

Low-Rank Extragradient Method for Nonsmooth and Low-Rank Matrix Optimization Problems*

Dan Garber

Technion - Israel Institute of Technology
dangar@technion.ac.il

Atara Kaplan

Technion - Israel Institute of Technology
ataragold@campus.technion.ac.il

Abstract

Low-rank and nonsmooth matrix optimization problems capture many fundamental tasks in statistics and machine learning. While significant progress has been made in recent years in developing efficient methods for *smooth* low-rank optimization problems that avoid maintaining high-rank matrices and computing expensive high-rank SVDs, advances for nonsmooth problems have been slow paced.

In this paper we consider standard convex relaxations for such problems. Mainly, we prove that under a natural *generalized strict complementarity* condition and under the relatively mild assumption that the nonsmooth objective can be written as a maximum of smooth functions, the *extragradient method*, when initialized with a “warm-start” point, converges to an optimal solution with rate $O(1/t)$ while requiring only two *low-rank* SVDs per iteration. We give a precise trade-off between the rank of the SVDs required and the radius of the ball in which we need to initialize the method. We support our theoretical results with empirical experiments on several nonsmooth low-rank matrix recovery tasks, demonstrating that using simple initializations, the extragradient method produces exactly the same iterates when full-rank SVDs are replaced with SVDs of rank that matches the rank of the (low-rank) ground-truth matrix to be recovered.

1 Introduction

Low-rank and nonsmooth matrix optimization problems have many important applications in statistics, machine learning, and related fields, such as *sparse PCA* [21, 35], *robust PCA* [29, 34, 2, 9, 39], *phase synchronization* [42, 6, 30], *community detection and stochastic block models* [1]¹, *low-rank and sparse covariance matrix recovery* [36], *robust matrix completion* [23, 10], and more. For many of these problems, convex relaxations, in which one replaces the nonconvex low-rank constraint with a trace-norm constraint, have been demonstrated in numerous papers to be highly effective both in theory (under suitable

*This version corrects an error in the original paper published in NeurIPS 2021 [22]: while the version [22] provides convergence rates w.r.t. the best iterate (which under the assumptions of the paper is guaranteed to be low-rank), this corrected version provides the same rates but for the ergodic sequence, i.e., the averaged iterate (which, under our assumptions, is the average of low-rank iterates).

¹in [42, 6, 30] and [1] the authors consider SDPs with linear objective function and affine constraints of the form $\mathcal{A}(\mathbf{X}) = \mathbf{b}$. By incorporating the linear constraints into the objective function via a ℓ_2 penalty term of the form $\lambda \|\mathcal{A}(\mathbf{X}) - \mathbf{b}\|_2$, $\lambda > 0$, we obtain a nonsmooth objective function.

assumptions) and empirically (see references above). These convex relaxations can be formulated as the following general nonsmooth optimization problem:

$$\min_{\mathbf{X} \in \mathcal{S}_n} g(\mathbf{X}), \quad (1)$$

where $g : \mathbb{S}^n \rightarrow \mathbb{R}$ is convex but nonsmooth, and $\mathcal{S}_n = \{\mathbf{X} \in \mathbb{S}^n \mid \text{Tr}(\mathbf{X}) = 1, \mathbf{X} \succeq 0\}$ is the spectrahedron in \mathbb{S}^n , \mathbb{S}^n being the space of $n \times n$ real symmetric matrices.

Problem (1), despite being convex, is notoriously difficult to solve in large scale. The simplest and most general approach applicable to it is the *projected subgradient method* [3, 7], which requires on each iteration to compute a Euclidean projection onto the spectrahedron \mathcal{S}_n , which in worst case amounts to $O(n^3)$ runtime per iteration. In many applications $g(\mathbf{X})$ follows a composite model, i.e., $g(\mathbf{X}) = h(\mathbf{X}) + w(\mathbf{X})$, where $h(\cdot)$ is convex and smooth and $w(\cdot)$ is convex and nonsmooth but admits a simple structure (e.g., nonsmooth regularizer). For such composite objectives, without the spectrahedron constraint, proximal methods such as FISTA [4] or splitting methods such as ADMM [33] are often very effective. However, with the spectrahedron constraint, all such methods require on each iteration to apply a subprocedure (e.g., computing the proximal mapping) which in worst case amounts to at least $O(n^3)$ runtime. A third type of off-the-shelf methods include those which are based on the *conditional gradient method* and adapted to nonsmooth problems, see for instance [32, 19, 37, 27]. The advantage of such methods is that no expensive high-rank SVD computations are needed. Instead, only a single leading eigenvector computation (i.e., a rank-one SVD) per iteration is required. However, similarly to the subgradient method, these suffer from slow convergence rates ($O(1/\sqrt{t})$ at best, when t is the iteration counter) even when the nonsmooth problem admits favorable structure (as we detail below). Thus, to conclude, standard first-order methods for Problem (1) require in worst case $\Omega(n^3)$ runtime per iteration or suffer from worst-case slow convergence rates.

