

# A self-consistent-field iteration for MAXBET with an application to multi-view feature extraction

Xijun Ma<sup>1</sup> · Chungen Shen<sup>2</sup> · Li Wang<sup>3</sup> · Lei-Hong Zhang<sup>4,1</sup> · Ren-Cang Li<sup>3,5</sup>

Received: 24 March 2021 / Accepted: 1 February 2022

© The Author(s), under exclusive licence to Springer Science+Business Media, LLC, part of Springer Nature 2022

#### Abstract

As an extension of the traditional principal component analysis, the multi-view canonical correlation analysis (MCCA) aims at reducing m high dimensional random variables  $\mathbf{s}_i \in \mathbb{R}^{n_i}$  (i = 1, 2, ..., m) by proper projection matrices  $X_i \in \mathbb{R}^{n_i \times \ell}$  so that the m reduced ones  $\mathbf{y}_i = X_i^{\mathrm{T}} \mathbf{s}_i \in \mathbb{R}^{\ell}$  have the "maximal correlation." Various measures of the correlation for  $y_i$  (i = 1, 2, ..., m) in MCCA have been proposed. One of the earliest criteria is the sum of all traces of pair-wise correlation matrices between  $y_i$  and  $y_i$  subject to the orthogonality constraints on  $X_i$ , i = 1, 2, ..., m. The resulting problem is to maximize a homogeneous quadratic function over the product of Stiefel manifolds and is referred to as the MAXBET problem. In this paper, the problem is first reformulated as a coupled nonlinear eigenvalue problem with eigenvector dependency (NEPv) and then solved by a novel self-consistent-field (SCF) iteration. Global and local convergences of the SCF iteration are studied and proven computational techniques in the standard eigenvalue problem are incorporated to yield more practical implementations. Besides the preliminary numerical evaluations on various types of synthetic problems, the efficiency of the SCF iteration is also demonstrated in an application to multi-view feature extraction for unsupervised learning.

**Keywords** MAXBET · Multi-view canonical correlation analysis · Multi-view feature extraction · Nonlinear eigenvalue problem · Self-consistent-field iteration · Stiefel manifold

Mathematics Subject Classification (2010)  $90C20 \cdot 90C06 \cdot 65F10 \cdot 65F15 \cdot 65F35$ 

Communicated by: Raymond H. Chan

Published online: 16 March 2022

□ Lei-Hong Zhang longzlh@suda.edu.cn

Extended author information available on the last page of the article.



13 Page 2 of 34 Adv Comput Math (2022) 48:13

## 1 Introduction

In this paper, we are concerned with the maximization of a homogeneous quadratic function  $f(X) = tr(X^T A X)$  over the Cartesian product of m Stiefel manifolds:

$$\max_{X \in \mathcal{M}} \left\{ f(X) := \operatorname{tr}(X^{\mathrm{T}} A X) \right\},\tag{1}$$

where  $A \in \mathbb{R}^{n \times n}$  is symmetric and  $tr(X^T A X)$  is the trace of  $X^T A X$ ,

$$\mathcal{M} := \mathbb{O}^{n_1 \times \ell} \times \ldots \times \mathbb{O}^{n_m \times \ell}$$

$$= \left\{ X = [X_1; \ldots; X_m] \in \mathbb{R}^{n \times \ell} \mid X_i^{\mathsf{T}} X_i = I_{\ell}, \ X_i \in \mathbb{R}^{n_i \times \ell} \right\},\,$$

 $\sum n_i = n$ , and the Stiefel manifold

$$\mathbb{O}^{n_i \times \ell} := \{ X_i \in \mathbb{R}^{n_i \times \ell} \mid X_i^{\mathsf{T}} X_i = I_\ell \}.$$

Necessarily,  $\ell \le n_i$  for all i. Here and in the rest of this paper, we will use notation  $[X_1; \ldots; X_m]$  as a matrix/vector constructor that stacks up the blocks  $X_i$  in order (as in MATLAB programming).

Our motivation for such a maximization is from the MAXBET problem [13] arising from the applied multivariate statistical analysis and data mining. The MAXBET problem is a generalization of the classical canonical correlation analysis (CCA) [10, 18, 20] for multi-view (also known as multi-set) situation, and the special case  $\ell=1$  is referred to as the maximal correlation problem (MCP) in the literature. In particular, the multi-view of a given random variable  $\mathbf{s} \in \mathbb{R}^n$  here means that  $\mathbf{s}$  can be split into multiple sub-variables, i.e.,  $\mathbf{s} = [\mathbf{s}_1; , \mathbf{s}_2; \dots; \mathbf{s}_m]$  with each  $\mathbf{s}_i \in \mathbb{R}^{n_i}$  representing a group of features to describe a certain statistical character. The MAXBET problem is to find projection matrices  $X_i \in \mathbb{R}^{n_i \times \ell}$ , one for each group of variables  $\mathbf{s}_i$ , to reduce the original  $n_i$  dimensional random variable  $\mathbf{s}_i$  to an  $\ell$ -dimensional one  $\mathbf{y}_i = X_i^T \mathbf{s}_i \in \mathbb{R}^{\ell}$  within this group. Suppose we now have p samples of  $\mathbf{s}$  collected in the data matrix  $S = [S_1; S_2; \dots; S_m] \in \mathbb{R}^{n \times p}$  with each  $S_i \in \mathbb{R}^{n_i \times p}$  and has been centralized, i.e.,  $S_i \mathbf{1}_p = 0$ . Under the MAXBET criterion (see [6, 13, 32]), the optimal projections for the sample data is formulated as the solution to the following maximization problem

$$\max_{\{X_i \in \mathbb{O}^{n_i \times \ell}\}} \sum_{i=1}^m \operatorname{tr}(X_i^{\mathsf{T}} S_i S_j^{\mathsf{T}} X_j). \tag{2}$$

Set  $A_{ij} = S_i S_i^{\mathrm{T}} \in \mathbb{R}^{n_i \times n_j}$ , and

$$A = \begin{bmatrix} A_{11} & A_{12} & \dots & A_{1m} \\ A_{21} & A_{22} & \dots & A_{2m} \\ \vdots & \vdots & & \vdots \\ A_{m1} & A_{m2} & \dots & A_{mm} \end{bmatrix}.$$
 (3)

Then the maximization problem (2) is the same as (1). Problem (2) has many applications. For example, it has been pointed out that MAXBET (2) coincides with another



Adv Comput Math (2022) 48:13 Page 3 of 34 13

least-squares criterion "MAXNEAR" [13]. In [31], an equivalent problem is considered in order to rotate the m configuration matrices for p objects to maximize certain agreement. Also, (2) is a generalization of the traditional multi-view CCA which sums over  $i \neq j$  and under constraints  $X_i^T S_i S_i^T X_i = I_\ell$  for all i (see [13]). Additionally, (2) is used to model the multi-view spectral clustering problem [49] in data mining. Many other applications can be found, for example, in [7, 9, 14, 17, 19, 44, 48].

MAXBET (1) has a very simple homogeneous quadratic objective function  $f(X) = \operatorname{tr}(X^T A X)$ , where the symmetric matrix A can be even assumed to be positive definite because f(X) and  $f(X) + c = \operatorname{tr}(X^T (A + \frac{c}{m\ell} I_n) X)$  over  $\mathcal{M}$  have the same maximizers. The difficulty, however, comes from the constraint  $X \in \mathcal{M}$ . In an extreme case when  $\ell = 1$  and  $n_i = 1$  for all  $1 \le i \le m$ , (1) reduces to an integer programming which is NP-hard (see, e.g., [14]). Another special case of m = 1 is equivalent to the classical eigenvalue problem of finding the eigenparis of A associated with the first  $\ell$  largest eigenvalues, and many sophisticated modern eigensolvers have been developed [15]. Motivated by the latter and a recent development on the unbalanced procrustes problem [56], in this paper, we will establish a connection of MAXBET to a particular nonlinear eigenvalue problem and develop efficient algorithms to tackle MAXBET (1). The efficiency of the new methods is demonstrated on various examples and on solving the model in the multi-view feature extraction problem [45] in data mining.

The rest of this paper is organized as follows. In Section 2, we will review the standard first-order optimality condition for (1) and explore new necessary conditions. Relying on these conditions, in Section 3, we will connect MAXBET to a coupled nonlinear eigenvalue system where each block  $X_i$  in X turns out to be an orthonormal eigenbasis matrix of a nonlinear eigenvector-dependent eigenvalue problem (NEPv). In Section 4, we will derive a simple self-consistent field (SCF) iteration based on our established characterization, and discuss the convergence behavior of this basic SCF iteration in Section 5. For a numerically efficient implementation of this SCF iteration, we will further integrate a practical strategy to accelerate the convergence. Numerical demonstrations of our proposed SCF iteration as well as its application to the multi-view feature extraction for unsupervised learning are reported in Section 7. Finally, concluding remarks are drawn in Section 8.

**Notation**  $\mathbb{R}^{m \times n}$  ( $\mathbb{C}^{m \times n}$ ) is the set of  $m \times n$  real (complex) matrices and  $\bullet^T$  and  $\bullet^H$  stand for the transpose and conjugate transpose of matrices/vectors, respectively.  $I_n = [e_1, \dots, e_n] \in \mathbb{R}^{n \times n}$  is the identity matrix, and  $\mathbf{1}_n \in \mathbb{R}^n$  is the vector of all ones. For  $B \in \mathbb{R}^{m \times n}$ ,  $\mathcal{R}(B)$  is the column subspace and its singular values are denoted by  $\sigma_i(B)$  for  $i = 1, \dots, \min(m, n)$  arranged in the nonincreasing order. For  $B \in \mathbb{R}^{n \times n}$ , sym $(B) = (B + B^T)/2$ ; if B is also symmetric, then  $\operatorname{eig}(B) = \{\lambda_i(B)\}_{i=1}^n$  denotes the set of its eigenvalues (counted by multiplicities) arranged in the nonincreasing order.  $B \succ 0 (\succeq 0)$  means that B is symmetric and positive definite (semi-definite), and  $B \prec 0 (\preceq 0)$  if  $-B \succ 0 (\succeq 0)$ .  $\|B\|_2$  and  $\|B\|_F$  are the spectral and Frobenius norm of matrix B, respectively. MATLAB-like notation is used to access the entries of a matrix or vector:  $X_{(i:j,k:l)}$  to denote the submatrix of a matrix X, consisting of the intersections of rows i to j and columns k to l, and when i:j is replaced by :,



13 Page 4 of 34 Adv Comput Math (2022) 48:13

it means all rows, similarly for columns;  $v_{(k)}$  refers the kth entry of a vector v and  $v_{(i:j)}$  is the subvector of v consisting of the ith to jth entries inclusive.

Given two  $X, Y \in \mathbb{O}^{n \times \ell}$ , the canonical angles between two subspaces  $\mathcal{R}(X)$  and  $\mathcal{R}(Y)$  of dimension  $\ell$  are given by

$$0 \le \theta_i(\mathcal{R}(X), \mathcal{R}(Y)) := \arccos \sigma_{\ell-i+1}(X^T Y) \le \frac{\pi}{2}, \quad \text{for } 1 \le i \le \ell,$$

and let

$$\Theta(\mathcal{R}(X), \mathcal{R}(Y)) = \operatorname{diag}(\theta_1(\mathcal{R}(X), \mathcal{R}(Y)), \dots, \theta_{\ell}(\mathcal{R}(X), \mathcal{R}(Y))). \tag{4}$$

For simplicity, we often write  $\Theta(X, Y)$  to mean  $\Theta(\mathcal{R}(X), \mathcal{R}(Y))$ .

# 2 Optimality conditions for MAXBET

We begin with the standard first- and second-order optimality conditions for MAX-BET (1). These are the extension of [50] for  $\ell=1$  to the general case  $\ell>1$ .

In what follows, we will adopt the convention that any matrix  $\mathbb{R}^{n \times \ell}$  will be implicitly partitioned into m blocks as  $X = [X_1; X_2; \dots; X_m]$  with  $X_i \in \mathbb{R}^{n_i \times \ell}$  as the ith block. Let  $A_i$  be the ith block rows of A in (3):

$$A_i = [A_{i1}, A_{i2}, \dots, A_{im}] \in \mathbb{R}^{n_i \times n}.$$
 (5)

The tangent space of  $\mathbb{O}^{n_i \times \ell}$  at  $X_i \in \mathbb{O}^{n_i \times \ell}$  is denoted by  $\mathcal{T}_{X_i} \mathbb{O}$ . It is known that (see, e.g., [2])

$$\mathcal{T}_{X_i} \mathbb{O}^{n_i \times \ell} = \left\{ Z_i \in \mathbb{R}^{n_i \times \ell} \middle| \begin{array}{l} Z_i = X_i K + (I_{n_i} - X_i X_i^{\mathsf{T}}) J \\ \forall K = -K^{\mathsf{T}} \in \mathbb{R}^{\ell \times \ell}, \ J \in \mathbb{R}^{n_i \times \ell} \end{array} \right\}.$$
 (6)

Consequently, the tangent space  $\mathcal{T}_X \mathcal{M}$  at  $X \in \mathcal{M} = \mathbb{O}^{n_1 \times \ell} \times \ldots \times \mathbb{O}^{n_m \times \ell}$  can be given by

$$\mathcal{T}_X \mathcal{M} = \{ Z = [Z_1; Z_2; \dots; Z_m] \in \mathbb{R}^{n \times \ell} \mid Z_i \in \mathcal{T}_{X_i} \mathbb{O}^{n_i \times \ell}, \ i = 1, 2, \dots, m \}.$$

For any  $Z \in \mathbb{R}^{n \times \ell}$ , the orthogonal projection onto  $\mathcal{T}_X \mathcal{M}$  is given by

$$\Pi_X(Z) = \begin{bmatrix} Z_1 - X_1 \cdot \operatorname{sym}(Z_1^{\mathsf{T}} X_1) \\ \vdots \\ Z_m - X_m \cdot \operatorname{sym}(Z_m^{\mathsf{T}} X_m) \end{bmatrix} \in \mathcal{T}_X \mathcal{M}.$$
 (7)

The first-order optimality condition in Lemma 2.1 has been developed in [32, Theorem 2.1] through the traditional Lagrangian multiplier theory [36]. Here we use the manifold structure of  $\mathcal{M}$  to give a simple proof.

**Lemma 2.1**  $X \in \mathcal{M}$  is a KKT point X of MAXBET (1) if and only if there are m symmetric  $\ell$ -by- $\ell$  matrices  $\{\Lambda_i\}_{i=1}^m$  such that

$$A_i X = X_i \Lambda_i, \quad i = 1, 2, \dots, m. \tag{8}$$

As a result,  $f(X) = \sum_{i=1}^{m} \operatorname{tr}(\Lambda_i)$ .



