

## Backward perturbation analysis and residual-based error bounds for the linear response eigenvalue problem

Lei-Hong Zhang · Wen-Wei Lin · Ren-Cang Li

Received: 9 February 2014 / Accepted: 25 August 2014  
© Springer Science+Business Media Dordrecht 2014

**Abstract** The numerical solution of a large scale linear response eigenvalue problem is often accomplished by computing a pair of deflating subspaces associated with the interesting part of the spectrum. This paper is concerned with the backward perturbation analysis for a given pair of approximate deflating subspaces or an approximate eigenquadruple. Various optimal backward perturbation bounds are obtained, as well as bounds for approximate eigenvalues computed through the pair of approximate deflating subspaces or approximate eigenquadruple. These results are reminiscent of many existing classical ones for the standard eigenvalue problem.

**Keywords** Linear response eigenvalue problem · Eigenvalue approximation · Rayleigh–Ritz approximation · Backward perturbation · Error bound · Deflating subspace

**Mathematics Subject Classification** 65L15 · 65F15 · 81Q15 · 15A18 · 15A42

---

Communicated by Daniel Kressner.

---

L.-H. Zhang  
School of Mathematics, Shanghai University of Finance and Economics, 777 Guoding Road,  
Shanghai 200433, People's Republic of China  
e-mail: longzlh@163.com

W.-W. Lin  
Department of Applied Mathematics, National Chiao Tung University, No.1001 University Road,  
Hsinchu 30013, Taiwan  
e-mail: wwlin@math.nctu.edu.tw

R.-C. Li (✉)  
Department of Mathematics, University of Texas at Arlington, P.O. Box 19408, Arlington,  
TX 76019-0408, USA  
e-mail: rcli@uta.edu

## 1 Introduction

In this paper, we are concerned with a backward perturbation analysis and residual-based error bounds for the linear response eigenvalue problem (LREP):

$$H\mathbf{z} := \begin{bmatrix} 0 & K \\ M & 0 \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix} = \lambda \begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix} =: \lambda\mathbf{z}, \quad (1.1)$$

where  $K$  and  $M$  are both  $n$ -by- $n$  real symmetric and positive definite. The matrix  $H$  in (1.1) is a special Hamiltonian matrix whose eigenvalues are real [1] and come in pairs  $\{\lambda, -\lambda\}$ . Therefore, we can order the  $2n$  eigenvalues of (1.1) as

$$-\lambda_n \leq \cdots \leq -\lambda_1 < \lambda_1 \leq \cdots \leq \lambda_n. \quad (1.2)$$

LREP (1.1) is mathematically equivalent to the so-called random phase approximation (RPA) eigenvalue problem in computational quantum chemistry and physics:

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

where  $A, B \in \mathbb{R}^{n \times n}$  are both symmetric matrices and  $\begin{bmatrix} A & B \\ B & A \end{bmatrix}$  is positive definite. The equivalent relationship is established through the orthogonal matrix  $J = \frac{1}{\sqrt{2}} \begin{bmatrix} I_n & I_n \\ I_n & -I_n \end{bmatrix}$  and the similarity transformation (see, e.g., [1,2])

$$J^T \begin{bmatrix} A & B \\ -B & -A \end{bmatrix} J = \begin{bmatrix} 0 & A - B \\ A + B & 0 \end{bmatrix} =: \begin{bmatrix} 0 & K \\ M & 0 \end{bmatrix}, \quad \text{and} \quad \begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix} := J^T \begin{bmatrix} \mathbf{u} \\ \mathbf{v} \end{bmatrix}. \quad (1.3)$$

RPA is one of the most widely used methods in studying the excitation states (energies) of physical systems in the study of collective motion of many-particle systems [1,41,42] which has applications in silicon nanoparticles, nanoscale materials, analysis of interstellar clouds [1,2], among others. The heart of RPA calculations is to compute a few eigenpairs associated with the smallest *positive* eigenvalues, which, by the equivalent relationship (1.3), are the eigenpairs associated with the eigenvalues  $\lambda_1 \leq \cdots \leq \lambda_k$  of (1.1).

In consistent with [1,2], throughout the rest of this paper, we relax the condition on  $K, M \in \mathbb{R}^{n \times n}$  to that they are symmetric positive semi-definite and one of them is definite, unless explicitly stated differently. This means that possibly  $\lambda_1 = 0$ . Also the assignments in (1.1) will be assumed.

As the dimension  $n$  is usually very large, LREP is generally solved by iterative methods. Roughly speaking, any large scale eigenvalue computation is about approximating certain invariant subspaces associated with the interested part of the spectrum, and the interested eigenvalues are then extracted from projecting the problem by approximate invariant subspaces into much smaller eigenvalue problems. In the case of the linear

response eigenvalue problem, it is the *pair of deflating subspaces* associated with the first few smallest  $\lambda_i$  that needs to be computed [2].

For two  $k$ -dimensional subspaces  $\mathcal{U}$  and  $\mathcal{V}$  in  $\mathbb{R}^n$ , we call  $\{\mathcal{U}, \mathcal{V}\}$  a *pair of deflating subspaces* of  $\{K, M\}$  if

$$K\mathcal{U} \subseteq \mathcal{V} \quad \text{and} \quad M\mathcal{V} \subseteq \mathcal{U}. \tag{1.4}$$

This notion of the pair of deflating subspaces is a generalization of the concept of the invariant subspace (or, eigenspace) in the standard eigenvalue problem upon considering the special structure in LREP (1.1) [1]. Whenever such a pair of deflating subspaces is available, we can project LREP (1.1) into a much smaller problem in the form of (1.1), an LREP by its own, whose spectrum are part of that of  $H$  (see more discussions in Sect. 2 and [1, 2]). Based on this fact, several efficient algorithms, including the Locally Optimal Block Preconditioned 4D Conjugate Gradient Method (LOBP4DCG) [2], the block Chebyshev-Davidson method [40], as well as the generalized Lanczos method [39, 43, 44], have been proposed. Each of these algorithms generates a sequence of approximate deflating subspace pairs that hopefully converge to or contain subspaces near the pair of deflating subspaces. The goal of this paper is to perform a backward perturbation analysis and to establish error bounds on the accuracy (in eigenvalue/eigenspace approximations) using proper residuals associated with any given approximate deflating subspace pair.

A related study is presented in [46]. The main difference among the results in [46] and those in the present paper is that the error bounds on eigenvalue/eigenspace approximations in [46] are characterized by the canonical angles between the approximate deflating subspace pair and the exact pair, whereas the error bounds in this paper use certain computable residuals. These two types of error bounds are well-established in the standard eigenvalue problem (see, e.g., [30, 33]), and both types are useful in analyzing the convergence and designing stopping criteria for iterative algorithms.

The rest of the paper is organized as follows. In Sect. 2, we will state some basic properties about LREP for use later. Section 3 gives a backward perturbation analysis for a given pair of approximate deflating subspaces or an approximate eigenquadruple, optimizes backward perturbation errors, and shows the near optimality of the so-called Rayleigh quotient pair. Section 4 derives several error bounds in terms of residuals for eigenvalue approximations. In Sect. 5, we review related results for the standard eigenvalue problem as a comparison. Finally in Sect. 6, we present our concluding remarks.

**Notation.**  $\mathbb{K}^{n \times m}$  is the set of all  $n \times m$  matrices whose entries belong to the number field  $\mathbb{K}$ ,  $\mathbb{K}^n = \mathbb{K}^{n \times 1}$ , and  $\mathbb{K} = \mathbb{K}^1$ , where  $\mathbb{K} = \mathbb{R}$  (the set of real numbers) or  $\mathbb{C}$  (the set of complex numbers).  $I_n$  (or simply  $I$  if its dimension is clear from the context) denotes the  $n \times n$  identity matrix. All vectors are column vectors and are in boldface. For a matrix  $Z$ ,

1.  $Z^T$  and  $Z^H$  denote its transpose and the conjugate transpose, respectively,
2.  $\mathcal{R}(Z)$  is the column space of  $Z$ , spanned by its column vectors,
3.  $Z^\dagger$  stands for the Moore-Penrose inverse and  $P_Z = ZZ^\dagger$  is the orthogonal projection onto  $\mathcal{R}(Z)$  and  $P_Z^\perp = I - P_Z$  [33],

4.  $\|Z\|_2$ ,  $\|Z\|_F$ , and  $\|Z\|_{\text{ui}}$  are the spectral norm, the Frobenius norm, and a general unitarily invariant norm, respectively,
5. The submatrices  $Z_{(k:\ell,i:j)}$ ,  $Z_{(k:\ell,:)}$ , and  $Z_{(:,i:j)}$  of  $Z$  consist of intersections of row  $k$  to row  $\ell$  and column  $i$  to column  $j$ , row  $k$  to row  $\ell$ , and column  $i$  to column  $j$ , respectively,
6. When  $Z$  is a square matrix, its trace is  $\text{trace}(Z)$ , its eigenvalue set is  $\text{eig}(Z)$ , and its spectral condition number is  $\kappa(Z) = \|Z\|_2 \|Z^{-1}\|_2$ .

The oplus-sum  $\mathcal{V} \oplus \mathcal{U}$  of two subspaces  $\mathcal{V}$  and  $\mathcal{U}$  in  $\mathbb{K}^n$  is a subspace of  $\mathbb{K}^{2n}$  and consists of all vectors  $[\mathbf{y}^T, 0]^T + [0, \mathbf{x}^T]^T$  for all  $\mathbf{y} \in \mathcal{V}$  and  $\mathbf{x} \in \mathcal{U}$ .

## 2 Preliminaries

Many theoretical properties of LREP have been established in [1,2]. In Theorem 2.1, we present certain decompositions on  $K$  and  $M$ , necessary for our developments later in this paper. The reader is referred to [1, section 2] for proofs and more.

**Theorem 2.1** *Suppose that  $K$  is semidefinite and  $M$  is definite. Then the following statements are true:*

- (i) *There exists a nonsingular  $\Phi \in \mathbb{R}^{n \times n}$  such that*

$$K = \Psi \Lambda^2 \Psi^T \quad \text{and} \quad M = \Phi \Phi^T,$$

where  $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$  with  $0 \leq \lambda_1 \leq \dots \leq \lambda_n$ , and  $\Psi = \Phi^{-T}$ .

- (ii) *If  $K$  is also definite, then all  $\lambda_i > 0$  and  $H$  is diagonalizable:*

$$H \begin{bmatrix} \Psi \Lambda & \Psi \Lambda \\ -\Phi & \Phi \end{bmatrix} = \begin{bmatrix} \Psi \Lambda & \Psi \Lambda \\ -\Phi & \Phi \end{bmatrix} \begin{bmatrix} -\Lambda & \\ & \Lambda \end{bmatrix}.$$

- (iii) *The eigen-decomposition of  $KM$  and  $MK$  are*

$$(KM)\Psi = \Psi \Lambda^2 \quad \text{and} \quad (MK)\Phi = \Phi \Lambda^2, \quad (2.1)$$

respectively.

As we have introduced in Sect. 1, for two given  $k$ -dimensional subspaces  $\mathcal{U} \subseteq \mathbb{R}^n$  and  $\mathcal{V} \subseteq \mathbb{R}^n$ , the pair  $\{\mathcal{U}, \mathcal{V}\}$  is called a *pair of deflating subspaces* of  $\{K, M\}$  if

$$K\mathcal{U} \subseteq \mathcal{V} \quad \text{and} \quad M\mathcal{V} \subseteq \mathcal{U} \quad (1.4)$$

hold. This definition is essentially the same as the existing ones for the product eigenvalue problem [3, 11, 26, 27]. It also closely relates to the deflating subspace for the generalized eigenvalue problem in [31] in that (1.4) can be equivalently restated as

$$\begin{bmatrix} M & \\ & K \end{bmatrix} (\mathcal{V} \oplus \mathcal{U}) \subset \begin{bmatrix} 0 & I_n \\ I_n & 0 \end{bmatrix} (\mathcal{V} \oplus \mathcal{U}),$$

noting that LREP (1.1) is equivalent to the generalized eigenvalue problem

$$\begin{bmatrix} M & \\ & K \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix} = \lambda \begin{bmatrix} 0 & I_n \\ I_n & 0 \end{bmatrix} \begin{bmatrix} \mathbf{y} \\ \mathbf{x} \end{bmatrix}. \tag{2.2}$$

Let  $U \in \mathbb{R}^{n \times k}$  and  $V \in \mathbb{R}^{n \times k}$  be the basis matrices for  $\mathcal{U}$  and  $\mathcal{V}$ , respectively. Alternatively, (1.4) can be restated as that there exist  $K_R \in \mathbb{R}^{k \times k}$  and  $M_R \in \mathbb{R}^{k \times k}$  such that

$$KU = VK_R \quad \text{and} \quad MV = UM_R, \tag{2.3}$$

and vice versa, or equivalently,

$$H \begin{bmatrix} V \\ U \end{bmatrix} = \begin{bmatrix} V \\ U \end{bmatrix} H_R \quad \text{with} \quad H_R := \begin{bmatrix} & K_R \\ M_R & \end{bmatrix},$$

i.e.,  $\mathcal{V} \oplus \mathcal{U}$  is an invariant subspace of  $H$  [1, Theorem 2.4]. We call  $\{U, V, K_R, M_R\}$  an *eigenquadruple* of  $\{K, M\}$ .

