

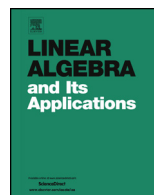


ELSEVIER

Contents lists available at ScienceDirect

Linear Algebra and its Applications

www.elsevier.com/locate/laa



On generalizing trace minimization principles

Xin Liang^{a,b,1}, Li Wang^{c,2}, Lei-Hong Zhang^{d,3}, Ren-Cang Li^{c,e,*,4}^a *Yau Mathematical Sciences Center, Tsinghua University, Beijing 100084, China*^b *Yanqi Lake Beijing Institute of Mathematical Sciences and Applications, Beijing 101408, China*^c *Department of Mathematics, University of Texas at Arlington, Arlington, TX 76019-0408, USA*^d *School of Mathematical Sciences, Soochow University, Suzhou 215006, Jiangsu, China*^e *Department of Mathematics, Hong Kong Baptist University, Kowloon Tong, Kowloon, Hong Kong*

ARTICLE INFO

Article history:

Received 20 May 2021

Accepted 11 October 2022

Available online 13 October 2022

Submitted by V. Mehrmann

MSC:

15A18

15A22

65F15

65K99

90C20

90C26

Keywords:

Trace minimization principle

Eigenvalue

Eigenvector

ABSTRACT

Various trace minimization principles have interplayed with numerical computations for the standard eigenvalue and generalized eigenvalue problems in general, as well as important applied eigenvalue problems including the linear response eigenvalue problem from electronic structure calculation and the symplectic eigenvalue problem of positive definite matrices that play important roles in classical Hamiltonian dynamics, quantum mechanics, and quantum information, among others. In this paper, Ky Fan's trace minimization principle is extended along the line of the Brockett cost function $\text{tr}(DX^HAX)$ in X on the Stiefel manifold, where D of an apt size is positive definite. Specifically, we investigate $\inf_X \text{tr}(DX^HAX)$ subject to $X^HBX = I_k$ (the $k \times k$ identity matrix) or $X^HBX = J_k$, where $J_k = \text{diag}(\pm 1)$. We establish conditions under which the infimum is finite and when it is finite, analytic solutions are obtained in terms of the eigen-

* Corresponding author.

E-mail addresses: liangxinslm@tsinghua.edu.cn (X. Liang), li.wang@uta.edu (L. Wang), longzh@suda.edu.cn (L.-H. Zhang), rcli@uta.edu (R.-C. Li).

¹ Supported in part by the National Natural Science Foundation of China NSFC-11901340 and by Tsinghua University Initiative Scientific Research Program.² Supported in part by NSF DMS-2009689.³ Supported in part by the National Natural Science Foundation of China NSFC-12071332.⁴ Supported in part by NSF grants DMS-1719620 and DMS-2009689.

Brockett cost function
 Linear response eigenvalue problem
 Symplectic eigenvalue problem of
 positive definite matrix

values and eigenvectors of the matrix pencil $A - \lambda B$, where B is possibly indefinite and possibly singular, and D is also possibly indefinite.

© 2022 Elsevier Inc. All rights reserved.

1. Introduction

Quadratic optimization problems with matrix arguments are drawing tremendous attentions lately in data science, and they often involve traces of certain quadratic forms, for example, the trace ratio maximization problem from the linear discriminant analysis (LDA) [1–6], the correlation maximization from the canonical correlation analysis (CCA) and its variants [7–11]. Trace minimizations are special quadratic optimization formulations that have played important interconnecting roles between theory and numerical computations. Some formulations admit analytical solutions in terms of matrix eigenvalue/singular value decompositions, but most don't. Among those that do admit analytical solutions, the most well-known one is perhaps Ky Fan's trace minimization principle [12] [13, p. 248]

$$\min_{X^H X = I_k} \text{tr}(X^H A X) = \sum_{i=1}^k \lambda_i, \tag{1.1}$$

where $\text{tr}(\cdot)$ is the trace of a square matrix, I_k is the $k \times k$ identity matrix, and $A \in \mathbb{C}^{n \times n}$ is Hermitian and its eigenvalues are denoted by λ_i ($i = 1, 2, \dots, n$) and arranged in the ascending order:

$$\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n. \tag{1.2}$$

Moreover for any minimizer X_{opt} of (1.1), i.e., $\text{tr}(X_{\text{opt}}^H A X_{\text{opt}}) = \sum_{i=1}^k \lambda_i$, its columns span A 's invariant subspace associated with the first k eigenvalues λ_i , $i = 1, 2, \dots, k$. The minimization principle (1.1) can be straightforwardly extended to

$$\min_{X^H B X = I_k} \text{tr}(X^H A X) = \sum_{i=1}^k \lambda_i, \tag{1.3}$$

where $A, B \in \mathbb{C}^{n \times n}$ are Hermitian and B is positive definite, and now λ_i are the eigenvalues of matrix pencil $A - \lambda B$ and ordered as in (1.2). Essentially, (1.3) is no more general than (1.1). In fact, upon substitutions: $X \leftarrow B^{1/2} X$ and $A \leftarrow B^{-1/2} A B^{-1/2}$, (1.3) reduces to (1.1). Both (1.1) and (1.3) have played roles of bridges between elegant mathematical theory and efficient numerical computations for large scale eigenvalue problems (see, e.g., [14–21] and references therein).

A nontrivial extension of (1.1) and (1.3) is [22,23]

$$\inf_{X^H B X = J_k} \text{tr}(X^H A X) = \sum_{i=1}^{k_+} \lambda_i^+ - \sum_{i=1}^{k_-} \lambda_i^- \tag{1.4}$$

for a positive semi-definite matrix pencil⁵ $A - \lambda B$, where λ_i^\pm are finite eigenvalues of $A - \lambda B$ and are arranged in the order as

$$\lambda_{n_-}^- \leq \dots \leq \lambda_1^- \leq \lambda_1^+ \leq \dots \leq \lambda_{n_+}^+, \tag{1.5}$$

n_- and n_+ are the numbers of negative and positive eigenvalues of B , respectively, and

$$0 \leq k_\pm \leq n_\pm, \quad k = k_+ + k_-, \quad J_k = \begin{bmatrix} I_{k_+} & \\ & -I_{k_-} \end{bmatrix} \in \mathbb{C}^{k \times k}. \tag{1.6}$$

The result (1.3) is immediately applicable to two important applied eigenvalue problems: one is the linear response eigenvalue problem from electronic structure calculation [24–26], and the other is the symplectic eigenvalues of positive definite matrices that play important roles in classical Hamiltonian dynamics, quantum mechanics, and quantum information, among others [27,28]. These applications will be outlined in appendix A to quickly derive two known trace minimization principles for the two eigenvalue problems, which have played foundational roles of bridges between theory and numerical computations for a few applied eigenvalue problems [24,25,29–32], respectively. Based on the principle (1.4), extensions of the eigensolver LOBPCG [16] to general definite pencils [33,34] were also worked out.

Recently, Liu, So, and Wu [35] laboriously analyzed how to solve

$$\min_{X^H X = I_k} \text{tr}(D X^H A X) \tag{1.7}$$

by numerical optimization techniques, where $A \in \mathbb{C}^{n \times n}$ and $D \in \mathbb{C}^{k \times k}$ are Hermitian but D may be indefinite. Its objective function $\text{tr}(D X^H A X)$ in X on the Stiefel manifold $\{X \in \mathbb{R}^{n \times k} : X^H X = I_k\}$ is known as the *Brockett cost function* in the case when D is diagonal and positive semi-definite, and optimizing it has often been used as an illustrative example for optimization on the Stiefel manifold [36, p. 80], [37,38].

Our goal in this paper is to go beyond the Brockett cost function to investigate, as an extension of (1.1),

$$\inf_X \text{tr}(D X^H A X), \tag{1.8}$$

subject to $X^H B X = I_k$ or $-I_k$ or J_k , where both B and D are possibly indefinite. Our first main result is an analytical solution to (1.8) for positive definite B (for which only

⁵ $A, B \in \mathbb{C}^{n \times n}$ are Hermitian and there exists $\lambda_0 \in \mathbb{R}$ such that $A - \lambda_0 B$ is positive semi-definite [23, Definition 1.1]. A brief review for the spectral properties of a positive semi-definite matrix pencil will be given at the beginning of section 3.

$X^H B X = I_k$ can be used as a constraint) in terms of the eigen-decompositions of D and matrix pencil $A - \lambda B$ and the solution lends itself to be computed by more efficient numerical linear algebra techniques [39–41,18]. In particular, this result yields an elegant solution to the widely studied (1.7) [36,35]. Our second main result is for a more general setting that B is indefinite and possibly singular and $A - \lambda B$ is a positive semi-definite matrix pencil. We show that the infimum in (1.8), subject to $X^H B X = I_k$ or $-I_k$, is finite if and only if D is positive semi-definite and establish analytical solutions to it when the infimum is finite.

Note that under the constraint either $X^H B X = I_k$ or $X^H B X = -I_k$, whether D is diagonal or not is inconsequential so long that it is Hermitian because we can always perform an eigen-decomposition $D = Q \Omega Q^H$ to get

$$\min_{X^H B X = \pm I_k} \operatorname{tr}(D X^H A X) = \min_{\tilde{X}^H B \tilde{X} = \pm I_k} \operatorname{tr}(\Omega \tilde{X}^H A \tilde{X}),$$

where X and \tilde{X} are related by $\tilde{X} = X Q$, Q is unitary and Ω is the diagonal matrix of the eigenvalues of D . So far, we have been focusing on “minimization”, but these formulations admit straightforward restatements for “maximization” by simply considering $-A$ instead.

The rest of this paper is organized as follows. We state our main results for (1.8) in section 2 for positive definite B (and with $X^H B X = I_k$) and in section 3 for the more general setting that B is genuinely indefinite. The proofs for the main results are presented in sections 4 and 5, respectively. We draw our conclusion in section 6. In appendix A, we demonstrate how easy it is to apply the trace minimization principle (1.4) to two applied eigenvalue problems: the linear response eigenvalue problem and the symplectic eigenvalue problem of positive definite matrices.

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. For a matrix $X \in \mathbb{C}^{m \times n}$, $\mathcal{N}(X) = \{x \in \mathbb{C}^n : Xx = 0\}$ and $\mathcal{R}(X) = \{Xx : x \in \mathbb{C}^n\}$ are the null space and the range of X , respectively. X^T and X^H are the transpose and the conjugate transpose of a vector or matrix, respectively. $A \succ 0$ ($A \succeq 0$) means that A is Hermitian positive (semi-)definite, and $A \prec 0$ ($A \preceq 0$) if $-A \succ 0$ ($-A \succeq 0$). $A^{1/2}$ is the unique positive semi-definite square root of a positive semi-definite matrix $A \succeq 0$.

