

AN EIGENVALUE-BASED METHOD FOR THE UNBALANCED PROCRUSTES PROBLEM*

LEI-HONG ZHANG[†], WEI HONG YANG[‡], CHUNGEN SHEN[§], AND JIAQI YING[¶]

Abstract. In this paper, we propose a novel eigenvalue-based approach to solving the unbalanced orthogonal Procrustes problem. By making effective use of the necessary condition for the global minimizer and the orthogonal constraint, we shall first show that the unbalanced Procrustes problem can be equivalently transformed into an eigenvalue minimization whose solution can be computed by solving a related eigenvector-dependent nonlinear eigenvalue problem. Through the exploitation of certain techniques in the nonlinear eigenvalue computations, we adapt the standard self-consistent field (SCF) iteration to solve the resulting nonlinear eigenvalue problem. Theoretical convergence analysis of this customized SCF iteration is performed, and practical strategies for a more efficient numerical implementation are discussed. Our numerical experience on preliminary tests indicates that the proposed eigenvalue-based SCF iteration is a promising method for the unbalanced orthogonal Procrustes problem.

Key words. Procrustes problem, nonlinear eigenvalue problem, self-consistent field iterations

AMS subject classifications. 65F15, 65H17, 90C30

DOI. 10.1137/19M1270872

1. Introduction. Given matrices $C \in \mathbb{R}^{m \times n}$ and $D \in \mathbb{R}^{m \times \ell}$ with $n \geq \ell$, the Procrustes problem minimizes the square of the Frobenius norm $\|CX - D\|_F^2$ over the Stiefel manifold $\mathbb{O}^{n \times \ell} := \{X \in \mathbb{R}^{n \times \ell} | X^T X = I_\ell\}$, i.e.,

$$(1.1) \quad \min_{X \in \mathbb{O}^{n \times \ell}} \|CX - D\|_F^2.$$

Equivalently, by expanding $\|CX - D\|_F^2$ and ignoring the constant $\|D\|_F^2$, (1.1) can be written as a problem to minimize a quadratic function over $\mathbb{O}^{n \times \ell}$, i.e.,

$$(1.2) \quad \min_{X \in \mathbb{O}^{n \times \ell}} \{f(X) := \text{trace}(X^T A X) + 2 \text{trace}(X^T B)\},$$

where $A = C^T C$ and $B = -C^T D$.

Originally, the terminology of the Procrustes problem (with $n = \ell$) is due to Hurley and Cattell [19], who tried to find an estimated transformation matrix X to relate a factor structure C from factor analysis to a specific hypothetical matrix D (see also [17]). The special case with $n = \ell$ is commonly called the balanced Procrustes

*Received by the editors June 27, 2019; accepted for publication (in revised form) by B. Vandereycken March 27, 2020; published electronically July 1, 2020.

<https://doi.org/10.1137/19M1270872>

Funding: The work of the first author was partially supported by the National Natural Science Foundation of China grant 11671246, and the National Key R & D Program of China (2018YFB0204404). The work of the second author was supported by the National Natural Science Foundation of China grant 11971118.

[†]School of Mathematical Sciences and Institute of Computational Science, Soochow University, Suzhou 215006, Jiangsu, China, and School of Mathematics, Shanghai University of Finance and Economics, Shanghai, 200433, China (longzlh@suda.edu.cn).

[‡]School of Mathematical Sciences, Fudan University, Shanghai, 200433, People's Republic of China (whyang@fudan.edu.cn).

[§]College of Science, University of Shanghai for Science and Technology, Shanghai 200093, China (shenchungen@gmail.com).

[¶]School of Mathematics, Shanghai University of Finance and Economics, Shanghai 200433, China (candy201314l@qq.com).

problem [14, 16, 30, 38, 49], and until now, there are several variants of (1.1) arising from statistics, psychometrika, biometrics, data mining, and other fields. For example, instead of the Frobenius norm, [41] considers the solutions of the Procrustes problem in a family of orthogonally invariant norms; [18, 3] discuss the solutions X subject to the set of symmetric matrices and certain closed cones (the cone of symmetric positive semidefinite and the cone of (symmetric) elementwise nonnegative matrices), respectively.

The minimization (1.1) with $n > \ell$ is referred to as the unbalanced Procrustes problem [8, 12, 14, 19, 17, 49]. Applications of the unbalanced Procrustes problem can arise, for example, in the orthogonal least squares regression (OLSR) for feature extraction [51, 30], the multidimensional similarity structure analysis (SSA) [6, Chapter 19], and the Maxbet problem [10, 26] in the canonical correlation analysis. The model OLSR is a supervised learning method in data mining, which aims at finding an orthogonal transformation matrix $X \in \mathbb{O}^{n \times \ell}$ to project high dimensional data (with dimension n) into a lower dimensional space (with dimension $\ell \ll n$). The mathematical formulation of OLSR is exactly an unbalanced Procrustes problem (1.1), and in section 6.2, we will give a detailed description of OLSR and apply our proposed method to solve it. In SSA, two configurations represented by matrices C and D are *similar* if they can be brought to a complete match by rigid motions and dilations [6, Chapter 19], and the Procrustes problem (1.1) is used to fit one configuration C to the other D . As pointed out in [6, section 19.8] two configurations are generally of the same dimensionality (i.e., $n = \ell$), but sometimes it is useful for them to have different dimensionality. A possible treatment of the latter leads to the unbalanced Procrustes problem, and in section 6.3, we will consider such an example. Another application is from solving a subproblem [26, section 4] of Maxbet in the form of (1.2), in which ℓ is the reduced dimension of a high dimensional random variable, and usually $\ell \ll n$.

For numerical methods, it is known that the balanced Procrustes problem admits a closed form solution which can be computed by the polar decomposition (e.g., [16]). Unfortunately, for the unbalanced case, there is no direct way to compute the solution X . In fact, the special case with $\ell = 1$ is widely known as the trust-region subproblem of the trust-region method in optimization (see, e.g., [31, 29]); even for this special case, it may admit a local but nonglobal minimizer [28], and computing the global solution is generally not an easy task. Due to the lack of closed form solution, necessary or sufficient conditions for the local and/or global minimizer of the unbalanced Procrustes problem (or the minimization (1.2)) have been developed; particularly, [8, 14, 50] establish a couple of useful and in-depth conditions that can be used to check if a computed one can be a local or a global solution.

For the general case $n > \ell > 1$, certain numerical approaches for computing an approximation to (1.2) or (1.1) are available. Besides resorting to generic optimization methods built on the Stiefel manifold (e.g., [1, 2, 12, 42]), several structure-exploiting algorithms have been proposed. For example, [5] presents an iterative method based on relaxation which produces a sequence that gradually reduces the objective function; in each iteration, the method finds a best plane rotation in a particularly chosen plane, and the resulting subproblem of computing the best plane rotation is called a planner Procrustes problem (i.e., $m = n = 2$ and $\ell = 1$ in (1.1)). In [49], a successive projection minimization approach is introduced in which, by alternatively fixing all but one column of X , each iteration solves a trust-region subproblem (i.e., (1.2) with $\ell = 1$) to update the approximation. Other numerical algorithms include the generalized power iteration in [30] and Newton-type iterations in [38]. However, none of them is an eigenvalue-based method so that the modern power of the state-

of-the-art eigensolvers can be exploited.

In this paper, we propose a self-consistent-field (SCF) iteration for the unbalanced Procrustes problem (1.1) or its equivalent minimization (1.2). The SCF iteration is an efficient method for solving a class of eigenvector-dependent nonlinear eigenvalue problems (NEPvs) [7], and can exploit the presently used relatively mature techniques in eigenvalue computations. It was originally used in electronic structure calculations (e.g., [25, 27, 33, 43]), and appears very useful in many applications arising recently from data science (e.g., [4, 24, 40, 44, 45, 48, 46, 47]). A general discussion on SCF together with some convergence results can be found in, e.g., [7]. For the unbalanced Procrustes problem (1.1) and/or (1.2), our contributions in this paper include the following three aspects: (i) by making use of the necessary conditions for the global minimizer and the orthogonal constraint $X \in \mathbb{O}^{n \times \ell}$, we construct an equivalent eigenvalue minimization of (1.2) for which a related NEPv can be established, (ii) by the general SCF iteration scheme and the special structure of (1.2), we propose a customized SCF iteration and establish the convergence, and (iii) we introduce a subspace acceleration procedure to improve the performance of the basic SCF. Numerical simulations are carried out and demonstrate that our eigenvalue-based SCF iteration is an efficient method for the unbalanced orthogonal Procrustes problem, particularly for the case $\ell \ll n$.

We organize the remainder of this paper as follows: In section 2, we first summarize some basic optimality conditions for the local and/or global minimizer of (1.2). Relying on the necessary conditions, in section 3, we then transform the problem (1.2) into an eigenvalue minimization, and also propose the eigenvector-dependent NEPv. The basic form of our SCF iteration is introduced in this section. In section 4, we shall present the convergence analysis of the SCF iteration. To make the basic SCF iteration more efficient, in section 5, we then discuss two practical treatments: the inexact SCF iteration and the subspace acceleration. Numerical results of our practical SCF iteration are reported in section 6, and final remarks are drawn in section 7.

Notation. Throughout this paper, all vectors are column vectors and are typeset in bold lowercase letters. For a matrix $C \in \mathbb{R}^{n \times m}$, C^T and $\mathcal{R}(C)$ denote its transpose and range space, respectively. The $n \times n$ identity matrix is I_n , and \mathbf{e}_j is the j th column of an identity matrix whose size is determined by the context. If C is square, we denote the set of its eigenvalues by $\text{eig}(C)$. When all of the eigenvalues of C are real, we order them as

$$\lambda_1(C) \leq \lambda_2(C) \leq \cdots \leq \lambda_n(C).$$

The eigendecomposition of A is given by $A = U\Theta U^T$ and $\Theta = \text{diag}(\theta_1, \dots, \theta_n)$ with the arranged eigenvalues

$$\theta_1 \leq \theta_2 \leq \cdots \leq \theta_n,$$

where $\theta_i = \lambda_i(A)$ is reserved in this paper. When A is positive semidefinite (positive definite), we denote it by $A \succcurlyeq 0$ ($\succ 0$). The notation $\|\cdot\|_{\text{ui}}$ stands for any unitary invariant norm, and it holds that

$$\|XYZ\|_{\text{ui}} \leq \|X\|_2 \|Y\|_{\text{ui}} \|Z\|_2.$$

To simplify our presentation, we shall also adopt MATLAB-like convention to access the (i, j) th entry of A as $A_{(i,j)}$.

2. Optimality conditions. For the trust-region subproblem (i.e., $\ell = 1$), the optimality conditions due to Gay [15] and Moré and Sorensen [29] (see also [35] and [31, Theorem 4.1]) are well known. The generalization for (1.2) or, equivalently, (1.1) has been discussed in [8, 14], and we summarize some basic facts below.

2.1. Local optimality conditions.

LEMMA 2.1 (first-order optimality [8, Theorem 3.8] and [14, Theorem 3.1]). *If $X \in \mathbb{O}^{n \times \ell}$ is a local minimizer, then there is a symmetric matrix $\Lambda \in \mathbb{R}^{\ell \times \ell}$ such that*

$$(2.1) \quad AX + X\Lambda = -B.$$

Note from (2.1) and the constraint $X^T X = I_\ell$ that $\Lambda = -X^T B - X^T A X$ is indeed determined by the local solution X , and also the matrix $X^T B = B^T X$ is symmetric. Moreover, we will further see in Theorem 2.4 that if X is a global minimizer, the matrix $X^T A X + \Lambda = -X^T B$ is positive semidefinite.

Note that the tangent space $\mathcal{T}_X \mathbb{O}^{n \times \ell}$ at any $X \in \mathbb{O}^{n \times \ell}$ can be expressed as, e.g., [2, p. 42]

$$(2.2a) \quad \mathcal{T}_X \mathbb{O}^{n \times \ell} := \{H \in \mathbb{R}^{n \times \ell} : H^T X + X^T H = 0\}$$

$$(2.2b) \quad = \{H = XK + (I_n - XX^T)J : K = -K^T \in \mathbb{R}^{\ell \times \ell}, J \in \mathbb{R}^{n \times \ell}\}.$$

With the standard inner product (or the Frobenius inner product) induced on $\mathcal{T}_X \mathbb{O}^{n \times \ell}$

$$\langle Z, Y \rangle = \text{trace}(Z^T Y) \quad \forall Z, Y \in \mathcal{T}_X \mathbb{O}^{n \times \ell},$$

it is known that the orthogonal projection of $Z \in \mathbb{R}^{n \times \ell}$ onto the tangent space $\mathcal{T}_X \mathbb{O}^{n \times \ell}$ at X is given by [2, eq. (3.35)]

$$(2.3) \quad \Pi_X(Z) := Z - X \frac{X^T Z + Z^T X}{2}.$$

LEMMA 2.2 (second-order optimality [8, Theorem 3.2] and [14, p. 609]). *If X is a local minimizer of (1.2), then*

$$\text{trace}(H^T A H) + \text{trace}(H \Lambda H^T) \geq 0 \quad \forall H \in \mathcal{T}_X \mathbb{O}^{n \times \ell};$$

moreover, if the strict inequality holds for any nonzero tangent vector $H \in \mathcal{T}_X \mathbb{O}^{n \times \ell}$ then X is a strict local minimizer.

It is worth mentioning that Lemma 2.2 is just the generalization of the second-order optimality condition [28, Lemma 2.1] for the trust-region subproblem (see also [35, 29] and [31, Theorem 4.1]), i.e., $\ell = 1$. In the general case $\ell > 1$, we can see that the elements of Λ are just the Lagrangian multipliers, which are uniquely given by X .

