

Extreme Point Pursuit—Part II: Further Error Bound Analysis and Applications

Junbin Liu, Ya Liu, Wing-Kin Ma, Mingjie Shao and Anthony Man-Cho So

Abstract— In the first part of this study, a convex-constrained penalized formulation was studied for a class of constant modulus (CM) problems. In particular, the error bound techniques were shown to play a vital role in providing exact penalization results. In this second part of the study, we continue our error bound analysis for the cases of partial permutation matrices, size-constrained assignment matrices and non-negative semi-orthogonal matrices. We develop new error bounds and penalized formulations for these three cases, and the new formulations possess good structures for building computationally efficient algorithms. Moreover, we provide numerical results to demonstrate our framework in a variety of applications such as the densest k -subgraph problem, graph matching, size-constrained clustering, non-negative orthogonal matrix factorization and sparse fair principal component analysis.

Index Terms— constant modulus optimization, non-convex optimization, error bound, densest subgraph problem, PCA, graph matching, clustering, ONMF

I. INTRODUCTION

In Part I of this paper [1], we considered a convex-constrained minimization framework for a class of constant modulus (CM) problems. Named extreme point pursuit (EXPP), this framework gives simple well-structured reformulations of the CM problems. This allows us to apply basic methods, such as the projected gradient (or subgradient) method, to build algorithms for CM problems; some underlying assumptions with the objective function are required, but they are considered reasonable in a wide variety of applications. As a requirement, EXPP chooses the constraint set as the convex hull of the CM set. When the projected gradient method is used, the computational efficiency will depend on whether the projection onto the convex hull of the CM set is easy to compute. We examined a number of CM examples that have such a benign property. But we also encountered some CM sets, namely, the partial permutation matrix (PPM) set and the size-constrained assignment matrix (SAM) set, whose convex hull projections may be inefficient to compute. Another difficult set is the non-negative semi-orthogonal matrix (NSOM) set, whose convex hull is not even known.

As Part II of this paper, we continue our analysis with the PPM set, the SAM set, and the NSOM set. Our focus is on the error bound principle for exact penalization, and we want to see if we can relax the EXPP constraint set for

the above three cases. We will consider some efficient-to-project constraint sets for the three cases, and we will show that, by suitably modifying the penalty function, we can once again achieve exact penalization. These case-specific results are more powerful versions of their counterparts in Part I of this paper. In the graph matching application, for example, we will numerically illustrate that the new method has much better computational efficiency than the state of the arts. Also, our results for NSOMs draw connections to some recently emerged results [2]–[4], as we shall see.

In addition to error bound analysis, we will numerically demonstrate EXPP in a variety of applications—such as the densest k -subgraph problem, graph matching, size-constrained clustering, ONMF, and PCA with sparsity and/or fairness. We will see that EXPP is a working method across different applications, offering reasonable performance and computational efficiency. The organization is as follows. Section II recapitulates some key concepts in Part I of this study. Section III provides further error bound analysis and devises new EXPP formulations. Section IV provides numerical results for different applications. Section V concludes this paper.

II. A RECAP OF PART I

We give a summary of Part I of this study, with a focus on the error bound principle. Let $\mathcal{D} \subseteq \mathbb{R}^n$ be a set which will be used to denote the domain of a problem. Let $\mathcal{X} \subseteq \mathcal{D}$ be a set which will often be chosen as a convex compact set. We consider a class of CM problems in the form of

$$\min_{\mathbf{x} \in \mathcal{V}} f(\mathbf{x}), \quad (1)$$

where $f : \mathcal{D} \rightarrow \mathbb{R}$ is K -Lipschitz continuous on \mathcal{X} ; $\mathcal{V} \subseteq \mathcal{X}$ is a non-empty closed CM set with modulus \sqrt{C} . Consider a general penalized formulation of (1):

$$\min_{\mathbf{x} \in \mathcal{X}} F_\lambda(\mathbf{x}) := f(\mathbf{x}) + \lambda h(\mathbf{x}), \quad (2)$$

where $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is a penalty function; $\lambda > 0$ is a given scalar. We want to find an \mathcal{X} and an h such that the penalized formulation (2) is “easy” to build an algorithm and has exact penalization guarantees. By exact penalization, we mean that the solution set of (2) is equal to that of (1). According to the error bound principle, formulation (2) is an exact penalization formulation of (1) if h is effectively an error bound function of \mathcal{V} relative to \mathcal{X} , i.e.,

$$\begin{aligned} \text{dist}(\mathbf{x}, \mathcal{V}) &\leq \nu h(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}, \\ 0 &= h(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{V}, \end{aligned}$$

The work of J. Liu, Y. Liu, W.-K. Ma and M. Shao was supported by the General Research Fund (GRF) of Hong Kong Research Grant Council (RGC) under Project ID CUHK 14208819. The work of A. M.-C. So was supported by the GRF of Hong Kong RGC under Project ID CUHK 14205421.

Junbin Liu and Ya Liu have equal contributions.

for some constant $\nu > 0$. More specifically, the exact penalization result holds when λ is sufficiently large such that $\lambda > K\nu$.

We studied a convex-constrained minimization formulation of problem (1), called EXPP, which is an instance of the formulation (2) with

$$\mathcal{X} = \text{conv}(\mathcal{V}), \quad h(\mathbf{x}) = C - \|\mathbf{x}\|_2^2. \quad (3)$$

We showed that the choice in (3) leads to an error bound, and consequently exact penalization, for many practical CM sets of interest. A merit with EXPP is that its penalty function h in (3) is simple. If the constraint set \mathcal{X} in (3) is friendly to handle, then we may build algorithms for EXPP without much difficulty. There are various possibilities for one to build algorithms for EXPP, and here we focus our attention on the projected gradient method or related methods. Assuming that f is differentiable, the projected gradient method for problem (2) is given by

$$\mathbf{x}^{l+1} = \Pi_{\mathcal{X}}(\mathbf{x}^l - \eta_l \nabla F_{\lambda}(\mathbf{x}^l)), \quad l = 0, 1, \dots, \quad (4)$$

where $\eta_l > 0$ is the step size. The computational efficiency of the projected gradient method depends on whether the projection $\Pi_{\mathcal{X}}$ is efficient to compute. We examined a collection of CM sets that have efficient-to-compute projections. But we also encounter CM sets whose projections can be expensive, or has no known way, to compute.

1. *partial permutation matrix (PPM) set:*

$$\mathcal{U}^{n,r} = \{\mathbf{X} \in \{0, 1\}^{n \times r} \mid \mathbf{X}^{\top} \mathbf{1} = \mathbf{1}, \mathbf{X} \mathbf{1} \leq \mathbf{1}\},$$

where $n \geq r$. The convex hull of $\mathcal{U}^{n,r}$ is

$$\text{conv}(\mathcal{U}^{n,r}) = \{\mathbf{X} \in [0, 1]^{n \times r} \mid \mathbf{X}^{\top} \mathbf{1} = \mathbf{1}, \mathbf{X} \mathbf{1} \leq \mathbf{1}\}. \quad (5)$$

There is no known easy way to compute the projection onto $\text{conv}(\mathcal{U}^{n,r})$. Solving the projection using a numerical solver can in practice be expensive for large n and r .

2. *size-constrained assignment matrix (SAM) set:*

$$\mathcal{U}_{\kappa}^{n,r} = \{\mathbf{X} \in \{0, 1\}^{n \times r} \mid \mathbf{X}^{\top} \mathbf{1} = \kappa, \mathbf{X} \mathbf{1} \leq \mathbf{1}\}, \quad (6)$$

where $n \geq r$, $\kappa \in \{1, \dots, n\}^r$, $\sum_{j=1}^r \kappa_j \leq n$. The convex hull of $\mathcal{U}_{\kappa}^{n,r}$ is

$$\text{conv}(\mathcal{U}_{\kappa}^{n,r}) = \{\mathbf{X} \in [0, 1]^{n \times r} \mid \mathbf{X}^{\top} \mathbf{1} = \kappa, \mathbf{X} \mathbf{1} \leq \mathbf{1}\}. \quad (7)$$

The projection onto $\text{conv}(\mathcal{U}_{\kappa}^{n,r})$ has the same issue as that in the PPM case.

3. *non-negative semi-orthogonal matrix (NSOM) set:*

$$\mathcal{S}_{+}^{n,r} = \mathcal{S}^{n,r} \cap \mathbb{R}_{+}^{n \times r},$$

where $n \geq r$, and $\mathcal{S}^{n,r} = \{\mathbf{X} \in \mathbb{R}^{n \times r} \mid \mathbf{X}^{\top} \mathbf{X} = \mathbf{I}\}$ is the semi-orthogonal matrix set. There are no known expressions with the convex hull of $\mathcal{S}_{+}^{n,r}$ and the associated projection.

In this second part of the study, we continue to study these CM sets.

III. FURTHER ERROR BOUND ANALYSIS

In this section we perform error bound analysis for the above three sets. The first and second subsections consider the PPM case. We will provide an application example to motivate the study, and then we will derive a new error bound and EXPP to overcome the computational issue described in the last section. Following the same genre, the SAM case is tackled in the third and fourth subsections, and the NSOM case studied in the fifth, sixth, and seventh subsections.

A. Example: Graph Matching

Consider the following problem for a given $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$:

$$\min_{\mathbf{X} \in \mathcal{U}^{n,n}} \|\mathbf{A} - \mathbf{X}^{\top} \mathbf{B} \mathbf{X}\|_{\mathbb{F}}^2.$$

This problem is called the graph matching (GM) problem in the context of computer vision. The goal is to match the nodes of two equal-size graphs by using the graph edge information. The GM problem does so by finding a set of one-to-one node associations, represented by \mathbf{X} , such that the associated Euclidean error between the two graph's adjacency matrices, represented by \mathbf{A} and \mathbf{B} , is minimized. Since a feasible \mathbf{X} is orthogonal, the GM problem can be rewritten as

$$\min_{\mathbf{X} \in \mathcal{U}^{n,n}} f(\mathbf{X}) := \|\mathbf{X} \mathbf{A} - \mathbf{B} \mathbf{X}\|_{\mathbb{F}}^2. \quad (8)$$

Following the review in the last section, the EXPP formulation of (8) is

$$\min_{\mathbf{X} \in \text{conv}(\mathcal{U}^{n,n})} F_{\lambda}(\mathbf{X}) = f(\mathbf{X}) - \lambda \|\mathbf{X}\|_{\mathbb{F}}^2. \quad (9)$$

It can be expensive to apply the projected gradient method (4) to the EXPP-GM problem (9). The projection operation $\Pi_{\text{conv}(\mathcal{U}^{n,n})}$ in the projected gradient method requires us to solve an optimization problem, namely, minimization of $\|\mathbf{Z} - \mathbf{X}\|_{\mathbb{F}}^2$ over $\mathbf{X} \in \text{conv}(\mathcal{U}^{n,n})$, where \mathbf{Z} is given. As mentioned, there is no easy-to-compute solution for this problem. The problem is nevertheless convex, and there is a specialized numerical solver for this problem based on gradient ascent of the dual problem [5]. Still, one would anticipate that solving an optimization problem at each step of the projected gradient method (4) would be computationally burdensome particularly when the problem size n is large.

At this point it is worthwhile to review an important prior study in GM [6]. The authors in that study developed a similar formulation as (8); they essentially put forth the same fundamental idea as EXPP under the concave minimization principle, though not under the error bound principle. They built an algorithm that exploits the problem structure for efficient computations. Specifically they considered the Frank-Wolfe method. The Frank-Wolfe method for (9) is given by

$$\mathbf{X}^{l+1} = \mathbf{X}^l + \alpha_l (\text{LO}_{\text{conv}(\mathcal{U}^{n,n})}(\nabla F_{\lambda}(\mathbf{X}^l)) - \mathbf{X}^l),$$

where $\alpha_l \in [0, 1]$ is a step size; $\text{LO}_{\mathcal{X}}(\mathbf{g}) \in \arg \min_{\mathbf{x} \in \mathcal{X}} \mathbf{g}^{\top} \mathbf{x}$ is the linear optimization (LO) oracle for \mathcal{X} . To efficiently compute each Frank-Wolfe iteration, we need an efficient way to compute $\text{LO}_{\text{conv}(\mathcal{U}^{n,n})}$. The authors of [6] did so by using the Hungarian algorithm [7], which is a specialized algorithm for solving linear optimization over the set of doubly stochastic

matrices and can efficiently compute $\text{LO}_{\text{conv}(\mathcal{U}^{n,n})}$ with a complexity of $\mathcal{O}(n^3)$.

We have a different proposition. To put into context, recall that $\Delta^n = \{\mathbf{x} \in \mathbb{R}_+^n \mid \mathbf{x}^\top \mathbf{1} = 1\}$ denotes the unit simplex. Define

$$\begin{aligned}\tilde{\Delta}^{n \times r} &= \{\mathbf{X} \in [0, 1]^{n \times r} \mid \mathbf{X}^\top \mathbf{1} = \mathbf{1}\} \\ &= \{\mathbf{X} \in \mathbb{R}^{n \times r} \mid \mathbf{x}_j \in \Delta^n, \forall j\}.\end{aligned}$$

This set is a relaxation of $\text{conv}(\mathcal{U}^{n,r})$ in (5) by taking out the row constraint $\mathbf{X}\mathbf{1} \leq \mathbf{1}$. Our idea is to derive an error bound function h of $\mathcal{U}^{n,r}$ relative to $\tilde{\Delta}^{n \times r}$. If we can do so, then we will have an exact penalization formulation (2) for the GM problem, and more generally, CM problems for the PPM case. Consequently there is a possibility for us to apply the projected gradient method (4) in a computationally efficient fashion. To be specific, we can compute the projection $\Pi_{\tilde{\Delta}^{n \times r}}(\mathbf{Z})$ of a given matrix $\mathbf{Z} \in \mathbb{R}^{n \times r}$ with a complexity of $\mathcal{O}(nr \log(n))$: the projection $\Pi_{\tilde{\Delta}^{n \times r}}(\mathbf{Z})$ corresponds to the unit simplex projections $\Pi_{\Delta^n}(\mathbf{z}_j)$'s for $j = 1, \dots, r$, and the unit simplex projection can be computed with a complexity of $\mathcal{O}(n \log(n))$ [8]. In the GM application, the complexity of the projection $\Pi_{\tilde{\Delta}^{n \times n}}$ is $\mathcal{O}(n^2 \log(n))$ —which seems more attractive than its Frank-Wolfe counterpart, $\mathcal{O}(n^3)$.

To illustrate the efficiencies of the above described projection and LO operations, we ran a numerical experiment. We tested the runtimes of (i) the projection onto $\text{conv}(\mathcal{U}^{n,n})$, implemented either by CVX or by the specialized dual gradient method [5]; (ii) the LO oracle for $\text{conv}(\mathcal{U}^{n,n})$, implemented by the Hungarian algorithm; and (iii) the projection onto $\tilde{\Delta}^{n \times n}$, done by column-wise unit-simplex projection. The results are shown in Fig. 1. It is seen that the projection onto $\tilde{\Delta}^{n \times n}$ is much faster than the other operations.

