

## ON AN EIGENVECTOR-DEPENDENT NONLINEAR EIGENVALUE PROBLEM\*

YUNFENG CAI<sup>†</sup>, LEI-HONG ZHANG<sup>‡</sup>, ZHAOJUN BAI<sup>§</sup>, AND REN-CANG LI<sup>¶</sup>

**Abstract.** We first provide existence and uniqueness conditions for the solvability of an algebraic eigenvalue problem with eigenvector nonlinearity. We then present a local and global convergence analysis for a self-consistent field (SCF) iteration for solving the problem. The well-known  $\sin \Theta$  theorem in the perturbation theory of Hermitian matrices plays a central role. The near-optimality of the local convergence rate of the SCF iteration revealed in this paper is demonstrated by examples from the discrete Kohn–Sham eigenvalue problem in electronic structure calculations and the maximization of the trace ratio in the linear discriminant analysis for dimension reduction.

**Key words.** nonlinear eigenvalue problem, self-consistent field iteration, convergence analysis

**AMS subject classifications.** 65F15, 65H17

**DOI.** 10.1137/17M115935X

**1. Introduction.** The following eigenvector-dependent nonlinear eigenvalue problem (NEPv) is to find  $V \in \mathbb{C}^{n \times k}$  with orthonormal columns and  $\Lambda \in \mathbb{C}^{k \times k}$  such that

$$(1.1) \quad H(V)V = V\Lambda,$$

where  $H(V) \in \mathbb{C}^{n \times n}$  is a Hermitian matrix-valued function of  $V \in \mathbb{C}^{n \times k}$  with orthonormal columns, i.e.,  $V^H V = I_k$ ,  $k \leq n$  (usually  $k \ll n$ ). Immediately, we infer from (1.1) that  $\Lambda = V^H H(V)V$ , necessarily Hermitian, and the eigenvalues of  $\Lambda$  are  $k$  of the  $n$  eigenvalues of  $H(V)$ . For the problem of practical interests, they are usually either the  $k$  smallest or the  $k$  largest eigenvalues of  $H(V)$ . We will state all our results for the case of the smallest eigenvalues. But they are equally valid if the word “smallest” is replaced by “largest.”

Often the dependency on  $V$  of  $H(V)$  satisfies

$$(1.2) \quad H(V) \equiv H(VQ) \quad \text{for any unitary } Q \in \mathbb{C}^{k \times k}.$$

The condition (1.2) implies that  $H(V)$  is a function of  $k$ -dimensional subspaces of  $\mathbb{C}^n$  or, equivalently, a function on the complex Grassmann manifold  $\mathcal{G}_k(\mathbb{C}^n)$ . In particular,

---

\*Received by the editors December 1, 2017; accepted for publication (in revised form) by W.-W. Lin June 4, 2018; published electronically September 13, 2018.

<http://www.siam.org/journals/simax/39-3/M115935.html>

**Funding:** The first author’s work was supported in part by National Natural Science Foundation of China grants NSFC-11671023 and NSFC-11301013. The second author’s work was supported in part by National Natural Science Foundation of China grants NSFC-11671246 and NSFC-91730303. The third author’s work was supported in part by NSF grants DMS-1522697 and CCF-1527091. The fourth author’s work was supported in part by NSF grants CCF-1527104 and DMS-1719620.

<sup>†</sup>LMAM & School of Mathematical Sciences, Peking University, Beijing, 100871, China (yfcai@math.pku.edu.cn).

<sup>‡</sup>School of Mathematics and Shanghai Key Laboratory of Financial Information Technology, Shanghai University of Finance and Economics, Shanghai 200433, China (zhang.leihong@mail.shufe.edu.cn).

<sup>§</sup>Department of Computer Science and Department of Mathematics, University of California, Davis, CA 95616 (zbai@ucdavis.edu).

<sup>¶</sup>School of Mathematical Sciences, Xiamen University, Xiamen, 361005 China, and Department of Mathematics, University of Texas at Arlington, Arlington, TX 76019-0408 (rccli@uta.edu).

if  $V$  is a solution, then so is  $VQ$  for any  $k \times k$  unitary matrix  $Q$ . Therefore, any solution  $V$  to (1.1) essentially represents a class  $\{VQ : Q \in \mathbb{C}^{k \times k}, Q^H Q = I_k\}$ , each of which solves (1.1). In light of this, we say that the solution to (1.1) is unique if  $V, \tilde{V}$  are two solutions to (1.1); then  $\mathcal{R}(V) = \mathcal{R}(\tilde{V})$ , where  $\mathcal{R}(V)$  and  $\mathcal{R}(\tilde{V})$  are the column subspaces of  $V$  and  $\tilde{V}$ , respectively.

The most well known origin of NEPv (1.1) is from Kohn–Sham density functional theory in electronic structure calculations; see [15, 20, 5] and references therein. NEPv (1.1) also arises from the discretized Gross–Pitaevskii equation for modeling particles in the state of matter called the Bose–Einstein condensate [1, 7, 8], optimization of the trace ratio in the linear discriminant analysis for dimension reduction [16], and balanced graph cuts [9].

In the first part of this paper, we present two sufficient conditions for the existence and uniqueness of the solution of NEPv (1.1). One is a Lipschitz-like condition on the matrix-value function  $H(V)$ . The other is a uniform gap between the  $k$ th and  $(k+1)$ st smallest eigenvalues of  $H(V)$ , known as the “uniform well-posedness” property for the Hartree–Fock equation in electronic structure calculations [4]. To the best of our knowledge, it is the first such kind of results on the existence and uniqueness of the solution of NEPv (1.1) from the linear algebraic point of view.

Self-consistent field (SCF) iteration is the most widely used algorithm to solve NEPv (1.1); see [15, 20] and references therein. It is conceivably a natural one to try. At the  $i$ th SCF iteration, one computes an approximation to the eigenvector matrix  $V_i$  associated with the  $k$  smallest eigenvalues of  $H(V_{i-1})$  evaluated at the previous approximation  $V_{i-1}$ , and then  $V_i$  is used as the next approximation to the solution of NEPv (1.1). When the iterative process converges, the computed eigenvectors are said to be self-consistent. In the second part of this paper, we provide a local and global convergence analysis of a plain SCF iteration for solving NEPv (1.1). We use two examples to show the near-optimality of the newly established local convergence rate. We closely examine applications of derived convergence results to electronic structure calculations and linear discriminant analysis for dimension reduction. In particular, with weaker assumptions, we can significantly improve previous convergence results in [14, 27] on the SCF iteration for solving the discrete Kohn–Sham NEPv.

We will begin the presentation in section 2 with a review of matrix norms, angles between subspaces, and perturbation bounds to be used in this paper. In section 3, we establish the existence and uniqueness of NEPv (1.1) under a Lipschitz-like condition and the uniform well-posedness of the eigenvalue gap of  $H(V)$ . In section 4, we start by stating a plain SCF iteration for solving NEPv (1.1) and then establish local and global convergence results for the SCF iteration. In section 5, we discuss two applications. Concluding remarks are given in section 6.

*Notation.*  $\mathbb{C}^{n \times m}$  is the set of all  $n \times m$  matrices with complex entries,  $\mathbb{C}^n = \mathbb{C}^{n \times 1}$ , and  $\mathbb{C} = \mathbb{C}^1$ . Correspondingly, we will use  $\mathbb{R}^{n \times m}$ ,  $\mathbb{R}^n$ , and  $\mathbb{R}$  for their counterparts for the real number case. The superscripts “ $\cdot^T$ ” and “ $\cdot^H$ ” take the transpose and the complex conjugate transpose of a matrix or vector, respectively.  $\mathbb{U}^{n \times k} = \{V \mid V \in \mathbb{C}^{n \times k}, V^H V = I_k\}$ , i.e., the set of all  $n \times k$  complex matrices with orthonormal columns, and  $\mathcal{G}_k(\mathbb{C}^n)$  denotes the complex Grassmann manifold of all  $k$ -dimensional subspaces of  $\mathbb{C}^n$ .  $I_n$  (or simply  $I$  if its dimension is clear from the context) is the  $n \times n$  identity matrix, and  $e_j$  is its  $j$ th column.  $\mathcal{R}(X)$  is the column space of matrix  $X$ . Denote by  $\lambda_j(H)$  for  $1 \leq j \leq n$  the eigenvalues of a Hermitian matrix  $H \in \mathbb{C}^{n \times n}$ . They are always arranged in nondecreasing order:  $\lambda_1(H) \leq \lambda_2(H) \leq \dots \leq \lambda_n(H)$ .  $\text{Diag}(x)$  denotes the diagonal matrix with the vector  $x$  on its diagonal.  $\text{diag}(A)$  stands for the

column vector containing the diagonal elements of the matrix  $A$ .

**2. Preliminaries.** For completeness, in this section, we review matrix norms, angles between subspaces, and perturbation bounds to be used later in this paper.

*Unitarily invariant norm.* A matrix norm  $\|\cdot\|_{\text{ui}}$  is called a *unitarily invariant norm* on  $\mathbb{C}^{m \times n}$  if it is a matrix norm and has the following two properties:

1.  $\|X^H A Y\|_{\text{ui}} = \|A\|_{\text{ui}}$  for all unitary matrices  $X$  and  $Y$ .
2.  $\|A\|_{\text{ui}} = \|A\|_2$  whenever  $A$  is of rank one, where  $\|\cdot\|_2$  is the spectral norm.

Two commonly used unitarily invariant norms are

$$\text{the spectral norm: } \|A\|_2 = \max_j \sigma_j,$$

$$\text{the Frobenius norm: } \|A\|_F = \left( \sum_j \sigma_j^2 \right)^{1/2},$$

where  $\sigma_1, \sigma_2, \dots, \sigma_{\min\{m,n\}}$  are the singular values of  $A$ ; see, e.g., [2, 22].

In this paper, for convenience, any  $\|\cdot\|_{\text{ui}}$  we use is generic to matrix sizes in the sense that it applies to matrices of all sizes [22, page 79]. Examples include the matrix spectral norm  $\|\cdot\|_2$  and the Frobenius norm  $\|\cdot\|_F$ . One important property of unitarily invariant norms is

$$(2.1) \quad \|ABC\|_{\text{ui}} \leq \|A\|_2 \cdot \|B\|_{\text{ui}} \cdot \|C\|_2$$

for any matrices  $A, B$ , and  $C$  of compatible sizes. Comparing  $\|\cdot\|_2$  with any  $\|\cdot\|_{\text{ui}}$ , we have

$$(2.2) \quad \|A\|_2 \leq \|A\|_{\text{ui}} \leq \min\{m, n\} \|A\|_2$$

for any  $A \in \mathbb{C}^{m \times n}$ . Sharper bounds than this are possible for a particular unitarily invariant norm. For example,

$$(2.3) \quad \|A\|_2 \leq \|A\|_F \leq \sqrt{\min\{m, n\}} \|A\|_2.$$

*Angles between subspaces.* Consider the complex Grassmann manifold  $\mathcal{G}_k(\mathbb{C}^n)$  consisting of all  $k$ -dimensional subspaces of  $\mathbb{C}^n$ , and let  $\mathcal{X}, \mathcal{Y} \in \mathcal{G}_k(\mathbb{C}^n)$ . Let  $X, Y \in \mathbb{C}^{n \times k}$  be the orthonormal basis matrices of  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively, i.e.,

$$\mathcal{R}(X) = \mathcal{X}, \quad X^H X = I_k \quad \text{and} \quad \mathcal{R}(Y) = \mathcal{Y}, \quad Y^H Y = I_k,$$

and let  $\sigma_j$  for  $1 \leq j \leq k$  be the singular values of  $Y^H X$  in ascending order, i.e.,  $\sigma_1 \leq \dots \leq \sigma_k$ ; then the  $k$  canonical angles  $\theta_j(\mathcal{X}, \mathcal{Y})$  between  $\mathcal{X}$  and  $\mathcal{Y}$  are defined by

$$(2.4) \quad 0 \leq \theta_j(\mathcal{X}, \mathcal{Y}) := \arccos \sigma_j \leq \frac{\pi}{2} \quad \text{for } 1 \leq j \leq k.$$

They are in descending order, i.e.,  $\theta_1(\mathcal{X}, \mathcal{Y}) \geq \dots \geq \theta_k(\mathcal{X}, \mathcal{Y})$ . Set

$$(2.5) \quad \Theta(\mathcal{X}, \mathcal{Y}) = \text{Diag}(\theta_1(\mathcal{X}, \mathcal{Y}), \dots, \theta_k(\mathcal{X}, \mathcal{Y})).$$

It can be seen that angles so defined are independent of the orthonormal basis matrices  $X$  and  $Y$ . A different way to define these angles is through the orthogonal projections onto  $\mathcal{X}$  and  $\mathcal{Y}$  [26]. Note that when  $k = 1$ , i.e.,  $X$  and  $Y$  are vectors, there is only one canonical angle between  $\mathcal{X}$  and  $\mathcal{Y}$ , and so we will simply write  $\theta(\mathcal{X}, \mathcal{Y})$ . With the

definition of canonical angles, Sun [23, page 95] proved that for any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$  on  $\mathbb{C}^{k \times k}$ ,  $\|\sin \Theta(\mathcal{X}, \mathcal{Y})\|_{\text{ui}}$  defines a unitarily invariant metric on  $\mathcal{G}_k(\mathbb{C}^n)$ . On a related note, Qiu, Zhang, and Li [19] proved that  $\|\Theta(\mathcal{X}, \mathcal{Y})\|_{\text{ui}}$  also defines a unitarily invariant metric on  $\mathcal{G}_k(\mathbb{C}^n)$ .