With the representation of the tangent vector given in (2.2b), we can further refine, as in [8] for the Procrustes problem (1.1), the second-order optimality condition of Lemma 2.2 as follows.

COROLLARY 2.3. *If X is a local minimizer of (1.2), then*

(i) *for any skew-symmetric $K \in \mathbb{R}^{\ell \times \ell}$ and any $J \in \mathbb{R}^{n \times \ell}$, it holds that*

$$(2.4) \quad \begin{aligned} & \text{trace}(J^T [A + X\Lambda X^T + (2I_n - XX^T)BX^T]J) - \text{trace}(X^T J \Lambda J^T X) \\ & + \text{trace}(J \Lambda J^T) + 2 \text{trace}(J^T (XB^T X + B)K) - \text{trace}(K^T X^T B K) \geq 0, \end{aligned}$$

and moreover, if the strict inequality holds for any skew-symmetric $K \in \mathbb{R}^{\ell \times \ell}$ and any $J \in \mathbb{R}^{n \times \ell}$ with $XK + (I_n - XX^T)J \neq 0$, then X is a strict local minimizer;

(ii) if $\ell \geq 2$, then $\text{trace}(X^T B) \leq 0$.

Proof. For item (i), it is a straightforward task to verify (2.4) by using the facts that (a) $X^T X = I_\ell$, (b) Lemma 2.1, (c) Lemma 2.2, and (d) the representation of $\mathcal{T}_X \mathbb{O}^{n \times \ell}$ in (2.2b). The detailed manipulations are omitted.

For (ii), we particularly choose $J = 0$ in (2.4), or equivalently, $H = XK$ in Lemma 2.2 to have

$$\text{trace}(K^T X^T B K) \leq 0 \quad \forall K = -K^T \in \mathbb{R}^{\ell \times \ell}.$$

By the symmetry of $X^T B$, let $X^T B = W^T \Sigma W$ be the spectral decomposition of $X^T B$, and $\Sigma = \text{diag}(\mu_1, \dots, \mu_\ell)$. Choose $K = K_{12} := W \text{diag}(\Xi, 0, \dots, 0) W^T$, where $\Xi = \begin{bmatrix} 0 & \zeta \\ -\zeta & 0 \end{bmatrix}$ for $\zeta > 0$; then it follows that

$$(\mu_1 + \mu_2)\zeta^2 \leq 0 \implies \mu_1 + \mu_2 \leq 0.$$

Similarly, we can define K_{ij} for $1 \leq i < j \leq \ell$ and add all of the resulting inequalities to have

$$(\ell - 1)(\mu_1 + \dots + \mu_\ell)\zeta^2 \leq 0 \implies \mu_1 + \dots + \mu_\ell = \text{trace}(X^T B) \leq 0.$$

This completes the proof. \square

As the special case of $\ell = 1$, we remark that Corollary 2.3(ii) may not hold. As an example with $A = \text{diag}(1, -1)$ and $B = [0, 1]^T$, the minimization

$$\min_{x_1^2 + x_2^2 = 1} x_1^2 - x_2^2 + 2x_2$$

has a local minimizer $X = [0, 1]^T$, but $B^T X = 1 > 0$.

2.2. Global optimality conditions. Suppose X_{opt} is a global minimizer. Then for any orthogonal $Q \in \mathbb{R}^{\ell \times \ell}$, we have

$$\begin{aligned} f(X_{\text{opt}}) &\leq f(X_{\text{opt}}Q) = \text{trace}(Q^T X_{\text{opt}}^T A X_{\text{opt}} Q) + 2 \text{trace}(Q^T X_{\text{opt}}^T B) \\ &= \text{trace}(X_{\text{opt}}^T A X_{\text{opt}}) + 2 \text{trace}(Q^T X_{\text{opt}}^T B). \end{aligned}$$

This leads to $\text{trace}(X_{\text{opt}}^T B) \leq \text{trace}(Q^T X_{\text{opt}}^T B)$ for any orthogonal $Q \in \mathbb{R}^{\ell \times \ell}$. Now, let $X_{\text{opt}}^T B = U \Gamma U^T$ be the spectral decomposition (note that $X_{\text{opt}}^T B$ is symmetric), and we can conclude $\Gamma = \text{diag}(\gamma_1, \dots, \gamma_\ell) \preceq 0$ because, otherwise, we can choose $Q = U \text{diag}(-\text{sign}(\gamma_1), \dots, -\text{sign}(\gamma_\ell)) U^T$ to have $\text{trace}(X_{\text{opt}}^T B) > \text{trace}(Q^T X_{\text{opt}}^T B)$. This then shows $-X_{\text{opt}}^T B \succcurlyeq 0$. Hence, using (2.1), one has that

$$(2.5) \quad -B^T X = X^T A X + \Lambda \succcurlyeq 0.$$

Consequently, the Cauchy interlacing inequalities [37]

$$\theta_i \leq \lambda_i(X^T A X) \leq \theta_{n-\ell+i}$$

and Weyl's theorem [37]

$$0 \leq \lambda_{\ell-r+1}(X^T A X + \Lambda) \leq \lambda_{\ell-r+j}(X^T A X) + \lambda_{\ell-j+1}(\Lambda) \leq \theta_{n-r+j} + \lambda_{\ell-j+1}(\Lambda)$$

$\forall 1 \leq j \leq r, 1 \leq r \leq \ell$ yield a necessary optimality condition (see, e.g., [14, Theorem 3.8] or [50, Theorem 3.2]).

THEOREM 2.4. *If X is a global minimizer of (1.2), then*

$$(2.6) \quad X^T A X + \Lambda = -X^T B = -B^T X \succcurlyeq 0,$$

which implies

$$\lambda_i(\Lambda) \geq -\theta_{n-i+1} \quad \text{for } 1 \leq i \leq \ell.$$

There are other types of necessary or sufficient optimality conditions for (1.2), and the reader is referred to [8, 14, 50] for more details.

3. A self-consistent field iteration. In this section, we will propose a novel eigenvalue-based algorithm for solving (1.2). Our method essentially targets the KKT condition given in (2.1) but makes effective use of the necessary condition (2.6) for the global minimizer.

3.1. A related NEPv. Relying on the following well-known von Neumann's trace inequality [39] (see also [34, eq. 6.81]), we will first, equivalently, transform (1.2) into an eigenvalue-related minimization, and then introduce an eigenvector-dependent NEPv for which an SCF iteration can apply.

LEMMA 3.1. *If C_1, C_2 are n -by- n complex matrices with singular values*

$$\sigma_1(C_1) \leq \cdots \leq \sigma_n(C_1), \quad \sigma_1(C_2) \leq \cdots \leq \sigma_n(C_2),$$

respectively, then

$$-\sum_{i=1}^n \sigma_i(C_1) \sigma_i(C_2) \leq \operatorname{Re}(\operatorname{trace}(C_1^H C_2)) \leq \sum_{i=1}^n \sigma_i(C_1) \sigma_i(C_2).$$

Equality holds on the right when the singular value decomposition of C_2 is $C_2 = U \Sigma_2 V^H$, and equality on the left holds when $C_2 = -U \Sigma_2 V^H$, where $C_1 = U \Sigma_1 V^H$ is the singular value decomposition of C_1 , and $\Sigma_i = \operatorname{diag}(\sigma_1(C_i), \sigma_2(C_i), \dots, \sigma_n(C_i))$ for $i = 1, 2$.

Define a symmetric matrix function

$$(3.1) \quad E(X) := A + X B^T + B X^T \in \mathbb{R}^{n \times n}.$$

By making use of the necessary optimality condition (2.1) and the orthogonal constraint $X^T X = I_\ell$, we have the following lemma.

LEMMA 3.2. *Let $X \in \mathbb{O}^{n \times \ell}$ satisfying $X^T B = B^T X$. Then X is a KKT point satisfying (2.1) if and only if X is an orthonormal eigenbasis matrix of $E(X)$.*

Proof. For a KKT point X , by $X^T X = I_\ell$, we can rewrite (2.1) as

$$(A + X B^T + B X^T) X = -X \Lambda + X B^T X \iff E(X) X = X \Psi_X$$

where $\Psi_X = -\Lambda + B^T X$ is symmetric. The converse also holds by the above relation and $X^T B = B^T X$. \square

The following theorem further ensures that any global minimizer X_{opt} of (1.2) is an orthonormal eigenbasis matrix of $E(X_{\text{opt}})$ associated with the first ℓ smallest eigenvalues.

THEOREM 3.3. For the matrix function $E(X)$ in (3.1), we have

$$(3.2) \quad \min_{X \in \mathbb{O}^{n \times \ell}} \sum_{i=1}^{\ell} \lambda_i(E(X)) = \min_{X \in \mathbb{O}^{n \times \ell}} f(X).$$

That is, the minimization (1.2) is equivalent to finding $X \in \mathbb{O}^{n \times \ell}$ so that the sum of the first ℓ smallest eigenvalues of $E(X)$ is minimized.¹ Moreover, any global minimizer X_{opt} of (1.2) is an orthonormal eigenbasis matrix of $E(X_{\text{opt}})$, i.e.,

$$(3.3) \quad E(X_{\text{opt}})X_{\text{opt}} = X_{\text{opt}}\Psi_{\text{opt}},$$

associated with the first ℓ smallest eigenvalues, and $\text{trace}(\Psi_{\text{opt}}) = f(X_{\text{opt}})$.

Proof. Let $X \in \mathbb{O}^{n \times \ell}$ be arbitrarily given, and let Z be an orthonormal eigenbasis matrix of $E(X)$ corresponding to the ℓ smallest eigenvalues. Note that

$$(3.4) \quad \sum_{i=1}^{\ell} \lambda_i(E(X)) = \text{trace}(Z^T E(X) Z) = \text{trace}(Z^T A Z) + 2 \text{trace}(Z^T B X^T Z).$$

Moreover, we can choose a suitable orthonormal basis $\tilde{Z} = ZQ$ for $\mathcal{R}(Z)$ to ensure $\text{trace}(\tilde{Z}^T E(X) \tilde{Z}) \geq f(X_{\text{opt}})$. The orthonormal matrix $Q \in \mathbb{R}^{\ell \times \ell}$ is obtained by minimizing $\text{trace}(Q^T Z^T B)$ over $\mathbb{O}^{\ell \times \ell}$. In fact,

$$\min_{Q^T Q = I_{\ell}} \text{trace}(Q^T Z^T B) = -\text{trace}(\Sigma) = -\sum_{i=1}^{\ell} \sigma_i(Z^T B),$$

where the minimum is achieved when $Q = -UV^T$ and $Z^T B = U\Sigma V^T$ is the SVD of $Z^T B$ and $\Sigma = \text{diag}(\sigma_1(Z^T B), \dots, \sigma_{\ell}(Z^T B))$ with $\sigma_1(Z^T B) \leq \dots \leq \sigma_{\ell}(Z^T B)$. Note that UV^T is the orthogonal polar factor [16, section 9.4.3] of $Z^T B$. On the other hand, by using Lemma 3.1, we have

$$\text{trace}(Z^T B X^T Z) \geq -\sum_{i=1}^{\ell} \sigma_i(Z^T B) \sigma_i(X^T Z),$$

where $\sigma_1(X^T Z) \leq \dots \leq \sigma_{\ell}(X^T Z)$ are the singular values of $X^T Z$, which are all less than 1 by the facts $X, Z \in \mathbb{O}^{n \times \ell}$. Setting $\tilde{Z} = ZQ$, we have

$$\begin{aligned} \text{trace}(\tilde{Z}^T B X^T \tilde{Z}) &= \text{trace}(Z^T B X^T Z) \\ &\geq -\sum_{i=1}^{\ell} \sigma_i(Z^T B) \sigma_i(X^T Z) \geq -\sum_{i=1}^{\ell} \sigma_i(Z^T B) = \text{trace}(\tilde{Z}^T B), \end{aligned}$$

and by (3.4)

$$\sum_{i=1}^{\ell} \lambda_i(E(X)) = \text{trace}(\tilde{Z}^T E(X) \tilde{Z}) = \text{trace}(\tilde{Z}^T A \tilde{Z}) + 2 \text{trace}(\tilde{Z}^T B X^T \tilde{Z})$$

¹We remark that (3.2) is equivalent to

$$\min_{X \in \mathbb{O}^{n \times \ell}} \min_{Y \in \mathbb{O}^{n \times \ell}} \text{trace}(Y^T E(X) Y) = \min_{X \in \mathbb{O}^{n \times \ell}} \text{trace}(X^T E(X) X) = \min_{X \in \mathbb{O}^{n \times \ell}} f(X).$$

$$\begin{aligned} &\geq \text{trace}(\tilde{Z}^T A \tilde{Z}) + 2 \text{trace}(\tilde{Z}^T B) \\ &= f(\tilde{Z}) \geq f(X_{\text{opt}}). \end{aligned}$$

Since $X \in \mathbb{O}^{n \times \ell}$ is arbitrary, the above implies that

$$\min_{X \in \mathbb{O}^{n \times \ell}} \sum_{i=1}^{\ell} \lambda_i(E(X)) \geq f(X_{\text{opt}}).$$

On the other hand, since

$$\min_{X \in \mathbb{O}^{n \times \ell}} \sum_{i=1}^{\ell} \lambda_i(E(X)) \leq \sum_{i=1}^{\ell} \lambda_i(E(X_{\text{opt}})) \leq \text{trace}(X_{\text{opt}}^T E(X_{\text{opt}}) X_{\text{opt}}) = f(X_{\text{opt}}),$$

we consequently have

$$\min_{X \in \mathbb{O}^{n \times \ell}} \sum_{i=1}^{\ell} \lambda_i(E(X)) = f(X_{\text{opt}}).$$

This proves (3.2).