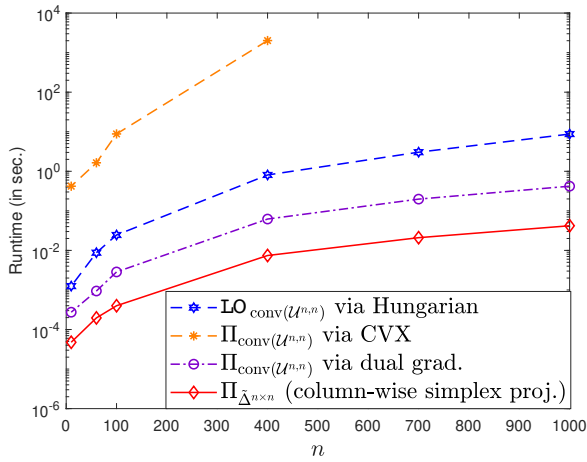


Fig. 1: Runtime comparison with the projections and LO oracle for the full PPM case ($n = r$). We used 200 randomly generated trials to evaluate the average runtimes, except for $\Pi_{\text{conv}(\mathcal{U}^{n,n})}$ via CVX which we tested only 20 trials due to the long runtime.

B. A New Error Bound and a New EXPP for PPMs

Motivated by the GM problem, we perform error bound analysis for the PPM set $\mathcal{U}^{n,r}$ relative to the column-wise unit simplex $\tilde{\Delta}^{n \times r}$. Our result is as follows.

Theorem 1 (error bound for partial permutation matrices)

For any $\mathbf{X} \in \tilde{\Delta}^{n \times r}$ we have the error bound

$$\text{dist}(\mathbf{X}, \mathcal{U}^{n,r}) \leq \nu \|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_{\ell_1} \quad (10a)$$

$$= \nu(r + \|\mathbf{X}\mathbf{1}\|_2^2 - 2\|\mathbf{X}\|_{\text{F}}^2), \quad (10b)$$

where $\nu = 11\sqrt{r}$.

The proof of Theorem 1 uses the same proof approach as that for the error bound for $\mathcal{U}_{\kappa}^{n,r}$ in Part I, Section IV.F, of this paper. But the latter is considered easier to show because of the presence of the row constraint $\mathbf{X}\mathbf{1} \leq \mathbf{1}$. To derive (10) we need to go further, analyzing the singular values of \mathbf{X} . In view of its complexity, the proof of Theorem 1 is relegated to Appendix B.

Theorem 1 gives rise to a new EXPP formulation for $\mathcal{U}^{n,r}$. By applying the error bound (10b) to the penalized formulation (2), we have an exact penalization formulation for the PPM case:

$$\min_{\mathbf{X} \in \tilde{\Delta}^{n \times r}} F_{\lambda}(\mathbf{X}) = f(\mathbf{X}) + \lambda(\|\mathbf{X}\mathbf{1}\|_2^2 - 2\|\mathbf{X}\|_{\text{F}}^2). \quad (11)$$

This new EXPP formulation resembles our previous EXPP, with addition of a friendly penalty term $\|\mathbf{X}\mathbf{1}\|_2^2$. Assuming a differentiable f , the new EXPP problem (11) can be efficiently handled by the projected gradient method (4). In particular, without counting the complexity of computing the gradient ∇F_{λ} , the per-iteration complexity of the projected gradient method is $\mathcal{O}(nr \log(n))$; see the discussion in the last subsection.

We should provide insight into the error bound in Theorem 1. The PPM set can be characterized as

$$\mathcal{U}^{n,r} = \tilde{\Delta}^{n \times r} \cap \mathcal{S}^{n,r}.$$

The error bound (10a) reveals that, by using $\tilde{\Delta}^{n \times r}$ as the constraint set and $\|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_{\ell_1}$ as the penalty, we can achieve exact penalization results. Particularly, $\|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_{\ell_1}$ appears as a penalty for promoting \mathbf{X} to be semi-orthogonal. Moreover, the error bound (10b) is an equivalent form of $\|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_{\ell_1}$. Consider the following result.

Lemma 1 Let $\mathbf{d} \in \mathbb{R}_{++}^r$ be given. Let $\mathbf{D} = \text{Diag}(\mathbf{d})$. For any $\mathbf{X} \in \mathbb{R}_+^{n \times r}$ with $\|\mathbf{x}_j\|_2^2 \leq d_j$ for all j , it holds that

$$\|\mathbf{X}^\top \mathbf{X} - \mathbf{D}\|_{\ell_1} = \mathbf{1}^\top \mathbf{d} + \|\mathbf{X}\mathbf{1}\|_2^2 - 2\|\mathbf{X}\|_{\text{F}}^2.$$

Proof of Lemma 1: For any $\mathbf{X} \in \mathbb{R}_+^{n \times r}$ with $\|\mathbf{x}_j\|_2^2 \leq d_j$ for all j ,

$$\begin{aligned}\|\mathbf{X}^\top \mathbf{X} - \mathbf{D}\|_{\ell_1} &= \sum_{j=1}^r (d_j - \|\mathbf{x}_j\|_2^2) + \sum_{j=1}^r \sum_{\substack{k=1 \\ j \neq k}}^r \mathbf{x}_j^\top \mathbf{x}_k \\ &= \sum_{j=1}^r (d_j - 2\|\mathbf{x}_j\|_2^2) + \left(\sum_{j=1}^r \mathbf{x}_j \right)^\top \left(\sum_{j=1}^r \mathbf{x}_j \right) \\ &= \mathbf{1}^\top \mathbf{d} - 2\|\mathbf{X}\|_{\text{F}}^2 + \|\mathbf{X}\mathbf{1}\|_2^2.\end{aligned}$$

The proof is complete. \blacksquare

Applying Lemma 1 to (10a) gives the error bound (10b); note that any $\mathbf{x} \in \Delta^n$ has $\|\mathbf{x}\|_2 \leq \|\mathbf{x}\|_1 = \mathbf{1}^\top \mathbf{x} = 1$. The error bound (10b) also reveals interesting insight. For any $\mathbf{X} \in \mathbb{R}_+^{n \times r}$ with $\|\mathbf{x}_j\|_2 \leq 1$ for all j , we can rewrite (10b) as

$$r + \|\mathbf{X}\mathbf{1}\|_2^2 - 2\|\mathbf{X}\|_F^2 = \sum_{j=1}^r c_2(\mathbf{x}_j) + \sum_{i=1}^n \rho_2(\bar{\mathbf{x}}_i), \quad (12)$$

where

$$c_2(\mathbf{x}) = 1 - \|\mathbf{x}\|_2^2, \quad \rho_2(\mathbf{x}) = \|\mathbf{x}\|_1^2 - \|\mathbf{x}\|_2^2$$

appear as penalties for columns and rows, respectively. According to Part I of this paper, c_2 is effectively an error bound function for the unit vector set $\mathcal{U}^n = \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ relative to the unit simplex Δ^n ; this means that c_2 promotes every column \mathbf{x}_j to lie in \mathcal{U}^n . As for ρ_2 , we note the basic norm result that $\|\mathbf{x}\|_1 \geq \|\mathbf{x}\|_2$ for any \mathbf{x} , and that $\|\mathbf{x}\|_1 = \|\mathbf{x}\|_2$ if and only if \mathbf{x} is a scaled unit vector, i.e., $\mathbf{x} = \alpha \mathbf{e}_i$ for some α and i . This means that ρ_2 promotes every row $\bar{\mathbf{x}}_i$ to be a scaled unit vector. Putting the column and row penalties together, we have the interpretation that (12) promotes \mathbf{X} to be a PPM.

C. Example: Size-Constrained Clustering

We turn our interest to the SAM set $\mathcal{U}_\kappa^{n,r}$ in (6). Let $\kappa \in \{1, \dots, n\}^r$, $\sum_{j=1}^r \kappa_j = n$, be given. Given a matrix $\mathbf{Y} \in \mathbb{R}^{m \times n}$, we consider the following problem

$$\min_{\mathbf{A} \in \mathbb{R}^{m \times r}, \mathbf{X} \in \mathcal{U}_\kappa^{n,r}} \|\mathbf{Y} - \mathbf{A}\mathbf{X}^\top\|_F^2. \quad (13)$$

This problem appears in size-constrained clustering and was used, for instance, in paper-to-session assignment [9]. In size-constrained clustering, we want to cluster a given set of data points $\mathbf{y}_1, \dots, \mathbf{y}_n$ into r clusters, and the constraint is that each cluster has size κ_j . In the size-constrained clustering problem (13), the j th column \mathbf{a}_j of \mathbf{A} describes the cluster center of cluster j . From (6) it is seen that \mathbf{X} takes the form

$$\mathbf{X}^\top = [\mathbf{e}_{l_1}, \dots, \mathbf{e}_{l_n}],$$

for some $l_i \in \{1, \dots, r\}$. In particular, \mathbf{e}_j is constrained to appear in the rows of \mathbf{X} for κ_j times.

A natural way to handle the size-constrained clustering problem (13) is to apply alternating minimization; see, e.g., [9]. In this study, we use the following reformulation of the size-constrained clustering problem

$$\min_{\mathbf{X} \in \mathcal{U}_\kappa^{n,r}} f(\mathbf{X}) := -\text{tr}(\mathbf{D}^{-1} \mathbf{X}^\top \mathbf{R} \mathbf{X}), \quad (14)$$

where $\mathbf{D} = \text{Diag}(\kappa)$; $\mathbf{R} = \mathbf{Y}^\top \mathbf{Y}$. Problem (14) is obtained by substituting the solution to \mathbf{A} given an $\mathbf{X} \in \mathcal{U}_\kappa^{n,r}$, specifically, $\mathbf{A} = \mathbf{Y}(\mathbf{X}^\top)^\dagger = \mathbf{Y} \mathbf{X} \mathbf{D}^{-1}$, to problem (13). Here, \dagger denotes the pseudo inverse. We encounter the same difficulty as the PPM case in Section III-A: The EXPP formulation can be applied to (14), but the projected gradient method for EXPP can be expensive to implement due to the computational cost

of numerically solving $\Pi_{\text{conv}(\mathcal{U}_\kappa^{n,r})}$. We want to replace the EXPP constraint set $\text{conv}(\mathcal{U}_\kappa^{n,r})$ with

$$\begin{aligned} \tilde{\mathcal{U}}_\kappa^{n,r} &= \{\mathbf{X} \in [0, 1]^{n \times r} \mid \mathbf{X}^\top \mathbf{1} = \kappa\} \\ &= \{\mathbf{X} \in \mathbb{R}^{n \times r} \mid \mathbf{x}_j \in \text{conv}(\mathcal{U}_{\kappa_j}^n) \forall j\}. \end{aligned}$$

Recall that $\mathcal{U}_\kappa^n = \{\mathbf{x} \in \{0, 1\}^n \mid \mathbf{1}^\top \mathbf{x} = \kappa\}$ and $\text{conv}(\mathcal{U}_\kappa^n) = \{\mathbf{x} \in [0, 1]^n \mid \mathbf{1}^\top \mathbf{x} = \kappa\}$, where $\kappa \in \{1, \dots, n\}$; and that $\Pi_{\text{conv}(\mathcal{U}_\kappa^n)}$ can be efficiently computed by a bisection algorithm with a complexity of $\mathcal{O}(n \log(n))$, given a solution precision [10, Algorithm 1].

D. A New Error Bound and a New EXPP for SAMs

We have the following result.

Theorem 2 (error bound for size-constrained assignment matrices) For any $\mathbf{X} \in \tilde{\mathcal{U}}_\kappa^{n,r}$, we have the error bound

$$\text{dist}(\mathbf{X}, \mathcal{U}_\kappa^{n,r}) \leq \nu \|\mathbf{X}^\top \mathbf{X} - \text{Diag}(\kappa)\|_{\ell_1} \quad (15a)$$

$$= \nu(\mathbf{1}^\top \kappa + \|\mathbf{X}\mathbf{1}\|_2^2 - 2\|\mathbf{X}\|_F^2), \quad (15b)$$

where $\nu = 3 \sum_{i=1}^r (1 + 2\kappa_i)(1 + \sqrt{\kappa_i})\sqrt{\mathbf{1}^\top \kappa}$. Note that $\nu \leq 18\kappa_{\max}^2 r^{3/2}$, where $\kappa_{\max} = \max\{\kappa_1, \dots, \kappa_r\}$.

The idea behind the proof of Theorem 2 is the same as that of the PPM case in Theorem 1. The actual proof of Theorem 2 is however more tedious. The reader can find the proof in Appendix C. The error bound (15) shares the same insight as its PPM counterpart (see Section III-B), and we shall not repeat.

Applying the error bound (15b) to the penalized formulation (2) gives rise to the following new EXPP formulation for the SAM case:

$$\min_{\mathbf{X} \in \tilde{\mathcal{U}}_\kappa^{n,r}} F_\lambda(\mathbf{X}) = f(\mathbf{X}) + \lambda(\|\mathbf{X}\mathbf{1}\|_2^2 - 2\|\mathbf{X}\|_F^2). \quad (16)$$

This new EXPP is identical to the new EXPP for the partial permutation matrix case in (11). Assuming a differentiable f , we can use the projected gradient method (4) to efficiently handle problem (16). Specifically, the projection $\Pi_{\tilde{\mathcal{U}}_\kappa^{n,r}}$, which contributes to the bulk of complexity with the projected gradient method, can be done with a complexity of $\mathcal{O}(nr \log(n))$; the idea is the same as that of the PPM case (see the last paragraph of Section III-A).

E. Example: Orthogonal Non-Negative Matrix Factorization

As an application example for the NSOM set $\mathcal{S}_+^{n,r}$, consider the orthogonal non-negative matrix factorization (ONMF) problem

$$\min_{\mathbf{A} \in \mathbb{R}_+^{m \times r}, \mathbf{X} \in \mathcal{S}_+^{n,r}} \|\mathbf{Y} - \mathbf{A}\mathbf{X}^\top\|_F^2, \quad (17)$$

where the given matrix \mathbf{Y} is non-negative. The ONMF problem is, in essence, a non-negatively scaled clustering problem. It is known by researchers that any point \mathbf{X} in $\mathcal{S}_+^{n,r}$ can be characterized as

$$\mathbf{X}^\top = [\alpha_1 \mathbf{e}_{l_1}, \dots, \alpha_n \mathbf{e}_{l_n}], \quad (18)$$

where $\alpha_i \geq 0$ for all i , $l_i \in \{1, \dots, r\}$, $\sum_{i:l_i=j} \alpha_i^2 = 1$ for all j ; see, e.g., [2]–[4], [11]. From this characterization, we

see that, fixing an assignment l_1, \dots, l_n , each \mathbf{a}_j is a cluster center, obtained by minimizing $\sum_{i:l_i=j} \|\mathbf{y}_i - \alpha_i \mathbf{a}_j\|_2^2$ over the non-negative \mathbf{a}_j and the non-negative scale-compensating scalars α_i 's. We consider the following reformulation of the ONMF problem

$$\min_{\mathbf{X} \in \mathcal{S}_+^{n,r}} f(\mathbf{X}) := -\text{tr}(\mathbf{X}^\top \mathbf{R} \mathbf{X}), \quad (19)$$

where $\mathbf{R} = \mathbf{Y}^\top \mathbf{Y}$. Problem (19) is obtained by putting the solution to \mathbf{A} given an $\mathbf{X} \in \mathcal{S}_+^{n,r}$ and a $\mathbf{Y} \in \mathbb{R}_+^{m \times n}$, i.e., $\mathbf{A} = \mathbf{Y}(\mathbf{X}^\top)^\dagger = \mathbf{Y} \mathbf{X}$, to (17). As discussed in Part I of this paper, there is no known expression for the convex hull of $\mathcal{S}_+^{n,r}$.