Adv Comput Math (2022) 48:13 Page 5 of 34 13

*Proof* Note that the gradient g(X) of  $f: \mathcal{M} \to \mathbb{R}$  at  $X \in \mathcal{M}$  is given by

$$g(X) = \Pi_X(\nabla f(X)) = 2AX - 2 \begin{bmatrix} X_1 \cdot \operatorname{sym}(X_1^{\mathsf{T}} A_1 X) \\ \vdots \\ X_m \cdot \operatorname{sym}(X_m^{\mathsf{T}} A_m X) \end{bmatrix}. \tag{9}$$

If  $X \in \mathcal{M}$  is a KKT point, then g(X) = 0 yielding (8) with  $\Lambda_i = \operatorname{sym}(X_i^T A_i X)$ . Suppose (8) with symmetric  $\Lambda_i$ . Pre-multiply (8) by  $X_i^T$  to get  $\Lambda_i = X_i^T A_i X$ . Since  $\Lambda_i$  is symmetric,  $\Lambda_i = \operatorname{sym}(\Lambda_i) = \operatorname{sym}(X_i^T A_i X)$  and hence  $A_i X = X_i \operatorname{sym}(X_i^T A_i X)$ , implying g(X) = 0, i.e., X is a KKT point. Finally,  $f(X) = \sum_{i=1}^m \operatorname{tr}(X_i^T A_i X) = \sum_{i=1}^m \operatorname{tr}(\Lambda_i)$ .

**Corollary 2.1** Let X be a KKT point. Then  $X_i^T A_i X$  and  $X_i^T G_i(X)$  for  $1 \le i \le m$  are symmetric, where

$$\check{A}_i = [A_{i1}, \dots, A_{i(i-1)}, 0, A_{i(i+1)}, \dots, A_{im}] \in \mathbb{R}^{n_i \times n},$$
(10a)

$$G_i(X) = \sum_{j \neq i}^m A_{ij} X_j = \check{A}_i X \in \mathbb{R}^{n_i \times \ell}, \quad 1 \le i \le m.$$
 (10b)

*Proof*  $X_i^T A_i X$  is symmetric because by (8)  $X_i^T A_i X = \Lambda_i$  which is symmetric. Expanding the left-hand side of (8) yields

$$A_{ii}X_i + \sum_{j \neq i}^m A_{ij}X_j = X_i\Lambda_i, \quad i = 1, 2, \dots, m.$$
 (11)

Hence, 
$$X_i^T G_i(X) = \Lambda_i - X_i^T A_{ii} X_i \in \mathbb{R}^{\ell \times \ell}$$
 is symmetric.

The result in Lemma 2.2 is a consequence of the second-order optimality condition (see, e.g., [36, 47]) on the manifold  $\mathcal{M}$ .

**Lemma 2.2** Let  $X \in \mathcal{M}$  be any local maximizer of (1). Then

$$\operatorname{tr}(Z^{\mathrm{T}}AZ) - \sum_{i=1}^{m} \operatorname{tr}(Z_{i}^{\mathrm{T}}Z_{i}\Lambda_{i}) \leq 0 \quad \text{for } Z \in \mathcal{T}_{X}\mathcal{M}, \tag{12}$$

where  $\Lambda_i = X_i^T A_i X$  for i = 1, 2, ..., m. If  $X \in \mathcal{M}$  is a KKT point and (12) holds strictly for any nonzero  $Z \in \mathcal{T}_X \mathcal{M}$ , then X is a strict local maximizer.

*Proof* We first compute the Hessian operator of  $f: \mathcal{M} \to \mathbb{R}$  acting on a tangent vector Z:

$$\operatorname{Hess} f(X)[Z] = \Pi_X(\mathbf{D}g(X)[Z]), \tag{13a}$$

where g(X) is as in (9) and  $\mathbf{D}g(X)[Z]$  is the directional derivative of g(X) along Z. We have

$$\frac{1}{2}\operatorname{Hess} f(X)[Z] = AZ - \begin{bmatrix} Z_1\Lambda_1 \\ \vdots \\ Z_m\Lambda_m \end{bmatrix} - \begin{bmatrix} X_1 \operatorname{sym}(X_1^{\mathsf{T}}A_1Z - X_1^{\mathsf{T}}Z_1\Lambda_1) \\ \vdots \\ X_m \operatorname{sym}(X_m^{\mathsf{T}}A_mZ - X_m^{\mathsf{T}}Z_m\Lambda_m) \end{bmatrix}. (13b)$$

13 Page 6 of 34 Adv Comput Math (2022) 48:13

Now, by  $X_i^T Z_i + Z_i^T X_i = 0$  for any  $Z_i \in \mathcal{T}_{X_i} \mathbb{O}^{n_i \times \ell}$ , it holds  $\operatorname{tr}(Z_i^T X_i S) = 0$  for any symmetric matrix  $S \in \mathbb{R}^{\ell \times \ell}$ , and hence

$$\operatorname{tr}\left(Z_{i}^{\mathrm{T}}X_{i} \operatorname{sym}(X_{i}^{\mathrm{T}}A_{i}Z - X_{i}^{\mathrm{T}}Z_{i}\Lambda_{i})\right) = 0 \quad \text{for } 1 \leq i \leq m.$$

The second-order necessary optimality condition  $\operatorname{tr}(Z^{\operatorname{T}}\operatorname{Hess} f(X)[Z]) \leq 0$  for any  $Z \in \mathcal{T}_X \mathcal{M}$  (see, e.g., [47]) together with (13) leads to (12). Furthermore, X is a strict local maximizer if  $\operatorname{tr}(Z^{\operatorname{T}}\operatorname{Hess} f(X)[Z]) < 0$  for any nonzero  $Z \in \mathcal{T}_X \mathcal{M}$ , i.e., (12) holds strictly.

**Corollary 2.2** If  $X \in \mathcal{M}$  be a local maximizer of (1), then  $\operatorname{tr}(J_i^T A_{ii} J_i) \leq \operatorname{tr}(\Lambda_i)$  for any  $J_i \in \mathbb{O}^{n_i \times \ell}$  such that  $J_i^T X_i = 0$ . If also  $\ell \geq 2$ , then  $\operatorname{tr}(X_i^T G_i(X)) \geq 0$ .

*Proof* Note that any such  $J_i \in \mathcal{T}_{X_i} \mathbb{O}^{n_i \times \ell}$ , and thus,  $Z \in \mathbb{R}^{n \times \ell}$  with  $Z_j = 0$  for all  $j \neq i$  and  $Z_i = J_i$  is in  $\mathcal{T}_X \mathcal{M}$ . Plug it into (12) to get  $\operatorname{tr}(J_i^T A_{ii} J_i) \leq \operatorname{tr}(\Lambda_i)$ .

Suppose  $\ell \geq 2$ . According to (6), construct  $Z \in \mathcal{T}_X \mathcal{M}$  by setting  $Z_j = 0$  for  $j \neq i$  and  $Z_i = X_i K$  for some skew-symmetric matrix  $K \in \mathbb{R}^{\ell \times \ell}$  to be specified. Plug it into (12) to get (keep in mind that  $X_i^T G_i(X)$  is symmetric)

$$\operatorname{tr}(K^{\mathsf{T}}X_{i}^{\mathsf{T}}A_{ii}X_{i}K) \leq \operatorname{tr}(K^{\mathsf{T}}K\Lambda_{i})$$

$$= \operatorname{tr}(K^{\mathsf{T}}KX_{i}^{\mathsf{T}}A_{ii}X_{i}) + \operatorname{tr}(K^{\mathsf{T}}KX_{i}^{\mathsf{T}}G_{i}(X))$$

$$= \operatorname{tr}(K^{\mathsf{T}}X_{i}^{\mathsf{T}}A_{ii}X_{i}K) + \operatorname{tr}(K^{\mathsf{T}}X_{i}^{\mathsf{T}}G_{i}(X)K),$$

yielding  $\operatorname{tr}(K^{\mathrm{T}}X_{i}^{\mathrm{T}}G_{i}(X)K) \geq 0$ . We now specify K as follows. Let the spectral decomposition  $X_{i}^{\mathrm{T}}G_{i}(X) = U\operatorname{diag}(\mu_{1},\ldots,\mu_{\ell})U^{\mathrm{T}}$ . Given any fixed (j,k) with  $1 \leq j < k \leq \ell$  and  $\xi \neq 0$ , let  $K = K_{jk} = \xi U(\boldsymbol{e}_{j}\boldsymbol{e}_{k}^{\mathrm{T}} - \boldsymbol{e}_{k}\boldsymbol{e}_{j}^{\mathrm{T}})U^{\mathrm{T}}$ . It is skew-symmetric, and thus

$$0 \le \operatorname{tr}(K_{jk}^{\mathsf{T}} X_i^{\mathsf{T}} G_i(X) K_{jk}) = \xi^2 (\mu_k + \mu_j).$$

Summing all of them up over  $1 \le j < k \le \ell$ , we have

$$0 \le \xi^2 \sum_{1 \le j < k \le \ell} (\mu_j + \mu_k) = \xi^2 (\ell - 1) \sum_{i=1}^{\ell} \mu_i$$
$$= \xi^2 (\ell - 1) \cdot \operatorname{tr}(X_i^{\mathrm{T}} G_i(X)),$$

yielding  $\operatorname{tr}(X_i^{\mathrm{T}}G_i(X)) \geq 0$ .

Besides these standard optimality conditions for a local maximizer, we will make use of the special structure of MAXBET (1) to establish the following necessary condition for a global maximizer.

**Theorem 1** Let  $X_{\text{opt}} = [X_{\text{opt},1}; \dots; X_{\text{opt},m}] \in \mathcal{M}$  be a global maximizer of MAXBET (1). Then  $X_{\text{opt},i}^T G_i(X_{\text{opt}}) \succeq 0$  for  $i = 1, 2, \dots, m$ .



Adv Comput Math (2022) 48:13 Page 7 of 34 13

*Proof* For any fixed  $1 \le i \le m$ , consider  $Y \in \mathcal{M}$  with  $Y_j = X_{\text{opt},j}$  for  $j \ne i$  and  $Y_i = X_{\text{opt},i}Q$ , where  $Q \in \mathbb{O}^{\ell \times \ell}$  is arbitrary. Since  $f(X_{\text{opt}}) - f(Y) \ge 0$  for any  $Q \in \mathbb{O}^{\ell \times \ell}$ , we have

$$\begin{split} 0 &\leq f(X_{\text{opt}}) - f(Y) \\ &= \sum_{j=1,k=1}^{m} \text{tr}(X_{\text{opt},j}^{\text{T}} A_{jk} X_{\text{opt},k}) - \sum_{j=1,k=1}^{m} \text{tr}(Y_{j}^{\text{T}} A_{jk} Y_{k}) \\ &= \text{tr}(X_{\text{opt},i}^{\text{T}} A_{ii} X_{\text{opt},i}) + 2 \sum_{j=1,\ j \neq i}^{m} \text{tr}(X_{\text{opt},i}^{\text{T}} A_{ij} X_{\text{opt},j}) \\ &- \text{tr}(Q^{\text{T}} X_{\text{opt},i}^{\text{T}} A_{ii} X_{\text{opt},i} Q) - 2 \sum_{j=1,\ j \neq i}^{m} \text{tr}(Q^{\text{T}} X_{\text{opt},i}^{\text{T}} A_{ij} X_{\text{opt},j}) \\ &= 2 \big[ \text{tr}(X_{\text{opt},i}^{\text{T}} G_{i}(X_{\text{opt}})) - \text{tr}(Q^{\text{T}} X_{\text{opt},i}^{\text{T}} G_{i}(X_{\text{opt}})) \big], \end{split}$$

which implies that

$$\operatorname{tr}(X_{\operatorname{opt},i}^{\operatorname{T}}G_{i}(X_{\operatorname{opt}})) \ge \max_{Q \in \mathbb{Q}^{\ell \times \ell}} \operatorname{tr}(Q^{\operatorname{T}}X_{\operatorname{opt},i}^{\operatorname{T}}G_{i}(X_{\operatorname{opt}})). \tag{14}$$

Now applying [55, Lemma 3], we conclude  $X_{\text{opt }i}^{\text{T}}G_i(X_{\text{opt}}) \geq 0$ .

# 3 Connect to a coupled nonlinear eigenvalue system

To connect the solution of MAXBET (1) with a nonlinear eigenvalue problem, we rewrite the KKT condition (11) as

$$[A_{ii} + G_i(X)X_i^{\mathrm{T}}]X_i = X_i\Lambda_i, i = 1, 2, ..., m$$

where  $G_i(X)$  is defined in (10). This implies that at a KKT point X, each  $X_i$  is an orthonormal eigenbasis matrix of  $A_{ii} + G_i(X)X_i^{\mathrm{T}}$ . However, the matrix  $G_i(X)X_i^{\mathrm{T}} \in \mathbb{R}^{n \times n}$  is not necessarily symmetric, even at KKT points. A further step is then to symmetrize it by considering the following matrices:

$$E_i(X) := A_{ii} + G_i(X)X_i^{\mathrm{T}} + X_iG_i(X)^{\mathrm{T}}, \quad i = 1, 2, \dots, m,$$
 (15)

which are symmetric and thus have real eigensystems at any given X.

The following theorem connects MAXBET with a coupled system of NEPv (16).

**Theorem 2**  $X \in \mathcal{M}$  is a KKT point (i.e., satisfying the system (8)) if and only if  $X_i^T G_i(X) \in \mathbb{R}^{n_i \times n_i}$  for  $1 \le i \le m$  are symmetric, and

$$E_i(X)X_i = X_i\Psi_i, \quad i = 1, 2, \dots, m,$$
 (16)

for some  $\Psi_i \in \mathbb{R}^{\ell \times \ell}$ . As a result,

$$f(X) = \text{tr}(\Psi_i) + \sum_{j \neq i, k \neq i}^{m} \text{tr}(X_j^{\mathsf{T}} A_{jk} X_k), \quad \forall i = 1, 2, \dots, m.$$
 (17)



13 Page 8 of 34 Adv Comput Math (2022) 48:13

*Proof* Let *X* be a KKT point. Then  $A_i X = X_i \Lambda_i$  by (8) and  $X_i^T G_i(X)$  is symmetric by Corollary 2.1. We then have

$$E_{i}(X)X_{i} = \left(A_{ii} + G_{i}(X)X_{i}^{T} + X_{i}G_{i}(X)^{T}\right)X_{i} = X_{i}(\Lambda_{i} + G_{i}(X)^{T}X_{i}) =: X_{i}\Psi_{i},$$

where  $\Psi_i = \Lambda_i + G_i(X)^T X_i$ . On the other hand, if  $E_i(X) X_i = X_i \Psi_i$  for a symmetric  $\Psi_i$ , then  $\Psi_i = X_i^T E_i(X) X_i$  and thus symmetric. Also,  $A_i X = X_i (\Psi_i - G_i(X)^T X_i) =: X_i \Lambda_i$ , where  $\Lambda_i = \Psi_i - G_i(X)^T X_i$  is symmetric because both  $\Psi_i$  and  $G_i(X)^T X_i$  are symmetric.