Furthermore, if X_{opt} is a global solution, then by

$$\begin{aligned} \sum_{i=1}^{\ell} \lambda_i(E(X_{\text{opt}})) &= f(X_{\text{opt}}) = \text{trace}(X_{\text{opt}}^T A X_{\text{opt}}) + 2 \text{trace}(X_{\text{opt}}^T B) \\ &= \text{trace}(X_{\text{opt}}^T E(X_{\text{opt}}) X_{\text{opt}}), \end{aligned}$$

it follows that X_{opt} is an orthogonal eigenbasis associated with the first ℓ smallest eigenvalues of $E(X_{\text{opt}})$, and thus (3.3) is true. \square

As a counterpart of (3.3), for a local minimizer X , we can also ensure that X is an orthonormal eigenbasis matrix of $E(X)$.

THEOREM 3.4. *Any local minimizer X is an orthonormal eigenbasis matrix of $E(X)$ satisfying $E(X)X = X\Psi_X$ for a symmetric $\Psi_X \in \mathbb{R}^{\ell \times \ell}$. Moreover, for $\ell \geq 2$, it holds $\lambda_1(\Psi_X) \leq \lambda_{\ell}(E(X))$.*

Proof. According to Lemma 3.2, we have known that any local minimizer X satisfies $E(X)X = X\Psi_X$ and $\text{trace}(\Psi_X) = f(X)$. Thus, all eigenvalues of Ψ_X are eigenvalues of $E(X)$. Let these eigenvalues be

$$\lambda_{\pi_1}(E(X)) \leq \cdots \leq \lambda_{\pi_{\ell}}(E(X)),$$

where $1 \leq \pi_1 \leq \cdots \leq \pi_{\ell} \leq n$, and $\lambda_i(\Psi_X) = \lambda_{\pi_i}(E(X))$ for $i = 1, \dots, \ell$. To show the last part of the theorem, we assume by contradiction that $\lambda_1(\Psi_X) > \lambda_{\ell}(E(X))$. Let $J \in \mathbb{O}^{n \times \ell}$ be an orthonormal eigenbasis matrix of $E(X)$ corresponding to the first ℓ smallest eigenvalues; then

$$\begin{aligned} \text{trace}(J^T E(X) J) &= \sum_{i=1}^{\ell} \lambda_i(E(X)) \\ &< \sum_{i=1}^{\ell} \lambda_i(\Psi_X) = \sum_{i=1}^{\ell} \lambda_{\pi_i}(E(X)) \\ &= \text{trace}(X^T E(X) X) \end{aligned}$$

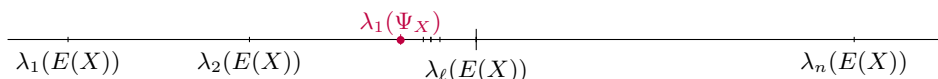
$$= \text{trace}(X^T A X) + 2 \text{trace}(X^T B),$$

and moreover, $J^T X = 0$. Now applying such a J to the second order necessary condition (2.4) and by $\text{trace}(X^T B) \leq 0$, we have

$$\begin{aligned} 0 &\leq \text{trace}(J^T [A + X \Lambda X^T + (2I_n - X X^T) B X^T] J) + \text{trace}(\Lambda) \\ &= \text{trace}(J^T A J) + \text{trace}(\Lambda) \\ &< \text{trace}(X^T A X) + 2 \text{trace}(X^T B) + \text{trace}(\Lambda) \\ &= \text{trace}(X^T B) \leq 0, \end{aligned}$$

a contradiction, and the proof is complete. \square

Comparing with the global minimizer X_{opt} where the associated eigenvalues of $\Psi_{X_{\text{opt}}}$ are just the first ℓ smallest eigenvalues of $E(X_{\text{opt}})$, Theorem 3.4 guarantees that only the smallest eigenvalue $\lambda_1(\Psi_X)$ of Ψ_X must be one of the first ℓ smallest ones of $E(X)$. The following figure illustrates this necessary condition of the location of $\lambda_1(\Psi_X)$ for the local minimizer X .



3.2. An SCF iteration for the associated NEPv. Note that we have characterized the global minimizer X_{opt} of (1.2) as the solution to an eigenvector-dependent NEPv in (3.3). Such a system involves the coefficient matrix $E(X)$ which depends on the pursued eigenvectors X , instead of the eigenvalue in the traditional setting. As has been mentioned in section 1, there are an increasing number of cases and applications of such NEPvs arising from electronic structure calculations (e.g., [25, 27, 33, 43]) and data science (e.g., [4, 24, 40, 44, 45, 48, 46, 47]) recently. For our case, we next propose an SCF iteration for solving the NEPv (3.3). For the k th iteration, the algorithm first finds an arbitrarily orthogonal eigenbasis \hat{X}_k for the current matrix $E(X_{k-1})$, and refines \hat{X}_k to obtain a better orthonormal eigenbasis matrix. The following is a technical lemma for this refinement step.

LEMMA 3.5. *For any given $Z \in \mathbb{O}^{n \times \ell}$, we have*

$$\min_{X \in \mathbb{O}^{n \times \ell}} \text{trace}(Z^T E(X) Z) = f(-ZP),$$

where the solution is $X = -ZP \in \mathbb{O}^{n \times \ell}$ with P being the orthogonal polar factor of $Z^T B$.

Proof. For the given Z , according to Lemma 3.1, we know that

$$\text{trace}(Z^T B X^T Z) \geq - \sum_{i=1}^{\ell} \sigma_i(Z^T B) \sigma_i(X^T Z) \geq - \sum_{i=1}^{\ell} \sigma_i(Z^T B),$$

where the minimum can be achieved when $X = -ZP$ with P being the orthogonal polar factor of $Z^T B$. Thus, we know by (3.4) that

$$\begin{aligned} \min_{X \in \mathbb{O}^{n \times \ell}} \text{trace}(Z^T E(X) Z) &= \text{trace}(Z^T A Z) + 2 \text{trace}(-Z^T B P^T) \\ &= \text{trace}(P^T Z^T A Z P) + 2 \text{trace}((-ZP)^T B) = f(-ZP), \end{aligned}$$

and the proof is complete. \square

Algorithm 3.1 An SCF iteration for the NEPv (3.2).

Given a symmetric $A \in \mathbb{R}^{n \times n}$, a matrix $B \in \mathbb{R}^{n \times \ell}$, and $X_0 \in \mathbb{O}^{n \times \ell}$, the iteration solves (3.2).

- 1: **for** $k = 1, \dots$, until convergence **do**
 - 2: Find an orthonormal eigenbasis matrix $\hat{X}_k \in \mathbb{O}^{n \times \ell}$ associated with the ℓ smallest eigenvalues of $E(X_{k-1})$, and we have
- $$(3.5) \quad E(X_{k-1})\hat{X}_k = \hat{X}_k\hat{\Psi}_k.$$
- 3: Compute orthogonal polar factor P of $\hat{X}_k^T B$ and set $X_k = -\hat{X}_k P$.
 - 4: **end for**
-

We now present the simplest form of our SCF iteration for the NEPv (3.2) in Algorithm 3.1.

For Algorithm 3.1, we remark that our SCF iteration can be understood as an alternating scheme. In fact, for the current X_{k-1} , in line 2 of Algorithm 3.1, we fix the variable $X = X_{k-1}$ inside $E(X)$ to solve

$$\min_{X \in \mathbb{O}^{n \times \ell}} \text{trace}(X^T E(X_{k-1})X)$$

to have \hat{X}_k ; the next update is to fix the $X = \hat{X}_k$ outside and solve

$$(3.6) \quad \min_{X \in \mathbb{O}^{n \times \ell}} \text{trace}(\hat{X}_k^T E(X)\hat{X}_k),$$

whose solution by Lemma 3.5 is $X_k = -\hat{X}_k P$, where P is the orthogonal polar factor of $\hat{X}_k^T B$. This is just the update in line 3 of Algorithm 3.1, and we therefore have

$$(3.7) \quad \min_{X \in \mathbb{O}^{n \times \ell}} \text{trace}(\hat{X}_k^T E(X)\hat{X}_k) = f(X_k).$$

Moreover, it should be noted that X_k is also an orthonormal eigenbasis matrix of $E(X_{k-1})$ satisfying

$$E(X_{k-1})X_k = X_k\Psi_k, \quad \text{where } \Psi_k = P^T\hat{\Psi}_k P.$$

It is this additional step 3 of Algorithm 3.1 that separates our SCF from the usual SCF iteration [7]. The reason is that the objective function $f(X)$ and the coefficient matrix $E(X)$ are not invariant under orthogonal transformation, i.e., $f(X) \neq f(XP)$ for $P \in \mathbb{O}^{\ell \times \ell}$. Thus, when we have an orthonormal eigenbasis matrix \hat{X}_k of $E(X_{k-1})$, it is still possible to improve the objective function by refining \hat{X}_k to a new eigenbasis matrix $X_k = -\hat{X}_k P$. Indeed, the associated P satisfies

$$P \in \operatorname{argmin}_{Q \in \mathbb{O}^{\ell \times \ell}} f(\hat{X}_k Q).$$

It is interesting to note that even if \hat{X}_k in line 2 of Algorithm 3.1 can be arbitrary, the next iteration X_k in line 3 is uniquely determined under certain conditions.

PROPOSITION 3.6. *For Algorithm 3.1, suppose the eigenvalue gap*

$$(3.8) \quad \eta_{k-1} = \lambda_{\ell+1}(E(X_{k-1})) - \lambda_{\ell}(E(X_{k-1})) > 0;$$

then any two orthogonal eigenbases \hat{X}_k and \tilde{X}_k associated with the ℓ smallest eigenvalues of $E(X_{k-1})$ satisfy $\tilde{X}_k = \hat{X}_k Q$ for some orthogonal matrix $Q \in \mathbb{R}^{\ell \times \ell}$. Furthermore, if additionally $\text{rank}(B^T \hat{X}_k) = \ell$, then X_k is uniquely determined.

Proof. If the eigenvalue gap $\eta_{k-1} > 0$, then the eigenspace associated with the ℓ smallest eigenvalues of $E(X_{k-1})$ is unique (see, e.g., [36, p. 244]), and the first part follows. If $\text{rank}(B^T \hat{X}_k) = \ell$, we know that the next iteration in line 3 is just $X_k = -\hat{X}_k P$, where P is the orthogonal polar factor of $\hat{X}_k^T B$ which is $\hat{X}_k^T B (B^T \hat{X}_k \hat{X}_k^T B)^{-\frac{1}{2}}$. Now since for any other choice \tilde{X}_k we have $\tilde{X}_k = \hat{X}_k Q$, it holds that

$$X_k = \hat{X}_k \hat{X}_k^T B (B^T \hat{X}_k \hat{X}_k^T B)^{-\frac{1}{2}} = \tilde{X}_k \tilde{X}_k^T B (B^T \tilde{X}_k \tilde{X}_k^T B)^{-\frac{1}{2}},$$

implying the uniqueness of X_k . \square

4. Convergence analysis. We next discuss some convergence behaviors of the basic SCF iteration (Algorithm 3.1). The following proposition summarizes some basic facts about the sequence $\{X_k\}_{k=0}^\infty$ from Algorithm 3.1.

PROPOSITION 4.1. *Let $\{X_k\}_{k=0}^\infty$ be the sequence generated by Algorithm 3.1; then the following hold:*

- (i) *For each $k \geq 1$, $-B^T X_k = -X_k^T B \succcurlyeq 0$ and $\text{trace}(X_k^T B) = -\sum_{i=1}^\ell \sigma_i(X_k^T B)$.*
- (ii) *The sequence $\{f(X_k)\}_{k=0}^\infty$ is nonincreasing and convergent.*
- (iii) *The sum of the ℓ smallest eigenvalues of $E(X_k)$ is nonincreasing and convergent; that is,*

$$(4.1) \quad \sum_{i=1}^\ell \lambda_i(E(X_{k-1})) \geq \sum_{i=1}^\ell \lambda_i(E(X_k)) \rightarrow \nu \quad \text{as } k \rightarrow \infty.$$

Proof. The term (i) is based on our design of SCF iteration in line 3 of Algorithm 3.1, which leads to $X_k^T B = -V \Sigma V^T \preccurlyeq 0$, where $U \Sigma V^T$ is the SVD of $\hat{X}_k^T B$. Since $X_k = -\hat{X}_k U V^T$, the singular values of $\hat{X}_k^T B^T$ and $X_k^T B$ are the same, and (i) is true.

For (ii), since \hat{X}_{k-1} is the orthogonal eigenbasis of $E(X_{k-1})$ and X_k is the solution to (3.6), it holds that

$$(4.2) \quad \begin{aligned} f(X_{k-1}) &\geq \text{trace}(\hat{X}_k^T E(X_{k-1}) \hat{X}_k) \geq \text{trace}(\hat{X}_k^T E(X_k) \hat{X}_k) \\ &= \text{trace}(X_k^T E(X_k) X_k) = f(X_k), \end{aligned}$$

i.e., $\{f(X_k)\}_{k=0}^\infty$ is nonincreasing. The convergence of $\{f(X_k)\}_{k=0}^\infty$ is then evident by the boundedness over $\mathbb{O}^{n \times \ell}$.