We turn to the error bound principle, seeing if we can obtain an error bound of $\mathcal{S}_+^{n,r}$ relative to a friendly constraint set \mathcal{X} . We should mention a recently-emerged related work [3]. Consider the following result.

Theorem 3 (a special case of Lemma 3.1 in [3]) Define $\tilde{\mathcal{S}}_+^{n \times r} = \{\mathbf{X} \in \mathbb{R}_+^{n \times r} \mid \|\mathbf{x}_j\|_2 = 1 \ \forall j\}$ as the column-wise non-negative unit sphere. For any $\mathbf{X} \in \tilde{\mathcal{S}}_+^{n \times r}$, we have the error bound

$$\text{dist}(\mathbf{X}, \mathcal{S}_+^{n,r}) \leq \sqrt{2r} \sqrt{\|\mathbf{X} \mathbf{1}\|_2^2 - r}. \quad (20)$$

The authors of [3] actually derived a more general result than the above, but the above result is considered the most representative and was the standard choice in the authors' numerical demonstrations. From the perspective of this study, the most interesting question lies in how the above error bound was shown. We will come back to this later. Using the error bound principle (applying (20) to the penalized formulation (2)), the authors of [3] gave the following exact penalization formulation

$$\min_{\mathbf{X} \in \tilde{\mathcal{S}}_+^{n,r}} f(\mathbf{X}) + \lambda \sqrt{\|\mathbf{X} \mathbf{1}\|_2^2 - r} \quad (21)$$

for the ONMF problem (19) or for CM problems for the NSOM case. The presence of a square root in the penalty of (21) makes the objective function non-smooth, which adds some challenge from the viewpoint of building algorithms. The authors of [3] handled problem (21) by using manifold optimization to deal with the manifold $\tilde{\mathcal{S}}_+^{n \times r}$ and by applying smooth approximation to deal with the penalty term. Taking inspiration from the above error bound result, another group of researchers derived a general error bound result as follows.

Theorem 4 (Theorem 5 in [4]) For any $\mathbf{X} \in \mathbb{R}^{n \times r}$, we have the error bound

$$\text{dist}(\mathbf{X}, \mathcal{S}_+^{n,r}) \leq 5r^{\frac{3}{4}} (\|\mathbf{X}_- \|_{\text{F}}^{\frac{1}{2}} + \|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_{\text{F}}^{\frac{1}{2}}), \quad (22)$$

where $\mathbf{X}_- = \max\{-\mathbf{X}, \mathbf{0}\}$.

F. A Modified Error Bound and a New EXPP for NSOMs

We derive a modified error bound result for $\mathcal{S}_+^{n,r}$. Let

$$\tilde{\mathcal{B}}_+^{n \times r} = \{\mathbf{X} \in \mathbb{R}_+^{n \times r} \mid \|\mathbf{x}_j\|_2 \leq 1 \ \forall j\}$$

define the non-negative column-wise unit ℓ_2 -norm ball. Let

$$\psi_p(\mathbf{X}) = \left[\sum_{j=1}^r c_1(\mathbf{x})^p + \sum_{i=1}^n \rho_1(\bar{\mathbf{x}}_i)^p \right]^{\frac{1}{p}} \quad (23)$$

be a penalty function, where $p \in \{1, 2\}$,

$$c_1(\mathbf{x}) = 1 - \|\mathbf{x}\|_2, \quad \rho_1(\bar{\mathbf{x}}) = \|\mathbf{x}\|_1 - \|\mathbf{x}\|_\infty.$$

The penalty ψ_p shares a similar rationale as the one in (12): for $\mathbf{X} \in \mathbb{R}^{n \times r}$ with $\|\mathbf{x}_j\|_2 \leq 1$ for all j , $c_1(\mathbf{x}_j)$ promotes \mathbf{x}_j to have unit ℓ_2 norm, while $\rho_1(\bar{\mathbf{x}}_i)$ promotes $\bar{\mathbf{x}}_i$ to be a scaled unit vector. We should note that $\|\mathbf{x}\|_1 \geq \|\mathbf{x}\|_\infty$, and that $\|\mathbf{x}\|_1 = \|\mathbf{x}\|_\infty$ if and only if \mathbf{x} is a scaled unit vector. We should also recall from (18) that an NSOM has rows being scaled unit vectors. Our result is as follows.

Theorem 5 (error bound for non-negative semi-orthogonal matrices) For any $\mathbf{X} \in \tilde{\mathcal{B}}_+^{n \times r}$, we have the error bounds

$$\text{dist}(\mathbf{X}, \mathcal{S}_+^{n,r}) \leq \nu \psi_2(\mathbf{X}) \leq \nu \psi_1(\mathbf{X}), \quad (24)$$

where $\nu = \max\{\sqrt{6}, 2\sqrt{r}\}$.

The proof of Theorem 5 is provided in Appendix D, and we will give insight into these error bounds in the next subsection. Theorem 5 can be used to build a new EXPP formulation. From (24) we have a further error bound

$$\nu \psi_1(\mathbf{X}) \leq \nu \left[r - \|\mathbf{X}\|_{\text{F}}^2 + \sum_{i=1}^n (\mathbf{1}^\top \bar{\mathbf{x}}_i - s_1(\bar{\mathbf{x}}_i)) \right],$$

where we use the fact that $\|\mathbf{x}_j\|_2 \geq \|\mathbf{x}_j\|_2^2$ whenever $\|\mathbf{x}_j\|_2 \leq 1$; we recall $s_1(\mathbf{x}) = x_{[1]}$; $\mathbf{1}^\top \bar{\mathbf{x}}_i - s_1(\bar{\mathbf{x}}_i)$ is an alternate form of $\rho_1(\bar{\mathbf{x}}_i)$ for non-negative $\bar{\mathbf{x}}_i$. The above error bound leads us to the following new EXPP formulation for the NSOM case:

$$\min_{\mathbf{X} \in \tilde{\mathcal{B}}_+^{n \times r}} F_\lambda(\mathbf{X}) = f(\mathbf{X}) + \lambda \left[\sum_{i=1}^n (\mathbf{1}^\top \bar{\mathbf{x}}_i - s_1(\bar{\mathbf{x}}_i)) - \|\mathbf{X}\|_{\text{F}}^2 \right] \quad (25)$$

Unlike the previous EXPP formulations, the above formulation has non-smooth penalty terms $-s_1(\bar{\mathbf{x}}_i)$'s. But since they are concave, they are considered not difficult to handle when we use majorization minimization techniques to build algorithms. Specifically, by Jensen's inequality, we have

$$-s_1(\mathbf{x}) \leq -s_1(\mathbf{x}') + x'_{l'} - x_{l'}, \quad \text{for any } \mathbf{x}, \mathbf{x}',$$

where l' is such that $x'_{l'} = s_1(\mathbf{x}')$. The above inequality can be used to conveniently construct a majorant for the penalty function. Moreover, it should be mentioned that the projection onto $\tilde{\mathcal{B}}_+^{n \times r}$ has a closed form. The projection $\Pi_{\tilde{\mathcal{B}}_+^{n \times r}}(\mathbf{Z})$ for a given matrix $\mathbf{Z} \in \mathbb{R}^{n \times r}$ corresponds to the projections $\Pi_{\mathcal{B}_+^n}(\mathbf{z}_j)$'s, where $\mathcal{B}_+^n = \mathcal{B}^n \cap \mathbb{R}_+^n$. The projection $\Pi_{\mathcal{B}_+^n}(\mathbf{z})$ for a given vector $\mathbf{z} \in \mathbb{R}^n$ equals \mathbf{z}_+ if $\|\mathbf{z}_+\|_2 \leq 1$, and $\mathbf{z}_+ / \|\mathbf{z}_+\|_2$ if $\|\mathbf{z}_+\|_2 > 1$; here, $\mathbf{z}_+ = \max\{\mathbf{0}, \mathbf{z}\}$.

G. Insights and Further Remarks for NSOMs

We should describe the insight behind the proof of Theorem 5. Some of the key steps of our proof are actually the same as those of the previous results in Theorems 3 and 4. We however change one important ingredient. Recall from (18) that every NSOM has rows being scaled unit vectors. This motivates us to consider the scaled unit vector set

$$\mathcal{W}^n := \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} = \alpha \mathbf{e}_i, \alpha \in \mathbb{R}, i \in \{1, \dots, n\}\}.$$

In particular we consider the error bound for \mathcal{W}^n :

Lemma 2 *For any $\mathbf{x} \in \mathbb{R}^n$ we have the error bound*

$$\text{dist}(\mathbf{x}, \mathcal{W}^n) \leq \|\mathbf{x}\|_1 - \|\mathbf{x}\|_\infty = \rho_1(\mathbf{x}). \quad (26)$$

Proof of Lemma 2: Let $\mathbf{x} \in \mathbb{R}^n$ be given. Let $\mathbf{y} = x_l \mathbf{e}_l$, where l is such that $|x_l| = \max\{|x_1|, \dots, |x_n|\}$. It can be verified that $\mathbf{y} = \Pi_{\mathcal{W}^n}(\mathbf{x})$ and $\text{dist}(\mathbf{x}, \mathcal{W}^n) = \|\mathbf{x} - \mathbf{y}\|_2 \leq \|\mathbf{x} - \mathbf{y}\|_1 = \rho_1(\mathbf{x})$. ■

Our proof differs in that we use the error bound in Lemma 2, while the prior results used a different bound. Interested readers are referred to Appendix E, wherein we delineate the difference.

In fact, from the proof of Theorem 5, we notice the following generalization.

Corollary 1 *Let $\mathcal{S}_e^{n,r} = \{\mathbf{X} \in \mathbb{R}^{n \times r} \mid \|\mathbf{x}_j\|_2 = 1, \bar{\mathbf{x}}_i \in \mathcal{W}^r, \forall i, j\}$. For any $\mathbf{X} \in \mathbb{R}^{n \times r}$ with $\|\mathbf{x}_j\|_2 \leq 1$ for all j , we have error bounds $\text{dist}(\mathbf{X}, \mathcal{S}_e^{n,r}) \leq \nu \psi_2(\mathbf{X}) \leq \nu \psi_1(\mathbf{X})$, where ν is given by the one in Theorem 5.*

We omit the proof as it is a trivial variation of the proof of Theorem 5. Note that $\mathcal{S}_+^{n,r} \subset \mathcal{S}_e^{n,r} \subset \mathcal{S}^{n,r}$, and $\mathcal{S}_e^{n,r}$ permits negative components.

It is also worth noting that Theorem 5 has connections with the previous results in Theorems 3 and 4. To describe it, we first derive the following result.

Lemma 3 *For any $\mathbf{x} \in \mathbb{R}^n$, it holds that $\rho_1(\mathbf{x})^2 \leq \|\mathbf{x}\|_1^2 - \|\mathbf{x}\|_2^2 = \rho_2(\mathbf{x})$.*

Proof of Lemma 3: Without loss of generality, assume $|x_1| \geq |x_i|$ for all i . We have

$$\begin{aligned} \rho_1(\mathbf{x})^2 &= (\|\mathbf{x}\|_1 - |x_1|)^2 = \|\mathbf{x}\|_1^2 - 2|x_1|\|\mathbf{x}\|_1 + |x_1|^2 \\ &\leq \|\mathbf{x}\|_1^2 - 2(|x_1|^2 + |x_2|^2 + \dots + |x_n|^2) + |x_1|^2 \\ &\leq \|\mathbf{x}\|_1^2 - \|\mathbf{x}\|_2^2. \end{aligned}$$

The proof is complete. ■

From Lemma 3 we have the following further result.

Lemma 4 *For any $\mathbf{X} \in \tilde{\mathcal{B}}_+^{n \times r}$, it holds that*

$$\psi_2(\mathbf{X}) \leq \sqrt{r + \|\mathbf{X}\mathbf{1}\|_2^2 - 2\|\mathbf{X}\|_F^2} \quad (27a)$$

$$= \|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_{\ell_1}^{\frac{1}{2}}. \quad (27b)$$

$$\leq r\|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_F^{\frac{1}{2}} \quad (27c)$$

Proof of Lemma 4: Let $\mathbf{X} \in \tilde{\mathcal{B}}_+^{n \times r}$ be given. Since $\|\mathbf{x}_j\|_2 \leq 1$, we have $c_1(\mathbf{x}_j)^2 \leq 1 - \|\mathbf{x}_j\|_2 \leq 1 - \|\mathbf{x}_j\|_2^2 = c_2(\mathbf{x}_j)$. By Lemma 3 we have $\rho_1(\bar{\mathbf{x}}_i)^2 \leq \rho_2(\bar{\mathbf{x}}_i)$. Applying these results to (23), and observing (12), we have (27a). Eq. (27b) is a direct consequence of Lemma 1. Eq. (27c) is the consequence of the norm result $\|\mathbf{a}\|_1 \leq \sqrt{n}\|\mathbf{a}\|_2$ for $\mathbf{a} \in \mathbb{R}^n$. ■

By Theorem 5, all the functions in (27) are effective error bound functions of $\mathcal{S}_+^{n,r}$ relative to $\tilde{\mathcal{B}}_+^{n \times r}$. If we limit \mathbf{X} to lie in $\tilde{\mathcal{B}}_+^{n \times r}$, then the effective error bound function in (27a) becomes that in (20) in Theorem 3. The effective error bound function in (27c) is identical to that in (22) in Theorem 4 for the special case of $\mathbf{X} \in \tilde{\mathcal{B}}_+^{n \times r}$. Hence, simply speaking, Theorem 5 can produce the error bound results in the prior studies.

We end with a further comment. The penalty functions ρ_1 and ρ_2 for scaled unit vectors were considered in [2]. The authors of that work studied a variation of the ONMF problem (17) wherein they remove the unit ℓ_2 -norm column constraints with \mathbf{X} . They showed stationarity results with the use of ρ_1 and ρ_2 . Their analysis is not based on error bounds, and the problem nature with $\mathcal{S}_+^{n,r}$ is considered more challenging than that with the scaled unit vectors.

IV. NUMERICAL RESULTS IN DIFFERENT APPLICATIONS

In this section we numerically demonstrate EXPP in different applications.

A. The Algorithm for EXPP

We should specify the algorithm used to implement EXPP in our experiments. The algorithm used is an extrapolated variant of the projected gradient (PG) method, which was numerically found to have faster convergence; see, e.g., [12]. It is also the same algorithm used in our prior studies [13], [14]. The algorithm is shown in Algorithm 1. To explain it, let us write down the EXPP formulation

$$\min_{\mathbf{x} \in \mathcal{X}} F_\lambda(\mathbf{x}) = f(\mathbf{x}) + \lambda h(\mathbf{x}), \quad (28)$$

where $\mathcal{X} = \text{conv}(\mathcal{V})$; $h(\mathbf{x}) = -\|\mathbf{x}\|_2^2$; f is assumed to be differentiable and have Lipschitz continuous gradient on \mathcal{X} . Line 6–7 of Algorithm 1 is the extrapolated PG step. If we replace \mathbf{z} by $\tilde{\mathbf{x}}^l$ and $\nabla G_{\lambda_k}(\mathbf{z}|\tilde{\mathbf{x}}^l)$ by $\nabla F_{\lambda_k}(\tilde{\mathbf{x}}^l)$, we return to the baseline PG step. The point \mathbf{z} is an extrapolated point, and the extrapolation sequence $\{\alpha_l\}$ is chosen as the FISTA sequence [15]. The function $G_\lambda(\mathbf{x}|\mathbf{x}')$ is a majorant of $F_\lambda(\mathbf{x})$ at \mathbf{x}' , and we apply majorization before the PG step. Specifically we majorize the concave h by $u(\mathbf{x}|\mathbf{x}') = -\|\mathbf{x}'\|_2^2 - 2(\mathbf{x} - \mathbf{x}')^\top \mathbf{x}'$ (it is the result of Jensen's inequality), and then we construct the majorant $G_\lambda(\mathbf{x}|\mathbf{x}') = f(\mathbf{x}) + \lambda u(\mathbf{x}|\mathbf{x}')$ (for $\lambda > 0$). With this choice we choose the step size η as the reciprocal of the Lipschitz constant of $\nabla G_\lambda(\mathbf{x}|\mathbf{x}')$. Furthermore, we apply a homotopy strategy wherein we gradually increase penalty parameter λ so that we start with a possibly convex problem and end with the target EXPP problem (with exact penalization).