2. Positive definite B

Throughout this section and section 4, $A, B \in \mathbb{C}^{n \times n}$ and $D \in \mathbb{C}^{k \times k}$ are Hermitian and B is positive definite. Then $A - \lambda B$ admits the following eigen-decomposition

$$U^H A U = \Lambda \equiv \operatorname{diag}(\lambda_1, \lambda_2, \dots, \lambda_n), \quad U^H B U = I_n, \tag{2.1}$$

where λ_i are the eigenvalues of $A - \lambda B$ and are, without loss of generality, arranged in the ascending order as in (1.2), and $U = [u_1, u_2, \dots, u_n]$ is the eigenvector matrix and B -unitary: $Au_i = \lambda_i Bu_i$ for all i and $U^H B U = I_n$. Let the eigen-decomposition of D be

$$Q^H D Q = \Omega \equiv \text{diag}(\omega_1, \omega_2, \dots, \omega_k), \tag{2.2a}$$

where $Q \in \mathbb{C}^{k \times k}$ is unitary and, without loss of generality,

$$\omega_1 \geq \dots \geq \omega_\ell \geq 0 \geq \omega_{\ell+1} \geq \dots \geq \omega_k. \tag{2.2b}$$

The case $\ell = 0$ or $\ell = k$ corresponds to when D has no positive eigenvalues, i.e., $D \preceq 0$, or no nonnegative eigenvalues, i.e., $D \succeq 0$, respectively.

Our main result of this section is Theorem 2.1 below.

Theorem 2.1. *Suppose that $A, B \in \mathbb{C}^{n \times n}$ and $D \in \mathbb{C}^{k \times k}$ are Hermitian and B is positive definite, admitting the eigen-decompositions in (2.1) and (2.2). Then*

$$\min_{X^H B X = I_k} \text{tr}(D X^H A X) = \sum_{i=1}^{\ell} \omega_i \lambda_i + \sum_{i=\ell+1}^k \omega_i \lambda_{i+n-k}. \tag{2.3}$$

Furthermore, any minimizer X_{opt} has the following characterizations:

- (a) *If D is nonsingular, then $\mathcal{R}(X_{\text{opt}} Q)$ is the eigenspace of $A - \lambda B$ [42, p. 303], associated with the ℓ smallest and $k - \ell$ largest eigenvalues of $A - \lambda B$. If also all ω_i are distinct, then*

$$(X_{\text{opt}} Q)^H A X_{\text{opt}} Q = \text{diag}(\underbrace{\lambda_1, \lambda_2, \dots, \lambda_\ell}_{\ell}, \underbrace{\lambda_{n-k+\ell+1}, \dots, \lambda_n}_{k-\ell}).$$

- (b) *Suppose that D is possibly singular and has ℓ_+ positive eigenvalues and ℓ_- negative eigenvalues, and let $\widehat{Q} \in \mathbb{C}^{k \times (\ell_+ + \ell_-)}$ be the one obtained from Q by keeping its first ℓ_+ and last ℓ_- columns. Then $\mathcal{R}(X_{\text{opt}} \widehat{Q})$ is the eigenspace of $A - \lambda B$ associated with its ℓ_+ smallest and ℓ_- largest eigenvalues. If also the nonzero eigenvalues of D are distinct, then*

$$(X_{\text{opt}} \widehat{Q})^H A X_{\text{opt}} \widehat{Q} = \text{diag}(\underbrace{\lambda_1, \lambda_2, \dots, \lambda_{\ell_+}}_{\ell_+}, \underbrace{\lambda_{n-\ell_-+1}, \dots, \lambda_n}_{\ell_-}).$$

There are a couple of remarks in order. Firstly, the minimization extracts out the extreme eigenvalues of $A - \lambda B$ from both ends. Secondly, if all ω_i are distinct and nonzero, then the columns of $X_{\text{opt}} Q$ are the associated eigenvectors. Thirdly, if D does have 0 as some of its eigenvalues, then in the notation of Theorem 2.1(b), those eigenvalues 0 can be matched to any λ_i ($\ell_+ + 1 \leq i \leq n - \ell_-$), other things being equal, to still yield

the same objective value as the optimal one in the right-hand side of (2.3). Fourthly, upon replacing A by $-A$, we can obtain immediately parallel results on maximizing $\text{tr}(DX^HAX)$ subject to $X^HBX = I_k$. The proof of Theorem 2.1 occupies a few pages and is deferred to section 4.

3. Genuinely indefinite B

Throughout this section and section 5, $A - \lambda B \in \mathbb{C}^{n \times n}$ is a positive semi-definite matrix pencil, i.e., A and B are Hermitian and there exists $\lambda_0 \in \mathbb{R}$ such that $A - \lambda_0 B \succeq 0$, and B is genuinely indefinite in the sense that B has both positive and negative eigenvalues. We are interested in (1.8):

$$\inf_X \text{tr}(DX^HAX) \quad \text{subject to } X^HBX = I_k \text{ or } -I_k \text{ or } J_k, \tag{1.8}$$

where J_k is as given in (1.6).

Before we investigate (1.8), we review some of the related concepts and results about a positive semi-definite matrix pencil $A - \lambda B$ [23]. Let the integer triplet (n_+, n_0, n_-) be the *inertia* of B , meaning B has n_+ positive, n_0 zero, and n_- negative eigenvalues, respectively. Necessarily

$$r := \text{rank}(B) = n_+ + n_-. \tag{3.1}$$

We say $\mu \neq \infty$ is a *finite eigenvalue* of $A - \lambda B$ if

$$\text{rank}(A - \mu B) < \max_{\lambda \in \mathbb{C}} \text{rank}(A - \lambda B), \tag{3.2}$$

and $x \in \mathbb{C}^n$ is a corresponding *eigenvector* if $x \notin \mathcal{N}(A) \cap \mathcal{N}(B)$ satisfies

$$Ax = \mu Bx, \tag{3.3}$$

or equivalently, $x \in \mathcal{N}(A - \mu B) \setminus (\mathcal{N}(A) \cap \mathcal{N}(B))$. It is known [23] that a *positive semi-definite pencil* $A - \lambda B$ has only $r = \text{rank}(B)$ finite eigenvalues all of which are real. Denote these finite eigenvalues by λ_i^\pm ordered as (1.5). It has been proved that for all i, j

$$\lambda_i^- \leq \lambda_0 \leq \lambda_j^+. \tag{3.4}$$

As in section 2, let D have its eigen-decomposition given by (2.2):

$$Q^H D Q = \Omega \equiv \text{diag}(\omega_1, \omega_2, \dots, \omega_k), \quad \omega_1 \geq \omega_2 \geq \dots \geq \omega_k.$$

Our first main result of the section is Theorem 3.1 below.

Theorem 3.1. *Suppose that $A, B \in \mathbb{C}^{n \times n}$ and $D \in \mathbb{C}^{k \times k}$ are Hermitian, $A \neq 0$ and B is genuinely indefinite, $k \leq n_+$, and the matrix pencil $A - \lambda B$ is positive semi-definite. Then*

$$\inf_{X^H B X = I_k} \operatorname{tr}(D X^H A X) > -\infty$$

if and only if $D \succeq 0$, in which case

$$\inf_{X^H B X = I_k} \operatorname{tr}(D X^H A X) = \sum_{i=1}^k \omega_i \lambda_i^+. \tag{3.5}$$

The infimum can be attained, when $A - \lambda B$ is diagonalizable, by X such that the columns of XQ are the eigenvectors of $A - \lambda B$ associated with its eigenvalues λ_i^+ for $1 \leq i \leq k$, respectively.

Our proof of this theorem is rather involved and will be given in section 5. Apply Theorem 3.1 to the matrix pencil $A - (-\lambda)(-B)$, we immediately conclude a parallel result on minimizing $\operatorname{tr}(D X^H A X)$ subject to $X^H B X = -I_k$. Our second main result stated in Theorem 3.2 below is for the more general constraint $X^H B X = J_k$, whose proof is deferred to section 5 as well.

Theorem 3.2. *Suppose that $A, B \in \mathbb{C}^{n \times n}$ and $D_{\pm} \in \mathbb{C}^{k_{\pm} \times k_{\pm}}$ are Hermitian, $A \neq 0$ and B is genuinely indefinite, $k_{\pm} \leq n_{\pm}$, and the matrix pencil $A - \lambda B$ is positive semi-definite. Let*

$$J_k = \begin{bmatrix} I_{k_+} & \\ & -I_{k_-} \end{bmatrix}, \quad D = \begin{matrix} k_+ & k_- \\ k_- & \end{matrix} \begin{bmatrix} D_+ & \\ & D_- \end{bmatrix},$$

and denote by $\omega_1^+ \geq \dots \geq \omega_{k_+}^+$ and $\omega_1^- \geq \dots \geq \omega_{k_-}^-$ the eigenvalues of D_+ and D_- , respectively. If $D_{\pm} \succeq 0$, then

$$\inf_{X^H B X = J_k} \operatorname{tr}(D X^H A X) = \sum_{i=1}^{k_+} \omega_i^+ \lambda_i^+ - \sum_{i=1}^{k_-} \omega_i^- \lambda_i^-. \tag{3.6}$$

The infimum can be attained when $A - \lambda B$ is diagonalizable.

Equation (3.6) for $k_- = 0$ reduces to (3.5). Comparing Theorem 3.2 with Theorem 3.1, one may be tempted to conjecture that in Theorem 3.2 if both $D_{\pm} \succeq 0$ is also a necessary condition for $\inf_{X^H B X = J_k} \operatorname{tr}(D X^H A X) > -\infty$, a question whose answer eludes us.

One important comment that we would like to emphasize about the conditions of Theorem 3.2 is that matrix D has to take the same block-diagonal structure as J_k . In

fact, Example 3.1 below shows that if D doesn't have the same block-diagonal structure as J_k , the infimum may not be able to be expressed simply as some sum of the products of the eigenvalues between D and $A - \lambda B$.

Example 3.1. Let $\mu = 2$ and $\delta = 1/4$, and

$$\gamma = \frac{1 - \delta}{1 + \delta} = \frac{3}{5}, \quad \nu = \frac{1 - \mu}{1 + \mu} = -\frac{1}{3},$$

$$\sigma = \sqrt{\frac{1}{2} \left(\frac{\nu}{\gamma} \sqrt{\frac{1 - \gamma^2}{1 - \nu^2}} + 1 \right)} = \frac{\sqrt{18 - 6\sqrt{2}}}{6} \approx .5141.$$

Consider

$$A = \begin{bmatrix} 1 & \\ & \mu \end{bmatrix}, \quad B = \begin{bmatrix} 1 & \\ & -1 \end{bmatrix}, \quad J_2 = B, \tag{3.7a}$$

$$\Omega = \begin{bmatrix} 1 & \\ & \delta \end{bmatrix} \succ 0, \quad Q = \begin{bmatrix} \sqrt{1 - \sigma^2} & -\sigma \\ \sigma & \sqrt{1 - \sigma^2} \end{bmatrix}, \quad D = Q^H \Omega Q \succ 0. \tag{3.7b}$$

$A - \lambda B$ is positive definite pencil because $A - 0 \cdot B = A \succ 0$. The two eigenvalues of $A - \lambda B$ are $\lambda_1^- = -\mu$, $\lambda_1^+ = 1$. D doesn't have the same block structure as J_2 since $\sigma \neq 0$. We claim that

$$\inf_{X^H B X = J_2} \text{tr}(D X^H A X) < \min\{1 + \delta\mu, \mu + \delta\} = 1 + \delta\mu = \frac{3}{2}, \tag{3.8}$$

implying that the infimum cannot be simply expressed as any of the two possible sums of products: $1 + \delta\mu$ and $\mu + \delta$, between the eigenvalues of $A - \lambda B$ and of D . To see (3.8), we consider

$$t = \sqrt{\frac{1}{2} \left(\sqrt{\frac{1 - \nu^2}{1 - \gamma^2}} - 1 \right)} = \frac{15\sqrt{2} - 18}{6} \approx .2988, \quad Y = \begin{bmatrix} \sqrt{1 + t^2} & \\ t & \sqrt{1 + t^2} \end{bmatrix}.$$

It can be verified that $Y^H B Y = J_2$, and hence

$$\inf_{X^H B X = J_2} \text{tr}(D X^H A X) \leq \text{tr}(D Y^H A Y) = \sqrt{2} \approx 1.4142 < \frac{3}{2},$$

as expected.

