C} = V^H\Lambda V$  and  $\tilde{\mathbf{C}} = U^H\tilde{\Lambda}U$ , both of which are Hermitian. It can be verified that

$$\begin{aligned} \tilde{\Lambda}\tilde{W}^{-1}W - \tilde{W}^{-1}W\Lambda &= U(\tilde{\mathbf{C}}\Sigma - \Sigma\mathbf{C})V^H, \\ \Lambda W^{-1}\tilde{W} - W^{-1}\tilde{W}\tilde{\Lambda} &= V(\mathbf{C}\Sigma^{-1} - \Sigma^{-1}\tilde{\mathbf{C}})U. \end{aligned}$$

Theorem 2.1 of [7] yields

$$\|\tilde{\mathbf{C}} - \mathbf{C}\|_{\text{ui}}^2 \leq \|\tilde{\mathbf{C}}\Sigma - \Sigma\mathbf{C}\|_{\text{ui}} \|\mathbf{C}\Sigma^{-1} - \Sigma^{-1}\tilde{\mathbf{C}}\|_{\text{ui}}. \quad (\text{A.12})$$

Mirsky's theorem [56, p.204] says

$$\|\tilde{\Lambda} - \Lambda\|_{\text{ui}} \leq \|\tilde{\mathbf{C}} - \mathbf{C}\|_{\text{ui}}. \quad (\text{A.13})$$

The main result (A.7) is now a consequence of (A.11) – (A.13).  $\square$

In Theorem A.2, the upper bound by (A.7) contains  $\|W\|_2$  and  $\|\tilde{W}\|_2$ . They can be bounded, too, in terms of extreme pos- and neg-type eigenvalues.

**Theorem A.3.** *Let  $\mathbf{A} - \lambda\mathbf{B}$  be a positive definite Hermitian pencil of order  $n$ , with eigenvalues given by and ordered as in (A.4), and let its eigen-decomposition be given by (A.6a). Then for any  $\lambda_0 \in (\lambda_{n-}^-, \lambda_1^+)$*

$$\|W\|_2 \leq \sqrt{\max\{\lambda_{n+}^+ - \lambda_0, \lambda_0 - \lambda_1^-\} \|(\mathbf{A} - \lambda_0\mathbf{B})^{-1}\|_2}, \quad (\text{A.14a})$$

$$\|W^{-1}\|_2 \leq \sqrt{\frac{1}{\min\{\lambda_1^+ - \lambda_0, \lambda_0 - \lambda_{n-}^-\}} \|\mathbf{A} - \lambda_0\mathbf{B}\|_2}. \quad (\text{A.14b})$$

In particular, taking  $\lambda_0 = (\lambda_{n_-}^- + \lambda_1^+)/2$  gives

$$\|W\|_2 \leq \sqrt{(\lambda_{n_+}^+ - \lambda_1^-) \|(\mathbf{A} - \lambda_0 \mathbf{B})^{-1}\|_2}, \quad (\text{A.15a})$$

$$\|W^{-1}\|_2 \leq \sqrt{\frac{2}{\lambda_1^+ - \lambda_{n_-}^-} \|\mathbf{A} - \lambda_0 \mathbf{B}\|_2}. \quad (\text{A.15b})$$

*Proof.* For  $\lambda_0 \in (\lambda_{n_-}^-, \lambda_1^+)$ ,  $\mathbf{A} - \lambda_0 \mathbf{B} \succ 0$ . We have  $\mathbf{A} - \lambda_0 \mathbf{B} \succeq \lambda_{\min}(\mathbf{A} - \lambda_0 \mathbf{B}) I_n$  and thus

$$\lambda_{\min}(\mathbf{A} - \lambda_0 \mathbf{B}) W^H W \preceq W^H (\mathbf{A} - \lambda_0 \mathbf{B}) W = J(\Lambda - \lambda_0 I) \preceq \max\{\lambda_{n_+}^+ - \lambda_0, \lambda_0 - \lambda_1^-\} I$$

which gives (A.14a). We also have

$$W^H (\mathbf{A} - \lambda_0 \mathbf{B}) W = J(\Lambda - \lambda_0 I) \succeq \min\{\lambda_1^+ - \lambda_0, \lambda_0 - \lambda_{n_-}^-\} I$$