Lastly, by the definition of  $E_i(X)$  in (15) we have

$$f(X) = \sum_{j \neq i, k \neq i}^{m} \operatorname{tr}(X_{j}^{T} A_{jk} X_{k}) + \operatorname{tr}(X_{i}^{T} A_{ii} X_{i}) + 2 \sum_{j=1, j \neq i}^{m} \operatorname{tr}(X_{i}^{T} A_{ij} X_{j})$$

$$= \sum_{j \neq i, k \neq i}^{m} \operatorname{tr}(X_{j}^{T} A_{jk} X_{k}) + \operatorname{tr}(X_{i}^{T} E_{i}(X) X_{i})$$

$$= \sum_{j \neq i, k \neq i}^{m} \operatorname{tr}(X_{j}^{T} A_{jk} X_{k}) + \operatorname{tr}(\Psi_{i}), \quad \text{(by (16))}$$

as expected.

# 3.1 Eigenspaces associated with a local maximizer

When X is a local maximizer, by Theorem 2, it satisfies (16), a coupled system of m nonlinear eigenvalue problems. Each  $E_i(X)$  depends on all eigenvector matrices  $X_j$  for  $1 \le j \le m$ , and  $X_i$  is an orthonormal eigenbasis matrix of  $E_i(X)$  associated with  $\ell$  of its eigenvalues in  $\operatorname{eig}(\Psi_i)$ . Let

$$\operatorname{eig}(\Psi_i) = \{\lambda_{\pi_{i,1}}(E_i(X)) \ge \dots \ge \lambda_{\pi_{i,\ell}}(E_i(X))\},\tag{18}$$

where  $1 \leq \pi_{i,1} \leq \cdots \leq \pi_{i,\ell} \leq n_i$ .

**Theorem 3** Suppose  $X \in \mathcal{M}$  is a local maximizer of (1) and  $2 \leq \ell \leq n$ . Then  $\lambda_{\pi_{i,1}}(E_i(X)) = \lambda_1(\Psi_i) \geq \lambda_\ell(E_i(X))$ .

*Proof* Suppose, to the contrary, that  $\lambda_{\pi_{i,1}}(E_i(X)) < \lambda_{\ell}(E_i(X))$ . Then we can choose an orthonormal eigenbasis matrix  $J_i \in \mathbb{O}^{n_i \times \ell}$  associated with the first  $\ell$  largest eigenvalues of  $E_i(X)$ . This implies that  $J_i^T X_i = 0$  and

$$\operatorname{tr}(J_{i}^{T}E_{i}(X)J_{i}) = \operatorname{tr}(J_{i}^{T}A_{ii}J_{i}) > \operatorname{tr}(X_{i}^{T}E_{i}(X)X_{i})$$

$$= \operatorname{tr}(\Psi_{i})$$

$$= \operatorname{tr}(X_{i}^{T}A_{ii}X_{i}) + 2\sum_{j \neq i}\operatorname{tr}(X_{i}^{T}A_{ij}X_{j})$$
(19)

$$= \operatorname{tr}(\Lambda_i) + \sum_{j \neq i} \operatorname{tr}(X_i^{\mathsf{T}} A_{ij} X_j). \tag{20}$$



Adv Comput Math (2022) 48:13 Page 9 of 34 13

On the other hand, by Corollary 2.2, we know that  $\operatorname{tr}(J_i^T A_{ii} J_i) \leq \operatorname{tr}(\Lambda_i)$ , which leads to  $\operatorname{tr}(X_i^T G_i(X)) = \sum_{j \neq i} \operatorname{tr}(X_i^T A_{ij} X_j) < 0$ . This contradicts  $\operatorname{tr}(X_i^T G_i(X)) \geq 0$  again by Corollary 2.2, and the conclusion follows.

## 3.2 Eigenspaces associated with a global maximizer

To further explore the eignspaces  $\mathcal{R}(X_i)$  of  $E_i(X)$  associated with a global maximizer, we will use the following well-known von Neumann's trace inequality [43] (see also [40]).

**Lemma 3.1** [43] For  $C_1$ ,  $C_2 \in \mathbb{R}^{m \times n}$ , we have

$$\operatorname{tr}(C_1^{\mathrm{T}}C_2) \leq \sum_{i=1}^n \sigma_i(C_1)\sigma_i(C_2).$$

**Theorem 4** Let  $X_{\text{opt}} = [X_{\text{opt},1}; \dots; X_{\text{opt},m}] \in \mathcal{M}$  be a global maximizer of MAXBET (1). Then, each  $X_{\text{opt},i}$  is an orthonormal eigenbasis matrix of  $E_i(X_{\text{opt}})$  associated with its  $\ell$  largest eigenvalues.

*Proof* Suppose there is an i  $(1 \le i \le m)$  for which the statement is not true. Then there exists a  $J_i \in \mathbb{O}^{n_i \times \ell}$  such that

$$\operatorname{tr}(\Psi_{\operatorname{opt},i}) = \operatorname{tr}(X_{\operatorname{opt},i}^{\operatorname{T}} E_i(X_{\operatorname{opt}}) X_{\operatorname{opt},i}) < \operatorname{tr}(J_i^{\operatorname{T}} E_i(X_{\operatorname{opt}}) J_i). \tag{21}$$

Now, let  $J_i^{\mathrm{T}}G_i(X_{\mathrm{opt}}) = U\Sigma V^{\mathrm{T}}$  be the SVD, and set  $Q = UV^{\mathrm{T}} \in \mathbb{O}^{\ell \times \ell}$ . Consider  $Y \in \mathcal{M}$  with  $Y_j = X_j$  for all  $j \neq i$  but  $Y_i = J_i Q$ . We have by (17)

$$f(Y) - f(X) = \operatorname{tr}(Q^{T} J_{i}^{T} A_{ii} J_{i} Q) + 2 \sum_{j \neq i} \operatorname{tr}(Q^{T} J_{i}^{T} A_{ij} X_{\operatorname{opt}, j}) - \operatorname{tr}(\Psi_{\operatorname{opt}, i})$$

$$> \operatorname{tr}(J_{i}^{T} A_{ii} J_{i}) + 2 \sum_{j \neq i} \operatorname{tr}(Q^{T} J_{i}^{T} A_{ij} X_{\operatorname{opt}, j}) \qquad (22)$$

$$- \operatorname{tr}(J_{i}^{T} E_{i}(X_{\operatorname{opt}}) J_{i}) \qquad (by (21))$$

$$= 2 \Big[ \operatorname{tr}(Q^{T} J_{i}^{T} \sum_{j \neq i} A_{ij} X_{\operatorname{opt}, j}) - \operatorname{tr}(J_{i}^{T} \sum_{j \neq i} A_{ij} X_{\operatorname{opt}, j} X_{\operatorname{opt}, i}^{T} J_{i}) \Big]$$

$$= 2 \Big[ \operatorname{tr}(Q^{T} J_{i}^{T} G_{i}(X_{\operatorname{opt}})) - \operatorname{tr}(J_{i}^{T} \sum_{j \neq i} A_{ij} X_{\operatorname{opt}, j} X_{\operatorname{opt}, i}^{T} J_{i}) \Big]$$

$$\geq 0, \qquad (23)$$



13 Page 10 of 34 Adv Comput Math (2022) 48:13

where the last inequality holds because by Lemma 3.1

$$\operatorname{tr}\left(J_{i}^{\mathsf{T}}\sum_{j\neq i}A_{ij}X_{\mathrm{opt},j}X_{\mathrm{opt},i}^{\mathsf{T}}J_{i}\right) \leq \sum_{j=1}^{\ell}\sigma_{j}(J_{i}^{\mathsf{T}}\sum_{j\neq i}A_{ij}X_{\mathrm{opt},j})\sigma_{j}(X_{\mathrm{opt},i}^{\mathsf{T}}J_{i})$$

$$\leq \sum_{j=1}^{\ell}\sigma_{j}(J_{i}^{\mathsf{T}}\sum_{j\neq i}A_{ij}X_{\mathrm{opt},j})$$

$$= \sum_{j=1}^{\ell}\sigma_{j}(J_{i}^{\mathsf{T}}G_{i}(X_{\mathrm{opt}}))$$

$$= \operatorname{tr}(Q^{\mathsf{T}}J_{i}^{\mathsf{T}}A_{ii}J_{i}Q),$$

where we have used  $0 \le \sigma_j(J_i^T X_i) \le ||J_i||_2 ||X_i||_2 = 1$  for all j. As a consequence of (23), we have  $f(Y) > f(X_{\text{opt}})$ , a contradiction, because  $X_{\text{opt}}$  is a global maximizer, and the proof is complete.

Both Theorems 1 and 4 establish necessary conditions for a global maximizer. Unfortunately, the converse of Theorem 4 may not be true; that is, X may not be a global maximizer even if X is a KKT point satisfying that each  $X_i$  is an orthonormal eigenbasis matrix of  $E_i(X)$  associated with its  $\ell$  largest eigenvalues. The following is such an example that numerically substantiates this statement.

Example 1 Consider m = 2,  $\ell = 2$ ,  $n_1 = n_2 = 3$ , and

$$A = \begin{bmatrix} 2 & 0 & 1 & 1 & 2 & 2 \\ 0 & 0 & 2 & 1 & 1 & 1 \\ 1 & 2 & 2 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 & 1 \\ 2 & 1 & 1 & 0 & 0 & 1 \\ 2 & 1 & 1 & 1 & 1 & 0 \end{bmatrix} \in \mathbb{R}^{6 \times 6}.$$

Using the Riemannian trust-region (RTR) method<sup>1</sup> [1, 2], we found two approximated KKT points

$$X = \begin{bmatrix} -0.9587 & 0.2769 \\ -0.1173 & -0.5924 \\ -0.2590 & -0.7566 \\ -0.0653 & -0.9263 \\ -0.7841 & 0.2776 \\ -0.6171 & -0.2547 \end{bmatrix} \text{ and } \widehat{X} = \begin{bmatrix} -0.9415 & 0.3334 \\ -0.1586 & -0.5691 \\ -0.2974 & -0.7517 \\ -0.5916 & 0.4355 \\ -0.3557 & -0.8962 \\ -0.7236 & 0.0845 \end{bmatrix}$$

for which  $||g(X)||_2 \approx 3.7 \times 10^{-10}$ ,  $||g(\widehat{X})||_2 \approx 9.7 \times 10^{-10}$ , and

$$f(X) = 15.1940 > 13.4575 = f(\widehat{X}),$$

showing at least that  $\widehat{X}$  is not a global maximizer. Numerically, it has been checked that for either  $Y \in \{X, \widehat{X}\}$ , each  $Y_i$  is an orthonormal eigenbasis matrix of  $E_i(Y)$ 

<sup>&</sup>lt;sup>1</sup>RTR is available at www.manopt.org.



Adv Comput Math (2022) 48:13 Page 11 of 34 13

associated with its  $\ell$  largest eigenvalues. Also, we sampled  $10^7$  random tangent "vectors" at X and  $\widehat{X}$  and found that the second-order sufficient inequality condition (12) strictly hold, indicating that both X and  $\widehat{X}$  are likely strictly local maxima. This numerical observation suggests that, unlike the traditional eigenvalue problem (i.e., m = 1), MAXBET (1) (with m > 1) may admit local but non-global maximizers.

## 4 An SCF iteration for MAXBET

We would like to mention that the SCF iteration is commonly used to solve NEPv [8] from the Kohn–Sham density functional theory in electronic structure calculations [33, 39]. Lately, it has been attracting a great deal of attention in data science (e.g., [4, 8, 30, 46, 52, 53, 55, 56]).

Our eigen-based method for MAXBET (1) is an inner-outer iteration scheme that alternatingly solves the coupled system (16) of m nonlinear eigenvalue problems in a self-consistent manner. Specifically, during the kth outer-loop iteration an approximate solution  $X^{(k)} = \left[X_1^{(k)}; \ldots; X_m^{(k)}\right]$  is updated one block at a time in the Gauss-Seidel style inner loop to produce m intermediate approximations

$$X^{(k+\frac{i}{m})} = \left[X_1^{(k+1)}; \dots; X_i^{(k+1)}; X_{i+1}^{(k)}; \dots; X_m^{(k)}\right], \ i = 1, 2, \dots, m,$$

the last one of which is the next approximation  $X^{(k+1)}$ . The updating of each block is based on the necessary conditions in Theorem 4 for a global maximizer. Algorithm 1 outlines the main algorithmic idea.

Remark 1 There are a few comments about Algorithm 1 in order.

## Algorithm 1 Basic SCF iteration for MAXBET (1).

**Require:** symmetric  $A \in \mathbb{R}^{n \times n}$ ,  $X^{(0)} \in \mathcal{M}$ ;

**Ensure:** a maximizer of MAXBET (1).

- 1: **for**  $k = 1, \ldots$ , until convergence **do**
- 2: **for** i = 1, ..., m **do**
- 3: Compute partial eigen-decomposition:

$$E_i(X^{(k+\frac{i-1}{m})})\widehat{X}_i^{(k+1)} = \widehat{X}_i^{(k+1)}\widehat{\Psi}_i^{(k+1)}, \quad \widehat{X}_i^{(k+1)} \in \mathbb{O}^{n_i \times \ell}$$

associated with the  $\ell$  largest eigenvalues of  $E_i(X^{(k+\frac{i-1}{m})})$ ;

- 4: Compute SVD:  $(\widehat{X}_{i}^{(k+1)})^{\mathrm{T}}G_{i}(X^{(k+\frac{i-1}{m})}) = U_{i}\Sigma_{i}V_{i}^{\mathrm{T}}$  and set  $X_{i}^{(k+1)} = \widehat{X}_{i}^{(k+1)}U_{i}V_{i}^{\mathrm{T}};$
- 5: end for
- 6: **end forreturn** the last  $X^{(k)}$ .