In the recent works [17, 18] it was established that for smooth objective functions, the high-rank SVD computations required for Euclidean projections onto the spectrahedron in standard gradient methods, can be replaced with low-rank SVDs in the close proximity of a low-rank optimal solution. This is significant since the runtime to compute a rank- r SVD of a given $n \times n$ matrix using efficient iterative methods typically scales with rn^2 (and further improves when the matrix is sparse), instead of n^3 for a full-rank SVD. These results depend on the existence of eigen-gaps in the gradient of the optimal solution, which we refer to as a *generalized strict complementarity condition*. These results also hinge on a unique property of the Euclidean projection onto the spectrahedron. The projection onto the spectrahedron of a matrix $\mathbf{X} \in \mathbb{S}^n$, which admits an eigen-decomposition $\mathbf{X} = \sum_{i=1}^n \lambda_i \mathbf{v}_i \mathbf{v}_i^\top$, is given by

$$\Pi_{\mathcal{S}_n}[\mathbf{X}] = \sum_{i=1}^n \max\{0, \lambda_i - \lambda\} \mathbf{v}_i \mathbf{v}_i^\top, \quad (2)$$

where $\lambda \in \mathbb{R}$ is the unique scalar satisfying $\sum_{i=1}^n \max\{0, \lambda_i - \lambda\} = 1$. This operation thus truncates all eigenvalues that are smaller than λ , while leaving the eigenvectors unchanged, thereby returning a matrix with rank equal to the number of eigenvalues greater than λ . Importantly, when the projection of \mathbf{X} onto \mathcal{S}_n is of rank r , only the first r components in the eigen-decomposition of \mathbf{X} are required to compute it in the first place, and thus, only a rank- r SVD of \mathbf{X} is required. In other words and simplifying,

[17, 18] show that under strict complementarity, at the proximity of an optimal solution of rank r , the exact Euclidean projection equals the rank- r truncated projection given by:

$$\hat{\Pi}_{\mathcal{S}_n}^r[\mathbf{X}] := \Pi_{\mathcal{S}_n} \left[\sum_{i=1}^r \lambda_i \mathbf{v}_i \mathbf{v}_i^\top \right]. \quad (3)$$

Extending the results of [17, 18] to the nonsmooth setting is difficult since the smoothness assumption is critical to the analysis. Moreover, while [17, 18] rely on certain eigen-gaps in the gradients at optimal points, for nonsmooth problems, since the subdifferential set is often not a singleton, it is not likely that a similar eigen-gap property holds for all subgradients of an optimal solution.

In this paper we show that under the mild assumption that Problem (1) can be formulated as a smooth convex-concave saddle-point problem, i.e., the nonsmooth term can be written as a maximum over (possibly infinite number of) smooth convex functions, we can obtain results in the spirit of [17, 18]. Concretely, we show that if a generalized strict complementarity (GSC) assumption holds for a low-rank optimal solution (see Assumption 1 in the sequel), the *extragradient method* for smooth convex-concave saddle-point problems [24, 31] (see Algorithm 1 below), when initialized in the proximity of the optimal solution, converges with its original convergence rate of $O(1/t)$, while requiring only two low-rank SVDs per iteration². It is important to recall that while the extragradient method requires two SVDs per iteration, it has the benefit of a fast $O(1/t)$ convergence rate, while simpler saddle-point methods such as mirror-descent-based only achieve a $O(1/\sqrt{t})$ rate [7].