For (iii), from (3.5) and (4.2), we have

$$\sum_{i=1}^\ell \lambda_i(E(X_{k-1})) = \text{trace}(\Psi_k) = \text{trace}(X_k^T E(X_{k-1}) X_k) \geq f(X_k),$$

and also

$$\sum_{i=1}^\ell \lambda_i(E(X_k)) = \text{trace}(X_{k+1}^T E(X_k) X_{k+1}) \leq \text{trace}(X_k^T E(X_k) X_k) = f(X_k).$$

This shows

$$\sum_{i=1}^\ell \lambda_i(E(X_{k-1})) \geq \sum_{i=1}^\ell \lambda_i(E(X_k)),$$

and by the boundedness of $\text{trace}(X^T E(X) X)$ over $\mathbb{O}^{n \times \ell}$, the claim (iii) is then established. \square

For the convergence of Algorithm 3.1, we next introduce the metric between two subspaces. Given two $X, Y \in \mathbb{O}^{n \times \ell}$, we define the distance between two subspaces $\mathcal{R}(X)$ and $\mathcal{R}(Y)$ of dimension ℓ by

$$(4.3) \quad \|\sin \Theta(\mathcal{R}(X), \mathcal{R}(Y))\|_2 = \|\mathcal{P}_X - \mathcal{P}_Y\|_2,$$

where \mathcal{P}_X denotes the orthogonal projection onto $\mathcal{R}(X)$, and

$$\Theta(\mathcal{R}(X), \mathcal{R}(Y)) = \text{diag}(\theta_1(\mathcal{R}(X), \mathcal{R}(Y)), \dots, \theta_\ell(\mathcal{R}(X), \mathcal{R}(Y)))$$

stand for the canonical angles between $\mathcal{R}(X)$ and $\mathcal{R}(Y)$, which according to the singular values $\sigma_i(X^T Y)$ are defined by

$$0 \leq \theta_i(\mathcal{R}(X), \mathcal{R}(Y)) := \arccos \sigma_i(X^T Y) \leq \frac{\pi}{2} \quad \text{for } 1 \leq i \leq \ell.$$

We remark that $\|\sin \Theta(\mathcal{R}(X), \mathcal{R}(Y))\|_2$ in (4.3) is a particular unitarily invariant metric [32] in the Grassmann manifold $\text{Grass}(n, \ell)$, in which each point $[X] \in \text{Grass}(n, \ell)$ can be represented by any basis X of $\mathcal{R}(X)$. In particular, if $X, Y \in \mathbb{O}^{n \times \ell}$, we have $\mathcal{P}_X = XX^T$ and $\mathcal{P}_Y = YY^T$, and thus $\|\sin \Theta(\mathcal{R}(X), \mathcal{R}(Y))\|_2 = \|XX^T - YY^T\|_2$. For simplicity purposes, we denote

$$(4.4) \quad \|[X] - [Y]\|_2 := \|\sin \Theta(\mathcal{R}(X), \mathcal{R}(Y))\|_2.$$

The following lemma is a modification of the one [29, Lemma 4.10] (see also [21, Proposition 7]) in which the standard Euclidean distance is replaced by (4.4).

LEMMA 4.2. *Assume that $[W] \in \text{Grass}(n, \ell)$ is an isolated accumulation point (in the metric (4.4)) of a sequence $\{[W_k]\}_{k=0}^\infty \subseteq \text{Grass}(n, \ell)$ such that, for every subsequence $\{[W_k]\}_{k \in \mathcal{K}}$ converging to $[W]$, there is an infinite subset $\hat{\mathcal{K}} \subseteq \mathcal{K}$ such that $\{\|[W_{k+1}] - [W_k]\|_2\}_{k \in \hat{\mathcal{K}}} \rightarrow 0$. Then the whole sequence $\{[W_k]\}_{k=0}^\infty$ converges to $[W]$.*

We remark that Lemma 4.2 is a sufficient condition for proving the convergence of $\{[W_k]\}_{k=0}^\infty$ to $[W]$. The condition that $[W] \in \text{Grass}(n, \ell)$ is an isolated accumulation point that does not exclude the existence of finitely many or even infinitely many accumulation points. Indeed, when we know that $\{[W_k]\}_{k=0}^\infty$ only contains finitely many accumulation points, then each is isolated, and the condition of Lemma 4.2 ensures that $[W_k] \rightarrow [W]$. Based on this lemma, we have the following convergence results for Algorithm 3.1.

THEOREM 4.3. *Let $\{X_k\}_{k=0}^\infty$ be the sequence generated by Algorithm 3.1. Let $\{X_k\}_{k \in \mathcal{K}}$ be any convergent subsequence with the accumulation point X satisfying*

$$(4.5) \quad \eta = \lambda_{\ell+1}(E(X)) - \lambda_\ell(E(X)) > 0.$$

Then

- (i) *X satisfies the first order optimality in Lemma 2.1 and also the necessary condition (2.6) for a global minimizer;*
- (ii) *if, additionally, $[X]$ is an isolated accumulation point of $\{[X_k]\}_{k=0}^\infty$, then $\{[X_k]\}_{k=0}^\infty$ converges to $[X]$;*
- (iii) *under the conditions of (ii), and if $\text{rank}(B^T X) = \ell$, then $\{X_k\}_{k=0}^\infty$ converges to X .*

Proof. For (i), we can choose a subsequence $\{X_k\}_{k \in \widehat{\mathcal{K}}}$ of $\{X_k\}_{k \in \mathcal{K}}$ such that $\{X_{k-1}\}_{k \in \widehat{\mathcal{K}}}$ and $\{X_{k+1}\}_{k \in \widehat{\mathcal{K}}}$ converge to $Y \in \mathbb{O}^{n \times \ell}$ and $Z \in \mathbb{O}^{n \times \ell}$, respectively. Thus, from

$$E(X_{k-1})X_k = X_k\Psi_k \quad \text{and} \quad E(X_k)X_{k+1} = X_{k+1}\Psi_{k+1}, \quad k \in \widehat{\mathcal{K}},$$

we have

$$(4.6) \quad E(Y)X = X\Psi \quad \text{and} \quad E(X)Z = Z\tilde{\Psi}.$$

Also, by (i) of Proposition 4.1, we know that $-H^T B = -B^T H \succcurlyeq 0$ for $H \in \{X, Y, Z\}$, and

$$\lim_{k \rightarrow \infty, k \in \widehat{\mathcal{K}}} \sum_{i=1}^{\ell} \lambda_i(E(X_{k-1})) = \text{trace}(\Psi) = \lim_{k \rightarrow \infty, k \in \widehat{\mathcal{K}}} \sum_{i=1}^{\ell} \lambda_i(E(X_k)) = \text{trace}(\tilde{\Psi}) = \nu.$$

Note that $0 \leq \sigma_i(X_{k-1}^T X_k) \leq 1$ and

$$(4.7) \quad \text{trace}(X_k^T B X_{k-1}^T X_k) - \text{trace}(X_k^T B) \geq \sum_{i=1}^{\ell} \sigma_i(X_k^T B) (1 - \sigma_i(X_{k-1}^T X_k)) \geq 0.$$

Applying $k \rightarrow \infty, k \in \widehat{\mathcal{K}}$ to (4.7), it holds that $\text{trace}(X^T B Y^T X) \geq \text{trace}(X^T B)$. Thus,

$$(4.8) \quad \text{trace}(X^T E(Y)X) - \text{trace}(X^T E(X)X) = 2(\text{trace}(X^T B Y^T X) - \text{trace}(X^T B)) \geq 0.$$

Moreover, using the fact that Z is an orthogonal eigenbasis of $E(X)$ associated with the first ℓ smallest eigenvalues and (4.6), we conclude that

$$(4.9) \quad \begin{aligned} \text{trace}(X^T E(Y)X) - \text{trace}(X^T E(X)X) &\leq \text{trace}(X^T E(Y)X) - \text{trace}(Z^T E(X)Z) \\ &= \text{trace}(\Psi) - \text{trace}(\tilde{\Psi}) = 0, \end{aligned}$$

which, combined with (4.8), yields $\text{trace}(X^T B Y^T X) = \text{trace}(X^T B)$. Therefore, using (4.6) and

$$\text{trace}(X^T B Y^T X) = \text{trace}(X^T B),$$

we get

$$\begin{aligned} \nu = \text{trace}(\tilde{\Psi}) &\leq \text{trace}(X^T E(X)X) \\ &= \text{trace}(X^T A X) + 2 \text{trace}(X^T B) \\ &= \text{trace}(X^T A X) + 2 \text{trace}(X^T B Y^T X) \\ &= \text{trace}(X^T E(Y)X) = \text{trace}(\Psi) = \nu. \end{aligned}$$

From $\text{trace}(\tilde{\Psi}) = \text{trace}(\Psi)$ and $E(X)Z = Z\tilde{\Psi}$, we know that X is also an orthogonal eigenbasis associated with the first ℓ smallest eigenvalues of $E(X)$. The eigenvalue gap condition (4.5) ensures that $\lfloor X \rfloor = \lfloor Z \rfloor$ (see, e.g., [36, p. 244]). In other words, $Z = XQ$ for an orthonormal $Q \in \mathbb{R}^{\ell \times \ell}$. Thus, the condition $E(X)Z = E(X)XQ = XQ\tilde{\Psi}$ yields $E(X)X = E(X)X = XQ\tilde{\Psi}Q^T$, or equivalently,

$$AX + X\Lambda = -B, \quad \text{with} \quad \Lambda = B^T X - Q\tilde{\Psi}Q^T.$$

This proves (i).

For (ii), as we have shown, there is a subsequence $\widehat{\mathcal{K}} \subseteq \mathcal{K}$ such that $\{\| \lfloor X_{k+1} \rfloor - \lfloor X_k \rfloor \|_2\}_{k \in \widehat{\mathcal{K}}} \rightarrow 0$ as $k \rightarrow \infty$. According to Lemma 4.2, the condition that $\lfloor X \rfloor$ is an isolated accumulation point of $\{\lfloor X_k \rfloor\}_{k=0}^\infty$ leads to (ii).

For (iii), according to (ii), the limit of any convergent subsequence of $\{X_k\}_{k=0}^\infty$ can be expressed as XQ for some $Q \in \mathbb{O}^{\ell \times \ell}$. Moreover, Proposition 4.1(ii) implies that $f(X) = f(XQ)$, leading to $\text{trace}(B^T X Q) = \text{trace}(B^T X)$. Note also that Proposition 4.1(i) implies that $B^T X$ is symmetric and $\text{trace}(B^T X) = -\sum_{i=1}^\ell \sigma_i(B^T X)$. Let the eigenvalue decomposition of $B^T X$ be $-V\Sigma V^T$ with $\Sigma \succcurlyeq 0$ and $\widehat{Q} = V^T Q V \in \mathbb{O}^{\ell \times \ell}$. Thus, we have

$$-\sum_{i=1}^\ell \sigma_i(B^T X) = \text{trace}(B^T X) = \text{trace}(B^T X Q) = \text{trace}(-\Sigma \widehat{Q}) = -\sum_{i=1}^\ell \sigma_i(B^T X) \widehat{Q}_{(i,i)}$$

implying $\sigma_i(B^T X)(1 - \widehat{Q}_{(i,i)}) = 0 \ \forall i = 1, 2, \dots, \ell$, where $\widehat{Q}_{(i,i)}$ are the diagonal entries of \widehat{Q} . Thus, if $\text{rank}(B^T X) = \ell$, we have $\widehat{Q}_{(i,i)} = 1 \ \forall i = 1, \dots, \ell$, which implies $\widehat{Q} = Q = I_\ell$. This proves (iii). \square

We next discuss the local convergence behavior for the SCF iteration. As Theorem 4.3(iii) implies that the sequence $\{X_k\}$ converges to a solution X , we will establish an estimate of the local convergence rate, which reflects the factors for the speed of the SCF iteration. Our convergence analysis for this part basically follows from the general local convergence analysis [7] for the SCF iteration in solving the standard eigenvector-dependent NEPv. The following lemma [45, Lemma 3.5] is obvious.

LEMMA 4.4. *Let $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ be of unit norm and $\|\mathbf{x} - \mathbf{y}\|_2 < 1$. Then*

$$(4.10) \quad \frac{\sqrt{2}}{2} \|\mathbf{x} - \mathbf{y}\|_2 \leq \sin \theta(\mathbf{x}, \mathbf{y}) = \|\lfloor \mathbf{x} \rfloor - \lfloor \mathbf{y} \rfloor\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_2,$$

where $\sin \theta(\mathbf{x}, \mathbf{y})$ is the sine of angle between \mathbf{x} and \mathbf{y} .