Algorithm 1 is also used to implement the new EXPP formulations (11), (16) and (25) for the PPM, SAM and

Algorithm 1 A PG-type algorithm for (28), with homotopy

- 1: **given:** a sequence $\{\lambda_k\}$, an extrapolation sequence $\{\alpha_l\}$, and a starting point \mathbf{x}^0
 - 2: $k \leftarrow 0$
 - 3: **repeat**
 - 4: $\tilde{\mathbf{x}}^0 \leftarrow \mathbf{x}^k, \mathbf{z} \leftarrow \mathbf{x}^k, l \leftarrow 0$
 - 5: **repeat**
 - 6: $\tilde{\mathbf{x}}^{l+1} \leftarrow \Pi_{\mathcal{X}}(\mathbf{z} - \eta \nabla G_{\lambda_k}(\mathbf{z} | \tilde{\mathbf{x}}^l))$ for some $\eta > 0$
 - 7: $\mathbf{z} \leftarrow \tilde{\mathbf{x}}^{l+1} + \alpha_l(\tilde{\mathbf{x}}^{l+1} - \tilde{\mathbf{x}}^l)$
 - 8: $l \leftarrow l + 1$
 - 9: **until** a stopping rule is met
 - 10: $\mathbf{x}^{k+1} \leftarrow \tilde{\mathbf{x}}^l$
 - 11: $k \leftarrow k + 1$
 - 12: **until** a stopping rule is met
 - 13: **output:** \mathbf{x}^k
-

NSOM cases, respectively; we will call them “EXPP-II” to avoid confusion with the previous EXPP. For the PPM and SAM cases we have $h(\mathbf{X}) = \|\mathbf{X}\mathbf{1}\|_2^2 - 2\|\mathbf{X}\|_F^2$. We choose $u(\mathbf{X}|\mathbf{X}') = \|\mathbf{X}\mathbf{1}\|_2^2 - 4\text{tr}(\mathbf{X}^\top \mathbf{X}') + 2\|\mathbf{X}'\|_F^2$; we keep the convex term $\|\mathbf{X}\mathbf{1}\|_2^2$ and majorize the concave term $-2\|\mathbf{X}\|_F^2$. For the NSOM case we have $h(\mathbf{X}) = -\|\mathbf{X}\|_F^2 + \sum_{i=1}^n (\mathbf{1}^\top \tilde{\mathbf{x}}_i - s_1(\tilde{\mathbf{x}}_i))$. We choose $u(\mathbf{X}|\mathbf{X}') = -2\text{tr}(\mathbf{X}^\top \mathbf{X}') + \|\mathbf{X}'\|_F^2 + \sum_{i=1}^n (\mathbf{1}^\top \tilde{\mathbf{x}}_i - x_{i,l'_i})$, where l'_i is such that $x'_{i,l'_i} = \max\{x'_{i,1}, \dots, x'_{i,r}\}$; again we keep the convex terms and majorize the concave terms. The rest of the operations are identical.

We should mention that Algorithm 1 is equipped with guarantees of convergence to a stationary point. To be specific, the majorization and extrapolated PG loop in Line 5–9 of Algorithm 1 was shown to be able to converge to a stationary point of problem (28) (with $\lambda = \lambda_k$) [13]. The premise of this is that h is differentiable and has Lipschitz continuous gradient. This premise is satisfied in most of the above described cases, with the NSOM case being the only exception. The critical-point convergence for the NSOM case (nonsmooth h) is a subject of future study.

We may deal with applications that have non-differentiable objective function f . For such cases we replace the extrapolated PG step in Line 6–7 of Algorithm 1 with the projected subgradient method.

Some details with the stopping rules of Algorithm 1 are as follows. Unless specified otherwise, the stopping rule for the inner loop is either $\|\tilde{\mathbf{x}}^{l+1} - \tilde{\mathbf{x}}^l\|_2 / \|\tilde{\mathbf{x}}^l\|_2 < \varepsilon_1$ or $l > \bar{L}$ for some given ε_1 and \bar{L} . Unless specified otherwise, the stopping rule for the outer loop is $\|\mathbf{x}^{k+1} - \mathbf{x}^k\|_2 / \|\mathbf{x}^k\|_2 < \varepsilon_2$, $\text{dist}(\mathbf{x}^k, \mathcal{V}) \leq \varepsilon_3$, or $\lambda_k > \bar{\lambda}$ for some given $\varepsilon_2, \varepsilon_3$ and $\bar{\lambda}$.

B. MIMO detection with MPSK Constellations

We consider MIMO detection. In this problem, we have a received signal model $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{v}$, where $\mathbf{y} \in \mathbb{C}^m$ is the received signal; $\mathbf{H} \in \mathbb{C}^{m \times n}$ is the channel; $\mathbf{x} \in \mathbb{C}^n$ is the transmitted symbol vector; $\mathbf{v} \in \mathbb{C}^m$ is noise. The goal is to detect \mathbf{x} from \mathbf{y} . Assuming that every x_i is drawn from an M -ary phase shift keying (PSK) constellation $\Theta_M = \{x \in$

$\mathbb{C} \mid x = e^{j\frac{2\pi l}{M} + j\frac{\pi}{M}}, l \in \{0, 1, \dots, M-1\}\}$, we consider the maximum-likelihood (ML) detector

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \Theta_M^n} \|\mathbf{y} - \mathbf{H}\mathbf{x}\|_2^2.$$

We use EXPP to handle the ML detector.

The benchmarked algorithms are (i) MMSE: the minimum mean square detector; (ii) SDR: the semidefinite relaxation detector [16]; (iii) LAMA: the approximate message passing method in [17] with incorporation of damping [18] (with damping factor 0.7). The parameter settings of EXPP are $\lambda_0 = 0.01$, $\lambda_{k+1} = 5\lambda_k$, $\varepsilon_1 = 10^{-4}$, $\bar{L} = 100$, $\bar{\lambda} = 10^4$.

The simulation settings are as follows. A number of 10,000 Monte-Carlo trials were run. The channel \mathbf{H} was generated based on a correlated MIMO channel model $\mathbf{H} = \mathbf{R}_r^{\frac{1}{2}} \tilde{\mathbf{H}} \mathbf{R}_t^{\frac{1}{2}}$, where the components of $\tilde{\mathbf{H}}$ are independent and identically distributed (i.i.d.) and follow a circular Gaussian distribution with mean zero and variance 1, and \mathbf{R}_t and \mathbf{R}_r follow the exponential model with parameter $\rho = 0.2$ [19]. The vector \mathbf{v} was generated as a component-wise i.i.d. circular Gaussian noise with mean 0 and variance σ^2 .

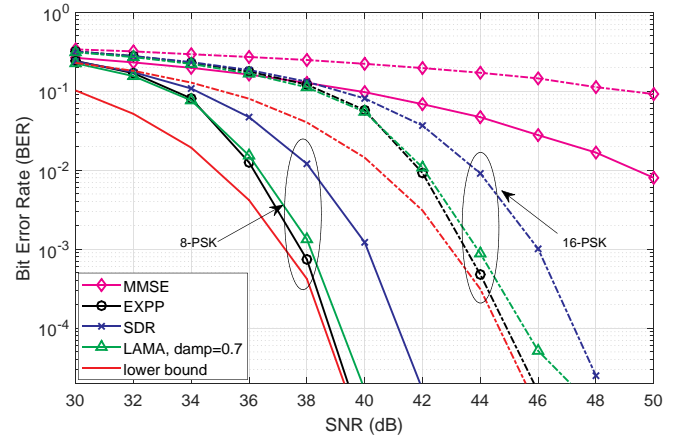


Fig. 2: MIMO detection BER performance. Solid lines: 8-PSK; dashed lines: 16-PSK.

Algorithms	SDR	EXPP	LAMA
time (in sec.)	1.017	0.026	0.021

TABLE I: Average runtime performance in MIMO detection. 8-PSK, SNR= 38dB.

Fig. 2 and Table I display the bit-error-rate (BER) and runtime performances of the various detectors, respectively. The problem size is $m = n = 80$. “Lower bound” is the performance baseline when there is no MIMO interference. EXPP gives the best BER performance and is closely followed by LAMA. The runtime performance of EXPP is slightly worse than that of LAMA, but is still competitive.

C. Densest k -Subgraph Problem

We consider the densest k -subgraph (DkS) problem. The problem is to identify the most connected subgraph from a

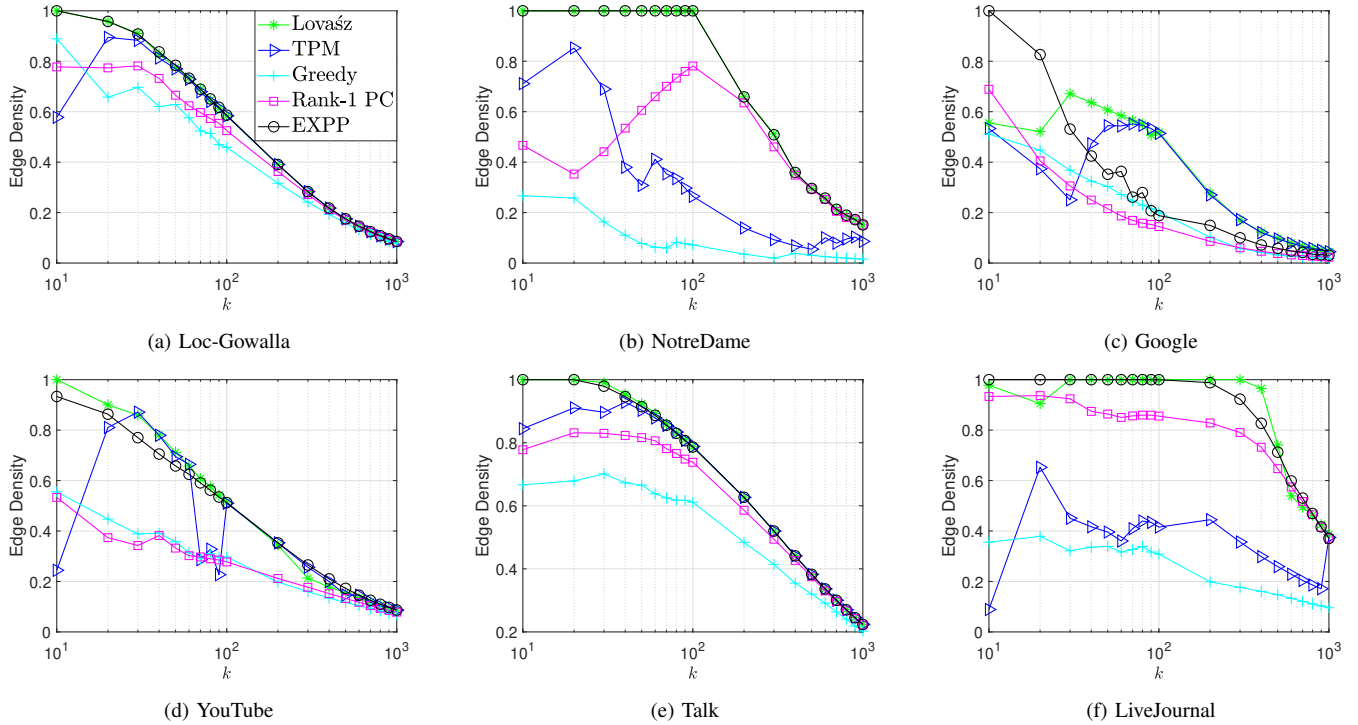


Fig. 3: Performance with the DkS problem.

graph. Let $\mathbf{W} \in \mathbb{R}^{n \times n}$ be the adjacency matrix of a given graph. Given a subgraph size k , the DkS problem is given by

$$\min_{\mathbf{x} \in \mathcal{U}_k^n} -\mathbf{x}^\top \mathbf{W} \mathbf{x},$$

where $\mathcal{U}_k^n = \{\mathbf{x} \in \{0, 1\}^n \mid \mathbf{1}^\top \mathbf{x} = k\}$. We use EXPP to handle this problem.

The benchmarked algorithms are (i) Greedy: the greedy algorithm in [20]; (ii) TPM: the truncated power method [21]; (iii) Rank-1 PC: the rank-one binary principal component approximation method [22]; (iv) Lovašz: Lovašz relaxation, implemented by the ADMM method and a Frank-Wolfe post-processing step [10]. The parameter settings of EXPP are $\lambda_0 = -\sigma_1(\mathbf{W})$, $\lambda_{k+1} = \lambda_k + 0.2\sigma_1(\mathbf{W})$, $\varepsilon_1 = 10^{-3}$, $\bar{L} = 50$, $\varepsilon_3 = 10^{-2}$, $\bar{\lambda} = 1.1\sigma_1(\mathbf{W})$. Note that the above choice of λ_0 is to make the corresponding EXPP problem convex; and that if $\lambda_k < 0$ then we choose $u(\mathbf{x}|\mathbf{x}') = h(\mathbf{x})$ (no majorization).

We performed our test on large real-world graphs [23]. Table II describes the graph sizes of the tested datasets, which range from nearly 200,000 nodes to nearly 4,000,000 nodes. We adopt the same data pre-processing as that in [10]. Fig. 3 shows the edge densities, $\mathbf{x}^\top \mathbf{W} \mathbf{x} / (k^2 - k)$, achieved by the various algorithms. Lovašz gives the best performance in all the datasets. EXPP is the second best. Except for the Google dataset, EXPP achieves similar performance as Lovašz. Table II shows the runtimes. We see that EXPP runs faster than Lovašz.

D. Graph Matching

We consider the GM problem described in Section III-A. The benchmarked algorithms are (i) PATH [6]; (ii) GNCCP

[24]; (iii) LAGSA [25]. PATH was concisely reviewed in Section III-A, and GNCCP and LAGSA are related methods. We employ EXPP-II in (11). The parameter settings of EXPP are $\lambda_0 = 10^{-5}$, $\lambda_{k+1} = 4\lambda_k$, $\varepsilon_1 = 10^{-4}$, $\varepsilon_2 = 10^{-5}$, $\bar{L} = 100$, $\bar{\lambda} = 10^3$.

We performed our test on graphs constructed from various image datasets. We follow a standard protocol to extract graphs from a pair of images; see e.g., [26]. We tested several image sets from 5 datasets: CMU-House [27], PASCAL [28], DTU [29], DAISY [30], and SUIRD [31]. CMU-House and PASCAL are small graphs and are commonly used in GM papers; DTU, DAISY, and SUIRD are large graphs. For each pair of images, we performed 10 independent graph constructions and used them for experiments.

