

# Non-smooth Optimization over Stiefel Manifolds

Fariba Zohrizadeh, Mohsen Kheirandishfard, Farhad Kamangar, and Ramtin Madani

**Abstract**—This paper is concerned with the class of non-convex optimization problems with orthogonality constraints. We develop computationally efficient relaxations that transform non-convex orthogonality constrained problems into polynomial-time solvable surrogates. A novel penalization technique is used to enforce feasibility and derive certain conditions under which the constraints of the original non-convex problem are guaranteed to be satisfied. Moreover, we extend our approach to a feasibility-preserving algorithm that solves a sequence of penalized relaxations to obtain feasible and near optimal points. Experimental results on synthetic and real datasets demonstrate the effectiveness of the proposed approach on the two practical applications of discriminative dimensionality reduction and graph clustering.

## I. INTRODUCTION

Consider the following optimization problem

$$\underset{\mathbf{P} \in \mathbb{R}^{n \times m}}{\text{minimize}} \quad \bar{f}_0(\mathbf{P}) + g_0(\mathbf{P}) \quad (1a)$$

$$\text{subject to} \quad \bar{f}_k(\mathbf{P}) \leq 0, \quad k \in \{1, \dots, p\}, \quad (1b)$$

$$\mathbf{P}^\top \mathbf{P} = \mathbf{I}_m, \quad (1c)$$

where  $g_0: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$  is a convex piecewise linear function and each  $\bar{f}_k: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$  is an arbitrary quadratic function of the form  $\bar{f}_k(\mathbf{P}) \triangleq \langle \mathbf{M}_k, \mathbf{P}\mathbf{P}^\top \rangle + \langle \mathbf{N}_k, \mathbf{P} \rangle + q_k$ , for every  $k \in \{0, 1, \dots, p\}$ , and  $\{\mathbf{M}_k \in \mathbb{S}_n\}_{k=0}^p$ ,  $\{\mathbf{N}_k \in \mathbb{R}^{n \times m}\}_{k=0}^p$  and  $\{q_k \in \mathbb{R}\}_{k=0}^p$  are given. With no loss of generality, we assume that  $q_0 = 0$  and write  $g_0$  in the form of  $g_0(\mathbf{P}) = \|\alpha(\mathbf{P}) + \mathbf{b}\|_1$ , where  $\mathbf{b} \in \mathbb{R}^w$  is a given vector,  $\alpha: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}^w$  is a linear matrix function defined as  $\alpha(\mathbf{Y}) \triangleq \sum_{i=1}^w \langle \mathbf{A}_i, \mathbf{Y} \rangle \mathbf{e}_i$ , the matrices  $\{\mathbf{A}_i \in \mathbb{R}^{n \times m}\}_{i=1}^w$  are given, and  $\{\mathbf{e}_i \in \mathbb{R}^w\}_{i=1}^w$  represent the standard basis for  $\mathbb{R}^w$ . The formulation (1a)–(1c) encompasses a broad class of computationally-hard optimization problems with a variety of practical applications in discriminative dimensionality reduction [2], graph matching [3], feature selection [4], [5], compressed modes [6], [7], among other areas of machine learning.

The majority of methods in the literature are focused on a special case of (1a)–(1c) that involves the minimization of a convex and smooth objective function over non-convex sets of the form  $\mathcal{S}_{n,m} \triangleq \{\mathbf{P} \in \mathbb{R}^{n \times m} \mid \mathbf{P}^\top \mathbf{P} = \mathbf{I}_m\}$ , known as the Stiefel manifolds. There are various iterative local search algorithms which preserve the structure of Stiefel manifolds via geodesics steps [8] or retractions [9]. Although these algorithms exhibit satisfactory performance in dealing with orthogonality constraints, they mostly restrict the objective function to the class of smooth functions and are not compatible with additional constraints [10]. To overcome these limitations,

general algorithms are proposed that work with either smooth or non-smooth objective functions [2], [7]. The paper [2] uses a family of semidefinite programming (SDP) problems to generate a converging sequence of points on Stiefel manifolds. The paper [7] introduces an inner-outer iteration scheme for solving  $\ell_1$ -regularized optimization problems with orthogonality constraints based on the augmented Lagrangian method from [11] and the proximal alternating minimization technique from [12]. Moreover, a series of splitting techniques are proposed in [6] and [13] that can efficiently handle non-smooth objective functions. They partition the problem into multiple sub-problems with analytical solutions and employ Bregman iterations [14] or its variants to obtain optimal solutions for orthogonality-constrained problems. In the more recent paper [15], an extended proximal alternating linearized minimization method is introduced to minimize convex functions subject to linear constraints and generalized orthogonality constraints.

The success of related sequential frameworks and penalized relaxations for nonconvex optimization is demonstrated in [16]–[18], and in [19] for quadratically-constrained quadratic programming. In [16], a sequential framework is introduced for solving BMIs without theoretical guarantees. In [17], [18], this approach is further investigated and theoretical results are offered through the notion of generalized Mangasarian-Fromovitz regularity condition. Another sequential SDP-based algorithm for pattern recognition is introduced in [2] that is not feasibility preserving.

### A. Contributions

Differentiated from the existing literature, we propose a computational approach with theoretical analysis for solving problems of the form (1a)–(1c), that guarantees the recovery of feasible points. The proposed approach generalizes the existing literature by including additional quadratic inequality constraints. The core of our approach is based on a novel and computationally efficient parabolic relaxation which transforms the non-convex problem (1a)–(1c) into a convex quadratically-constrained quadratic program (QCQP). To ensure that the solution of the relaxed problem is feasible for (1a)–(1c), we incorporate a penalty term into the objective function and derive certain conditions that guarantee the recovery of feasible points. Moreover, under certain conditions, we prove that by starting from any arbitrary initial point on a Stiefel manifold (not necessarily feasible), a sequence of penalized relaxations can be solved to find a feasible and near-optimal point. Unlike the existing algorithms, if mild assumptions are satisfied, the proposed sequential scheme is feasibility-preserving and improves the objective monotonically at every step. To corroborate the effectiveness of our method, we perform experiments on two practical applications

Parts of this paper have appeared in the conference paper [1]. Compared with the conference version, the new additions to this paper are detailed proofs and major theoretical results that guarantee the convergence of the proposed algorithm.

with both synthetic and real datasets. The experimental results demonstrate that the proposed approach exhibits comparable results for both applications.

## B. Notation

Throughout this paper, the scalars, vectors, and matrices are shown by italic, bold lower-case and bold upper-case letters, respectively. The symbols  $\mathbb{R}^n$ ,  $\mathbb{R}^{n \times m}$ ,  $\mathbb{S}_n$ , and  $\mathbb{S}_n^+$  denote the set of real  $n$ -dimensional vectors, real  $n \times m$  matrices, real symmetric  $n \times n$  matrices, and real positive semidefinite matrices, respectively. The symbols  $\text{tr}\{\cdot\}$  and  $(\cdot)^\top$  are indicative of the trace and transpose operators, respectively. Given a vector  $\mathbf{a}$  and a matrix  $\mathbf{A}$ , the symbols  $a_i$  and  $A_{ij}$ , respectively, refer to the  $i^{\text{th}}$  element of  $\mathbf{a}$  and the  $(i, j)^{\text{th}}$  element of  $\mathbf{A}$ . The notation  $\mathbf{A} \succeq 0$  states that  $\mathbf{A}$  is symmetric positive semidefinite. Given matrices  $\mathbf{A}$  and  $\mathbf{B}$  of the same size,  $\langle \mathbf{A}, \mathbf{B} \rangle \triangleq \text{tr}\{\mathbf{A}^\top \mathbf{B}\}$  and  $\mathbf{A} \circ \mathbf{B}$ , respectively, denote the Frobenius inner-product and the Hadamard product of  $\mathbf{A}$  and  $\mathbf{B}$ . The operator  $\text{diag}(\cdot)$  gets a vector and forms a diagonal matrix with its input on the diagonal elements. The notation  $\|\cdot\|_p$  refers to either matrix norm or vector norm depending on the context,  $\|\cdot\|_F$  shows the Frobenius norm, and  $|\cdot|$  indicates the absolute value or the cardinality of a set depending on the context. The symbol  $\mathbf{I}_m$  denotes the identity matrix of size  $m$  and the letter  $\mathcal{K}$  is used as a shorthand for the set  $\{1, \dots, p\}$ . The symbol  $\mathcal{S}_{n,m}$  as the set of real  $n \times m$  matrices with orthonormal columns, i.e.,  $\mathcal{S}_{n,m} \triangleq \{\mathbf{P} \in \mathbb{R}^{n \times m} \mid \mathbf{P}^\top \mathbf{P} = \mathbf{I}_m\}$ . The projection operator  $\text{proj}_{\mathcal{S}_{n,m}} : \mathbb{R}^{n \times m} \rightarrow \mathcal{S}_{n,m}$  is defined as  $\text{proj}_{\mathcal{S}_{n,m}} \mathbf{H} = \arg \min \{\|\mathbf{P} - \mathbf{H}\|_F \mid \mathbf{P} \in \mathcal{S}_{n,m}\}$ .

## II. PROBLEM FORMULATION

Optimization problems of the form (1a)–(1c) can be computationally challenging due to the non-convexities of the objective function and constraints. In order to derive convex relaxations, we first lift the problem into a higher dimensional space by introducing an auxiliary variable  $\mathbf{X} \in \mathbb{S}_n$ , accounting for the quadratic term  $\mathbf{P}\mathbf{P}^\top$ . For every  $k \in \{0\} \cup \mathcal{K}$ , define  $f_k : \mathbb{R}^{n \times m} \times \mathbb{S}_n \rightarrow \mathbb{R}$  as:

$$f_k(\mathbf{P}, \mathbf{X}) \triangleq \langle \mathbf{M}_k, \mathbf{X} \rangle + \langle \mathbf{N}_k, \mathbf{P} \rangle + q_k. \quad (2)$$

Using the auxiliary variable  $\mathbf{X}$ , the optimization problem (1a)–(1c) can be equivalently reformulated as

$$\underset{\substack{\mathbf{P} \in \mathbb{R}^{n \times m} \\ \mathbf{X} \in \mathbb{S}_n}}{\text{minimize}} \quad f_0(\mathbf{P}, \mathbf{X}) + g_0(\mathbf{P}) \quad (3a)$$

$$\text{subject to} \quad f_k(\mathbf{P}, \mathbf{X}) \leq 0 \quad k \in \mathcal{K}, \quad (3b)$$

$$\mathbf{P}^\top \mathbf{P} = \mathbf{I}_m, \quad (3c)$$

$$\mathbf{P} \mathbf{P}^\top = \mathbf{X}, \quad (3d)$$

with a convex objective function and convex linear inequality constraints (3b). The above formulation is still not convex due to the presence of the constraints (3c) and (3d) that capture all non-convexities of the problem.

## A. Convex Relaxation

In order to convexify the lifted problem (3a)–(3d), we relax the constraints (3c) and (3d) to

$$\mathbf{I}_m - \mathbf{P}^\top \mathbf{P} \in \mathcal{C} \quad \wedge \quad \mathbf{X} - \mathbf{P}\mathbf{P}^\top \in \mathcal{D} \quad \wedge \quad \text{tr}\{\mathbf{X}\} = m, \quad (4)$$

where  $\mathcal{C} \subseteq \mathbb{S}_m$  and  $\mathcal{D} \subseteq \mathbb{S}_n$  are convex cones to be defined. In this work, we consider the common-practice semidefinite programming (SDP) relaxation and introduce a novel convex relaxation that transforms the problem (3a)–(3d) into a convex quadratically-constrained quadratic program (QCQP).