13 Page 12 of 34 Adv Comput Math (2022) 48:13

1. Both  $\widehat{X}_i^{(k+1)}$  and  $X_i^{(k+1)}$  at lines 3 and 4 maximize  $\operatorname{tr}\left(X_i^{\mathrm{T}} E_i(X^{(k+\frac{i-1}{m})}) X_i\right)$  over  $X_i \in \mathbb{O}^{n_i \times \ell}$ . Hence,

$$\operatorname{tr}\left((\widehat{X}_{i}^{(k+1)})^{\mathrm{T}}E_{i}(X^{(k+\frac{i-1}{m})})\widehat{X}_{i}^{(k+1)}\right) = \operatorname{tr}\left((X_{i}^{(k+1)})^{\mathrm{T}}E_{i}(X^{(k+\frac{i-1}{m})})X_{i}^{(k+1)}\right) \\ \geq \operatorname{tr}\left((X_{i}^{(k)})^{\mathrm{T}}E_{i}(X^{(k+\frac{i-1}{m})})X_{i}^{(k)}\right). \tag{24}$$

When the inequality in (24) is strict, later in Theorem 6 we will prove that  $f(X^{(k+\frac{i}{m})}) > f(X^{(k+\frac{i-1}{m})})$ , i.e., the objective function value in the ith inner step strictly increases. An implication from line 4 is  $(X_i^{(k+1)})^T G_i(X^{(k+\frac{i-1}{m})}) = V_i \Sigma_i V_i^T \succeq 0$ , and thus

$$\operatorname{tr}\left((X_i^{(k+1)})^{\mathrm{T}}G_i(X^{(k+\frac{i-1}{m})})\right) = \sum_{i=1}^{\ell} \sigma_j\left((X_i^{(k+1)})^{\mathrm{T}}G_i(X^{(k+\frac{i-1}{m})})\right). \tag{25}$$

Moreover, the following relation holds

$$E_{i}(X^{(k+\frac{i-1}{m})})X_{i}^{(k+1)} = X_{i}^{(k+1)}\Psi_{i}^{(k+1)}, \ \Psi_{i}^{(k+1)} = (U_{i}V_{i}^{T})^{T}\widehat{\Psi}_{i}^{(k+1)}(U_{i}V_{i}^{T}).$$
(26)

2. As there are infinitely many choices of orthonormal eigenbasis matrix (at line 3) associated with the first  $\ell$  largest eigenvalues of  $E_i(X^{(k+\frac{i-1}{m})})$ , we are particularly interested in the one that maximizes the objective function f(X) over  $X_i$  while with  $X_j = X_j^{(k+1)}$  (j < i) and  $X_j = X_j^{(k)}$  (j > i) fixed. Note

$$f(X) = \operatorname{tr}(X_i^{\mathsf{T}} A_{ii} X_i) + 2 \operatorname{tr}(X_i^{\mathsf{T}} G_i (X^{(k + \frac{i - 1}{m})}) + c_k$$
 (27)

for some constant  $c_k$ . Given the subspace  $\mathcal{R}(\widehat{X}_i^{(k+1)})$ , any other orthonormal eigenbasis matrix takes the form  $X_i = \widehat{X}_i^{(k+1)} P$  for some  $P \in \mathbb{O}^{\ell \times \ell}$ . We would like to maximize the right-hand side of (27) over such  $X_i$ , which leads to  $P = U_i V_i^T$  as in line 4.

3. For a stopping criterion at line 1, we can use the relative difference of the two successive objective function values and/or the scaled norm of the gradient:

$$\frac{|f(X^{(k+1)}) - f(X^{(k)})|}{|f(X^{(k)})|} \le \varepsilon_{\text{scf}_g} \quad \text{and/or} \quad \frac{\|g(X^{(k)})\|_1}{\|A\|_1} \le \varepsilon_{\text{scf}_g}. \tag{28}$$

The scaling factor  $||A||_1$  is introduced in the latter because a global maximizer  $X_{\text{opt}}$  is invariant under arbitrary positive scaling on A, but the gradient is dependent on such scaling. The  $\ell_1$ -matrix norm is used for numerical convenience.

We commented before that  $\widehat{X}_i^{(k+1)}$  at line 3 as an orthonormal eigenbasis matrix is not unique even if the eigenvalue gap

$$\xi_i^{(k)} := \lambda_\ell \left( E_i(X^{(k + \frac{i-1}{m})}) \right) - \lambda_{\ell+1} \left( E_i(X^{(k + \frac{i-1}{m})}) \right) > 0.$$
 (29)

Nonetheless, the next proposition shows that  $X_i^{(k+1)}$  is still unique under (29) and a full rank assumption.



Adv Comput Math (2022) 48:13 Page 13 of 34 13

**Proposition 5** In Algorithm 1, if (29) holds and if rank  $((\widehat{X}_i^{(k+1)})^T G_i(X^{(k+\frac{i-1}{m})})) = \ell$ , then  $X_i^{(k+1)}$  at line 4 is uniquely determined.

Proof Under (29), it is known that the eigenspace corresponding to the  $\ell$  largest eigenvalues of  $E_i(X^{(k+\frac{i-1}{m})})$  is unique [41, p. 244], and thus any two orthonormal eigenbases matrices  $\widehat{X}_i^{(k+1)}$  and  $\widetilde{X}_i^{(k+1)}$  for the eigenspace are related by  $\widetilde{X}_i^{(k+1)} = \widehat{X}_i^{(k+1)}Q$  for some  $Q \in \mathbb{O}^{\ell \times \ell}$ . In what follows, we will show that with  $\widetilde{X}_i^{(k+1)}$  instead of  $\widehat{X}_i^{(k+1)}$ , the resulting  $X_i^{(k+1)}$  at the end of line 4 is independent of Q. For convenience, we simplify notation by denoting  $\widehat{X}_i = \widehat{X}_i^{(k+1)}$ ,  $\widetilde{X}_i = \widetilde{X}_i^{(k+1)}$  and  $G_i = G_i(X^{(k+\frac{i-1}{m})})$ . Now  $X_i^{(k+1)} = \widetilde{X}_i P$ , where P is the orthogonal polar factor of  $\widetilde{X}_i^T G_i$ , which under  $\operatorname{rank}(\widetilde{X}_i^T G_i) = \ell$  is orthogonal and unique [27], and can be written as  $P = \widetilde{X}_i^T G_i(G_i^T \widetilde{X}_i \widetilde{X}_i^T G_i)^{-1/2}$ . Thus with  $\widetilde{X}_i$ , the output of line 4 is  $X_i^{(k+1)} = \widetilde{X}_i \widetilde{X}_i^T G_i(G_i^T \widetilde{X}_i \widetilde{X}_i^T G_i)^{-1/2} = \widehat{X}_i \widehat{X}_i^T G_i(G_i^T \widehat{X}_i \widehat{X}_i^T G_i)^{-1/2}$ , independent of Q.

Remark 2 In practice, the generic condition (29) will be always true. Nevertheless, we point out that this generic condition is not a must for the execution of Algorithm 1. Indeed, when  $\xi_i^{(k)} = 0$ , we can choose  $\widehat{X}_i^{(k+1)}$  as an arbitrary orthonormal eigenbasis matrix associated with the  $\ell$  largest eigenvalues of  $E_i(X^{(k+\frac{i-1}{m})})$ , and we will see in Theorem 7 that a particular limit point of  $\{X^{(k)}\}$  can still be a KKT point satisfying the necessary global optimality conditions in Theorems 1 and 4.

# 5 Convergence analysis of the SCF iteration

We next establish some convergence results for Algorithm 1 in order to better understand the behavior of this basic SCF iteration in Algorithm 1 so that a more practical algorithm can be designed in the next section.

**Lemma 5.1** [23, 24, 28] Let  $U \in \mathbb{O}^{n \times \ell}$  such that  $\mathcal{R}(U)$  is an eigenspace of a symmetric matrix  $H \in \mathbb{R}^{n \times n}$  associated with its  $\ell$  largest eigenvalues  $\lambda_j(H)$ ,  $j = 1, 2, \ldots, \ell$ , and let  $V \in \mathbb{O}^{n \times \ell}$ . If  $\lambda_{\ell+1}(H) < \lambda_{\ell}(H)$ , then

$$\frac{\sum_{j=1}^{\ell} \left( \lambda_j(H) - \lambda_j(V^{H}HV) \right)}{\lambda_1(H) - \lambda_n(H)} \leq \|\sin\Theta(V, U)\|_F^2 
\leq \frac{\sum_{j=1}^{\ell} \left( \lambda_j(H) - \lambda_j(V^{H}HV) \right)}{\lambda_{\ell}(H) - \lambda_{\ell+1}(H)}.$$

**Theorem 6** Let  $\{X^{(k)}\}$  be the sequence from Algorithm 1. Then

1. The sequence  $\{f(X^{(k)})\}\$  converges monotonically, and

$$f(X^{(k+1)}) - f(X^{(k)}) \ge \sum_{i=1}^{m} \xi_i^{(k)} \cdot \|\sin\Theta(X_i^{(k)}, X_i^{(k+1)})\|_F^2$$
 (30)



13 Page 14 of 34 Adv Comput Math (2022) 48:13

where  $\xi_i^{(k)}$  is defined in (29). Therefore for any  $1 \leq i \leq m$ , if  $\liminf_{k \to \infty} \xi_i^{(k)} > 0$ , then  $\lim_{k \to \infty} \sin \Theta(X_i^{(k)}, X_i^{(k+1)}) = 0$ ;

2. If  $f(X^{(k)}) = f(X^{(k+1)})$  and rank  $((X_i^{(k+1)})^T G_i(X^{(k+\frac{i-1}{m})})) = \ell$  for  $i = 1, \ldots, m$ , then  $X^{(k)} = X^{(k+1)}$ , and  $X^{(k)}$  is a KKT point satisfying the necessary global optimality conditions in Theorems 1 and 4.

Proof We have

$$\operatorname{tr}\left((X_{i}^{(k+1)})^{\mathrm{T}}G_{i}(X^{(k+\frac{i-1}{m})})(X_{i}^{(k)})^{\mathrm{T}}X_{i}^{(k+1)}\right) \\
\leq \sum_{j=1}^{\ell} \sigma_{j}\left((X_{i}^{(k+1)})^{\mathrm{T}}G_{i}(X^{(k+\frac{i-1}{m})})\right) \cdot \sigma_{j}\left((X_{i}^{(k)})^{\mathrm{T}}X_{i}^{(k+1)}\right) \quad \text{(by Lemma 3.1)} \\
\leq \sum_{j=1}^{\ell} \sigma_{j}\left((X_{i}^{(k+1)})^{\mathrm{T}}G_{i}(X^{(k+\frac{i-1}{m})})\right) \quad (\operatorname{since } 0 \leq \sigma_{j}\left((X_{i}^{(k)})^{\mathrm{T}}X_{i}^{(k+1)}\right) \leq 1) \\
= \operatorname{tr}\left((X_{i}^{(k+1)})^{\mathrm{T}}G_{i}(X^{(k+\frac{i-1}{m})})\right). \quad (\operatorname{by (25)}) \quad (31)$$

Noticing that  $X^{(k+\frac{i-1}{m})}$  differs from  $X^{(k+\frac{i}{m})}$  only in their *i*th blocks, we get

$$\begin{split} &f(X^{(k+\frac{i}{m})}) - f(X^{(k+\frac{i-1}{m})}) \\ &= \operatorname{tr}\left((X_{i}^{(k+1)})^{\mathrm{T}} A_{ii} X_{i}^{(k+1)}\right) + 2 \operatorname{tr}\left((X_{i}^{(k+1)})^{\mathrm{T}} G_{i} (X^{(k+\frac{i-1}{m})})\right) \\ &- \operatorname{tr}\left((X_{i}^{(k)})^{\mathrm{T}} E_{i} (X^{(k+\frac{i-1}{m})}) X_{i}^{(k)}\right) \\ &\geq \operatorname{tr}\left((X_{i}^{(k+1)})^{\mathrm{T}} A_{ii} X_{i}^{(k+1)}\right) + 2 \operatorname{tr}\left((X_{i}^{(k+1)})^{\mathrm{T}} G_{i} (X^{(k+\frac{i-1}{m})}) (X_{i}^{(k)})^{\mathrm{T}} X_{i}^{(k+1)}\right) \\ &- \operatorname{tr}\left((X_{i}^{(k)})^{\mathrm{T}} E_{i} (X^{(k+\frac{i-1}{m})}) X_{i}^{(k)}\right) \quad \text{(by (31))} \\ &= \operatorname{tr}\left((X_{i}^{(k+1)})^{\mathrm{T}} E_{i} (X^{(k+\frac{i-1}{m})}) X_{i}^{(k+1)}\right) - \operatorname{tr}\left((X_{i}^{(k)})^{\mathrm{T}} E_{i} (X^{(k+\frac{i-1}{m})}) X_{i}^{(k)}\right) \quad \text{(32)} \\ &\geq 0, \quad \text{(by (24))} \end{split}$$

where the last inequality is strict if (24) holds strictly. Hence

$$f(X^{(k+1)}) = f(X^{(k+\frac{m}{m})}) \ge f(X^{(k+\frac{m-1}{m})}) \ge \dots \ge f(X^{(k+\frac{1}{m})}) \ge f(X^{(k)}).$$
 (34)

Also,  $f(X^{(k+1)}) > f(X^{(k)})$  if one of the inequalities in (34) holds strictly. Since  $X_i^{(k+1)}$  is an orthonormal eigenbasis matrix of  $E_i(X^{(k+\frac{i-1}{m})})$  associated with its first  $\ell$  largest eigenvalues while  $X_i^{(k)}$  is an approximation, by Lemma 5.1 we get

$$\begin{split} & \xi_i^{(k)} \cdot \| \sin \Theta(X_i^{(k)}, X_i^{(k+1)}) \|_{\mathrm{F}}^2 \\ & \leq \sum_{j=1}^{\ell} \left( \lambda_j (E_i(X^{(k+\frac{i-1}{m})})) - \lambda_j ((X^{(k)})^{\mathrm{T}} E_i(X^{(k+\frac{i-1}{m})}) X^{(k)}) \right) \\ & = \operatorname{tr} \left( (X_i^{(k+1)})^{\mathrm{T}} E_i(X^{(k+\frac{i-1}{m})}) X_i^{(k+1)} \right) - \operatorname{tr} \left( (X_i^{(k)})^{\mathrm{T}} E_i(X^{(k+\frac{i-1}{m})}) X_i^{(k)} \right). \end{split}$$



Adv Comput Math (2022) 48:13 Page 15 of 34 13

Therefore, together with (32), we have<sup>2</sup>

$$f(X^{(k+\frac{i}{m})}) - f(X^{(k+\frac{i-1}{m})}) \geq \xi_i^{(k)} \cdot \|\sin\Theta(X_i^{(k)}, X_i^{(k+1)})\|_{\mathrm{F}}^2,$$

and (30) follows by summing it up from i = 1 to m.