Table III shows the average matching accuracies and runtimes of the different algorithms; the matching accuracy is defined as the ratio of the number of correctly matched nodes to the node size n . We did not test PATH and LAGSA on large-size graphs because they run too slowly. EXPP-II gives the best matching accuracy performance for all the datasets. More importantly, EXPP-II runs much faster than the other algorithms. This is particularly so for large-size graphs.

We also want to examine sensitivities with outliers. We used the same protocol to construct graphs, except that we added 10 outliers to the graphs. Table IV shows the average matching accuracies and runtimes. EXPP-II is seen to show better matching accuracies than the benchmarked algorithms.

For visual illustration, we provide matching examples in Fig. 4–5. An example for the case of $n = 500$ is enlarged and shown in Fig. 6 for clear visualization.

Datasets	n	m	Runtime (in sec.)				
			Lovař	TPM	Greedy	Rank-1 PC	EXPP
Lcc-Gowalla	196,591	950,327	41.32	0.10	0.02	0.56	2.40
NotreDame	325,729	1,497,134	18.95	0.07	0.03	0.91	3.11
Google	875,713	5,105,039	274.13	0.50	0.13	5.91	25.36
YouTube	1,134,890	2,987,624	1.05e3	1.04	0.13	2.94	18.27
Talk	2,394,385	5,021,410	1.21e3	1.50	0.33	5.02	36.92
LiveJournal	3,997,962	34,681,189	1.29e3	3.82	0.66	22.78	124.68

TABLE II: Average runtime with the DkS problem. n = number of nodes, m = number of edges.



Fig. 4: Illustration of graph matching instances. First row to third row: CMU-House, Car, Motorbike, with $n = 20$; fourth row to fifth row: DTU-House, HerzJesu, with $n = 100$; sixth row to eighth row: Fountain, Semper and Stadium, $n = 110$, 10 outliers. Red lines: correctly matched nodes, green lines: mismatched nodes. Blue dots and magenta dots represent outliers in the two figures respectively.

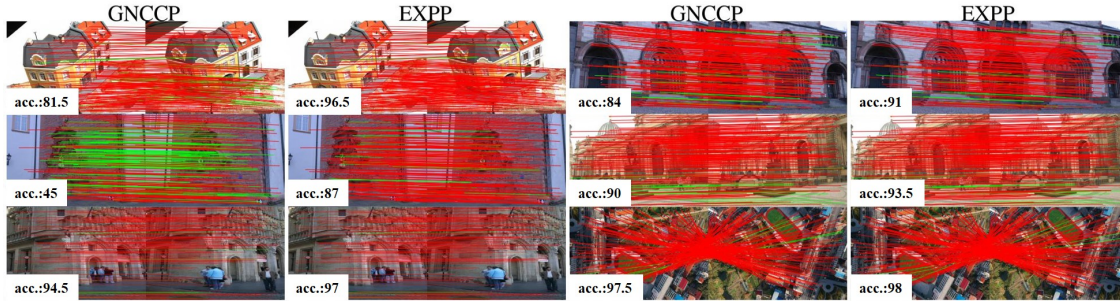


Fig. 5: Illustration of graph matching instances with $n = 200$. First row: DTU-House, HerzJesu; second row: Fountain, Semper; third row: Brussels, Stadium. Red lines: correctly matched nodes, green lines: mismatched nodes.

E. Size-Constrained Clustering

We consider the size-constrained clustering problem described in Section III-C. The benchmarked algorithms are (i) AM: an alternating minimization method [9]; (ii) SC: K -means followed by a post-processing based on mixed integer

linear programming [32]; (iii) E-Kmeans: an extension of K -means [33]. We employ EXPP-II in (16). We revise the objective function (cf. (14)) as $f(x) + \gamma \|X\|_F^2$, with $\gamma = 1.1\sigma_1(R)/\kappa_{[r]}$, such that the revised objective function is convex and EXPP-II starts with a convex problem. The parameter

Data	n	Metric	EXPP-II	GNCCP	PATH	LAGSA
CMU-House	20	acc.	77.00	74.00	54.50	43.00
		time	0.03	0.09	1.21	0.41
Car in PASCAL	20	acc.	82.50	65.00	59.50	61.00
		time	0.02	0.04	1.22	0.76
Motorbike in PASCAL	20	acc.	96.00	80.50	70.00	73.50
		time	0.02	0.08	1.16	0.64
DTU-House	500	acc.	83.62	60.64	-	-
		time	16	711	-	-
HerzJesu in DAISY	500	acc.	82.56	49.18	-	-
		time	11	786	-	-
Semper in DAISY	500	acc.	92.04	54.46	-	-
		time	13	619	-	-
Stadium in SUIRD	500	acc.	91.68	40.46	-	-
		time	10	435	-	-

TABLE III: Graph matching results. n = number of nodes, acc. = matching accuracy in percent, time = runtime in seconds.

Data	Metric	EXPP-II	GNCCP	PATH	LAGSA
DTU-House	acc.	75.70	60.10	29.10	41.60
	time	0.51	3.57	234.55	283.32
HerzJesu in DAISY	acc.	78.50	58.90	22.90	35.50
	time	0.50	3.71	209.89	307.61
Fountain in DAISY	acc.	79.50	59.00	30.90	45.10
	time	0.43	2.85	325.34	278.73
Semper in DAISY	acc.	79.50	59.00	30.90	45.10
	time	0.50	4.42	236.76	371.07
Brussels in DAISY	acc.	77.80	48.40	37.80	40.40
	time	0.48	3.08	207.81	302.79
Stadium in SUIRD	acc.	82.10	63.10	43.00	58.20
	time	0.49	3.79	213.74	308.04

TABLE IV: Graph matching results in the presence of outliers. n = 110 nodes, 10 outliers, acc. = matching accuracy in percent, time = runtime in seconds.

settings of EXPP-II are $\lambda_0 = 10^{-3}K\nu$, $\lambda_{k+1} = \lambda_k + 0.1K\nu$, $\varepsilon_1 = 10^{-4}$, $\varepsilon_3 = 10^{-2}$, $\bar{L} = 200$, $\bar{\lambda} = K\nu$, where ν is specified in Theorem 2 and K is the Lipschitz constant of $f(\mathbf{X})$ on $\tilde{\mathcal{U}}_{\kappa}^{n,r}$.

We evaluate the algorithms on several real-world datasets, some of which are used by the prior work [32]. The results are shown in Table V. We see that the clustering accuracies, defined as the number of correctly clustered data points normalized by the total number of data points, of EXPP-II are generally good compared to the other algorithms. Note that SC fails to work for the Fashion MNIST dataset, due possibly to the large data size. EXPP-II cannot compete with E-Kmeans in terms of runtimes, but otherwise EXPP-II is much faster than the other algorithms.

F. Orthogonal Non-Negative Matrix Factorization

We consider the ONMF problem described in Section III-E. The benchmarked algorithms are (i) ONPMF [11]; (ii) NSNCP [2]: an alternating minimization algorithm for a penalized ONMF formulation; (iii) EP4ORTH [3]: a manifold optimization algorithm for the formulation reviewed in (21). We consider both the EXPP-II formulation in (25) and the EXPP formulation in Part I of this paper. In EXPP we have $h(\mathbf{X}) = -\|\mathbf{X}\|_F^2$ and $\mathcal{X} = \mathcal{B}^{n,r} \cap \mathbb{R}_+^{n \times r}$, and we employ Dykstra's projection algorithm [38] to perform projection onto \mathcal{X} . The parameter settings of EXPP-II (respectively [resp.] EXPP) are $\lambda_0 = 10^{-15}$ (resp. 10^{-5}), $\lambda_{k+1} = 10\lambda_k$ (resp. $5\lambda_k$), $\varepsilon_1 = 10^{-9}$, $\bar{L} = 50$ (resp. 100), $\bar{\lambda} = K\nu$, where ν is specified in Theorem 5 and K is the Lipschitz constant of f on

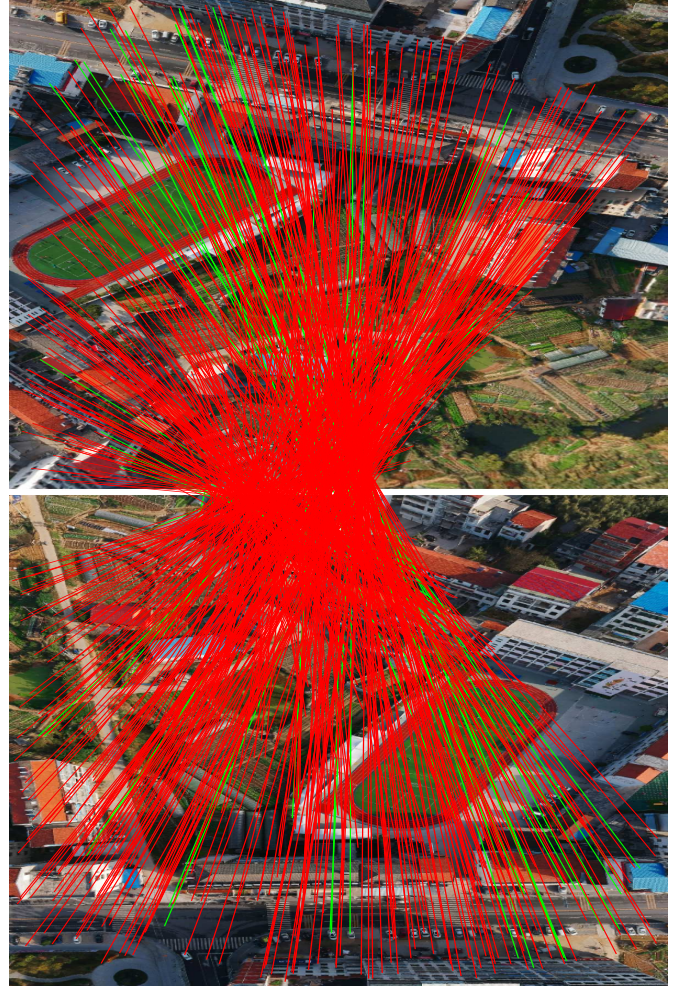


Fig. 6: A graph matching instance of Stadium with $n = 500$. Red lines: correctly matched nodes, green lines: mismatched nodes.

$\tilde{\mathcal{B}}_+^{n,r}$. Unlike the previous applications, we found that EXPP and EXPP-II do not work if we start with a λ_0 such that the problem is convex (the corresponding solution is $\mathbf{0}$ which is not a good starting point). We initialize EXPP and EXPP-II with the NNDSVD method [39] which is commonly adopted by other ONMF methods such as ONPMF and EP4ORTH.

We evaluate the algorithms on various real-world datasets: face images [40], hyperspectral images [41], and text [42]. Table VI describes the dimensions of these datasets. The accuracies and runtimes of the various algorithms are shown in Table VII and Table VIII, respectively. In general, EXPP and EXPP-II perform comparably with the other methods. EXPP-II has faster runtimes while EXPP has higher clustering accuracies.

G. Sparse and Fair Principal Component Analysis

We consider sparse and fair PCAs. Let $\mathbf{Y} \in \mathbb{R}^{m \times n}$ be a data matrix. Sparse PCA considers

$$\min_{\mathbf{X} \in \mathcal{S}^{m \times r}} f(\mathbf{X}) = -\text{tr}(\mathbf{X}^\top \mathbf{R} \mathbf{X}) + \mu \|\mathbf{X}\|_{\ell_1},$$

Datasets	m	n	κ	Metric	AM	EXPP-II	SC	E-Kmeans
Fashion MNIST [34]	784	70,000	$r = 10$; 7K samples/cluster	acc.	0.485	0.525	-	0.450
				time	7792	125	-	51
Iris [35]	4	150	(50, 50, 50)	acc.	92.00	93.33	81.33	93.33
				time	0.153	0.100	0.283	0.006
Glass Identification [36]	9	214	(70, 76, 17, 13, 9, 29)	acc.	37.38	48.60	52.34	45.33
				time	0.390	0.025	0.286	0.007
Balance Scale [37]	4	625	(288, 288, 49)	acc.	60.00	67.52	65.60	61.44
				time	0.226	0.077	0.295	0.012

TABLE V: Size-constrained clustering results. acc. = clustering accuracy in percent, time = runtime in seconds.

	TDT2	Reuters	YALE	ORL	Jasper Ridge	Samson
m	34,642	5,423	1024	1,024	198	156
n	7,809	748	165	150	10,000	9,025
r	25	10	15	15	4	3

TABLE VI: Datasets for ONMF.

Algorithms	TDT2	Reuters	YALE	ORL	Jasper Ridge	Samson
DTPP	67.20	49.20	43.64	48.00	80.28	89.86
ONPMF	69.47	50.40	35.15	42.67	79.71	94.81
NSNCP	67.44	49.20	42.42	56.00	81.49	94.23
EP4ORTH	70.48	50.27	37.58	40.67	82.86	93.65
EXPP-II	67.38	52.14	37.58	46.67	90.27	92.72
EXPP	74.30	64.84	37.58	42.67	94.46	95.86

TABLE VII: ONMF clustering accuracies in percent.

Algorithms	TDT2	Reuters	YALE	ORL	Jasper Ridge	Samson
DTPP	87.46	2.19	0.20	0.11	0.63	0.63
ONPMF	1008.67	31.94	5.02	4.38	25.10	15.90
NSNCP	79.95	1.96	1.79	0.77	5.38	2.25
EP4ORTH	40.66	1.37	1.10	1.08	2.21	2.14
EXPP-II	10.11	0.36	0.25	0.16	2.68	2.34
EXPP	63.60	2.60	0.88	0.68	8.68	6.57

TABLE VIII: ONMF runtimes in seconds.

where $\mu > 0$ is given; $\mathbf{R} = \mathbf{Y}\mathbf{Y}^\top$ is the data correlation matrix. The aim is to recover sparse principal components. Let $\mathbf{Y}_t \in \mathbb{R}^{m \times n_t}$, $t = 1, \dots, T$, be data matrices of different groups. Fair PCA considers

$$\min_{\mathbf{X} \in \mathbb{S}^{m \times r}} f(\mathbf{X}) = \max_{t=1, \dots, T} -\text{tr}(\mathbf{X}^\top \mathbf{R}_t \mathbf{X}),$$

where $\mathbf{R}_t = \mathbf{Y}_t \mathbf{Y}_t^\top$ is the data correlation matrix associated with the t th group. The aim of fair PCA is to reduce possible bias to individual groups. We also consider sparse fair PCA

$$\min_{\mathbf{X} \in \mathbb{S}^{m \times r}} f(\mathbf{X}) = \max_{t=1, \dots, T} -\text{tr}(\mathbf{X}^\top \mathbf{R}_t \mathbf{X}) + \mu \|\mathbf{X}\|_{\ell_1}.$$

We use EXPP to handle all the above PCAs.

The benchmarked algorithms for sparse PCA are: (i) IMRP: iterative minimization of rectangular Procrustes [43]; (ii) Gpower- ℓ_1 : generalized power method with ℓ_1 penalty [44]. The benchmarked algorithms for fair PCA are (i) ARPGDA: the alternating Riemannian projected gradient descent ascent algorithm [45]; (ii) F-FPCA: a subgradient-type algorithm [46]. There is no algorithm to benchmark for sparse fair PCA.