1) *Semidefinite Programming Relaxation*: This relaxation provides a powerful method for tackling non-convex polynomial optimization problems [20]. The SDP relaxation of the problem (3a)–(3d) can be derived by having  $\mathcal{C} = \mathbb{S}_m^+$  and  $\mathcal{D} = \mathbb{S}_n^+$ . Despite the effectiveness of this relaxation in providing high-quality solutions, its applicability is limited to the problems of moderate size due to the computational cost of imposing high-dimensional conic constraints.

2) *Parabolic Relaxation*: We propose a computationally efficient convex relaxation as an alternative to the SDP relaxation. In order to formulate the proposed relaxation for the problem (3a)–(3d), we need to set  $\mathcal{C} = \mathcal{V}_m$  and  $\mathcal{D} = \mathcal{V}_n$ , where for every positive integer  $o$ , set  $\mathcal{V}_o \subseteq \mathbb{S}_o$  is defined as follows

$$\mathcal{V}_o \triangleq \{\mathbf{H} \in \mathbb{S}_o \mid H_{ii} + H_{jj} \geq 2|H_{ij}|, \forall i, j \in \{1, \dots, o\}\}.$$

**Remark 1.** It can be easily observed that if  $(\mathcal{C}, \mathcal{D}) = (\mathcal{V}_m, \mathcal{V}_n)$ , the constraints (3c) and (3d) are equivalent to the following convex quadratic inequalities:

$$\|\mathbf{P}(\dot{\mathbf{e}}_i - \dot{\mathbf{e}}_j)\|_2^2 \leq 2, \quad \forall i, j \in \{1, \dots, m\}, \quad (5a)$$

$$\|\mathbf{P}(\dot{\mathbf{e}}_i + \dot{\mathbf{e}}_j)\|_2^2 \leq 2, \quad \forall i, j \in \{1, \dots, m\}, \quad (5b)$$

$$\|\mathbf{P}^\top(\ddot{\mathbf{e}}_i - \ddot{\mathbf{e}}_j)\|_2^2 \leq X_{ii} + X_{jj} - 2X_{ij}, \quad \forall i, j \in \{1, \dots, n\}, \quad (5c)$$

$$\|\mathbf{P}^\top(\ddot{\mathbf{e}}_i + \ddot{\mathbf{e}}_j)\|_2^2 \leq X_{ii} + X_{jj} + 2X_{ij}, \quad \forall i, j \in \{1, \dots, n\}, \quad (5d)$$

$$\text{tr}\{\mathbf{X}\} = m. \quad (5e)$$

where  $\{\dot{\mathbf{e}}_k \in \mathbb{R}^m\}_{k=1}^m$  and  $\{\ddot{\mathbf{e}}_k \in \mathbb{R}^n\}_{k=1}^n$  denote the standard basis for  $\mathbb{R}^m$  and  $\mathbb{R}^n$ , respectively. Hence, the proposed relaxation reduces (3a)–(3d) to a convex QCQP.

Notice that either of the aforementioned relaxations may fail to produce a feasible point for (1a)–(1c), because in general, an optimal solution to a convex relaxation does not necessarily satisfy the constraints (3c) and (3d). In what follows, we propose a penalization technique that guarantees the recovery of feasible points for (1a)–(1c) under certain conditions.

## III. PENALIZATION

In this section, we show that by including a penalty term in the objective, one can obtain feasible points for the non-convex problem (3a)–(3d). Given an arbitrary initial point  $\tilde{\mathbf{P}} \in \mathcal{S}_{n,m}$ , that is not necessarily feasible, we transform the problem (3a)–(3d) into the following convex relaxation with

revised objective function:

$$\begin{aligned} & \underset{\substack{\mathbf{P} \in \mathbb{R}^{n \times m} \\ \mathbf{X} \in \mathbb{S}_n}}{\text{minimize}} & f_0(\mathbf{P}, \mathbf{X}) + g_0(\mathbf{P}) - \mu \langle \mathbf{P}, \check{\mathbf{P}} \rangle \end{aligned} \quad (6a)$$

$$\text{subject to} \quad f_k(\mathbf{P}, \mathbf{X}) \leq 0 \quad k \in \mathcal{K}, \quad (6b)$$

$$\mathbf{I}_m - \mathbf{P}^\top \mathbf{P} \in \mathcal{C}, \quad (6c)$$

$$\mathbf{X} - \mathbf{P} \mathbf{P}^\top \in \mathcal{D}, \quad (6d)$$

$$\text{tr}\{\mathbf{X}\} = m, \quad (6e)$$

where  $(\mathcal{C}, \mathcal{D}) \in \{(\mathbb{S}_m^+, \mathbb{S}_n^+), (\mathcal{V}_m, \mathcal{V}_n)\}$ , and the fixed parameter  $\mu > 0$  sets a trade-off between the original objective function and the linear penalty term  $\langle \mathbf{P}, \check{\mathbf{P}} \rangle$ .

**Remark 2.** If an optimal solution  $(\check{\mathbf{P}}, \check{\mathbf{X}})$  of the problem (6a)–(6e) satisfies the constraints (3c) and (3d), then  $\check{\mathbf{P}}$  is feasible for (1a)–(1c).

In the remainder of this section, certain conditions are introduced to guarantee that the penalized relaxation (6a)–(6e) produces feasible points for the non-convex problem (3a)–(3d).

**Definition 1.** Define feasibility distance  $d_{\mathcal{F}}: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$  as

$$d_{\mathcal{F}}(\mathbf{P}) \triangleq \inf\{\|\mathbf{C} - \mathbf{P}\|_{\mathbb{F}} \mid \mathbf{C} \in \mathcal{F}\}, \quad (7)$$

where  $\mathcal{F}$  denotes the feasible set of the problem (1a)–(1c).

**Definition 2.** Define the singularity function  $s: \mathcal{S}_{n,m} \rightarrow \mathbb{R}$  as:

$$s(\mathbf{P}) \triangleq \sup_{\mathbf{D} \in \mathcal{Z}_{\mathbf{P}}} \left\{ \min_{k \in \mathcal{K}} \{-\langle 2\mathbf{M}_k \mathbf{P} + \mathbf{N}_k, \mathbf{D} \rangle\} \right\}, \quad (8)$$

where  $\mathcal{Z}_{\mathbf{P}} \triangleq \{\mathbf{D} \in \mathbb{R}^{n \times m} \mid \mathbf{P}^\top \mathbf{D} = \mathbf{0} \wedge \|\mathbf{D}\|_{\mathbb{F}} \leq 1\}$ . A point  $\mathbf{P} \in \mathcal{S}_{n,m}$  is said to satisfy the Mangasarian-Fromovitz constraint qualification (MFCQ) condition if it is feasible for the problem (1a)–(1c) and  $s(\mathbf{P}) > 0$ .

**Theorem 1.** Define the constants

$$\beta \triangleq \max_{\mathbf{P} \in \mathcal{S}_{m,n}} \{ |g_0(\mathbf{P}) + \langle \mathbf{M}_0, \mathbf{P} \mathbf{P}^\top \rangle + \langle \mathbf{N}_0, \mathbf{P} \rangle| \}, \quad (9a)$$

$$\psi \triangleq 2\|\mathbf{M}_0\|_{\mathbb{F}} + \|\mathbf{N}_0\|_{\mathbb{F}} + \sum_{i=1}^w \|\mathbf{A}_i\|_{\mathbb{F}}, \quad (9b)$$

$$\kappa \triangleq 4 \max_{k \in \mathcal{K}} \{\|\mathbf{M}_k\|_{\mathbb{F}}\} + \max_{k \in \mathcal{K}} \{\|\mathbf{N}_k\|_{\mathbb{F}}\} \quad (9c)$$

and let  $\check{\mathbf{P}} \in \mathcal{F}$  be a feasible point for the problem (1a)–(1c) that satisfies the MFCQ condition. If

$$\mu > \max\{\beta^{-1}\psi^2, \beta(26\kappa)^2 s(\check{\mathbf{P}})^{-2}, 144\beta\}, \quad (10)$$

then the penalized relaxation (6a)–(6e) has a unique optimal solution  $(\check{\mathbf{P}}, \check{\mathbf{X}})$ , that satisfies (3c) and (3d). Moreover,  $\check{\mathbf{P}}$  is feasible for (1a)–(1c) and  $\bar{f}_0(\check{\mathbf{P}}) + g_0(\check{\mathbf{P}}) \leq \bar{f}_0(\check{\mathbf{P}}) + g_0(\check{\mathbf{P}})$ .

*Proof.* See Section V for the proof.  $\square$

**Remark 3.** For every point  $\mathbf{P} \in \mathcal{S}_{m,n}$ , it is straightforward

---

### Algorithm 1 Sequential Penalized Relaxation

---

**Input:**  $\check{\mathbf{P}} \in \mathcal{S}_{n,m}$ , a fixed parameter  $\mu > 0$ , and  $k = 0$ ,  
1: **repeat**  
2:    $k \leftarrow k + 1$   
3:    $\mathbf{P}^k \leftarrow$  solve (6a)–(6e) with the penalty  $\mu \langle \mathbf{P}, \check{\mathbf{P}} \rangle$   
4:    $\check{\mathbf{P}} \leftarrow \text{proj}_{\mathcal{S}_{n,m}} \mathbf{P}^k$   
5: **until** stopping criteria is met  
**Output:**  $\mathbf{P}^k$

---

to calculate  $s(\mathbf{P})$  by solving the following convex problem:

$$\begin{aligned} & \underset{t \in \mathbb{R}, \mathbf{D} \in \mathcal{Z}_{\mathbf{P}}}{\text{maximize}} & t \\ & \text{subject to} & t \leq -\langle 2\mathbf{M}_k \mathbf{P} + \mathbf{N}_k, \mathbf{D} \rangle, \quad k \in \mathcal{K}. \end{aligned}$$

Notice that  $\beta$  is upper bounded by  $\psi$  and it can be simply lower-bounded by any arbitrary member of the set  $\mathcal{S}_{m,n}$ . This certifies the existence of a bounded  $\mu$  that satisfies (10). In practice, there is no need to compute  $s(\mathbf{P})$  for fine-tuning parameter  $\mu$ , since (10) offers a conservative sufficient condition and usually, there exists a smaller  $\mu$  that satisfies Theorem 1. In Section IV, we assess the sensitivity of our approach with respect to different choices of  $\mu$ .

Theorem 1 is concerned with the case where the initial point  $\check{\mathbf{P}}$  is feasible for the original problem (1a)–(1c). However, finding a feasible starting point can be difficult due to the presence of the non-convex quadratic inequality constraints (1b). The next theorem states that even if  $\check{\mathbf{P}}$  is not feasible, the proposed penalized relaxation can still result in a feasible point for the non-convex problem (1a)–(1c).

**Theorem 2.** Consider an initial  $\check{\mathbf{P}} \in \mathcal{S}_{n,m}$  that satisfies

$$d_{\mathcal{F}}(\check{\mathbf{P}}) < 1, \quad (12a)$$

$$s(\check{\mathbf{P}}) > \kappa d_{\mathcal{F}}(\check{\mathbf{P}}) [1 + (1 - d_{\mathcal{F}}(\check{\mathbf{P}}))^{-1}], \quad (12b)$$

where  $\kappa$  is defined in (9c). If  $\mu$  is sufficiently large, then the penalized convex relaxation (6a)–(6e) has a unique optimal solution  $(\check{\mathbf{P}}, \check{\mathbf{X}})$  that satisfies (3c) and (3d). Moreover,  $\check{\mathbf{P}}$  is feasible for (1a)–(1c).