Whenever a pair of deflating subspaces  $\{\mathcal{R}(U), \mathcal{R}(V)\}$  is at hand, part of the eigenpairs of  $H$  can be obtained via solving the smaller eigenvalue problem [1, Theorem 2.5]: if

$$H_R \hat{\mathbf{z}} := \begin{bmatrix} & K_R \\ M_R & \end{bmatrix} \begin{bmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{x}} \end{bmatrix} = \lambda \begin{bmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{x}} \end{bmatrix} =: \lambda \hat{\mathbf{z}}, \tag{2.4}$$

then  $\left(\lambda, \begin{bmatrix} V\hat{\mathbf{y}} \\ U\hat{\mathbf{x}} \end{bmatrix}\right)$  is an eigenpair of  $H$ . The matrix  $H_R$  is the restriction of  $H$  onto  $\mathcal{V} \oplus \mathcal{U}$  with respect to the basis matrices  $V$  and  $U$  of  $\mathcal{V}$  and  $\mathcal{U}$ , respectively. Moreover, the eigenvalues of  $H_R$  are uniquely determined by the pair of deflating subspaces  $\{\mathcal{U}, \mathcal{V}\}$ ; in the other word, different choices of the basis matrices for  $\mathcal{U}$  and  $\mathcal{V}$  result in the same eigenvalues. In fact, if  $\hat{U} = UD_1 \in \mathbb{R}^{n \times k}$  and  $\hat{V} = VD_2 \in \mathbb{R}^{n \times k}$  are new basis matrices for  $\mathcal{R}(U)$  and  $\mathcal{R}(V)$ , respectively, and

$$K\hat{U} = \hat{V}\hat{K}_R \quad \text{and} \quad M\hat{V} = \hat{U}\hat{M}_R, \tag{2.5}$$

then  $\hat{K}_R = D_2^{-1}K_RD_1$  and  $\hat{M}_R = D_1^{-1}K_RD_2$  by comparing (2.3) to (2.5) after substituting in  $\hat{U} = UD_1$  and  $\hat{V} = VD_2$ . Thus

$$\hat{H}_R := \begin{bmatrix} & \hat{K}_R \\ \hat{M}_R & \end{bmatrix} = \begin{bmatrix} D_2 & \\ & D_1 \end{bmatrix}^{-1} H_R \begin{bmatrix} D_2 & \\ & D_1 \end{bmatrix}. \tag{2.6}$$

Evidently,  $\hat{H}_R$  and  $H_R$  must have the same eigenvalues.

Two particular choices of  $\{K_R, M_R\}$  to satisfy (2.3) are

$$K_R = (U^T V)^{-1} U^T K U, \quad M_R = (V^T U)^{-1} V^T M V; \tag{2.7a}$$

$$K_R = (V^T V)^{-1} V^T K U, \quad M_R = (U^T U)^{-1} U^T M V. \tag{2.7b}$$

In (2.7a),  $U^T V$  has to be assumed invertible, which is guaranteed since one of  $K$  and  $M$  is definite [1, Lemma 2.7]. By what we just proved, the associated  $H_R$  with either (2.7a) or (2.7b) must have the same eigenvalues.

In practical computations, however,  $\{\mathcal{U}, \mathcal{V}\}$  is usually a pair of approximate deflating spaces, i.e., no  $\{K_R, M_R\}$  that satisfies (2.3) exists. Dependent on how good  $\{\mathcal{U}, \mathcal{V}\}$  is as a pair of approximate deflating spaces, the equations in (2.3) is satisfied approximately to an appropriate level for some  $\{K_R, M_R\}$  like the ones given in (2.7). Regardless whether  $\{\mathcal{U}, \mathcal{V}\}$  is a pair of exact deflating spaces or approximate ones,  $\{K_R, M_R\}$  by (2.7b) is always well-defined, but for (2.7a), it is well-defined only if  $U^T V$  is nonsingular. This requirement is automatically satisfied if  $\{\mathcal{U}, \mathcal{V}\}$  is exact. It must also be true when  $\{\mathcal{U}, \mathcal{V}\}$  is a reasonably accurate approximation to an exact pair. Therefore it is quite reasonable to assume that  $U^T V$  is nonsingular from now on.

We note that  $\{K_R, M_R\}$  by (2.7a) relates to the structure-preserving projection  $H_{SR}$  of  $H$  in [2, (2.2)] that plays an important role numerically there. To highlight this particular pair, we will call  $\{K_R, M_R\}$  by (2.7a) a *Rayleigh quotient pair* of LREP (1.1) associated with  $\{\mathcal{R}(U), \mathcal{R}(V)\}$  and introduce

$$K_{RQ} := (U^T V)^{-1} U^T K U, \quad M_{RQ} := (V^T U)^{-1} V^T M V \quad (2.8)$$

for the ease of future references. Both  $K_{RQ}$  and  $M_{RQ}$  vary with different selections of  $U$  and  $V$  as the basis matrices of  $\mathcal{R}(U)$  and  $\mathcal{R}(V)$ , respectively. But the eigenvalues of the induced

$$H_{RQ} = \begin{bmatrix} & K_{RQ} \\ M_{RQ} & \end{bmatrix}. \quad (2.9)$$

do not. In fact, with new basis matrices  $\widehat{U} = U D_1$  and  $\widehat{V} = V D_2$  and, accordingly, new  $\widehat{K}_{RQ}$  and  $\widehat{M}_{RQ}$ ,  $\widehat{H}_{RQ}$  is similar to  $H_{RQ}$  [an equation like (2.6) holds].

For the definition and properties of unitarily invariant norms, the reader is referred to [5, 33] for details. In this article, for convenience, any  $\|\cdot\|_{ui}$  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  $\|\cdot\|_F$ . Two important properties of unitarily invariant norms are

$$\|X\|_2 \leq \|X\|_{ui}, \quad \|XYZ\|_{ui} \leq \|X\|_2 \cdot \|Y\|_{ui} \cdot \|Z\|_2$$

for any matrices  $X$ ,  $Y$ , and  $Z$  of compatible sizes.

### 3 Backward errors and optimal residuals

We recall our default assumption on  $K$ ,  $M \in \mathbb{R}^{n \times n}$ : both are symmetric and positive semi-definite and one of them is definite.

Let  $\mathcal{U} \subseteq \mathbb{R}^n$  and  $\mathcal{V} \subseteq \mathbb{R}^n$  be two  $k$ -dimensional subspaces, and let  $U \in \mathbb{R}^{n \times k}$  and  $V \in \mathbb{R}^{n \times k}$  be the basis matrices for  $\mathcal{U}$  and  $\mathcal{V}$ , respectively. As discussed in Sect. 2,  $\{\mathcal{U}, \mathcal{V}\}$  is a pair of deflating subspaces of  $\{K, M\}$  if and only if the equations in (2.3)

hold for some  $K_R \in \mathbb{R}^{k \times k}$  and  $M_R \in \mathbb{R}^{k \times k}$ . In this case,  $\{U, V, K_R, M_R\}$  is an *eigenquadruple*.

But in practice,  $\{U, V\}$  is likely a pair of approximate deflating subspaces in the sense that the *residuals*

$$\mathcal{R}_K(K_R) := KU - VK_R, \quad \mathcal{R}_M(M_R) := MV - UM_R \tag{3.1}$$

are tiny in norm for some  $K_R \in \mathbb{R}^{k \times k}$  and  $M_R \in \mathbb{R}^{k \times k}$ . In this case,  $\{U, V, K_R, M_R\}$  is an *approximate eigenquadruple*. Set

$$H_R = \begin{bmatrix} & K_R \\ M_R & \end{bmatrix}, \tag{3.2}$$

associated with such  $K_R$  and  $M_R$ . Different from  $\{U, V\}$  being exact, now  $\text{eig}(H_R) \not\subset \text{eig}(H)$  but hopefully some or all eigenvalues of  $H_R$  are good approximations to some eigenvalues of  $H$ . Naturally if

$$H_R \hat{z} := \begin{bmatrix} & K_R \\ M_R & \end{bmatrix} \begin{bmatrix} \hat{y} \\ \hat{x} \end{bmatrix} = \lambda \begin{bmatrix} \hat{y} \\ \hat{x} \end{bmatrix} =: \lambda \hat{z},$$

we may take  $\left(\lambda, \begin{bmatrix} V\hat{y} \\ U\hat{x} \end{bmatrix}\right)$  as an approximate eigenpair of  $H$  [1,2] in view of our discussions in the previous section.

In this section, we are interested in answering the following three questions:

1. Given an approximate eigenquadruple, what are the smallest symmetric perturbations  $\Delta K$  and  $\Delta M$  (to  $K$  and  $M$ , respectively) in norm such that the given eigenquadruple is an exact eigenquadruple of  $\{K + \Delta K, M + \Delta M\}$ ?
2. Given a pair of approximate deflating subspaces  $\{U, V\}$ , what are the smallest residuals  $\mathcal{R}_K(K_R)$  and  $\mathcal{R}_M(M_R)$  in norm optimizing among all possible  $K_R$  and  $M_R$ ?
3. It turns out that the so-called Rayleigh quotient pair  $\{K_{RQ}, M_{RQ}\}$  is not the one that minimizes  $\mathcal{R}_K(K_R)$  and  $\mathcal{R}_M(M_R)$  in norm. But how far are  $K_{RQ}$  and  $M_{RQ}$  from their optimal counterparts?

*Remark 3.1* The residuals defined by (3.1) for LREP (1.1) can also be recasted into the residual for the equivalent generalized eigenvalue problem (2.2):

$$\begin{bmatrix} M & \\ & K \end{bmatrix} \begin{bmatrix} V \\ U \end{bmatrix} - \begin{bmatrix} 0 & I_n \\ I_n & 0 \end{bmatrix} \begin{bmatrix} V \\ U \end{bmatrix} \begin{bmatrix} & K_R \\ M_R & \end{bmatrix} = \begin{bmatrix} \mathcal{R}_M(M_R) & \\ & \mathcal{R}_K(K_R) \end{bmatrix}.$$

But this connection appears of no use to us in answering the three questions we just raised because there seems no existing results that can be simply applied to (2.2) with its block structure taken advantage of.

### 3.1 Optimal backward errors

In this subsection, we shall investigate the first question raised at the beginning of the section.

Throughout this subsection,  $\{U, V, K_R, M_R\}$  is assumed an approximate eigenquadruple of  $\{K, M\}$  with  $U, V \in \mathbb{R}^{n \times k}$  satisfying, for convenience,

$$U^T U = V^T V = I_k, \quad \text{and} \quad \text{rank}(U^T V) = k, \quad (3.3)$$

$K_R, M_R \in \mathbb{R}^{k \times k}$ . Define  $\mathcal{R}_K(K_R)$  and  $\mathcal{R}_M(M_R)$  by (3.1),  $H_R$  as in (3.2), and

$$S_K := (U^T V) K_R, \quad S_M := (V^T U) M_R. \quad (3.4)$$

We note that the first condition  $U^T U = V^T V = I_k$  is simply about normalizing the basis matrices for the associated approximate deflating subspaces in question. While it is not essential as far as the approximate deflating subspaces are concerned and not required in (3.1), it removes possible ill-conditioningness in the basis matrices and make optimal backward errors reflect better how good the subspaces are as approximate deflating subspaces. In the case of (2.7), the eigenvalues of the associated  $H_R$  are not affected by this normalization.