In what follows, we sometimes place a vector or matrix in one or both arguments of  $\theta_j(\cdot, \cdot)$ ,  $\theta(\cdot, \cdot)$ , and  $\Theta(\cdot, \cdot)$  with the understanding that it is about the subspace spanned by the vector or the columns of the matrix argument. The following lemma provides a convenient way to compute  $\|\sin \Theta(\mathcal{X}, \mathcal{Y})\|_{\text{ui}}$ .

LEMMA 2.1. *Let  $[X, X_c]$  and  $[Y, Y_c]$  be two unitary matrices with  $X, Y \in \mathbb{C}^{n \times k}$ . Then*

$$\|\sin \Theta(X, Y)\|_{\text{ui}} = \|X_c^H Y\|_{\text{ui}} = \|X^H Y_c\|_{\text{ui}}$$

for any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$  on  $\mathbb{C}^{k \times k}$ .

Because orthonormal bases for subspaces are not unique, two subspaces  $\mathcal{X}$  and  $\mathcal{Y}$  of dimension  $k$  are close in terms of their canonical angles, or equivalently some norm  $\|\sin \Theta(\mathcal{X}, \mathcal{Y})\|_{\text{ui}}$  can have orthonormal basis matrices  $X, Y \in \mathbb{C}^{n \times k}$  that are far apart in the sense that  $\|X - Y\|_{\text{ui}} \gg \|\sin \Theta(\mathcal{X}, \mathcal{Y})\|_{\text{ui}}$ . The next lemma, whose proof can be found in [29, Lemma 4.1], says that one can always choose the basis matrices that differ from each other by  $O(\|\sin \Theta(\mathcal{X}, \mathcal{Y})\|_{\text{ui}})$ .

LEMMA 2.2 (see [29, Lemma 4.1]). *Suppose  $X, Y \in \mathbb{U}^{n \times k}$ . Then there exists a unitary matrix  $Q \in \mathbb{R}^{k \times k}$  such that*

$$(2.6) \quad \|\sin \Theta(X, Y)\|_{\text{ui}} \leq \|XQ - Y\|_{\text{ui}} \leq \sqrt{2} \|\sin \Theta(X, Y)\|_{\text{ui}}$$

for any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ .

Each  $\mathcal{X} \in \mathcal{G}_k(\mathbb{C}^n)$  can be represented uniquely by the orthogonal projector  $P_{\mathcal{X}}$  onto the subspace  $\mathcal{X}$ . Given  $X^H X = I_k$  such that  $\mathcal{X} = \mathcal{R}(X)$ , we have

$$P_{\mathcal{X}} = P_X := X X^H.$$

Note that even though  $P_{\mathcal{X}}$  is explicitly defined by  $X$ , it is independent of the choice of  $X$  so long as  $\mathcal{R}(X) = \mathcal{X}$ . Therefore, any norm on the differences between the orthogonal projectors induces a metric on  $\mathcal{G}_k(\mathbb{C}^n)$ . Naturally, we would ask whether there is any relation between  $\|\sin \Theta(X, Y)\|_{\text{ui}}$  and  $\|P_X - P_Y\|_{\text{ui}}$ . Indeed, for any  $X, Y \in \mathbb{U}^{n \times k}$ , we have

$$(2.7) \quad \|\sin \Theta(X, Y)\|_2 = \|P_X - P_Y\|_2, \quad \|\sin \Theta(X, Y)\|_{\text{F}} = \frac{1}{\sqrt{2}} \|P_X - P_Y\|_{\text{F}}.$$

Both equalities in (2.7) are the simple consequences of the fact that the singular values of  $P_X - P_Y$  consist of each  $\sin \theta_i(X, Y)$  repeated twice and  $n - 2k$  zeros [22, page 43]. In addition, we have

$$(2.8) \quad \frac{1}{2} \|P_X - P_Y\|_{\text{ui}} \leq \|\sin \Theta(X, Y)\|_{\text{ui}} \leq \|P_X - P_Y\|_{\text{ui}}$$

for any unitarily invariant norm  $\|\cdot\|_{\text{ui}}$ . Closely related, the singular values of  $P_X(I - P_Y)$  consist of all  $\sin \theta_i(X, Y)$  and  $n - k$  zeros [22, page 43], and therefore

$$(2.9) \quad \|\sin \Theta(X, Y)\|_{\text{ui}} = \|P_X(I - P_Y)\|_{\text{ui}} = \|P_Y(I - P_X)\|_{\text{ui}}.$$

In the rest of this paper, we will treat  $\mathcal{X}$ ,  $P_{\mathcal{X}}$ , and  $P_X$  indistinguishably whenever convenient.

*Perturbation of Hermitian matrices.* A well-known theorem of Weyl is the following.

LEMMA 2.3 (see [22, page 203]). *For two Hermitian matrices  $A, \tilde{A} \in \mathbb{C}^{n \times n}$ , we have*

$$|\lambda_j(A) - \lambda_j(\tilde{A})| \leq \|A - \tilde{A}\|_2 \quad \text{for } 1 \leq j \leq n.$$

The next lemma is essentially [6, Theorem 5.1] due to Davis and Kahan.

LEMMA 2.4 (see [6]). *Let  $H$  and  $M$  be two Hermitian matrices, and let  $S$  be a matrix of a compatible size as determined by the Sylvester equation*

$$HY - YM = S.$$

*If either all eigenvalues of  $H$  are contained in a closed interval that contains no eigenvalue of  $M$ , or vice versa, then the Sylvester equation has a unique solution  $Y$ , and, moreover,*

$$\|Y\|_{\text{ui}} \leq \frac{1}{\delta} \|S\|_{\text{ui}},$$

where  $\delta = \min |\mu - \omega|$  over all eigenvalues  $\mu$  of  $M$  and all eigenvalues  $\omega$  of  $H$ .

Finally, the well-known Davis–Kahan  $\sin \Theta$  theorem in [6] (see also [22, 13]) will play a central role in our later analysis. However, we will not explicitly state the theorem here for two reasons. The first reason is that it can be inferred from Lemma 2.4, and the second one is that we can derive a better locally convergent rate of the SCF iteration by going through the actual proof of the theorem, as we will later in this paper.

**3. Existence and uniqueness.** Recall (1.2) about the dependency of  $H(V)$  on  $V$ , which makes  $H(\cdot)$  a Hermitian matrix-valued function on the complex Grassmann manifold  $\mathcal{G}_k(\mathbb{C}^n)$ . For convenience, we will treat  $H(V)$  and  $H(\mathcal{V})$  indistinguishably, where  $\mathcal{V} = \mathcal{R}(V)$ . As a convention,  $\lambda_j(H(V))$  for  $1 \leq j \leq n$  are the eigenvalues of  $H(V)$ , arranged in nondecreasing order. The following theorem gives sufficient conditions for the existence and uniqueness of the solution of NEP<sub>V</sub> (1.1).

THEOREM 3.1. *Assume that for the given unitarily invariant norm  $\|\cdot\|_{\text{ui}}$  and the spectral norm  $\|\cdot\|_2$ , there exist positive constants  $\xi_{\text{ui}}$  and  $\xi_2$  such that for any  $V, \tilde{V} \in \mathbb{U}^{n \times k}$ ,*

$$(3.1a) \quad \|H(V) - H(\tilde{V})\|_{\text{ui}} \leq \xi_{\text{ui}} \|\sin \Theta(V, \tilde{V})\|_{\text{ui}},$$

$$(3.1b) \quad \|H(V) - H(\tilde{V})\|_2 \leq \xi_2 \|\sin \Theta(V, \tilde{V})\|_2,$$

and also assume that there exists a positive constant  $\delta$  such that for any  $V \in \mathbb{U}^{n \times k}$ ,

$$(3.2) \quad \lambda_{k+1}(H(V)) - \lambda_k(H(V)) \geq \delta.$$

*If  $\delta > \xi_{\text{ui}} + \xi_2$ , then NEP<sub>V</sub> (1.1) has a unique solution.*

*Remark 1.* Before we provide a proof, three comments are in order.

(a) The conditions in (3.1) are Lipschitz-like conditions. As we have pointed out, Weyl’s lemma (Lemma 2.3) and the Davis–Kahan  $\sin \Theta$  theorem play central roles in our analysis. The former requires the 2-norm inequality (3.1b), and the latter works for the general unitary-invariant norm inequality (3.1b). The proof below needs (3.1b), but at the place where (3.1a) is used, it can be simply replaced by using (3.1b) instead. Thus it may seem that the theorem is made unnecessarily more complicated

by including both conditions in (3.1) rather than just (3.1b) alone. But we argue that there are situations where the theorem is stronger, namely if and when  $\xi_{\text{ui}} < \xi_2$  for some  $\|\cdot\|_{\text{ui}}$  other than the spectral norm and hence the condition on  $\delta > \xi_{\text{ui}} + \xi_2$  is weaker than  $\delta > 2\xi_2$ . By the same logic, if  $\xi_{\text{ui}} \geq \xi_2$ , then we should just use the version of this theorem with  $\|\cdot\|_{\text{ui}}$  also being the spectral norm.

(b) Any one of the conditions in (3.1) yields one for the other by using (2.2). For example, (3.1a) leads to (3.1b) with  $\xi_2 = k\xi_{\text{ui}}$ , and likewise (3.1b) leads to (3.1a) with  $\xi_{\text{ui}} = n\xi_2$ . Conceivably, the resulting  $\xi$ -constant is likely worse than being obtained through a direct estimation.

(c) The assumption (3.2) requires a uniform gap between the  $k$ th and  $(k + 1)$ st eigenvalues of every  $H(V)$  for  $V \in \mathbb{U}^{n \times k}$ . This is known as the “uniform well-posedness” property for using the SCF iteration to solve the Hartree–Fock equation in electronic structure calculations [4]. It is undoubtedly strong and may be hard to verify.

*Proof of Theorem 3.1.* We prove the theorem by constructing a mapping on  $\mathcal{G}_k(\mathbb{C}^n)$  whose fixed-point is a solution to NEPv (1.1), and vice versa. For any  $\mathcal{V} \in \mathcal{G}_k(\mathbb{C}^n)$ , let  $V \in \mathbb{U}^{n \times k}$  such that  $\mathcal{R}(V) = \mathcal{V}$ . Because of (3.2),  $H(V)$  has a unique invariant subspace associated with its  $k$  smallest eigenvalues. We define  $\phi(\mathcal{V})$  to be that subspace. Any solution  $V$  to NEPv (1.1) satisfies  $\mathcal{R}(V) = \phi(\mathcal{R}(V))$ ; i.e.,  $\mathcal{R}(V)$  is a fixed-point of  $\phi$ , and vice versa.