*Proof.* See Section V for the proof.  $\square$

#### A. Sequential Penalized Relaxation

Motivated by Theorems 1 and 2, this section presents a sequential approach that solves a sequence of penalized relaxations of the form (6a)–(6e) to infer high-quality feasible points for the non-convex problem (1a)–(1c). The proposed scheme starts from an initial point  $\check{\mathbf{P}}$  on the Stiefel manifold. In each round, the solution of the penalized relaxation (6a)–(6e) is projected onto the Stiefel manifold and then the projected point is employed as an initialization for the next round. Once a feasible point for (1a)–(1c) is obtained, according to Theorem 1, the proposed scheme preserves feasibility and generates a sequence of points whose objective values monotonically improves. The details of the sequential scheme are delineated in Algorithm 1.

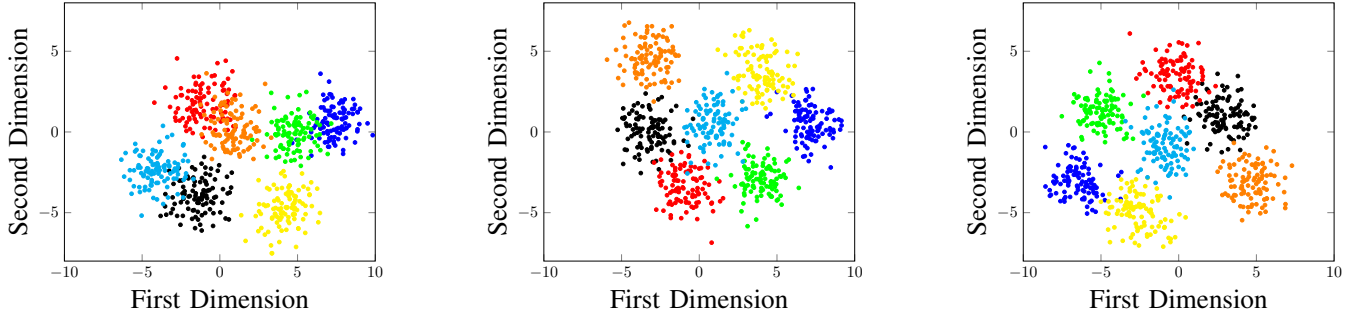


Fig. 1: Two dimensional representation on a training set from the synthetic data set. **Left:** MMDA [2], **middle:** SPR-S, **right:** SPR-Q. The results show that the SPR-S and SPR-Q algorithms have provided more discriminative 2D representations compared to the MMDA method.

The following theorem guarantees the convergence of Algorithm 1 to at least a locally optimal solution.

**Theorem 3.** Let  $\check{\mathcal{F}} \triangleq \{\mathbf{P} \in \mathcal{F} | \check{f}_0(\mathbf{P}) + g_0(\mathbf{P}) \leq h\}$  denote an epigraph of the problem (1a)–(1c) such that  $s(\mathbf{P}) > 0$  for every  $\mathbf{P} \in \check{\mathcal{F}}$ . If  $\check{\mathbf{P}} \in \check{\mathcal{F}}$  and

$$\mu > \max \left\{ \beta^{-1} \psi^2, \beta (26\kappa)^2 \min_{\mathbf{P} \in \check{\mathcal{F}}} \{s(\mathbf{P})\}^{-2}, 144\beta \right\}, \quad (13)$$

then the sequence generated by Algorithm (1) converges to a local minimizer of the problem (1a)–(1c).

#### IV. EXPERIMENTAL RESULTS

In this section, we conduct numerical experiments on real and synthetic datasets to verify the effectiveness of the proposed sequential approach, termed SPR, in solving non-convex optimization problems with orthogonality constraints. In Subsections IV-A and IV-B, we apply SPR on two practical problems involving orthogonality constraints. We use SPR-S and SPR-P to refer to the combination of Algorithm 1 with the SDP relaxation and the proposed parabolic relaxation, respectively. To solve the penalized relaxations in each round of the algorithm, we use MOSEK version 7.0. Through the experiments, we leverage the inherent sparsity patterns of the problems to reduce the computational cost of solving large-scale semidefinite programs. This enables us to break down large-scale conic constraints into a set of smaller ones [21]. Since finding a feasible point for (1a)–(1c) can be computationally demanding, we initialize Algorithm 1 with an arbitrary starting point on the Stiefel manifold and aim to improve the quality of the point. If the algorithm can recover a feasible point for (1a)–(1c), according to Theorem 1, it can generate a sequence of feasible points whose objective values monotonically improve. To measure the level of infeasibility, define  $\text{tr}\{\bar{\mathbf{X}} - \bar{\mathbf{P}}\bar{\mathbf{P}}^\top\}$  as the feasibility violation of an arbitrary feasible point  $(\bar{\mathbf{P}}, \bar{\mathbf{X}})$  of the problem (6a)–(6e). We terminate the sequential algorithm once the feasibility violation and objective value improvement are less than  $10^{-5}$  or if the round number exceeds 100. Notice that the Nesterov acceleration technique can be employed to improve the convergence behaviour of the SPR algorithm. However, in this case, the algorithm may fail to preserve the monotonically decreasing order of the objective values even if the initial point is feasible.

We apply the sequential algorithm on two fundamental machine learning problems of discriminative dimensionality reduction and graph clustering. Notice that each of these problems are well-studied in the literature and several approaches have been developed to efficiently target these applications. Therefore, it is not the intent of this work to compete with these state-of-the-art problem-specific approaches, but rather to demonstrate the potential of Algorithm 1 in solving the problems of form (1a)–(1c) that widely arise in different areas of machine learning.

##### A. Experiment I: Discriminative Dimensionality Reduction

Given a collection of high-dimensional data points from  $c$  different classes, the problem of discriminative dimensionality reduction aims to learn a low-dimensional subspace on which the projection of different classes are well-separated. To this end, [2] proposed a max-min distance analysis (MMDA) that maximizes the minimum distance between all class pairs. This problem can be cast as a non-convex and non-smooth optimization problem of form

$$\underset{\mathbf{P} \in \mathbb{R}^{n \times m}}{\text{maximize}} \quad \min_{1 \leq i < j \leq c} \langle \mathbf{A}^{ij}, \mathbf{P}\mathbf{P}^\top \rangle \quad (14a)$$

$$\text{subject to} \quad \mathbf{P}^\top \mathbf{P} = \mathbf{I}_m, \quad (14b)$$

where each  $\mathbf{A}^{ij} \in \mathbb{S}_n$  is a given weighted distance matrix between the  $i^{\text{th}}$  and  $j^{\text{th}}$  classes. In this experiment, we evaluate the performance of the SPR algorithm for solving the problems of form (14a)–(14b). Closely related to our work, [2] uses a sequence of local SDP relaxations to find the solution of problem (14a)–(14b). We benchmark the SPR method against the MMDA on both real and synthetic datasets. To ensure the comparison is fair, both methods use the same initial point and the same distance matrices  $\mathbf{A}^{ij}$  which are computed based on [2]. Other parameter settings of the MMDA are set to their default values. Following [2], we conduct 100 independent experiments on 10-dimensional synthetic data from seven classes. For each class  $i$ , a mean vector  $\boldsymbol{\eta}_i \in \mathbb{R}^{10}$  is sampled from 10-dimensional zero mean Gaussian distribution with co-variance matrix  $2\mathbf{I}_{10}$  and then a pair of training and testing sets, each with 100 members, is generated based on the Gaussian distribution  $\mathcal{N}(\boldsymbol{\eta}_i, \mathbf{I}_{10})$ .

To compare the classification error rate, we project each test set into subspaces with varying dimensions, learned on

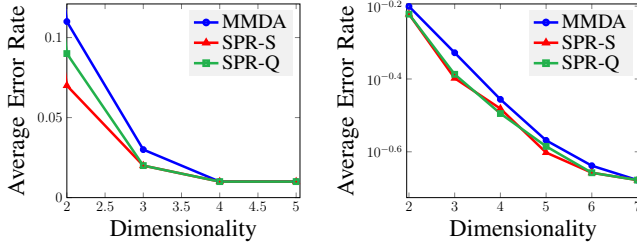


Fig. 2: Performance of SPR comparing to MMDA [2] on **left**: synthetic dataset, **right** YALE dataset [23]. *Best viewed in color.*

its corresponding training set. The projected instances are then classified using the nearest mean classifier. Figure 2 (left) shows the average classification error rate with respect to the reduced dimensionality on the synthetic datasets. To run the experiment on the synthetic datasets, we set  $\mu$  to 100 and 200 for SPR-S and SPR-P, respectively. Moreover, we conduct this experiment on the YALE dataset consisting of 165 frontal face images of 15 individuals under different illumination and lightening conditions [22]. Each image is of size  $32 \times 32$  pixels. The results of this experiment are illustrated in Figure 2 (right). According to Figure 2, SPR-S and SPR-P perform on par or better than the MMDA algorithm on both real and synthetic datasets in the problem of discriminative dimensionality reduction. In the experiment on the YALE dataset, we set  $\mu$  to 5000 and 10000 for SPR-S and SPR-P, respectively. To qualitatively compare the methods, Figure 1 visualizes the results of projecting a randomly chosen training set from the synthetic dataset on the 2D space. Observe that comparing to the MMDA method, the SPR-based algorithms learn more discriminative 2D representations that are suitable for classification tasks.

To assess the sensitivity of the SPR algorithm with respect to the parameter  $\mu$ , we perform the discriminative dimensionality reduction experiment with  $m = 2$  on YALE dataset and report the results in Figure 3 for various choices of  $\mu$ . Observe that the final solution obtained by the proposed algorithm is not very sensitive to the choice of  $\mu$ . According to the figure, the SPR-S requires smaller values of  $\mu$  to recover feasible points, e.g.  $\mu = 5000$ , while SPR-P fails to find feasible points for such choice of  $\mu$ . Moreover, it can be seen that if  $\mu$  exceeds a certain threshold, both SPR-S and SPR-P provide the same sequence of feasible points.