For item (b), by assumption, we conclude from (34) that

$$f(X^{(k+\frac{i}{m})}) = f(X^{(k+\frac{i-1}{m})})$$

for all i = 1, 2, ..., m, which implies that (31) must become an equality. That is, for  $1 \le i \le m$ ,

$$\operatorname{tr}\left((X_i^{(k+1)})^{\mathrm{T}}G_i(X^{(k+\frac{i-1}{m})})(X_i^{(k)})^{\mathrm{T}}X_i^{(k+1)}\right) = \sum_{j=1}^{\ell} \sigma_j\left((X_i^{(k+1)})^{\mathrm{T}}G_i(X^{(k+\frac{i-1}{m})})\right).$$
(35)

Since  $(X_i^{(k+1)})^T G_i(X^{(k+\frac{i-1}{m})}) \geq 0$ , its spectral decomposition takes the form  $U \Sigma U^T$  where  $U \in \mathbb{O}^{\ell \times \ell}$  and  $\Sigma = \operatorname{diag}(\sigma_1, \ldots, \sigma_\ell)$  with  $\sigma_j = \sigma_j((X_i^{(k+1)})^T G_i(X^{(k+\frac{i-1}{m})})$ . We know all  $\sigma_j > 0$  because  $\operatorname{rank} \left( (X_i^{(k+1)})^T G_i(X^{(k+\frac{i-1}{m})}) \right) = \ell$ . Denote by  $Q = U^T(X_i^{(k)})^T X_i^{(k+1)} U$ . Equation (35) leads to

$$\sum_{j=1}^{\ell} \sigma_j \cdot (1 - Q_{(j,j)}) = 0 \quad \Rightarrow \quad Q_{(j,j)} = 1, \quad \forall j = 1, 2, \dots, \ell.$$

This shows that  $Q=I_\ell$  by the fact that  $X_i^{(k)}, X_i^{(k+1)} \in \mathbb{O}^{n_i \times \ell}$  and  $U \in \mathbb{O}^{\ell \times \ell}$  and thus  $1 \leq Q_{(j,j)} \leq 1$ . Hence  $X_i^{(k)} = X_i^{(k+1)}$  for  $1 \leq i \leq m$ , i.e.,  $X^{(k)} = X^{(k+1)}$ .

For the claim that  $X^{(k)}$  is a KKT point satisfying the necessary conditions in Theorems 1 and 4, we note that  $X^{(k+\frac{i-1}{m})} = X^{(k)} = X^{(k+1)}$  for  $1 \le i \le m$ . Because  $X_i^{(k+1)}$  is an orthonormal eigenbasis matrix associated with the first  $\ell$  largest eigenvalues of  $E_i(X^{(k+\frac{i-1}{m})}) = E_i(X^{(k)})$  and  $(X_i^{(k)})^T G_i(X^{(k)}) = (X_i^{(k+1)})^T G_i(X^{(k+\frac{i-1}{m})}) \ge 0$ , Theorem 2 guarantees that  $X^{(k)}$  is a KKT point, which also fulfills the necessary global optimality conditions in Theorems of 1 and 4.

Remark 3 In the proof of Theorem 6, we find that  $f(X^{(k+\frac{i}{m})}) - f(X^{(k+\frac{i-1}{m})}) \geq 0$  as long as  $\widehat{X}_i^{(k+1)}$  at line 3 of Algorithm 1 is a good enough approximation such that (33) holds true, and as a consequence the monotonic convergence of  $\{f(X^{(k)})\}$  can be ensured. This opens up the door to using modern iterative eigensolvers [3, 29, 37]. In our numerical results, we use LOBPCG³ for the task at line 3 of Algorithm 1, in which the (j+1)st iterate  $U_i^{(j+1)}$  (as an approximation of  $\widehat{X}_i^{(k+1)}$ ) of LOBPCG [22]



<sup>&</sup>lt;sup>2</sup>This relation also holds when  $\xi_i^{(k)} = 0$  because of (33).

<sup>&</sup>lt;sup>3</sup>https://cn.mathworks.com/matlabcentral/fileexchange/48-lobpcg-m.

13 Page 16 of 34 Adv Comput Math (2022) 48:13

for an approximate eigenbasis matrix corresponding to the first  $\ell$  largest eigenvalues of  $E_i(X^{(k+\frac{i-1}{m})})$  is obtained by solving

$$U_{i}^{(j+1)} \in \underset{U_{i} \in \mathcal{R}\left([U_{i}^{(j)}, U_{i}^{(j-1)}, PR^{(j)}]\right), \ U_{i} \in \mathbb{O}^{n_{i} \times \ell}}{\operatorname{argmax}} \operatorname{tr}(U_{i}^{T} E_{i}(X^{(k+\frac{i-1}{m})})U_{i})$$
(36)

where  $P \in \mathbb{R}^{n_i \times n_i}$  is a preconditioner, and  $R^{(j)} = E_i(X^{(k+\frac{i-1}{m})})\widehat{X}_k^{(j)} - \widehat{X}_k^{(j)}\Theta_j$  denotes the residual corresponding to the Ritz eigenpair  $(\widehat{X}_k^{(j)}, \Theta_j)$ . With initially  $U_i^{(0)} = X_i^{(k)}$ , (24) is obviously true, and thus the monotonic convergence of  $\{f(X^{(k)})\}$  follows.

For analyzing the convergence of  $\{X^{(k)}\}$ , the following lemma [16, Lemma 2.7] (see also [34, Lemma 4.10]) provides a sufficient condition that will be used.

**Lemma 5.2** Assume that X is an isolated accumulation point of a sequence  $\{X^{(k)}\}$  such that, for every subsequence  $\{X^{(k)}\}_{k\in\mathcal{K}}$  converging to X, there is an infinite subset  $\widehat{\mathcal{K}}\subseteq\mathcal{K}$  such that the sequence  $\{\|X^{(k)}-X^{(k+1)}\|_{\mathrm{F}}\}_{k\in\widehat{\mathcal{K}}}$  converges to 0. Then the whole sequence  $\{X^{(k)}\}_{k=0}^{\infty}$  converges to X.

Note that *X* being an isolated accumulation point does not exclude the existence of finitely many or even infinitely many accumulation points. If there are only finitely many accumulation points, then each must be isolated.

**Theorem 7** Let  $\{X^{(k)}\}$  be the sequence from Algorithm 1, and  $X = [X_1; ...; X_m] \in \mathcal{M}$  be an accumulation point of  $\{X^{(k)}\}$ . Suppose that

$$\xi_i := \lambda_\ell \big( E_i(X) \big) - \lambda_{\ell+1} \big( E_i(X) \big) > 0, \tag{37}$$

and rank  $(X_i^T \check{A}_i X) = \ell$  for i = 1, 2, ..., m, where  $\check{A}_i$  is defined by (10a).

- X is a KKT point satisfying the necessary global optimality conditions in Theorems 1 and 4;
- 2. If X is an isolated accumulation point  $\{X^{(k)}\}$ , then  $\lim_{k\to\infty} X^{(k)} = X$ .

*Proof* For item (a), suppose  $\{X^{(k)}\}_{k\in\mathcal{K}}$  is the subsequence that converges to X, and let  $\{X^{(k)}\}_{k\in\widehat{\mathcal{K}}}$  be a subsequence of  $\{X^{(k)}\}_{k\in\mathcal{K}}$  so that  $\{X^{(k+1)}\}_{k\in\widehat{\mathcal{K}}}$  converges to Z, i.e.,  $\widehat{\mathcal{K}}\subseteq\mathcal{K}$  and

$$\lim_{k \to \infty, \ k \in \mathcal{K}} X^{(k)} = X, \quad \lim_{k \to \infty, \ k \in \widehat{\mathcal{K}}} X^{(k+1)} = Z.$$

We have f(X) = f(Z). Now, for i = 1 in (26), Algorithm 1 says

$$E_1(X^{(k)})X_1^{(k+1)} = X_1^{(k+1)}\Psi_1^{(k+1)}$$
(38)



Adv Comput Math (2022) 48:13 Page 17 of 34 13

and in the limit, we have  $E_1(X)Z_1=Z_1\Psi_1$  and, because of  $\xi_1>0$ ,  $Z_1$  is an orthonormal eigenbasis matrix of  $E_1(X)$  associated with its  $\ell$  largest eigenvalues. Moreover, for any i,

$$\lim_{k \to \infty} f(X^{(k+\frac{i}{m})}) = \lim_{k \to \infty} f(X^{(k+\frac{i-1}{m})}) = f(X),$$

and by (33), we have

$$\lim_{k \to \infty, k \in \mathcal{K}} \left[ \text{tr}\left( (X_i^{(k+1)})^{\text{T}} E_i(X^{(k+\frac{i-1}{m})}) X_i^{(k+1)} \right) - \text{tr}\left( (X_i^{(k)})^{\text{T}} E_i(X^{(k+\frac{i-1}{m})}) X_i^{(k)} \right) \right] = 0, (39)$$

which for i=1 leads to  $\operatorname{tr}\left(Z_1^{\operatorname{T}}E_1(X)Z_1\right)=\operatorname{tr}\left(X_1^{\operatorname{T}}E_1(X)X_1\right)$ . This, together with  $\xi_1>0$ , ensures that  $X_1$  is an orthonormal eigenbasis matrix of  $E_1(X)$  associated with its  $\ell$  largest eigenvalues as well, and hence  $X_1=Z_1Q_1$  for some  $Q_1\in\mathbb{O}^{\ell\times\ell}$ . Similarly to the argument in concluding  $X^{(k)}=X^{(k+1)}$  in the proof of Theorem 6(b), we can use  $\operatorname{rank}\left(X_1^{\operatorname{T}}G_1(X)\right)=\operatorname{rank}\left(Z_1^{\operatorname{T}}G_1(X)\right)=\ell$  to prove  $X_1=Z_1$ . Applying the same argument for i=2 with  $\lim_{k\to\infty,\ k\in\mathcal{K}}X^{(k+\frac{1}{m})}\to X$ , we can obtain  $Z_2=X_2$ . Continue this procedure until i=m to get Z=X and (16). As  $X_i^{\operatorname{T}}G_i(X)\succeq 0$  for all  $i=1,2,\ldots,m$ , by Theorem 2 we conclude that X is a KKT point satisfying the necessary conditions for the global maximizer established in Theorems 1 and 4.

For item (b), since we have shown that any convergent subsequence  $\{X^{(k)}\}_{k\in\widehat{\mathcal{K}}}$  of  $\{X^{(k)}\}_{k\in\mathcal{K}}$  satisfies  $\lim_{k\to\infty,\,k\in\widehat{\mathcal{K}}}\|X^{(k)}-X^{(k+1)}\|_{\mathrm{F}}=0$ , Lemma 5.2 guarantees that the whole sequence  $\{X^{(k)}\}$  converges to X.

Let X be an accumulation point of  $\{X^{(k)}\}$ . In what follows, we will establish a bound on  $\|\sin\Theta(X_i^{(k)},X_i)\|_{\mathrm{F}}$  to reveal how the convergence of the ith block  $X_i^{(k)}$  depends on other blocks. The bound in Theorem 8 below is by no means tightest for  $\|\sin\Theta(X_i^{(k)},X_i)\|_{\mathrm{F}}$  but it is informative enough to guide practical eigen-computations at line 3 of Algorithm 1. Let

$$c_0 = ||A||_2 + \sqrt{m} \max_{1 \le i \le m} ||A_i||_2.$$

**Lemma 5.3** Let X be a KKT point of MAXBET (1). Then for any  $Y \in \mathcal{M}$ , we have

$$|f(Y) - f(X)| \le c_0 ||X - Y||_F^2.$$
 (40)

*Proof* Denote by  $R = Y - X = [R_1; ...; R_m]$ . We have

$$f(Y) - f(X) = \operatorname{tr}(R^{\mathsf{T}}AR) + 2\operatorname{tr}(R^{\mathsf{T}}AX) = \operatorname{tr}(R^{\mathsf{T}}AR) + 2\sum_{i=1}^{m} \operatorname{tr}(R_{i}^{\mathsf{T}}A_{i}X). \tag{41}$$

Since X is a KKT point, by Lemma 2.1 we have  $A_i X = X_i \Lambda_i$  and thus,  $\operatorname{tr}(R_i^{\mathrm{T}} A_i X) = \operatorname{tr}(R_i^{\mathrm{T}} X_i \Lambda_i)$ . Also, from the symmetry of  $\Lambda_i$  and

$$I_{\ell} = Y_i^{\mathrm{T}} Y_i = (X_i + R_i)^{\mathrm{T}} (X_i + R_i) = I_{\ell} + X_i^{\mathrm{T}} R_i + R_i^{\mathrm{T}} X_i + R_i^{\mathrm{T}} R_i,$$

13 Page 18 of 34 Adv Comput Math (2022) 48:13

we know that  $\operatorname{tr}(R_i^T X_i \Lambda_i) = -\operatorname{tr}(R_i^T R_i \Lambda_i)/2$ . With this, (41) leads to

$$f(Y) - f(X) = \operatorname{tr}(R^{\mathrm{T}}AR) - \sum_{i=1}^{m} \operatorname{tr}(R_{i}^{\mathrm{T}}R_{i}\Lambda_{i}).$$

The conclusion (40) hence follows by further noticing that (see Lemma 3.1)

$$|\operatorname{tr}(R^{T}AR)| \le ||A||_{2} \cdot ||R||_{F}^{2}, \quad ||\operatorname{tr}(R_{i}^{T}R_{i}\Lambda_{i})|| \le ||\Lambda_{i}||_{2} \cdot ||R_{i}||_{F}^{2}, \quad \sum_{i=1}^{m} ||R_{i}||_{F}^{2} = ||R||_{F}^{2},$$

and 
$$\|\Lambda_i\|_2 = \|X_i^T A_i X\|_2 \le \sqrt{m} \|A_i\|_2$$
.

**Theorem 8** *Under the assumptions of Theorem 7(b), for sufficiently large k, it holds that for i* = 1, 2, ..., m

$$\|\sin\Theta(X_{i}^{(k+1)}, X_{i})\|_{F} \le \frac{4\|\check{A}_{i}\|_{2} \left(\sqrt{m}\|X_{i} - X_{i}^{(k)}\|_{F} + \sqrt{\sum_{j < i} \|X_{j} - X_{j}^{(k+1)}\|_{F}^{2} + \sum_{j > i} \|X_{j} - X_{j}^{(k)}\|_{F}^{2}}\right)}{\xi_{i}}, \tag{42}$$

and

$$\|X^{(k)} - X\|_{F}^{2} \ge \sum_{i=1}^{m} \frac{\xi_{i}^{(k)}}{c_{0}} \cdot \|\sin\Theta(X_{i}^{(k)}, X_{i}^{(k+1)})\|_{F}^{2},\tag{43}$$

where  $\check{A}_i$  is defined by (10a),  $\xi_i$  and  $\xi_i^{(k)}$  are given by (37) and (29), respectively.