In what follows, we will prove  $\phi$  is strictly contractive on  $\mathcal{G}_k(\mathbb{C}^n)$  endowed with the distance metric

$$\text{dist}(\mathcal{V}, \tilde{\mathcal{V}}) := \|\sin \Theta(\mathcal{V}, \tilde{\mathcal{V}})\|_{\text{ui}}.$$

To this end, we consider  $\mathcal{V}, \tilde{\mathcal{V}} \in \mathcal{G}_k(\mathbb{C}^n)$  and let  $V, \tilde{V} \in \mathbb{U}^{n \times k}$  such that  $\mathcal{R}(V) = \mathcal{V}$  and  $\mathcal{R}(\tilde{V}) = \tilde{\mathcal{V}}$ , respectively. Write the eigendecompositions of  $H(V)$  and  $H(\tilde{V})$  as

$$H(V) = [U, U_c] \text{Diag}(\Lambda, \Lambda_c) [U, U_c]^H \quad \text{and} \quad H(\tilde{V}) = [\tilde{U}, \tilde{U}_c] \text{Diag}(\tilde{\Lambda}, \tilde{\Lambda}_c) [\tilde{U}, \tilde{U}_c]^H,$$

where  $[U, U_c], [\tilde{U}, \tilde{U}_c] \in \mathbb{C}^{n \times n}$  are unitary,  $U, \tilde{U} \in \mathbb{U}^{n \times k}$ , and

$$\begin{aligned} \Lambda &= \text{Diag}(\lambda_1(H(V)), \dots, \lambda_k(H(V))), & \Lambda_c &= \text{Diag}(\lambda_{k+1}(H(V)), \dots, \lambda_n(H(V))), \\ \tilde{\Lambda} &= \text{Diag}(\lambda_1(H(\tilde{V})), \dots, \lambda_k(H(\tilde{V}))), & \tilde{\Lambda}_c &= \text{Diag}(\lambda_{k+1}(H(\tilde{V})), \dots, \lambda_n(H(\tilde{V}))). \end{aligned}$$

By Lemma 2.3 and the Lipschitz-like conditions (3.1), we have

$$|\lambda_j(H(V)) - \lambda_j(H(\tilde{V}))| \leq \|H(V) - H(\tilde{V})\|_2 \leq \xi_2 \|\sin \Theta(V, \tilde{V})\|_2 \quad \text{for } 1 \leq j \leq n,$$

which, together with (3.2) and  $\delta > \xi_{\text{ui}} + \xi_2$ , leads to

$$\begin{aligned} \lambda_{k+1}(H(V)) - \lambda_k(H(\tilde{V})) &= \lambda_{k+1}(H(V)) - \lambda_k(H(V)) + \lambda_k(H(V)) - \lambda_k(H(\tilde{V})) \\ &\geq \delta - \xi_2 \|\sin \Theta(V, \tilde{V})\|_2 \\ (3.3) \qquad \qquad \qquad &\geq \delta - \xi_2 > 0 \end{aligned}$$

since  $\|\sin \Theta(V, \tilde{V})\|_2 \leq 1$  always. Now define  $R = H(V)\tilde{U} - \tilde{U}\tilde{\Lambda}$ . We have

$$U_c^H R = \Lambda_c U_c^H \tilde{U} - U_c^H \tilde{U} \tilde{\Lambda}.$$

On the other hand, it can be seen that  $R = [H(V) - H(\tilde{V})]\tilde{U}$ . Therefore,

$$(3.4) \qquad \Lambda_c U_c^H \tilde{U} - U_c^H \tilde{U} \tilde{\Lambda} = U_c^H [H(V) - H(\tilde{V})]\tilde{U}.$$

Next we apply Lemmas 2.1 and 2.4 to get

$$(3.5a) \quad \text{dist}(\phi(\mathcal{V}), \phi(\tilde{\mathcal{V}})) = \|\sin \Theta(U, \tilde{U})\|_{\text{ui}} = \|U_c^H \tilde{U}\|_{\text{ui}}$$

$$(3.5b) \quad \leq \frac{1}{\lambda_{k+1}(H(V)) - \lambda_k(H(\tilde{V}))} \|U_c^H [H(V) - H(\tilde{V})] \tilde{U}\|_{\text{ui}}$$

$$(3.5c) \quad \leq \frac{\xi_{\text{ui}}}{\delta - \xi_2} \|\sin \Theta(V, \tilde{V})\|_{\text{ui}} \\ = \frac{\xi_{\text{ui}}}{\delta - \xi_2} \text{dist}(\mathcal{V}, \tilde{\mathcal{V}}),$$

where we have used Lemma 2.1 for (3.5a), applied Lemma 2.4 to (3.4) for (3.5b), and used the assumption (3.1a) and the inequality (3.3) for (3.5c). This completes the proof that the mapping  $\phi$  is strictly contractive on  $\mathcal{G}_k(\mathbb{C}^n)$  since the factor  $\xi_{\text{ui}}/(\delta - \xi_2) < 1$ . By the Banach fixed-point theorem [10],  $\phi$  has a unique fixed-point in  $\mathcal{G}_k(\mathbb{C}^n)$ , or equivalently NEPv (1.1) has a unique solution.  $\square$

In section 5, we will verify the satisfiability of the Lipschitz-like conditions (3.1) for two NEPvs arising in electronic structure calculations and linear discriminant analysis.

#### 4. SCF iteration and convergence analysis.

**4.1. SCF iteration.** A natural and most widely used method to solve NEPv (1.1) is the SCF iteration shown in Algorithm 1; see [15, 20] and references therein for its usage and variants in electronic structure calculations.

---

**Algorithm 1** SCF iteration for solving NEPv (1.1).

---

**Require:**  $V_0 \in \mathbb{C}^{n \times k}$  with orthonormal columns, i.e.,  $V_0^H V_0 = I_k$ ;

**Ensure:** a solution to NEPv (1.1).

- 1: **for**  $i = 1, 2, \dots$  until convergence **do**
  - 2:   construct  $H_i = H(V_{i-1})$
  - 3:   compute the partial eigenvalue decomposition  $H_i V_i = V_i \Lambda_i$ , where  $V_i \in \mathbb{U}^{n \times k}$  and  $\Lambda_i = \text{Diag}(\lambda_1(H_i), \dots, \lambda_k(H_i))$ .
  - 4: **end for**
  - 5: **return** the last  $V_i$  as a solution to NEPv (1.1).
- 

At each iterative step of SCF, a linear eigenvalue problem for  $H_i = H(V_{i-1})$  is partially solved. It is hoped that eventually  $\mathcal{R}(V_i)$  converges to some subspace  $\mathcal{V}_* \in \mathcal{G}_k(\mathbb{C}^n)$ . When it does, the orthonormal basis matrix  $V_*$  of  $\mathcal{V}_*$  will satisfy NEPv (1.1), provided  $H(V)$  is continuous at  $V_*$ . We will state all our convergence analysis explicitly for the case of the smallest eigenvalues in the next sections. However, in some other applications, such as the one to be discussed in section 5.2, the solutions  $V$  of interest to (1.1) are those such that the eigenvalues of  $\Lambda = V^H H(V) V$  correspond to the  $k$  largest eigenvalues of  $H(V)$ . The results of convergence analysis are equally valid if the word “smallest” is simply replaced by “largest,”

Note that at line 2 of Algorithm 1, we use the word “construct” to mean that sometimes  $H_i$  may not be explicitly computed but rather exists in some form in such a way that matrix-vector products by  $H_i$  can be efficiently performed.

To monitor the progress of convergence, we can compute the normalized residual

$$(4.1) \quad \text{NRes}_i = \frac{\|H_{i+1} V_i - V_i (V_i^H H_{i+1} V_i)\|}{\|H_{i+1}\| + \|\Lambda_i\|},$$

where  $\|\cdot\|$  is some matrix norm that is easy to compute, e.g., the Frobenius norm. But some of the quantities in defining  $\text{NRes}_i$  may not be necessarily needed in the SCF iteration, e.g.,  $H_{i+1}V_i$  and  $\|H_{i+1}\|$ , and they may not be cheap to compute. Therefore, we should not compute  $\text{NRes}_i$  too often, especially for the first many iterations of SCF when convergence has not yet happened. Also, only a very rough estimate of  $\|H_{i+1}\|$  is enough. There are metrics other than (4.1) that have been used to monitor the convergence of the SCF iteration when it comes to a particular application, e.g., the use of  $\rho(V_i) = \text{Diag}(V_i V_i^H)$  that corresponds to the charge density of electrons in electronic structure calculations (see [3, 4] and references therein). The idea is to examine the difference of  $\|\rho(V_i) - \rho(V_{i-1})\|$ . When it falls below a prescribed tolerance, convergence is claimed.

**4.2. Local convergence of SCF.** Let  $V_*$  be a solution to NEPv (1.1). The three assumptions we will make are as follows:

(A1) The eigenvalue gap

$$(4.2) \quad \delta_* = \lambda_{k+1}(H(V_*)) - \lambda_k(H(V_*)) > 0.$$

(A2) The matrix-valued function  $H(V)$  is continuous over  $\mathbb{U}^{n \times k}$  at  $V = V_*$ .

(A3) There exists a nonnegative constant  $\chi < \infty$  such that for some  $q \geq 1$ ,

$$(4.3) \quad \limsup_{\|\sin \Theta(V, V_*)\|_{\text{ui}} \rightarrow 0} \frac{\|(I - P_*)[H(V) - H(V_*)]P_*\|_{\text{ui}}}{\|\sin \Theta(V, V_*)\|_{\text{ui}}^q} \leq \chi,$$

where  $P_* = V_* V_*^H$  is the orthogonal projector onto  $\mathcal{R}(V_*)$ .

The following theorem is the main result on the local convergence of the SCF iteration.

**THEOREM 4.1.** *Assume (A1), (A2), (A3) and that  $\chi < \delta_*$  if  $q = 1$  (not necessary to assume  $\chi < \delta_*$  if  $q > 1$ ), and let  $\{V_i\}_i$  be the sequence generated by the SCF iteration (Algorithm 1) with initial guess  $V_0$ . If  $V_0$  is sufficiently close to  $V_*$  in the sense that  $\|\sin \Theta(V_0, V_*)\|_{\text{ui}}$  is sufficiently small, then there exists a sequence  $\{\tau_i\}_i$  such that*

$$(4.4) \quad \|\sin \Theta(V_i, V_*)\|_{\text{ui}} \leq \tau_{i-1} \|\sin \Theta(V_{i-1}, V_*)\|_{\text{ui}}^q$$

and

$$(4.5) \quad \lim_{i \rightarrow \infty} \tau_i = \frac{\chi}{\delta_*},$$

and the following hold:

1. for  $q = 1$ , all  $\tau_i < 1$ , and thus the SCF iteration is locally linearly convergent to  $\mathcal{R}(V_*)$  with the linear convergence rate no larger than  $\chi/\delta_*$ ;
2. for  $q > 1$ , the SCF iteration is locally convergent to  $\mathcal{R}(V_*)$  of order  $q$ .

*Proof.* For  $q = 1$ , as  $\chi < \delta_*$ , we can pick two positive constants  $\epsilon_1 < \delta_*/3$  and  $\epsilon_2$  such that

$$(4.6) \quad \tau := \frac{\chi + \epsilon_2}{\delta_* - 3\epsilon_1} < 1;$$

otherwise, for  $q > 1$ , any two positive constants  $\epsilon_1 < \delta_*/3$  and  $\epsilon_2$  will do. By (A2) and (A3), there exists a positive number  $\Delta$  with

$$(4.7) \quad \Delta \leq 1 \text{ for } q = 1 \quad \text{or} \quad \Delta < \min \left\{ 1, \left( \frac{\delta_* - 3\epsilon_1}{\chi + \epsilon_2} \right)^{1/(q-1)} \right\} \text{ for } q > 1$$

such that, whenever  $\|\sin \Theta(V, V_*)\|_{\text{ui}} \leq \Delta$ ,

$$(4.8a) \quad \|H(V) - H(V_*)\|_2 \leq \epsilon_1,$$

$$(4.8b) \quad \frac{\|(I - P_*)[H(V) - H(V_*)]P_*\|_{\text{ui}}}{\|\sin \Theta(V, V_*)\|_{\text{ui}}^q} \leq \chi + \epsilon_2.$$

Now suppose that  $\|\sin \Theta(V_0, V_*)\|_{\text{ui}} \leq \Delta$  and set

$$(4.9a) \quad \epsilon_{i1} = \|H(V_i) - H(V_*)\|_2,$$

$$(4.9b) \quad \epsilon_{i2} = \max \left\{ \frac{\|(I - P_*)[H(V_i) - H(V_*)]P_*\|_{\text{ui}}}{\|\sin \Theta(V_i, V_*)\|_{\text{ui}}^q} - \chi, 0 \right\},$$

$$(4.9c) \quad \tau_i = \frac{\chi + \epsilon_{i2}}{\delta_* - 3\epsilon_{i1}}.$$

Then  $\epsilon_{01} \leq \epsilon_1 < \delta_*/3$  and  $\epsilon_{02} \leq \epsilon_2$ , and hence

$$(4.10) \quad \tau_0 \Delta^{q-1} \leq \tau \Delta^{q-1} < 1.$$

To see this, for  $q = 1$ ,  $\tau_0 \leq \tau < 1$  by (4.6). If  $q > 1$ , (4.10) is a consequence of (4.7).