**Lemma 3.1** *Factorize  $U^T V$  as  $U^T V = W_1^T W_2$ , where  $W_i \in \mathbb{R}^{k \times k}$  are nonsingular. Then*

$$H_R = \text{diag}(W_2, W_1)^{-1} \begin{bmatrix} 0 & W_1^{-T} S_K W_1^{-1} \\ W_2^{-T} S_M W_2^{-1} & 0 \end{bmatrix} \text{diag}(W_2, W_1). \quad (3.5)$$

*In the case when  $\{U, V, K_R, M_R\}$  is an exact eigenquadruple of  $\{K, M\}$ ,*

$$S_K = U^T K U, \quad S_M = V^T M V. \quad (3.6)$$

*Proof* The Eq. (3.5) can be verified straightforwardly after substituting in  $S_K$  and  $S_M$  as given by (3.4). When  $\{U, V, K_R, M_R\}$  is exact, we can take  $K_R$  and  $M_R$  as in (2.7a). Now use (3.4) to see (3.6).  $\square$

Perturbations  $\Delta K$  and  $\Delta M$  (to  $K$  and  $M$ , respectively) such that the given eigenquadruple  $\{U, V, K_R, M_R\}$  is an exact eigenquadruple of  $\{K + \Delta K, M + \Delta M\}$  are the ones that satisfy

$$(K + \Delta K)U = V K_R, \quad (M + \Delta M)V = U M_R. \quad (3.7)$$

Since  $K$  and  $M$  are symmetric, we further restrict  $\Delta K$  and  $\Delta M$  to be symmetric, too. The first and foremost question is, naturally, if such perturbations  $\Delta K$  and  $\Delta M$  exist, and then if they do, what the smallest perturbations in norm are. For this purpose, we need the following lemma:



**Lemma 3.2** ([36, Lemma 1.4]) *Given  $Z_1, Z_2 \in \mathbb{C}^{n \times k}$ , define*

$$\mathbb{S} = \{S \in \mathbb{C}^{n \times n} : S^H = S, SZ_1 = Z_2\}.$$

1.  $\mathbb{S} \neq \emptyset$  if and only if  $Z_1$  and  $Z_2$  satisfy

$$Z_2 P_{Z_1^H} = Z_2 \quad \text{and} \quad (P_{Z_1} Z_2 Z_1^\dagger)^H = P_{Z_1} Z_2 Z_1^\dagger.$$

2. In the case of  $\mathbb{S} \neq \emptyset$ , any  $S \in \mathbb{S}$  can be expressed by

$$S = Z_2 Z_1^\dagger + (Z_1^\dagger)^H Z_2^H - (Z_1^\dagger)^H Z_2^H P_{Z_1} + P_{Z_1}^\perp T P_{Z_1}^\perp,$$

where  $T \in \mathbb{C}^{n \times n}$  is Hermitian and arbitrary. Moreover,

$$S_{\text{opt}} = Z_2 Z_1^\dagger + (Z_1^\dagger)^H Z_2^H - (Z_1^\dagger)^H Z_2^H P_{Z_1} \in \mathbb{S}$$

is the unique matrix such that

$$\|S_{\text{opt}}\|_F = \min_{S \in \mathbb{S}} \|S\|_F.$$

**Lemma 3.3** *Given approximate eigenquadruple  $\{U, V, K_R, M_R\}$  satisfying (3.3), define*

$$\mathbb{L} := \{(\Delta K, \Delta M) : \Delta K^T = \Delta K \in \mathbb{R}^{n \times n}, \Delta M^T = \Delta M \in \mathbb{R}^{n \times n} \text{ satisfying (3.7)}\}.$$

$\mathbb{L} \neq \emptyset$  if and only if  $S_K$  and  $S_M$  defined by (3.4) are symmetric.

*Proof* We first apply Lemma 3.2 with

$$S = \Delta K, \quad Z_1 = U, \quad Z_2 = V K_R - K U.$$

Notice  $P_{Z_1^H} = U^T U = I$  and

$$P_{Z_1} Z_2 Z_1^\dagger = U U^T (V K_R - K U) U^T = U \left[ \underbrace{(U^T V) K_R}_{S_K} - U^T K U \right] U^T$$

to conclude that  $\Delta K$  exists if and only if  $S_K$  is symmetric, as was to be shown. Next we apply Lemma 3.2 again but with  $S = \Delta M$ ,  $Z_1 = V$ , and  $Z_2 = U M_R - M V$ .  $\square$

As we pointed out in Sect. 2, for any given pair of deflating subspaces, the associated  $\{K_R, M_R\}$  may be expressed in different ways, e.g., the ones in (2.7a) and (2.7b). Now, by Lemma 3.3, it becomes clear that (2.7a) is a good choice for any given  $\{U, V\}$  because it ensures that  $\mathbb{L} \neq \emptyset$  due to (3.3) and (3.4).

In the case of  $\mathbb{L} \neq \emptyset$ , we define the *optimal backward error* by<sup>1</sup>

$$\zeta(U, V, K_R, M_R) := \min_{(\Delta K, \Delta M) \in \mathbb{L}} (\|\Delta K\|_{\text{ui}} + \|\Delta M\|_{\text{ui}}) \quad (3.8)$$

for the given unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ . For any particular unitarily invariant norm, we will attach a suggestive subscript to  $\zeta$  to indicate the norm used, e.g.,  $\zeta_2(U, V, K_R, M_R)$  and  $\zeta_F(U, V, K_R, M_R)$  defined under the spectral norm and the Frobenius norm, respectively.

**Theorem 3.1** *Suppose  $S_K$  and  $S_M$  defined by (3.4) are symmetric. Then*

$$\begin{aligned} \zeta_F(U, V, K_R, M_R) &= \sqrt{2\|\mathcal{R}_K(K_R)\|_F^2 - \|U^T \mathcal{R}_K(K_R)\|_F^2} \\ &\quad + \sqrt{2\|\mathcal{R}_M(M_R)\|_F^2 - \|V^T \mathcal{R}_M(M_R)\|_F^2}, \\ \zeta_2(U, V, K_R, M_R) &= \|\mathcal{R}_K(K_R)\|_2 + \|\mathcal{R}_M(M_R)\|_2, \end{aligned}$$

and for a general unitarily invariant norm,

$$\begin{aligned} \|\mathcal{R}_K(K_R)\|_{\text{ui}} + \|\mathcal{R}_M(M_R)\|_{\text{ui}} &\leq \zeta(U, V, K_R, M_R) \\ &\leq 2 \left[ \|\mathcal{R}_K(K_R)\|_{\text{ui}} + \|\mathcal{R}_M(M_R)\|_{\text{ui}} \right]. \end{aligned} \quad (3.9)$$

*Proof* Note that the minimization for  $\zeta(U, V, K_R, M_R)$  can be separated into the  $K$ -part and  $M$ -part. For the Frobenius norm, we can apply directly Lemma 3.2 with  $Z_1 = U$  and  $Z_2 = -\mathcal{R}_K(K_R)$  to get the optimal  $\Delta K$  as

$$\Delta K_{\text{opt}}(K_R) = [-\mathcal{R}_K(K_R)]U^T + U[-\mathcal{R}_K(K_R)]^T - U[-\mathcal{R}_K(K_R)]^T U U^T, \quad (3.10)$$

whose Frobenius norm is  $\sqrt{2\|\mathcal{R}_K(K_R)\|_F^2 - \|U^T \mathcal{R}_K(K_R)\|_F^2}$ , and similarly for the optimal  $\Delta M$  in the Frobenius norm.

Expand  $U$  to an orthogonal matrix  $[U, U_\perp] \in \mathbb{R}^{n \times n}$  and write

$$\Delta K = [U, U_\perp] \begin{bmatrix} T_{11} & T_{12} \\ T_{21} & T_{22} \end{bmatrix} [U, U_\perp]^T.$$

Since  $\Delta K U = -\mathcal{R}_K(K_R)$  by (3.7), we have  $\begin{bmatrix} T_{11} \\ T_{21} \end{bmatrix} = -\begin{bmatrix} U^T \\ U_\perp^T \end{bmatrix} \mathcal{R}_K(K_R)$ , and thus

$$\Delta K = -[U, U_\perp] \begin{bmatrix} U^T \mathcal{R}_K(K_R) & \mathcal{R}_K(K_R)^T U_\perp \\ U_\perp^T \mathcal{R}_K(K_R) & T_{22} \end{bmatrix} [U, U_\perp]^T,$$

<sup>1</sup> Conceivably, there are other possible ones that are equally appropriate. For example, one may define another *optimal backward error* by replacing  $\|\Delta K\|_{\text{ui}} + \|\Delta M\|_{\text{ui}}$  in (3.8) by  $\sqrt{\|\Delta K\|_{\text{ui}}^2 + \|\Delta M\|_{\text{ui}}^2}$ . Such  $\zeta$  differs from (3.8) within a constant factor.

where  $T_{22}$  is symmetric and arbitrary. Therefore

$$\begin{aligned} \|\Delta K\|_{\text{ui}} &= \left\| \begin{bmatrix} U^T \mathcal{R}_K(K_R) & \mathcal{R}_K(K_R)^T U_\perp \\ U_\perp^T \mathcal{R}_K(K_R) & T_{22} \end{bmatrix} \right\|_{\text{ui}} \geq \left\| \begin{bmatrix} U^T \mathcal{R}_K(K_R) \\ U_\perp^T \mathcal{R}_K(K_R) \end{bmatrix} \right\|_{\text{ui}} \\ &= \|\mathcal{R}_K(K_R)\|_{\text{ui}} \end{aligned}$$

for any  $T_{22}$ . Setting  $T_{22} = 0$ , we have

$$\|\Delta K\|_{\text{ui}} \leq \left\| \begin{bmatrix} U^T \mathcal{R}_K(K_R) \\ U_\perp^T \mathcal{R}_K(K_R) \end{bmatrix} \right\|_{\text{ui}} + \left\| \begin{bmatrix} \mathcal{R}_K(K_R)^T U_\perp \\ 0 \end{bmatrix} \right\|_{\text{ui}} \leq 2\|\mathcal{R}_K(K_R)\|_{\text{ui}}.$$

Similar inequalities hold for the optimal  $\Delta M$ . Together, they yield (3.9).

Finally for the spectral norm, by the dilation theorem of Kreĭn and Kahan (see, e.g., [14, 18] and [38, Theorem 1.2.3]), the optimal  $\Delta K_{\text{opt}}$  in the sense that  $\|\Delta K\|_2$  is smallest as  $T_{22}$  varies among all possible symmetric matrices is

$$\|\Delta K_{\text{opt}}\|_2 = \left\| \begin{bmatrix} U^T \mathcal{R}_K(K_R) \\ U_\perp^T \mathcal{R}_K(K_R) \end{bmatrix} \right\|_2 = \|\mathcal{R}_K(K_R)\|_2,$$

and similarly for the optimal  $\Delta M$  in the spectral norm. □

We remark that the optimal backward perturbation matrices  $\Delta K$  and  $\Delta M$  for the spectral norm and for the Frobenius norm may be different. In particular, the optimal  $\{\Delta K, \Delta M\}$  for the Frobenius norm is unique and can be explicitly stated as by (3.10), while the optimal  $\{\Delta K, \Delta M\}$  for the spectral norm, in general, is not unique and we do not have an explicit expression for it.

### 3.2 Optimal residuals

Theorem 3.1 gives the minimal spectral norm and Frobenius norm for an approximate eigenquadruple of  $\{K, M\}$ . In this subsection, we shall investigate the second question raised at the beginning of the section.