## B. Experiment II: Graph Clustering

Given a weighted graph  $\mathcal{G}$  with  $n$  vertices, the graph clustering problem aims to partition  $\mathcal{G}$  into a set of sub-graphs such that the vertices within each one are more densely connected to each other than those belonging to different sub-graphs. Inspired by the well-known spectral clustering technique [24], this experiment incorporates a set of non-negative constraints to formulate the graph clustering problem as the following

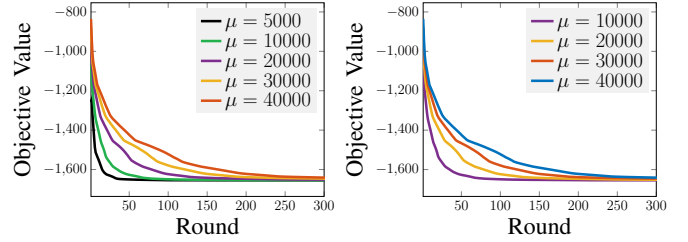


Fig. 3: Sensitivity analysis of SPR-S (**left**) and SPR-Q (**right**) with respect to different choices of parameter  $\mu$  for the discriminative dimensionality reduction problem, where  $m = 2$ . This experiment is performed on the YALE dataset. *Best viewed in color.*

Dataset	$n$	Dim.	$m$	ONGR	SPR-S	SPR-Q
Iris	150	4	3	79.84	<b>86.71</b>	81.23
Spiral	312	2	3	87.44	<b>95.76</b>	94.15
Jain	373	2	2	88.42	<b>92.33</b>	90.26
Compound	399	2	6	74.57	74.25	<b>76.48</b>
R15	600	2	15	86.07	85.36	<b>86.94</b>
Aggregation	788	2	7	<b>87.84</b>	86.39	84.66

TABLE I: Clustering performance (%) on the UCI datasets [26] and shape sets [27]–[29].

optimization [25]:

$$\underset{P \in \mathbb{R}^{n \times m}}{\text{minimize}} \quad \langle L, PP^\top \rangle \quad (15a)$$

$$\text{subject to} \quad P^\top P = I_m, \quad (15b)$$

$$P \geq 0, \quad (15c)$$

where  $L$  denotes the Laplacian matrix of the weighted graph  $\mathcal{G}$  and  $\geq$  is the element-wise inequality operator. Comparing to the spectral clustering, formulation (15a)–(15c) offers a more interpretable clustering framework which requires no further post-processing steps to identify the cluster members. Given  $\hat{P}$ , the optimal solution of the above problem, each vertex  $i$  is assigned to a cluster with label  $\text{argmax}_j \hat{P}_{ij}$ . [25] proposed a fast and scalable heuristic, denoted by ONGR, to solve large-scale instances of the form (15a)–(15c). Due to the fact that this problem is a special case of (1a)–(1c), we apply the SPR algorithm to find the solution of (15a)–(15c) and use the same procedure as [25] to create the Laplacian matrix  $L$ . To make a fair comparison between the ONGR and SPR, we use the same initialization for both methods. Table I reports the clustering performance of the SPR against [25] on well-known datasets from the UCI machine learning repository [26] and shape sets [27], [28]. For each dataset,  $n$ , Dim, and  $m$  refer to the number of sample points, dimension of each point, and the number of classes, respectively. The scores for each method is computed by averaging over 30 independent runs for each dataset. As the results indicate, SPR-S and SPR-P exhibit better performance compared to [25] on most of the datasets. Through this experiment, we set  $\mu = 1000$  in the SPR algorithm and use the default parameter settings for the ONGR algorithm.

## V. PROOFS

This section presents the proof of Theorems 1, 2, and 3. Before proceeding with the proofs, we provide some prerequisite lemmas.

Using the well-known epigraph technique [20], the non-smooth term  $g_0(\mathbf{P})$  in (3a) can be removed by adding a pair of linear constraints and incorporating an additional term into the objective function. This reformulation of (3a)–(3d) leads to the following penalized non-convex problem:

$$\underset{\substack{\mathbf{P} \in \mathbb{R}^{n \times m} \\ \mathbf{t} \in \mathbb{R}^w}}{\text{minimize}} \quad \mathbf{1}^\top \mathbf{t} + \langle \mathbf{M}_0, \mathbf{P}\mathbf{P}^\top \rangle + \langle \mathbf{N}_0, \mathbf{P} \rangle - \mu \langle \check{\mathbf{P}}, \mathbf{P} \rangle \quad (16a)$$

$$\text{subject to} \quad \bar{\gamma} : +\alpha(\mathbf{P}) + \mathbf{b} \leq \mathbf{t}, \quad (16b)$$

$$\gamma : -\alpha(\mathbf{P}) - \mathbf{b} \leq \mathbf{t}, \quad (16c)$$

$$\lambda : \langle \mathbf{M}_k, \mathbf{P}\mathbf{P}^\top \rangle + \langle \mathbf{N}_k, \mathbf{P} \rangle + q_k \leq 0, \quad k \in \mathcal{K}, \quad (16d)$$

$$\Omega : \mathbf{P}^\top \mathbf{P} = \mathbf{I}_m, \quad (16e)$$

with  $\bar{\gamma} \in \mathbb{R}^w$ ,  $\gamma \in \mathbb{R}^w$ ,  $\lambda \in \mathbb{R}^{|\mathcal{K}|}$ , and  $\Omega \in \mathbb{S}_m$  as the dual variables associated with the constraints (16b), (16c), (16d), and (16e), respectively. Observe that the problems (16a)–(16e) and (1a)–(1c) are equivalent, if  $\mu = 0$ . In what follows, we show that under certain conditions, the optimal solution of (16a)–(16e) can be obtained in polynomial time via convex relaxation.

**Lemma 1.** *Consider an arbitrary point  $\check{\mathbf{P}} \in \mathbb{R}^{n \times m}$ . Every optimal solution  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  of the problem (16a)–(16e) satisfies,*

$$0 \leq \|\check{\mathbf{P}} - \check{\mathbf{P}}\|_{\text{F}} - d_{\mathcal{F}}(\check{\mathbf{P}}) \leq 2\sqrt{\beta\mu^{-1}} \quad (17)$$

where  $\beta$  is defined in (9a).

*Proof.* According to Definition 1, the distance between an arbitrary point  $\check{\mathbf{P}}$  and any points in  $\mathcal{F}$  is greater than or equal to  $d_{\mathcal{F}}(\check{\mathbf{P}})$ . This implies that  $\|\check{\mathbf{P}} - \check{\mathbf{P}}\|_{\text{F}} - d_{\mathcal{F}}(\check{\mathbf{P}})$  is lower bounded by zero. To prove the validity of the upper bound, let  $\bar{\mathbf{P}}$  be an arbitrary member of  $\{\mathbf{P} \in \mathcal{F} \mid \|\mathbf{P} - \check{\mathbf{P}}\|_{\text{F}} = d_{\mathcal{F}}(\check{\mathbf{P}})\}$ . Since  $\check{\mathbf{P}}$  is the minimizer of the optimization problem (16a)–(16e), one can write:

$$\begin{aligned} & \|\alpha(\check{\mathbf{P}}) + \mathbf{b}\|_1 + \langle \mathbf{M}_0, \check{\mathbf{P}}\check{\mathbf{P}}^\top \rangle + \langle \mathbf{N}_0, \check{\mathbf{P}} \rangle - \mu \langle \check{\mathbf{P}}, \check{\mathbf{P}} \rangle \\ & \leq \|\alpha(\bar{\mathbf{P}}) + \mathbf{b}\|_1 + \langle \mathbf{M}_0, \bar{\mathbf{P}}\bar{\mathbf{P}}^\top \rangle + \langle \mathbf{N}_0, \bar{\mathbf{P}} \rangle - \mu \langle \bar{\mathbf{P}}, \bar{\mathbf{P}} \rangle. \end{aligned}$$

and due to feasibility of  $\check{\mathbf{P}}$  and  $\bar{\mathbf{P}}$  we have:

$$\begin{aligned} & \frac{\mu}{2} \|\check{\mathbf{P}} - \bar{\mathbf{P}}\|_{\text{F}}^2 - \beta \\ & \leq \|\alpha(\check{\mathbf{P}}) + \mathbf{b}\|_1 + \langle \mathbf{M}_0, \check{\mathbf{P}}\check{\mathbf{P}}^\top \rangle + \langle \mathbf{N}_0, \check{\mathbf{P}} \rangle + \frac{\mu}{2} \|\check{\mathbf{P}} - \bar{\mathbf{P}}\|_{\text{F}}^2 \\ & \leq \|\alpha(\bar{\mathbf{P}}) + \mathbf{b}\|_1 + \langle \mathbf{M}_0, \bar{\mathbf{P}}\bar{\mathbf{P}}^\top \rangle + \langle \mathbf{N}_0, \bar{\mathbf{P}} \rangle + \frac{\mu}{2} \|\bar{\mathbf{P}} - \check{\mathbf{P}}\|_{\text{F}}^2 \\ & \leq \frac{\mu}{2} \|\bar{\mathbf{P}} - \check{\mathbf{P}}\|_{\text{F}}^2 + \beta = \frac{\mu}{2} \times d_{\mathcal{F}}(\check{\mathbf{P}})^2 + \beta. \end{aligned}$$

where  $\beta$  is defined in (9a). Therefore,

$$\|\check{\mathbf{P}} - \bar{\mathbf{P}}\|_{\text{F}} - d_{\mathcal{F}}(\check{\mathbf{P}}) \leq \sqrt{\|\check{\mathbf{P}} - \bar{\mathbf{P}}\|_{\text{F}}^2 - d_{\mathcal{F}}(\check{\mathbf{P}})^2} \leq 2\sqrt{\beta\mu^{-1}}$$

which proves the right side of (17).  $\square$

Based on Lemma 1, the next lemma guarantees that if the initial point  $\check{\mathbf{P}}$  satisfies the MFCQ regularity condition and if it is close to the feasible set, then under some assumptions,

MFCQ is satisfied by every optimal point  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  of the problem (16a)–(16e) as well.

**Lemma 2.** *Consider an arbitrary  $\check{\mathbf{P}} \in \mathcal{S}_{n,m}$ . Every optimal solution  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  of the problem (16a)–(16e) satisfies*

$$s(\check{\mathbf{P}}) \geq s(\check{\mathbf{P}}) - \kappa(d_{\mathcal{F}}(\check{\mathbf{P}}) + 2\sqrt{\beta\mu^{-1}}) \quad (19)$$

where  $\beta$  and  $\kappa$  are defined in (9a) and (9c).

*Proof.* Due to compactness of the set  $\mathcal{Z}_{\check{\mathbf{P}}}$ , the supremum in (8) is attainable. As a result, there exists  $\check{\mathbf{D}} \in \mathcal{Z}_{\check{\mathbf{P}}}$  such that:

$$s(\check{\mathbf{P}}) = \min_{\check{\mathbf{D}} \in \mathcal{K}} \{-\langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \check{\mathbf{D}} \rangle\}. \quad (20)$$

and hence,

$$s(\check{\mathbf{P}}) \leq -\langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \check{\mathbf{D}} \rangle \quad \forall k \in \mathcal{K}. \quad (21)$$

On the other hand, we have  $\mathbf{D}' \in \mathcal{Z}_{\check{\mathbf{P}}}$ , where

$$\mathbf{D}' \triangleq (\mathbf{I}_n - \check{\mathbf{P}}\check{\mathbf{P}}^\top)\check{\mathbf{D}}. \quad (22)$$