Our contributions can be summarized as follows:

- We prove that even under (standard) strict complementarity, the projected sub-gradient method, when initialized with a “warm-start”, may produce iterates with rank higher than that of the optimal solution. This phenomena further motivates our saddle-point approach. See Lemma 5.
- We suggest a generalized strict complementarity (GSC) condition for saddle-point problems and prove that when $g(\cdot)$ — the objective function in Problem (1), admits a highly popular saddle-point structure (one which captures all applications we mentioned in this paper), GSC w.r.t. an optimal solution to Problem (1) implies GSC (with the same parameters) w.r.t. a corresponding optimal solution of the equivalent saddle-point problem (the other direction always holds). See Section 3.
- **Main result:** we prove that for a smooth convex-concave saddle-point problem and an optimal solution which satisfies GSC, the extragradient method, when initialized with a “warm-start”, converges with its original rate of $O(1/t)$ while requiring only two low-rank SVDs per iteration. Moreover, we prove GSC facilitates a precise and powerful tradeoff: increasing the rank of SVD computations (beyond the rank of the optimal solution) can significantly increase the radius of the ball in which the method needs to be initialized. See Theorem 1.
- We present extensive numerical evidence that demonstrate both the plausibility of the GSC assumption in various tasks, and more importantly, demonstrate that

²note that the extragradient method computes two projected-gradient steps on each iteration, and thus two SVDs are needed per iteration.

indeed the extragradient method with simple initialization converges correctly (i.e., produces exactly the same sequences of iterates) when the rank of the SVDs used to compute the (truncated) projections matches the rank of the (low-rank) ground-truth matrix to be recovered, instead of naively using full-rank SVDs (as suggested by (2)). See Section 5.

1.1 Additional related work

Since, as in the works [17, 18] mentioned before which deal with smooth objectives, strict complementarity plays a key role in our analysis, we refer the interested reader to the recent works [16, 40, 13, 20] which also exploit this property for efficient smooth and convex optimization over the spectrahedron. Strict complementarity has also played an instrumental role in two recent and very influential works which used it to prove linear convergence rates for proximal gradient methods [43, 14].

Besides convex relaxations such as Problem (1), considerable advances have been made in the past several years in developing efficient *nonconvex* methods with global convergence guarantees for low-rank matrix problems. In [38] the authors consider semidefinite programs and prove that under a smooth manifold assumption on the constraints, such methods converge to the optimal global solution. In [25] the authors prove global convergence of factorized nonconvex gradient descent from a “warm-start” initialization point for non-linear smooth minimization on the positive semidefinite cone. Very recently, [8] has established, under statistical conditions, fast convergence results from “warm-start” initialization of nonconvex first-order methods, when applied to nonsmooth nonconvex matrix recovery problems which are based on the explicit factorization of the low-rank matrix. A result of similar flavor concerning nonsmooth and nonconvex formulation of robust recovery of low-rank matrices from random linear measurements was presented in [26]. Finally, several recent works have considered nonconvex low-rank regularizers which result in nonconvex nonsmooth optimization problems, but guarantee convergence only to a stationary point [28, 41].

2 Strict Complementarity for Nonsmooth Optimization and Difficulty of Applying Low-Rank Projected Subgradient Steps

Our analysis of the nonsmooth Problem (1) naturally depends on certain subgradients of an optimal solution which, in many aspects, behave like the gradients of smooth functions. The existence of such a subgradient is guaranteed from the first-order optimality condition for constrained convex minimization problems:

Lemma 1 (first-order optimality condition, see [3]). *Let $g : \mathbb{S}^n \rightarrow \mathbb{R}$ be a convex function. Then $\mathbf{X}^* \in \mathcal{S}_n$ minimizes g over \mathcal{S}_n if and only if there exists a subgradient $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ such that $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$.*

For some $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ which satisfies the first-order optimality condition for an optimal solution \mathbf{X}^* , if the multiplicity of the smallest eigenvalue equals $r^* = \text{rank}(\mathbf{X}^*)$, then

it can be shown that the optimal solution satisfies a strict complementarity assumption. The equivalence between a standard strict complementarity assumption on some low-rank optimal solution of a *smooth* optimization problem over the spectrahedron and an eigen-gap in the gradient of the optimal solution was established in [40]. We generalize this equivalence to also include nonsmooth problems. The proof follows similar arguments and is given in Appendix A.