Given a pair of approximate deflating subspaces  $\{\mathcal{U}, \mathcal{V}\}$ , there are many  $K_R \in \mathbb{R}^{k \times k}$  and  $M_R \in \mathbb{R}^{k \times k}$ , e.g., the ones in the form of (3.11) below, which, combined with the basis matrices  $U$  and  $V$  for  $\mathcal{U}$  and  $\mathcal{V}$ , lead to approximate eigenquadruple of  $\{K, M\}$ . Each approximate eigenquadruple gives rise to an optimal backward error  $\zeta(U, V, K_R, M_R)$  as defined by (3.8). A natural question then is how small  $\zeta(U, V, K_R, M_R)$  can get by varying  $K_R$  and  $M_R$ . Because of Lemma 3.3, we will consider only these  $K_R$  and  $M_R$ :

$$K_R = (U^T V)^{-1} S_K \quad \text{and} \quad M_R = (V^T U)^{-1} S_M, \tag{3.11}$$

where  $S_K, S_M \in \mathbb{R}^{k \times k}$  are symmetric.

Our investigation reveals a similar conclusion to that for the standard nonsymmetric eigenvalue problem in [15]: the Rayleigh quotient pair  $\{K_{RQ}, M_{RQ}\}$  in (2.8) does not achieve the minimum in general, but is a reasonably good and computable choice.

We begin by the case for  $k = 1$ . In this case,  $K_R = S_K/(\mathbf{u}^T \mathbf{v})$  is a scalar, and by calculation, the optimal  $K_R$  in the spectral norm (using  $\|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1$ ) is

$$K_R = \mathbf{u}^T K \mathbf{v} \quad (3.12)$$

which is different from

$$K_{RQ} = \mathbf{u}^T K \mathbf{u} / (\mathbf{u}^T \mathbf{v})$$

unless  $\mathbf{u} = \pm \mathbf{v}$  or  $\{\mathcal{R}(\mathbf{u}), \mathcal{R}(\mathbf{v})\}$  is already a pair of deflating subspaces. For the Frobenius norm, simple calculations yield the optimal  $K_R$  as

$$K_R = \frac{2\mathbf{u}^T K \mathbf{v} - \mathbf{u}^T K \mathbf{u} (\mathbf{u}^T \mathbf{v})}{2 - (\mathbf{u}^T \mathbf{v})^2} \quad (3.13)$$

which is not equal to  $K_{RQ}$ , either. It is also noticed that the optimal  $K_R$  in (3.12) for the spectral norm differs from the one in (3.13) for the Frobenius norm.

In general for  $k > 1$ , it seems not easy to derive closed formulas for the optimal  $K_R$  and  $M_R$  with respect to any unitarily invariant norm. But for the Frobenius norm, the closed solutions for the optimal  $K_R$  and  $M_R$  can be derived as Theorem 3.2 below shows.

**Theorem 3.2** *For the Frobenius norm, there is a unique  $\{K_R, M_R\}$  in the form of (3.11) that minimizes  $\zeta_F(U, V, K_R, M_R)$ , and the corresponding  $S_K$  and  $S_M$  are*

$$S_K = 2 \int_0^\infty e^{(I_k - 2Q^T Q)t} (Q^T V^T K U + U^T K V Q - U^T K U) e^{(I_k - 2Q^T Q)t} dt, \quad (3.14a)$$

$$S_M = 2 \int_0^\infty e^{(I_k - 2Q Q^T)t} (Q U^T M V + V^T M U Q^T - V^T M V) e^{(I_k - 2Q Q^T)t} dt, \quad (3.14b)$$

where  $Q = (U^T V)^{-1}$ .

*Proof* First, upon using the facts<sup>2</sup> [12, (6.7) in Chapter 15]:

$$\|Z\|_F^2 = \text{trace}(Z^T Z) \quad \text{and} \quad \frac{\partial \text{trace}(ZS)}{\partial S} = Z + Z^T - \text{Diag}(Z) \quad \text{for} \quad S = S^T,$$

we see that the first order optimality conditions for minimizing

$$2\|\mathcal{R}_K(K_R)\|_F^2 - \|U^T \mathcal{R}_K(K_R)\|_F^2 \quad \text{and} \quad 2\|\mathcal{R}_M(M_R)\|_F^2 - \|V^T \mathcal{R}_M(M_R)\|_F^2$$

<sup>2</sup> For a function  $f$  with a matrix argument  $X \in \mathbb{R}^{n \times m}$ , its partial derivative  $\frac{\partial f(X)}{\partial X} \in \mathbb{R}^{n \times m}$  with its  $(i, j)$ th entry  $\frac{\partial f(X)}{\partial X_{(i,j)}}$  [12, Chapter 15].

are

$$\begin{aligned} & S_K(2Q^T Q - I_k) + (2Q^T Q - I_k)S_K - \text{Diag}(S_K[2Q^T Q - I_k]) \\ &= -2U^T K U + 2(Q^T V^T K U + U^T K V Q) + \text{Diag}(U^T K U - 2Q^T V^T K U), \end{aligned} \tag{3.15a}$$

$$\begin{aligned} & S_M(2Q Q^T - I_k) + (2Q Q^T - I_k)S_M - \text{Diag}(S_M[2Q Q^T - I_k]) \\ &= -2V^T M V + 2(Q U^T M V + V^T M U Q^T) + \text{Diag}(V^T M V - 2Q U^T M V), \end{aligned} \tag{3.15b}$$

where  $\text{Diag}(Z)$  stands for the diagonal matrix whose diagonal entries are those of  $Z$ . Equations in (3.15) are linear matrix equations. Denoting by<sup>3</sup>

$$X = S_K(2Q^T Q - I_k), \quad \widehat{X} = -U^T K U + 2Q^T V^T K U,$$

and noting that

$$\text{Diag}(X) = \text{Diag}\left(\frac{X + X^T}{2}\right) \quad \text{and} \quad \text{Diag}(\widehat{X}) = \text{Diag}\left(\frac{\widehat{X} + \widehat{X}^T}{2}\right),$$

we can rewrite (3.15a) as

$$X + X^T - \text{Diag}(X) = \widehat{X} + \widehat{X}^T - \text{Diag}(\widehat{X}).$$

By comparing the diagonals on both sides, we know  $\text{Diag}(X) = \text{Diag}(\widehat{X})$ , and thus (3.15a) reduces to  $X + X^T = \widehat{X} + \widehat{X}^T$ , or

$$S_K(I_k - 2Q^T Q) + (I_k - 2Q^T Q)S_K = 2(U^T K U - Q^T V^T K U - U^T K V Q), \tag{3.16}$$

which is a Lyapunov equation. Since  $I_k - 2Q^T Q$  is negative definite due to (3.3), we know that (3.16) admits a unique solution which is given by (3.14a). The same argument leads to (3.14b) for  $S_M$ . The proof is completed.  $\square$

Even though (3.14a) and (3.14b) give closed-form solutions, they still involve the integrations over  $t \in [0, \infty)$ . For the special case  $\mathcal{R}(U) = \mathcal{R}(V)$ , we can take  $U = V$  and thus  $Q = I$ , (3.14) yields  $S_K = U^T K U$  and  $S_M = V^T M V$ .

### 3.3 Near optimality of the Rayleigh quotient pair

Previously, we introduced Rayleigh quotient pair  $\{K_{\text{RQ}}, M_{\text{RQ}}\}$ :

$$K_{\text{RQ}} = (U^T V)^{-1} U^T K U \quad \text{and} \quad M_{\text{RQ}} = (V^T U)^{-1} V^T M V. \tag{2.8}$$

<sup>3</sup> This idea of turning (3.15a) into the more “friendly” (3.16) is due to one of the referees.

It is in general not the optimal pair that minimizes  $\zeta(U, V, K_R, M_R)$  in the Frobenius norm and the spectral norm. So for a given pair of approximate deflating subspaces  $\{\mathcal{U}, \mathcal{V}\}$ , there are better pairs  $\{K_R, M_R\}$ , in the sense of giving smaller  $\zeta(U, V, K_R, M_R)$ , than the Rayleigh quotient pair to extract partial spectral information for  $H$  from.

On the other hand, consider

$$H_{\text{RQ}} = \begin{bmatrix} & K_{\text{RQ}} \\ M_{\text{RQ}} & \end{bmatrix}. \quad (2.9)$$

Factorize  $U^T V$  as  $U^T V = W_1^T W_2$ , where  $W_i \in \mathbb{R}^{k \times k}$  are nonsingular. Recall the structure-preserving restriction

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

introduced in [1, 2]. It can be verified that

$$H_{\text{RQ}} = [\text{diag}(W_2, W_1)]^{-1} H_{\text{SR}} [\text{diag}(W_2, W_1)],$$

and thus  $H_{\text{RQ}}$  and  $H_{\text{SR}}$  have the same eigenvalues. In [1, 2], it was the eigenvalues of  $H_{\text{SR}}$ , and thus of  $H_{\text{RQ}}$ , too, that were used to approximate part of the eigenvalues of  $H$ , given the pair of approximate deflating subspaces  $\{\mathcal{U}, \mathcal{V}\}$ . There, it was also proved that such eigenvalue approximation is optimal in the sense of the trace minimization principle obtained there. Therefore the Rayleigh quotient pair must be a reasonably good pair and cannot be too far from the optimal one  $\{K_R, M_R\}$  in the form of (3.11) that minimizes  $\zeta(U, V, K_R, M_R)$ . In this subsection, we will justify such a claim.

Recall our assumptions (3.3), and let  $K_R$  and  $M_R$  be in the form of (3.11), where  $S_K$  and  $S_M$  are symmetric. The angle between  $\mathcal{U} = \mathcal{R}(U)$  and  $\mathcal{V} = \mathcal{R}(V)$  is defined by

$$\theta_{\max}(\mathcal{U}, \mathcal{V}) := \arccos \sigma_{\min}(U^T V),$$

where  $\sigma_{\min}(U^T V)$  is the smallest singular value of  $U^T V$ .

Owing to the definition of  $\zeta(U, V, K_R, M_R)$  and Theorem 3.1, we will focus on minimizing the norms of  $\mathcal{R}_K(K_R)$  and  $\mathcal{R}_M(M_R)$ , separately.

**Lemma 3.4** For any  $K_R, M_R$ , and unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ ,

$$\|K_{\text{RQ}} - K_R\|_{\text{ui}} \leq \alpha \cdot \|\mathcal{R}_K(K_R)\|_{\text{ui}}, \quad (3.17a)$$

$$\|M_{\text{RQ}} - M_R\|_{\text{ui}} \leq \alpha \cdot \|\mathcal{R}_M(M_R)\|_{\text{ui}}, \quad (3.17b)$$

where

$$\alpha = \frac{\sqrt{1 + \sin \theta_{\max}(\mathcal{U}, \mathcal{V})}}{\cos \theta_{\max}(\mathcal{U}, \mathcal{V})}. \quad (3.18)$$

*Proof* We will prove (3.17a) only since (3.17b) can be proved in the same way. Let  $V_{\perp} \in \mathbb{R}^{n \times (n-k)}$  that makes  $[V, V_{\perp}]$  an orthogonal matrix and set  $P = [U, V_{\perp}]$ . Write

$$\begin{aligned} \mathcal{R}_K(K_R) &= KU - VK_R \\ &= (KU - VK_{RQ}) + V(K_{RQ} - K_R) \\ &= \mathcal{R}_K(K_{RQ}) + V(K_{RQ} - K_R). \end{aligned} \tag{3.19}$$

Now using  $U^T \mathcal{R}_K(K_{RQ}) = 0$  and  $V_{\perp}^T V = 0$ , we get

$$P^T \mathcal{R}_K(K_R) = \begin{bmatrix} U^T V(K_{RQ} - K_R) \\ V_{\perp}^T \mathcal{R}_K(K_{RQ}) \end{bmatrix}.$$

Therefore

$$\begin{aligned} \|P^T \mathcal{R}_K(K_R)\|_{\text{ui}} &\geq \|(U^T V)(K_{RQ} - K_R)\|_{\text{ui}} \\ &\geq \sigma_{\min}(U^T V) \cdot \|K_{RQ} - K_R\|_{\text{ui}}, \end{aligned} \tag{3.20}$$

$$\begin{aligned} \|P^T \mathcal{R}_K(K_R)\|_{\text{ui}} &\leq \|P\|_2 \|\mathcal{R}_K(K_R)\|_{\text{ui}} \\ &= \sqrt{1 + \sin \theta_{\max}(\mathcal{U}, \mathcal{V})} \|\mathcal{R}_K(K_R)\|_{\text{ui}}. \end{aligned} \tag{3.21}$$