As a result, we have

$$\begin{aligned} & -\langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \mathbf{D}' \rangle \\ & \geq s(\check{\mathbf{P}}) + \langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \check{\mathbf{D}} \rangle - \langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \mathbf{D}' \rangle \\ & = s(\check{\mathbf{P}}) + \langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \mathbf{D}' + \check{\mathbf{P}}\check{\mathbf{P}}^\top \check{\mathbf{D}} \rangle - \langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \mathbf{D}' \rangle \\ & = s(\check{\mathbf{P}}) + \langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \mathbf{D}' + \check{\mathbf{P}}(\check{\mathbf{P}} - \check{\mathbf{P}})^\top \check{\mathbf{D}} \rangle - \langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \mathbf{D}' \rangle \\ & = s(\check{\mathbf{P}}) - \langle 2\mathbf{M}_k(\check{\mathbf{P}} - \check{\mathbf{P}}), \mathbf{D}' \rangle + \langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \check{\mathbf{P}}(\check{\mathbf{P}} - \check{\mathbf{P}})^\top \check{\mathbf{D}} \rangle \\ & = s(\check{\mathbf{P}}) - \langle 2\mathbf{M}_k^\top \mathbf{D}' - \check{\mathbf{D}}(2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k)^\top \check{\mathbf{P}}, \check{\mathbf{P}} - \check{\mathbf{P}} \rangle \\ & \geq s(\check{\mathbf{P}}) - \|2\mathbf{M}_k^\top \mathbf{D}' - \check{\mathbf{D}}(2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k)^\top \check{\mathbf{P}}\|_{\text{F}} \|\check{\mathbf{P}} - \check{\mathbf{P}}\|_{\text{F}} \\ & = s(\check{\mathbf{P}}) - (4\|\mathbf{M}_k\|_{\text{F}} + \|\mathbf{N}_k\|_{\text{F}}) \|\check{\mathbf{P}} - \check{\mathbf{P}}\|_{\text{F}} \\ & \geq s(\check{\mathbf{P}}) - \kappa \|\check{\mathbf{P}} - \check{\mathbf{P}}\|_{\text{F}}. \end{aligned}$$

Now, according to Lemma 1, we can write

$$\begin{aligned} -\langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \mathbf{D}' \rangle & \geq s(\check{\mathbf{P}}) - \kappa \|\check{\mathbf{P}} - \check{\mathbf{P}}\|_{\text{F}} \\ & \geq s(\check{\mathbf{P}}) - \kappa(d_{\mathcal{F}}(\check{\mathbf{P}}) + 2\sqrt{\beta\mu^{-1}}). \end{aligned}$$

Additionally, based on Definition 2, for every  $k \in \mathcal{K}$  we have:

$$s(\check{\mathbf{P}}) \geq -\langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \mathbf{D}' \rangle,$$

which together with the above relations conclude (19).  $\square$

The next lemma guarantees the existence of Lagrange multipliers corresponding to optimal solutions of the problem (16a)–(16e).

**Lemma 3.** *Consider an arbitrary  $\check{\mathbf{P}} \in \mathcal{S}_{n,m}$  that satisfies*

$$s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}}) > 0. \quad (25)$$

*If the following inequality holds true,*

$$\mu > 4\beta[\kappa^{-1}s(\check{\mathbf{P}}) - d_{\mathcal{F}}(\check{\mathbf{P}})]^{-2}, \quad (26)$$

*then for every primal optimal pair  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  of (16a)–(16e), there exists Lagrange multipliers  $(\check{\gamma}, \check{\gamma}, \check{\lambda}, \check{\Omega}) \in \mathbb{R}^w \times \mathbb{R}^w \times \mathbb{R}^{|\mathcal{K}|} \times \mathbb{S}_m$  that satisfy the following Karush–Kuhn–Tucker (KKT) condi-*

tions

$$\nabla_{\mathbf{P}} \mathcal{L}(\check{\mathbf{P}}, \check{\mathbf{t}}, \check{\boldsymbol{\gamma}}, \check{\boldsymbol{\lambda}}, \check{\boldsymbol{\Omega}}) = \mu \check{\mathbf{P}}, \quad (27a)$$

$$\mathbf{1} + \check{\boldsymbol{\gamma}} + \check{\boldsymbol{\gamma}} = \mathbf{0}, \quad (27b)$$

$$\check{\boldsymbol{\gamma}} \circ (\alpha(\check{\mathbf{P}}) + \mathbf{b} - \check{\mathbf{t}}) = \mathbf{0}, \quad (27c)$$

$$\check{\boldsymbol{\gamma}} \circ (-\alpha(\check{\mathbf{P}}) - \mathbf{b} - \check{\mathbf{t}}) = \mathbf{0}, \quad (27d)$$

$$\check{\lambda}_k (\langle \mathbf{M}_k, \check{\mathbf{P}} \check{\mathbf{P}}^\top \rangle + \langle \mathbf{N}_k, \check{\mathbf{P}} \rangle + q_k) = 0 \quad k \in \mathcal{K}, \quad (27e)$$

$$\check{\boldsymbol{\gamma}} \leq \mathbf{0}, \quad \check{\boldsymbol{\gamma}} \leq \mathbf{0}, \quad \check{\boldsymbol{\lambda}} \leq \mathbf{0}, \quad (27f)$$

where  $\mathcal{L}(\mathbf{P}, \mathbf{t}, \boldsymbol{\gamma}, \boldsymbol{\lambda}, \boldsymbol{\Omega})$  represents the Lagrangian function of (16a)–(16e), defined as

$$\begin{aligned} \mathcal{L}(\mathbf{P}, \mathbf{t}, \boldsymbol{\gamma}, \boldsymbol{\lambda}, \boldsymbol{\Omega}) &\triangleq \mathbf{1}^\top \mathbf{t} + \langle \mathbf{M}_0, \mathbf{P} \mathbf{P}^\top \rangle + \langle \mathbf{N}_0, \mathbf{P} \rangle \\ &- \bar{\boldsymbol{\gamma}}^\top (\alpha(\mathbf{P}) + \mathbf{b} - \mathbf{t}) + \boldsymbol{\gamma}^\top (\alpha(\mathbf{P}) + \mathbf{b} + \mathbf{t}) \\ &- \sum_{k \in \mathcal{K}} \lambda_k (\langle \mathbf{M}_k, \mathbf{P} \mathbf{P}^\top \rangle + \langle \mathbf{N}_k, \mathbf{P} \rangle + q_k) - \langle \boldsymbol{\Omega}, \mathbf{P}^\top \mathbf{P} - \mathbf{I}_m \rangle. \end{aligned} \quad (28)$$

and,  $\beta$  and  $\kappa$  are defined in (9a) and (9c).

*Proof.* According to Lemma 2, and due to the assumptions (25) and (26) we have  $s(\check{\mathbf{P}}) > 0$ . Therefore, according to Definition 2, there exist  $\check{\mathbf{D}} \in \mathbb{R}^{n \times m}$  such that the following Mangasarian-Fromovitz constraint qualification conditions are satisfied:

$$+\alpha(\check{\mathbf{D}}) - \check{\mathbf{d}} < \mathbf{0} \quad (29a)$$

$$-\alpha(\check{\mathbf{D}}) - \check{\mathbf{d}} < \mathbf{0} \quad (29b)$$

$$\langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \check{\mathbf{D}} \rangle < 0, \quad k \in \mathcal{K} \quad (29c)$$

$$\check{\mathbf{P}}^\top \check{\mathbf{D}} = \mathbf{0}, \quad (29d)$$

where  $\check{\mathbf{d}} = 2|\alpha(\check{\mathbf{D}})|$ . According to (29a)–(29d) and due to the optimality of  $\check{\mathbf{P}}$ , we can conclude that there exists a dual point  $(\check{\boldsymbol{\gamma}}, \check{\boldsymbol{\gamma}}, \check{\boldsymbol{\lambda}}, \check{\boldsymbol{\Omega}}) \in \mathbb{R}^w \times \mathbb{R}^w \times \mathbb{R}^p \times \mathbb{S}_m$  for which the KKT conditions are satisfied.  $\square$

**Lemma 4.** Let  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  be a primal optimal solution of (16a)–(16e) with the corresponding Lagrange multipliers  $(\check{\boldsymbol{\gamma}}, \check{\boldsymbol{\gamma}}, \check{\boldsymbol{\lambda}}, \check{\boldsymbol{\Omega}})$  that satisfy the KKT conditions (27a)–(27e). The point  $(\check{\mathbf{P}}, \check{\mathbf{P}} \check{\mathbf{P}}^\top)$  is the unique primal solution of the penalized convex relaxation (6a)–(6e), if

$$-\mathbf{M}_0 + \sum_{k \in \mathcal{K}} \check{\lambda}_k \mathbf{M}_k + \theta \mathbf{I}_n \prec_{\check{\mathcal{D}}}^* \mathbf{0}, \quad (30a)$$

$$\check{\boldsymbol{\Omega}} - \theta \mathbf{I}_m \prec_{\check{\mathcal{C}}}^* \mathbf{0}, \quad (30b)$$

where  $\check{\mathcal{C}}$  and  $\check{\mathcal{D}}$  denote the dual cones of  $\mathcal{C}$  and  $\mathcal{D}$ , respectively.

*Proof.* Consider the following equivalent formulation of the

penalized convex relaxation (6a)–(6e)

$$\begin{aligned} \text{minimize} \quad & \mathbf{1}^\top \mathbf{t} + \langle \mathbf{M}_0, \mathbf{X} \rangle + \langle \mathbf{N}_0, \mathbf{P} \rangle - \mu \langle \check{\mathbf{P}}, \mathbf{P} \rangle \quad (31a) \\ & \mathbf{P} \in \mathbb{R}^{n \times m} \\ & \mathbf{t} \in \mathbb{R}^w \times \mathbb{R}^1 \\ & \mathbf{X} \in \mathbb{S}_n \end{aligned}$$

$$\text{subject to} \quad \bar{\boldsymbol{\gamma}} : +\alpha(\mathbf{P}) + \mathbf{b} \leq \mathbf{t}, \quad (31b)$$

$$\boldsymbol{\gamma} : -\alpha(\mathbf{P}) - \mathbf{b} \leq \mathbf{t}, \quad (31c)$$

$$\boldsymbol{\lambda} : \langle \mathbf{M}_k, \mathbf{X} \rangle + \langle \mathbf{N}_k, \mathbf{P} \rangle + q_k \leq 0, \quad k \in \mathcal{K}, \quad (31d)$$

$$\boldsymbol{\Psi} : \mathbf{I}_m - \mathbf{P}^\top \mathbf{P} \preceq_{\mathcal{C}} \mathbf{0}, \quad (31e)$$

$$\boldsymbol{\Lambda} : \mathbf{X} - \mathbf{P} \mathbf{P}^\top \preceq_{\mathcal{D}} \mathbf{0}, \quad (31f)$$

$$\theta : \text{tr}\{\mathbf{X}\} = m, \quad (31g)$$

with  $\bar{\boldsymbol{\gamma}} \in \mathbb{R}^w$ ,  $\boldsymbol{\gamma} \in \mathbb{R}^w$ ,  $\boldsymbol{\lambda} \triangleq [\lambda_1, \dots, \lambda_k]^\top \in \mathbb{R}^{|\mathcal{K}|}$ ,  $\boldsymbol{\Psi} \in \mathbb{S}_m$ ,  $\boldsymbol{\Lambda} \in \mathbb{S}_n$ , and  $\theta \in \mathbb{R}$  as the dual variables associated with the constraints (31b), (31c), (31d), (31e), (31f), (31g), respectively.