4. Proof of Theorem 2.1

We start with three lemmas as preparation. The first lemma is about a result from majorization [43,13]. Given two multisets of real numbers $\{\alpha_i\}_{i=1}^m$ and $\{\beta_i\}_{i=1}^m$, we say that $\{\beta_i\}_{i=1}^m$ majorizes $\{\alpha_i\}_{i=1}^m$ if

$$\sum_{i=1}^j \alpha_i^\downarrow \leq \sum_{i=1}^j \beta_i^\downarrow, \quad \text{for } j = 1, 2, \dots, m$$

with equality holds for $j = m$, where $\{\alpha_i^\downarrow\}_{i=1}^m$ is from re-ordering $\{\alpha_i\}_{i=1}^m$ in the decreasing order, i.e.,

$$\alpha_1^\downarrow \geq \alpha_2^\downarrow \geq \dots \geq \alpha_m^\downarrow$$

(similarly for $\{\beta_i^\downarrow\}_{i=1}^m$). We also use notation α_i^\uparrow obtained from re-ordering $\{\alpha_i\}_{i=1}^m$ as well but in the increasing order.

Lemma 4.1. *Let $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_m$. If $\{\beta_i\}_{i=1}^m$ majorizes $\{\alpha_i\}_{i=1}^m$, then*

$$\sum_{i=1}^m \gamma_i \beta_i^\uparrow \leq \sum_{i=1}^m \gamma_i \alpha_i \leq \sum_{i=1}^m \gamma_i \beta_i^\downarrow. \tag{4.1}$$

Furthermore, if all γ_i are distinct, then the first inequality becomes an equality if and only if $\alpha_i = \beta_i^\uparrow$ for all i . Similarly, if all γ_i are distinct, then the second inequality becomes an equality if and only if $\alpha_i = \beta_i^\downarrow$ for all i .

The first part of the lemma is exactly the same as [17, Lemma 2.3], except that here it is not required that all $\gamma_i \geq 0$. The second part on the inequalities becoming equalities was not explicitly stated there, but it follows from the proof there straightforwardly. This lemma likely appeared elsewhere⁶ but an explicit reference is hard to find. As a corollary, we have

$$\sum_{i=1}^m \gamma_i \beta_i^\uparrow \leq \sum_{i=1}^m \gamma_i \beta_i \leq \sum_{i=1}^m \gamma_i \beta_i^\downarrow$$

because clearly $\{\beta_i\}_{i=1}^m$ majorizes $\{\beta_i\}_{i=1}^m$ itself.

Proof of Lemma 4.1. Without loss of generality, we may assume $\gamma_m > 0$; otherwise we can always pick a scalar ξ such that $\gamma_m + \xi \geq 0$, and let

$$\tilde{\gamma}_i := \gamma_i + \xi > 0 \quad \text{for } 1 \leq i \leq m.$$

By assumption, we have $\sum_{i=1}^m \alpha_i = \sum_{i=1}^m \beta_i = \sum_{i=1}^m \beta_i^\uparrow = \sum_{i=1}^m \beta_i^\downarrow =: \eta$, and thus

$$\sum_{i=1}^m \gamma_i \beta_i^\uparrow = -\xi \eta + \sum_{i=1}^m \tilde{\gamma}_i \beta_i^\uparrow, \quad \sum_{i=1}^m \gamma_i \alpha_i = -\xi \eta + \sum_{i=1}^m \tilde{\gamma}_i \alpha_i, \quad \sum_{i=1}^m \gamma_i \beta_i^\downarrow = -\xi \eta + \sum_{i=1}^m \tilde{\gamma}_i \beta_i^\downarrow.$$

It suffices to prove the lemma for $\tilde{\gamma}_1 \geq \tilde{\gamma}_2 \geq \dots \geq \tilde{\gamma}_m > 0$, instead.

⁶ After the paper was accepted, we found that Lemma 4.1 in fact had appeared in [13, Lemma 4.3.51 on p. 255] and Lemmas 4.2 and 4.3 had appeared in [13, Theorem 4.3.45 on p. 249].

The argument below up to (4.3) appears in the proof of [17, Lemma 2.3]. It is repeated here for the purpose of arguing when the equality signs in (4.1) are attained. Suppose $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_m > 0$ and set

$$p_j = \sum_{i=1}^j \beta_i^\uparrow, \quad s_j = \sum_{i=1}^j \alpha_i, \quad t_j = \sum_{i=1}^j \beta_i^\downarrow, \quad p_0 = s_0 = t_0 = 0.$$

Since $\{\beta_i\}_{i=1}^m$ majorizes $\{\alpha_i\}_{i=1}^m$, we have

$$p_j \leq s_j \leq t_j, \quad p_m = s_m = t_m$$

and thus

$$\begin{aligned} \sum_{i=1}^m \gamma_i \alpha_i &= \sum_{i=1}^m (s_i - s_{i-1}) \gamma_i \\ &= \sum_{i=1}^m s_i \gamma_i - \sum_{i=2}^m s_{i-1} \gamma_i \\ &= s_m \gamma_m + \sum_{i=1}^{m-1} s_i (\gamma_i - \gamma_{i+1}) \\ &\leq t_m \gamma_m + \sum_{i=1}^{m-1} t_i (\gamma_i - \gamma_{i+1}) \\ &= \sum_{i=1}^m \gamma_i \beta_i^\downarrow, \end{aligned} \tag{4.2}$$

$$\begin{aligned} \sum_{i=1}^m \gamma_i \alpha_i &= s_m \gamma_m + \sum_{i=1}^{m-1} s_i (\gamma_i - \gamma_{i+1}) \\ &\geq p_m \gamma_m + \sum_{i=1}^{m-1} p_i (\gamma_i - \gamma_{i+1}) \\ &= \sum_{i=1}^m \gamma_i \beta_i^\uparrow, \end{aligned} \tag{4.3}$$

as required. To figure out when any of the inequalities above is an equality, we look at (4.2), for an example. We notice that there is only one inequality sign during the derivation in (4.2). In order for the inequality to become an equality, assuming all γ_i are distinct, we will have to have $s_i = t_i$ for all i and consequently, $\alpha_i = \beta_i^\downarrow$ for all i . \square

The next two lemmas relate the diagonal entries of a Hermitian matrix with its eigenvalues.

Lemma 4.2 ([43, Exercise II.1.12, p. 35]). *The multiset of the diagonal entries of a Hermitian matrix is majorized by the multiset of its eigenvalues.*

Lemma 4.3. *For a Hermitian matrix, if the multiset of its diagonal entries is the same as the multiset of its eigenvalues, then it is diagonal.*

Proof. This lemma is probably known, but we could not find a reference to it. For completeness, we provide a quick proof. Let $A = [a_{ij}] \in \mathbb{C}^{n \times n}$ be such a Hermitian matrix with eigenvalues $\{\lambda_i\}_{i=1}^n$. By the assumption,

$$\|A\|_F^2 = \sum_{i,j=1}^n |a_{ij}|^2 = \sum_{i=1}^n |\lambda_i|^2 = \sum_{i=1}^n |a_{ii}|^2,$$

where $\|A\|_F$ denotes the Frobenius norm of A . Hence $|a_{ij}|^2 = 0$ for all $i \neq j$, as expected. \square

Now we are ready to prove Theorem 2.1.

Proof of Theorem 2.1. Recall the eigen-decomposition (2.1) with (1.2) for $A - \lambda B$ and the eigen-decomposition (2.2) for D . Consider first that D is nonsingular, i.e., all $\omega_i \neq 0$.

Introducing

$$Y = U^{-1}XQ \quad \Rightarrow \quad X = UYQ^H, \tag{4.4}$$

we find that the left-hand side of (2.3) can be transformed to

$$\min_{X^H B X = I_k} \text{tr}(DX^H A X) = \min_{Y^H Y = I_k} \text{tr}(\Omega Y^H \Lambda Y),$$

and any minimizer of one yields a minimizer of the other according to (4.4).

For any given $Y \in \mathbb{C}^{n \times k}$ with $Y^H Y = I_k$, denote the eigenvalues of $Y^H \Lambda Y$ by

$$\mu_1 \leq \mu_2 \leq \dots \leq \mu_k,$$

where we suppress the dependency of μ_i on Y for clarity. Cauchy’s interlacing inequalities say that

$$\lambda_{i+n-k} \geq \mu_i \geq \lambda_i \text{ for all } 1 \leq i \leq k. \tag{4.5}$$

Denote the diagonal entries of $Y^H \Lambda Y$ by $(Y^H \Lambda Y)_{(i,i)}$ for $i = 1, 2, \dots, k$, and let α_i be the reordering of $(Y^H \Lambda Y)_{(i,i)}$ in the increasing order, i.e.,

$$\alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_k.$$

Evidently, $\{(Y^H \Lambda Y)_{(i,i)}\}_{i=1}^k$ is majorized by $\{\alpha_i\}_{i=1}^k$ because they are the same up to a permutation. By Lemma 4.2, $\{\alpha_i\}_{i=1}^k$ is majorized by $\{\mu_i\}_{i=1}^k$. We have

$$\begin{aligned}
 \text{tr}(\Omega Y^H \Lambda Y) &= \sum_{i=1}^k \omega_i (Y^H \Lambda Y)_{(i,i)} \quad (\text{use Lemma 4.1}) \\
 &\geq \sum_{i=1}^k \omega_i \alpha_i \quad (\text{use Lemma 4.1}) \\
 &\geq \sum_{i=1}^k \omega_i \mu_i \\
 &= \sum_{i=1}^{\ell} \omega_i \mu_i + \sum_{i=\ell+1}^k \omega_i \mu_i \quad (\text{use (4.5)}) \\
 &\geq \sum_{i=1}^{\ell} \omega_i \lambda_i + \sum_{i=\ell+1}^k \omega_i \lambda_{i+n-k}. \tag{4.6}
 \end{aligned}$$

Since Y is arbitrary, we have

$$\min_{Y^H Y = I_k} \text{tr}(\Omega Y^H \Lambda Y) \geq \sum_{i=1}^{\ell} \omega_i \lambda_i + \sum_{i=\ell+1}^k \omega_i \lambda_{i+n-k}. \tag{4.7}$$

Taking $Y = [e_1, e_2, \dots, e_{\ell}, e_{n-k+\ell+1}, e_{n-k+\ell+2}, \dots, e_n]$ where e_i is the i th column of I_n , we see that $\text{tr}(\Omega Y^H \Lambda Y)$ is equal to the right-hand side of (4.7). Therefore we have (2.3).