In deriving (3.21), we have used

$$\begin{aligned} P^T P = \begin{bmatrix} I_k & U^T V_{\perp} \\ V_{\perp}^T U & I_{n-k} \end{bmatrix} &\Rightarrow \|P\|_2^2 = \|P^T P\|_2 = 1 + \|U^T V_{\perp}\|_2 \\ &= 1 + \sin \theta_{\max}(\mathcal{U}, \mathcal{V}). \end{aligned}$$

Combining (3.20) and (3.21), we get (3.17a). □

**Theorem 3.3** For any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ ,

$$\min \|K_{RQ} - K_R\|_{\text{ui}} \leq \alpha \cdot \min \|\mathcal{R}_K(K_R)\|_{\text{ui}}, \tag{3.22a}$$

$$\min \|M_{RQ} - M_R\|_{\text{ui}} \leq \alpha \cdot \min \|\mathcal{R}_M(M_R)\|_{\text{ui}}, \tag{3.22b}$$

and

$$\min \|\mathcal{R}_K(K_R)\|_{\text{ui}} \leq \|\mathcal{R}_K(K_{RQ})\|_{\text{ui}} \leq (1 + \alpha) \cdot \min \|\mathcal{R}_K(K_R)\|_{\text{ui}}, \tag{3.23a}$$

$$\min \|\mathcal{R}_M(M_R)\|_{\text{ui}} \leq \|\mathcal{R}_M(M_{RQ})\|_{\text{ui}} \leq (1 + \alpha) \cdot \min \|\mathcal{R}_M(M_R)\|_{\text{ui}}, \tag{3.23b}$$

where the “min” in (3.22a) and (3.23a) are taken over all  $K_R$  in the form of (3.11) with symmetric  $S_K$ , and the ones in (3.22b) and (3.23b) are taken over all  $M_R$  in the form of (3.11) with symmetric  $S_M$ , and  $\alpha$  is given by (3.18).

*Proof* The inequalities in (3.22) are direct consequences of Lemma 3.4. In what follows, we will prove (3.23a) only since (3.23b) can be proved in the same way. The first inequality in (3.23a) is evident. For the second inequality there, we note

$$\mathcal{R}_K(K_{RQ}) = \mathcal{R}_K(K_R) - V(K_{RQ} - K_R)$$

by (3.19) and thus

$$\begin{aligned} \|\mathcal{R}_K(K_{RQ})\|_{\text{ui}} &\leq \|\mathcal{R}_K(K_R)\|_{\text{ui}} + \|V(K_{RQ} - K_R)\|_{\text{ui}} \\ &= \|\mathcal{R}_K(K_R)\|_{\text{ui}} + \|K_{RQ} - K_R\|_{\text{ui}} \\ &\leq (1 + \alpha)\|\mathcal{R}_K(K_R)\|_{\text{ui}} \end{aligned} \quad (3.24)$$

for all  $K_R$  in the form of (3.11) with symmetric  $S_K$ . Minimizing the right-hand side of (3.24) over all  $S_K$  leads to the second inequality in (3.23a).  $\square$

Although the Rayleigh quotient pair  $\{K_{RQ}, M_{RQ}\}$  may not be optimal in the sense of achieving  $\min \|\mathcal{R}_K(K_R)\|_{\text{ui}}$  and  $\min \|\mathcal{R}_M(K_M)\|_{\text{ui}}$ , respectively, Theorem 3.3 says that it is not too far from the optimal, provided  $\alpha$  is not too big: the smaller  $\alpha$  is, the closer to the optimal the Rayleigh quotient pair  $\{K_{RQ}, M_{RQ}\}$  will be. Note  $\alpha$  is proportional to  $\theta_{\max}(\mathcal{U}, \mathcal{V})$  which, in the limit, approaches  $\theta_{\max}(\mathcal{U}_{\text{exact}}, \mathcal{V}_{\text{exact}})$ , where  $\{\mathcal{U}_{\text{exact}}, \mathcal{V}_{\text{exact}}\}$  is the pair that  $\{\mathcal{U}, \mathcal{V}\}$  is supposed to approximate. So it is an intrinsic quantity of the targeted deflating subspaces.

## 4 Residual-based error bounds for eigenvalues

As preparation, we first cite an eigenvalue perturbation result for a positive definite pencil in Subsect. 4.1 and apply it to LREP (1.1) in Subsect. 4.2, and then come to develop residual based error bounds in Subsect. 4.3. Results in both Subsects. 4.1 and 4.2 are of independent interests on their own from the rest of this article.

In what follows,  $A \succ 0$  means that  $A$  is Hermitian and positive definite.

### 4.1 A perturbation bound for positive definite pencil

Consider a Hermitian matrix pencil  $A - \lambda B$ , where  $A, B \in \mathbb{C}^{n \times n}$  are Hermitian. It is called a *positive definite pencil* if there is a  $\lambda_0 \in \mathbb{R}$  such that  $A - \lambda_0 B \succ 0$  [16, 19, 25].

Suppose that  $A - \lambda B$  is a positive definite pencil and  $B$  is nonsingular. Let  $n_+$  and  $n_-$  be the numbers of positive and negative eigenvalues of  $B$ , respectively. Note  $n_+ + n_- = n$ . It is known [25] that  $A - \lambda B$  has only real eigenvalues which we will divide into two groups  $\{\lambda_i^-\}_{i=1}^{n_-}$  and  $\{\lambda_i^+\}_{i=1}^{n_+}$  and which can be arranged in the order as

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



Moreover,  $A - \lambda B$  is diagonalizable [10,25]: there exists nonsingular  $Z \in \mathbb{C}^{n \times n}$  such that

$$Z^H A Z = \text{diag}(-\Lambda_-, \Lambda_+), \quad Z^H B Z = J := \text{diag}(-I_{n_-}, I_{n_+}), \quad (4.1)$$

where  $\Lambda_{\pm} = \text{diag}(\lambda_1^{\pm}, \lambda_2^{\pm}, \dots, \lambda_{n_{\pm}}^{\pm})$ .

**Lemma 4.1** ([24, Theorem A.2]) *Let  $A - \lambda B$  be a positive definite pencil with nonsingular  $B$  and with the eigen-decomposition (4.1). Suppose it is perturbed to another positive definite pencil  $\tilde{A} - \lambda \tilde{B}$  with nonsingular  $\tilde{B}$ , and adopt the same notations for this perturbed pencil as those for  $A - \lambda B$  except with a tilde on each symbol. Then for any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ ,*

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

where  $\Lambda = \text{diag}(\Lambda_-, \Lambda_+)$  and  $\xi = \max\{\|\Lambda\|_2, \|\tilde{\Lambda}\|_2\}$ .

The concept of positive definite pencil is closely related to that of the so-called *definite pencil* in the past literature [32,34,35] which requires some real linear combination of  $A$  and  $B$  to be positive definite. The latter is more general, encompassing the former. In general,  $B$  may be singular, but Lemma 4.1 excludes the case. When  $B$  is singular, infinite eigenvalues occur. In order to be able to deal with both finite and infinite eigenvalues at the same time, in the literature number pairs  $(\alpha, \beta)$  were used to represent eigenvalues  $\alpha/\beta$  which is finite if  $\beta \neq 0$  and infinite otherwise and the chordal distance was used to measure the difference between two eigenvalues, finite or not. Lemma 4.1 resembles various perturbation bounds in [8,22,23,32,34] for the definite pencils.

#### 4.2 Perturbation bounds for LREP

Consider LREP (1.1) with  $K > 0$  and  $M > 0$ . It is equivalent to the generalized eigenvalue problem for the matrix pencil [1]

$$\mathbf{A} - \lambda \mathbf{B} \equiv \begin{bmatrix} M & \\ & K \end{bmatrix} - \lambda \begin{bmatrix} 0 & I_n \\ I_n & 0 \end{bmatrix}. \quad (4.2)$$

$\mathbf{A} - \lambda \mathbf{B}$  is a positive definite pencil because  $\mathbf{A} - 0 \cdot \mathbf{B} = \mathbf{A} \succ 0$ . Recall Theorem 2.1. We find the eigen-decomposition for  $\mathbf{A} - \lambda \mathbf{B}$ :

$$Z^T \mathbf{A} Z = \text{diag}(\Lambda, \Lambda), \quad Z^T \mathbf{B} Z = \text{diag}(-I_n, I_n), \quad (4.3)$$

where

$$Z = \begin{bmatrix} \Psi \Lambda^{1/2} & \Psi \Lambda^{1/2} \\ -\Phi \Lambda^{-1/2} & \Phi \Lambda^{-1/2} \end{bmatrix}. \quad (4.4)$$

The next theorem bounds  $Z$  and its inverse from above and below.

**Theorem 4.1** For  $Z$  in (4.4) with  $\Lambda$ ,  $\Phi$  and  $\Psi$  defined in Theorem 2.1,

$$2\gamma_1 \leq \|Z\|_2^2 \leq 2\gamma_2, \quad \frac{1}{2}\gamma_1 \leq \|Z^{-1}\|_2^2 \leq \frac{1}{2}\gamma_2, \quad (4.5)$$

where

$$\gamma_1 = \max \left\{ \|M^{-1}\|_2 \lambda_1, \frac{\|M\|_2}{\lambda_n} \right\}, \quad \gamma_2 = \max \left\{ \|M^{-1}\|_2 \lambda_n, \frac{\|M\|_2}{\lambda_1} \right\}.$$

They are also valid if all occurrences of  $M$  are replaced by  $K$ .

*Proof* It can be verified that

$$ZZ^T = 2 \begin{bmatrix} \Psi \Lambda \Psi^T & 0 \\ 0 & \Phi \Lambda^{-1} \Phi^T \end{bmatrix}, \quad Z^{-T} Z^{-1} = \frac{1}{2} \begin{bmatrix} \Phi \Lambda^{-1} \Phi^T & 0 \\ 0 & \Psi \Lambda \Psi^T \end{bmatrix}.$$

Therefore

$$\begin{aligned} \|Z\|_2^2 &\leq 2 \max \left\{ \|\Psi\|_2^2 \lambda_n, \frac{\|\Phi\|_2^2}{\lambda_1} \right\} = 2 \max \left\{ \|M^{-1}\|_2 \lambda_n, \frac{\|M\|_2}{\lambda_1} \right\}, \\ \|Z\|_2^2 &\geq 2 \max \left\{ \|\Psi\|_2^2 \lambda_1, \frac{\|\Phi\|_2^2}{\lambda_n} \right\} = 2 \max \left\{ \|M^{-1}\|_2 \lambda_1, \frac{\|M\|_2}{\lambda_n} \right\}, \\ \|Z^{-1}\|_2^2 &\leq \frac{1}{2} \max \left\{ \frac{\|\Phi\|_2^2}{\lambda_1}, \|\Psi\|_2^2 \lambda_n \right\} = \frac{1}{2} \max \left\{ \frac{\|M\|_2}{\lambda_1}, \|M^{-1}\|_2 \lambda_n \right\}, \\ \|Z^{-1}\|_2^2 &\geq \frac{1}{2} \max \left\{ \frac{\|\Phi\|_2^2}{\lambda_n}, \|\Psi\|_2^2 \lambda_1 \right\} = \frac{1}{2} \max \left\{ \frac{\|M\|_2}{\lambda_n}, \|M^{-1}\|_2 \lambda_1 \right\}, \end{aligned}$$

where we have used  $M = \Phi \Phi^T$  and  $M^{-1} = \Psi \Psi^T$ . Together, they yield (4.5). To see the last claim of this theorem, we let  $\widehat{\Psi} = \Psi \Lambda$  and  $\widehat{\Phi} = \Phi \Lambda^{-1}$ . It can be verified that

$$K = \Psi \Lambda^2 \Psi^T = \widehat{\Psi} \widehat{\Psi}^T, \quad M = \Phi \Phi^T = \widehat{\Phi} \Lambda^2 \widehat{\Phi}^T, \quad Z = \begin{bmatrix} \widehat{\Psi} \Lambda^{-1/2} & \widehat{\Psi} \Lambda^{-1/2} \\ -\widehat{\Phi} \Lambda^{1/2} & \widehat{\Phi} \Lambda^{1/2} \end{bmatrix},$$

and  $K^{-1} = \widehat{\Phi} \widehat{\Phi}^T$ . Following the same lines of argument as above, we see all inequalities in (4.5) are valid if all occurrences of  $M$  are replaced by  $K$ .  $\square$