Let  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  and  $(\check{\boldsymbol{\gamma}}, \check{\boldsymbol{\gamma}}, \check{\boldsymbol{\lambda}}, \check{\boldsymbol{\Omega}})$  denote a pair of primal and dual optimal solutions for the nonconvex problem (16a)–(16e) and define

$$\check{\boldsymbol{\Lambda}} \triangleq -\mathbf{M}_0 + \sum_{k \in \mathcal{K}} \check{\lambda}_k \mathbf{M}_k.$$

In order to prove the lemma, we show that the following pair satisfies the KKT condition for (31a)–(31g):

$$(\mathbf{P}, \mathbf{X}, \mathbf{t}) = (\check{\mathbf{P}}, \check{\mathbf{P}} \check{\mathbf{P}}^\top, \check{\mathbf{t}}), \quad (32a)$$

$$(\bar{\boldsymbol{\gamma}}, \boldsymbol{\gamma}, \boldsymbol{\lambda}, \boldsymbol{\Psi}, \boldsymbol{\Lambda}, \theta) = (\check{\boldsymbol{\gamma}}, \check{\boldsymbol{\gamma}}, \check{\boldsymbol{\lambda}}, \check{\boldsymbol{\Omega}} - \theta \mathbf{I}_m, \check{\boldsymbol{\Lambda}} + \theta \mathbf{I}_m, \theta). \quad (32b)$$

These conditions can be formulated as:

- Stationarity with respect to  $\mathbf{X}$ ,  $\mathbf{P}$  and  $\mathbf{t}$ , respectively:

$$\boldsymbol{\Lambda} + \mathbf{M}_0 - \sum_{k \in \mathcal{K}} \lambda_k \mathbf{M}_k - \theta \mathbf{I}_n = \mathbf{0}, \quad (33a)$$

$$2\mathbf{P} \boldsymbol{\Psi} + 2\boldsymbol{\Lambda} \mathbf{P} - \mathbf{N}_0 + \sum_{k \in \mathcal{K}} \lambda_k \mathbf{N}_k + \sum_{i=1}^w (\bar{\gamma}_i - \gamma_i) \mathbf{A}_i + \mu \check{\mathbf{P}} = \mathbf{0}, \quad (33b)$$

$$\mathbf{1} + \bar{\boldsymbol{\gamma}} + \boldsymbol{\gamma} = \mathbf{0}, \quad (33c)$$

- Complementary slackness:

$$\bar{\boldsymbol{\gamma}} \circ (\alpha(\mathbf{P}) + \mathbf{b} - \mathbf{t}) = \mathbf{0}, \quad (33d)$$

$$\boldsymbol{\gamma} \circ (-\alpha(\mathbf{P}) - \mathbf{b} - \mathbf{t}) = \mathbf{0}, \quad (33e)$$

$$\lambda_k (\langle \mathbf{M}_k, \mathbf{X} \rangle + \langle \mathbf{N}_k, \mathbf{P} \rangle + q_k) = \mathbf{0}, \quad k \in \mathcal{K}, \quad (33f)$$

$$\boldsymbol{\Psi} (\mathbf{P}^\top \mathbf{P} - \mathbf{I}_m) = \mathbf{0}, \quad (33g)$$

$$\boldsymbol{\Lambda} (\mathbf{P} \mathbf{P}^\top - \mathbf{X}) = \mathbf{0}, \quad (33h)$$

- Dual feasibility:

$$\bar{\boldsymbol{\gamma}} \leq \mathbf{0}, \quad \boldsymbol{\gamma} \leq \mathbf{0}, \quad \boldsymbol{\lambda} \leq \mathbf{0}, \quad (33i)$$

$$\boldsymbol{\Psi} \prec_{\check{\mathcal{C}}}^* \mathbf{0}, \quad \boldsymbol{\Lambda} \prec_{\check{\mathcal{D}}}^* \mathbf{0}. \quad (33j)$$

The equations (33a)–(33c) are followed from the definition of  $\check{\boldsymbol{\Lambda}}$  and  $\check{\boldsymbol{\Psi}}$ , along with the stationarity conditions (27a)–(27b). Complementary slackness conditions are immediate consequences of primal feasibility. Finally, the dual feasibility conditions are resulted from (27f) and the assumptions (30a) and (30b).  $\square$

The next lemma provides an upper bound on the Lagrange multipliers of the problem (16a)–(16e), that will be used to show that this problem can be relaxed to (6a)–(6e) with no effect on the solution.

**Lemma 5.** Consider an arbitrary  $\check{\mathbf{P}} \in \mathcal{S}_{n,m}$  that satisfies (25) and let  $\mu$  satisfy (26). For every solution  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  of (16a)–

(16e), there exist Lagrange multipliers  $(\check{\gamma}, \check{\gamma}, \check{\lambda}, \check{\Omega}) \in \mathbb{R}^w \times \mathbb{R}^w \times \mathbb{R}^p \times \mathbb{S}_m$  that satisfy the KKT conditions (27a)–(27f) as well as the inequalities:

$$-\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \leq \frac{d_{\mathcal{F}}(\check{\mathbf{P}}) + \mu^{-1}\psi + 2\sqrt{\beta\mu^{-1}}}{s(\check{\mathbf{P}}) - \kappa(d_{\mathcal{F}}(\check{\mathbf{P}}) + 2\sqrt{\beta\mu^{-1}})} \quad (34a)$$

$$\left\| \frac{2}{\mu} \check{\Omega} + \mathbf{I}_m \right\|_{\mathbb{F}} \leq \kappa_2 \left( -\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \right) + d_{\mathcal{F}}(\check{\mathbf{P}}) + \mu^{-1}\psi + 2\sqrt{\beta\mu^{-1}} \quad (34b)$$

where constant  $\kappa_2$  is given by

$$\kappa_2 \triangleq 2 \max_{k \in \mathcal{K}} \{ \| \mathbf{M}_k \|_{\mathbb{F}} \} + \max_{k \in \mathcal{K}} \{ \| \mathbf{N}_k \|_{\mathbb{F}} \}, \quad (35)$$

and  $\beta$ ,  $\psi$  and  $\kappa$  are defined in (9a)–(9c).

*Proof.* According to Lemma 2, and due to the assumptions (25) and (26), we have  $s(\check{\mathbf{P}}) > 0$ . Hence, there exists  $\check{\mathbf{D}} \in \mathcal{Z}_{\check{\mathbf{P}}}$  such that:

$$-\langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \check{\mathbf{D}} \rangle > 0, \quad k \in \mathcal{K}. \quad (36)$$

According to Lemma 3 and due to optimality of  $\check{\mathbf{D}}$ , the condition (27a) yields

$$\begin{aligned} \mathbf{0} &= \langle \nabla_{\mathcal{P}} \mathcal{L}(\check{\mathbf{P}}, \check{\mathbf{t}}, \check{\gamma}, \check{\gamma}, \check{\lambda}, \check{\Omega}), \check{\mathbf{D}} \rangle \\ &= \langle 2\mathbf{M}_0 \check{\mathbf{P}} + \mathbf{N}_0, \check{\mathbf{D}} \rangle - \sum_{k \in \mathcal{K}} \check{\lambda}_k \langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \check{\mathbf{D}} \rangle + \langle 2\check{\mathbf{P}} \check{\Omega}, \check{\mathbf{D}} \rangle \\ &\quad - (\check{\gamma} - \check{\gamma})^\top \alpha(\check{\mathbf{D}}) - \mu \langle \check{\mathbf{P}}, \check{\mathbf{D}} \rangle. \end{aligned}$$

The above equation together with Definition 2 and the fact that  $\check{\mathbf{D}}^\top \check{\mathbf{P}} = \mathbf{0}$  give rise to the following inequality

$$\begin{aligned} -(\mathbf{1}^\top \check{\lambda}) s(\check{\mathbf{P}}) &\leq \sum_{k \in \mathcal{K}} \check{\lambda}_k \langle 2\mathbf{M}_k \check{\mathbf{P}} + \mathbf{N}_k, \check{\mathbf{D}} \rangle \\ &= \langle 2\mathbf{M}_0 \check{\mathbf{P}} + \mathbf{N}_0, \check{\mathbf{D}} \rangle - (\check{\gamma} - \check{\gamma})^\top \alpha(\check{\mathbf{D}}) + \mu \langle \check{\mathbf{P}} - \check{\mathbf{P}}, \check{\mathbf{D}} \rangle \\ &\leq 2 \| \mathbf{M}_0 \|_{\mathbb{F}} + \| \mathbf{N}_0 \|_{\mathbb{F}} + \sum_{i=1}^w \| \mathbf{A}_i \|_{\mathbb{F}} + \mu \| \check{\mathbf{P}} - \check{\mathbf{P}} \|_{\mathbb{F}}, \end{aligned}$$

which leads to

$$-\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \leq \frac{\mu^{-1}(2 \| \mathbf{M}_0 \|_{\mathbb{F}} + \| \mathbf{N}_0 \|_{\mathbb{F}} + \sum_{i=1}^w \| \mathbf{A}_i \|_{\mathbb{F}}) + \| \check{\mathbf{P}} - \check{\mathbf{P}} \|_{\mathbb{F}}}{s(\check{\mathbf{P}})}.$$

The above inequality along with Lemmas 1 and 2, concludes (34a).

Now, (34b) can be proven by pre-multiplying  $\check{\mathbf{P}}^\top$  to the KKT stationarity condition (27a) resulting in

$$\begin{aligned} (2\check{\Omega} + \mu \mathbf{I}_m) &\triangleq 2\check{\mathbf{P}}^\top \mathbf{M}_0 \check{\mathbf{P}} + \check{\mathbf{P}}^\top \mathbf{N}_0 - \sum_{i=1}^w (\check{\gamma}_i - \check{\gamma}_i) \check{\mathbf{P}}^\top \mathbf{A}_i \\ &\quad - \sum_{k \in \mathcal{K}} \check{\lambda}_k (2\check{\mathbf{P}}^\top \mathbf{M}_k \check{\mathbf{P}} + \check{\mathbf{P}}^\top \mathbf{N}_k) - \mu \check{\mathbf{P}}^\top (\check{\mathbf{P}} - \check{\mathbf{P}}). \end{aligned}$$

Hence, according to (27b) and (27f), we have

$$\begin{aligned} \| 2\check{\Omega} + \mu \mathbf{I}_m \|_{\mathbb{F}} &\leq 2 \| \mathbf{M}_0 \|_{\mathbb{F}} + \| \mathbf{N}_0 \|_{\mathbb{F}} + \sum_{i=1}^w |\check{\gamma}_i - \check{\gamma}_i| \times \| \mathbf{A}_i \|_{\mathbb{F}} \\ &\quad - \sum_{k \in \mathcal{K}} \check{\lambda}_k (2 \| \mathbf{M}_k \|_{\mathbb{F}} + \| \mathbf{N}_k \|_{\mathbb{F}}) + \mu \| \check{\mathbf{P}} - \check{\mathbf{P}} \|_{\mathbb{F}} \\ &\leq \psi - \kappa_2 (\mathbf{1}^\top \check{\lambda}) + \mu d_{\mathcal{F}}(\check{\mathbf{P}}) + 2\sqrt{\beta\mu}, \end{aligned}$$

which concludes (34b).  $\square$

Using Lemma 5, the next lemma offers conditions to guarantee that penalized relaxations give feasible points for (16a)–(16e).

**Lemma 6.** Consider an initial point  $\check{\mathbf{P}} \in \mathcal{S}_{n,m}$  and  $\mu > 0$ . Let  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  be a primal optimal solution of (16a)–(16e) with the corresponding Lagrange multipliers  $(\check{\gamma}, \check{\gamma}, \check{\lambda}, \check{\Omega})$  that satisfy the KKT conditions (27a)–(27e). Define

$$\varepsilon \triangleq \frac{1}{4} \left( 1 - d_{\mathcal{F}}(\check{\mathbf{P}}) - \frac{\kappa d_{\mathcal{F}}(\check{\mathbf{P}})}{s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}})} \right). \quad (38)$$

If the following inequalities hold true

$$2\mu^{-1} \| \mathbf{M}_0 \| \leq \varepsilon, \quad (39a)$$

$$-\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \leq \frac{d_{\mathcal{F}}(\check{\mathbf{P}})}{s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}})} + \frac{\varepsilon}{2\kappa_2}, \quad (39b)$$

$$\left\| \frac{2}{\mu} \check{\Omega} + \mathbf{I}_m \right\|_{\mathbb{F}} \leq \kappa_2 \left( -\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \right) + d_{\mathcal{F}}(\check{\mathbf{P}}) + \varepsilon, \quad (39c)$$

then the pair  $(\check{\mathbf{P}}, \check{\mathbf{P}} \check{\mathbf{P}}^\top)$  is the unique primal solution of the penalized convex relaxation (6a)–(6d), where  $\kappa$  and  $\kappa_2$  are defined in (9c) and (35), respectively.