THEOREM 4.5. *Under the assumptions of Theorem 4.3(iii), then $\{X_k\}$ converges to X , and moreover, for sufficiently large k , we have*

$$(4.11) \quad \|\lfloor X \rfloor - \lfloor X_{k+1} \rfloor\|_2 \leq \frac{3\|B\|_2}{\eta} \|X_k - X\|_2,$$

where η is given by (4.5). If $\ell = 1$, then

$$(4.12) \quad \|\lfloor X \rfloor - \lfloor X_{k+1} \rfloor\|_2 \leq \frac{3\sqrt{2}\|B\|_2}{\eta} \|\lfloor X \rfloor - \lfloor X_k \rfloor\|_2,$$

$$(4.13) \quad \|X - X_{k+1}\|_2 \leq \frac{3\sqrt{2}\|B\|_2}{\eta} \|X - X_k\|_2.$$

Proof. For k , we denote $Z_k = X_k - X$ and $\Delta_k = BZ_k^T + Z_k B^T$, and thus, $E(X_k) = E(X) + \Delta_k$. Also, let

$$(4.14) \quad X_{k+1} = XC + X_\perp S, \quad C \in \mathbb{R}^{\ell \times \ell}, \quad S \in \mathbb{R}^{(n-\ell) \times (n-\ell)}.$$

It is then known that for any unitarily invariant norm $\|\cdot\|_{\text{ui}}$, it follows that

$$(4.15) \quad \|\sin \Theta(\mathcal{R}(X_{k+1}), \mathcal{R}(X))\|_{\text{ui}} = \|S\|_{\text{ui}}.$$

By premultiplying X_{\perp}^T on both sides of $E(X_k)X_{k+1} = (E(X) + \Delta_k)X_{k+1} = X_{k+1}\Psi_{k+1}$, we have

$$(4.16) \quad X_{\perp}^T E(X)X_{k+1} - X_{\perp}^T X_{k+1}\Psi_{k+1} = \Psi_{\perp}S - S\Psi_{k+1} = -X_{\perp}^T \Delta_k X_{k+1},$$

where $\Psi_{\perp} = X_{\perp}^T E(X)X_{\perp} \in \mathbb{R}^{(n-\ell) \times (n-\ell)}$ whose eigenvalues are $\text{eig}(E(X))$ with $\text{eig}(\Psi) = \text{eig}(X^T E(X)X)$ excluded. Note that the solution S in the Sylvester equation (4.16) satisfies (see [9])

$$(4.17) \quad \|S\|_{\text{ui}} = \|\sin \Theta(\mathcal{R}(X_{k+1}), \mathcal{R}(X))\|_{\text{ui}} \leq \frac{\|X_{\perp}^T \Delta_k X_{k+1}\|_{\text{ui}}}{\eta_k} \leq \frac{3\|X_{\perp}^T \Delta_k X_{k+1}\|_{\text{ui}}}{2\eta},$$

where $\eta_k = \lambda_{\ell+1}(E(X_k)) - \lambda_{\ell}(E(X_k))$ which can be assumed to be positive and $\eta_k \geq 2\eta/3$ for sufficiently large k . Therefore, (4.11) follows by $\|X_{\perp}^T \Delta_k X_{k+1}\|_2 \leq 2\|B\|_2\|Z_k\|_2$, and (4.12)–(4.13) follow by (4.10). \square

Remark 4.6. As an implication of Theorem 4.5, we know that the convergence of SCF is fast whenever $\|B\|_2$ is small and the eigenvalue gap η in (4.5) is large. Indeed, the result (4.11) reflects the nature of the problem (1.2): When B is small in the norm (relative to that of A), then (1.2) is much like an eigenvalue problem of finding the eigenpairs associated with the ℓ smallest eigenvalues of A , while when $\|B\|$ gets bigger (relative to that of A), the problem (1.2) gets closer to maximizing $\text{trace}(B^T X)$ over $\mathbb{O}^{n \times \ell}$, whose solution is just the orthogonal polar factor of B . This useful observation will be used in our implementation for the eigensolver in line 2 of Algorithm 3.1; furthermore, we will illustrate the convergence behavior with a different norm $\|B\|_2$ in our numerical testing.

5. Efficient implementation of the SCF iteration. In this section, we provide several techniques to make the simple SCF iteration more efficient.

5.1. The inexact SCF iteration. Note that the main computational effort of the pure SCF iteration lies in computing an orthonormal eigenbasis matrix of $E(X_{k-1})$ associated with the ℓ smallest eigenvalues (line 2 in Algorithm 3.1); we will discuss efficient iterative eigensolvers for this task, instead of resorting the full eigendecomposition of $E(X_{k-1})$. The following proposition first justifies the monotonic convergence of the SCF iteration even with an approximate orthonormal eigenbasis matrix \hat{X}_k of $E(X_{k-1})$ in line 2 in Algorithm 3.1.

PROPOSITION 5.1. *For the $(k-1)$ th iterate $X_{k-1} \in \mathbb{O}^{n \times \ell}$ in Algorithm 3.1, let $\hat{X}_k \in \mathbb{O}^{n \times \ell}$ be an approximate orthonormal eigenbasis matrix associated with the first ℓ smallest eigenvalues of $E(X_{k-1})$ satisfying*

$$(5.1) \quad \text{trace}(\hat{X}_k^T E(X_{k-1})\hat{X}_k) \leq \text{trace}(X_{k-1}^T E(X_{k-1})X_{k-1}) = f(X_{k-1});$$

then the next iterate X_k of the SCF iteration Algorithm 3.1 satisfies $f(X_k) \leq f(X_{k-1})$.

Proof. Recall from Algorithm 3.1 that $X_k = -\hat{X}_k P$, where P is the orthogonal polar factor of $\hat{X}_k^T B$. By Lemma 3.5, X_k minimizes $\text{trace}(\hat{X}_k^T E(X)\hat{X}_k)$ over all $X \in \mathbb{O}^{n \times \ell}$. Thus,

$$(5.2) \quad \begin{aligned} f(X_k) &= \text{trace}(X_k^T E(X_k)X_k) = \text{trace}(\hat{X}_k^T E(X_k)\hat{X}_k) \\ &\leq \text{trace}(\hat{X}_k^T E(X_{k-1})\hat{X}_k) \leq f(X_{k-1}), \end{aligned}$$

and the proof is complete. \square

We remark, by (5.2), that as long as the inequality (5.1) is strict, we have $f(X_k) < f(X_{k-1})$. If X_{k-1} is not the solution of the NEPv (3.3), then (5.1) holds strictly as long as the approximate \hat{X}_k is better than X_{k-1} in the sense of Rayleigh quotient minimization. The latter could be generically true for any efficient eigensolver. Our choice here is the locally optimal block preconditioned conjugate gradient (LOBPCG) method² (see [22, 23]). It is known that the $(j+1)$ th iterate $\hat{X}_k^{(j+1)}$ of LOBPCG for an approximate eigenbasis corresponding to the first ℓ smallest eigenvalues of $E(X_{k-1})$ is obtained by solving

$$(5.3) \quad \hat{X}_k^{(j+1)} := \arg \min_{X \in \mathcal{R}([\hat{X}_k^{(j-1)}, \hat{X}_k^{(j)}, NR^{(j)}])} \text{trace}(X^T E(X_{k-1}) X),$$

where $N \in \mathbb{R}^{n \times n}$ is a symmetric preconditioner for the matrix $E(X_{k-1})$, and $R^{(j)} = E(X_{k-1})\hat{X}_k^{(j)} - \hat{X}_k^{(j)}\Theta_j$ is the residual associated with the Ritz eigenpair $(\hat{X}_k^{(j)}, \Theta_j)$. Thus, if we choose the initial $\hat{X}_k^{(0)} = X_{k-1}$, by (5.3), for the approximate $\hat{X}_k = \hat{X}_k^{(j+1)}$ we have

$$\begin{aligned} \text{trace}(\hat{X}_k^T E(X_{k-1}) \hat{X}_k) &= \text{trace}((\hat{X}_k^{(j+1)})^T E(X_{k-1}) \hat{X}_k^{(j+1)}) \\ &\leq \text{trace}((\hat{X}_k^{(0)})^T E(X_{k-1}) \hat{X}_k^{(0)}) \\ &= \text{trace}(X_{k-1}^T E(X_{k-1}) X_{k-1}) = f(X_{k-1}), \end{aligned}$$

validating the condition (5.1).

The termination of the LOBPCG for the eigenpairs of $E(X_{k-1})$ is controlled by the residual $\|R^{(j)}\| \leq \text{tol_lobpcg}$ and the maximal number `maxit_lobpcg` of iterations. By Remark 4.6, we know that when $\|B\|$ is small relative to $\|A\|$ (implying that the problem is close to an eigenvalue problem associated with A), we attempt to choose a relatively tight stopping criterion or use a preconditioner for LOBPCG. Detailed parameters of `tol_lobpcg` and `maxit_lobpcg` will be specified in section 6.

5.2. Inexact SCF iteration with subspace acceleration. In the traditional eigenvalue problem computation, the Rayleigh–Ritz (RR) procedure [11, 16] is one of most widely used and efficient approaches. It serves as a basic framework for most state-of-the-art eigensolvers for large scale and sparse eigenvalue problems. In general, the RR procedure seeks a good subspace (usually the Krylov subspace) containing sufficient information of desired eigenpairs, onto which the original eigenproblem is projected to compute the Ritz pairs as approximate eigenpairs of the original problem.

For the NEPv (3.3), we observe that the iterates $\{X_k\}$ gradually converge to X_{opt} and naturally collect increasingly useful information of X_{opt} . Suppose \mathcal{S} is a subspace which can be formed by the previous t iterates X_{k-t}, \dots, X_{k-1} at the current $(k-1)$ th iteration. In the framework of RR, we attempt to refine the current X_{k-1} by seeking a better approximate \check{X}_{k-1} than X_{k-1} in the next SCF iteration: $E(\check{X}_{k-1})X_k = X_k\Lambda_k$. In particular, we let Q_k be the orthogonal basis of \mathcal{S} , and we choose to find $\check{X}_{k-1} \in \mathcal{R}(Q_k)$. Assume $\text{rank}(\mathcal{S}) = r_k$, i.e., $Q_k \in \mathbb{O}^{n \times r_k}$, and thus, $\check{X}_{k-1} = Q_k Y_k$ for some $Y_k \in \mathbb{O}^{r_k \times \ell}$. In the spirit of RR, the parameter matrix Y_k can be naturally obtained via solving

$$(5.4) \quad \min_{\check{X} \in \mathcal{R}(Q_k)} f(\check{X}) \iff \min_{Y \in \mathbb{O}^{r_k \times \ell}} \text{trace}(Y^T A_k Y) + 2 \text{trace}(B_k^T Y),$$

²LOBPCG (in MATLAB) is available at <http://cn.mathworks.com/matlabcentral/fileexchange/48-lobpcg-m>.

where $A_k = Q_k^T A Q_k \in \mathbb{R}^{r_k \times r_k}$ and $B_k = Q_k^T B \in \mathbb{R}^{r_k \times \ell}$. It can be seen by Theorem 3.3 that (5.4) is equivalent to

$$\min_{Y \in \mathbb{O}^{r_k \times \ell}} \sum_{i=1}^{\ell} \lambda_i(A_k + B_k Y^T + Y B_k^T),$$

and the solution Y_k is an orthogonal eigenbasis of $A_k + B_k Y^T + Y B_k^T$ associated with the first ℓ smallest eigenvalues.

Notice that (5.4) is again a minimization of the quadratic function over the Stiefel manifold, but the dimension r_k is generally small. The Riemannian Trust-Region (RTR) method³ in [1, 2] is suitable for such a task, where we mainly need to specify the gradient and Hessian operator of the objective function $f(X)$ restricted on $\mathbb{O}^{n \times \ell}$ for its implementation. By calculation, we have the following proposition.

PROPOSITION 5.2. *Let $f|_{\mathbb{O}^{n \times \ell}(X)}$ be the restriction of the objective function $f(X)$ of (1.2) on the Stiefel manifold $\mathbb{O}^{n \times \ell}$. Then for any $X \in \mathbb{O}^{n \times \ell}$, the gradient $g(X)$ of $f|_{\mathbb{O}^{n \times \ell}(X)}$ is $g(X) = AX + X\Lambda + B$, where $\Lambda = -X^T A X - \frac{X^T B + B^T X}{2}$. Moreover, the Hessian operation $h(X) : \mathcal{T}_X \mathbb{O}^{n \times \ell} \rightarrow \mathcal{T}_X \mathbb{O}^{n \times \ell}$ of $f|_{\mathbb{O}^{n \times \ell}(X)}$ at $X \in \mathbb{O}^{n \times \ell}$ acting on an arbitrary tangent vector $H \in \mathcal{T}_X \mathbb{O}^{n \times \ell}$ is given by*

$$h(X)[H] = \Pi_X(AH + H\Lambda),$$

where $\Pi_X(\bullet)$ is the projection defined by (2.3).

Similar to the monotonic property in Proposition 5.1, we next show that the solution \check{X}_{k-1} to (5.4) obtained from RTR can be computed inexactly.

PROPOSITION 5.3. *For $t \geq 1$, suppose the approximate solution \check{X}_{k-1} to (5.4) obtained from RTR satisfies $f(\check{X}_{k-1}) \leq f(X_{k-1})$ and $\hat{X}_k \in \mathbb{O}^{n \times \ell}$ is an approximate orthonormal eigenbasis matrix associated with the ℓ smallest eigenvalues of $E(\check{X}_{k-1})$ satisfying*

$$(5.5) \quad \text{trace}(\hat{X}_k^T E(\check{X}_{k-1}) \hat{X}_k) \leq \text{trace}(\check{X}_{k-1}^T E(\check{X}_{k-1}) \check{X}_{k-1}) = f(\check{X}_{k-1});$$

then the next iterate $X_k = -\hat{X}_k P$ satisfies $f(X_k) \leq f(X_{k-1})$, where P is the orthogonal polar factor of $\hat{X}_k^T B$.

Proof. We note by assumptions and Lemma 3.5 that

$$\begin{aligned} f(X_{k-1}) &\geq f(\check{X}_{k-1}) \geq \text{trace}(\hat{X}_k^T E(\check{X}_{k-1}) \hat{X}_k) \geq \min_{X \in \mathbb{O}^{n \times \ell}} \text{trace}(\hat{X}_k^T E(X) \hat{X}_k) \\ &= f(-\hat{X}_k P) = f(X_k), \end{aligned}$$

and the assertion follows. \square

To efficiently apply this subspace acceleration strategy, in our numerical testing, we restart the subspace acceleration by introducing the maximal number of blocks, say ς_m , of the subspace $\mathcal{S} = \mathcal{R}([X_{k-t}, \dots, X_{k-1}])$; that is, we require $t \leq \varsigma_m$, and whenever $t > \varsigma_m$, we discard the first block by setting $\mathcal{S}(:, 1 : \ell) = [\]$ and then restart expanding \mathcal{S} (in our numerical testing, we set $\varsigma_m = 5$). Incorporating such subspace acceleration strategy, we arrive at the following practical SCF iteration, named SCFRTR (Algorithm 5.1).

³The MATLAB solver of the generic RTR method is available at <https://www.manopt.org>.

Algorithm 5.1 SCFRTR: The SCF iteration with subspace acceleration for (1.2).