In what follows, we will prove that  $\|\sin \Theta(V_i, V_*)\|_{\text{ui}} \leq \Delta$  and (4.4) hold for all  $i \geq 1$ . As a consequence, the inequalities in (4.8) hold for  $V = V_i$ , and  $\epsilon_{i1} \leq \epsilon_1 < \delta_*/3$  and  $\epsilon_{i2} \leq \epsilon_2$ , and hence  $\tau_i \Delta^{q-1} \leq \tau \Delta^{q-1} < 1$ . These, in particular, imply that  $\|\sin \Theta(V_i, V_*)\|_{\text{ui}} \rightarrow 0$  because

$$\|\sin \Theta(V_i, V_*)\|_{\text{ui}} \leq \tau \Delta^{q-1} \|\sin \Theta(V_{i-1}, V_*)\|_{\text{ui}},$$

and the limiting equality in (4.4) holds.

Due to similarity, it suffices to show that  $\|\sin \Theta(V_1, V_*)\|_{\text{ui}} \leq \Delta$  and (4.4) hold for  $i = 1$ . Let the eigendecompositions of  $H(V_0)$  and  $H(V_*)$  be

$$H(V_0) = [V_1, V_{1c}] \text{Diag}(\Lambda_1, \Lambda_{1c}) [V_1, V_{1c}]^H$$

and

$$H(V_*) = [V_*, V_{*c}] \text{Diag}(\Lambda_*, \Lambda_{*c}) [V_*, V_{*c}]^H,$$

respectively, where  $[V_1, V_{1c}], [V_*, V_{*c}] \in \mathbb{C}^{n \times n}$  are unitary,  $V_1, V_* \in \mathbb{U}^{n \times k}$ , and

$$\begin{aligned} \Lambda_1 &= \text{Diag}(\lambda_1(H(V_0)), \dots, \lambda_k(H(V_0))), & \Lambda_{1c} &= \text{Diag}(\lambda_{k+1}(H(V_0)), \dots, \lambda_n(H(V_0))), \\ \Lambda_* &= \text{Diag}(\lambda_1(H(V_*)), \dots, \lambda_k(H(V_*))), & \Lambda_{*c} &= \text{Diag}(\lambda_{k+1}(H(V_*)), \dots, \lambda_n(H(V_*))). \end{aligned}$$

By Lemma 2.3, we know that

$$|\lambda_j(H(V_0)) - \lambda_j(H(V_*))| \leq \|H(V_0) - H(V_*)\|_2 = \epsilon_{01}$$

for  $1 \leq j \leq n$ . It follows that

$$(4.11) \quad \begin{aligned} \lambda_{k+1}(H(V_0)) - \lambda_k(H(V_*)) &= \lambda_{k+1}(H(V_0)) - \lambda_k(H(V_0)) + \lambda_k(H(V_0)) - \lambda_k(H(V_*)) \\ &\geq \delta_* - \epsilon_{01} > \frac{2\delta_*}{3} > 0. \end{aligned}$$

Now define  $R_1 = H(V_0)V_* - V_*\Lambda_*$ . We have

$$(4.12) \quad (V_{1c})^H R_1 = \Lambda_{1c} (V_{1c})^H V_* - (V_{1c})^H V_* \Lambda_*.$$

On the other hand, it can be seen that  $R_1 = [H(V_0) - H(V_*)]V_*$ . Therefore,

$$\Lambda_{1c}(V_{1c})^H V_* - (V_{1c})^H V_* \Lambda_* = (V_{1c})^H [H(V_0) - H(V_*)] V_*$$

Next we apply Lemmas 2.1 and 2.4, (2.1), and (2.8) to get

$$\begin{aligned} \|\sin \Theta(V_1, V_*)\|_{ui} &= \|(V_{1c})^H V_*\|_{ui} \\ &\leq \frac{1}{\lambda_{k+1}(H(V_0)) - \lambda_k(H(V_*))} \|(V_{1c})^H [H(V_0) - H(V_*)] V_*\|_{ui} \\ &= \frac{1}{\lambda_{k+1}(H(V_0)) - \lambda_k(H(V_*))} \|(I - P_1)[H(V_0) - H(V_*)] P_*\|_{ui} \\ &\leq \frac{1}{\delta_* - \varepsilon_{01}} \|(I - P_1)[H(V_0) - H(V_*)] P_*\|_{ui} \\ &\leq \frac{1}{\delta_* - \varepsilon_{01}} (\|(P_1 - P_*)[H(V_0) - H(V_*)] P_*\|_{ui} \\ &\quad + \|(I - P_*)[H(V_0) - H(V_*)] P_*\|_{ui}) \\ &\leq \frac{1}{\delta_* - \varepsilon_{01}} (2\|\sin \Theta(V_1, V_*)\|_{ui} \varepsilon_{01} + (\chi + \varepsilon_{02}) \|\sin \Theta(V_0, V_*)\|_{ui}^q). \end{aligned}$$

Solving the above inequality for  $\|\sin \Theta(V_1, V_*)\|_{ui}$ , we obtain

$$(4.13) \quad \|\sin \Theta(V_1, V_*)\|_{ui} \leq \tau_0 \|\sin \Theta(V_0, V_*)\|_{ui}^q \leq \tau \Delta^{q-1} \|\sin \Theta(V_0, V_*)\|_{ui},$$

where we have used  $\|\sin \Theta(V_0, V_*)\|_{ui} \leq \Delta$  for the second inequality. The first inequality in (4.13) is (4.4) for  $i = 1$ , and the second inequality there implies that  $\|\sin \Theta(V_1, V_*)\|_{ui} \leq \Delta \leq 1$  because  $\tau \Delta^{q-1} < 1$ .  $\square$

*Remark 2.* (a) The assumption (A1) is similar to the “uniform well-posedness” assumption (3.2) on the eigenvalue gap of  $H(V)$  in Theorem 3.1.

(b) Let  $[V_*, V_{*c}] \in \mathbb{C}^{n \times n}$  be unitary. Notice that

$$\|(I - P_*)[H(V) - H(V_*)] P_*\|_{ui} = \|V_{*c}^H [H(V) - H(V_*)] V_*\|_{ui}.$$

The assumption (A3) is about the closeness of the (2, 1)-block of

$$[V_*, V_{*c}]^H H(V) [V_*, V_{*c}]$$

in limit to the (2, 1)-block of  $[V_*, V_{*c}]^H H(V_*) [V_*, V_{*c}]$ , which is 0.

(c) The assumption (A3) with  $q = 1$  is weaker than the Lipschitz-like condition (3.1a). This is because (3.1a) implies

$$\begin{aligned} \limsup_{\|\sin \Theta(V, V_*)\|_{ui} \rightarrow 0} \frac{\|(I - P_*)[H(V) - H(V_*)] P_*\|_{ui}}{\|\sin \Theta(V, V_*)\|_{ui}} \\ \leq \limsup_{\|\sin \Theta(V, V_*)\|_{ui} \rightarrow 0} \frac{\|H(V) - H(V_*)\|_{ui}}{\|\sin \Theta(V, V_*)\|_{ui}} \leq \xi_{ui}. \end{aligned}$$

In other words, the Lipschitz-like condition (3.1a) implies the assumption (A3) with  $q = 1$  and  $\chi = \xi_{ui}$ .

(d) From the proof of Theorem 4.1, we can see that the assumption (A3) can be relaxed to

$$(4.14) \quad \limsup_{i \rightarrow \infty} \frac{\|(I - P_*)[H(V_i) - H(V_*)] P_*\|_{ui}}{\|\sin \Theta(V_i, V_*)\|_{ui}} \leq \chi,$$

instead for all  $V \in \mathbb{U}^{n \times k}$  that go to  $V_*$ .

*Example 1.* We give an example to show that the local convergence rate revealed in Theorem 4.1 is nearly achievable, which implies its near-optimality. Consider the following single-particle Hamiltonian in electronic structure calculations studied in [14, 27, 32]:

$$(4.15) \quad H(V) = L + \alpha \cdot \text{Diag}(L^{-1}\rho(V)),$$

where  $L = \text{tridiag}(-1, 2, -1)$  is a discrete 1-D Laplacian,  $\alpha$  is some real constant,  $\rho(V) = \text{diag}(VV^T)$ , and  $V^T V = I_k$ . This is a good place for us to point out again that all our developments in this paper are valid for real NEPv (1.1), i.e.,  $H(V)$  is an  $n \times n$  real symmetric matrix-valued function of real  $V$  with  $V^T V = I_k$ , and at the same time  $Q$  in (1.2) is restricted to any orthogonal matrix.

To numerically demonstrate the local convergence rate, we use the following approach to compute an estimated convergence rate and the corresponding observed convergence rate for solving NEPv (1.1) with  $H(V)$  here by the SCF iteration (Algorithm 1). For a given  $\alpha$ , we compute an “exact” solution  $\widehat{V}_*$  by setting a small tolerance of  $10^{-14}$  in SCF. At convergence, the eigenvalue gap  $\delta_*$  of the assumption (A1) is estimated by

$$\widehat{\delta}_* = \lambda_{k+1}(H(\widehat{V}_*)) - \lambda_k(H(\widehat{V}_*)).$$

We approximate the quantity  $\chi$  in the assumption (A3) by using the quantity

$$\frac{\|(I - \widehat{P}_*)[H(V_i) - H(\widehat{V}_*)]\widehat{P}_*\|_2}{\|\sin \Theta(V_i, \widehat{V}_*)\|_2}$$

near the end of the SCF iteration when it stays almost unchanged at a constant  $\widehat{\chi}$ , where  $\widehat{P}_* = \widehat{V}_* \widehat{V}_*^T$ . Consequently, an estimated convergence rate in Theorem 4.1 is the quantity  $\widehat{\chi}/\widehat{\delta}_*$ . The corresponding observed convergence rate  $\widehat{\tau}$  is the numerical limit of the sequence

$$\widehat{\tau}_i = \frac{\|\sin \Theta(V_i, \widehat{V}_*)\|_2}{\|\sin \Theta(V_{i-1}, \widehat{V}_*)\|_2}.$$

Figure 1 shows the estimated convergence rates  $\widehat{\chi}/\widehat{\delta}_*$  and observed convergence rates  $\widehat{\tau}$  for different values of  $\alpha$ . We can see that the estimated convergence rates are *tight* upper bounds for the observed convergence rates for all tested values of  $\alpha$ . In particular, for  $\alpha = 0.05$ , the bound is essentially reached.

**4.3. Global convergence of SCF.** Our main results for the global convergence of the SCF iteration are based on the following two inequalities to be established:

$$(4.16) \quad \|\sin \Theta(V_i, V_{i+1})\|_2 \leq \tau_2 \|\sin \Theta(V_{i-1}, V_i)\|_2,$$

$$(4.17) \quad \|\sin \Theta(V_i, V_{i+1})\|_{\text{ui}} \leq \tau_{\text{ui}} \|\sin \Theta(V_{i-1}, V_i)\|_{\text{ui}}$$

for some constants  $\tau_2, \tau_{\text{ui}} < 1$  to be specified in the theorem below.