*Proof.* According to Lemma 4, it suffices to find  $\theta \in \mathbb{R}$  that satisfies (30a) and (30b). Define:

$$\theta \triangleq -\frac{2^{-1}\mu(\kappa - \kappa_2) d_{\mathcal{F}}(\check{\mathbf{P}})}{s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}})} - \mu\varepsilon$$

The conic inequality (30a) can be proven as follows:

$$\begin{aligned} 2\mu^{-1} \| \mathbf{M}_0 - \sum_{k \in \mathcal{K}} \check{\lambda}_k \mathbf{M}_k \|_{\mathbb{F}} &\leq 2\mu^{-1} \| \mathbf{M}_0 \| - (\kappa - \kappa_2) \frac{\mathbf{1}^\top \check{\lambda}}{\mu} \\ &\leq \varepsilon + (\kappa - \kappa_2) \left( \frac{d_{\mathcal{F}}(\check{\mathbf{P}})}{s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}})} + \frac{\varepsilon}{2\kappa_2} \right) \\ &= \frac{(\kappa - \kappa_2) d_{\mathcal{F}}(\check{\mathbf{P}})}{s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}})} + 2\varepsilon - \frac{(3\kappa_2 - \kappa)\varepsilon}{2\kappa_2} < -2\mu^{-1}\theta \\ &\Rightarrow -\mathbf{M}_0 + \sum_{k \in \mathcal{K}} \check{\lambda}_k \mathbf{M}_k \prec_{\check{\mathbf{C}}} -\theta \mathbf{I}_n. \end{aligned}$$

Moreover, the conic inequality (30b) can be proven as:

$$\begin{aligned} \left\| \frac{2}{\mu} \check{\Omega} + \mathbf{I}_m \right\|_{\mathbb{F}} &\leq \kappa_2 \left( -\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \right) + d_{\mathcal{F}}(\check{\mathbf{P}}) + \varepsilon \\ &= \kappa_2 \left( -\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \right) + \varepsilon + 1 - 4\varepsilon - \frac{\kappa d_{\mathcal{F}}(\check{\mathbf{P}})}{s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}})} \\ &\leq \kappa_2 \left( \frac{d_{\mathcal{F}}(\check{\mathbf{P}})}{s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}})} + \frac{\varepsilon}{2\kappa_2} \right) \\ &\quad + \varepsilon + 1 - 4\varepsilon - \frac{\kappa d_{\mathcal{F}}(\check{\mathbf{P}})}{s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}})} \\ &= 1 - \frac{(\kappa - \kappa_2) d_{\mathcal{F}}(\check{\mathbf{P}})}{s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}})} - \frac{5\varepsilon}{2} \\ &= 1 + 2\mu^{-1}\theta - \frac{\varepsilon}{2} < 1 + 2\mu^{-1}\theta \\ &\Rightarrow \check{\Omega} + \frac{\mu}{2} \mathbf{I}_m \prec_{\check{\mathcal{D}}} \left( \frac{\mu}{2} + \theta \right) \mathbf{I}_m \Rightarrow \check{\Omega} \prec_{\check{\mathcal{D}}} \theta \mathbf{I}_m. \end{aligned}$$

Hence, according to Lemma 4, the pair  $(\check{\mathbf{P}}, \check{\mathbf{P}} \check{\mathbf{P}}^\top)$  is the unique primal solution of the penalized convex relaxation



(6a)–(6d). □

*Proof of Theorem 1.* Due to the main assumption, it is straightforward to verify the following three inequalities:

$$\mu^{-1}\psi < \sqrt{\beta\mu^{-1}}, \quad (42a)$$

$$2\kappa\sqrt{\beta\mu^{-1}} < 13^{-1}s(\check{\mathbf{P}}), \quad (42b)$$

$$\sqrt{\beta\mu^{-1}} < 12^{-1}. \quad (42c)$$

Consider an arbitrary optimal solution  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  of (16a)–(16e). The point  $\check{\mathbf{P}}$  is consequently feasible for (1a)–(1c). Therefore  $d_{\mathcal{F}}(\check{\mathbf{P}}) = 0$  and the inequalities (25) and (26) are satisfied. According to Lemma 5, there exist Lagrange multipliers  $(\check{\gamma}, \check{\gamma}, \check{\lambda}, \check{\Omega}) \in \mathbb{R}^w \times \mathbb{R}^w \times \mathbb{R}^p \times \mathbb{S}_m$  corresponding to  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  that satisfy the KKT conditions (27a)–(27f) as well as the inequalities (34a) and (34b). Based on Lemma 6 and since  $d_{\mathcal{F}}(\check{\mathbf{P}}) = 0$ , in order to prove the theorem, it suffices to show that:

$$2\mu^{-1}\|\mathbf{M}_0\| \leq 4^{-1} \quad (43a)$$

$$-\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \leq \frac{4^{-1}}{2\kappa_2} \quad (43b)$$

$$\left\| \frac{2}{\mu} \check{\Omega} + \mathbf{I}_m \right\|_{\text{F}} \leq \kappa_2 \left( -\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \right) + 4^{-1}. \quad (43c)$$

- (43a) is the direct consequence of (42a):

$$2\mu^{-1}\|\mathbf{M}_0\| \leq \mu^{-1}\psi \leq \sqrt{\beta\mu^{-1}} \leq 12^{-1} < 4^{-1}. \quad (44)$$

- (43b) is the direct consequence of (34a), (42b), and (42c):

$$-\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \leq \frac{\mu^{-1}\psi + 2\sqrt{\beta\mu^{-1}}}{s(\check{\mathbf{P}}) - 2\kappa\sqrt{\beta\mu^{-1}}} \leq \frac{\sqrt{\beta\mu^{-1}} + 2\sqrt{\beta\mu^{-1}}}{s(\check{\mathbf{P}}) - 2\kappa\sqrt{\beta\mu^{-1}}} \quad (45a)$$

$$\leq \frac{\sqrt{\beta\mu^{-1}} + 2\sqrt{\beta\mu^{-1}}}{s(\check{\mathbf{P}}) - 13^{-1}s(\check{\mathbf{P}})} = \frac{3\sqrt{\beta\mu^{-1}}}{(1 - 13^{-1})s(\check{\mathbf{P}})} \quad (45b)$$

$$\leq \frac{3 \times 13^{-1}(2\kappa)^{-1}s(\check{\mathbf{P}})}{(1 - 13^{-1})s(\check{\mathbf{P}})} = \frac{4^{-1}}{2\kappa} < \frac{4^{-1}}{2\kappa_2}. \quad (45c)$$

- (43c) can be concluded from (34b), (42a), and (42c):

$$\left\| \frac{2}{\mu} \check{\Omega} + \mathbf{I}_m \right\|_{\text{F}} \leq \kappa_2 \left( -\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \right) + \mu^{-1}\psi + 2\sqrt{\beta\mu^{-1}} \quad (46a)$$

$$\leq \kappa_2 \left( -\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \right) + 3\sqrt{\beta\mu^{-1}} \quad (46b)$$

$$\leq \kappa_2 \left( -\frac{\mathbf{1}^\top \check{\lambda}}{\mu} \right) + 4^{-1}. \quad (46c)$$

Hence, according to Lemma 6, the point  $(\check{\mathbf{P}}, \check{\mathbf{P}}\check{\mathbf{P}}^\top)$  is the unique optimal solution for the penalized relaxation (6a)–(6e), for which the relaxed constraints (3c) and (3d) are satisfied. Finally, due to feasibility of the pair  $(\check{\mathbf{P}}, \check{\mathbf{P}}\check{\mathbf{P}}^\top)$ , we have:

$$\bar{f}_0(\check{\mathbf{P}}) + g_0(\check{\mathbf{P}}) - \mu m = f_0(\check{\mathbf{P}}, \check{\mathbf{P}}\check{\mathbf{P}}^\top) + g_0(\check{\mathbf{P}}) - \mu \langle \check{\mathbf{P}}, \check{\mathbf{P}} \rangle \quad (47a)$$

$$\geq f_0(\check{\mathbf{P}}, \check{\mathbf{P}}\check{\mathbf{P}}^\top) + g_0(\check{\mathbf{P}}) - \mu \langle \check{\mathbf{P}}, \check{\mathbf{P}} \rangle \quad (47b)$$

$$\geq \bar{f}_0(\check{\mathbf{P}}) + g_0(\check{\mathbf{P}}) - \mu m \quad (47c)$$

and the proof is completed. □

*Proof of Theorem 2.* Consider an arbitrary optimal solution  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  of (16a)–(16e). Due to the main assumption, (25) is satisfied and if  $\mu$  is large, then (26) is satisfied as well. Moreover, according to Lemma 5, there exist Lagrange multipliers  $(\check{\gamma}, \check{\gamma}, \check{\lambda}, \check{\Omega}) \in \mathbb{R}^w \times \mathbb{R}^w \times \mathbb{R}^p \times \mathbb{S}_m$  corresponding to  $(\check{\mathbf{P}}, \check{\mathbf{t}})$  that satisfy the KKT conditions (27a)–(27f) as well as the inequalities (34a) and (34b). According to Lemma 6, the proof follows directly from the fact that

$$\varepsilon = \frac{1}{4} \left( 1 - d_{\mathcal{F}}(\check{\mathbf{P}}) - \frac{\kappa d_{\mathcal{F}}(\check{\mathbf{P}})}{s(\check{\mathbf{P}}) - \kappa d_{\mathcal{F}}(\check{\mathbf{P}})} \right) > 0, \quad (48)$$

and therefore, if  $\mu$  is sufficiently large, the inequalities (34a) and (34b) conclude (39a)–(39c). As a result, if  $\mu$  is sufficiently large,  $(\check{\mathbf{P}}, \check{\mathbf{P}}\check{\mathbf{P}}^\top)$  is the unique primal solution of the penalized convex relaxation (6a)–(6d). □

**Lemma 7.** Consider a set  $\mathcal{Q} \in \mathcal{F}$  for which there exists  $\nu > 0$  such that  $s(\mathbf{P}) > \nu$  for every  $\mathbf{P} \in \mathcal{Q}$ . For every

$$\mu > \max\{\beta^{-1}\psi^2, \beta(26\kappa)^2\nu^{-2}, 144\beta\}, \quad (49)$$

define  $h_{\mathcal{Q},\mu} : \mathcal{Q} \rightarrow \mathcal{F}$  as the function mapping any initial point  $\check{\mathbf{P}} \in \mathcal{Q}$  in problem (16a)–(16e) to its unique solution  $\check{\mathbf{P}}^*$  (whose existence and uniqueness is guaranteed by Theorem (1)). Then  $h_{\mathcal{Q},\mu}$  is continuous throughout  $\mathcal{Q}$ .