EXPP is implemented by the projected subgradient method since the objective functions are non-smooth. The parameter settings of EXPP are $\lambda_0 = 0$, $\lambda_{k+1} = \lambda_k + 0.1K$, $\varepsilon_1 = 10^{-3}$, $\varepsilon_3 = 10^{-6}$, $\bar{L} = 300$, $\bar{\lambda} = K$, where K is the Lipschitz constant of f over $\mathcal{B}^{n,r}$. The step size of the projected subgradient method is set as $c/\sqrt{l+1}$ for some $c > 0$. We

initialize EXPP with PCA. The EXPP for sparse PCA, fair PCA, and sparse fair PCA are called EXPP-S, EXPP-F, and EXPP-SF, respectively.

We use the following performance metrics to evaluate the performance of the various algorithms: explained variance and cardinality, which indicate the trade-off between sparsity and variance; and minimum variance, which measures fairness. Explained variance is measured as

$$\text{tr}(\mathbf{X}^\top \mathbf{R} \mathbf{X}) / \text{tr}(\hat{\mathbf{X}}^\top \mathbf{R} \hat{\mathbf{X}}),$$

where \mathbf{X} is the principal component matrix recovered by an algorithm; $\hat{\mathbf{X}}$ is that of PCA. Minimum variance is measured as

$$\min_{t=1, \dots, T} \text{tr}(\mathbf{X}^\top \mathbf{R}_t \mathbf{X}).$$

As a minor step, we project \mathbf{X} onto $\mathcal{S}^{n,r}$ before we evaluate the above performance metrics; some algorithms such as Gpower- ℓ_1 do not guarantee semi-orthogonality exactly. The cardinality is measured as the number of elements that are larger than $0.01 \max_{i,j} |x_{ij}|$. We normalize minimum variance and cardinality to $[0, 1]$.

We test the algorithms on three real-world datasets: MNIST, Fashion MNIST, and CIFAR10. Each dataset consists of 10 different classes. The information about the datasets and the corresponding EXPP parameters are provided in Table IX. To create a data matrix, we randomly selected 5 classes. Then, for each class, we randomly selected data points to form one data matrix \mathbf{Y}_t . The data lengths of the different classes are $\{n_t\} = \{5, 10, 500, 1000, 5000\}$; we consider unbalanced data groups. We conducted 100 Monte Carlo runs for each dataset.

The results are shown in Table X. For sparse PCA, we see that EXPP-S performs better than Gpower- ℓ_1 while IMRP outperforms EXPP-S. We note that IMRP adopts a different sparsity penalty. For fair PCA, EXPP-F performs comparably with ARPGDA and works better than F-FPCA. EXPP-SF is seen to be able to strike a balance between sparsity and fairness. The numerical results demonstrate the versatility of EXPP to handle different formulations.

V. CONCLUSION

As the second part of our study, we developed new EXPP formulations for the cases of PPMs, SAMs and NSOMs. They have lower requirements with constraints and can be efficiently handled from the viewpoint of building algorithms. We also demonstrated the utility of EXPP by performing numerical experiments on a variety of applications.

Datasets	Image Size	Number of Images	Data Type	EXPP-S	EXPP-SF	EXPP-F
MNIST [47]	28×28	70000	Digits	$c = 0.025, \mu = 0.3$	$c = 10, \mu = 0.2$	$c = 0.5$
Fashion MNIST [34]	28×28	70000	Clothing	$c = 0.05, \mu = 0.25$	$c = 100, \mu = 0.2$	$c = 20$
CIFAR10 [48]	32×32	60000	real-world objects	$c = 0.1, \mu = 0.25$	$c = 20, \mu = 0.35$	$c = 20$

TABLE IX: Datasets information and EXPP parameter setups. c : the constant in stepsize; μ : sparsity penalty parameter.

Algorithms	MNIST, $r = 10$.				Fashion MNIST, $r = 15$.				CIFAR10, $r = 20$.			
	Expl. Var.	Min. Var.	Card.	Time /s	Expl. Var.	Min. Var.	Card.	Time /s	Expl. Var.	Min. Var.	Card.	Time /s
IMRP	0.7733	0.3149	0.1037	0.6130	0.7890	0.3093	0.0642	1.6011	0.7126	0.5665	0.0462	4.0329
Gpower- ℓ_1	0.6252	0.3010	0.0998	0.4039	0.7070	0.3141	0.1056	1.8938	0.4876	0.3796	0.0445	4.0616
EXPP-S	0.6439	0.2427	0.0834	0.7227	0.7579	0.2696	0.0769	0.4915	0.6250	0.4222	0.0517	0.6302
EXPP-SF	0.4117	0.5639	0.1505	4.9375	0.3336	0.5164	0.0813	3.9784	0.4676	0.5221	0.1868	3.3913
EXPP-F	0.7967	1	0.9912	1.9908	0.7685	1	0.9755	1.9273	0.9354	0.9996	0.9938	1.0820
ARPGDA	0.8143	0.9895	0.9947	0.6747	0.8460	0.9895	0.9948	0.7944	0.9494	1	0.9981	1.5184
F-FPCA	0.9560	0.8189	1	6.3386	0.9490	0.8133	1	6.9849	0.9781	0.9484	1	4.2289

TABLE X: Sparse and/or fair PCA results.

VI. ACKNOWLEDGMENT

The authors would like to thank Prof. Aritra Konar, Prof. Nikos Sidiropoulos, and Prof. Jiang Bo for sharing their source codes and data sets for the DkS, size-constrained clustering, and ONMF problems, respectively.

APPENDIX

A. Some Basic Results

We describe some basic results which will be used in this Appendix. The following lemma encapsulates the proof procedure of Theorems 1, 2 and 5 in a generic form.

Lemma 5 Let $\mathcal{X} \subseteq \mathbb{R}^n$ be a compact set. Let $\mathcal{V} \subseteq \mathcal{X}$ a set. Suppose that there exists a mapping $\mathbf{y} : \mathcal{X} \rightarrow \mathbb{R}^n$, a function $h : \mathcal{X} \rightarrow \mathbb{R}$, and a scalar $\delta > 0$ such that, for any $\mathbf{x} \in \mathcal{X}$,

$$\|\mathbf{x} - \mathbf{y}(\mathbf{x})\|_2 \leq h(\mathbf{x}), \quad (29)$$

$$h(\mathbf{x}) < \delta \implies \mathbf{y}(\mathbf{x}) \in \mathcal{V}. \quad (30)$$

Then we have the inequality

$$\text{dist}(\mathbf{x}, \mathcal{V}) \leq \max \left\{ 1, \frac{B}{\delta} \right\} h(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}, \quad (31)$$

where B is any constant such that

$$B \geq \sup \{ \|\mathbf{x} - \mathbf{z}\|_2 \mid \mathbf{x} \in \mathcal{X}, \mathbf{z} \in \mathcal{V} \}. \quad (32)$$

Proof of Lemma 5: Let $\mathbf{x} \in \mathcal{X}$ be given. Suppose that $h(\mathbf{x}) \leq \delta$. Then $\text{dist}(\mathbf{x}, \mathcal{V}) \leq \|\mathbf{x} - \mathbf{y}(\mathbf{x})\|_2 \leq h(\mathbf{x})$. On the other hand, suppose that $h(\mathbf{x}) > \delta$. Then $\text{dist}(\mathbf{x}, \mathcal{V}) \leq B \leq Bh(\mathbf{x})/\delta$. Combining the two cases leads to (31). ■

It should be noted that the error bound analyses in [3], [4] essentially introduced the procedure in Lemma 5 in a case-specific manner for $\mathcal{S}_+^{n,r}$. Our error bound analysis for $\mathcal{U}_\kappa^{n,r}$ in Part I of this paper also used this procedure. The following results were previously shown and will be used.

1. unit sphere $\mathcal{S}^n = \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\|_2 = 1\}$: Let $\mathbf{x} \in \mathbb{R}^n$, $\|\mathbf{x}\|_2 \leq 1$, be given. Let $\mathbf{y} = \Pi_{\mathcal{S}^n}(\mathbf{x})$. It holds that

$$\text{dist}(\mathbf{x}, \mathcal{S}^n) = \|\mathbf{x} - \mathbf{y}\|_2 \leq 1 - \|\mathbf{x}\|_2 \quad (33)$$

2. unit vector set $\mathcal{U}^n = \{\mathbf{e}_1, \dots, \mathbf{e}_n\} \subseteq \mathbb{R}^n$: Let $\mathbf{x} \in \Delta^n$ be given. Let $\mathbf{y} = \Pi_{\mathcal{U}^n}(\mathbf{x})$. It holds that

$$\text{dist}(\mathbf{x}, \mathcal{U}^n) = \|\mathbf{x} - \mathbf{y}\|_2 \leq 2(1 - \|\mathbf{x}\|_2^2). \quad (34)$$

3. selection vector set $\mathcal{U}_\kappa^n = \{\mathbf{x} \in \{0, 1\}^n \mid \mathbf{1}^\top \mathbf{x} = \kappa\}$, $\kappa \in \{1, \dots, n\}$: Let $\mathbf{x} \in \text{conv}(\mathcal{U}_\kappa^n)$ be given. Let $\mathbf{y} = \Pi_{\mathcal{U}_\kappa^n}(\mathbf{x})$. It holds that

$$\text{dist}(\mathbf{x}, \mathcal{U}_\kappa^n) = \|\mathbf{x} - \mathbf{y}\|_2 \leq 2(\kappa - \|\mathbf{x}\|_2^2). \quad (35)$$

4. scaled unit vector set $\mathcal{W}^n = \{\mathbf{x} \in \mathbb{R}^n \mid \mathbf{x} = \alpha \mathbf{u}, \alpha \in \mathbb{R}, \mathbf{u} \in \mathcal{U}^n\}$: Let $\mathbf{x} \in \mathbb{R}^n$ be given. Let $\mathbf{y} = \Pi_{\mathcal{W}^n}(\mathbf{x}) = x_l \mathbf{e}_l$, where l is such that $|x_l| = \|\mathbf{x}\|_\infty$. It holds that

$$\text{dist}(\mathbf{x}, \mathcal{W}^n) = \|\mathbf{x} - \mathbf{y}\|_2 \leq \|\mathbf{x}\|_1 - \|\mathbf{x}\|_\infty. \quad (36)$$

B. Proof of Theorem 1

Let $\mathbf{X} \in \mathbb{R}^{n \times r}$, with $\mathbf{x}_j \in \Delta^n$ for all j , be given. Let $\mathbf{Y} \in \mathbb{R}^{n \times r}$ be given by

$$\mathbf{y}_j = \Pi_{\mathcal{U}^n}(\mathbf{x}_j), \quad j = 1, \dots, r. \quad (37)$$

Let

$$h(\mathbf{X}) = 2\|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_{\ell_1}. \quad (38)$$

We want to apply Lemma 5. We first show the first condition (29) of Lemma 5. Using the error bound result (34) for \mathcal{U}^n , the error $\|\mathbf{X} - \mathbf{Y}\|_F$ is bounded as

$$\|\mathbf{X} - \mathbf{Y}\|_F \leq \sum_{j=1}^r \|\mathbf{x}_j - \mathbf{y}_j\|_2 \leq 2(r - \|\mathbf{X}\|_F^2).$$

From the proof of Lemma 1 (see Section III-B), it can be seen that $\|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_{\ell_1} \geq r - \|\mathbf{X}\|_F^2$. It follows that the choice in (37)–(38) satisfies the first condition (29) of Lemma 5.

Second we show that the second condition (30) of Lemma 5 is satisfied for some δ . Consider the following lemma.

Lemma 6 Let $\mathbf{Y} \in \{0, 1\}^{n \times r}$ be a matrix with $\mathbf{y}_j \in \{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ for all j , and with $n \geq r$. The singular values of \mathbf{Y} satisfy $\sigma_i(\mathbf{Y})^2 \in \{0, 1, \dots, r\}$ for all i and $\sum_{i=1}^r \sigma_i(\mathbf{Y})^2 = r$.

Proof of Lemma 6: Represent each \mathbf{y}_j by $\mathbf{y}_j = \mathbf{e}_{l_j}$ for some $l_j \in \{1, \dots, n\}$. We have

$$\mathbf{Y}\mathbf{Y}^\top = \sum_{j=1}^r \mathbf{e}_{l_j} \mathbf{e}_{l_j}^\top = \text{Diag}(\mathbf{d}),$$

where $d_i \in \{0, 1, \dots, r\}$ counts the number of times \mathbf{e}_i appears in $\mathbf{y}_1, \dots, \mathbf{y}_r$. As a basic SVD result, the above equation implies that $\sigma_i(\mathbf{Y})^2 = d_{[i]}$ for all i . Since $\sum_{i=1}^r d_i = r$, it follows that $\sum_{i=1}^r \sigma_i(\mathbf{Y})^2 = r$. ■

Lemma 6 implies that, if $\sigma_i(\mathbf{Y}) < \sqrt{2}$ for all i , then $\sigma_1(\mathbf{Y}) = \dots = \sigma_r(\mathbf{Y}) = 1$ and consequently \mathbf{Y} lies in $\mathcal{U}^{n,r}$. Suppose that $h(\mathbf{X}) < \delta$ for some $\delta > 0$. By the Weyl inequality,

$$\begin{aligned} \sigma_i(\mathbf{Y}) &\leq \sigma_i(\mathbf{X}) + \sigma_1(\mathbf{Y} - \mathbf{X}) \leq \sigma_i(\mathbf{X}) + \|\mathbf{Y} - \mathbf{X}\|_F \\ &\leq \sigma_i(\mathbf{X}) + h(\mathbf{X}) \\ &< \sigma_i(\mathbf{X}) + \delta. \end{aligned} \quad (39)$$

Also, it is seen from (38) that

$$\begin{aligned} \frac{1}{2}h(\mathbf{X}) &\geq \|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_F = \|\boldsymbol{\sigma}(\mathbf{X})^2 - \mathbf{1}\|_2 \\ &\geq |\sigma_i(\mathbf{X})^2 - 1|, \end{aligned}$$

for any i . The above equation implies $\sigma_i(\mathbf{X}) < \sqrt{1 + \delta/2}$. Applying this to (39) yields

$$\sigma_i(\mathbf{Y}) < \sqrt{1 + \delta/2} + \delta \leq 1 + \delta/2 + \delta.$$

By choosing $\delta = 2(\sqrt{2} - 1)/3$ such that $\sigma_i(\mathbf{Y}) < \sqrt{2}$ for all i , the matrix \mathbf{Y} must lie in $\mathcal{U}^{n,r}$. Hence, the second condition (30) of Lemma 5 is satisfied for the choice of h and \mathbf{y} in (37)–(38) and for $\delta = 2(\sqrt{2} - 1)/3$.

With the two conditions of Lemma 5 satisfied, we are ready to apply Lemma 5. It can be verified that $B = \sqrt{2r}$ satisfies the requirement (32) in Lemma 5 (see (36) in Part I of this paper). Applying the above h , δ and B to (31) in Lemma 5, we get

$$\text{dist}(\mathbf{X}, \mathcal{U}^{n,r}) \leq \frac{2\sqrt{2}}{2(\sqrt{2} - 1)/3} \sqrt{r} \|\mathbf{X}^\top \mathbf{X} - \mathbf{I}\|_{\ell_1}.$$

The above inequality leads to the final result (10a); note $(2\sqrt{2})/[2(\sqrt{2} - 1)/3] = 10.2426$. The accompanied result (10b) is obtained by applying Lemma 1 to (10a).