**THEOREM 4.2.** *Assume the Lipschitz-like conditions (3.1) hold and  $\xi_{\text{ui}}$  and  $\xi_2$  are the corresponding positive constants, and let  $\{V_i\}_i$  be generated by the SCF iteration (Algorithm 1). Suppose that there exists a positive constant  $\delta$  such that*

$$(4.18) \quad \lambda_{k+1}(H(V_i)) - \lambda_k(H(V_i)) \geq \delta \quad \text{for all } i = 1, 2, \dots$$

(a) *If  $\delta > \xi_{\text{ui}} + \xi_2$ , then the inequality (4.17) holds for all  $i$  with*

$$(4.19) \quad \tau_{\text{ui}} = \frac{\xi_{\text{ui}}}{\delta - \xi_2} < 1.$$

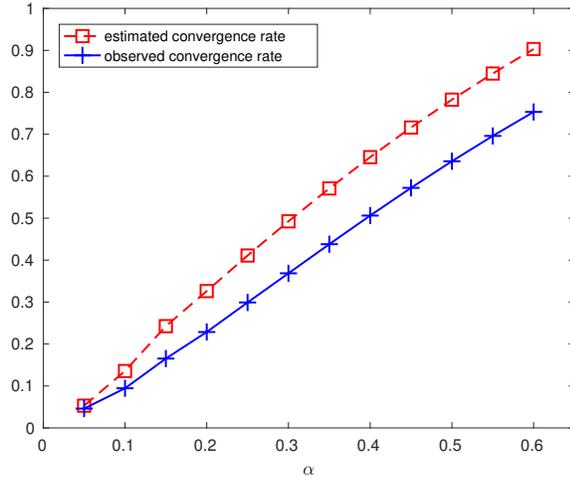


FIG. 1. Convergence rate: estimated  $\widehat{\chi}/\widehat{\delta}_*$  and observed  $\widehat{\tau}$  for different values of  $\alpha$  for the single-particle Hamiltonian  $H(V)$  defined in (4.15) with  $n = 10$  and  $k = 2$ .

(b) If  $\delta > \xi_2 + \|\sin \Theta(V_0, V_1)\|_2 \xi_2$ , then the inequality (4.16) holds for all  $i$  with

$$(4.20) \quad \tau_2 = \frac{\xi_2}{\delta - \xi_2 \|\sin \Theta(V_0, V_1)\|_2} < 1.$$

(c) If  $\delta > \max\{\xi_{ui}, \xi_2\} + \|\sin \Theta(V_0, V_1)\|_2 \xi_2$ , then both inequalities (4.16) and (4.17) hold for all  $i$  with

$$(4.21) \quad \tau_2 = \frac{\xi_2}{\delta - \xi_2 \|\sin \Theta(V_0, V_1)\|_2} < 1 \quad \text{and} \quad \tau_{ui} = \frac{\xi_{ui}}{\delta - \xi_2 \|\sin \Theta(V_0, V_1)\|_2} < 1.$$

These inequalities imply that, in their respective cases,  $\mathcal{R}(V_i)$  converges and that the limit, denoted by  $V_* \in \mathcal{G}_k(\mathbb{C}^n)$ , is the solution of NEPv (1.1). In other words,  $\sin \Theta(V_{i-1}, V_i) \rightarrow 0$  as  $i \rightarrow \infty$  and the SCF iteration is globally linearly convergent.

*Proof.* Let the eigendecomposition of  $H(V_i)$  be

$$H(V_i) = [V_{i+1}, V_{i+1,c}] \text{Diag}(\Lambda_{i+1}, \Lambda_{i+1,c}) [V_{i+1}, V_{i+1,c}]^H,$$

where  $[V_{i+1}, V_{i+1,c}] \in \mathbb{C}^{n \times n}$  is unitary,  $V_{i+1} \in \mathbb{U}^{n \times k}$ ,

$$\Lambda_{i+1} = \text{Diag}(\lambda_1(H(V_i)), \dots, \lambda_k(H(V_i))),$$

and

$$\Lambda_{i+1,c} = \text{Diag}(\lambda_{k+1}(H(V_i)), \dots, \lambda_n(H(V_i))).$$

For convenience, introduce  $\eta_i = \|\sin \Theta(V_{i-1}, V_i)\|_2$ . By Lemma 2.3 and (3.1), it holds that

$$(4.22) \quad |\lambda_j(H(V_{i-1})) - \lambda_j(H(V_i))| \leq \|H(V_{i-1}) - H(V_i)\|_2 \leq \xi_2 \eta_i$$

for all  $j$ . Combine (4.18) and (4.22) to get

$$(4.23) \quad \begin{aligned} \lambda_{k+1}(H(V_{i-1})) - \lambda_k(H(V_i)) &= \lambda_{k+1}(H(V_{i-1})) - \lambda_k(H(V_{i-1})) \\ &\quad + \lambda_k(H(V_{i-1})) - \lambda_k(H(V_i)) \\ &\geq \delta - \xi_2 \eta_i. \end{aligned}$$

Define  $R_i = H(V_{i-1})V_{i+1} - V_{i+1}\Lambda_{i+1}$ . We have

$$V_{i,c}^H R_i = \Lambda_{i,c} V_{i,c}^H V_{i+1} - V_{i,c}^H V_{i+1} \Lambda_{i+1}.$$

On the other hand, it can be verified that  $R_i = [H(V_{i-1}) - H(V_i)]V_{i+1}$ . Therefore,

$$(4.24) \quad \Lambda_{i,c} V_{i,c}^H V_{i+1} - V_{i,c}^H V_{i+1} \Lambda_{i+1} = V_{i,c}^H [H(V_{i-1}) - H(V_i)] V_{i+1}.$$

Apply Lemmas 2.1 and 2.4 to get, provided that we can prove  $\delta - \xi_2 \eta_i > 0$ ,

(4.25a)

$$(4.25b) \quad \begin{aligned} \|\sin \Theta(V_i, V_{i+1})\|_{\text{ui}} &= \|V_{i,c}^H V_{i+1}\|_{\text{ui}} \\ &\leq \frac{1}{\lambda_{k+1}(H(V_{i-1})) - \lambda_k(H(V_i))} \|V_{i,c}^H [H(V_{i-1}) - H(V_i)] V_{i+1}\|_{\text{ui}} \end{aligned}$$

$$(4.25c) \quad \begin{aligned} &\leq \frac{1}{\delta - \xi_2 \eta_i} \|H(V_{i-1}) - H(V_i)\|_{\text{ui}} \\ &\leq \frac{\xi_{\text{ui}}}{\delta - \xi_2 \eta_i} \|\sin \Theta(V_{i-1}, V_i)\|_{\text{ui}}, \end{aligned}$$

where we have used Lemma 2.1 for (4.25a) and Lemma 2.4 for (4.25b).

Item (a) is an immediate consequence of (4.25c) because all

$$\eta_i = \|\sin \Theta(V_{i-1}, V_i)\|_2 \leq 1,$$

and thus  $\delta - \xi_2 \eta_i \geq \delta - \xi_2 > 0$  by assumption.

For item (b), specialize (4.25c) to the case  $\|\cdot\|_{\text{ui}} = \|\cdot\|_2$  to get

$$(4.26) \quad \eta_{i+1} = \|\sin \Theta(V_i, V_{i+1})\|_2 \leq \frac{\xi_2}{\delta - \xi_2 \eta_i} \|\sin \Theta(V_{i-1}, V_i)\|_2 = \frac{\xi_2}{\delta - \xi_2 \eta_i} \eta_i.$$

By assumption,  $\delta > \xi_2 + \xi_2 \eta_1$ ,  $\delta - \xi_2 \eta_1 > \xi_2 > 0$ , and thus  $\eta_2 \leq \xi_2 \eta_1 / (\delta - \xi_2 \eta_1) < \eta_1$ . Now assume  $\eta_{i+1} < \eta_i$  for all  $i \leq \ell - 1$  ( $\ell \geq 2$ ). Using (4.26), we get

$$\eta_{\ell+1} \leq \frac{\xi_2 \eta_\ell}{\delta - \xi_2 \eta_\ell} < \frac{\xi_2 \eta_\ell}{\delta - \xi_2 \eta_1} < \eta_\ell.$$

Thus, by mathematical induction, we conclude that  $\eta_{i+1} < \eta_i$  for all  $i \geq 1$ . Finally, by (4.26), we obtain

$$(4.27) \quad \delta - \xi_2 \eta_i \geq \delta - \xi_2 \eta_1 > 0$$

and

$$\eta_{i+1} \leq \frac{\xi_2 \eta_i}{\delta - \xi_2 \eta_i} < \frac{\xi_2}{\delta - \xi_2 \eta_1} \eta_i = \tau_2 \eta_i,$$

where  $\tau_2$  is given by (4.20).

Finally, for item (c), the assumption on  $\delta$  is stronger than the one in item (b). Hence (4.16) holds for all  $i$  with  $\tau_2$  given by (4.20), which is the same as the one in (4.21). By combining (4.25) and (4.27), we get (4.17) with  $\tau_{\text{ui}}$  given in (4.21).  $\square$

Theorem 4.2 looks similar to Theorem 3.1 on the existence and uniqueness of the solution of NEPv (1.1). Both use the Lipschitz-like conditions in (3.1) but differ on gap assumptions between the  $k$ th and  $(k+1)$ st eigenvalues. Theorem 4.2 only requires the uniform gap assumption (4.18) with three different assumptions on the size of the

gap  $\delta$  on  $H(V_i)$  for all  $V_i$  generated by the SCF iteration, whereas the gap assumption in Theorem 3.1 is for all  $V \in \mathbb{U}^{n \times k}$ . This seems weaker, but it is not clear whether (4.18) is any easier to use than (3.2). Depending on how large  $\|\sin \Theta(V_0, V_1)\|_2$  is,  $\delta > \xi_{\text{ui}} + \xi_2$  needed for item (a) can be a significantly stronger assumption than the ones for items (b) and (c). Another difference is that under the conditions of Theorem 3.1, NEPv (1.1) has a unique solution, whereas for Theorem 4.2, it only guarantees that NEPv (1.1) has a solution which is the limit of  $\mathcal{R}(V_i)$ .

**5. Applications.** In this section, we apply the previous convergence analysis to the discretized Kohn–Sham NEPv in electronic structure calculations and the NEPv arising from linear discriminant analysis for dimension reduction. When applicable, we compare with the existing results. We note that both examples take the form (1.1) but for real numbers, i.e.,  $H(V) \in \mathbb{R}^{n \times n}$  are symmetric and  $V \in \mathbb{R}^{n \times k}$  has orthonormal columns. As we commented before, the general theory so far remains valid after replacing all  $\mathbb{C}$  by  $\mathbb{R}$  and  $\mathbb{U}^{n \times k}$  by  $\mathbb{O}^{n \times k} := \{V \mid V \in \mathbb{R}^{n \times k}, V^T V = I_k\}$ .

**5.1. The discretized Kohn–Sham NEPv.** Consider the following discretized Kohn–Sham NEPv studied in [17, 28, 14, 12, 21] and references therein:

$$(5.1a) \quad H(V)V = V\Lambda,$$

where the matrix-valued function

$$(5.1b) \quad H(V) = \frac{1}{2}L + V_{\text{ion}} + \sum_{\ell} w_{\ell} w_{\ell}^T + \text{Diag}(L^{\dagger} \rho) + \text{Diag}(\mu_{\text{xc}}^T(\rho) \mathbf{1})$$

is the plane-wave discretized Hamiltonian of the total energy functional. The first term corresponds to the kinetic energy, and  $L$  is a finite dimensional representation of the Laplacian operator. The second term  $V_{\text{ion}}$  is for the ionic pseudopotential sampled on the suitably chosen Cartesian grid in the local ionic potential energy. The third term represents a discretized pseudopotential reference projection function in the nonlocal ionic potential energy. The fourth term is for the Hartree potential energy, where  $\rho \equiv \rho(V) := \text{Diag}(V V^T) \in \mathbb{R}^n$  and  $L^{\dagger}$  is the pseudoinverse of  $L$ . The last term is for the exchange correlation energy, where  $\mu_{\text{xc}}(\rho) = \frac{\partial \epsilon_{\text{xc}}(\rho)}{\partial \rho} \in \mathbb{R}^{n \times n}$ ,  $\epsilon_{\text{xc}}(\rho)$  is an exchange correlation functional, and  $\mathbf{1}$  is a vector of all ones.