Given a symmetric $A \in \mathbb{R}^{n \times n}$, $B \in \mathbb{R}^{n \times \ell}$, and $X_0 \in \mathbb{O}^{n \times \ell}$, SCFRTR returns an approximate solution of (1.2).

```

1: Set  $\mathcal{S} = \emptyset$ ,  $t = 0$ ;
2: for  $k = 1, \dots$ , until convergence do
3:   if  $t \geq 1$  then
4:     Find an orthonormal basis matrix  $Q_k \in \mathbb{R}^{n \times r_k}$  of  $\mathcal{S}$ .
5:     Solve (5.4) by RTR for  $Y_k$  and set  $X_{k-1} = Q_k Y_k$ .
6:   end if
7:   Call LOBPCG (or other eigensolver) to find an approximate orthonormal eigen-
     basis matrix  $\hat{X}_k \in \mathbb{R}^{n \times \ell}$  associated with the first  $\ell$  smallest eigenvalues of
      $E(X_{k-1})$ .
8:   Compute the orthogonal polar factor  $P$  of  $\hat{X}_k^T B$ , and set  $X_k = -\hat{X}_k P$ .
9:   if  $t = \varsigma_m$  then
10:     $\mathcal{S}(:, 1 : \ell) = [ ]$  and  $t = t - 1$ 
11:   else
12:     $\mathcal{S} = [\mathcal{S}, X_k]$  and  $t = t + 1$ 
13:   end if
14: end for

```

6. Numerical experiments. In this section, we will conduct numerical testing to show basic behavior of our method; also we will evaluate the performance on a preliminary test problems and report the numerical experiments together with those from other generic optimization methods designed upon the Stiefel manifold. All of these experiments were conducted on a PC using Windows 7 (64bit) system with Intel Core i5-3230M CPU (2.6GHz) and 4GB memory on MATLAB (2010b).

There are two goals in our numerical report: (i) to demonstrate the improvements of the inexact eigenvalue computation and our subspace acceleration strategy used in Algorithm 5.1 over the pure SCF iteration Algorithm 3.1, and (ii) to show the efficiency of the algorithm SCFRTR by comparing it with four other solvers (the RTR method of Absil, Baker, and Gallivan [1], the WYBB algorithm by Wen and Yin [42], the JDCP algorithm of Jiang and Dai [20], and the GPI iteration of Nie, Zhang, and Li [30]). Following the stopping criterion used in [42, 20], for our method SCFRTR we terminate the iteration in step 2 of Algorithm 5.1 when $k > 60$ or if any of the following three conditions is satisfied:

$$(6.1) \quad \frac{f(X_k) - f(X_{k+1})}{|f(X_k)| + 1} \leq \text{tol}_f, \quad \frac{\|X_k - X_{k+1}\|_F}{\sqrt{n}} \leq \text{tol}_x, \quad \|AX_k + B + X_k \Lambda_k\|_F \leq \text{tol}_g.$$

Note that the last condition of (6.1) is the absolute norm of the gradient at the computed point X_k , which might be affected by the scaling of the data matrices A and B . Therefore, in order to eliminate the impact from the scaling and in order to have a relative norm of the gradient at the computed point, for a given problem with data matrices A and B in our testing, we conduct a procedure

$$A \leftarrow \frac{A}{\|A\|_1 + \|B\|_1}, \quad B \leftarrow \frac{B}{\|A\|_1 + \|B\|_1}.$$

We remark that, without modifying the stopping criteria used in [42, 20], the absolute norm of the gradient (6.1) on the scaled problem now serves as an approximation of

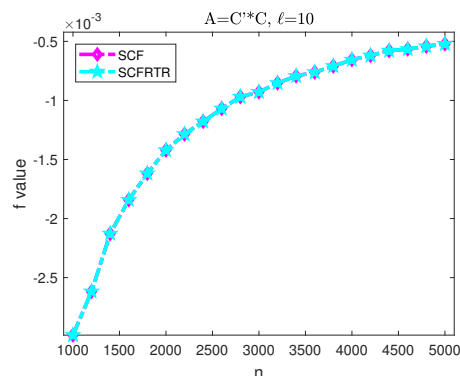


FIG. 1. Optimal values.

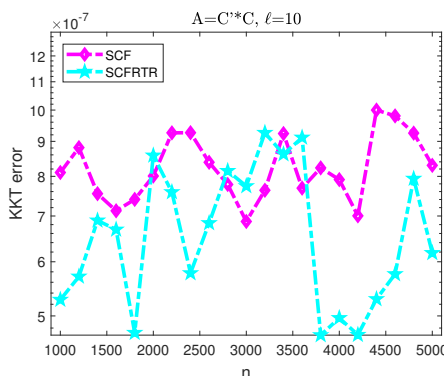


FIG. 2. KKT error.

the relative norm of the gradient at X_k for the unscaled one. In our testing, we set $\text{tol}_f = 10^{-12}$ and $\text{tol}_x = \text{tol}_g = 10^{-6}$; moreover, unless otherwise specified, we choose X_0 as the orthonormal basis of $\mathcal{R}(B)$ by the modified Gram–Schmidt process if $\text{rank}(B) = \ell$ and $X_0 = [I_\ell, 0]^T$ otherwise, for all methods in the tests of sections 6.1, 6.2, and 6.3.

6.1. The pure SCF versus the practical SCF. In this subsection, we will demonstrate the behavior of our practical SCFRTR (Algorithm 5.1) by comparing that of the pure SCF algorithm (Algorithm 3.1). For this purpose, we choose the data matrices A and B generated by MATLAB command `randn`: $C = \text{randn}(n)$; $A = C' * C$; $B = \text{randn}(n, \ell)$ for (1.2). With an increment 200, we choose various n in the interval $[1000, 5000]$, and compare the results of optimal values, KKT errors $\|AX_k + B + X_k\Lambda_k\|_F$, the number of iterations, and running CPU time (in seconds) from SCFRTR (Algorithm 5.1) and SCF (Algorithm 3.1), in Figures 1–4, respectively. In our implementation of the baseline SCF in Algorithm 3.1, we use the MATLAB function `eigs` to compute an orthonormal eigenbasis matrix \tilde{X}_k in step 2. For SCFRTR, we call LOBPCG by setting the inner tolerance $\text{tol_lobpcg} = 10^{-7}$ and activate the Cholesky factor of A (or the incomplete Cholesky for sparse case) as the preconditioner whenever $\|A\|_1 / \|B\|_1 \geq 100$ and the estimated condition number of A is larger than 10^6 ; moreover, we set the maximal inner iterations $\text{maxit_lobpcg} = 200$ in the first three outer iterations when no preconditioner is used, and set $\text{maxit_lobpcg} \leq 20$ for others.

Figures 1 and 2 show that, in terms of the optimal values and KKT errors, SCFRTR and SCF nearly are of the same performance. However, for the efficiency, we observed from Figures 3 and 4 that SCFRTR saves about 2/3 fraction of the number of outer-loop iterations, and reduces the running times in CPU significantly from those used in the pure SCF iteration. This demonstrates the improvement of our practical SCF iteration over the pure SCF.

6.2. Comparison with other solvers. We now compare our algorithm with the RTR algorithm of Absil, Baker, and Gallivan [1], the WYBB algorithm of Wen and Yin [42], the JDCP algorithm of Jiang and Dai [20], and the GPI iteration of Nie, Zhang, and Li [30]. Here, RTR [1] is a generic trust-region method for minimization over a general Riemannian manifold (we terminate RTR if the norm of Riemannian gradient is no greater than 10^{-4}), while WYBB and JDCP are two generic optimization methods built on the Stiefel manifold. In particular, the WYBB algorithm

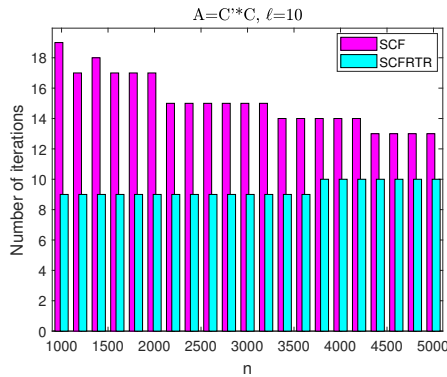


FIG. 3. Number of iterations.

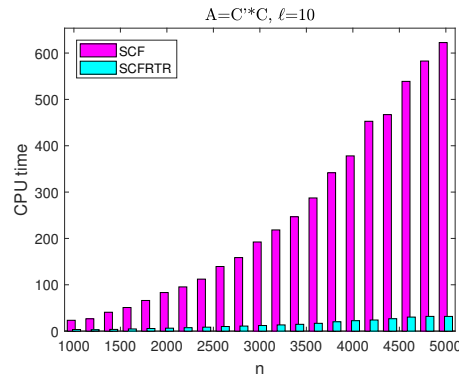


FIG. 4. CPU time.

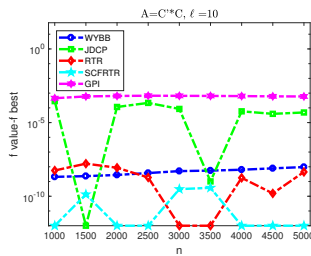
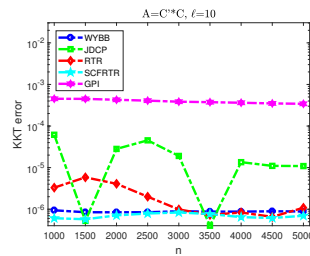
FIG. 5. Difference of f .

FIG. 6. KKT error.

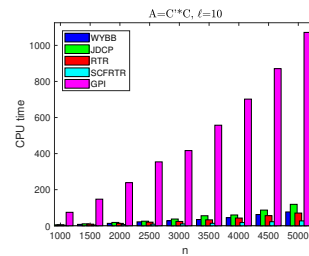


FIG. 7. CPU time.

is a feasible gradient method which is based on the Cayley transform and Barzilai–Borwein stepsizes, and JDCP is a uniform framework of feasible gradient methods which decompose each feasible point X into the range space of X and the null space of X^T . GPI is the generalized power iteration especially designed for the unbalanced Procrustes problem (1.1), and the same stopping criteria (6.1) as SCFRTR are set for GPI in our testing.

6.2.1. Test on synthetic problems. We first report the numerical results of (1.2) obtained by solving the randomly generated problems in section 6.1 but with $1000 \leq n \leq 5000$. In Figures 5–7 we plot the averaged quantities of the difference of optimal values, KKT errors, and CPU time over ten random tests for each n . The notion of “f best” in Figure 5 stands for the smallest computed optimal values among all solvers. These preliminary numerical results demonstrate the efficiency of SCFRTR for our generated minimization problems (1.2).

We next extend our numerical evaluation. As has been pointed out in Remark 4.6, the factor of $\frac{\|B\|^2}{\eta}$ has an impact on the convergence of the SCF iteration. Note that the eigenvalue gap $\eta = \lambda_{\ell+1}(E(X_{\text{opt}})) - \lambda_{\ell}(E(X_{\text{opt}}))$ cannot be estimated a priori; in order to evaluate behavior of the iteration with respect to that factor, we therefore only choose to vary the relative norm $\|A\|/\|B\|$ between A and B . In particular, we test both the minimization (1.2) and the Procrustes problem (1.1) by selecting various parameters $a > 0$ in

$$(6.2) \quad \min_{X \in \mathbb{O}^{n \times \ell}} \text{trace}(X^T A_a X) + 2 \text{trace}(X^T B), \quad A_a = aA,$$

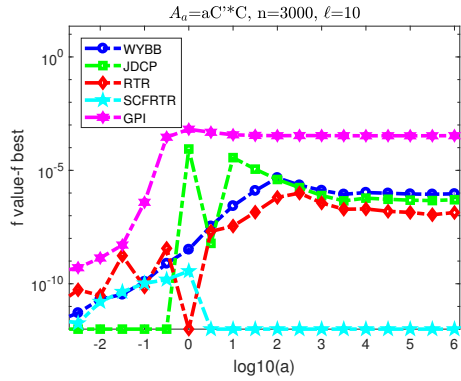


FIG. 8. Difference of optimal values for (6.2).

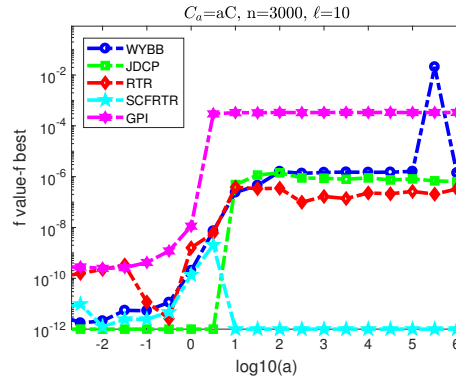


FIG. 9. Difference of optimal values for (6.3).

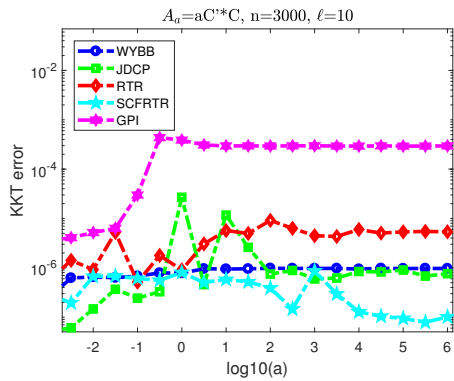


FIG. 10. KKT error for (6.2).