Suppose now all ω_i are distinct and Y_{opt} is a minimizer. We must have

$$\begin{aligned}
 \text{tr}(\Omega Y_{\text{opt}}^H \Lambda Y_{\text{opt}}) &= \sum_{i=1}^k \omega_i (Y_{\text{opt}}^H \Lambda Y_{\text{opt}})_{(i,i)} \\
 &= \sum_{i=1}^k \omega_i \alpha_i \tag{4.8}
 \end{aligned}$$

$$= \sum_{i=1}^k \omega_i \mu_i \tag{4.9}$$

$$= \sum_{i=1}^{\ell} \omega_i \lambda_i + \sum_{i=\ell+1}^k \omega_i \lambda_{i+n-k}, \tag{4.10}$$

where $\alpha_1 \leq \alpha_2 \leq \dots \leq \alpha_k$ are the reordering of $(Y_{\text{opt}}^H \Lambda Y_{\text{opt}})_{(i,i)}$, and $\mu_1 \leq \mu_2 \leq \dots \leq \mu_k$ are the eigenvalues of $Y_{\text{opt}}^H \Lambda Y_{\text{opt}}$. For the equalities in (4.8) – (4.10) to hold, we must have for all i

$$\begin{aligned}
 (Y_{\text{opt}}^H \Lambda Y_{\text{opt}})_{(i,i)} &= \alpha_i = \mu_i = \lambda_i \quad \text{for } 1 \leq i \leq \ell, \\
 (Y_{\text{opt}}^H \Lambda Y_{\text{opt}})_{(i,i)} &= \alpha_i = \mu_i = \lambda_{n-k+i} \quad \text{for } \ell + 1 \leq i \leq k,
 \end{aligned}$$

and $Y_{\text{opt}}^H \Lambda Y_{\text{opt}} = \text{diag}(\alpha_1, \alpha_2, \dots, \alpha_k) = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_\ell, \lambda_{n-k+\ell+1}, \dots, \lambda_n)$. Now use the relation (4.4) to conclude the proof for the case when D is nonsingular.

Return to the case when D is singular, i.e., some of its eigenvalues $\omega_i = 0$. Let \widehat{Q} be as the one defined in item (b) and \widehat{Q}_\perp be the columns of Q not in \widehat{Q} . The eigen-decomposition (2.2) of D can be rewritten as

$$D = [\widehat{Q}, \widehat{Q}_\perp] \begin{bmatrix} \widehat{\Omega} & 0 \\ 0 & 0 \end{bmatrix} [\widehat{Q}, \widehat{Q}_\perp]^H,$$

where $\widehat{\Omega} = \text{diag}(\omega_1, \dots, \omega_{\ell_+}, \omega_{k-\ell_-+1}, \dots, \omega_k)$ composed of all nonzero eigenvalues of D . It can be verified that

$$\text{tr}(DX^HAX) = \text{tr}(\widehat{\Omega}\widehat{Y}^H A\widehat{Y}), \tag{4.11}$$

where $\widehat{Y} = X\widehat{Q}$. If $X^H B X = I_k$, then $\widehat{Y}^H B \widehat{Y} = \widehat{Q}^H X^H B X \widehat{Q} = \widehat{Q}^H \widehat{Q} = I_{\ell_+ + \ell_-}$. On the other hand, given $\widehat{Y} \in \mathbb{C}^{n \times (\ell_+ + \ell_-)}$ such that $\widehat{Y}^H B \widehat{Y} = I_{\ell_+ + \ell_-}$, we can expand it to $Y = [\widehat{Y}, \widehat{Y}_c] \in \mathbb{C}^{n \times k}$ such that $Y^H B Y = I_k$ and then let $X = Y[\widehat{Q}, \widehat{Q}_\perp]^H$ for which it can be seen that (4.11) holds. This proves

$$\min_{X^H B X = I_k} \text{tr}(DX^HAX) = \min_{\widehat{Y}^H B \widehat{Y} = I_{\ell_+ + \ell_-}} \text{tr}(\widehat{\Omega}\widehat{Y}^H A\widehat{Y}), \tag{4.12}$$

and a minimizer for one leads to a minimizer for the other. The right-hand side of (4.12) is a minimization problem belonging to the case of nonsingular D that we just dealt with. \square

5. Proof of Theorems 3.1 and 3.2

The next two lemmas will be needed in our later proofs.

Lemma 5.1 ([44, Corollary 5.12]). *Let $J_m = \text{diag}(I_{n_+}, -I_{n_-})$ and $m = n_+ + n_-$. A set of vectors u_1, \dots, u_k in \mathbb{C}^m satisfying $u_i^H J_n u_j = \pm \delta_{ij}$ for $i, j = 1, \dots, k$ can be complemented to a basis $\{u_1, \dots, u_m\}$ of \mathbb{C}^m satisfying $u_i^H J_m u_j = \pm \delta_{ij}$ for $i, j = 1, \dots, m$, where δ_{ij} is the Kronecker delta which is 1 for $i = j$ and 0 otherwise, and the numbers of 1 and -1 among $u_i^H J_m u_i$ for $1 \leq i \leq n$ are n_+ and n_- , respectively.*

Lemma 5.2 ([44, Example 6.3]). *Let $J_m = \text{diag}(I_{n_+}, -I_{n_-})$ and $m = n_+ + n_-$. A matrix $X \in \mathbb{C}^{m \times m}$ satisfies $X^H J_m X = J_m$ if and only if it is of the form*

$$X = \begin{bmatrix} (I_{n_+} + WW^H)^{1/2} & W \\ W^H & (I_{n_-} + W^H W)^{1/2} \end{bmatrix} \begin{bmatrix} V_+ \\ V_- \end{bmatrix}, \tag{5.1}$$

where $V_+ \in \mathbb{C}^{n_+ \times n_+}$ and $V_- \in \mathbb{C}^{n_- \times n_-}$ are unitary, and $W \in \mathbb{C}^{n_+ \times n_-}$.

Lemma 5.1 can also be found in many classical monographs, e.g., [45,46], and Lemma 5.2 can be found in [47,22], where (5.1) is called a (hyperbolic) polar decomposition of X .

Proof of Theorem 3.1. We may assume, without loss of generality, that $A \succeq 0$; Otherwise, noticing

$$\text{tr}(DX^HAX) = \text{tr}(DX^H(A - \lambda_0B)X) + \lambda_0 \text{tr}(D), \tag{5.2}$$

we may consider $\text{tr}(DX^H(A - \lambda_0B)X)$, instead.

In what follows, we will suppose $A \succeq 0$ and also $\lambda_0 = 0$.

By [23, Lemma 3.8], $A - \lambda B$ admits the following eigen-decomposition: there exists a nonsingular matrix $U \in \mathbb{C}^{n \times n}$ such that

$$U^H A U = \begin{matrix} & \begin{matrix} n_+ - m_0 & n_- - m_0 & 2m_0 & n_0 \end{matrix} \\ \begin{matrix} n_+ - m_0 \\ n_- - m_0 \\ 2m_0 \\ n_0 \end{matrix} & \begin{bmatrix} \Lambda_+ & & & \\ & -\Lambda_- & & \\ & & \Lambda_b & \\ & & & \Lambda_\infty \end{bmatrix} \end{matrix} =: \begin{matrix} r & n_0 \\ n_0 & \end{matrix} \begin{bmatrix} \Lambda_r & \\ & \Lambda_\infty \end{bmatrix} =: \Lambda \succeq 0, \tag{5.3a}$$

$$U^H B U = \begin{matrix} & \begin{matrix} n_+ - m_0 & n_- - m_0 & 2m_0 & n_0 \end{matrix} \\ \begin{matrix} n_+ - m_0 \\ n_- - m_0 \\ 2m_0 \\ n_0 \end{matrix} & \begin{bmatrix} I_{n_+ - m_0} & & & \\ & -I_{n_- - m_0} & & \\ & & J_b & \\ & & & 0 \end{bmatrix} \end{matrix} =: \begin{matrix} r & n_0 \\ n_0 & \end{matrix} \begin{bmatrix} J_r & \\ & J_\infty \end{bmatrix} =: J_n, \tag{5.3b}$$

where $0 \leq m_0 \leq \min\{n_+, n_-\}$, and⁷

$$U = \begin{bmatrix} U_+ & U_- & U_b & U_\infty \end{bmatrix} =: \begin{bmatrix} U_r & U_\infty \end{bmatrix}, \tag{5.3c}$$

$$\Lambda_+ = \text{diag}(\lambda_{m_0+1}^+, \dots, \lambda_{n_+}^+), \quad \Lambda_- = \text{diag}(\lambda_{n_-}^-, \dots, \lambda_{m_0+1}^-), \tag{5.3d}$$

$$\Lambda_0 = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}, \quad F_2 = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}, \tag{5.3e}$$

$$J_b = \text{diag}(\underbrace{F_2, \dots, F_2}_{m_0}), \quad \Lambda_b = \text{diag}(\underbrace{\Lambda_0, \dots, \Lambda_0}_{m_0}), \quad \Lambda_\infty \succeq 0. \tag{5.3f}$$

Both Λ in (5.3a) and J_n in (5.3b) are diagonal if $m_0 = 0$, i.e., in the absence of blocks Λ_b , J_b , and U_b . For the case, we say that $A - \lambda B$ is diagonalizable. It can be seen from (5.3) that the finite eigenvalues of $A - \lambda B$ are given by

⁷ Recall the simplification due to (5.2). In general, Λ_0 in (5.3e) takes the form $\begin{bmatrix} 0 & \lambda_0 \\ \lambda_0 & 1 \end{bmatrix}$, and thus $\lambda_{m_0}^- = \dots = \lambda_1^- = \lambda_0 = \lambda_1^+ = \dots = \lambda_{m_0}^+$.

$$\lambda_{n_-}^- \leq \dots \leq \lambda_{m_0+1}^- \leq \underbrace{0 = \dots = 0}_{m_0} = \underbrace{0 = \dots = 0}_{m_0} \leq \lambda_{m_0+1}^+ \leq \dots \leq \lambda_{n_+}^+,$$

which, compared to (3.4), implies $\lambda_{m_0}^- = \dots = \lambda_1^- = 0 = \lambda_1^+ = \dots = \lambda_{m_0}^+$, and they come from $\Lambda_b - \lambda J_b$.

Letting $Y = U^{-1}XQ$, we can transform (1.8) for the case $X^H B X = I_k$ into

$$\inf_{X^H B X = I_k} \text{tr}(DX^H A X) = \inf_{Y^H J_n Y = I_k} \text{tr}(\Omega Y^H \Lambda Y), \tag{5.4}$$

where $k \leq n_+$.

First we deal with the case when **the matrix B is singular**, i.e., $n_0 > 0$ in (5.3).

Partition $Y = \begin{matrix} r \\ n_0 \end{matrix} \begin{matrix} k \\ Y_r \\ Y_\infty \end{matrix}$, and then

$$\begin{aligned} \inf_{Y^H J_n Y = I_k} \text{tr}(\Omega Y^H \Lambda Y) &= \inf_{Y_r^H J_r Y_r = I_k} [\text{tr}(\Omega Y_r^H \Lambda_r Y_r) + \text{tr}(\Omega Y_\infty^H \Lambda_\infty Y_\infty)] \\ &= \inf_{Y_r^H J_r Y_r = I_k} \text{tr}(\Omega Y_r^H \Lambda_r Y_r) + \inf_{Y_\infty} \text{tr}(\Omega Y_\infty^H \Lambda_\infty Y_\infty). \end{aligned} \tag{5.5}$$

We will examine the two terms in (5.5) separately. Constraint $Y^H J_n Y = I_k$ yields $Y_r^H J_r Y_r = I_k$, leaving $Y_\infty \in \mathbb{C}^{n_0 \times k}$ arbitrary. Restricting Y_∞ to a rank-1 matrix xy^H , we find

$$\begin{aligned} \inf_{Y_\infty} \text{tr}(\Omega Y_\infty^H \Lambda_\infty Y_\infty) &\leq \inf_{\text{rank}(Y_\infty) \leq 1} \text{tr}(\Omega Y_\infty^H \Lambda_\infty Y_\infty) \\ &= \inf_{y,x} \text{tr}(\Omega x y^H \Lambda_\infty y x^H) \\ &= \inf_{y,x} (x^H \Omega x)(y^H \Lambda_\infty y). \end{aligned}$$

There are three cases.