The discretized Kohn–Sham NEPv (5.1) is of the NEPv form (1.1). To apply the results in the previous sections, we first have to estimate how  $H(V)$  changes with respect to  $V$ . For this purpose, it suffices to know how  $\mu_{\text{xc}}(\rho)$  changes with respect to  $\rho \equiv \rho(V)$  since the first three terms in  $H(V)$  are independent of  $V$ . We assume that there exist positive constants  $\sigma_2$  and  $\sigma_{\text{F}}$  such that

$$(5.2a) \quad \|\text{Diag}(\mu_{\text{xc}}(\rho)^T \mathbf{1}) - \text{Diag}(\mu_{\text{xc}}(\tilde{\rho})^T \mathbf{1})\|_2 = \|[\mu_{\text{xc}}(\rho) - \mu_{\text{xc}}(\tilde{\rho})]^T \mathbf{1}\|_{\infty} \leq \sigma_2 \|\rho - \tilde{\rho}\|_{\infty}$$

and

$$(5.2b) \quad \|\text{Diag}(\mu_{\text{xc}}(\rho)^T \mathbf{1}) - \text{Diag}(\mu_{\text{xc}}(\tilde{\rho})^T \mathbf{1})\|_{\text{F}} = \|[\mu_{\text{xc}}(\rho) - \mu_{\text{xc}}(\tilde{\rho})]^T \mathbf{1}\|_2 \leq \sigma_{\text{F}} \|\rho - \tilde{\rho}\|_2$$

for all  $\rho \equiv \rho(V)$  and  $\tilde{\rho} \equiv \rho(\tilde{V})$ , where  $\|\cdot\|_{\infty}$  is either the  $\ell_{\infty}$ -norm of a vector or the  $\ell_{\infty}$ -norm of a matrix. With these assumptions (5.2), we can verify that  $H(V)$  satisfy

the Lipschitz-like conditions (3.1):

$$\begin{aligned}
 \|H(V) - H(\tilde{V})\|_2 &\leq \| \text{Diag}(L^\dagger(\rho - \tilde{\rho})) \|_2 + \| \text{Diag}(\mu_{\text{xc}}^\text{T}(\rho)\mathbf{1} - \mu_{\text{xc}}^\text{T}(\tilde{\rho})\mathbf{1}) \|_2 \\
 &= \|L^\dagger(\rho - \tilde{\rho})\|_\infty + \|[\mu_{\text{xc}}(\rho) - \mu_{\text{xc}}(\tilde{\rho})]^\text{T}\mathbf{1}\|_\infty \\
 &\leq \|L^\dagger\|_\infty \|\rho - \tilde{\rho}\|_\infty + \sigma_2 \|\rho - \tilde{\rho}\|_\infty \\
 &= (\|L^\dagger\|_\infty + \sigma_2) \max_i |e_i^\text{T}(VV^\text{T} - \tilde{V}\tilde{V}^\text{T})e_i| \\
 &\leq (\|L^\dagger\|_\infty + \sigma_2) \|\sin \Theta(V, \tilde{V})\|_2 \\
 (5.3a) \quad &\equiv \xi_2^{\text{ks}} \|\sin \Theta(V, \tilde{V})\|_2
 \end{aligned}$$

and

$$\begin{aligned}
 \|H(V) - H(\tilde{V})\|_\text{F} &\leq \| \text{Diag}(L^\dagger(\rho - \tilde{\rho})) \|_\text{F} + \| \text{Diag}(\mu_{\text{xc}}^\text{T}(\rho)\mathbf{1} - \mu_{\text{xc}}^\text{T}(\tilde{\rho})\mathbf{1}) \|_\text{F} \\
 &= \|L^\dagger(\rho - \tilde{\rho})\|_2 + \|[\mu_{\text{xc}}(\rho) - \mu_{\text{xc}}(\tilde{\rho})]^\text{T}\mathbf{1}\|_2 \\
 &\leq \|L^\dagger\|_2 \|\rho - \tilde{\rho}\|_2 + \sigma_\text{F} \|\rho - \tilde{\rho}\|_2 \\
 &\leq (\|L^\dagger\|_2 + \sigma_\text{F}) \|VV^\text{T} - \tilde{V}\tilde{V}^\text{T}\|_\text{F} \\
 (5.3b) \quad &\equiv \xi_\text{F}^{\text{ks}} \|\sin \Theta(V, \tilde{V})\|_\text{F},
 \end{aligned}$$

where  $\xi_2^{\text{ks}} = \|L^\dagger\|_\infty + \sigma_2$  and  $\xi_\text{F}^{\text{ks}} = \sqrt{2}(\|L^\dagger\|_2 + \sigma_\text{F})$ . By Theorem 3.1, we have the following sufficient condition on the existence and uniqueness of (5.1).

**THEOREM 5.1.** *Under the assumption (5.2), if for any  $V \in \mathbb{O}^{n \times k}$*

$$(5.4) \quad \lambda_{k+1}(H(V)) - \lambda_k(H(V)) > \min\{\xi_2^{\text{ks}} + \xi_\text{F}^{\text{ks}}, 2\xi_2^{\text{ks}}\},$$

*then the discretized Kohn–Sham NEPv (5.1) has a unique solution.*

It should be emphasized that for many realistic physical systems in electronic structure calculations, the eigenvalue gap condition (5.4) is not satisfied, and yet the solution to the discretized Kohn–Sham NEPv (5.1) does appear to be unique. It is a subject of further investigation.

Next we consider the convergence of the SCF iteration for solving the NEPv (5.1). For applying the local and global convergence results of Theorems 4.1 and 4.2, we note that the assumption (A3) in (4.3) becomes

$$(5.5) \quad \limsup_{\|\sin \Theta(V, V_*)\|_{\text{ui}} \rightarrow 0} \frac{\|(I - P_*)[\text{Diag}(L^\dagger[\rho - \rho_*]) + \text{Diag}([\mu_{\text{xc}}(\rho) - \mu_{\text{xc}}(\rho_*)]^\text{T}\mathbf{1})]P_*\|_{\text{ui}}}{\|\sin \Theta(V, V_*)\|_{\text{ui}}^q} \leq \chi,$$

where  $\rho_* := \rho(V_*)$ . Evidently, by Lipschitz-like conditions (5.3), we can take  $\chi = \xi_{\text{ui}}^{\text{ks}}$  and  $q = 1$  for  $\text{ui} \in \{2, \text{F}\}$ . The following theorem summarizes the local and global convergence of the SCF iteration.

**THEOREM 5.2.** *Let  $V_*$  be a solution of the discretized Kohn–Sham NEPv (5.1), and let  $\{V_i\}_i$  be the sequence generated by the SCF iteration (Algorithm 1). Then, under the assumption (5.2), the following hold:*

- (i) *If, for  $\text{ui} \in \{2, \text{F}\}$ ,  $\delta_* := \lambda_{k+1}(H(V_*)) - \lambda_k(H(V_*)) > \xi_{\text{ui}}^{\text{ks}}$  and  $\mathcal{R}(V_0)$  is sufficiently close to  $\mathcal{R}(V_*)$ , then  $\mathcal{R}(V_i)$  linearly converges to  $\mathcal{R}(V_*)$ . Moreover,*

$$(5.6) \quad \|\sin \Theta(V_{i+1}, V_*)\|_{\text{ui}} \leq \tau_i^{\text{ks}} \|\sin \Theta(V_i, V_*)\|_{\text{ui}},$$

*where  $\tau_i^{\text{ks}} < 1$  and  $\lim_{i \rightarrow \infty} \tau_i^{\text{ks}} = \frac{\xi_{\text{ui}}^{\text{ks}}}{\delta_*}$ .*

(ii) *Theorem 4.2 is valid for  $u_i = F$  and with  $\xi_2$  and  $\xi_{ui}$  replaced by  $\xi_2^{ks}$  and  $\xi_{ui}^{ks}$ , respectively.*

*Proof.* Result (i) follows from Theorem 4.1 since the subspace approximation condition (4.3) holds for the constants  $\chi = \xi_{ui}^{ks}$ . Result (ii) immediately follows from Theorem 4.2.  $\square$

Let us compare Theorem 5.2 with the previous convergence results of the SCF iteration on NEPv (5.1) obtained by Liu et al. [14]. We first restate the following main results of [14].

**THEOREM 5.3** (see [14, Theorem 4.2]). *For NEPv (5.1), suppose that there exists a constant  $\sigma$  such that for all  $\rho, \tilde{\rho} \in \mathbb{R}^n$ ,*

$$(5.7a) \quad \|\text{Diag}(\mu_{xc}(\rho)^T \mathbf{1}) - \text{Diag}(\mu_{xc}(\tilde{\rho})^T \mathbf{1})\|_F \leq \sigma \|\rho - \tilde{\rho}\|_2,$$

$$(5.7b) \quad \left\| \frac{\partial^2 \epsilon_{xc}}{\partial \rho^2} \mathbf{1} \right\|_2 \leq \sigma.$$

*Let  $\{V_i\}_i$  be the sequence generated by the SCF iteration (Algorithm 1), let  $V_*$  be a solution of NEPv (5.1), and let  $\delta_* = \lambda_{k+1}(H(V_*)) - \lambda_k(H(V_*)) > 0$ . If  $V_i$  is sufficiently close to  $V_*$ , i.e.,  $\|\sin \Theta(V_i, V_*)\|_2$  is sufficiently small, then*

$$(5.8) \quad \|\sin \Theta(V_{i+1}, V_*)\|_2 \leq \frac{2\sqrt{n}(\|L^\dagger\|_2 + \sigma)}{\delta_*} \|\sin \Theta(V_i, V_*)\|_2 + \mathcal{O}(\|\sin \Theta(V_i, V_*)\|_2^2).$$

**THEOREM 5.4** (see [14, Theorem 3.3]). *Assume (5.7), and suppose there exists a constant  $\delta > 0$  such that  $\lambda_{k+1}(H(V_i)) - \lambda_k(H(V_i)) \geq \delta > 0$  for all  $i$ , where  $\{V_i\}_i$  is the sequence generated by the SCF iteration (Algorithm 1). If  $\delta > 12k\sqrt{n}(\|L^\dagger\|_2 + \sigma)$ , then  $\mathcal{R}(V_i)$  converges to a solution  $\mathcal{R}(V_*)$  of NEPv (5.1).*

First we note that the assumption (5.7b) on the twice differentiability of the exchange correlation functional  $\epsilon_{xc}$  is not necessary in Theorem 5.2. On the local convergence, Theorem 5.3 requires the eigenvalue gap  $\delta_* > 2\sqrt{n}(\|L^\dagger\|_2 + \sigma)$ . In contrast, for  $u_i = 2$ , Theorem 5.2(i) only requires  $\delta_* > \|L^\dagger\|_\infty + \sigma_2$ , which is a much weaker condition. This can be verified as follows. By the assumption (5.2) of Theorem 5.2(i), let

$$\hat{\sigma}_2 = \sup_{\rho \neq \tilde{\rho}} \frac{\|[\mu_{xc}(\rho) - \mu_{xc}(\tilde{\rho})]^T \mathbf{1}\|_\infty}{\|\rho - \tilde{\rho}\|_\infty} \quad \text{and} \quad \hat{\sigma}_F = \sup_{\rho \neq \tilde{\rho}} \frac{\|[\mu_{xc}(\rho) - \mu_{xc}(\tilde{\rho})]^T \mathbf{1}\|_2}{\|\rho - \tilde{\rho}\|_2}.$$

Then  $\frac{1}{\sqrt{n}}\hat{\sigma}_F \leq \hat{\sigma}_2 \leq \sqrt{n}\hat{\sigma}_F$  and  $\hat{\sigma}_F \leq \sigma$ . Since Theorem 5.2 also holds for  $\sigma_2 = \hat{\sigma}_2$  and  $\sigma_F = \hat{\sigma}_F$ , we have

$$(5.9) \quad \|L^\dagger\|_\infty + \hat{\sigma}_2 \leq \sqrt{n}(\|L^\dagger\|_2 + \hat{\sigma}_F) < 2\sqrt{n}(\|L^\dagger\|_2 + \sigma).$$

Therefore, Theorem 5.2(i) has a weaker condition on the eigenvalue gap  $\delta_*$  than the one required by Theorem 5.3. In fact, since the first inequality in (5.9) is overestimated, Theorem 5.2(i) has a significantly weaker condition on the eigenvalue gap by removing the factor  $\sqrt{n}$ . By the inequalities (5.9), we can also see that the new bound (5.6) on  $\|\sin \Theta(V_{i+1}, V_*)\|_2$  is much sharper than the first-order bound (5.8) of Theorem 5.3. In addition, we note that for  $u_i = F$ , Theorem 5.2(i) provides the convergence rate  $\xi_F^{ks}/\delta_*$  of the SCF iteration, which is absent in [14].