to give

$$W^{-H} W^{-1} \preceq \frac{1}{\min\{\lambda_1^+ - \lambda_0, \lambda_0 - \lambda_{n_-}^-\}} (\mathbf{A} - \lambda_0 \mathbf{B})$$

which yields (A.14b). □

## References

- [1] A.R. Amir-Moéz. Extreme properties of eigenvalues of a Hermitian transformation and singular values of the sum and product of linear transformations. *Duke Math. J.*, 23:463–476, 1956.
- [2] E. Anderson, Z. Bai, C. Bischof, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, S. Ostrouchov, and D. Sorensen. *LAPACK Users' Guide*. SIAM, Philadelphia, 3rd edition, 1999.
- [3] Zhaojun Bai, Ren-Cang Li, and Yangfeng Su. Lecture notes on matrix eigenvalue computations. Prepared for 2009 Summer School on Numerical Linear Algebra, Chinese Academy of Science, July 2009.
- [4] Lawrence Barkwell and Peter Lancaster. Overdamped and gyroscopic vibrating systems. *J. Appl. Mech.*, 59(1):176–181, 1992.
- [5] T. Betcke, N. J. Higham, V. Mehrmann, C. Schröder, and F. Tisseur. NLEVP: A collection of nonlinear eigenvalue problems. *ACM Trans. Math. Software*, 39(2):7:1–7:28, 2013.
- [6] R. Bhatia. *Matrix Analysis*. Graduate Texts in Mathematics, vol. 169. Springer, New York, 1996.
- [7] R. Bhatia, F. Kittaneh, and Ren-Cang Li. Some inequalities for commutators and an application to spectral variation. II. *Lin. Multilin. Alg.*, 43(1-3):207–220, 1997.
- [8] R. Bhatia and Ren-Cang Li. On perturbations of matrix pencils with real spectra. II. *Math. Comp.*, 65(214):637–645, 1996.
- [9] Moody T. Chu and Shufang Xu. Spectral decomposition of real symmetric quadratic  $\lambda$ -matrices and its applications. *Math. Comp.*, 78:293–313, 2009.
- [10] G. Davis. Numerical solution of a quadratic matrix equation. *SIAM J. Sci. Statist. Comput.*, 2(2):164–175, 1981.
- [11] J. Demmel. *Applied Numerical Linear Algebra*. SIAM, Philadelphia, PA, 1997.
- [12] R. Duffin. A minimax theory for overdamped networks. *Indiana Univ. Math. J.*, 4:221–233, 1955.
- [13] D. C. Dzens and W. W. Lin. Homotopy continuation method for the numerical solutions of generalised symmetric eigenvalue problems. *J. Austral. Math. Soc. Ser. B*, 32:437–456, 4 1991.
- [14] D. K. Faddeev and V. N. Faddeeva. *Computational Methods of Linear Algebra*. Undergraduate Mathematics Books. W.H.Freeman & Co Ltd, San Francisco, 1963. Translated by R. C. Williams.
- [15] Ky Fan. On a theorem of Weyl concerning eigenvalues of linear transformations. I. *Proc. Natl. Acad. Sci. USA*, 35(11):pp. 652–655, 1949.
- [16] I. Gohberg, P. Lancaster, and L. Rodman. *Matrix Polynomials*, volume 58 of *Classics in Applied Mathematics*. SIAM, Philadelphia, 2009.
- [17] G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, Maryland, 3rd edition, 1996.
- [18] Gene H. Golub and Qiang Ye. An inverse free preconditioned Krylov subspace method for symmetric generalized eigenvalue problems. *SIAM J. Sci. Comput.*, 24(1):312–334, 2002.
- [19] Anne Greenbaum. *Iterative Methods for Solving Linear Systems*. SIAM, Philadelphia, 1997.