1. $\Omega \not\geq 0$ and $\Lambda_\infty \neq 0$: we have $\inf_{y,x} (x^H \Omega x)(y^H \Lambda_\infty y) = -\infty$, which leads to that the second infimum in (5.5) is $-\infty$.
2. $\Omega \not\geq 0$ and $\Lambda_\infty = 0$: we have $Y_\infty^H \Lambda_\infty Y_\infty = 0$, which leads to that the second infimum in (5.5) is 0. But our later proof for nonsingular B shows that for the case the first infimum in (5.5) is $-\infty$.
3. $\Omega \succeq 0$: we have $Y_\infty^H \Lambda_\infty Y_\infty \succeq 0$, and $\text{tr}(\Omega Y_\infty^H \Lambda_\infty Y_\infty) \geq 0$ and $\text{tr}(\Omega Y_\infty^H \Lambda_\infty Y_\infty) = 0$ for $Y_\infty = 0$, which leads to that the second infimum in (5.5) is 0.

The first infimum in (5.5): $\inf_{Y_r^H J_r Y_r = I_k} \text{tr}(\Omega Y_r^H \Lambda_r Y_r)$, falls into the case when **the matrix B is nonsingular**, which we are about to investigate.

Suppose now that B is nonsingular, i.e., $n_0 = 0$ in (5.3).

Consider first that $m_0 = 0$, namely **the pencil $A - \lambda B$ is also diagonalizable**. Then $J_n = \text{diag}(I_{n_+}, -I_{n_-})$. Since $Y^H J_n Y = I_k$, by Lemma 5.1 we can complement Y to $\tilde{Y} = [Y \quad Y_c] \in \mathbb{C}^{n \times n}$ such that $\tilde{Y}^H J_n \tilde{Y} = J_n$. By Lemma 5.2, \tilde{Y} has a hyperbolic polar decomposition

$$\tilde{Y} = \begin{bmatrix} (I_{n_+} + \tilde{\Sigma} \tilde{\Sigma}^H)^{1/2} & \tilde{\Sigma} \\ \tilde{\Sigma}^H & (I_{n_-} + \tilde{\Sigma}^H \tilde{\Sigma})^{1/2} \end{bmatrix} \begin{bmatrix} \tilde{V}_+ \\ \tilde{V}_- \end{bmatrix}, \tag{5.6}$$

where $\tilde{V}_+ \in \mathbb{C}^{n_+ \times n_+}$ and $\tilde{V}_- \in \mathbb{C}^{n_- \times n_-}$ are unitary, and $\tilde{\Sigma} \in \mathbb{C}^{n_+ \times n_-}$. Let $\tilde{\Sigma} = W_+ \Sigma W_-^H$ be the singular value decomposition of $\tilde{\Sigma}$, where

$$\Sigma = \begin{bmatrix} \Sigma_0 \\ 0 \end{bmatrix} \text{ if } n_+ \geq n_-, \text{ or } \Sigma = [0 \quad \Sigma_0] \text{ if } n_+ < n_-.$$

Hence plug $\tilde{\Sigma} = W_+ \Sigma W_-^H$ into (5.6) to turn $\tilde{Y} = W S V^H$, where

$$W = \begin{matrix} n_+ & n_- \\ n_- \end{matrix} \begin{bmatrix} W_+ & \\ & W_- \end{bmatrix}, \quad V = \begin{matrix} n_+ & n_- \\ n_- \end{matrix} \begin{bmatrix} V_+ & \\ & V_- \end{bmatrix} := \begin{matrix} n_+ & n_- \\ n_- \end{matrix} \begin{bmatrix} \tilde{V}_+^H W_+ & \\ & \tilde{V}_-^H W_- \end{bmatrix}, \tag{5.7a}$$

and

$$\begin{aligned} S &= \begin{bmatrix} (I_{n_+} + \Sigma \Sigma^H)^{1/2} & \Sigma \\ \Sigma^H & (I_{n_-} + \Sigma^H \Sigma)^{1/2} \end{bmatrix} \\ &= \begin{bmatrix} (I + \Sigma_0^2)^{1/2} & 0 & \Sigma_0 \\ 0 & I_{|n_+ - n_-|} & 0 \\ \Sigma_0 & 0 & (I + \Sigma_0^2)^{1/2} \end{bmatrix}. \end{aligned} \tag{5.7b}$$

Noticing $Y = \tilde{Y} I_{n;k}$ where $I_{n;k} = \begin{bmatrix} I_k \\ 0 \end{bmatrix} \in \mathbb{C}^{n \times k}$, we have from (5.4)

$$\begin{aligned} \inf_{Y^H J_n Y = I_k} \text{tr}(\Omega Y^H \Lambda Y) &= \inf_{\tilde{Y}^H J_n \tilde{Y} = J_n} \text{tr}(\Omega I_{n;k}^H \tilde{Y}^H \Lambda \tilde{Y} I_{n;k}) \\ &= \inf_{\substack{\Sigma_0 \geq 0 \text{ diagonal} \\ V_+, V_-, \tilde{W}_+, \tilde{W}_- \text{ unitary}}} \text{tr}(I_{n;k} \Omega I_{n;k}^H V S W^H \Lambda W S V^H) \\ &= \inf_{\substack{\Sigma_0 \geq 0 \text{ diagonal} \\ V_+, V_-, \tilde{W}_+, \tilde{W}_- \text{ unitary}}} \text{tr}(\tilde{\Omega}_V S \Lambda W S), \end{aligned} \tag{5.8}$$

where $\tilde{\Omega}_V = V^H I_{n;k} \Omega I_{n;k}^H V$ and $\Lambda_W = W^H \Lambda W$. Use (5.7a) to see

$$\tilde{\Omega}_V = \begin{matrix} n_+ & n_- \\ n_- \end{matrix} \begin{bmatrix} \tilde{\Omega}_{+,V} & \\ & 0 \end{bmatrix}, \quad \Lambda_W = \begin{matrix} n_+ & n_- \\ n_- \end{matrix} \begin{bmatrix} \Lambda_{+,W} & \\ & -\Lambda_{-,W} \end{bmatrix},$$

where $\tilde{\Omega}_+ = \begin{matrix} k & n_+ - k \\ \Omega & \\ & 0 \end{matrix}$, $\tilde{\Omega}_{+,V} = V_+^H \tilde{\Omega}_+ V_+$, $\Lambda_{+,W} = W_+^H \Lambda_+ W_+$, and $\Lambda_{-,W} = W_-^H \Lambda_- W_-$. As a result, $\text{tr}(\tilde{\Omega}_V S \Lambda_W S)$ can be given by

$$\begin{aligned} & \text{tr} \left(\begin{bmatrix} \tilde{\Omega}_{+,V} & \\ & 0 \end{bmatrix} \begin{bmatrix} (I_{n_+} + \Sigma \Sigma^H)^{1/2} & \Sigma \\ \Sigma^H & (I_{n_-} + \Sigma^H \Sigma)^{1/2} \end{bmatrix} \right. \\ & \quad \left. \times \begin{bmatrix} \Lambda_{+,W} & \\ & \Lambda_{-,W} \end{bmatrix} \begin{bmatrix} (I_{n_+} + \Sigma \Sigma^H)^{1/2} & \Sigma \\ \Sigma^H & (I_{n_-} + \Sigma^H \Sigma)^{1/2} \end{bmatrix} \right) \\ & = \text{tr} \left(\begin{bmatrix} \tilde{\Omega}_{+,V} (I_{n_+} + \Sigma \Sigma^H)^{1/2} & \tilde{\Omega}_{+,V} \Sigma \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \Lambda_{+,W} (I_{n_+} + \Sigma \Sigma^H)^{1/2} & \Lambda_{+,W} \Sigma \\ \Lambda_{-,W} \Sigma^H & \Lambda_{-,W} (I_{n_-} + \Sigma^H \Sigma)^{1/2} \end{bmatrix} \right) \\ & = \text{tr}(\tilde{\Omega}_{+,V} [(I_{n_+} + \Sigma \Sigma^H)^{1/2} \Lambda_{+,W} (I_{n_+} + \Sigma \Sigma^H)^{1/2} - \Sigma \Lambda_{-,W} \Sigma^H]). \end{aligned}$$

The last infimum in (5.8) becomes

$$\begin{aligned} & \inf_{\substack{\Sigma_0 \geq 0 \text{ diagonal} \\ V_+, V_-, \tilde{W}_+, W_- \text{ unitary}}} \text{tr}(\tilde{\Omega}_V S \Lambda_W S) \\ & = \inf_{\substack{\Sigma_0 \geq 0 \text{ diagonal} \\ V_+, \tilde{W}_+, W_- \text{ unitary}}} \text{tr}(\tilde{\Omega}_{+,V} [(I_{n_+} + \Sigma \Sigma^H)^{1/2} \Lambda_{+,W} (I_{n_+} + \Sigma \Sigma^H)^{1/2} - \Sigma \Lambda_{-,W} \Sigma^H]). \end{aligned} \tag{5.9}$$

This infimum is $-\infty$ if $\Omega \not\leq 0$. In fact, suppose $\tilde{\Omega}_+ x_+ = \omega_k x_+$ where $\omega_k < 0$, and x_+ is a unit eigenvector. Construct $\hat{V}_+ = [x_+ \quad V_{+,c}] \in \mathbb{C}^{n_+ \times n_+}$ that is unitary. Thus, upon restrictions $V_+ = \hat{V}_+$, $W_+ = I$, $W_- = I$, $\Sigma_0 = \text{diag}(\sigma, 0, \dots, 0)$, we have by (5.9)

$$\begin{aligned} & \inf_{\substack{\Sigma_0 \geq 0 \text{ diagonal} \\ V_+, V_-, \tilde{W}_+, W_- \text{ unitary}}} \text{tr}(\tilde{\Omega}_V S \Lambda_W S) \leq \inf_{\sigma} \text{tr}(\hat{V}_+^H \tilde{\Omega}_+ \hat{V}_+ [\Lambda_+ + \sigma^2 (\lambda_{n_+}^+ - \lambda_{n_-}^-) e_1 e_1^H]) \\ & = \inf_{\sigma} \sigma^2 (\lambda_{n_+}^+ - \lambda_{n_-}^-) \text{tr}(\hat{V}_+^H \tilde{\Omega}_+ \hat{V}_+ e_1 e_1^H) + \text{tr}(\hat{V}_+^H \tilde{\Omega}_+ \hat{V}_+ \Lambda_+) \\ & = \inf_{\sigma} \sigma^2 (\lambda_{n_+}^+ - \lambda_{n_-}^-) (e_1^H \hat{V}_+^H \tilde{\Omega}_+ \hat{V}_+ e_1) + \text{tr}(\hat{V}_+^H \tilde{\Omega}_+ \hat{V}_+ \Lambda_+) \\ & = \inf_{\sigma} \sigma^2 (\lambda_{n_+}^+ - \lambda_{n_-}^-) (x_+^H \tilde{\Omega}_+ x_+) + \text{tr}(\hat{V}_+^H \tilde{\Omega}_+ \hat{V}_+ \Lambda_+) \\ & = \inf_{\sigma} \sigma^2 (\lambda_{n_+}^+ - \lambda_{n_-}^-) \omega_k + \text{tr}(\hat{V}_+^H \tilde{\Omega}_+ \hat{V}_+ \Lambda_+) \\ & = -\infty, \end{aligned}$$

as long as $\Lambda_+ \neq 0$ or $\Lambda_- \neq 0$, which is equivalent to $A \neq 0$, where e_1 is the first column of the identity matrix.