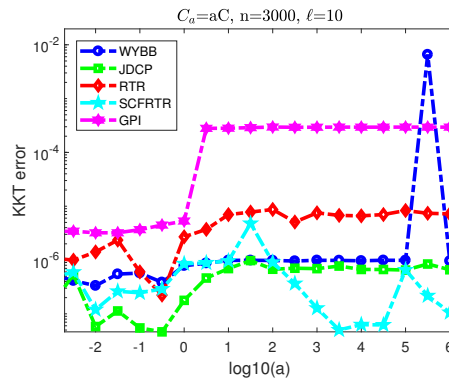


FIG. 11. KKT error for (6.3).

and

$$(6.3) \quad \min_{X \in \mathbb{Q}^{n \times \ell}} \|C_a X - D\|_F^2, \quad C_a = aC,$$

where B, C , and D are generated with entries randomly drawn from standard normal distribution, and $A = C^T C$ in (6.2). By fixing $n = 3000$ and $\ell = 10$, in Figures 8–13, we plot the averaged quantities of the difference of optimal values, KKT errors, and CPU time, for the minimization problems (6.2) and (6.3), over ten random tests for each a .

We observed that as a gets large, the quadratic term $\text{trace}(X^T A X)$ dominates the function $f(X)$ and thus (6.2) gradually becomes an eigenvalue problem, for which our eigenvalue-based solver SCFRTR manifests its efficiency. This can be seen by the consuming CPU times (we also noticed that the number of outer-loop iterations of SCFRTR does not change much when $a \geq 10^3$ and $a \geq 10^2$ for (6.2) and (6.3), respectively). On the other hand, for a small $a \leq 1$, even though all solvers efficiently handle (6.2) and (6.3), our method SCFRTR generally needs a little bit more CPU time to reach the stopping rule than RTR, WYBB, and JDCP. We think that in this case, the linear term $\text{trace}(X^T B)$ dominates $f(X)$, and gradient-based methods are able to solve it efficiently; the consuming CPU time of our SCFRTR is due to the calling of the eigensolver LOBPCG. Overall, this type of numerical testing also illustrates that our eigenvalue-based method SCFRTR is an efficient candidate to

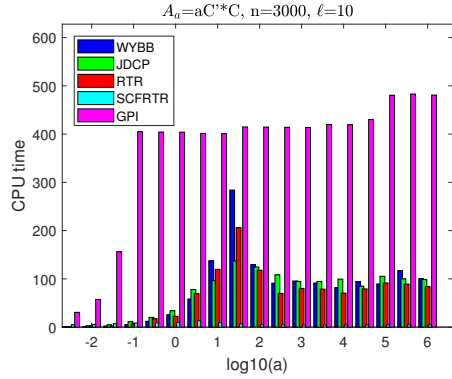


FIG. 12. CPU time for (6.2).

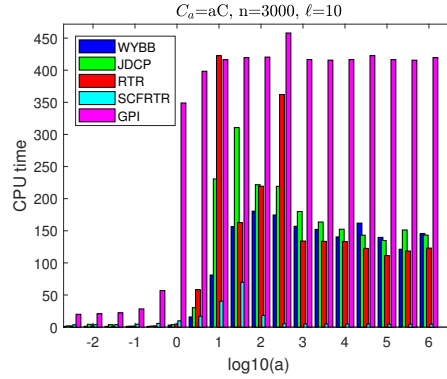


FIG. 13. CPU time for (6.3).

solve the minimization (1.2) or the Procrustes problem (1.1).

6.2.2. Test on the orthogonal least squares regression for feature extraction. We now evaluate the efficiency of our method on a practical model in data mining. The orthogonal least squares regression (OLSR) proposed in [51] (see also [30]) is a supervised learning method in linear discriminant analysis. Suppose we are given a dataset $S = [\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_m] \in \mathbb{R}^{n \times m}$ containing m samples with n features drawn from ℓ classes. Let $K = [\mathbf{k}_1, \dots, \mathbf{k}_m] \in \mathbb{R}^{\ell \times m}$ be the corresponding class indicator matrix; in particular, if the sample \mathbf{s}_i belongs to the j th class, then $\mathbf{k}_i = \mathbf{e}_j = [0, \dots, 0, 1, 0, \dots, 0]^T \in \mathbb{R}^\ell$. OLSR is based on the least squares regression to establish a dimension reduction procedure, and the model is to find a transformation matrix $X \in \mathbb{O}^{n \times \ell}$ with usually $\ell \ll n$ and an associated bias $\mathbf{b} \in \mathbb{R}^\ell$ by solving

$$(6.4) \quad \min_{X \in \mathbb{O}^{n \times \ell}, \mathbf{b} \in \mathbb{R}^\ell} \{J(X, \mathbf{b}) := \|S^T X + \mathbf{e}\mathbf{b}^T - K^T\|_F^2\},$$

where $\mathbf{e} = [1, \dots, 1]^T \in \mathbb{R}^m$. It is claimed that the orthogonal constraint of the transformation matrix X is able to preserve more local structure and discriminant information [51]. Take the partial derivative of $J(X, \mathbf{b})$ with respect to \mathbf{b} and set it to zero to have $\mathbf{b} = (K\mathbf{e} - X^T S \mathbf{e})/m$, (6.4) can be solved by the following unbalanced Procrustes problem:

$$(6.5) \quad \min_{X \in \mathbb{O}^{n \times \ell}} \|CX - D\|_F^2, \quad \text{where } C = \left(I_m - \frac{\mathbf{e}\mathbf{e}^T}{m}\right) S^T, \quad D = \left(I_m - \frac{\mathbf{e}\mathbf{e}^T}{m}\right) K^T.$$

For our testing in this case, we choose image datasets and text datasets⁴ summarized in Table 1. With the same stopping criterion used in subsection 6.2.1, we report the numerical results of the five methods in Table 2. Our experiments on the nine datasets in Table 1 demonstrate the efficiency of SCFRTR for OLSR model (6.5).

TABLE 1
Summary of test datasets.

	Dataset	Feature (n)	Number of samples (m)	Number of classes (ℓ)
Image	ORL	10304	400	40
	Mpeg-7	6000	1400	70
Text	Text-1	7511	1946	2
	20Newsgroups-4	8014	3970	4
	Cora-DS	6234	751	9
	Cora-HA	3989	400	7
	Cora-ML	8329	1617	4
	Cora-OS	6737	1246	4
	Cora-PL	7949	1575	9

TABLE 2
Numerical results on OLSR model for the datasets in Table 1.

Dataset	KKT error					CPU time (in seconds)				
	RTR	WYBB	JDCP	GPI	SCFRTR	RTR	WYBB	JDCP	GPI	SCFRTR
ORL	4.96e-6	9.66e-7	2.25e-6	3.98e-5	3.12e-7	5.95e+2	7.34e+2	8.26e+2	1.42e+4	9.28e+1
Mpeg-7	8.28e-5	9.02e-7	8.57e-7	3.47e-4	9.53e-7	7.55e+2	8.54e+2	8.05e+2	5.54e+3	3.82e+2
Text-1	2.97e-5	9.32e-7	3.44e-7	1.17e-5	8.17e-7	2.94e+1	3.48e+1	4.39e+1	9.28e+2	1.02e+1
20Newsgroups-4	7.52e-5	9.79e-7	8.76e-7	5.61e-5	9.82e-7	6.62e+1	7.40e+2	5.00e+2	2.08e+3	7.94e+1
Cora-DS	9.11e-5	9.38e-7	3.21e-7	3.64e-5	5.26e-7	4.29e+1	7.20e+1	1.09e+2	1.93e+3	5.45e+1
Cora-HA	6.88e-5	7.00e-7	4.14e-7	3.14e-5	4.56e-7	3.06e+1	2.75e+1	3.49e+1	7.41e+2	2.07e+1
Cora-ML	2.23e-5	9.14e-7	1.20e-6	9.20e-6	9.29e-7	1.46e+2	8.83e+1	1.18e+2	2.99e+3	7.35e+1
Cora-OS	6.40e-6	8.83e-7	2.69e-7	7.85e-6	4.40e-7	6.80e+1	4.61e+1	6.90e+1	1.51e+3	3.28e+1
Cora-PL	3.21e-5	9.60e-7	4.53e-7	3.03e-5	4.80e-7	1.00e+2	8.95e+1	1.03e+2	3.31e+3	8.39e+1

6.3. The effectiveness in finding a better KKT point. Besides the computation efficiency, we finally remark that our SCF iteration, which is based on a necessary optimality condition (2.4) for the global minimizer, seems more prone to find the global minimizer than a generic optimization method. Using the same stopping criterion for each method as in section 6.2, we have two types of numerical evaluation for this claim.

In the first test, for a given pair (n, ℓ) , we generate 10^4 problems (1.2) each with $\mathbf{C}=\text{randn}(n)$, $\mathbf{A}=\mathbf{C}'\mathbf{C}$, and $\mathbf{B}=\text{randn}(n, \ell)$. For each particular test problem, we compare the final computed objective function values from RTR, WYBB, JDCP, GPI,

⁴All datasets are available at <http://www.escience.cn/people/fpnie/papers.html>, except for the dataset ORL. The ORL database (<http://featureselection.asu.edu/datasets.php>) consists of 400 face images taken from 40 distinct subjects, each with 10 images taken at different times, varying the lighting, facial expressions (open/closed eyes, smiling/not smiling), and facial details (glasses/no glasses). The description of other datasets in Table 1 can be found in [13, Chapter 5].

and SCF⁵ starting from the same random initial guess. In order to evaluate whether a method M_1 returns a better computed solution than the other method M_2 (denoted by $M_1 \triangleright M_2$), we count the numbers of cases for M_1 when $f_{M_2} - f_{M_1} \geq 0.1 \cdot |f_{M_1}|$ (such a criterion, in general, protects the effect from the rounding-off errors), where f_{M_1} and f_{M_2} represent the computed objective function values from M_1 and M_2 , respectively. Table 3 summarizes the number of problems for which the objective function values from SCF are superior (denoted by SCF \triangleright) and inferior to (denoted by SCF \triangleleft) others for various pairs (n, ℓ) . Also reported are the mean and standard deviation of the objective function values for each method. As the computed solution from the SCF iteration satisfies a global necessary condition (2.4), it can be seen that, in many cases, our customized SCF iteration can achieve a KKT point with a smaller objective value than others.

TABLE 3

Number of cases (among 10^4) in which SCF is superior and inferior to others, and the mean and standard deviation of objective values for each method.

(n, ℓ)	Number of SCF \triangleright M (SCF \triangleleft M)				
	Mean (std) of the computed f_M				
	RTR	WYBB	JDCP	GPI	SCF
(5, 2)	1765 (134) -2.053 (1.1)	1717 (168) -2.064 (1.1)	1719 (142) -2.060 (1.1)	1815 (178) -2.060 (1.1)	- -2.214 (1.0)
(5, 3)	3239 (169) -2.128 (1.5)	3197 (170) -2.135 (1.5)	3076 (166) -2.146 (1.5)	3334 (217) -2.120 (1.5)	- -2.416 (1.5)
(10, 3)	954 (98) -4.386 (1.3)	930 (117) -4.390 (1.3)	893 (103) -4.393 (1.3)	1025 (164) -4.384 (1.3)	- -4.471 (1.2)
(10, 5)	1940 (99) -5.009 (2.0)	1945 (104) -5.009 (2.0)	1827 (100) -5.021 (2.0)	1972 (126) -5.008 (2.0)	- -5.218 (2.0)
(10, 7)	3031 (175) -1.170 (3.4)	3065 (204) -1.170 (3.4)	2894 (190) -1.183 (3.4)	2983 (180) -1.174 (3.4)	- -1.395 (3.4)

In our second numerical evaluation, the data matrices $F, M \in \mathbb{R}^{10 \times 4}$ (see [6, Table 19.1]) given in (6.6) are from the real application of the Procrustes problem (1.1) in the multidimensional similarity structure analysis (SSA) [6].

(6.6)

$$F = \begin{bmatrix} 0.08 & 0.59 & 2.54 & 0.23 \\ 2.05 & 0.31 & 0.20 & 0.03 \\ 1.22 & 0.73 & 0.33 & 1.63 \\ 2.11 & 0.25 & 0.79 & 1.15 \\ 1.75 & -0.14 & 1.56 & -0.93 \\ 1.84 & -0.12 & 1.17 & 1.09 \\ 0.27 & 1.31 & 1.82 & -0.68 \\ 0.41 & 2.17 & 0.04 & 0.91 \\ 0.49 & 2.46 & -0.22 & -0.07 \\ 0.13 & 1.96 & 0.14 & -0.39 \end{bmatrix}, \quad M = \begin{bmatrix} 1.89 & -0.76 & 1.50 & -0.76 \\ 2.30 & -1.30 & 0.27 & 1.41 \\ 2.26 & -0.34 & -1.15 & -1.38 \\ 2.53 & -1.21 & -0.92 & 0.29 \\ 1.65 & -1.05 & 1.33 & -0.38 \\ 2.20 & -1.42 & -0.71 & 0.28 \\ 2.06 & 0.88 & 1.46 & -0.39 \\ 2.54 & 1.63 & -1.22 & 0.10 \\ 1.83 & 2.01 & -0.23 & -0.34 \\ 1.74 & 2.01 & 0.77 & 1.00 \end{bmatrix}.$$

The matrix F represents the coordinates (configuration) of a vector representation of ten emotions, and M is the corresponding solution for a replication study [6, Chapter 19]. In SSA, two configurations are similar if they can be brought to a complete match by rigid motions and dilations [6]. As pointed out in [6, section 19.8], two configurations generally have the same dimensionality, but sometimes it is useful for

⁵Due to the small dimension n , in this subsection, we inactivated our subspace acceleration procedure in subsection 5.2 and just used the basic SCF iteration.

them not to have the same dimensionality. A technical treatment [6, section 19.8] for the latter situation is to append zero columns on one configuration until it matches the other, and then solve the balanced Procrustes problem. Another possible way is just to solve the unbalanced Procrustes problem (1.1). Therefore, in our experiments, we choose either F or M with certain columns deleted as the target matrix D to be used in (1.1), and the other is for C . Note that $F, M \in \mathbb{R}^{10 \times 4}$ and we can delete one and two columns to have $\ell = 3$ and $\ell = 2$, respectively. For each problem's result, we randomly generate 10^4 initial points for the five methods. Excluding cases where all methods return the same objective function value, in Tables 4 and 5, we report the numbers of problems for which the objective values from SCF are superior and inferior to others among these 10^4 tests. The results again indicate that, in many cases, SCF can converge to a KKT point with a smaller objective value than others.