*Proof.* According to Berge's maximum theorem,  $h_{\mathcal{Q},\mu}$  is upper hemicontinuous. However, according to Theorem (1) it is a function and therefore, it is continuous. □

*Proof of Theorem 3.* Let  $\{\mathbf{P}^k\}_{k=0}^\infty$  represent the sequence generated by Algorithm (1). Assume by induction that  $\mathbf{P}^k \in \check{\mathcal{F}}$  and let  $(\check{\mathbf{P}}, \check{\mathbf{P}}\check{\mathbf{P}}^\top)$  be the solution to the problem (6a)–(6e) with  $\check{\mathbf{P}} = \mathbf{P}^k$ . According to the optimality of  $\check{\mathbf{P}}$  and feasibility of  $\mathbf{P}^k$  we have:

$$\begin{aligned} \bar{f}_0(\check{\mathbf{P}}) + g_0(\check{\mathbf{P}}) + \frac{\mu}{2} \|\check{\mathbf{P}} - \mathbf{P}^k\|_{\text{F}}^2 &\leq \\ \bar{f}_0(\check{\mathbf{P}}) + g_0(\check{\mathbf{P}}) - \mu \langle \check{\mathbf{P}}, \mathbf{P}^k \rangle + \mu m &\leq \\ \bar{f}_0(\check{\mathbf{P}}) + g_0(\check{\mathbf{P}}) - \mu \langle \mathbf{P}^k, \mathbf{P}^k \rangle + \mu m &\leq \bar{f}_0(\mathbf{P}^k) + g_0(\mathbf{P}^k). \end{aligned} \quad (50)$$

Hence, according to Theorem (1),  $\mathbf{P}^{k+1} = \check{\mathbf{P}} \in \check{\mathcal{F}}$  and the sequence  $\{\bar{f}_0(\mathbf{P}^k) + g_0(\mathbf{P}^k)\}_{k=0}^\infty$  is monotonically non-increasing and convergent. On the other hand, according to (50), we have

$$\begin{aligned} \|\mathbf{P}^{k+1} - \mathbf{P}^k\|_{\text{F}}^2 &\leq \\ 2\mu^{-1} [\bar{f}_0(\mathbf{P}^k) + g_0(\mathbf{P}^k) - \bar{f}_0(\mathbf{P}^{k+1}) - g_0(\mathbf{P}^{k+1})] &\quad (51) \end{aligned}$$

which implies that the sequence  $\{\mathbf{P}^k\}_{k=0}^\infty$  is convergent to a point  $\mathbf{P}^\infty \in \check{\mathcal{F}}$ .

Define  $h_{\check{\mathcal{F}},\mu} : \check{\mathcal{F}} \rightarrow \check{\mathcal{F}}$  as the function mapping any initial point  $\check{\mathbf{P}} \in \check{\mathcal{F}}$  in problem (16a)–(16e) to its unique solution. According to Lemma 7,  $h_{\check{\mathcal{F}},\mu}$  is continuous, and therefore:

$$h_{\check{\mathcal{F}},\mu}(\mathbf{P}^\infty) = \mathbf{P}^\infty \quad (52)$$

Now, according to Lemma 3, there exists Lagrange multipliers  $(\check{\gamma}, \check{\gamma}, \check{\lambda}, \check{\Omega}) \in \mathbb{R}^w \times \mathbb{R}^w \times \mathbb{R}^{|\mathcal{K}|} \times \mathbb{S}_m$  that satisfy the following

Karush–Kuhn–Tucker (KKT) conditions

$$\nabla_P \mathcal{L}(\mathbf{P}^\infty, \mathbf{t}^*, \tilde{\gamma}^*, \tilde{\gamma}^*, \tilde{\lambda}^*, \tilde{\Omega}^*) = \mu \mathbf{P}^\infty, \quad (53a)$$

$$\mathbf{1} + \tilde{\gamma}^* + \tilde{\gamma}^* = \mathbf{0}, \quad (53b)$$

$$\tilde{\gamma}^* \circ (+\alpha(\mathbf{P}^\infty) + \mathbf{b} - \mathbf{t}^*) = \mathbf{0}, \quad (53c)$$

$$\tilde{\gamma}^* \circ (-\alpha(\mathbf{P}^\infty) - \mathbf{b} - \mathbf{t}^*) = \mathbf{0}, \quad (53d)$$

$$\tilde{\lambda}_k (\langle \mathbf{M}_k, \mathbf{P}^\infty \mathbf{P}^{\infty \top} \rangle + \langle \mathbf{N}_k, \mathbf{P}^\infty \rangle + q_k) = 0 \quad k \in \mathcal{K}, \quad (53e)$$

$$\tilde{\gamma}^* \leq \mathbf{0}, \quad \tilde{\gamma}^* \leq \mathbf{0}, \quad \tilde{\lambda}^* \leq \mathbf{0}, \quad (53f)$$

and the pair of primal and dual solutions  $\mathbf{P}^\infty$  and  $(\tilde{\gamma}^*, \tilde{\gamma}^*, \tilde{\lambda}^*, \tilde{\Omega}^* - \mathbf{I}/2)$  satisfy KKT optimality conditions for the problem (1a) – (1c).  $\square$

## VI. CONCLUSIONS

This work introduces convex relaxations for solving a broad class of non-convex and non-smooth optimization problems involving orthogonality constraints. The proposed approach relies on solving a sequence of penalized convex relaxations to find feasible and near optimal points for a given non-convex orthogonality-constrained problem. Experimental results on two fundamental problems in machine learning demonstrate the potential and effectiveness of the proposed approach in solving practical problems.

## REFERENCES

- [1] F. Zohrizadeh, M. Kheirandishfard, F. Kamangar, and R. Madani, “Non-smooth optimization over stiefel manifolds with applications to dimensionality reduction and graph clustering,” in *IJCAI*, 2019, pp. 1319–1326.
- [2] W. Bian and D. Tao, “Max-min distance analysis by using sequential SDP relaxation for dimension reduction,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, no. 5, pp. 1037–1050, 2011.
- [3] B. Jiang, J. Tang, C. H. Ding, and B. Luo, “Nonnegative orthogonal graph matching,” in *AAAI*, 2017.
- [4] J. Tang and H. Liu, “Unsupervised feature selection for linked social media data,” in *SIGKDD*, 2012.
- [5] Y. Yang, H. T. Shen, Z. Ma, Z. Huang, and X. Zhou, “ $\ell_{2,1}$ -norm regularized discriminative feature selection for unsupervised learning,” in *IJCAI*, 2011.
- [6] V. Ozoliņš, R. Lai, R. Caflich, and S. Osher, “Compressed modes for variational problems in mathematics and physics,” *PNAS*, vol. 110, no. 46, pp. 18 368–18 373, 2013.
- [7] W. Chen, H. Ji, and Y. You, “An augmented Lagrangian method for  $\ell_1$ -regularized optimization problems with orthogonality constraints,” *SIAM J SCI COMPUT*, vol. 38, no. 4, pp. B570–B592, 2016.
- [8] T. E. Abrandan, J. Eriksson, and V. Koivunen, “Steepest descent algorithms for optimization under unitary matrix constraint,” *IEEE T SIGNAL PROCES*, vol. 56, no. 3, pp. 1134–1147, 2008.
- [9] Z. Wen and W. Yin, “A feasible method for optimization with orthogonality constraints,” *MATH PROGRAM*, vol. 142, no. 1-2, pp. 397–434, 2013.
- [10] B. Gao, X. Liu, X. Chen, and Y.-x. Yuan, “A new first-order algorithmic framework for optimization problems with orthogonality constraints,” *SIAM J OPTIMIZ*, vol. 28, no. 1, pp. 302–332, 2018.
- [11] M. Fortin and R. Glowinski, *Augmented Lagrangian methods: applications to the numerical solution of boundary-value problems*. Elsevier, 2000, vol. 15.
- [12] H. Attouch, J. Bolte, and B. F. Svaiter, “Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized Gauss–Seidel methods,” *MATH PROGRAM*, vol. 137, no. 1-2, pp. 91–129, 2013.
- [13] R. Lai and S. Osher, “A splitting method for orthogonality constrained problems,” *SIAM J SCI COMPUT*, vol. 58, no. 2, pp. 431–449, 2014.
- [14] W. Yin, S. Osher, D. Goldfarb, and J. Darbon, “Bregman iterative algorithms for  $\ell_1$ -minimization with applications to compressed sensing,” *SIAM J IMAGING SCI*, vol. 1, no. 1, pp. 143–168, 2008.
- [15] H. Zhu, X. Zhang, D. Chu, and L.-Z. Liao, “Nonconvex and nonsmooth optimization with generalized orthogonality constraints: An approximate augmented Lagrangian method,” *SIAM J SCI COMPUT*, pp. 1–42, 2017.
- [16] S. Ibaraki and M. Tomizuka, “Rank minimization approach for solving BMI problems with random search,” in *Proceedings of the 2001 American Control Conference.(Cat. No. 01CH37148)*, vol. 3. IEEE, 2001, pp. 1870–1875.
- [17] M. Kheirandishfard, F. Zohrizadeh, and R. Madani, “Convex relaxation of bilinear matrix inequalities part I: Theoretical results,” in *IEEE 57th Annual Conference on Decision and Control (CDC)*, 2018.
- [18] M. Kheirandishfard, F. Zohrizadeh, M. Adil, and R. Madani, “Convex relaxation of bilinear matrix inequalities part II: Applications to optimal control synthesis,” in *IEEE 57th Annual Conference on Decision and Control (CDC)*, 2018.
- [19] R. Madani, M. Kheirandishfard, J. Lavaei, and A. Atamturk, “Penalized conic relaxations for quadratically-constrained quadratic programming,” *Preprint*, 2019.
- [20] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [21] K. Nakata, K. Fujisawa, M. Fukuda, M. Kojima, and K. Murota, “Exploiting sparsity in semidefinite programming via matrix completion II: Implementation and numerical results,” *Mathematical Programming*, vol. 95, no. 2, pp. 303–327, 2003.
- [22] P. N. Belhumeur, J. P. Hespanha, and D. J. Kriegman, “Eigenfaces vs. fisherfaces: Recognition using class specific linear projection,” *IEEE TPAMI*, no. 7, pp. 711–720, 1997.
- [23] A. S. Georghiadis, P. N. Belhumeur, and D. J. Kriegman, “From few to many: Illumination cone models for face recognition under variable lighting and pose,” *IEEE T PATTERN ANAL*, vol. 23, no. 6, pp. 643–660, 2001.
- [24] A. Y. Ng, M. I. Jordan, and Y. Weiss, “On spectral clustering: Analysis and an algorithm,” in *NIPS*, 2002, pp. 849–856.
- [25] J. Han, K. Xiong, and F. Nie, “Orthogonal and nonnegative graph reconstruction for large scale clustering,” in *IJCAI*, 2017.
- [26] D. Dua and C. Graff, “UCI machine learning repository,” 2017. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [27] A. Gionis, H. Mannila, and P. Tsaparas, “Clustering aggregation,” *ACM T KNOWL DISCOV D*, vol. 1, no. 1, p. 4, 2007.
- [28] A. K. Jain and M. H. Law, “Data clustering: A user’s dilemma,” in *ICPRAI*, 2005.
- [29] H. Chang and D.-Y. Yeung, “Robust path-based spectral clustering,” *PATTERN RECOGN*, vol. 41, no. 1, pp. 191–203, 2008.