*Proof* The proof for (42) relies on a matrix perturbation result and follows closely to that of [8, Theorem 4.1]. For any  $1 \le i \le m$ , decompose  $X_i^{(k+1)} = X_i K_i + X_{i,\perp} L_i$ , where  $X_{i,\perp} \in \mathbb{O}^{n_i \times (n_i - \ell)}$  such that  $[X_i, X_{i,\perp}] \in \mathbb{O}^{n_i \times n_i}$ , and  $K_i \in \mathbb{R}^{\ell \times \ell}$ ,  $L_i \in \mathbb{R}^{(n_i - \ell) \times \ell}$ . Then  $\|L_i\|_F = \|\sin \Theta(X_i^{(k+1)}, X_i)\|_F$ . By calculations, for any  $Y \in \mathcal{M}$ , we have

$$\begin{split} \|E_{i}(X) - E_{i}(Y)\|_{F} &\leq \|X_{i} - Y_{i}\|_{F} \cdot (\|G_{i}(X)\|_{2} + \|G_{i}(Y)\|_{2}) + 2\|G_{i}(X) - G_{i}(Y)\|_{F} \\ &\leq \|X_{i} - Y_{i}\|_{F} \cdot (\|\check{A}_{i}X\|_{2} + \|\check{A}_{i}Y\|_{2}) + 2\|\check{A}_{i}(X - Y)\|_{F} \\ &\leq 2\sqrt{m}\|X_{i} - Y_{i}\|_{F} \cdot \|\check{A}_{i}\|_{2} + 2\|\check{A}_{i}\|_{2} \cdot \|X - Y\|_{F} \\ &= 2\|\check{A}_{i}\|_{2}(\sqrt{m}\|X_{i} - Y_{i}\|_{F} + \|X - Y\|_{F}). \end{split}$$

Let  $\Delta E_i = E_i(X^{(k+\frac{i-1}{m})}) - E_i(X)$ . Setting  $Y = X^{(k+\frac{i-1}{m})}$  gives

$$\begin{split} \|\Delta E_i\|_{\mathrm{F}} &= \|E_i(X^{(k+\frac{i-1}{m})}) - E_i(X)\|_{\mathrm{F}} \\ &\leq 2\|\check{A}_i\|_2 \left(\sqrt{m}\|X_i - X_i^{(k)}\|_{\mathrm{F}} + \sqrt{\sum_{j < i} \|X_j - X_j^{(k+1)}\|_{\mathrm{F}}^2 + \sum_{j > i} \|X_j - X_j^{(k)}\|_{\mathrm{F}}^2}\right). \tag{44} \end{split}$$

Since  $E_i(X)X_i = X_i\Psi_i$  for the limit X and  $X_i^{(k+1)}$  satisfies the first equation in (26), we get

$$E_i(X)X_i^{(k+1)} - X_i^{(k+1)}\Psi_i^{(k+1)} = \left[E_i(X) - E_i(X^{(k+\frac{i-1}{m})})\right]X_i^{(k+1)}.$$

Pre-multiplying both sides of this equation by  $X_{i,\perp}^{\mathrm{T}}$ , we get

$$\Psi_{i,\perp} L_i - L_i \Psi_i^{(k+1)} = -X_{i,\perp}^{\mathrm{T}} \Delta E_i X_i^{(k+1)}, \tag{45}$$



Adv Comput Math (2022) 48:13 Page 19 of 34 13

where  $\Psi_{i,\perp} = X_{i,\perp}^T E_i(X) X_{i,\perp}$  whose eigenvalues are given by  $\operatorname{eig}(E_i(X)) \setminus \operatorname{eig}(\Psi_i)$ . Observe that  $L_i$  in (45) is a solution to Sylvester (45), and by [12, (5.1)], we find

$$||L_i||_{\mathcal{F}} = ||\sin\Theta(X_i^{(k+1)}, X_i)||_{\mathcal{F}} \le \frac{||X_{i, \perp}^{\mathsf{T}} \Delta E_i X_i^{(k+1)}||_{\mathcal{F}}}{\hat{\xi}_i^{(k)}} \le \frac{2||\Delta E_i||_{\mathcal{F}}}{\xi_i},$$

where  $\hat{\xi}_i^{(k)} = \|\lambda_\ell(\Psi_i^{(k+1)}) - \lambda_1(\Psi_{i,\perp})\|$  is the gap between the eigenvalues of  $\Psi_i^{(k+1)}$  and  $\Psi_{i,\perp}$ , which satisfies  $\hat{\xi}_i^{(k)} \geq \hat{\xi}_i/2$  for sufficiently large k. Consequently, (42) follows according to (44).

For (43), by applying Lemma 5.3 with  $Y = X^{(k)}$  and (30), we get

$$\sum_{i=1}^{m} \xi_{i}^{(k)} \| \sin \Theta(X_{i}^{(k)}, X_{i}^{(k+1)}) \|_{F}^{2} \le f(X^{(k+1)}) - f(X^{(k)})$$

$$\le f(X) - f(X^{(k)}) \le c_{0} \|X - X^{(k)}\|_{F}^{2}.$$

as was to be shown.

We observe from (42) that the convergence of the ith block  $X_i^{(k)}$  depends on (i)  $\|\check{A}_i\|_2$ , (ii) the eigenvalue gap  $\xi_i$ , and (iii) the convergence of other blocks  $X_j^{(k)}$ . For (i), it reflects the fact that when  $A = \operatorname{diag}(A_{11}, \ldots, A_{mm})$  (i.e.,  $\check{A}_i = 0$ ), the convergence is achieved in one outer iteration step. The factor  $\xi_i$  in (ii) says that the larger  $\xi_i$  is, the faster the convergence will be. For (iii), it implies that the slow convergence of  $X_j^{(k)}$  for  $j \neq i$  can spread to  $X_i^{(k)}$ , and therefore, if  $\|\check{A}_j\|_2$  is relatively large, we need to set a reasonably high accuracy requirement for the eigencomputation for  $X_i^{(k+1)}$ .

Note that the right-hand side of (43) can be computed after the (k + 1)st outer iteration step, providing a lower bound for the uncomputable error  $||X^{(k)} - X||_F$ .

# 6 Subspace accelerated SCF

For a practical implementation, we accelerate the SCF iteration (Algorithm 1) by a subspace acceleration technique. Similarly to [56], we collect the information from the previous t ( $t \le \varsigma$  and  $\varsigma > 0$ ) outer-loop iterates and  $X^{(k+1)}$  immediately after the inner loop to form a subspace  $\mathcal C$  for acceleration. Specifically, immediately after line 5 but before line 6 of Algorithm 1, define

$$C = [X^{(k-t+1)}, \dots, X^{(k)}, X^{(k+1)}] \in \mathbb{R}^{n \times (t+1)\ell},$$
  

$$C_i = [X_i^{(k-t+1)}, \dots, X_i^{(k)}, X_i^{(k+1)}],$$

and let

$$C_i = \mathcal{R}(C_i)$$
 and  $C = C_1 \oplus \ldots \oplus C_m$ .

We then refine the current approximation  $X^{(k+1)}$  by maximizing f(X) over  $X \in \mathcal{M}$  subject to  $\mathcal{R}(X_i) \subseteq \mathcal{C}_i$  for all i:

$$X^{(k+1)} \leftarrow \underset{X \in \mathcal{M}, \, \mathcal{R}(X_i) \subseteq \mathcal{C}_i \, \forall i}{\operatorname{argmax}} f(X). \tag{46}$$

13 Page 20 of 34 Adv Comput Math (2022) 48:13

To solve (46), we first compute orthonormal basis matrix  $Q_i^{(k+1)} \in \mathbb{O}^{n_i \times r_i}$  of  $C_i$ , and represent any  $X_i \in \mathbb{O}^{n_i \times \ell}$  such that  $\mathcal{R}(X_i) \subseteq C_i$  by  $Q_i^{(k+1)} V_i$  for  $V_i \in \mathbb{O}^{r_i \times \ell}$ . We have

$$f(X) = \sum_{i=1}^{m} \sum_{j=1}^{m} \operatorname{tr}\left(V_{i}^{T} (Q_{i}^{(k+1)})^{T} A_{ij} Q_{j}^{(k+1)} V_{j}\right) = \operatorname{tr}(V^{T} \mathring{A} V) := \mathring{f}(V)$$

where  $\mathring{A} = [\mathring{A}_{ij}] \in \mathbb{R}^{r \times r}$  with  $\mathring{A}_{ij} = (Q_i^{(k+1)})^T A Q_j^{(k+1)} \in \mathbb{R}^{r_i \times r_j}$  and  $r = \sum_{i=1}^m r_i$ . The optimal  $V = [V_1; \ldots; V_m]$  can be found by solving

$$\max_{\{V_i \in \mathbb{O}^{r_i \times \ell}\}} \mathring{f}(V). \tag{47}$$

Problem (47) is still a MAXBET problem in the form of (1) but of a much smaller size, as  $r_i$  is generally far less than  $n_i$ . Because of that, we can use the Riemannian trust-region method [2, 51] without inflicting much computational burden. Once optimal  $\{V_i\}$  of (47) is found, we construct a new refined iterate  $X^{(k+1)}$  with block  $X_i^{(k+1)} = Q_i^{(k+1)} V_i$  for  $i = 1, \ldots, m$ .

Equipped with this subspace acceleration, we propose our practical version of the SCF iteration (named as Scf\_Maxbet) for (1) as in Algorithm 2.

# Algorithm 2 Scf\_Maxbet: subspace accelerated SCF for MAXBET (1).

**Require:** symmetric  $A \in \mathbb{R}^{n \times n}$ ,  $X^{(0)} \in \mathcal{M}$ , and  $\zeta > 0$ ;

**Ensure:** a maximizer of MAXBET (1).

- 1: Set C = [] and t = 0.
- 2: **for**  $k = 1, \ldots$ , until convergence **do**
- 3: **for** i = 1, ..., m **do**
- Compute (by e.g., LOBPCG) an approximate orthonormal eigenbasis matrix  $\widehat{X}_i^{(k+1)}$  associated with the  $\ell$  largest eigenvalues of  $E_i(X^{(k+\frac{i-1}{m})})$ , accurate enough so that

$$\operatorname{tr}\left((\widehat{X}_{i}^{(k+1)})^{\mathrm{T}}E_{i}(X^{(k+\frac{i-1}{m})})\widehat{X}_{i}^{(k+1)}\right) > \operatorname{tr}\left((X_{i}^{(k)})^{\mathrm{T}}E_{i}(X^{(k+\frac{i-1}{m})})X_{i}^{(k)}\right);$$

- 5: Compute SVD:  $(\widehat{X}_{i}^{(k+1)})^{T}G_{i}(X^{(k+\frac{i-1}{m})})$  and set  $X_{i}^{(k+1)} = \widehat{X}_{i}^{(k+1)}U_{i}V_{i}^{T};$
- 6: end for
- 7: **if** t > 1 **then**
- 8: Set  $C = [C, X^{(k+1)}]$  and solve (47) to get a refined  $X^{(k+1)}$ ;
- 9: Replace the last  $\ell$  columns of C by the refined  $X^{(k+1)}$ ;
- 10: **end if**
- if  $t = \zeta$  then delete the first  $\ell$  columns of C else t = t + 1 end if;
- 12: **end forreturn** the last  $X^{(k)}$ .



Adv Comput Math (2022) 48:13 Page 21 of 34 13

# 7 Numerical experiments

In this section, we present numerical experiments on synthetic problems and real-world applications to evaluate the performance of the proposed algorithm Scf\_Maxbet and compare it with the basic SCF (Algorithm 1) and the Riemannian trust-region method (RTR) [2, 51], a generic solver for optimization on a general Riemannian manifold. All tests were conducted by MATLAB (R2016a) on a PC under Windows 7 (64bit) system with Intel(R) Core(TM) i7-4510U CPU (2.6 GHz) and 4 GB memory.

To manage the cost for partial eigen-decompositions in Algorithm 2, we invoke LOBPCG (see [22]) at line 4 to compute  $\widehat{X}_i^{(k+1)}$  if  $n_i \geq 500$ , otherwise we simply call the MATLAB built-in function eig. We terminate the basic SCF and the accelerated SCF iteration Scf\_Maxbet when one of

$$\begin{split} \operatorname{res} &= \frac{\|AX^{(k)} - X^{(k)}\Lambda^{(k)}\|_1}{\|A\|_1} \, \leq \, \varepsilon_{\operatorname{res}} = 10^{-6}, \\ &\frac{|\, f(X^{(k-1)}) - f(X^{(k)})\,|}{|\, f(X^{(k-1)})\,|} \, \leq \, \varepsilon_{\operatorname{fun}} = 10^{-15} \end{split}$$

is satisfied or the number of outer-loop iterations reaches 3000.

## 7.1 Evaluation on synthetic problems

We first report numerical results of the basic SCF, Scf\_Maxbet and RTR for solving (1) on synthetic correlation matrices A. For RTR, we note that it stops when the Frobenius norm of the Riemannian gradient is less than a tolerance  $\varepsilon_{\text{rtr}}$ . To be consistent with our stopping criteria, we choose a relaxed tolerance  $\varepsilon_{\text{rtr}} = 5 \times 10^{-5}$  for RTR and also set the maximum number of RTR iterations to 3000; all other parameters in RTR are the default ones.

We use the following four ways to generate the random correlation matrix A:

- corr. The MATLAB built-in function corr(n) produces a random  $n \times n$  correlation matrix;
- gall. The MATLAB built-in function gallery ('randcorr', n) generates a random  $n \times n$  correlation matrix with random eigenvalues from a uniform distribution [5, 11];
- ksih. The function ksih (n) generates a random correlation matrix having the distribution  $\Psi_n(n)$  [42];
- rand. The function<sup>4</sup> randcorr (n) [38] generates a random  $n \times n$  correlation matrix.

#### 7.1.1 Different feature sizes n

In this subsection, we use functions corr(n), gallery('randcorr', n), ksih(n), and randcorr(n) to generate a number of instances with n ranging



<sup>&</sup>lt;sup>4</sup>https://www.mathworks.com/matlabcentral/fileexchange/68810-randcorr.

13 Page 22 of 34 Adv Comput Math (2022) 48:13

from 2000 to 7000,  $\ell = 10$ , m = 2, and  $n_i = \frac{n}{m}$  for i = 1, 2. Numerical results in terms of CPU time and the scaled norm of the gradient (denoted by res) averaged over 5 random correlation matrices are plotted in Figs. 1 and 2.