In the rest of this section, we shall adopt a notational convention: any perturbed quantity is denoted by the same symbol but with a *tilde*. A straightforward application of Lemma 4.1 leads to

**Theorem 4.2** For LREP (1.1) with  $K \succ 0$  and  $M \succ 0$  admitting the decompositions in Theorem 2.1, let  $Z$  be defined by (4.4). Suppose also  $\widetilde{K} \succ 0$  and  $\widetilde{M} \succ 0$ . Then for any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ ,

$$\|\text{diag}(\widetilde{\Lambda}, \widetilde{\Lambda}) - \text{diag}(\Lambda, \Lambda)\|_{\text{ui}} \leq \|Z\|_2 \|\widetilde{Z}\|_2 \|\text{diag}(\widetilde{M}, \widetilde{K}) - \text{diag}(M, K)\|_{\text{ui}}. \quad (4.6)$$

In particular,

$$\max_{1 \leq i \leq n} |\tilde{\lambda}_i - \lambda_i| \leq \|Z\|_2 \|\tilde{Z}\|_2 \max\{\|\tilde{M} - M\|_2, \|\tilde{K} - K\|_2\}, \tag{4.7a}$$

$$\sqrt{\sum_{i=1}^n |\tilde{\lambda}_i - \lambda_i|^2} \leq \frac{1}{\sqrt{2}} \|Z\|_2 \|\tilde{Z}\|_2 \sqrt{\|\tilde{M} - M\|_F^2 + \|\tilde{K} - K\|_F^2}. \tag{4.7b}$$

In the left-hand side of (4.6), the difference between  $\tilde{\Lambda}$  and  $\Lambda$  appears twice. This repetition is handily removed for the spectral and Frobenius norm in (4.7). In general, it is not so easy to remove the repetition without weakening the inequality a little bit. In the corollary below, we show one way of doing it.

**Corollary 4.1** *Under the conditions of Theorem 4.2,*

$$\|\tilde{\Lambda} - \Lambda\|_{\text{ui}} \leq \|Z\|_2 \|\tilde{Z}\|_2 \left[ \|\tilde{M} - M\|_{\text{ui}} + \|\tilde{K} - K\|_{\text{ui}} \right]. \tag{4.8}$$

*Proof* It suffices to show that (4.8) holds for all Ky Fan  $\ell$ -norm  $\|\cdot\|_{(\ell)}$  which is the sum of the  $\ell$  largest singular values of its argument [5, 13, 33].

Let  $\{i_1, i_2, \dots, i_n\}$  be the permutation of  $\{1, 2, \dots, n\}$  such that

$$|\tilde{\lambda}_{i_1} - \lambda_{i_1}| \geq \dots \geq |\tilde{\lambda}_{i_n} - \lambda_{i_n}|.$$

For  $1 \leq \ell \leq n$ , we have

$$\begin{aligned} \|\text{diag}(\tilde{\Lambda}, \tilde{\Lambda}) - \text{diag}(\Lambda, \Lambda)\|_{(2\ell)} &= 2 \sum_{j=1}^{\ell} |\tilde{\lambda}_{i_j} - \lambda_{i_j}| \\ &= 2 \|\tilde{\Lambda} - \Lambda\|_{(\ell)}, \\ \|\text{diag}(\tilde{M}, \tilde{K}) - \text{diag}(M, K)\|_{(2\ell)} &\leq \|\tilde{M} - M\|_{(2\ell)} + \|\tilde{K} - K\|_{(2\ell)} \\ &\leq 2 \left[ \|\tilde{M} - M\|_{(\ell)} + \|\tilde{K} - K\|_{(\ell)} \right]. \end{aligned}$$

By Theorem 4.2, we have

$$\|\tilde{\Lambda} - \Lambda\|_{(\ell)} \leq \|Z\|_2 \|\tilde{Z}\|_2 \left[ \|\tilde{M} - M\|_{(\ell)} + \|\tilde{K} - K\|_{(\ell)} \right],$$

as expected. □

Another way to develop perturbation bounds for  $\lambda_i$  is through noticing the fact that<sup>4</sup> the singular values of  $K^{1/2}M^{1/2}$  are  $0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n$ . This makes a theorem of Mirsky [33, p.204] applicable.

---

<sup>4</sup> This is actually true even for  $K \geq 0$  and  $M \geq 0$ . But in stating Theorem 4.3, we stick to our default assumption on  $K$  and  $M$  just for consistency.

**Theorem 4.3** For LREP (1.1) with  $K \geq 0$  and  $M > 0$ , suppose  $\tilde{K} \geq 0$  and  $\tilde{M} > 0$ . Then for any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ ,

$$\|\tilde{\Lambda} - \Lambda\|_{\text{ui}} \leq \|\tilde{K}^{1/2}\tilde{M}^{1/2} - K^{1/2}M^{1/2}\|_{\text{ui}}. \quad (4.9)$$

One unsatisfactory part of this theorem is that the right-hand side of (4.9) is not explicitly expressed in terms of the norms of  $\tilde{K} - K$  and  $\tilde{M} - M$  whose bounds are usually known. But this can be overcome by writing, e.g.,

$$\tilde{K}^{1/2}\tilde{M}^{1/2} - K^{1/2}M^{1/2} = \tilde{K}^{1/2}(\tilde{M}^{1/2} - M^{1/2}) + (\tilde{K}^{1/2} - K^{1/2})M^{1/2}$$

and then bound the norms of  $\tilde{K}^{1/2} - K^{1/2}$  and  $\tilde{M}^{1/2} - M^{1/2}$  in terms of the norms of  $\tilde{K} - K$  and  $\tilde{M} - M$ , respectively. For example, if  $\tilde{K}^{1/2} + K^{1/2} \geq 2\zeta I_n > 0$ , then [4, 21, 45]

$$\|\tilde{K}^{1/2} - K^{1/2}\|_{\text{ui}} \leq \frac{1}{2\zeta} \|\tilde{K} - K\|_{\text{ui}}.$$

We shall omit the detail of bounding the right-hand side of (4.9) along this line to save space.

We note that Theorem 4.3 requires a weaker condition on  $K$  and  $M$  than Theorem 4.2 does, but there is a consequence: the right-hand side of (4.9) cannot be bounded in terms of the norms of  $\tilde{K} - K$  and  $\tilde{M} - M$  unless both  $\tilde{K}^{1/2} + K^{1/2} > 0$  and  $\tilde{M}^{1/2} + M^{1/2} > 0$ .

*Remark 4.1* We point out that the perturbation results in [6, 7, 28] can be used to derive bounds on the differences  $\tilde{\lambda}_i^2 - \lambda_i^2$  by noticing item (iii) of Theorem 2.1:  $KM$  and  $\tilde{K}\tilde{M}$  are diagonalizable and have real spectra. As an example, by [6, Theorem 3.1], we have the following: For LREP (1.1) with  $K \geq 0$  and  $M > 0$ , suppose  $\tilde{K} \geq 0$  and  $\tilde{M} > 0$ . Then for any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ ,

$$\|\tilde{\Lambda}^2 - \Lambda^2\|_{\text{ui}} \leq [\kappa(M)\kappa(\tilde{M})]^{1/4} \|\tilde{M}\tilde{K} - MK\|_{\text{ui}}. \quad (4.10)$$

Based on this result, different residual based error bounds from those in the next subsection are readily established. But for this article, we shall limit ourselves on bounding the differences  $\tilde{\lambda}_i - \lambda_i$  only, and if needed, residual based error bounds on  $\tilde{\lambda}_i^2 - \lambda_i^2$  can be obtained in the similar way.

### 4.3 Residual based error bounds for LREP

In this subsection, we shall focus on residual based error bounds derivable through Theorem 4.2 and its Corollary 4.1 only and omit those derivable through Theorem 4.3 and (4.10) in Remark 4.1 for two reasons. The first reason is similarity in technicality and the second one is that the bounds may be in different forms but are comparable in sharpness.

Consider an approximate eigenquadruple  $\{U, V, K_R, M_R\}$  of  $\{K, M\}$ , where  $U, V \in \mathbb{R}^{n \times k}$  satisfying, as before,

$$U^T U = V^T V = I_k, \quad \text{and} \quad \text{rank}(U^T V) = k, \tag{3.3}$$

and  $K_R, M_R \in \mathbb{R}^{k \times k}$ , and define  $\mathcal{R}_K(K_R)$  and  $\mathcal{R}_M(M_R)$  by (3.1) and  $H_R$  by (3.2). In Subsect. 3.1, we showed that  $\{U, V, K_R, M_R\}$  of  $\{K, M\}$  is an exact eigenquadruple of

$$\{\tilde{K}, \tilde{M}\} := \{K + \Delta K, M + \Delta M\} \tag{4.11}$$

with bounds in norm on  $\Delta K$  and  $\Delta M$  in terms of the residuals  $\mathcal{R}_K(K_R)$  and  $\mathcal{R}_M(M_R)$ . If the two residuals are sufficiently small, then  $\tilde{K} > 0$  and  $\tilde{M} > 0$  and the eigenvalue problem for the corresponding  $\tilde{H}$  is again an LREP, making all results in Subsect. 4.2 applicable.

**Lemma 4.2** *Suppose  $\|\mathcal{R}_K(K_R)\|_2 < \sigma_{\min}(K)$  and  $\|\mathcal{R}_M(M_R)\|_2 < \sigma_{\min}(M)$ . Then  $H_R$  given in (3.2) is similar to an LREP of  $2k \times 2k$ . Consequently, all eigenvalues of  $H_R$  are real and they come in  $\{\pm\lambda\}$  pairs.*

*Proof* By Theorem 3.1, the approximate exact eigenquadruple  $\{U, V, K_R, M_R\}$  of  $\{K, M\}$  is an exact eigenquadruple of  $\{\tilde{K}, \tilde{M}\}$  as in (4.11) with

$$\|\Delta K\|_2 = \|\mathcal{R}_K(K_R)\|_2 < \sigma_{\min}(K), \quad \|\Delta M\|_2 = \|\mathcal{R}_M(M_R)\|_2 < \sigma_{\min}(M).$$

Now apply Lemma 3.1 to conclude that  $H_R$  is similar to

$$\begin{bmatrix} 0 & W_1^{-T}(U^T \tilde{K} U)W_1^{-1} \\ W_2^{-T}(V^T \tilde{M} V)W_2^{-1} & 0 \end{bmatrix}$$

whose eigenvalue problem is an LREP, where  $W_i$  are as defined in Lemma 3.1.  $\square$

In what follows, whenever  $H_R$  is similar to an LREP of  $2k \times 2k$ , we will denote its eigenvalues by

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

Let  $Z \in \mathbb{R}^{2n \times 2n}$  be the one that diagonalizes  $\mathbf{A} - \lambda \mathbf{B}$  defined in (4.2) and (4.3) and similarly  $\tilde{Z}$  diagonalizes  $\tilde{\mathbf{A}} - \lambda \tilde{\mathbf{B}}$  which is similarly defined in terms of  $\tilde{K}$  and  $\tilde{M}$ .

The appearance of  $\tilde{Z}$  is the unsatisfactory part of the results below since  $\Delta K$  and  $\Delta M$  are usually unknown. But we argue that it does not necessarily invalidate the usefulness of these results. Because for sufficiently small  $\mathcal{R}_K(K_R)$  and  $\mathcal{R}_M(M_R)$  in norm, it is reasonable to expect  $\|\tilde{Z}\|_2 \approx \|Z\|_2$  and the latter can be bounded as in Theorem 4.1.