On the global convergence, the condition  $\delta > 12k\sqrt{n}(\|L^\dagger\|_2 + \sigma)$  in Theorem 5.4 is a much more stringent condition than the one required by Theorem 5.2(ii). This is due to the fact that

$$(5.10) \quad \xi_{\text{ui}}^{\text{ks}} + \xi_2^{\text{ks}} \leq (\sqrt{n} + 1)(\|L^\dagger\|_2 + \hat{\sigma}_F) < 12k\sqrt{n}(\|L^\dagger\|_2 + \sigma).$$

Now let us examine the implications of these results for the simple single-particle Hamiltonian (4.15) with the nonlinearity controlled by the parameter  $\alpha$ . To ensure the local convergence of the SCF iteration, the sufficient condition from the analysis in [14] is that the parameter  $\alpha$  must satisfy

$$(5.11) \quad \alpha < \alpha_L := \frac{\delta_*}{2\sqrt{n}\|L^\dagger\|_2}.$$

In contrast, by Theorem 5.2(ii), the upper bound is

$$(5.12) \quad \alpha < \tilde{\alpha}_L := \max \left\{ \frac{\delta_*}{\|L^\dagger\|_1}, \frac{\delta_*}{\sqrt{2}\|L^\dagger\|_2} \right\}.$$

Since  $\|L^\dagger\|_1 \leq \sqrt{n}\|L^\dagger\|_2$ ,<sup>1</sup>  $\tilde{\alpha}_L$  is always larger than  $\alpha_L$  and does not explicitly depend on  $n$  (though  $\|L^\dagger\|_2$  depends on  $n$ ), it implies that the sufficient condition (5.12) is less stringent than (5.11). This is also confirmed by numerical results seen in Example 1 and [27, Table 1].

For the global convergence, an earlier one of Yang, Gao, and Meza [27] requires

$$(5.13) \quad \alpha < \alpha_F := \frac{\delta}{\ln \frac{1-\gamma}{\gamma} \cdot n^4 \|L^\dagger\|_1},$$

where  $\gamma$  is a constant and  $\gamma \ll 1$ , and  $\delta$  is the one as in Theorems 4.2 and 5.4. Liu et al. [14] improved the upper bound (5.13) to

$$(5.14) \quad \alpha < \alpha_G := \frac{\delta}{12k\sqrt{n}\|L^\dagger\|_2}.$$

In contrast, the result in Theorem 5.2(ii) requires

$$(5.15) \quad \alpha < \tilde{\alpha}_G := \max \left\{ \frac{\delta}{(1 + \|\sin \Theta(V_0, V_1)\|_2)\|L^\dagger\|_1}, \frac{\delta}{\|L^\dagger\|_1 + \sqrt{2}\|L^\dagger\|_2} \right\}.$$

As we can see, unlike the previous bounds,  $\alpha_F$  and  $\alpha_G$ ,  $\tilde{\alpha}_G$  does not explicitly depend on  $n$ . Furthermore,  $\tilde{\alpha}_G$  is always larger than  $\alpha_G$ , which in turn is larger than  $\alpha_F$  for

$$\gamma < \left[ 1 + \exp \left( \frac{12k}{n^{7/2}} \cdot \frac{\|L^\dagger\|_2}{\|L^\dagger\|_1} \right) \right]^{-1},$$

i.e.  $\tilde{\alpha}_G > \alpha_G > \alpha_F$ . This means the result (5.15) predicts a much larger range of  $\alpha$  than (5.13) and (5.14) could, within which the SCF iteration converges. This is again confirmed by numerical experiments reported in Example 1 and [27, Table 1].

---

<sup>1</sup>Numerical observation suggests that  $\|L^\dagger\|_1/\|L^\dagger\|_2 \leq 1.7072$ . However, we do not have a rigorous proof.

**5.2. The trace ratio problem.** In this section, we discuss an application to a trace ratio maximization problem (TRP) arising from the linear discriminant analysis for dimension reduction [16, 30, 31]. Given symmetric matrices  $A, B \in \mathbb{R}^{n \times n}$  and  $B \succ 0$  (positive definite), TRP solves the following optimization problem:

$$(5.16) \quad \max_{V \in \mathbb{O}^{n \times k}} \frac{\text{tr}(V^T A V)}{\text{tr}(V^T B V)},$$

where  $\text{tr}(\cdot)$  denotes the trace of a square matrix. Employing the first-order optimality condition (i.e., the KKT condition) yields that any critical point  $V \in \mathbb{O}^{n \times k}$  of (5.16) is a solution of the following NEP<sub>v</sub>:

$$(5.17a) \quad H(V)V = V\Lambda,$$

where

$$(5.17b) \quad H(V) = A - \psi(V)B \in \mathbb{R}^{n \times n} \quad \text{and} \quad \psi(V) = \frac{\text{tr}(V^T A V)}{\text{tr}(V^T B V)}.$$

Necessarily,  $\Lambda = V^T H(V)V \in \mathbb{R}^{k \times k}$  is symmetric. Evidently,  $H(VQ) \equiv H(V)$  for any orthogonal  $Q \in \mathbb{R}^{k \times k}$ . NEP<sub>v</sub> (5.17) takes the form of NEP<sub>v</sub> (1.1). We have the following theorem that characterizes the relation between any global solution  $V_*$  of (5.16) and solutions of NEP<sub>v</sub> (5.17).

**THEOREM 5.5** (see [31, Theorem 2.1]).  *$V \in \mathbb{O}^{n \times k}$  is a global maximizer of (5.16) if and only if it solves NEP<sub>v</sub> (5.17) and the eigenvalues of  $\Lambda \equiv V^T H(V)V$  correspond to the  $k$  largest eigenvalues of  $H(V)$ .*

Theorem 5.5 transforms TRP (5.16) into NEP<sub>v</sub> (5.17), and naturally it leads to an SCF iteration for finding the desired solution. The SCF iteration is the same as Algorithm 1, except a simple modification of  $\Lambda_i$  at line 3 to

$$\Lambda_i = \text{Diag}(\lambda_{n-k+1}(H_i), \dots, \lambda_n(H_i)),$$

namely consisting of the  $k$  largest eigenvalues of  $H_i$ .

In [30, Theorem 5.1], it is shown that an SCF iteration is globally convergent to a global maximizer  $V_*$  of (5.16) for any given initial guess  $V_0$ . In what follows, we will apply the convergence results in section 4 to estimate the local convergence rate of the SCF iteration. To that end, we need to establish the assumption (A3) at the beginning of section 4.2. The next lemma is similar to [30, Theorem 5.2] but with tighter constants.

**LEMMA 5.6.** *Let  $V_* \in \mathbb{O}^{n \times k}$  solve NEP<sub>v</sub> (5.17). For any  $V \in \mathbb{O}^{n \times k}$ , we have*

$$(5.18a) \quad |\psi(V_*) - \psi(V)| \leq \kappa_F \|\sin \Theta(V, V_*)\|_F^2,$$

$$(5.18b) \quad |\psi(V_*) - \psi(V)| \leq \kappa_2 \|\sin \Theta(V, V_*)\|_2^2,$$

where

$$(5.19) \quad \kappa_F = \frac{\lambda_n(H(V_*)) - \lambda_1(H(V_*))}{\sum_{i=1}^k \lambda_i(B)} \quad \text{and} \quad \kappa_2 = \frac{\sum_{i=1}^k [\lambda_{n-i+1}(H(V_*)) - \lambda_i(H(V_*))]}{\sum_{i=1}^k \lambda_i(B)}.$$

*Proof.* We note that  $\mathcal{R}(V_*)$  is the invariant subspace of  $H(V_*)$  and

$$\text{tr}(V_*^T H(V_*)V_*) = 0,$$

and thus  $\text{tr}(V^T H(V_*)V) = \text{tr}(V^T H(V_*)V) - \text{tr}(V_*^T H(V_*)V_*)$  for any  $V \in \mathbb{O}^{n \times k}$ . Viewing  $\mathcal{R}(V)$  as an approximate invariant subspace of  $H(V_*)$ , by [11, Theorem 2.2] (see also [24, item 2 of Theorem 3.1]), we have that

$$\begin{aligned} 0 &\leq \sum_{i=1}^k (\lambda_i(V^T H(V_*)V) - \lambda_i(V_*^T H(V_*)V_*)) \\ &\leq [\lambda_n(H(V_*)) - \lambda_1(H(V_*))] \sum_{i=1}^k \sin \theta_i^2(V, V_*) \end{aligned}$$

and

$$\begin{aligned} 0 &\leq \sum_{i=1}^k (\lambda_i(V^T H(V_*)V) - \lambda_i(V_*^T H(V_*)V_*)) \\ &\leq \sum_{i=1}^k [\lambda_{n-i+1}(H(V_*)) - \lambda_i(H(V_*))] \sin \theta_i^2(V, V_*). \end{aligned}$$

Consequently, we have

$$\begin{aligned} (5.20a) \quad |\text{tr}(V^T H(V_*)V)| &= |\text{tr}(V^T H(V_*)V) - \text{tr}(V_*^T H(V_*)V_*)| \\ &\leq [\lambda_n(H(V_*)) - \lambda_1(H(V_*))] \|\sin \Theta(V, V_*)\|_F^2 \end{aligned}$$

and

$$(5.20b) \quad |\text{tr}(V^T H(V_*)V)| \leq \left( \sum_{i=1}^k [\lambda_{n-i+1}(H(V_*)) - \lambda_i(H(V_*))] \right) \|\sin \Theta(V, V_*)\|_2^2,$$

respectively. On the other hand, since

$$|\text{tr}(V^T H(V_*)V)| = |\text{tr}(V^T AV) - \psi(V_*) \text{tr}(V^T BV)|,$$

we get

$$|\psi(V_*) - \psi(V)| = \frac{|\text{tr}(V^T H(V_*)V)|}{\text{tr}(V^T BV)} \leq \frac{|\text{tr}(V^T H(V_*)V)|}{\sum_{i=1}^k \lambda_i(B)},$$

which, combined with (5.20), yields (5.18). □

*Remark 3.* We remark that the constant  $\kappa_2$  in (5.18b) for the 2-norm case is smaller than the following constant given in [30, Theorem 5.2]:

$$k \left( \max_{i=1, \dots, k} |\lambda_{n-i+1}(H(V_*))| + \max_{i=1, \dots, k} |\lambda_i(H(V_*))| \right) / \left( \sum_{i=1}^k \lambda_i(B) \right)$$

because for any  $i = 1, 2, \dots, k$ ,

$$0 < \lambda_{n-i+1}(H(V_*)) - \lambda_i(H(V_*)) \leq \max_{i=1, \dots, k} |\lambda_{n-i+1}(H(V_*))| + \max_{i=1, \dots, k} |\lambda_i(H(V_*))|.$$

We now are able to establish the quadratic convergence of the SCF iteration for NEP<sub>V</sub> (5.17).

THEOREM 5.7. Let  $A, B \in \mathbb{R}^{n \times n}$  be symmetric with  $B \succ 0$ , and let  $V_*$  be any global solution to TRP (5.16). If  $\delta_* := \lambda_{n-k+1}(H(V_*)) - \lambda_{n-k}(H(V_*)) > 0$ , then for any given  $V_0 \in \mathbb{O}^{n \times k}$ , the SCF iteration (Algorithm 1 with a modification  $\Lambda_i = \text{Diag}(\lambda_{n-k+1}(H_i), \dots, \lambda_n(H_i))$  at line 3) converges quadratically to  $\mathcal{R}(V_*)$ . Moreover,

$$(5.21) \quad \limsup_{i \rightarrow \infty} \frac{\|\sin \Theta(V_i, V_*)\|_{\text{ui}}}{\|\sin \Theta(V_{i-1}, V_*)\|_{\text{ui}}^2} \leq \frac{\chi_{\text{ui}}}{\delta_*},$$

where  $\text{ui} \in \{2, \mathbb{F}\}$ ,  $\chi_{\text{ui}} = \kappa_{\text{ui}} \|B\|_{\text{ui}}$  with  $\kappa_{\text{ui}}$  given by (5.19).