TABLE 4

Number of cases (among 10^4) in which SCF is superior and inferior to others, and the mean and standard deviation of objective values. M (with certain columns deleted) is the D in (1.1).

		Number of SCF \triangleright M (SCF \triangleleft M) Mean (std) of the computed f_M				
	$D = M$ with	RTR	WYBB	JDCP	GPI	SCF
$\ell = 3$	$M(:, 1) = []$	3819 (0)	3676 (0)	3722 (0)	4573 (0)	-
		5.403 (3.0)	5.315 (3.0)	5.344 (3.0)	5.866 (3.0)	3.057 (0.0)
$\ell = 3$	$M(:, 2) = []$	3106 (0)	3135 (0)	3193 (0)	2601 (0)	-
		5.699 (2.9)	5.716 (2.9)	5.752 (2.9)	5.388 (2.7)	3.786 (0.0)
$\ell = 2$	$M(:, [1, 2]) = []$	3337 (0)	3234 (0)	3308 (0)	3890 (0)	-
		4.661 (2.9)	4.598 (2.9)	4.643 (2.9)	5.001 (3.0)	2.609 (0.0)

TABLE 5

Number of cases (among 10^4) in which SCF is superior and inferior to others, and the mean and standard deviation of objective values. F (with certain columns deleted) is the target D in (1.1).

		Number of SCF \triangleright M (SCF \triangleleft M) Mean (std) of the computed f_M				
	$D = F$ with	RTR	WYBB	JDCP	GPI	SCF
$\ell = 3$	$F(:, 1) = []$	2562 (2068)	2171 (2173)	2366 (2224)	1822 (2507)	-
		4.078 (2.1)	3.862 (2.1)	3.929 (2.2)	3.565 (2.0)	3.862 (2.0)
$\ell = 3$	$F(:, 2) = []$	4100 (0)	4091 (0)	4226 (0)	4348 (0)	-
		5.712 (3.0)	5.706 (3.0)	5.789 (3.0)	5.864 (3.1)	3.185 (0.0)

7. Conclusions. In this paper, we treat the numerical solution of the unbalanced orthogonal Procrustes problem or the minimization of a quadratic function over a Stiefel manifold from the eigenvalue point of view. With a motivation of exploiting the modern techniques in eigenvalue computations, we first connect the problem (1.2) with an eigenvalue minimization (3.2), for which an eigenvector-dependent nonlinear eigenvalue problem (3.3) can be established. Our proposed algorithm is a succinct SCF iteration in which each iteration involves computing a particular eigenbasis matrix of a symmetric matrix. The convergence behavior of this basic SCF iteration is analyzed and two practical implementation treatments are further built in to make the pure SCF iteration more efficient. Preliminary numerical testings are conducted to indicate that such an eigenvalue-based approach is an efficient way to solve the unbalanced Procrustes problem or the minimization (1.2), particularly for $\ell \ll n$.

Acknowledgments. The authors are grateful to the associate editor and two anonymous referees for their helpful comments and suggestions which significantly improved the presentation of this paper.

REFERENCES

- [1] P.-A. ABSIL, C. G. BAKER, AND K. A. GALLIVAN, *Trust-region methods on Riemannian manifolds*, Found. Comput. Math., 7 (2007), pp. 303–330.
- [2] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization Algorithms on Matrix Manifolds*, Princeton University Press, Princeton, NJ, 2008.
- [3] L.-E. ANDERSSON AND T. ELFVING, *A constrained Procrustes problem*, SIAM J. Matrix Anal. Appl., 18 (1997), pp. 124–139, <https://doi.org/10.1137/S0895479894277545>.
- [4] Z. BAI, D. LU, AND B. VANDEREYCKEN, *Robust Rayleigh quotient minimization and nonlinear eigenvalue problems*, SIAM J. Sci. Comput., 40 (2018), pp. A3495–A3522, <https://doi.org/10.1137/18M1167681>.
- [5] A. W. BOJANCZYK AND A. LUTOBORSKI, *The Procrustes problem for orthogonal Stiefel matrices*, SIAM J. Sci. Comput., 21 (1999), pp. 1291–1304, <https://doi.org/10.1137/S106482759630992X>.
- [6] I. BORG AND J. LINGOES, *Multidimensional Similarity Structure Analysis*, Springer-Verlag, New York, 1987.
- [7] Y. CAI, L.-H. ZHANG, Z. BAI, AND R.-C. LI, *On an eigenvector-dependent nonlinear eigenvalue problem*, SIAM J. Matrix Anal. Appl., 39 (2018), pp. 1360–1382, <https://doi.org/10.1137/17M115935X>.
- [8] M. T. CHU AND N. T. TRENDAFILOV, *The orthogonally constrained regression revisited*, J. Comput. Graph. Stat., 10 (2001), pp. 746–771.
- [9] C. DAVIS AND W. KAHAN, *The rotation of eigenvectors by a perturbation*, III, SIAM J. Numer. Anal., 7 (1970), pp. 1–46, <https://doi.org/10.1137/0707001>.
- [10] J. P. VAN DE GEER, *Linear relations among k sets of variables*, Psychometrika, 49 (1984), pp. 70–94.
- [11] J. W. DEMMEL, *Applied Numerical Linear Algebra*, SIAM, Philadelphia, 1997, <https://doi.org/10.1137/1.9781611971446>.
- [12] A. EDELMAN, T. A. ARIAS, AND S. T. SMITH, *The geometry of algorithms with orthogonality constraints*, SIAM J. Matrix Anal. Appl., 20 (1999), pp. 303–353, <https://doi.org/10.1137/S0895479895290954>.
- [13] T. HOKIMOTO, ED., *Advances in Statistical Methodologies and Their Application to Real Problems*, InTech, Rijeka, Croatia, 2017.
- [14] L. ELDÉN AND H. PARK, *A Procrustes problem on the Stiefel manifold*, Numer. Math., 82 (1999), pp. 599–619.
- [15] D. M. GAY, *Computing optimal locally constrained steps*, SIAM J. Sci. Statist. Comput., 2 (1981), pp. 186–197, <https://doi.org/10.1137/0902016>.
- [16] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 4th ed., Johns Hopkins University Press, Baltimore, MD, 2013.
- [17] J. C. GOWER AND G. B. DIJKSTERHUIS, *Procrustes Problems*, Oxford University Press, New York, 2004.
- [18] N. J. HIGHAM, *The symmetric Procrustes problem*, BIT, 28 (1998), pp. 133–143.
- [19] J. R. HURLEY AND R. B. CATTELL, *The Procrustes program: Producing direct rotation to test a hypothesized factor structure*, Beh. Sci., 7 (1962), pp. 258–262.
- [20] B. JIANG AND Y.-H. DAI, *A framework of constraint preserving update schemes for optimization on Stiefel manifold*, Math. Program., 153 (2015), pp. 535–575.
- [21] C. KANZOW AND H.-D. QI, *A QP-free constrained Newton-type method for variational inequality problems*, Math. Program., 85 (1999), pp. 81–106.
- [22] A. V. KNYAZEV, *Toward the optimal preconditioned eigensolver: Locally optimal block preconditioned conjugate gradient method*, SIAM J. Sci. Comput., 23 (2001), pp. 517–541, <https://doi.org/10.1137/S1064827500366124>.
- [23] A. V. KNYAZEV AND K. NEYMEYR, *Efficient solution of symmetric eigenvalue problems using multigrid preconditioners in the locally optimal block conjugate gradient method*, Electron. Trans. Numer. Anal., 15 (2003), pp. 38–55.
- [24] L. LI AND Z. ZHANG, *Semi-supervised domain adaptation by covariance matching*, IEEE Trans. Pattern Anal. Mach. Intell., 41 (2019), pp. 2724–2739.
- [25] X. LIU, X. WANG, Z. WEN, AND Y. YUAN, *On the convergence of the self-consistent field iteration in Kohn–Sham density functional theory*, SIAM J. Matrix Anal. Appl., 35 (2014),

- pp. 546–558, <https://doi.org/10.1137/130911032>.
- [26] X.-G. LIU, X.-F. WANG, AND W.-G. WANG, *Maximization of matrix trace function of product Stiefel manifolds*, SIAM J. Matrix Anal. Appl., 36 (2015), pp. 1489–1506, <https://doi.org/10.1137/15M100883X>.
 - [27] R. M. MARTIN, *Electronic Structure: Basic Theory and Practical Methods*, Cambridge University Press, Cambridge, UK, 2004.
 - [28] J. M. MARTÍNEZ, *Local minimizers of quadratic functions on Euclidean balls and spheres*, SIAM J. Optim., 4 (1994), pp. 159–176, <https://doi.org/10.1137/0804009>.
 - [29] J. J. MORÉ AND D. C. SORENSEN, *Computing a trust region step*, SIAM J. Sci. Statist. Comput., 4 (1983), pp. 553–572, <https://doi.org/10.1137/0904038>.
 - [30] F. NIE, R. ZHANG, AND X. LI, *A generalized power iteration method for solving quadratic problem on the Stiefel manifold*, Sci. China Inf. Sci., 60 (2017), 112101.
 - [31] J. NOCEDAL AND S. WRIGHT, *Numerical Optimization*, 2nd ed., Springer, New York, 2006.
 - [32] L. QIU, Y. ZHANG, AND C.-K. LI, *Unitarily invariant metrics on the Grassmann space*, SIAM J. Matrix Anal. Appl., 27 (2005), pp. 507–531, <https://doi.org/10.1137/040607605>.
 - [33] Y. SAAD, J. R. CHELIKOWSKY, AND S. M. SHONTZ, *Numerical methods for electronic structure calculations of materials*, SIAM Rev., 52 (2010), pp. 3–54, <https://doi.org/10.1137/060651653>.
 - [34] G. A. F. SEBER, *A Matrix Handbook for Statisticians*, John Wiley & Sons, Hoboken, NJ, 2007.
 - [35] D. C. SORENSEN, *Newton's method with a model trust region modification*, SIAM J. Numer. Anal., 19 (1982), pp. 409–426, <https://doi.org/10.1137/0719026>.
 - [36] G. W. STEWART, *Matrix Algorithms: Volume II: Eigensystems*, SIAM, Philadelphia, 2001, <https://doi.org/10.1137/1.9780898718058>.
 - [37] G. W. STEWART AND J.-G. SUN, *Matrix Perturbation Theory*, Academic Press, Boston, 1990.
 - [38] T. VIKLANDS, *Algorithms for the Weighted Orthogonal Procrustes Problem and Other Least Squares Problems*, Ph.D. thesis, Umeå University, Umeå, Sweden, 2006.
 - [39] J. VON NEUMANN, *Some matrix-inequalities and metrization of matrix-space*, Tomck. Univ. Rev., 1 (1937), pp. 286–300.
 - [40] Z. WANG, Q. RUAN, AND G. AN, *Projection-optimal local Fisher discriminant analysis for feature extraction*, Neural Comput. & Applic., 26 (2015), pp. 589–601.
 - [41] G. A. WATSON, *The solution of orthogonal Procrustes problems for a family of orthogonally invariant norms*, Adv. Comput. Math., 2 (1994), pp. 393–405.
 - [42] Z. WEN AND W. YIN, *A feasible method for optimization with orthogonality constraints*, Math. Program., 142 (2013), pp. 397–434.
 - [43] C. YANG, W. GAO, AND J. C. MEZA, *On the convergence of the self-consistent field iteration for a class of nonlinear eigenvalue problems*, SIAM J. Matrix Anal. Appl., 30 (2009), pp. 1773–1788, <https://doi.org/10.1137/080716293>.
 - [44] L.-H. ZHANG, *Uncorrelated trace ratio LDA for undersampled problems*, Patt. Recog. Lett., 32 (2011), pp. 476–484.
 - [45] L.-H. ZHANG, *On a self-consistent-field-like iteration for maximizing the sum of the Rayleigh quotients*, J. Comput. Appl. Math., 257 (2014), pp. 14–28.
 - [46] L.-H. ZHANG AND R.-C. LI, *Maximization of the sum of the trace ratio on the Stiefel manifold, I: Theory*, Sci. China Math., 57 (2014), pp. 2495–2508.
 - [47] L.-H. ZHANG AND R.-C. LI, *Maximization of the sum of the trace ratio on the Stiefel manifold, II: Computation*, Sci. China Math., 58 (2015), pp. 1549–1566.
 - [48] L.-H. ZHANG, L.-Z. LIAO, AND M. K. NG, *Fast algorithms for the generalized Foley–Sammon discriminant analysis*, SIAM J. Matrix Anal. Appl., 31 (2010), pp. 1584–1605, <https://doi.org/10.1137/080720863>.
 - [49] Z. Y. ZHANG AND K. Q. DU, *Successive projection method for solving the unbalanced procrustes problem*, Sci. China Math., 49 (2006), pp. 971–986.
 - [50] Z. Y. ZHANG, Y. QIU, AND K. Q. DU, *Conditions for optimal solutions of unbalanced Procrustes problem on Stiefel manifold*, J. Comput. Math., 25 (2007), pp. 661–671.
 - [51] H. ZHAO, Z. WANG, AND F. NIE, *Orthogonal least squares regression for feature extraction*, Neurocomputing, 216 (2016), pp. 200–207.