C. Proof of Theorem 2

Assume $\kappa_1 \geq \dots \geq \kappa_r$ without loss of generality. Let $\mathbf{X} \in \mathbb{R}^{n \times r}$, with $\mathbf{x}_j \in \text{conv}(\mathcal{U}_{\kappa_j}^n)$ for all j , be given. Let $\mathbf{Y} \in \mathbb{R}^{n \times r}$ be given by

$$\mathbf{y}_j = \Pi_{\mathcal{U}_{\kappa_j}^n}(\mathbf{x}_j), \quad j = 1, \dots, r. \quad (40)$$

Let

$$h(\mathbf{X}) = 2\|\mathbf{X}^\top \mathbf{X} - \text{Diag}(\boldsymbol{\kappa})\|_{\ell_1}. \quad (41)$$

Once again we apply Lemma 5. Following the same proof as Theorem 1, it can be shown that

$$\|\mathbf{X} - \mathbf{Y}\|_F \leq 2(\mathbf{1}^\top \boldsymbol{\kappa} - \|\mathbf{X}\|_F^2) \leq h(\mathbf{X}),$$

where the first inequality is due to the error bound result (35) for \mathcal{U}_{κ}^n ; the second inequality can be seen from the proof of Lemma 1. The first condition (29) of Lemma 5 is thereby

satisfied. The second condition (30) of Lemma 5 is more challenging to show. Consider the following lemma.

Lemma 7 Let $\mathbf{Y} \in \{0, 1\}^{n \times r}$ be a matrix with $\mathbf{y}_j \in \mathcal{U}_{\kappa_j}^n$ for all j , and with $n \geq r$. If $\sum_{i=1}^r \sigma_i(\mathbf{Y})^4 < \mathbf{1}^\top (\boldsymbol{\kappa}^2) + 2$, then it must be true that $\mathbf{y}_i^\top \mathbf{y}_j = 0$ for all $i \neq j$. Consequently, the matrix \mathbf{Y} lies in $\mathcal{U}_{\boldsymbol{\kappa}}^{n,r}$.

Proof of Lemma 7: Consider the matrix product $\mathbf{Y}^\top \mathbf{Y}$. Let $i, j \in \{1, \dots, r\}$, $i \neq j$, be any two distinct indices. On the one hand,

$$\begin{aligned} \|\mathbf{Y}^\top \mathbf{Y}\|_F^2 &\geq \sum_{l=1}^r \|\mathbf{y}_l\|_2^4 + 2(\mathbf{y}_i^\top \mathbf{y}_j)^2 \\ &= \mathbf{1}^\top (\boldsymbol{\kappa}^2) + 2(\mathbf{y}_i^\top \mathbf{y}_j)^2. \end{aligned}$$

On the other hand,

$$\|\mathbf{Y}^\top \mathbf{Y}\|_F^2 = \|\boldsymbol{\sigma}(\mathbf{Y})^2\|_2^2 = \sum_{i=1}^r \sigma_i(\mathbf{Y})^4.$$

The above two equations imply that

$$(\mathbf{y}_i^\top \mathbf{y}_j)^2 \leq \frac{1}{2} [\sum_{i=1}^r \sigma_i(\mathbf{Y})^4 - \mathbf{1}^\top (\boldsymbol{\kappa}^2)].$$

Since $\mathbf{Y} \in \{0, 1\}^{n \times r}$, we have $\mathbf{y}_i^\top \mathbf{y}_j \in \{0, 1, \dots, n\}$. If $\sum_{i=1}^r \sigma_i(\mathbf{Y})^4 < \mathbf{1}^\top (\boldsymbol{\kappa}^2) + 2$, then we are left with $\mathbf{y}_i^\top \mathbf{y}_j = 0$. The proof is complete. ■

Suppose that $h(\mathbf{X}) < \delta$ for some $\delta > 0$. By the same proof as before (cf. (39)),

$$\sigma_i(\mathbf{Y}) < \sigma_i(\mathbf{X}) + \delta. \quad (42)$$

Also, from (41),

$$\begin{aligned} \frac{1}{2}h(\mathbf{X}) &\geq \|\mathbf{X}^\top \mathbf{X} - \text{Diag}(\boldsymbol{\kappa})\|_F \geq \|\boldsymbol{\sigma}(\mathbf{X})^2 - \boldsymbol{\kappa}\|_2 \\ &\geq |\sigma_i(\mathbf{X})^2 - \kappa_i|, \end{aligned}$$

for any i . Here, the second inequality is due to the von Neumann trace inequality result $\|\mathbf{A} - \mathbf{B}\|_F \geq \|\boldsymbol{\sigma}(\mathbf{A}) - \boldsymbol{\sigma}(\mathbf{B})\|_2$. The above equation implies $\sigma_i(\mathbf{X}) < \sqrt{\kappa_i + \delta/2}$. Putting it into (42) leads to

$$\begin{aligned} \sigma_i(\mathbf{Y})^2 &< \kappa_i + \delta/2 + 2\delta\sqrt{\kappa_i + \delta/2} + \delta^2 \\ &\leq \kappa_i + \delta/2 + 2\delta[\sqrt{\kappa_i} + \delta/(2\sqrt{\kappa_i})] + \delta^2 \\ &= \kappa_i + (1/2 + 2\sqrt{\kappa_i})\delta + (1 + 1/\sqrt{\kappa_i})\delta^2, \end{aligned} \quad (43)$$

where the second equation is due to the inequality $\sqrt{a+b} \leq \sqrt{a} + b/\sqrt{a}$ for $a > 0$ and $b \geq 0$. Suppose that

$$(1 + 1/\sqrt{\kappa_i})^{1/2} \delta \leq 1, \quad \forall i. \quad (44)$$

Eq. (43) can be further bounded as

$$\sigma_i(\mathbf{Y})^2 < \kappa_i + \underbrace{[1/2 + 2\sqrt{\kappa_i} + (1 + 1/\sqrt{\kappa_i})^{1/2}]}_{:=\alpha_i} \delta, \quad (45)$$

Taking square on (45) gives

$$\sigma_i(\mathbf{Y})^4 < \kappa_i^2 + 2\kappa_i\alpha_i\delta + \alpha_i^2\delta^2. \quad (46)$$

Suppose that

$$\alpha_i\delta \leq 1, \quad \forall i. \quad (47)$$

Then we further bound (46) as

$$\sigma_i(\mathbf{Y})^4 < \kappa_i^2 + \underbrace{(2\kappa_i\alpha_i + \alpha_i^2)}_{:=\beta_i} \delta, \quad (48)$$

and consequently,

$$\sum_{i=1}^r \sigma_i(\mathbf{Y})^4 < \mathbf{1}^\top (\boldsymbol{\kappa}^2) + (\mathbf{1}^\top \boldsymbol{\beta}) \delta.$$

Let us choose

$$\delta = 2/(\mathbf{1}^\top \boldsymbol{\beta}).$$

It can be verified that this choice of δ satisfies the requirements (44) and (47). By invoking Lemma 7, we see that \mathbf{Y} lies in $\mathcal{U}_{\boldsymbol{\kappa}}^{n,r}$. Hence we have the second condition (30) of Lemma 5 satisfied.

Finally we assemble the components together to obtain the error bound. Let $B = \sqrt{2\mathbf{1}^\top \boldsymbol{\kappa}}$, which can be verified to satisfy (32). We express β_i as

$$\begin{aligned} \beta_i &= (1 + 2\kappa_i) [1/2 + 2\sqrt{\kappa_i} + \underbrace{(1 + 1/\sqrt{\kappa_i})^{\frac{1}{2}}}_{\leq \sqrt{2} \leq 1.5}] \\ &\leq 2(1 + 2\kappa_i)(1 + \sqrt{\kappa_i}). \end{aligned}$$

We have

$$\frac{B}{\delta} \leq \sqrt{2\mathbf{1}^\top \boldsymbol{\kappa}} \underbrace{\left(\sum_{j=1}^r (1 + 2\kappa_j)(1 + \sqrt{\kappa_j}) \right)}_{:=\gamma}.$$

Note that the right-hand side of the above equation is greater than 1. Applying the above h , δ and B to (31) in Lemma 5 leads to

$$\text{dist}(\mathbf{X}, \mathcal{U}^{n,r}) \leq 2\sqrt{2\mathbf{1}^\top \boldsymbol{\kappa}} \gamma \|\mathbf{X}^\top \mathbf{X} - \text{Diag}(\boldsymbol{\kappa})\|_{\ell_1},$$

and consequently the final result (15a). In addition, applying Lemma 1 to (15a) gives (15b).

D. Proof of Theorem 5

Let $\mathbf{X} \in \mathbb{R}_+^{n \times r}$, with $\|\mathbf{x}_j\|_2 \leq 1$ for all j , be given. Let $\mathbf{W} \in \mathbb{R}^{n \times r}$ be given by $\mathbf{w}_i = x_{i,l_i} \mathbf{e}_{l_i}$ for all i , where l_i is such that $x_{i,l_i} = \max\{x_{i,1}, \dots, x_{i,r}\}$. Let $\mathbf{Y} \in \mathbb{R}^{n \times r}$ be given by

$$\mathbf{y}_j = \begin{cases} \mathbf{w}_j / \|\mathbf{w}_j\|_2, & \mathbf{w}_j \neq \mathbf{0} \\ \mathbf{u}, & \mathbf{w}_j = \mathbf{0} \end{cases}$$

for all j , where \mathbf{u} denotes any non-negative unit ℓ_2 norm vector. It is worth noting that

$$\bar{\mathbf{w}}_i = \Pi_{\mathcal{W}^r}(\bar{\mathbf{x}}_i), \quad \mathbf{y}_j = \Pi_{\mathcal{S}^n}(\mathbf{w}_j),$$

for all i, j . Let

$$h(\mathbf{X}) = \sqrt{6} \left[\sum_{j=1}^r c_1(\mathbf{x}_j)^2 + \sum_{i=1}^n \rho_1(\bar{\mathbf{x}}_i)^2 \right]^{\frac{1}{2}}, \quad (49)$$

where $c_1(\mathbf{x}) = 1 - \|\mathbf{x}\|_2$, $\rho_1(\mathbf{x}) = \|\mathbf{x}\|_1 - \|\mathbf{x}\|_\infty$. As before we apply Lemma 5. We begin by showing the first condition (29) of Lemma 5. Using the error bounds (36) and (33) for \mathcal{W}^r and \mathcal{S}^n , respectively, it holds that

$$\|\bar{\mathbf{x}}_i - \bar{\mathbf{w}}_i\|_2 \leq \rho_1(\bar{\mathbf{x}}_i), \quad (50)$$

$$\begin{aligned} \|\mathbf{w}_j - \mathbf{y}_j\|_2 &\leq 1 - \|\mathbf{w}_j\|_2 \\ &\leq 1 - \|\mathbf{x}_j\|_2 + \|\mathbf{w}_j - \mathbf{x}_j\|_2 \end{aligned} \quad (51)$$

$$\leq \sqrt{2} [c_1(\mathbf{x}_j)^2 + \|\mathbf{w}_j - \mathbf{x}_j\|_2^2]^{\frac{1}{2}}, \quad (52)$$

where (51) is due to the triangle inequality; (52) is due to $(|a| + |b|)^2 \leq 2(|a|^2 + |b|^2)$. This gives rise to

$$\begin{aligned} \|\mathbf{X} - \mathbf{Y}\|_F &\leq \|\mathbf{X} - \mathbf{W}\|_F + \|\mathbf{W} - \mathbf{Y}\|_F \\ &\leq \sqrt{2} (\|\mathbf{X} - \mathbf{W}\|_F^2 + \|\mathbf{W} - \mathbf{Y}\|_F^2)^{\frac{1}{2}} \\ &\leq \sqrt{2} \left[3\|\mathbf{X} - \mathbf{W}\|_F^2 + 2\sum_{j=1}^r c_1(\mathbf{x}_j)^2 \right]^{\frac{1}{2}} \\ &\leq \sqrt{6} \left[\sum_{i=1}^n \rho_1(\bar{\mathbf{x}}_i)^2 + \sum_{j=1}^r c_1(\mathbf{x}_j)^2 \right]^{\frac{1}{2}} \\ &= h(\mathbf{X}), \end{aligned} \quad (53)$$

where the third inequality is due to (52); the fourth inequality is due to (50). The first condition (29) of Lemma 5 is obtained.

The proof of the second condition (30) of Lemma 5 is as follows. From the way \mathbf{Y} is constructed, we notice that \mathbf{Y} is a non-negative semi-orthogonal matrix if $\|\mathbf{w}_j\|_2 > 0$ for all j . Suppose that $h(\mathbf{X}) < \delta$ for some $\delta > 0$. As a reverse version of (53), we have, for any j ,

$$\begin{aligned} \delta &> \sqrt{6} \left[\|\mathbf{X} - \mathbf{W}\|_F^2 + \sum_{j=1}^r c_1(\mathbf{x}_j)^2 \right]^{\frac{1}{2}} \\ &\geq \sqrt{6} [\|\mathbf{x}_j - \mathbf{w}_j\|_2^2 + c_1(\mathbf{x}_j)^2]^{\frac{1}{2}} \\ &\geq \sqrt{3} [\|\mathbf{x}_j - \mathbf{w}_j\|_2 + c_1(\mathbf{x}_j)] \\ &= \sqrt{3} (\|\mathbf{x}_j - \mathbf{w}_j\|_2 + 1 - \|\mathbf{x}_j\|_2) \\ &\geq \sqrt{3} (1 - \|\mathbf{w}_j\|_2). \end{aligned}$$

Here, the third inequality is due to $(|a| + |b|)^2 \leq 2(|a|^2 + |b|^2)$ and $c_1(\mathbf{x}_j) \geq 0$ for any $\|\mathbf{x}_j\|_2 \leq 1$; the last inequality is due to the triangle inequality. The above equation implies that, if $\delta = \sqrt{3}$, then $\|\mathbf{w}_j\|_2 > 0$. The second condition (30) of Lemma 5 is shown.

Finally we put together the components. We choose $B = \sqrt{2r}$, which can be verified to satisfy (32). The inequality (31) associated with the above δ , h and B is

$$\text{dist}(\mathbf{X}, \mathcal{S}_+^{n,r}) \leq \nu \underbrace{\left[\sum_{j=1}^r c_1(\mathbf{x}_j)^2 + \sum_{i=1}^n \rho_1(\bar{\mathbf{x}}_i)^2 \right]^{\frac{1}{2}}}_{=\psi_2(\mathbf{X})},$$

where $\nu = \max\{\sqrt{6}, 2\sqrt{r}\}$. Furthermore, as a basic norm result, we have $\psi_2(\mathbf{X}) \leq \sum_{j=1}^r c_1(\mathbf{x}_j) + \sum_{i=1}^n \rho_1(\bar{\mathbf{x}}_i) = \psi_1(\mathbf{X})$. The proof of Theorem 5 is complete.