From Figs. 1 and 2, we have the following observations: (i) the basic SCF, Scf\_Maxbet and RTR require more CPU time as n increases, as expected, but Scf\_Maxbet outperforms the basic SCF and RTR for each given n. (ii) Scf\_Maxbet achieves more accurate approximations in terms of res for all instances generated by corr(n) and gallery('randcorr',n), and for those generated by ksih(n) and randcorr(n), RTR and Scf\_Maxbet have roughly the same accuracy (residuals res from RTR in Fig. 2 are calculated at their computed solutions). (iii) The accelerated SCF Scf\_Maxbet improves significantly both the accuracy and the efficiency over the basic SCF; in particular, it is observed that in most cases, the basic SCF iteration cannot achieve res  $\leq 10^{-6}$  within the maximal number of iterations. Therefore, in our subsequent testing, we only report the results of Scf\_Maxbet and RTR.

#### 7.1.2 Different reduced dimensions $\ell$

With n=6000, m=2 and  $n_i=\frac{n}{m}$  for i=1,2, Scf\_Maxbet and RTR are now evaluated with  $\ell$  varying from 10 to 20 for correlation matrices A generated by

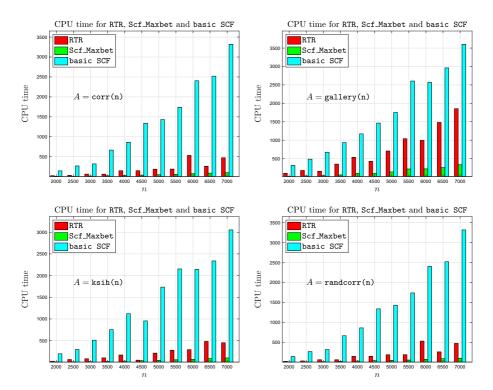


Fig. 1 CPU time for RTR, Scf\_Maxbet and basic SCF with varying n



Adv Comput Math (2022) 48:13 Page 23 of 34 13

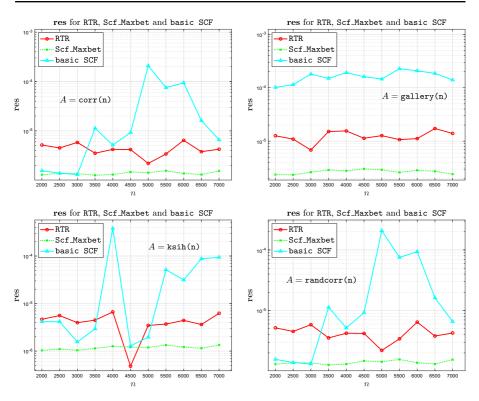


Fig. 2 Residual res for RTR, Scf\_Maxbet, and basic SCF with varying n

gallery, rand, ksih, and corr. Numerical results on these random instances are plotted in Figs. 3 and 4. The results show that Scf\_Maxbet consumes less CPU time and yet achieves more accurate approximate solutions.

#### 7.1.3 Different numbers of views m

Next, we evaluate the performances of Scf\_Maxbet and RTR with different numbers of views. For  $m \in \{2, 3, 4, 5, 6, 8, 10\}$ , n = 6000,  $\ell = 10$ , and  $n_i = \frac{n}{m}$ , for i = 1, 2, ..., m, we generate A as before and report the numerical results by Scf\_Maxbet and RTR in Figs. 5 and 6. It shows that as m increases, especially for the instances generated by gallery ('randcorr', n), the efficiency of Scf\_Maxbet slightly dwindles, partially due to the increase of the number of for-loops: lines 3-6 in Algorithm 2. Nevertheless, Scf\_Maxbet overall is still more efficient than RTR and is suitable for solving (1) particularly when  $n \gg m\ell$ .

#### 7.2 An application to multi-view feature extraction

In many recent applications, real-world data points frequently are represented by multiple sets of features which usually reflect various characteristics of the targeted object. In this subsection, we will apply Scf\_Maxbet for unsupervised multi-view



13 Page 24 of 34 Adv Comput Math (2022) 48:13

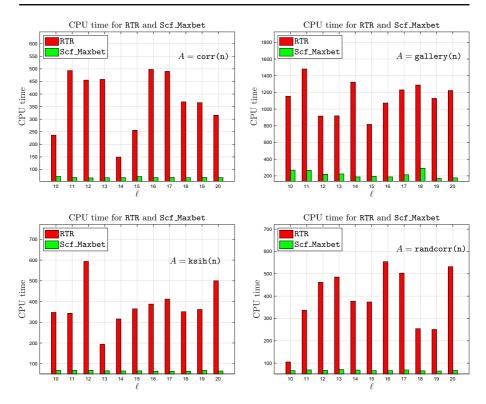


Fig. 3 CPU time for RTR and Scf\_Maxbet with various  $\ell$ 

feature extraction [45]. The number of features is dramatically reduced much fewer and the extracted features are often used later for following up learning tasks much more efficiently.

Suppose  $S \in \mathbb{R}^{n \times q}$  is a centralized multi-view data matrix that consists of q random samples. The matrix A in (1) is one of the following:

- $A = SS^{T}$  (the corresponding MAXBET (1) is known as SUMCOR [21, 54]);
- $A = SS^{T} \text{diag}(A_{11}, \dots, A_{mm})$  (the corresponding MAXBET (1) is known as MPLS [35]).

Note that if we ignore the multi-view structure embedded in the data, we can simply use the principal component analysis (PCA) to extract multiple features. In particular, PCA extracts first  $\ell$  principal components obtained by a proper orthogonal projection matrix  $X \in \mathbb{O}^{n \times \ell}$ . Note that MAXBET (1) reduces to PCA if m = 1. As a baseline, we will also report learning results by PCA.

Another closely related model in MCCA (multi-view canonical correlation analysis) is the successive deflation MCP [54] which serves as an approximate model of MAXBET (1). Specially, MCP first solves (1) with  $\ell = 1$  to obtain a vector  $\boldsymbol{x}_1$ , and then uses a successive deflation scheme to obtain the subsequent  $\boldsymbol{x}_2, \ldots, \boldsymbol{x}_\ell$ . The final  $X = [\boldsymbol{x}_1, \ldots, \boldsymbol{x}_\ell]$  collects these vectors and satisfies  $X \in \mathcal{M}$ . It should be pointed out that the projection matrix X from the successive deflation MCP is



Adv Comput Math (2022) 48:13 Page 25 of 34 13

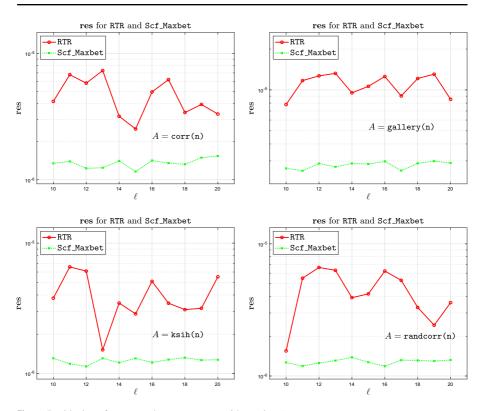


Fig. 4 Residual res for RTR and Scf\_Maxbet with varying  $\ell$ 

generally not the solution to the MAXBET problem when  $\ell > 1$ . We will also present the results from the successive deflation MCP [54] to compare with the performance of MAXBET.

Numerically, the projection matrix X from PCA can be simply obtained by the SVD of S. For the projection matrix X from the successive deflation MCP [54], we use two approaches, namely, the nested Lanczos-type iteration (named as nLMcp) [54] and the Riemannian trust-region method (named as RTRMCP) [2, 51] with tolerance  $10^{-4}$  for the associated residuals, to solve MCP resulting from each deflation step. For MAXBET (1), we apply our proposed algorithm Scf\_Maxbet with  $\varepsilon_{\text{res}} = 10^{-4}$ .

#### 7.2.1 Datasets

Datasets used in our experiment are summarized in Table 1 [10, 55]: the first one is a synthetic dataset while the others are publicly available image datasets<sup>5</sup>. We

<sup>&</sup>lt;sup>5</sup>The datasets Caltech101-7 and Caltech101-20 [26] are available at https://www.vision.caltech.edu/ Image\_Datasets/Caltech101/, and the dataset Scene15 [25] is available at figshare.com/articles/15-Scene\_Image\_Dataset/7007177.

13 Page 26 of 34 Adv Comput Math (2022) 48:13

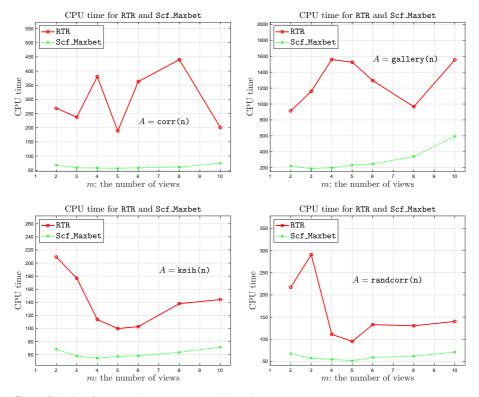


Fig. 5 CPU time for RTR and Scf\_Maxbet with varying m

generate the synthetic data in a similar way to [10, 55]. Specifically,  $S = [S_1; ...; S_m]$  with

- c = 5 (i.e., the number of classes), q = 2000 samples (i.e., each class has 400 samples), m = 5 (i.e., the number of multi-views), and  $n_i = 2000$  features for each view, and thus  $n = \sum_{i=1}^{m} n_i = 10000$ ;
- The ith view data matrix  $S_i = P_i Z_i + \sigma E \in \mathbb{R}^{n_i \times q}$ , where  $P_i \in \mathbb{R}^{n_i \times 100}$  whose entries are i.i.d. sampled from Gaussian distribution N(i/10, 1), noise level  $\sigma = 10^{-3}$ , and  $E \in \mathbb{R}^{n_i \times q}$  whose entries are i.i.d. sampled from normal distribution N(0, 1), matrix  $Z_i = [Z_i^{[1]}, \dots, Z_i^{[c]}] \in \mathbb{R}^{100 \times q}$ , and entries of  $Z_i^{[j]} \in \mathbb{R}^{100 \times 400}$  are i.i.d. and are sampled from a Gaussian distribution N((j-1)/5, 1) for  $1 \le j \le c = 5$ .

#### 7.2.2 Performances

To demonstrate the effectiveness of feature extraction by each approach, we rely on supervised learning, the 1-nearest neighbor (1NN) classifier, on the extracted features by the learned projection matrix X, as follows:



Adv Comput Math (2022) 48:13 Page 27 of 34 13

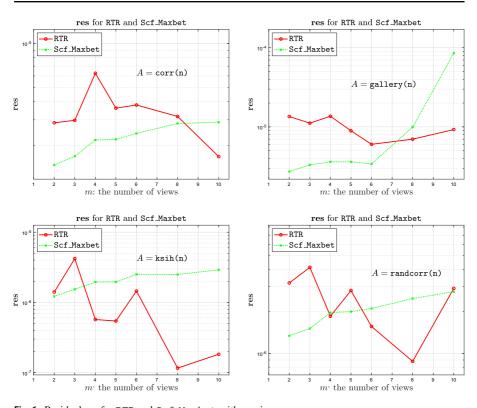


Fig. 6 Residual res for RTR and Scf\_Maxbet with varying m

- 1. Each dataset is split into training and testing data portions (30% training and 70% testing);
- 2. Projection matrices *X* are trained by the training data by PCA, Scf\_Maxbet, nLMcp and RTRMCP, respectively;
- 3. Data points of each view for both training and testing are projected onto a reduced common space via the projection matrices *X*;
- 4. Labels of testing data points are predicted via the 1-nearest neighbor (1NN) classifier on the projected testing data points.

Table 1 Multi-view datasets

Dataset	Samples (q)	m	Multiple views $(n_i)$	Features (n)	Classes (c)
Synthetic	2000	5	2000;2000;2000;2000;2000	10,000	5
Caltech101-7	1474	6	254;512;1180;1008;64;1000	4018	7
Caltech101-20	2386	6	254;512;1180;1008;64;1000	4018	20
Scene15	4310	6	254;512;531;360;64;1000	2721	15



 Table 2
 Means (std) of the classification accuracy on synthetic datasets

Dataset		accuracy_sumcor	лс		accuracy-plst			accuracy_pca
	в	Scf_Maxbet	RTRMCP	пГМср	Scf_Maxbet	RTRMCP	пГМср	PCA
Synthetic 5	5 10	77.09%(0.0307) 70.04%(0.0619)	59.57%(0.0094) 49.32%(0.0065)	59.64%(0.0093) 49.33%(0.0081)	70.93%(0.0606) 71.75%(0.0312)	59.64%(0.0048) 49.17%(0.0051)	59.62%(0.0050) 49.17%(0.0051)	74.18%(0.0051) 54.49%(0.0065)



Adv Comput Math (2022) 48:13 Page 29 of 34 13

Table 3 CPU time on synthetic datasets

Dataset		CPU_sumcor			CPU_plst		
Synthetic	<ul><li>ℓ</li><li>5</li><li>10</li></ul>	Scf_Maxbet 65.01 41.72	RTRMCP 607.47 1307.68	nLMcp 152.79 284.30	Scf_Maxbet 33.24 27.42	RTRMCP 557.68 1244.19	nLMcp 106.41 192.75

We report both classification accuracy (i.e., the percentage of correct predictions over the testing data) and CPU time for computing *X* as measures of performance. We use "accuracy\_sumcor" and "accuracy\_plst" to refer to the average accuracy corresponding to the two different choices of *A* mentioned at the beginning of sub-Section 7.2, and accordingly "CPU\_sumcor" and "CPU\_plst" are for the average CPU times. Label "accuracy\_pca" is for the classification accuracy by PCA.

For synthetic data, we report means and standard deviations (std) of classification accuracy and average CPU time over 10 random splittings for each dataset (30% training and 70% testing) in Tables 2 and 3. It can be seen that Scf\_Maxbet consistently outperforms other methods in terms of classification accuracy and CPU time. We did not show CPU time by PCA because it just calls MATLAB's  ${\tt svd}$  on  ${\tt S}$ .