So far, we have shown that if $\Omega \not\leq 0, A \neq 0$, the infimum is $-\infty$, for any positive semi-definite pencil $A - \lambda B$ with B genuinely indefinite, except the case $A - \lambda B$ is not diagonalizable, to which we will return.

In what follows, suppose that $\Omega \succeq 0$. Then

$$\begin{aligned} & \inf_{\substack{\Sigma_0 \succeq 0 \text{ diagonal} \\ V_+, V_-, W_+, W_- \text{ unitary}}} \text{tr} \left(\tilde{\Omega}_{+,V} [(I + \Sigma \Sigma^H)^{1/2} \Lambda_{+,W} (I + \Sigma \Sigma^H)^{1/2} - \Sigma \Lambda_{-,W} \Sigma^H] \right) \\ &= \inf_{V_\pm, W_\pm \text{ unitary}} \inf_{\Sigma_0 \succeq 0 \text{ diagonal}} \text{tr} \left(\tilde{\Omega}_{+,V} [(I + \Sigma \Sigma^H)^{1/2} \Lambda_{+,W} (I + \Sigma \Sigma^H)^{1/2} - \Sigma \Lambda_{-,W} \Sigma^H] \right). \end{aligned} \tag{5.10}$$

Note that $\tilde{\Omega}_{+,V}, \Lambda_{+,W}, -\Lambda_{-,W}$ are positive semi-definite. In the following, we will show that the “infimum” in (5.10) is $\sum_{i=1}^k \omega_i \lambda_i^+$ and is attained at $\Sigma = 0$.

Firstly, since $\tilde{\Omega}_{+,V} \succeq 0$ and $\Sigma(-\Lambda_{-,W})\Sigma^H \succeq 0$, we have by Theorem 2.1

$$\text{tr}(-\tilde{\Omega}_{+,V} \Sigma \Lambda_{-,W} \Sigma^H) = \text{tr}(\Sigma(-\Lambda_{-,W})\Sigma^H) \geq 0. \tag{5.11}$$

Secondly, we claim that

$$\text{tr} \left(\tilde{\Omega}_{+,V} (I + \Sigma \Sigma^H)^{1/2} \Lambda_{+,W} (I + \Sigma \Sigma^H)^{1/2} \right) \geq \sum_{i=1}^k \omega_i \lambda_i^+. \tag{5.12}$$

Denote the eigenvalues of $(I + \Sigma \Sigma^H)^{1/2} \Lambda_{+,W} (I + \Sigma \Sigma^H)^{1/2}$ by $\delta_1 \leq \delta_2 \leq \dots \leq \delta_{n_+}$. By Ostrowski’s theorem [13, p. 283], we know

$$\lambda_i^+ \leq \delta_i \leq (1 + \|\Sigma\|_2^2) \lambda_i^+ \quad \text{for } 1 \leq i \leq n_+, \tag{5.13}$$

where $\|\Sigma\|_2$ is the spectral norm of Σ . Let $I_{n_+,k} = \begin{bmatrix} I_k \\ 0 \end{bmatrix} \in \mathbb{C}^{n_+ \times k}$. We have

$$\begin{aligned} & \text{tr} \left(\tilde{\Omega}_{+,V} (I + \Sigma \Sigma^H)^{1/2} \Lambda_{+,W} (I + \Sigma \Sigma^H)^{1/2} \right) \\ &= \text{tr} \left(V_+^H I_{n_+,k} \Omega I_{n_+,k}^H V_+ (I + \Sigma \Sigma^H)^{1/2} \Lambda_{+,W} (I + \Sigma \Sigma^H)^{1/2} \right) \\ &= \text{tr} \left(\Omega (V_+^H I_{n_+,k})^H (I + \Sigma \Sigma^H)^{1/2} \Lambda_{+,W} (I + \Sigma \Sigma^H)^{1/2} (V_+^H I_{n_+,k}) \right) \\ &\geq \sum_{i=1}^k \omega_i \delta_i \quad (\text{by Theorem 2.1}) \\ &\geq \sum_{i=1}^k \omega_i \lambda_i^+, \end{aligned} \tag{5.14}$$

where we have used (5.13) in the last step. This is (5.12). It is not hard to see that the equality in (5.14) is attained at $\Sigma = 0$ and appropriately chosen V_+ and W_+ . Combining (5.9), (5.10), (5.11), and (5.12) completes the proof of the theorem for the case when $A - \lambda B$ is diagonalizable.

Consider now that $m_0 > 0$, namely **the pencil $A - \lambda B$ is not diagonalizable**. We perturb $A - \lambda B$ to $(A + \varepsilon E) - \lambda B$ with $\varepsilon > 0$, where

$$E = U^{-H} \text{diag}(0, 0, E_b, 0)U^{-1}, \quad E_b = \text{diag}(\underbrace{E_0, \dots, E_0}_{m_0}), \quad E_0 = \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix}.$$

Clearly $(A + \varepsilon E) \succeq 0$ and the pencil is diagonalizable. Letting $\varepsilon \rightarrow 0^+$ leads to the desired result, based on the case for diagonalizable $A - \lambda B$. \square

Applying Theorem 3.1 to the matrix pencil $A - (-\lambda)(-B)$, we immediately conclude the following corollary.

Corollary 5.3. *Suppose the conditions of Theorem 3.1, except now $k \leq n_-$. Then*

$$\inf_{X^H B X = -I_k} \text{tr}(D X^H A X) > -\infty$$

if and only if $D \succeq 0$, in which case

$$\inf_{X^H B X = -I_k} \text{tr}(D X^H A X) = -\sum_{i=1}^k \omega_i \lambda_i^-. \tag{5.15}$$

The infimum can be attained, when $A - \lambda B$ is diagonalizable, by X such that the columns of XQ are the eigenvectors of $A - \lambda B$ associated with its eigenvalues λ_i^- for $1 \leq i \leq k$, respectively.

With the help of Theorem 3.1 and Corollary 5.3, we are ready to prove Theorem 3.2.

Proof of Theorem 3.2. Partition $X = [X_+ \ X_-]$ with $X_+ \in \mathbb{C}^{n \times k_+}$, and let $D_{\pm} = Q_{\pm} \Omega_{\pm} Q_{\pm}^H$ be the eigen-decompositions of D_+ and D_- , respectively. First consider the case that $A - \lambda B$ is diagonalizable. We have

$$\begin{aligned} \min_{X^H B X = J_k} \text{tr}(D X^H A X) &= \min_{\substack{X_{\pm}^H B X_{\pm} = \pm I_{k_{\pm}} \\ X_+^H B X_- = 0}} \text{tr}(D_+ X_+^H A X_+ + D_- X_-^H A X_-) \\ &\geq \min_{X_+^H B X_+ = I_{k_+}} \text{tr}(D_+ X_+^H A X_+) + \min_{X_-^H B X_- = -I_{k_-}} \text{tr}(D_- X_-^H A X_-) \\ &= \sum_{i=1}^{k_+} \omega_i^+ \lambda_i^+ - \sum_{i=1}^{k_-} \omega_i^- \lambda_i^-, \end{aligned}$$

of which the last equality holds by making the columns of $X_{\pm} Q_{\pm}$ be the eigenvectors of $A - \lambda B$ associated with its eigenvalues λ_i^{\pm} for $1 \leq i \leq k_{\pm}$, respectively. This proves (3.6).

For the case that $A - \lambda B$ is not diagonalizable, (3.6) also holds, by using the same technique at the end of the proof of Theorem 3.1 above. \square

6. Conclusion

Previously, the classical Ky Fan's trace minimization principle on $\min_X \text{tr}(X^H A X)$ subject to $X^H X = I_k$ for a Hermitian matrix A has been extended to about

$$\begin{aligned} & \min_{X^H B X = I_k} \text{tr}(X^H A X) \quad \text{for positive definite } B, \text{ or more generally} \\ & \inf_{X^H B X = J_k} \text{tr}(X^H A X) \quad \text{for genuinely indefinite } B, \end{aligned}$$

where J_k is diagonal with diagonal entries ± 1 . The extension for a positive definite B is rather straightforward, but quite complicated when B is genuinely indefinite [22,23]. In fact, the infimum can be $-\infty$ for the last case.

Our extensions in this paper are along the line of the Brockett cost function: $\text{tr}(D X^H A X)$ in X satisfying $X^H X = I_k$, when D is Hermitian and positive semi-definite. Specifically, we present elegant analytic solutions, in terms of eigenvalues and eigenvectors of matrix pencil $A - \lambda B$, to

$$\min_{X^H B X = I_k} \text{tr}(D X^H A X) \quad \text{for positive definite } B, \quad (6.1a)$$

$$\inf_{X^H B X = J_k} \text{tr}(D X^H A X) \quad \text{for genuinely indefinite } B, \quad (6.1b)$$

where D is no longer assumed to be positive semi-definite. Our analytic solutions are concise and our algebraic technique compares favorably to previously laborious effort for the case $B = I$ via the usual optimization technique [35].

Each of earlier trace minimization principles (1.1), (1.3), (1.4), and the two principles in the appendices have inspired efficient numerical methods for corresponding large scale eigenvalue problems. The potential implications of the new ones established in this paper for (6.1) remain to be seen.

Declaration of competing interest

There is none.

Acknowledgements

The authors are grateful to the reviewers for their helpful comments and suggestions.

Appendix A. Two applications of (1.4)

In this section, we outline two applications of the trace minimization principle (1.4) [22,23] to two important applied eigenvalue problems. The applications will lead to two

known results but the goal is to demonstrate how easy (1.4) can be put into good use. Recall:

$$\inf_{X^H B X = J_k} \text{tr}(X^H A X) = \sum_{i=1}^{k_+} \lambda_i^+ - \sum_{i=1}^{k_-} \lambda_i^-, \tag{1.4}$$

where $A - \lambda B$ is a positive semi-definite matrix pencil with the inertia of B being (n_+, n_0, n_-) , λ_i^\pm are finite eigenvalues of $A - \lambda B$ (cf. section 3) and are arranged in the order as

$$\lambda_{n_-}^- \leq \dots \leq \lambda_1^- \leq \lambda_1^+ \leq \dots \leq \lambda_{n_+}^+, \tag{1.5}$$

and $0 \leq k_\pm \leq n_\pm$, $k = k_+ + k_-$.