E. Relationship of the Proof of Theorem 5 with Prior Work

A key step of the proof of Theorem 5 in Appendix D is the construction of \mathbf{W} and \mathbf{Y} . This step is identical to that in the prior studies [3], [4]. The important difference lies in the bound in (50). We use the error bound for scaled unit vectors, given in Lemma 2, to derive the bound in (50). The prior studies used another bound: For any $\mathbf{X} \in \mathbb{R}_+^{n \times r}$, it holds that

$$\|\mathbf{X} - \mathbf{W}\|_F^2 \leq \sum_{j=1}^r \sum_{l=1, l \neq j}^r \mathbf{x}_j^\top \mathbf{x}_l = \|\mathbf{X} \mathbf{1}\|_2^2 - \|\mathbf{X}\|_F^2; \quad (54)$$

see [3, Lemma 3.1] and [4, Lemma 7]. An oversimplified way to understand the prior studies is as follows: replace our bound

(50) with (54), and then proceed with the error bound proof in Appendix D. There are differences with the fine details and the reader is referred to the prior studies [3], [4] for the details. Our bound (50) may be seen as an improvement over (54). From (50) and the proof of Lemma 4, we have

$$\|\mathbf{X} - \mathbf{W}\|_F^2 \leq \sum_{i=1}^n \rho_1(\bar{\mathbf{x}}_i)^2 \leq \|\mathbf{X}\mathbf{1}\|_2^2 - \|\mathbf{X}\|_F^2.$$

REFERENCES

- [1] J. Liu, Y. Liu, W.-K. Ma, M. Shao, and A. M.-C. So, "Extreme point pursuit—Part I: A framework for constant modulus optimization," submitted to *IEEE Trans. Signal Process.*, 2024.
- [2] S. Wang, T.-H. Chang, Y. Cui, and J.-S. Pang, "Clustering by orthogonal NMF model and non-convex penalty optimization," *IEEE Trans. Signal Process.*, vol. 69, pp. 5273–5288, 2021.
- [3] B. Jiang, X. Meng, Z. Wen, and X. Chen, "An exact penalty approach for optimization with nonnegative orthogonality constraints," *Math Program.*, vol. 198, no. 1, pp. 855–897, 2023.
- [4] X. Chen, Y. He, and Z. Zhang, "Tight error bounds for the sign-constrained Stiefel manifold," *arXiv preprint arXiv:2210.05164*, 2022.
- [5] B. Jiang, Y.-F. Liu, and Z. Wen, " ℓ_p -norm regularization algorithms for optimization over permutation matrices," *SIAM J. Optim.*, vol. 26, no. 4, pp. 2284–2313, 2016.
- [6] M. Zaslavskiy, F. Bach, and J.-P. Vert, "A path following algorithm for the graph matching problem," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 31, no. 12, pp. 2227–2242, 2008.
- [7] L. F. McGinnis, "Implementation and testing of a primal-dual algorithm for the assignment problem," *Oper. Res.*, vol. 31, no. 2, pp. 277–291, 1983.
- [8] L. Condat, "Fast projection onto the simplex and the ℓ_1 ball," *Math. Program.*, vol. 158, no. 1-2, pp. 575–585, 2016.
- [9] N. D. Sidiropoulos and E. Tsakonas, "Signal processing and optimization tools for conference review and session assignment," *IEEE Signal Process. Mag.*, vol. 32, no. 3, pp. 141–155, 2015.
- [10] A. Konar and N. D. Sidiropoulos, "Exploring the subgraph density-size trade-off via the Lovász extension," in *Proc. Int. Conf. Web Search Data Min.*, pp. 743–751, 2021.
- [11] F. Pompili, N. Gillis, P.-A. Absil, and F. Glineur, "Two algorithms for orthogonal nonnegative matrix factorization with application to clustering," *Neurocomput.*, vol. 141, pp. 15–25, 2014.
- [12] Y. Xu and W. Yin, "A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion," *J. Imag. Sci.*, vol. 6, no. 3, pp. 1758–1789, 2013.
- [13] M. Shao, Q. Li, W.-K. Ma, and A. M.-C. So, "A framework for one-bit and constant-envelope precoding over multiuser massive MISO channels," *IEEE Trans. Signal Process.*, vol. 67, no. 20, pp. 5309–5324, 2019.
- [14] M. Shao and W.-K. Ma, "Binary MIMO detection via homotopy optimization and its deep adaptation," *IEEE Trans. Signal Process.*, vol. 69, pp. 781–796, 2020.
- [15] A. Beck, *First-order Methods in Optimization*. SIAM, 2017.
- [16] W.-K. Ma, P.-C. Ching, and Z. Ding, "Semidefinite relaxation based multiuser detection for M-ary PSK multiuser systems," *IEEE Trans. Signal Process.*, vol. 52, no. 10, pp. 2862–2872, 2004.
- [17] C. Jeon, R. Ghods, A. Maleki, and C. Studer, "Optimal data detection in large MIMO," *arXiv preprint arXiv:1811.01917*, 2018.
- [18] S. Rangan, P. Schniter, A. K. Fletcher, and S. Sarkar, "On the convergence of approximate message passing with arbitrary matrices," *IEEE Trans. Inf. Theory*, vol. 65, no. 9, pp. 5339–5351, 2019.
- [19] S. L. Loyka, "Channel capacity of MIMO architecture using the exponential correlation matrix," *IEEE Commun. Lett.*, vol. 5, no. 9, pp. 369–371, 2001.
- [20] U. Feige, D. Peleg, and G. Kortsarz, "The dense k -subgraph problem," *Algorithm.*, vol. 29, pp. 410–421, 2001.
- [21] X.-T. Yuan and T. Zhang, "Truncated power method for sparse eigenvalue problems," *J. Mach. Learn. Res.*, vol. 14, no. 4, 2013.
- [22] D. Papailiopoulos, I. Mitliagkas, A. Dimakis, and C. Caramanis, "Finding dense subgraphs via low-rank bilinear optimization," in *Proc. Int. Conf. Mach. Learn.*, pp. 1890–1898, 2014.
- [23] J. Leskovec and A. Krevl, "SNAP Datasets: Stanford large network dataset collection," 2014.
- [24] Z.-Y. Liu and H. Qiao, "GNCCP—graduated nonconvexity and concavity procedure," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 6, pp. 1258–1267, 2013.
- [25] Y. Xia, "An efficient continuation method for quadratic assignment problems," *Comput. Oper. Res.*, vol. 37, no. 6, pp. 1027–1032, 2010.
- [26] F. Zhou and F. De la Torre, "Factorized graph matching," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 9, pp. 1774–1789, 2015.
- [27] "CMU/VASC Image Database." <http://vasc.ri.cmu.edu/ldb/html/motion/house/index.html>.
- [28] M. Leordeanu, R. Sukthankar, and M. Hebert, "Unsupervised learning for graph matching," *Int. J. Comput. Vis.*, vol. 96, pp. 28–45, 2012.
- [29] H. Aanæs, R. R. Jensen, G. Vogiatzis, E. Tola, and A. B. Dahl, "Large-scale data for multiple-view stereopsis," *Int. J. Comput. Vis.*, vol. 120, pp. 153–168, 2016.
- [30] E. Tola, V. Lepetit, and P. Fua, "DAISY: An efficient dense descriptor applied to wide-baseline stereo," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 5, pp. 815–830, 2009.
- [31] Y. Liu, X. Gong, J. Chen, S. Chen, and Y. Yang, "Rotation-invariant siamese network for low-altitude remote-sensing image registration," *IEEE J. Sel. Top. Appl. Earth. Obs. Remote. Sens.*, vol. 13, pp. 5746–5758, 2020.
- [32] S. Zhu, D. Wang, and T. Li, "Data clustering with size constraints," *Knowl. Based Syst.*, vol. 23, no. 8, pp. 883–889, 2010.
- [33] N. Ganganath, C.-T. Cheng, and K. T. Chi, "Data clustering with cluster size constraints using a modified k -means algorithm," in *Proc. Int. Conf. Cyber-Enabled Distrib. Comput. Knowl. Discovery*, pp. 158–161, 2014.
- [34] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [35] R. A. Fisher, "Iris," UCI Machine Learning Repository, 1988.
- [36] B. German, "Glass Identification," UCI Machine Learning Repository, 1987.
- [37] R. Sieglar, "Balance Scale," UCI Machine Learning Repository, 1994.
- [38] J. P. Boyle and R. L. Dykstra, "A method for finding projections onto the intersection of convex sets in Hilbert spaces," in *Lecture Notes Statist.*, pp. 28–47, Springer, 1986.
- [39] C. Boutsidis and E. Gallopoulos, "SVD based initialization: A head start for nonnegative matrix factorization," *Pattern Recognit.*, vol. 41, no. 4, pp. 1350–1362, 2008.
- [40] D. Cai, X. He, J. Han, and H.-J. Zhang, "Orthogonal Laplacianfaces for face recognition," *IEEE Trans. Image Process.*, vol. 15, no. 11, pp. 3608–3614, 2006.
- [41] F. Zhu, Y. Wang, B. Fan, S. Xiang, G. Meng, and C. Pan, "Spectral unmixing via data-guided sparsity," *IEEE Trans. Image Process.*, vol. 23, no. 12, pp. 5412–5427, 2014.
- [42] D. Cai, X. He, and J. Han, "Locally consistent concept factorization for document clustering," *IEEE Trans. Knowl. Data Eng.*, vol. 23, no. 6, pp. 902–913, 2010.
- [43] K. Benidis, Y. Sun, P. Babu, and D. P. Palomar, "Orthogonal sparse PCA and covariance estimation via Procrustes reformulation," *IEEE Trans. Signal Process.*, vol. 64, no. 23, pp. 6211–6226, 2016.
- [44] M. Journée, Y. Nesterov, P. Richtárik, and R. Sepulchre, "Generalized power method for sparse principal component analysis," *J. Mach. Learn. Res.*, vol. 11, no. 2, pp. 517–553, 2010.
- [45] M. Xu, B. Jiang, W. Pu, Y.-F. Liu, and A. M.-C. So, "An efficient alternating Riemannian/projected gradient descent ascent algorithm for fair principal component analysis," in *Proc. Int. Conf. Acoust. Speech Signal Process.*, pp. 7195–7199, IEEE, 2024.
- [46] G. Zalcberg and A. Wiesel, "Fair principal component analysis and filter design," *IEEE Trans. Signal Process.*, vol. 69, pp. 4835–4842, 2021.
- [47] Y. LeCun, C. Cortes, and C. Burges, "MNIST handwritten digit database," *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, vol. 2, 2010.
- [48] A. Krizhevsky, G. Hinton, et al., "Learning multiple layers of features from tiny images," *MIT and NYU, Tech. Rep.*, 2009.



Junbin Liu received the B.S. degree from the South China University of Technology, Guangzhou, China, in 2017 and M.S. degree from the University of Chinese Academy of Sciences in 2020. He is currently pursuing the Ph.D. degree with the Department of Electronic Engineering, the Chinese University of Hong Kong, under the supervision of Professor Wing-Kin Ma.

His research interests encompass statistical signal processing methods, optimization theories, and their wide-ranging applications.



Ya Liu received the B.S. degree from the University of Electronic Science and Technology of China (UESTC), Chengdu, China, in 2018. She is currently pursuing a Ph.D. degree with the Department of Electronic Engineering at the Chinese University of Hong Kong (CUHK), under the supervision of Professor Wing-Kin Ma.

Her research focuses include matrix factorization, signal processing, and optimization, along with their diverse applications.



Wing-Kin Ma (Fellow, IEEE) is currently a Professor with the Department of Electronic Engineering, The Chinese University of Hong Kong (CUHK), Hong Kong. His research interests include signal processing, optimization and communications, with recent focus on i) optimization and statistical aspects with structured matrix factorization, with application to remote sensing and data science; and ii) coarsely quantized MIMO transceiver designs.

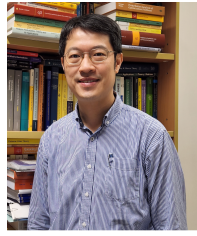
Dr. Ma has rich experience in editorial service, such as an Associate Editor, and then later, a Senior Area Editor, and then from 2021 to 2023, the Editor-in-Chief of the IEEE TRANSACTIONS ON SIGNAL PROCESSING, and many others. He was a Tutorial Speaker in EUSIPCO 2011 and ICASSP 2014, and was an IEEE Signal Processing Society (SPS) Distinguished Lecturer in 2018–2019. He was the recipient of the Research Excellence Award 2013–2014 by CUHK, the 2015 IEEE Signal Processing Magazine Best Paper Award, the 2016 IEEE Signal Processing Letters Best Paper Award, and the 2018 IEEE SPS Best Paper Award. He served as a Member of the Signal Processing for Communications and Networking Technical Committee (SPCOM-TC) in 2015–2020, a Member of Signal Processing Theory and Methods Technical Committee (SPTM-TC) in 2012–2017, the SPS Regional Director-at-Large for Region 10 in 2020–2021, and a Technical Program Co-Chair of ICASSP 2023. He co-founded and co-organized One World Signal Processing in 2020, a virtual seminar series for signal processing.



Mingjie Shao (S'16-M'20) received the B.S. degree from the Xidian University, Xi'an, China, in 2015 and Ph.D. degree from the Chinese University of Hong Kong (CUHK) in 2020. He was a Postdoctoral Fellow with the Department of Electronic Engineering, CUHK from 2020 to 2023. He is currently a Research Professor (Qilu Young Scholar) in the School of Information Science and Engineering, Shandong University, Qingdao, China. He was the recipient of the Hong Kong PhD Fellowship Scheme (HKPFS) from August 2015. He was listed in the

Student Best Paper Finalists in ICASSP 2017.

His research interests focus on convex and non-convex optimization, statistical signal processing and machine learning for wireless communication.



Anthony Man-Cho So (Fellow, IEEE) received the B.S.E. degree in computer science from Princeton University, Princeton, NJ, USA, with minors in applied and computational mathematics, engineering and management systems, and German language and culture; the M.Sc. degree in computer science and the Ph.D. degree in computer science with a Ph.D. minor in mathematics from Stanford University, Stanford, CA, USA. He is currently Dean of the Graduate School, Deputy Master of Morningside College, and a Professor with the Department of

Systems Engineering and Engineering Management at The Chinese University of Hong Kong (CUHK), Hong Kong SAR, China. His research interests include optimization theory and its applications in various areas of science and engineering, including computational geometry, machine learning, signal processing, and statistics.

Dr. So is a Fellow of the Hong Kong Institution of Engineers. He is the recipient of a number of research and teaching awards, including the 2024 INFORMS Computing Society Prize, the SIAM Review SIGEST Award in 2024, the 2022 University Grants Committee Teaching Award, the 2018 IEEE Signal Processing Society Best Paper Award, the 2015 IEEE Signal Processing Society Signal Processing Magazine Best Paper Award, the 2014 IEEE Communications Society Asia-Pacific Outstanding Paper Award, the 2013 CUHK Vice-Chancellor's Exemplary Teaching Award, and the 2010 INFORMS Optimization Society Optimization Prize for Young Researchers. He currently serves on the Editorial Boards of *Journal of Global Optimization*, *Mathematics of Operations Research*, *Mathematical Programming*, *Open Journal of Mathematical Optimization*, *Optimization Methods and Software*, and *SIAM Journal on Optimization*. He was also the Lead Guest Editor of the Special Issue on Non-Convex Optimization for Signal Processing and Machine Learning of the IEEE SIGNAL PROCESSING MAGAZINE and a Guest Editor of the Special Issue on Advanced Optimization Theory and Algorithms for Next Generation Wireless Communication Networks of the IEEE JOURNAL ON SELECTED AREAS IN COMMUNICATIONS.