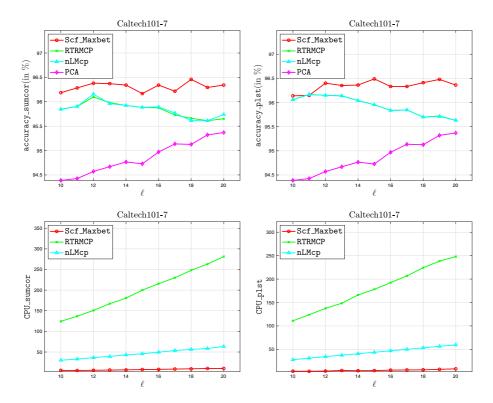
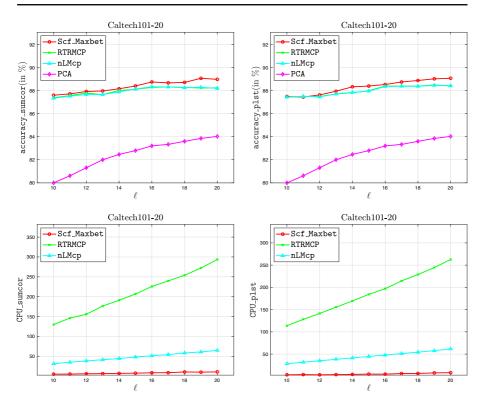


Fig. 7 Results for Caltech101-7 with varying  $\ell$ 

13 Page 30 of 34 Adv Comput Math (2022) 48:13



**Fig. 8** Results for Caltech101-20 with varying  $\ell$ 

For real-world data, we plot average CPU time and classification accuracy over 10 random draws in Figs. 7, 8, and 9. We observe that (i) Scf\_Maxbet achieves better accuracies on both SUMCOR and PLST, (ii) Scf\_Maxbet consumes much less CPU time, and (iii) in contrast to nLMcp and RTRMCP, the variation of CPU time with different  $\ell$  by Scf\_Maxbet is relatively moderate.

# 8 Concluding remarks

From a new perspective of maximizing a homogeneous quadratic function over the product of Stiefel manifolds, in this paper, we first reformulate the MAXBET problem as a coupled NEPv. The new formulation naturally leads itself to an alternating SCF iteration in which the blocks of projection matrix X are updated alternatingly in a Gauss-Seidel fashion. We have developed a theory for the relation between MAXBET and this coupled NEPv, and also established convergence results for the Gauss-Seidel-type updating scheme combined with the SCF iteration. The plain SCF iteration is then accelerated by a subspace acceleration strategy, and our preliminary numerical evaluations demonstrate that the subspace-accelerated SCF iteration is more efficient than the generic Riemannian trust-region solver RTR for problems



Adv Comput Math (2022) 48:13 Page 31 of 34 13

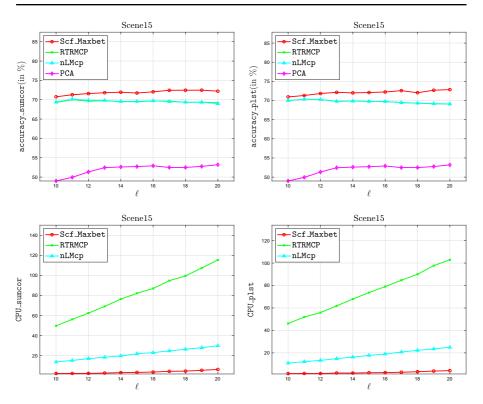


Fig. 9 Results for Scene 15 with varying  $\ell$ 

with large block sizes  $n_i$ . Finally, we use an application of MAXBET to unsupervised feature extraction learning to show the effectiveness of the MAXBET criterion in MCCA and the efficiency of the alternating SCF iteration for the MAXBET problem.

**Acknowledgements** The authors are grateful to the anonymous referees for their useful comments and suggestions to improve the presentation of this paper.

**Funding** The work of the third author was supported partially by NSF DMS-2009689 and NIH R01AG075582; the work of the fourth author was supported partially by the National Natural Science Foundation of China NSFC-12071332; and the work of the fifth author was supported partially by NSF DMS-1719620 and DMS-2009689.

#### **Declarations**

**Conflict of interest** The authors declare no competing interests.

#### References

 Absil, P.-A., Baker, C.G., Gallivan, K.A.: Trust-region methods on Riemannian manifolds. Found. Comput. Math. 7(3), 303–330 (2007)



13 Page 32 of 34 Adv Comput Math (2022) 48:13

 Absil, P.-A., Mahony, R., Sepulchre, R.: Optimization Algorithms on Matrix Manifolds. Princeton University Press, Princeton (2008)

- 3. Bai, Z., Demmel, J.W., Dongarra, J., Ruhe, A., Vorst, H., Van Der Vorst, H.: Templates for the Solution of Algebraic Eigenvalue Problems: A Practical Guide. SIAM, Philadelphia (2000)
- Bai, Z., Lu, D., Vandereycken, B.: Robust Rayleigh quotient minimization and nonlinear eigenvalue problems. SIAM J. Sci. Comput. 40, A3495–A3522 (2018)
- 5. Bendel, R.B., Mickey, M.R.: Population correlation matrices for sampling experiments. Comm Statist. Simul. Comput. **7**(2), 163–182 (1978)
- 6. Ten Berge, J.M.F.: Generalized approaches to the MAXBET problem and the MAXDIFF problem, with applications to canonical correlations. Psychometrika **53**(4), 487–494 (1988)
- 7. Ten Berge, J.M.F., Knol, D.L.: Orthogonal rotations to maximal agreement for two or more matrices of different column orders. Psychometrika **49**(4), 49–55 (1984)
- 8. Cai, Y., Zhang, L.-H., Bai, Z., Li, R.-C.: On an eigenvector-dependent nonlinear eigenvalue problem. SIAM J. Matrix Anal. Appl. **39**(3), 1360–1382 (2018)
- Chu, M.T., Watterson, J.L.: On a multivariate eigenvalue problem, Part I: Algebraic theory and a power method. SIAM J. Sci. Comput. 14(4), 1089–1106 (1993)
- Cunningham, J.P., Ghahramani, Z.: Linear dimensionality reduction: Survey, insights, and generalizations. J. Mach. Lear. Res. 16, 2859–2900 (2015)
- Davies, P.I., Higham, N.J.: Numerically stable generation of correlation matrices and their factors. BIT 40, 640–651 (2000)
- 12. Davis, C., Kahan, W.: The rotation of eigenvectors by a perturbation. III. SIAM J. Numer. Anal. 7, 1–46 (1970)
- 13. Van de Geer, J.P.: Linear relations among k sets of variables. Psychometrika 49, 70–94 (1984)
- Goemans, M.X., Williamson, D.P.: Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. J. ACM 42, 1115–1145 (1995)
- 15. Golub, G.H., Van Loan, C.F.: Matrix Computations, 4th edn. Johns Hopkins University Press, Baltimore (2013)
- Guan, Y., Chu, M.T., Chu, D.: SVD-based algorithms for the best rank-1 approximation of a symmetric tensor. SIAM J. Matrix Anal. Appl. 39, 1095–1115 (2018)
- 17. Hanafi, M., Ten Berge, J.M.F.: Global optimality of the successive Maxbet algorithm. Psychometrika 68, 97–103 (2003)
- 18. Hardoon, D.R., Szedmak, S., Shawe-Taylor, J.: Canonical correlation analysis: An overview with application to learning methods. Neural Comput. 16, 2639–2664 (2004)
- 19. Horst, P.: Relations among m sets of measures. Psychometrika 26, 129–149 (1961)
- 20. Hotelling, H.: Relations between two sets of variates. Biometrika 28, 321–377 (1936)
- 21. Kettenring, J.R.: Canonical analysis of several sets of variables. Biometrika 58(1), 433–451 (1971)
- Knyazev, A.V.: Toward the optimal preconditioned eigensolver: Locally optimal block preconditioned conjugate gradient method. SIAM J. Sci. Comput. 23(2), 517–541 (2001)
- Knyazev, A.V., Argentati, M.E.: Rayleigh-Ritz majorization error bounds with applications to FEM. SIAM J Matrix Anal Appl 31(3), 1521–1537 (2010)
- Kovač-Striko, J., Veselić, K.: Some remarks on the spectra of Hermitian matrices. Linear Algebra Appl. 145, 221–229 (1991)
- Lazebnik, S., Schmid, C., Ponce, J.: Beyond bags of features: Spatial pyramid matching for recognizing natural scene categories. In: Proceedings of the 2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06), pp. 2169–2178 (2006)
- Li, F.F., Fergus, R., Perona, P.: Learning generative visual models from few training examples: An incremental bayesian approach tested on 101 object categories. Comput. Vis. Image Underst. 106, 59–70 (2007)
- Li, R.-C.: New perturbation bounds for the unitary polar factor. SIAM J. Matrix Anal. Appl. 16, 327–332 (1995)
- 28. Li, R.-C.: Accuracy of computed eigenvectors via optimizing a Rayleigh quotient. BIT **44**(3), 585–593 (2004)
- Li, R.-C.: Rayleigh quotient based optimization methods for eigenvalue problems. In: Bai, Z., Gao, W., Su, Y. (eds.) Matrix Functions and Matrix Equations, volume 19 of Series in Contemporary Applied Mathematics, pp. 76–108. World Scientific, Singapore (2015)
- Li, Z., Nie, F., Chang, X., Yang, Y.: Beyond trace Weighted Harmonic mean of trace ratios for multiclass discriminant analysis. IEEE Trans. Knowl. Data Engrg. 29(10), 2100–2110 (2017)



Adv Comput Math (2022) 48:13 Page 33 of 34 13

31. Lingoes, J.C., Borg, I.: A direct approach to individual differences scaling using increasingly complex transformations. Psychometrika **43**, 491–519 (1978)

- Liu, X.-G., Wang, X.-F., Wang, W.-G.: Maximization of matrix trace function of product S,tiefel manifolds. SIAM J Matrix Anal. Appl. 36(4), 1489–1506 (2015)
- Martin, R.M.: Electronic structure: Basic theory and practical methods. Cambridge University Press, Cambridge (2004)
- 34. Moré, J.J., Sorensen, D.C.: Computing a trust region step. SIAM J. Sci. Statist. Comput. 4(3), 553–572 (1983)
- 35. Nielsen, A.A.: Multiset canonical correlations analysis and multispectral, truly multitemporal remote sensing data. IEEE Trans. Image Process. 11(3), 293–305 (2002)
- 36. Nocedal, J.: Numerical Optimization, 2nd edn. Springer, New York (2006)
- 37. Parlett, B.N.: The symmetric eigenvalue problem. SIAM. Philadelphia (1998)
- 38. Pourahmadi, M., Wang, X.: Distribution of random correlation matrices: Hyperspherical parameterization of the Cholesky factor. Statist. Probab. Lett. **106**(3), 5–12 (2015)
- Saad, Y., Chelikowsky, J.R., Shontz, S.M.: Numerical methods for electronic structure calculations of materials. SIAM Rev. 52(1), 3–54 (2010)
- 40. Seber, G.A.F.: A Matrix Handbook for Statisticians. Wiley, New Jersey (2007)
- 41. Stewart, G.W.: Matrix Algorithms, Vol. II: Eigensystems. SIAM, Philadelphia (2001)
- 42. Veleva, E.: Generation of correlation matrices. AIP Conf. Proc. 1895(120008), 1–7 (2017)
- von Neumann, J.: Some matrix-inequalities and metrization of matrix- space. Tomck. Univ Some Rev. 1, 286–300 (1937)
- 44. Wang, L., Li, R.-C.: A scalable algorithm for large-scale unsupervised multi-view partial least squares. IEEE Trans. Big Data. https://doi.org/10.1109/TBDATA.2020.3014937, to appear (2020)
- 45. Wang, L., Zhang, L.-H., Bai, Z., Li, R.-C.: Orthogonal canonical correlation analysis and applications. Optim. Methods Softw. **35**(4), 787–807 (2020)
- Wang, Z., Ruan, Q., An, G.: Projection-optimal local Fisher discriminant analysis for feature extraction. Neural Comput. Appl. 26, 589–601 (2015)
- 47. Yang, W.H., Zhang, L.-H., Song, R.Y.: Optimality conditions of the nonlinear programming on Riemannian, manifolds. Pac. J. Optim. 10, 415–434 (2014)
- 48. Yang, X., Liu, W., Liu, W., Tao, D.: A survey on canonical correlation analysis. IEEE Trans. Knowl. Data Engrg. 33(6), 2349–2368 (2020)
- 49. Yu, Y., Zhang, L.-H., Zhang, S.: Simultaneous clustering of multiview biomedical data using manifold optimization. Bioinformatics **35**(20), 4029–4037 (2019)
- 50. Zhang, L.-H.: Riemannian Newton method for the multivariate eigenvalue problem. SIAM J. Matrix Anal. Appl. 31(5), 2972–2996 (2010)
- 51. Zhang, L.-H.: Riemannian trust-region method for the maximal correlation problem. Numer. Funct. Anal. Optim. **33**(3), 338–362 (2012)
- 52. Zhang, L.-H., Li, R.-C.: Maximization of the sum of the trace ratio on the Stiefel, manifold, I Theory. Sci. China Math. 57(12), 2495–2508 (2014)
- Zhang, L.-H., Li, R.-C.: Maximization of the sum of the trace ratio on the Stiefel manifold, II Computation. Sci. China Math. 58, 1549–1566 (2015)
- Zhang, L.-H., Ma, X., Shen, C.: A structure-exploiting nested Lanczos-type iteration for the multiview canonical correlation analysis. SIAM J. Sci. Comput. 43(4), A2685–A2713 (2021)
- Zhang, L.-H., Li, W.ang., Bai, Z., Li, R.-C.: A self-consistent-field iteration for orthogonal canonical correlation analysis. IEEE Trans. Pattern Anal. Mach Intell. 44(2), 890–904 (2022)
- 56. Zhang, L.-H., Yang, W.H., Shen, C., Ying, J.: An eigenvalue-based method for the unbalanced Procrustes problem. SIAM J Matrix Anal. Appl. **41**(3), 957–983 (2020)

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

13 Page 34 of 34 Adv Comput Math (2022) 48:13

#### **Affiliations**

# Xijun Ma<sup>1</sup> · Chungen Shen<sup>2</sup> · Li Wang<sup>3</sup> · Lei-Hong Zhang<sup>4,1</sup> □ · Ren-Cang Li<sup>3,5</sup>

Xijun Ma 2018310147@live.sufe.edu.cn

Chungen Shen shenchungen@usst.edu.cn

Li Wang li.wang@uta.edu

Ren-Cang Li rcli@uta.edu

- School of Mathematics, Shanghai University of Finance and Economics, 777 Guoding Road, Shanghai, 200433, China
- College of Science, University of Shanghai for Science and Technology, Shanghai, 200093, China
- Department of Mathematics, University of Texas at Arlington, Arlington, 76019-0408, TX, USA
- School of Mathematical Sciences, Soochow University, Suzhou, 215006, Jiangsu, China
- Department of Mathematics, Hong Kong Baptist University, Kowloon Tong, Kowloon, Hong Kong, Hong Kong, China