A.1. Linear response eigenvalue problem

The linear response eigenvalue problem from electronic structure calculations [24,25, 31,26] (and references therein) is the eigenvalue problem for

$$\begin{bmatrix} 0 & K \\ M & 0 \end{bmatrix} - \lambda I_{2n}, \tag{A.1}$$

where $K, M \in \mathbb{C}^{n \times n}$ are Hermitian positive definite. This eigenvalue problem is equivalent to the one for

$$A - \lambda B := \begin{bmatrix} M & 0 \\ 0 & K \end{bmatrix} - \lambda \begin{bmatrix} 0 & I_n \\ I_n & 0 \end{bmatrix} \tag{A.2}$$

in the sense that both have the same eigenvalues and eigenvectors. It can be seen that B has n eigenvalues $+1$ and n eigenvalues -1 , and all eigenvalues of $A - \lambda B$ are finite and can be divided into two groups: n of the positive-type and n of the negative type as in (1.5). In fact, it is proved [24,25] that (A.1) (hence (A.2) by equivalency) has $2n$ real eigenvalues $\pm \lambda_i$ with

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

and thus in the notation of (1.5): $\lambda_i^- = -\lambda_i$ and $\lambda_i^+ = \lambda_i$ for $1 \leq i \leq n$. Furthermore, there exists a nonsingular $V \in \mathbb{C}^{n \times n}$ such that⁸

$$K = V \Lambda^2 V^H, \quad M = U U^H, \tag{A.3}$$

⁸ In [24,25], it was primarily stated in terms of the real number field \mathbb{R} , but was commented that all results hold for the complex number field \mathbb{C} after minor modifications [24, section 6].

where $\Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$ and $U = V^{-H}$. Now for any $X, Y \in \mathbb{C}^{n \times k}$ such that $X^H Y = I_k$, it can be verified that

$$Z = \begin{bmatrix} X & 0 \\ 0 & Y \end{bmatrix} \in \mathbb{C}^{2n \times 2k}, \quad Z^H B Z = \begin{bmatrix} 0 & X^H Y \\ Y^H X & 0 \end{bmatrix} = \begin{bmatrix} 0 & I_k \\ I_k & 0 \end{bmatrix}.$$

Because $Z^H B Z$ is symmetric and has eigenvalues ± 1 with each repeating k times, there exists an orthogonal matrix $Q \in \mathbb{R}^{2k \times 2k}$ such that

$$(ZQ)^H B (ZQ) = Q^H (Z^H B Z) Q = \begin{bmatrix} I_k & \\ & -I_k \end{bmatrix} := J_{2k}$$

(cf. (1.6) with $k_+ = k_- = k$). We have

$$\begin{aligned} \text{tr}(X^H M X + Y^H K Y) &= \text{tr}(Z^H A Z) = \text{tr}((ZQ)^H A (ZQ)) \\ &\geq \inf_{W^H B W = J_{2k}} \text{tr}(W^H A W) \end{aligned} \tag{A.4}$$

$$= \sum_{i=1}^k \lambda_i - \sum_{i=1}^k (-\lambda_i) = 2 \sum_{i=1}^k \lambda_i, \tag{A.5}$$

where the equality in (A.5) is due to (1.4) and the size of $W \in \mathbb{C}^{2n \times 2k}$ in (A.4) is implied by the context. Since $X, Y \in \mathbb{C}^{n \times k}$ are arbitrary, subject to $X^H Y = I_k$, we conclude

$$\inf_{X^H Y = I_k} \text{tr}(X^H M X + Y^H K Y) \geq 2 \sum_{i=1}^k \lambda_i. \tag{A.6}$$

On the other hand, for $X = V_{(:,1:k)} \Lambda_k^{1/2}$ and $Y = U_{(:,1:k)} \Lambda_k^{-1/2}$ with $\Lambda_k = \text{diag}(\lambda_1, \dots, \lambda_k)$, we have by (A.3) and $V^H U = I_n$ that

$$X^H Y = \Lambda_k^{1/2} V_{(:,1:k)}^H U_{(:,1:k)} \Lambda_k^{-1/2} = \Lambda_k^{1/2} (V^H U)_{(1:k,1:k)} \Lambda_k^{-1/2} = I_k,$$

where $V_{(:,1:k)}$ and $U_{(:,1:k)}$ are the submatrices consisting of the first k columns of V and U , respectively, and $(V^H U)_{(1:k,1:k)}$ is the leading principal submatrix of $V^H U$. It can be verified that $\text{tr}(X^H M X) = \text{tr}(Y^H K Y) = \text{tr}(\Lambda_k) = \sum_{i=1}^k \lambda_i$. Together with (A.6), we arrive at

$$\sum_{i=1}^k \lambda_i = \frac{1}{2} \min_{X^H Y = I_k} \text{tr}(X^H M X + Y^H K Y), \tag{A.7}$$

which is precisely the trace minimization principle obtained in [24] for the case when both K and M are positive definite.

When one of K and M or both are assumed only semi-definite, we may use the limiting argument by applying (A.7) to $K + \epsilon I_n$ and $M + \epsilon I_n$ and then letting $\epsilon \rightarrow 0^+$ to conclude that (A.7) with \min replaced by \inf remains valid:

$$\sum_{i=1}^k \lambda_i = \frac{1}{2} \inf_{X^H Y = I_k} \text{tr}(X^H M X + Y^H K Y), \tag{A.8}$$

for $K \succeq 0$ and $M \succeq 0$. We comment that the decompositions in (A.3) remain valid so long as $M \succ 0$ even if K is singular, but the decompositions no longer hold if also M is singular.

Similarly, one may apply (1.4) to the generalized linear response eigenvalue problem to arrive at the same trace minimization principle obtained in [29,30]. We omit the detail here. In [25,30], the authors went on to develop LOBPCG type eigensolvers to solve the first few eigenpairs for the standard and generalized linear response eigenvalue problems.

Remark A.1. The block-diagonal structure in A of (A.2) ensures that the eigenvalues of the matrix pencil $A - \lambda B$ appear in pairs $\pm \lambda_i$. But purely from the perspective of application, (1.4) is equally applicable to $A - \lambda B$ with B being the same as in (A.2) while $A \succeq 0$ but without admitting the block-diagonal structure. In the latter case, the eigenvalues no longer appear in pairs.

A.2. Symplectic eigenvalues of positive definite matrices

The symplectic eigenvalues of positive definite matrices play important roles in classical Hamiltonian dynamics, quantum mechanics, and quantum information, among others [27,28] (and references therein). Let (cf. J_{2n} in the previous subsection)

$$\mathcal{J}_{2n} = \begin{bmatrix} 0 & I_n \\ -I_n & 0 \end{bmatrix} \in \mathbb{R}^{2n \times 2n}. \tag{A.9}$$

A matrix $Z \in \mathbb{R}^{2n \times 2n}$ is called a *symplectic* matrix if $Z^T \mathcal{J}_{2n} Z = \mathcal{J}_{2n}$. It can be seen that $\mathcal{J}_{2n}^T = -\mathcal{J}_{2n}$ and $\mathcal{J}_{2n}^2 = -I_{2n}$. In 1936, Williamson [48] proved the following matrix decomposition results: Given a symmetric positive definite matrix $A \in \mathbb{R}^{2n \times 2n}$ (yes for real matrix only), there exists a symplectic $Z \in \mathbb{R}^{2n \times 2n}$, i.e., $Z^T \mathcal{J}_{2n} Z = \mathcal{J}_{2n}$, such that

$$Z^T A Z = \text{diag}(\Lambda, \Lambda) \quad \text{with} \quad \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_n), \tag{A.10a}$$

where all $\lambda_i > 0$ and can be arranged in the ascending order

$$0 < \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n. \tag{A.10b}$$

The decomposition is referred to as *Williamson’s diagonal form* or *Williamson’s normal form* of A . The values λ_i are called the *symplectic eigenvalues* of A . We caution that these values λ_i are not really eigenvalues of A .

It follows from $Z^T \mathcal{J}_{2n} Z = \mathcal{J}_{2n}$ that $Z^{-T} = -\mathcal{J}_{2n} Z \mathcal{J}_{2n}$. Hence by (A.10), we have

$$AZ = Z^{-T} \begin{bmatrix} \Lambda & \\ & \Lambda \end{bmatrix} = -\mathcal{J}_{2n} Z \mathcal{J}_{2n} \begin{bmatrix} \Lambda & \\ & \Lambda \end{bmatrix} = \mathcal{J}_{2n} Z \begin{bmatrix} 0 & -\Lambda \\ \Lambda & 0 \end{bmatrix}, \tag{A.11}$$

which is in the form of the eigen-decomposition of matrix pencil $A - \lambda \mathcal{J}_{2n}$, and from which it can be read off that the eigenvalues of $A - \lambda \mathcal{J}_{2n}$ are $\pm \iota \lambda_i$ for $1 \leq i \leq n$, where $\iota = \sqrt{-1}$ is the imaginary unit.

Now let $B = \iota \mathcal{J}_{2n}$ which is Hermitian and has n eigenvalues $+1$ and n eigenvalues -1 . Matrix pencil $A - \lambda B$ is positive definite because $A - \lambda_0 B = A \succ 0$ for $\lambda_0 = 0$, and has only real finite eigenvalues which can be divided into two groups: n of the positive type and n of the negative type as in (1.5): $\lambda_i^- = -\lambda_i$ and $\lambda_i^+ = \lambda_i$ for $1 \leq i \leq n$.

For any $X \in \mathbb{R}^{2n \times 2k}$ such that $X^H \mathcal{J}_{2n} X = \mathcal{J}_{2k}$, we have $X^T B X = \iota \mathcal{J}_{2k}$ which is Hermitian. There exists a unitary $Q \in \mathbb{C}^{2k \times 2k}$ such that

$$(XQ)^H B (XQ) = Q^H (X^T B X) Q = \begin{bmatrix} I_k & \\ & -I_k \end{bmatrix} =: J_{2k}$$

(cf. (1.6) with $k_+ = k_- = k$). Hence by (1.4), we have

$$\begin{aligned} \text{tr}(X^T A X) &= \text{tr}((XQ)^H A (XQ)) \geq \inf_{W^H B W = J_{2k}} \text{tr}(W^H A W) \\ &= \sum_{i=1}^k \lambda_i - \sum_{i=1}^k (-\lambda_i) = 2 \sum_{i=1}^k \lambda_i. \end{aligned}$$

Since X is arbitrary, subject to $X^T \mathcal{J}_{2n} X = \mathcal{J}_{2k}$, we conclude that

$$X^T \mathcal{J}_{2n} X = \mathcal{J}_{2k} \implies \inf_{X^T \mathcal{J}_{2n} X = \mathcal{J}_{2k}} \text{tr}(X^T A X) \geq 2 \sum_{i=1}^k \lambda_i. \tag{A.12}$$

On the other hand, for $X = [Z_{(:,1:k)}, Z_{(:,n+1:n+k)}]$ (the submatrix of Z consisting of its first k columns and its $(n + 1)$ st to $(n + k)$ th columns), we have $X^T \mathcal{J}_{2n} X = \mathcal{J}_{2k}$. It follows from (A.11) that

$$AZ_{(:,1:k)} = \mathcal{J}_{2n} Z_{(:,n+1:n+k)} \Lambda_k, \quad AZ_{(:,n+1:n+k)} = -\mathcal{J}_{2n} Z_{(:,1:k)} \Lambda_k,$$

where $\Lambda_k = \text{diag}(\lambda_1, \dots, \lambda_k)$, and thus $AX = \mathcal{J}_{2n} X \begin{bmatrix} 0 & -\Lambda_k \\ \Lambda_k & 0 \end{bmatrix}$. Therefore

$$X^T A X = X^T \mathcal{J}_{2n} X \begin{bmatrix} 0 & -\Lambda_k \\ \Lambda_k & 0 \end{bmatrix} = \mathcal{J}_{2k} \begin{bmatrix} 0 & -\Lambda_k \\ \Lambda_k & 0 \end{bmatrix} = \begin{bmatrix} \Lambda_k & \\ & \Lambda_k \end{bmatrix}$$

yielding $\text{tr}(X^T AX) = 2 \sum_{i=1}^k \lambda_i$, which together with (A.12) lead to

$$\sum_{i=1}^k \lambda_i = \frac{1}{2} \min_{X^T \mathcal{J}_{2n} X = \mathcal{J}_{2k}} \text{tr}(X^T AX). \quad (\text{A.13})$$

Equation (A.13) is precisely the trace minimization principle obtained in [27], and was recently used as the theoretical foundation in [32] for numerically computing the symplectic eigenvalues λ_i of A . Perturbation bounds can also be found in [27,49].