Proof. Note by (5.17b) that  $\|H(V) - H(V_*)\|_{\text{ui}} = |\psi(V) - \psi(V_*)| \cdot \|B\|_{\text{ui}}$ . Thus, for  $\text{ui} \in \{2, \mathbb{F}\}$ , we can use Lemma 5.6 to obtain

$$\|H(V) - H(V_*)\|_{\text{ui}} = |\psi(V) - \psi(V_*)| \cdot \|B\|_{\text{ui}} \leq \kappa_{\text{ui}} \|B\|_{\text{ui}} \cdot \|\sin \Theta(V, V_*)\|_{\text{ui}}^2.$$

Consequently, the assumption (A3) of Theorem 4.1 for local convergence is satisfied for  $q = 2$  due to the fact that

$$(5.22) \quad \begin{aligned} \limsup_{\|\sin \Theta(V, V_*)\|_{\text{ui}} \rightarrow 0} \frac{\|(I - P_*)[H(V) - H(V_*)]P_*\|_{\text{ui}}}{\|\sin \Theta(V, V_*)\|_{\text{ui}}^2} \\ \leq \limsup_{\|\sin \Theta(V, V_*)\|_{\text{ui}} \rightarrow 0} \frac{\|H(V) - H(V_*)\|_{\text{ui}}}{\|\sin \Theta(V, V_*)\|_{\text{ui}}^2} \\ \leq \kappa_{\text{ui}} \|B\|_{\text{ui}} = \chi_{\text{ui}}. \end{aligned}$$

Thus, by Theorem 4.1, the locally quadratic convergence of the SCF iteration immediately follows under the assumption of the eigenvalue gap  $\delta_* = \lambda_{n-k+1}(H(V_*)) - \lambda_{n-k}(H(V_*)) > 0$ .  $\square$

Example 2. We present an example to demonstrate the local quadratic convergence revealed in Theorem 5.7. Let  $A = Z \text{Diag}(1, 2, \dots, n)Z$  and  $B = L_m \otimes I_m + I_m \otimes L_m + \alpha I_n$ , where  $Z = I_n - \mathbf{211}^T/n$  is a Householder matrix,  $\mathbf{1}$  is a vector of all ones,  $B \in \mathbb{R}^{n \times n}$  is a regularized standard 2-D discrete Laplacian on the unit square based upon a 5-point stencil with equally spaced mesh points, and  $L_m = \text{tridiag}(-1, 2, -1)$ .  $\alpha > 0$  is a regularization parameter, usually determined by the cross-validation technique over a prescribed set [30].

To numerically demonstrate the quadratic convergence rate, we take  $n = 400$  (i.e.,  $m = 20$ ),  $k = 10$ , and  $\alpha = \alpha(t) = 2t$  for  $t = 0, 1, \dots, 10$ . For each  $\alpha(t)$ , the SCF iteration starts with  $V_0 = [e_1, \dots, e_k]$  and terminates and returns  $\widehat{V}_*$  whenever  $\text{NRes}_i$  in (4.1) is no larger than  $10^{-14}$ .  $\widehat{V}_*$  is then treated as an exact solution. The observed quadratic rate is taken to be

$$\widehat{\tau} = \frac{\|\sin \Theta(V_i, \widehat{V}_*)\|_2}{\|\sin \Theta(V_{i-1}, \widehat{V}_*)\|_2^2}$$

for  $i$  near the end of the SCF iteration. Correspondingly, the estimated quadratic rate is taken to be  $\widehat{\chi}/\widehat{\delta}_*$ , where

$$\widehat{\chi} = \frac{\|(I - \widehat{P}_*)[H(V_i) - H(\widehat{V}_*)]\widehat{P}_*\|_2}{\|\sin \Theta(V_i, \widehat{V}_*)\|_2^2} \quad \text{and} \quad \widehat{\delta}_* = \lambda_{n-k+1}(H(\widehat{V}_*)) - \lambda_{n-k}(H(\widehat{V}_*))$$

for  $i$  near the end of the SCF iteration. Figure 2 shows both  $\widehat{\tau}$  and  $\widehat{\chi}/\widehat{\delta}_*$  for different values of  $\alpha(t)$ . As we can see, the estimated quadratic convergence rates are generally tight as upper bounds for the observed ones.

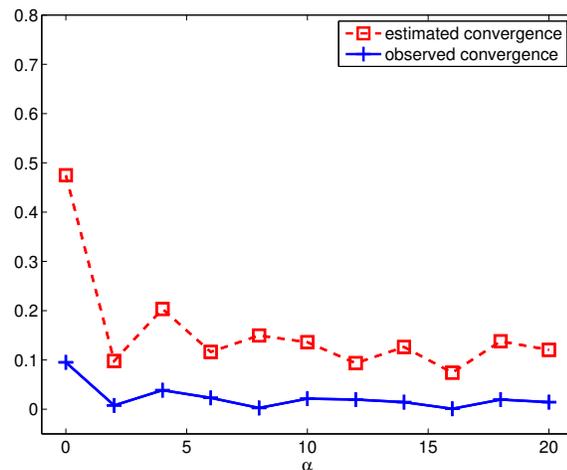


FIG. 2. Estimated quadratic convergence rate  $\hat{\chi}/\hat{\delta}_*$  and observed rate  $\hat{\tau}$  of the SCF iteration for solving NEPv (5.17) as  $\alpha = \alpha(t)$  varies with  $t = 0, 1, \dots, 10$ .

**6. Concluding remarks.** In this paper, we identified two sufficient conditions for the existence and uniqueness of the NEPv (1.1), namely Lipschitz-like conditions (3.1) and a uniform eigenvalue gap condition (3.2). The latter is undoubtedly strong and may be hard to verify in general unless the coefficient matrix  $H(V)$  is very special, such as the one for the Hartree–Fock integro-differential equations by Cancès and Le Bris [4].

Throughout the paper, we have assumed (1.2), i.e.,  $H(V) \equiv H(VQ)$  for any unitary  $Q \in \mathbb{C}^{k \times k}$  which makes  $H(V)$  a matrix-valued function on the Grassmann manifold of  $k$ -dimensional subspaces. As a result, Lipschitz-like conditions and the convergence results of the SCF iteration are stated in terms of the sine of the canonical angles between the subspaces. Looking beyond (1.2), we point out that most of our developments can still be adapted to the situations where (1.2) is no longer true. Possible modifications include, in general, replacing all  $\|\sin \Theta(V, \tilde{V})\|$  by  $\|V - \tilde{V}\|$ .

We presented local and global convergence analysis for the plain SCF iteration (Algorithm 1) for solving NEPv (1.1) and showed their applications to discrete Kohn–Sham NEPv (5.1) and the TRP (5.16). For these applications, we are able to demonstrate the near-optimality of the convergence rates revealed in this paper. Furthermore, for the instance of the Kohn–Sham problem (5.1), we have significantly improved the previous results in [27, 14]. Our analysis so far has been on the plain SCF iteration, i.e., without incorporating any accelerated schemes, such as the ones in [18, 25]. It would be an interesting topic to examine whether our analysis can be carried over to those accelerated SCF iterations.

**Acknowledgments.** We are grateful to the anonymous referees and the associate editor for their insightful comments and suggestions that significantly improved this paper. References [12, 21] were brought to our attention by one of the referees.

#### REFERENCES

- [1] W. BAO AND Q. DU, *Computing the ground state solution of Bose–Einstein condensates by a normalized gradient flow*, SIAM J. Sci. Comput., 25 (2004), pp. 1674–1697, <https://doi.org/10.1137/S1064827503422956>.

- [2] R. BHATIA, *Matrix Analysis*, Springer, New York, 1996.
- [3] Y. CAI, Z. BAI, J. E. PASK, AND N. SUKUMAR, *Hybrid preconditioning for iterative diagonalization of ill-conditioned generalized eigenvalue problems in electronic structure calculations*, *J. Comput. Phys.*, 255 (2013), pp. 16–30.
- [4] E. CANCÈS AND C. LE BRIS, *On the convergence of SCF algorithms for the Hartree–Fock equations*, *M2AN Math. Model. Numer. Anal.*, 34 (2000), pp. 749–774.
- [5] H. CHEN, X. DAI, X. GONG, L. HE, AND A. ZHOU, *Adaptive finite element approximations for Kohn–Sham models*, *Multiscale Model. Simul.*, 12 (2014), pp. 1828–1869, <https://doi.org/10.1137/130916096>.
- [6] C. DAVIS AND W. M. KAHAN, *The rotation of eigenvectors by a perturbation. III*, *SIAM J. Numer. Anal.*, 7 (1970), pp. 1–46, <https://doi.org/10.1137/0707001>.
- [7] E. JARLEBRING, S. KVAAL, AND W. MICHELIS, *An inverse iteration method for eigenvalue problems with eigenvector nonlinearities*, *SIAM J. Sci. Comput.*, 36 (2014), pp. A1978–A2001, <https://doi.org/10.1137/130910014>.
- [8] S.-H. JIA, H.-H. XIE, M.-T. XIE, AND F. XU, *A full multigrid method for nonlinear eigenvalue problems*, *Sci. China Math.*, 59 (2016), pp. 2037–2048.
- [9] L. JOST, S. SETZER, AND M. HEIN, *Nonlinear eigenproblems in data analysis: Balanced graph cuts and the RatioDCA-Prox*, in *Extraction of Quantifiable Information from Complex Systems*, Lect. Notes Comput. Sci. Eng. 102, S. Dahlke, W. Dahmen, M. Griebel, W. Hackbusch, R. Schneider, C. Schwab, and H. Yserentant, eds., Springer, Berlin, 2014, pp. 263–279.
- [10] M. A. KHAMSI, *Introduction to metric fixed point theory*, in *Topics in Fixed Point Theory*, S. Almezal, Q. H. Ansari, and M. A. Khamsi, eds., Springer, Cham, 2014, pp. 1–32.
- [11] A. V. KNYAZEV AND M. E. ARGENTATI, *Rayleigh–Ritz majorization error bounds with applications to FEM*, *SIAM J. Matrix Anal. Appl.*, 31 (2010), pp. 1521–1537, <https://doi.org/10.1137/08072574X>.
- [12] A. LEVITT, *Convergence of gradient-based algorithms for the Hartree–Fock equations*, *ESAIM Math Model Numer. Anal.*, 46 (2012), pp. 1321–1336.
- [13] R.-C. LI, *Matrix perturbation theory*, in *Handbook of Linear Algebra*, 2nd ed., L. Hogben, ed., CRC Press, Boca Raton, FL, 2014, pp. 21–1–21–30.
- [14] 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>.
- [15] R. M. MARTIN, *Electronic Structure: Basic Theory and Practical Methods*, Cambridge University Press, Cambridge, UK, 2004.
- [16] T. T. NGO, M. BELLALIJ, AND Y. SAAD, *The trace ratio optimization problem for dimensionality reduction*, *SIAM J. Matrix Anal. Appl.*, 31 (2010), pp. 2950–2971, <https://doi.org/10.1137/090776603>.
- [17] M. C. PAYNE, M. P. TETER, D. C. ALLEN, T. A. ARIAS, AND J. D. JOANNOPOULOS, *Iterative minimization techniques for ab initio total energy calculation: Molecular dynamics and conjugate gradients*, *Rev. Modern Phys.*, 64 (1992), pp. 1045–1097.
- [18] P. PULAY, *Improved SCF convergence*, *J. Comput. Chem.*, 3 (1982), pp. 556–560.
- [19] 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>.
- [20] 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>.
- [21] R. E. STANTON, *Intrinsic convergence in closed shell SCF calculations. A general criterion*, *J. Comput. Phys.*, 75 (1981), pp. 5416–5422.
- [22] G. W. STEWART AND J.-G. SUN, *Matrix Perturbation Theory*, Academic Press, Boston, 1990.
- [23] J.-G. SUN, *Matrix Perturbation Analysis*, Academic Press, Beijing, 1987 (in Chinese).
- [24] Z.-M. TENG, L.-Z. LU, AND R.-C. LI, *Cluster-robust accuracy bounds for Ritz subspaces*, *Linear Algebra Appl.*, 480 (2015), pp. 11–26.
- [25] H. F. WALKER AND P. NI, *Anderson acceleration for fixed-point iterations*, *SIAM J. Numer. Anal.*, 49 (2011), pp. 1715–1735, <https://doi.org/10.1137/10078356X>.
- [26] P.-Å. WEDIN, *On angles between subspaces*, in *Matrix Pencils*, B. Kågström and A. Ruhe, eds., Springer, New York, 1983, pp. 263–285.
- [27] 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>.
- [28] C. YANG, J. C. MEZA, AND L.-W. WANG, *A trust region direct constrained minimization algorithm for the Kohn–Sham equation*, *SIAM J. Sci. Comput.*, 29 (2007), pp. 1854–1875,

- <https://doi.org/10.1137/060661442>.
- [29] 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.
  - [30] 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>.
  - [31] L.-H. ZHANG, W. YANG, AND L.-Z. LIAO, *A note on the trace quotient problem*, Optim. Lett., 8 (2014), pp. 1637–1645.
  - [32] Z. ZHAO, Z.-J. BAI, AND X.-Q. JIN, *A Riemannian Newton algorithm for nonlinear eigenvalue problems*, SIAM J. Matrix Anal. Appl., 36 (2015), pp. 752–774, <https://doi.org/10.1137/140967994>.