- [20] C.-H. Guo and P. Lancaster. Algorithms for hyperbolic quadratic eigenvalue problems. *Math. Comp.*, 74:1777–1791, 2005.
- [21] Chun-Hua Guo. Numerical solution of a quadratic eigenvalue problem. *Linear Algebra Appl.*, 385(0):391–406, 2004.
- [22] Chun-Hua Guo, Nicholas J. Higham, and Françoise Tisseur. Detecting and solving hyperbolic quadratic eigenvalue problems. *SIAM J. Matrix Anal. Appl.*, 30(4):1593–1613, 2009.
- [23] N. J. Higham, D. Mackey, and F. Tisseur. Definite matrix polynomials and their linearization by definite pencils. *SIAM J. Matrix Anal. Appl.*, 31(2):478–502, 2009.
- [24] Nicholas J. Higham and Hyun-Min Kim. Numerical analysis of a quadratic matrix equation. *IMA J. Numer. Anal.*, 20(4):499–519, 2000.
- [25] Nicholas J. Higham, Ren-Cang Li, and Françoise Tisseur. Backward error of polynomial eigenproblems solved by linearization. *SIAM J. Matrix Anal. Appl.*, 29(4):1218–1241, 2007.
- [26] 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.
- [27] A. J. Hoffman and H. W. Wielandt. The variation of the spectrum of a normal matrix. *Duke Math. J.*, 20:37–39, 1953.
- [28] Andrew V. Knyazev. Toward the optimal preconditioned eigensolver: Locally optimal block preconditioned conjugate gradient method. *SIAM J. Sci. Comput.*, 23(2):517–541, 2001.
- [29] 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.
- [30] J. Kovač-Striko and K. Veselić. Trace minimization and definiteness of symmetric pencils. *Linear Algebra Appl.*, 216:139–158, 1995.
- [31] Peter Lancaster. Inverse spectral problems for semisimple damped vibrating systems. *SIAM J. Matrix Anal. Appl.*, 29(1):279–301, 2007.
- [32] Ren-Cang Li. On perturbations of matrix pencils with real spectra. *Math. Comp.*, 62:231–265, 1994.
- [33] Ren-Cang Li. On perturbations of matrix pencils with real spectra, a revisit. *Math. Comp.*, 72:715–728, 2003.
- [34] Ren-Cang Li. On Meinardus’ examples for the conjugate gradient method. *Math. Comp.*, 77(261):335–352, 2008. Electronically published on September 17, 2007.
- [35] Ren-Cang Li. Rayleigh quotient based numerical methods for eigenvalue problems. Gene Golub SIAM Summer School 2013, Fudan University, Shanghai, July 22–August 9 2013. SIAM. <http://math.fudan.edu.cn/g2s3/>.
- [36] Ren-Cang Li, Wen-Wei Lin, and Chern-Shuh Wang. Structured backward error for palindromic polynomial eigenvalue problems. *Numer. Math.*, 116(1):95–122, 2010.
- [37] Xin Liang and Ren-Cang Li. Extensions of Wielandt’s min-max principles for positive semi-definite pencils. *Lin. Multilin. Alg.*, 2013. DOI:10.1080/03081087.2013.803242.
- [38] Xin Liang, Ren-Cang Li, and Zhaojun Bai. Trace minimization principles for positive semi-definite pencils. *Linear Algebra Appl.*, 438:3085–3106, 2013.