So far $A \in \mathbb{R}^{2n \times 2n}$ is assumed to be symmetric positive definite. However, the trace minimization principle (A.13) can be extended to semi-definite A . Notice that $A - \lambda B$ with $B = \iota \mathcal{J}_{2n}$ is a positive semi-definite matrix pencil when $A \in \mathbb{R}^{2n \times 2n}$ is symmetric positive semi-definite. Regardless, $A - \lambda B$ still has $2n$ eigenvalues $\pm \lambda_i$ [23], but instead of (A.10b) we will have

$$0 \leq \lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_n, \quad (\text{A.14})$$

i.e., possibly $\lambda_1 = 0$. Indeed $\lambda_1 = 0$ if and only if A is singular. By applying (A.13) to $A + \epsilon I_{2n}$ and then letting $\epsilon \rightarrow 0^+$, we arrive at

$$\sum_{i=1}^k \lambda_i = \frac{1}{2} \inf_{X^T \mathcal{J}_{2n} X = \mathcal{J}_{2k}} \text{tr}(X^T AX) \quad (\text{A.15})$$

for generally $A \in \mathbb{R}^{2n \times 2n}$ being symmetric positive semi-definite.

References

- [1] E. Kokiopoulou, J. Chen, Y. Saad, Trace optimization and eigenproblems in dimension reduction methods, *Numer. Linear Algebra Appl.* 18 (3) (2011) 565–602.
- [2] T. Ngo, M. Bellalij, Y. Saad, The trace ratio optimization problem for dimensionality reduction, *SIAM J. Matrix Anal. Appl.* 31 (5) (2010) 2950–2971.
- [3] L.-H. Zhang, L.-Z. Liao, M.K. Ng, Fast algorithms for the generalized Foley-Sammon discriminant analysis, *SIAM J. Matrix Anal. Appl.* 31 (4) (2010) 1584–1605.
- [4] L.-H. Zhang, L.-Z. Liao, M.K. Ng, Superlinear convergence of a general algorithm for the generalized Foley-Sammon discriminant analysis, *J. Optim. Theory Appl.* 157 (3) (2013) 853–865.
- [5] L.-H. Zhang, R.-C. Li, Maximization of the sum of the trace ratio on the Stiefel manifold, I: theory, *Sci. China Math.* 57 (12) (2014) 2495–2508.
- [6] L.-H. Zhang, R.-C. Li, Maximization of the sum of the trace ratio on the Stiefel manifold, II: computation, *Sci. China Math.* 58 (7) (2015) 1549–1566.
- [7] D. Chu, L. Liao, M.K. Ng, X. Zhang, Sparse canonical correlation analysis: new formulation and algorithm, *IEEE Trans. Pattern Anal. Mach. Intell.* 35 (12) (2013) 3050–3065.
- [8] J.P. Cunningham, Z. Ghahramani, Linear dimensionality reduction: survey, insights, and generalizations, *J. Mach. Learn. Res.* 16 (2015) 2859–2900.
- [9] P. Horst, Generalized canonical correlations and their applications to experimental data, *J. Clin. Psychol.* 17 (4) (1961) 331–347.
- [10] L. Wang, L.-H. Zhang, Z. Bai, R.-C. Li, Orthogonal canonical correlation analysis and applications, *Optim. Methods Softw.* 35 (4) (2020) 787–807.
- [11] L.-H. Zhang, L. Wang, Z. Bai, R.-C. Li, A self-consistent-field iteration for orthogonal canonical correlation analysis, *IEEE Trans. Pattern Anal. Mach. Intell.* 44 (2) (2022) 890–904, <https://doi.org/10.1109/TPAMI.2020.3012541>.

- [12] K. Fan, On a theorem of Weyl concerning eigenvalues of linear transformations. I, Proc. Natl. Acad. Sci. USA 35 (11) (1949) 652–655.
- [13] R.A. Horn, C.R. Johnson, Matrix Analysis, 2nd edition, Cambridge University Press, New York, NY, 2013.
- [14] G.H. Golub, Q. Ye, An inverse free preconditioned Krylov subspace method for symmetric generalized eigenvalue problems, SIAM J. Sci. Comput. 24 (1) (2002) 312–334.
- [15] A. Klinvex, F. Saied, A. Sameh, Parallel implementations of the trace minimization scheme TraceMIN for the sparse symmetric eigenvalue problem, Comput. Math. Appl. 65 (3) (2013) 460–468, <https://doi.org/10.1016/j.camwa.2012.06.011>.
- [16] A.V. Knyazev, Toward the optimal preconditioned eigensolver: locally optimal block preconditioned conjugate gradient method, SIAM J. Sci. Comput. 23 (2) (2001) 517–541.
- [17] R.-C. Li, Accuracy of computed eigenvectors via optimizing a Rayleigh quotient, BIT Numer. Math. 44 (3) (2004) 585–593.
- [18] R.-C. Li, Rayleigh quotient based optimization methods for eigenvalue problems, in: Z. Bai, W. Gao, Y. Su (Eds.), Matrix Functions and Matrix Equations, in: Series in Contemporary Applied Mathematics, vol. 19, World Scientific, Singapore, 2015, pp. 76–108.
- [19] B.N. Parlett, The Symmetric Eigenvalue Problem, SIAM, Philadelphia, 1998.
- [20] P. Quillen, Q. Ye, A block inverse-free preconditioned Krylov subspace method for symmetric generalized eigenvalue problems, J. Comput. Appl. Math. 233 (5) (2010) 1298–1313.
- [21] A.H. Sameh, J.A. Wisniewski, A trace, minimization algorithm for the generalized eigenvalue problem, SIAM J. Numer. Anal. 19 (6) (1982) 1243–1259.
- [22] J. Kovač-Striko, K. Veselić, Trace minimization and definiteness of symmetric pencils, Linear Algebra Appl. 216 (1995) 139–158.
- [23] X. Liang, R.-C. Li, Z. Bai, Trace minimization principles for positive semi-definite pencils, Linear Algebra Appl. 438 (2013) 3085–3106.
- [24] Z. Bai, R.-C. Li, Minimization principle for linear response eigenvalue problem, I: theory, SIAM J. Matrix Anal. Appl. 33 (4) (2012) 1075–1100.
- [25] Z. Bai, R.-C. Li, Minimization principles for the linear response eigenvalue problem II: computation, SIAM J. Matrix Anal. Appl. 34 (2) (2013) 392–416.
- [26] D.J. Thouless, Vibrational states of nuclei in the random phase approximation, Nucl. Phys. 22 (1) (1961) 78–95.
- [27] R. Bhatia, T. Jain, On symplectic eigenvalues of positive definite matrices, J. Math. Phys. 56 (2015) 112201.
- [28] K.R. Parthasarathy, The symmetry group of Gaussian states in $L^2(\mathbb{R}^n)$, in: A.N. Shiryayev, S.R.S. Varadhan, E.L. Presman (Eds.), Prokhorov and Contemporary Probability Theory, Springer Berlin Heidelberg, Berlin, Heidelberg, 2013, pp. 349–369.
- [29] Z. Bai, R.-C. Li, Minimization principles and computation for the generalized linear response eigenvalue problem, BIT Numer. Math. 54 (1) (2014) 31–54.
- [30] Z. Bai, R.-C. Li, W.-W. Lin, Linear response eigenvalue problem solved by extended locally optimal preconditioned conjugate gradient methods, Sci. China Math. 59 (8) (2016) 1443–1460.
- [31] D. Rocca, Z. Bai, R.-C. Li, G. Galli, A block variational procedure for the iterative diagonalization of non-Hermitian random-phase approximation matrices, J. Chem. Phys. 136 (2012) 034111.
- [32] N.T. Son, P.-A. Absil, B. Gao, T. Stykel, Symplectic eigenvalue problem via trace minimization and Riemannian optimization, arXiv:2101.02618, 2021.
- [33] D. Kressner, M.M. Pandur, M. Shao, An indefinite variant of LOBPCG for definite matrix pencils, Numer. Algorithms (2013) 1–23.
- [34] M.M. Pandur, Preconditioned gradient iterations for the eigenproblem of definite matrix pairs, Electron. Trans. Numer. Anal. 51 (2019) 331–362.
- [35] H. Liu, A.M.-C. So, W. Wu, Quadratic optimization with orthogonality constraint: explicit Łojasiewicz exponent and linear convergence of retraction-based line-search and stochastic variance-reduced gradient methods, Math. Program., Ser. A 178 (1–2) (2019) 215–262.
- [36] P.-A. Absil, R. Mahony, R. Sepulchre, Optimization Algorithms on Matrix Manifolds, Princeton University Press, Princeton, NJ, 2008.
- [37] P. Birtea, I. Çaşu, D. Comănescu, First order optimality conditions and steepest descent algorithm on orthogonal Stiefel manifolds, Opt. Lett. 13 (2019) 1773–1791.
- [38] R. Brockett, Dynamical systems that sort lists, diagonalize matrices, and solve linear programming problems, Linear Algebra Appl. 146 (1991) 79–91.
- [39] E. Anderson, Z. Bai, C. Bischof, J. Demmel, J. Dongarra, J.D. Croz, A. Greenbaum, S. Hammarling, A. McKenney, S. Ostrouchov, D. Sorensen, LAPACK Users' Guide, 3rd edition, SIAM, Philadelphia, 1999.

- [40] Z. Bai, J. Demmel, J. Dongarra, A. Ruhe, H. van der Vorst (Eds.), *Templates for the Solution of Algebraic Eigenvalue Problems: A Practical Guide*, SIAM, Philadelphia, 2000.
- [41] J. Demmel, *Applied Numerical Linear Algebra*, SIAM, Philadelphia, PA, 1997.
- [42] G.W. Stewart, J.-G. Sun, *Matrix Perturbation Theory*, Academic Press, Boston, 1990.
- [43] R. Bhatia, *Matrix Analysis*, Graduate Texts in Mathematics, vol. 169, Springer, New York, 1996.
- [44] K. Veselić, *Damped Oscillations of Linear Systems*, Lecture Notes in Mathematics, vol. 2023, Springer, Berlin, 2011.
- [45] A. Mal'cev, *Foundation of Linear Algebra*, Freeman, 1963.
- [46] I. Gohberg, P. Lancaster, L. Rodman, *Indefinite Linear Algebra and Applications*, Birkhäuser, Basel, Switzerland, 2005.
- [47] K. Veselić, A Jacobi eigenreduction algorithm for definite matrix pairs, *Numer. Math.* 64 (1993) 241–269.
- [48] J. Williamson, On the algebraic problem concerning the normal forms of linear dynamical systems, *Am. J. Math.* 58 (1) (1936) 141–163.
- [49] M. Idel, S. Soto Gaona, M.M. Wolf, Perturbation bounds for Williamson's symplectic normal form, *Linear Algebra Appl.* 525 (2017) 45–58.