**Theorem 4.4** *If  $\|\mathcal{R}_K(K_R)\|_2 < \sigma_{\min}(K)$  and  $\|\mathcal{R}_M(M_R)\|_2 < \sigma_{\min}(M)$ , then there are  $k$  positive eigenvalues of  $H$ :*

$$\lambda_{i_1} \leq \cdots \leq \lambda_{i_k}$$

such that

$$\max_{1 \leq j \leq k} |\lambda_{i_j} - \mu_j| \leq \|Z\|_2 \|\tilde{Z}\|_2 \max\{\|\mathcal{R}_K(K_R)\|_2, \|\mathcal{R}_M(M_R)\|_2\}. \quad (4.12)$$

*Proof* The conditions of the theorem ensure that  $\{U, V, K_R, M_R\}$  is an exact eigenquadruple of  $\{\tilde{K}, \tilde{M}\}$  in (4.11) with  $\|\Delta K\|_2 = \|\mathcal{R}_K(K_R)\|_2$  and  $\|\Delta M\|_2 = \|\mathcal{R}_M(M_R)\|_2$ . Thus  $\mu_j$  for  $1 \leq j \leq k$  are among the positive eigenvalues of  $\tilde{H}$ . Let  $\lambda_{i_j}$  be the  $i_j$ th positive eigenvalue of  $\tilde{H}$ . The inequality (4.12) is now a consequence of (4.7a).  $\square$

In a similar way, we can prove

**Theorem 4.5** *If  $\sqrt{2}\|\mathcal{R}_K(K_R)\|_F < \sigma_{\min}(K)$  and  $\sqrt{2}\|\mathcal{R}_M(M_R)\|_F < \sigma_{\min}(M)$ , then there are  $k$  eigenvalues of  $H$ :*

$$\lambda_{i_1} \leq \cdots \leq \lambda_{i_k}$$

such that

$$\begin{aligned} \sqrt{\sum_{1 \leq j \leq k} |\lambda_{i_j} - \mu_j|^2} &\leq \frac{1}{\sqrt{2}} \|Z\|_2 \|\tilde{Z}\|_2 \left[ 2\|\mathcal{R}_K(K_R)\|_F^2 - \|U^T \mathcal{R}_K(K_R)\|_F^2 \right. \\ &\quad \left. + 2\|\mathcal{R}_M(M_R)\|_F^2 - \|V^T \mathcal{R}_M(M_R)\|_F^2 \right]^{1/2}. \end{aligned} \quad (4.13)$$

*Proof* The conditions on  $\|\mathcal{R}_K(K_R)\|_F$  and  $\|\mathcal{R}_M(M_R)\|_F$  ensure that  $\tilde{H}$  defined with the optimal  $\Delta K$  and  $\Delta M$  in the Frobenius norm is an LREP because, by the proof of Theorem 3.1,

$$\begin{aligned} \|\Delta K\|_2 &\leq \|\Delta K\|_F = \sqrt{2\|\mathcal{R}_K(K_R)\|_F^2 - \|U^T \mathcal{R}_K(K_R)\|_F^2} \\ &\leq \sqrt{2}\|\mathcal{R}_K(K_R)\|_F < \sigma_{\min}(K) \end{aligned}$$

and similarly  $\|\Delta M\|_2 < \sigma_{\min}(M)$ . The inequality (4.13) is now a consequence of (4.7b).  $\square$

The conditions on  $\mathcal{R}_K(K_R)$  and  $\mathcal{R}_M(M_R)$  in Theorem 4.5 seem to be stronger than necessary at first sight. It would be more natural to have the same conditions as stated in Theorem 4.4. The thing is that we don't know if  $\|\Delta K\|_2 \leq \|\mathcal{R}_K(K_R)\|_2$  for the optimal  $\Delta K$  in the Frobenius norm while we do know  $\|\Delta K\|_2 = \|\mathcal{R}_K(K_R)\|_2$  for the optimal  $\Delta K$  in the spectral norm. This same reasoning explains the seemingly stronger than necessary conditions in Theorem 4.6 below for any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ .

**Theorem 4.6** *If  $2\|\mathcal{R}_K(K_R)\|_{\text{ui}} < \sigma_{\min}(K)$  and  $2\|\mathcal{R}_M(M_R)\|_{\text{ui}} < \sigma_{\min}(M)$ , then there are  $k$  eigenvalues of  $H$ :*

$$\lambda_{i_1} \leq \dots \leq \lambda_{i_k}$$

such that

$$\|\tilde{\Omega} - \Omega\|_{\text{ui}} \leq 2\|Z\|_2\|\tilde{Z}\|_2\left[\|\mathcal{R}_K(K_R)\|_{\text{ui}} + \|\mathcal{R}_M(M_R)\|_{\text{ui}}\right],$$

where  $\Omega = \text{diag}(\lambda_{i_1}, \dots, \lambda_{i_k})$  and  $\tilde{\Omega} = \text{diag}(\mu_1, \dots, \mu_k)$ .

*Proof* The conditions on  $\|\mathcal{R}_K(K_R)\|_{\text{ui}}$  and  $\|\mathcal{R}_M(M_R)\|_{\text{ui}}$  ensure that  $\tilde{H}$  defined with the optimal  $\Delta K$  and  $\Delta M$  in the unitary invariant norm is an LREP because, by the proof of Theorem 3.1,

$$\|\Delta K\|_2 \leq \|\Delta K\|_{\text{ui}} \leq 2\|\mathcal{R}_K(K_R)\|_{\text{ui}} < \sigma_{\min}(K)$$

and similarly  $\|\Delta M\|_2 < \sigma_{\min}(M)$ . By Corollary 4.1, we have

$$\begin{aligned} \|\tilde{\Omega} - \Omega\|_{\text{ui}} &\leq \|\tilde{\Lambda} - \Lambda\|_{\text{ui}} \\ &\leq \|Z\|_2\|\tilde{Z}\|_2\left[\|\Delta K\|_{\text{ui}} + \|\Delta M\|_{\text{ui}}\right] \\ &\leq 2\|Z\|_2\|\tilde{Z}\|_2\left[\|\mathcal{R}_K(K_R)\|_{\text{ui}} + \|\mathcal{R}_M(M_R)\|_{\text{ui}}\right], \end{aligned}$$

as expected. □

*Remark 4.2* For a given pair of approximate deflating subspaces  $\{\mathcal{U}, \mathcal{V}\}$ , as in [1, 2], likely we will use the associated Rayleigh quotient pair  $\{K_{\text{RQ}}, M_{\text{RQ}}\}$  to make up an approximate eigenquadruple  $\{U, V, K_{\text{RQ}}, M_{\text{RQ}}\}$  of  $\{K, M\}$ . Theorems 4.4, 4.5, and 4.6 can be applied with  $K_R = K_{\text{RQ}}$  and  $M_R = M_{\text{RQ}}$  to arrive at corresponding error bounds on eigenvalue approximations by the eigenvalues of  $H_{\text{RQ}}$  in (2.9).

*Remark 4.3* The error bounds for eigenvalues we have so far are of linear order, i.e., proportional to norms of  $\tilde{K} - K$  and  $\tilde{M} - M$ , or norms of  $K_{\text{RQ}}$  and  $M_{\text{RQ}}$ . Improvements are conceivably possible to derive error bounds that are of the second order, i.e., quadratically dependent on the norms of  $\tilde{K} - K$  and  $\tilde{M} - M$ , or norms of  $K_{\text{RQ}}$  and  $M_{\text{RQ}}$  if certain gap information between those eigenvalues being approximated and the rest is known and positive. In part, this is because the close resemblance of LREP to the symmetric eigenvalue problem for which quadratic residual error bounds are abundant [9, 20, 29, 30, 33]. In fact, Kressner, Pandur, and Shao [17, (19)] presented a quadratic residual error bound for the eigenvalue problem we discussed in Sect. 4.1. Given the rich structure in LREP (1.1), we expect that many results in the aforementioned literature could be extended. But since this paper is already lengthy, we shall pursue a systematic study for quadratic residual error bounds elsewhere.

## 5 Compare with the standard eigenvalue problems

Analogous questions to what we have been investigating so far had been thoroughly studied for the standard eigenvalue problems. Our results here resemble those in the literature. In what follows, we give a brief review on the related results.

Consider the eigenvalue problem for  $C \in \mathbb{C}^{n \times n}$ . Suppose  $X \in \mathbb{C}^{n \times k}$  whose columns span an approximate invariant subspace, i.e.,  $CX \approx XD$  for some  $D \in \mathbb{C}^{k \times k}$  in the sense that

$$\mathcal{R}(D) = CX - XD$$

is relatively small in norm. Let

$$\mathbb{E}_1 := \{E \in \mathbb{C}^{n \times n} : (C + E)X = XD\}.$$

Each  $E$  in  $\mathbb{E}_1$  makes  $\mathcal{R}(X)$  an invariant subspace of  $C + E$  associated with its partial spectrum  $\text{eig}(D)$ . Define the optimal backward error by

$$\eta(X, D) := \min_{E \in \mathbb{E}_1} \|E\|_{\text{ui}}.$$

It is shown [38, Theorem 2.4.2] that

$$\eta(X, D) = \|\mathcal{R}(D)\|_{\text{ui}}.$$

In the other word, for any given approximate eigen-matrix pair  $\{D, X\}$ , the minimal norm  $\|E\|_{\text{ui}}$  among all possible backward perturbations is given by  $\|\mathcal{R}(D)\|_{\text{ui}}$ . The next question is to minimize  $\eta(X, D)$  as  $D$  varies. If also  $X^H X = I_k$ , we have (see, e.g., [38, Theorem 2.4.1], [36, Theorem 2.1], [33, Theorem IV.1.15])

$$\min_{D \in \mathbb{C}^{k \times k}} \|\mathcal{R}(D)\|_{\text{ui}} = \|CX - X(X^H CX)\|_{\text{ui}},$$

i.e., the Rayleigh quotient matrix  $D = X^H CX$  achieves the minimum of  $\|\mathcal{R}(D)\|_{\text{ui}}$  over  $D \in \mathbb{C}^{k \times k}$ , and it is unique for  $\|\cdot\|_{\text{ui}} = \|\cdot\|_{\text{F}}$ .

For Hermitian  $C \in \mathbb{C}^{n \times n}$ , it is often desirable to enforce that  $C + E$  remains Hermitian too. In this case, define [37]

$$\eta(X, D) := \min_{E \in \mathbb{E}_2} \|E\|_{\text{ui}},$$

where

$$\mathbb{E}_2 := \{E = E^H \in \mathbb{C}^{n \times n} : (C + E)X = XD\}.$$



Suppose  $X^H X = I_k$ . It is shown that (a special case of the main theorem in [15, p.478])

$$\eta_2(X, D) := \min_{E \in \mathbb{E}_2} \|E\|_2 = \|\mathcal{R}(D)\|_2,$$

$$\eta_F(X, D) := \min_{E \in \mathbb{E}_2} \|E\|_F = \sqrt{2\|\mathcal{R}(D)\|_F^2 - \|X^H \mathcal{R}(D)\|_F^2},$$

but no closed formula for  $\eta(X, D)$  for a general  $\|\cdot\|_{\text{ui}}$  is known. Moreover, among all Hermitian  $D \in \mathbb{C}^{k \times k}$ , the Rayleigh quotient  $D = X^H C X$  achieves the minimums for  $\eta(X, D)$  for all  $\|\cdot\|_{\text{ui}}$  [36, Theorem 2.2], i.e.,

$$X^H C X = \operatorname{argmin}_D \eta(X, D).$$

Let the eigenvalues of  $C$  and those of  $D = X^H C X$  be

$$\lambda_1 \leq \dots \leq \lambda_n, \quad \mu_1 \leq \dots \leq \mu_k,$$

respectively. As a consequence of well-known perturbation results for Hermitian matrices [33], there exist  $i_1 < i_2 < \dots < i_k$  such that

$$\max_{1 \leq j \leq k} |\lambda_{i_j} - \mu_j| \leq \|\mathcal{R}(D)\|_2,$$

$$\sqrt{\sum_{1 \leq j \leq k} |\lambda_{i_j} - \mu_j|^2} \leq \sqrt{2\|\mathcal{R}(D)\|_F^2 - \|X^H \mathcal{R}(D)\|_F^2}.$$

Similar inequalities in a unitarily invariant norm can be derived, too [36]. Quadratic residual error bounds are also available if certain eigenvalue gap information is known [9, 20, 29, 30, 33].