- [39] V. B. Lidskii. The proper values of the sum and product of symmetric matrices. *Dokl. Akad. Nauk SSSR*, 75:769–772, 1950. In Russian. Translation by C. Benster available from the National Translation Center of the Library of Congress.
- [40] D. E. Longsine and S. F. McCormick. Simultaneous Rayleigh-quotient minimization methods for  $Ax = \lambda Bx$ . *Linear Algebra Appl.*, 34:195–234, 1980.
- [41] D. Steven Mackey, Niloufer Mackey, Christian Mehl, and Volker Mehrmann. Structured polynomial eigenvalue problems: Good vibrations from good linearizations. *SIAM J. Matrix Anal. Appl.*, 28(4):1029–1051, 2006.
- [42] D. Steven Mackey, Niloufer Mackey, Christian Mehl, and Volker Mehrmann. Vector spaces of linearizations for matrix polynomials. *SIAM J. Matrix Anal. Appl.*, 28(4):971–1004, 2006.
- [43] A.S. Markus. *Introduction to the Spectral Theory of Polynomial Operator Pencils*. Translations of mathematical monographs, vol. 71. AMS, Providence, RI, 1988.
- [44] C. B. Moler and G. W. Stewart. An algorithm for generalized matrix eigenvalue problems. *SIAM J. Numer. Anal.*, 10(2):241–256, 1973.
- [45] Jorge Nocedal and Stephen Wright. *Numerical Optimization*. Springer, 2nd edition, 2006.
- [46] E. E. Ovtchinnikov. Sharp convergence estimates for the preconditioned steepest descent method for Hermitian eigenvalue problems. *SIAM J. Numer. Anal.*, 43(6):2668–2689, 2006.
- [47] B. N. Parlett. *The Symmetric Eigenvalue Problem*. SIAM, Philadelphia, 1998.
- [48] Boris Teodorovich Polyak. *Introduction to Optimization*. Optimization Software, New York, 1987.
- [49] E. H. Rogers. A mimmax theory for overdamped systems. *Arch. Rational Mech. Anal.*, 16:89–96, 1964.
- [50] E. H. Rogers. Variational properties of nonlinear spectra. *Indiana Univ. Math. J.*, 18:479–490, 1969.
- [51] Yousef Saad. *Iterative Methods for Sparse Linear Systems*. SIAM, Philadelphia, 2nd edition, 2003.
- [52] B. Samokish. The steepest descent method for an eigenvalue problem with semi-bounded operators. *Izv. Vyssh. Uchebn. Zaved. Mat.*, 5:105–114, 1958. in Russian.
- [53] Bernhard A. Schmitt. Perturbation bounds for matrix square roots and Pythagorean sums. *Linear Algebra Appl.*, 174:215–227, 1992.
- [54] G. W. Stewart. Perturbation bounds for the definite generalized eigenvalue problem. *Linear Algebra Appl.*, 23:69–86, 1979.
- [55] G. W. Stewart. *Matrix Algorithms, Vol. II: Eigensystems*. SIAM, Philadelphia, 2001.
- [56] G. W. Stewart and Ji-Guang Sun. *Matrix Perturbation Theory*. Academic Press, Boston, 1990.
- [57] Ji-Guang Sun. A note on Stewart’s theorem for definite matrix pairs. *Linear Algebra Appl.*, 48:331–339, 1982.
- [58] Ji-Guang Sun. Perturbation bounds for eigenspaces of a definite matrix pair. *Numer. Math.*, 41:321–343, 1983.
- [59] Wenyu Sun and Ya-Xiang Yuan. *Optimization Theory and Methods – Nonlinear Programming*. Springer, New York, 2006.

- [60] I. Takahashi. A note on the conjugate gradient method. *Inform. Process. Japan*, 5:45–49, 1965.
- [61] F. Tisseur and K. Meerbergen. The quadratic eigenvalue problem. *SIAM Rev.*, 43(2):235–386, 2001.
- [62] Françoise Tisseur. Backward error and condition of polynomial eigenvalue problems. *Linear Algebra Appl.*, 309(1-3):339–361, April 2000.
- [63] K. Veselić. A Jacobi eigenreduction algorithm for definite matrix pairs. *Numer. Math.*, 64:241–269, 1993.
- [64] K. Veselić. Note on interlacing for hyperbolic quadratic pencils. In Jussi Behrndt, Karl-Heinz Förster, and Carsten Trunk, editors, *Recent Advances in Operator Theory in Hilbert and Krein Spaces*, volume 198 of *Oper. Theory: Adv. Appl.*, pages 305–307. Birkhäuser, Boston, 2010.
- [65] K. Veselić. *Damped Oscillations of Linear Systems*, volume 2023 of *Lecture Notes in Mathematics*. Springer, Berlin, 2011.
- [66] 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.
- [67] Heinrich Voss. A minmax principle for nonlinear eigenproblems depending continuously on the eigenparameter. *Numer. Linear Algebra Appl.*, 16(11-12):899–913, 2009.
- [68] Songbin Wei and Imin Kao. Vibration analysis of wire and frequency response in the modern wiresaw manufacturing process. *Journal of Sound and Vibration*, 231(5):1383–1395, 2000.
- [69] H. Wielandt. An extremum property of sums of eigenvalues. *Proc. Amer. Math. Soc.*, 6:106–110, 1955.
- [70] H. Yang. Conjugate gradient methods for the Rayleigh quotient minimization of generalized eigenvalue problems. *Computing*, 51:79–94, 1993.