Finally, for (nonnormal)  $C \in \mathbb{C}^{n \times n}$  with the availability of both left and right approximate invariant subspaces, Kahan, Parlett, and Jiang [15] analyzed the backward perturbation and the residuals for a given quadruple  $\{X_L, X_R, D_L, D_R\}$ , where  $\{D_R, X_R\}$  and  $\{D_L, X_L\}$  are approximate right and left eigen-matrix pairs of  $C$ , respectively, and  $X_L, X_R \in \mathbb{C}^{n \times k}$  have orthonormal columns. Let

$$\mathbb{E}_3 := \{E \in \mathbb{C}^{n \times n} : (C + E)X_R = X_R D_R \text{ and } X_L^H (C + E) = D_L X_L^H\},$$

$$\mathcal{R}_R(D_R) = C X_R - X_R D_R, \quad \mathcal{R}_L(D_L) = X_L^H C - D_L X_L^H.$$

It is shown that [15]  $\mathbb{E}_3 \neq \emptyset$  if and only if  $D_R = (X_L^H X_R)^{-1} D_L (X_L^H X_R)$  and

$$\eta_2(X_L, X_R, D_L, D_R) := \min_{E \in \mathbb{E}_3} \|E\|_2 = \max\{\|\mathcal{R}_L(D_L)\|_2, \|\mathcal{R}_R(D_R)\|_2\},$$

$$\eta_F(X_L, X_R, D_L, D_R) := \min_{E \in \mathbb{E}_3} \|E\|_F$$

$$= \sqrt{\|\mathcal{R}_L(D_L)\|_F^2 + \|\mathcal{R}_R(D_R)\|_F^2 - \|X_L^H \mathcal{R}_R(D_R)\|_F^2}.$$

In this case, the Rayleigh quotient matrices

$$D_{R;RQ} = (X_L^H X_R)^{-1} X_L^H C X_R, \quad D_{L;RQ} = X_L^H C X_R (X_L^H X_R)^{-1}$$

in general do not achieve the minimum of either  $\eta_2$  or  $\eta_F$  any more over all possible  $D_R$  and  $D_L$ .

## 6 Concluding remarks

In approximations for the standard eigenvalue problem, in the past much attention was drawn to investigate the approximation accuracy by a given approximate invariant subspace. Numerous results some of which are reviewed in Sect. 5 have been obtained and can be found in, e.g., [30,33,36,38] and references therein. They are particularly important for today's large scale eigenvalue computations because often it is approximate invariant subspaces that get computed first and then the interested eigenvalues/eigenvectors are then extracted from projecting the problems by approximate invariant subspaces into much smaller eigenvalue problems.

While the linear response eigenvalue problem (1.1) is a standard eigenvalue problem, it has its own block and symmetry structures that are not exploited in the existing theory. Keeping these special structures in mind, in this paper, we have developed a backward perturbation analysis and error bounds for the approximation accuracy of eigenvalues generated by a pair of approximate deflating subspaces or eigenquadruple. Our results are specific for LREP and cannot be derived from the existing ones such as those in Sect. 5, and they are useful for convergence analysis and designing stopping criteria for iterative methods for LREP.

We have assumed so far that  $K$  and  $M$  in LREP (1.1) under investigation are real and symmetric, beside assumptions on their definiteness. We remark that all results are valid for complex Hermitian  $K$  and  $M$ , with the same assumptions on their definiteness, after minor changes: replacing all  $\mathbb{R}$  by  $\mathbb{C}$  and all superscripts  $(\cdot)^T$  by complex conjugate transposes  $(\cdot)^H$ .

**Acknowledgments** The authors are grateful to the editor and two anonymous referees for their careful reading and helpful comments and suggestions, which have improved the paper considerably. The explicit formulas of  $S_K$  and  $S_M$  and their derivations in Theorem 3.2 are in fact due to one of the referees. Zhang was supported in part by the National Natural Science Foundation of China NSFC-11101257, NSFC-11371102, and the Basic Academic Discipline Program, the 11th five year plan of 211 Project for Shanghai University of Finance and Economics. Part of this work is done while Zhang was a visiting scholar at the Department of Mathematics, University of Texas at Arlington from February 2013 to January 2014. Lin was supported in part by National Center of Theoretical Science, Taiwan and ST Yau Center, Chiao Tung University, Taiwan. Li was supported in part by NSF grants DMS-1115834 and DMS-1317330, by a Research Gift Grant from Intel Corporation, and by National Center of Theoretical Science, Taiwan while he visited in December 2013.

## References

1. Bai, Z., Li, R.C.: Minimization principle for linear response eigenvalue problem, I: theory. *SIAM J. Matrix Anal. Appl.* **33**(4), 1075–1100 (2012)

2. Bai, Z., Li, R.C.: Minimization principles for the linear response eigenvalue problem II: computation. *SIAM J. Matrix Anal. Appl.* **34**(2), 392–416 (2013)
3. Benner, P., Mehrmann, V., Xu, H.: Perturbation analysis for the eigenvalue problem of a formal product of matrices. *BIT* **42**(1), 1–43 (2002)
4. Bhatia, R.: Some inequalities for norm ideals. *Commun. Math. Phys.* **111**, 33–39 (1987)
5. Bhatia, R.: Matrix analysis. Graduate texts in mathematics, vol. 169. Springer, New York (1996)
6. Bhatia, R., Kittaneh, F., Li, R.C.: Some inequalities for commutators and an application to spectral variation. II. *Lin. Multilin. Alg.* **43**(1–3), 207–220 (1997)
7. Bhatia, R., Kittaneh, F., Li, R.C.: Eigenvalues of symmetrizable matrices. *BIT* **38**(1), 1–11 (1998)
8. Bhatia, R., Li, R.C.: On perturbations of matrix pencils with real spectra. II. *Math. Comp.* **65**(214), 637–645 (1996)
9. Cao, Z.H., Xie, J.J., Li, R.C.: A sharp version of Kahan’s theorem on clustered eigenvalues. *Linear Algebra Appl.* **245**, 147–155 (1996)
10. Dzensg, D.C., Lin, W.W.: Homotopy continuation method for the numerical solutions of generalised symmetric eigenvalue problems. *J. Austral. Math. Soc. Ser. B* **32**, 437–456 (1991)
11. Granat, R., Kågström, B., Kressner, D.: Computing periodic deflating subspaces associated with a specified set of eigenvalues. *BIT* **43**(1), 1–18 (2003)
12. Harville, D.A.: Matrix Algebra From a Statistician’s Perspective. Springer, New York (1997)
13. Horn, R.A., Johnson, C.R.: Matrix Analysis. Cambridge University Press, Cambridge (1985)
14. Kahan, W.: Inclusion theorems for clusters of eigenvalues of Hermitian matrices. Computer Science Department, University of Toronto, Technical report (1967)
15. Kahan, W., Parlett, B.N., Jiang, E.: Residual bounds on approximate eigensystems of nonnormal matrices. *SIAM J. Numer. Anal.* **19**, 470–484 (1982)
16. Kovač-Striko, J., Veselić, K.: Some remarks on the spectra of Hermitian matrices. *Linear Algebra Appl.* **145**, 221–229 (1991)
17. Kressner, D., Pandur, M.M., Shao, M.: An indefinite variant of LOBPCG for definite matrix pencils. *Numer. Alg.* **66**, 681–703 (2014)
18. Krein, M.G.: The theory of self-adjoint extensions of semi-bounded Hermitian transformations and its applications. *Mat. Sb.* **20**, 431–495 (1947)
19. Lancaster, P., Ye, Q.: Variational properties and Rayleigh quotient algorithms for symmetric matrix pencils. *Oper. Theory: Adv. Appl.* **40**, 247–278 (1989)
20. Li, C.K., Li, R.C.: A note on eigenvalues of perturbed Hermitian matrices. *Linear Algebra Appl.* **395**, 183–190 (2005)
21. Li, R.C.: A perturbation bound for definite pencils. *Linear Algebra Appl.* **179**, 191–202 (1993)
22. Li, R.C.: On perturbations of matrix pencils with real spectra. *Math. Comp.* **62**, 231–265 (1994)
23. Li, R.C.: On perturbations of matrix pencils with real spectra, a revisit. *Math. Comp.* **72**, 715–728 (2003)
24. Liang, X., Li, R.C.: The hyperbolic quadratic eigenvalue problem. Technical Report 2014–01, Department of Mathematics, University of Texas at Arlington. [www.uta.edu/math/preprint/](http://www.uta.edu/math/preprint/) (2014)
25. Liang, X., Li, R.C., Bai, Z.: Trace minimization principles for positive semi-definite pencils. *Linear Algebra Appl.* **438**, 3085–3106 (2013)
26. Lin, W.W., Sun, J.G.: Perturbation analysis for the eigenproblem of periodic matrix pairs. *Linear Algebra Appl.* **337**(13), 157–187 (2001)
27. Lin, W.W., van Dooren, P., Xu, Q.F.: Equivalent characterizations of periodical invariant subspaces. NCTS Preprints Series 1998–8, National Center for Theoretical Sciences, Math. Division, National Tsing Hua University, Hsinchu, Taiwan (1998)
28. Lu, T.X.: Perturbation bounds of eigenvalues of symmetrizable matrices. *Numer. Math. J. Chin. Univ.* **16**, 177–185 (1994). In Chinese
29. Mathias, R.: Quadratic residual bounds for the Hermitian eigenvalue problem. *SIAM J. Matrix Anal. Appl.* **19**, 541–550 (1998)
30. Parlett, B.N.: The Symmetric Eigenvalue Problem. SIAM, Philadelphia (1998)
31. Stewart, G.W.: On the sensitivity of the eigenvalue problem  $Ax = \lambda Bx$ . *SIAM J. Numer. Anal.* **4**, 669–686 (1972)
32. Stewart, G.W.: Perturbation bounds for the definite generalized eigenvalue problem. *Linear Algebra Appl.* **23**, 69–86 (1979)
33. Stewart, G.W., Sun, J.G.: Matrix Perturbation Theory. Academic Press, Boston (1990)

34. Sun, J.G.: A note on Stewart's theorem for definite matrix pairs. *Linear Algebra Appl.* **48**, 331–339 (1982)
35. Sun, J.G.: Perturbation bounds for eigenspaces of a definite matrix pair. *Numer. Math.* **41**, 321–343 (1983)
36. Sun, J.G.: Backward perturbation analysis of certain characteristic subspaces. *Numer. Math.* **65**, 357–382 (1993)
37. Sun, J.G.: A note on backward perturbations for the Hermitian eigenvalue problem. *BIT* **35**, 385–393 (1995)
38. Sun, J.G.: Stability and accuracy: perturbation analysis of algebraic eigenproblems. Technical report UMINF 1998–07, ISSN-0348-0542, Faculty of Science and Technology, Department of Computing Science, Umeå University (1998)
39. Teng, Z., Li, R.C.: Convergence analysis of Lanczos-type methods for the linear response eigenvalue problem. *J. Comput. Appl. Math.* **247**, 17–33 (2013)
40. Teng, Z., Zhou, Y., Li, R.C.: A block Chebyshev-Davidson method for linear response eigenvalue problems. Technical Report 2013–11, Department of Mathematics, University of Texas at Arlington. [www.uta.edu/math/preprint/](http://www.uta.edu/math/preprint/) (2013)
41. Thouless, D.J.: Vibrational states of nuclei in the random phase approximation. *Nucl. Phys.* **22**(1), 78–95 (1961)
42. Thouless, D.J.: *The quantum mechanics of many-body systems*. Academic Press, New York (1972)
43. Tsiper, E.V.: Variational procedure and generalized Lanczos recursion for small-amplitude classical oscillations. *JETP Lett.* **70**(11), 751–755 (1999)
44. Tsiper, E.V.: A classical mechanics technique for quantum linear response. *J. Phys. B At. Mol. Opt. Phys.* **34**(12), L401–L407 (2001)
45. van Hemmen, J.L., Ando, T.: An inequality for trace ideals. *Commun. Math. Phys.* **76**, 143–148 (1980)
46. Zhang, L.H., Xue, J., Li, R.C.: Rayleigh-Ritz approximation for the linear response eigenvalue problem. *SIAM J. Matrix Anal. Appl.* **35**, 765–782 (2014)