Definition 1 (strict complementarity). *An optimal solution $\mathbf{X}^* \in \mathcal{S}_n$ of rank r^* for Problem (1) satisfies the strict complementarity assumption with parameter $\delta > 0$, if there exists an optimal solution of the dual problem³ $(\mathbf{Z}^*, s^*) \in \mathbb{S}^n \times \mathbb{R}$ such that $\text{rank}(\mathbf{Z}^*) = n - r^*$, and $\lambda_{n-r^*}(\mathbf{Z}^*) \geq \delta$.*

Lemma 2. *Let $\mathbf{X}^* \in \mathcal{S}_n$ be a rank- r^* optimal solution to Problem (1). \mathbf{X}^* satisfies the (standard) strict complementarity assumption with parameter $\delta > 0$ if and only if there exists a subgradient $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ such that $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$ and $\lambda_{n-r^*}(\mathbf{G}^*) - \lambda_n(\mathbf{G}^*) \geq \delta$.*

Throughout this paper we assume a weaker and more general assumption than strict complementarity, namely generalized strict complementarity (GSC), which we present now.

Assumption 1 (generalized strict complementarity). *We say an optimal solution \mathbf{X}^* to Problem (1) satisfies the generalized strict complementarity assumption with parameters r, δ , if there exists a subgradient $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ such that $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$ and $\lambda_{n-r}(\mathbf{G}^*) - \lambda_n(\mathbf{G}^*) \geq \delta$.*

In [17] the author presents several characteristic properties of the gradient of the optimal solution in optimization problems over the spectrahedron. Using the existence of subgradients which satisfy the condition in Lemma 1, we can extend these properties also to the nonsmooth setting. The following lemma shows that GSC with parameters (r, δ) for some $\delta > 0$ (Assumption 1) is a sufficient condition for the optimal solution to be of rank at most r . The proof follows immediately from the proof of the analogous Lemma 7 in [17], by replacing the gradient of the optimal solution with a subgradient for which the first-order optimality condition holds.

Lemma 3. *Let \mathbf{X}^* be an optimal solution to Problem (1) and write its eigen-decomposition as $\mathbf{X}^* = \sum_{i=1}^{r^*} \lambda_i \mathbf{v}_i \mathbf{v}_i^T$. Then, any subgradient $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ which satisfies $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$, admits an eigen-decomposition such that the set of vectors $\{\mathbf{v}_i\}_{i=1}^{r^*}$ is a set of leading eigenvectors of $(-\mathbf{G}^*)$ which corresponds to the eigenvalue $\lambda_1(-\mathbf{G}^*) = -\lambda_n(\mathbf{G}^*)$. Furthermore, there exists at least one such subgradient.*

One motivation for assuming (standard) strict complementarity (Assumption 1 with parameters $r = \text{rank}(\mathbf{X}^*)$ and $\delta > 0$) is that it guarantees a certain notion of robustness of the problem to small perturbations in the parameters. It is well known (see for instance [3]) that a projected subgradient step from \mathbf{X}^* with respect to a subgradient $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ for which the first-order optimality condition holds, returns the optimal solution \mathbf{X}^* itself. This implies that $\text{rank}(\Pi_{\mathcal{S}_n}[\mathbf{X}^* - \eta \mathbf{G}^*]) = \text{rank}(\mathbf{X}^*)$ (here η is the step-size). Without (standard) strict complementarity however, a small change in the parameters could result

³Denote $q(\mathbf{Z}, s) = \min_{\mathbf{X} \in \mathbb{S}^n} \{g(\mathbf{X}) + s(1 - \text{Tr}(\mathbf{X})) - \langle \mathbf{Z}, \mathbf{X} \rangle\}$. The dual problem of Problem (1) can be written as: $\max_{\{\mathbf{Z} \succeq 0, s \in \mathbb{R}\}} \{q(\mathbf{Z}, s) \mid (\mathbf{Z}, s) \in \text{dom}(q)\}$.

in a higher rank matrix. This is captured in the following lemma which is analogous to Lemma 3 in [18], where again the proof is straightforward from the proof in [18] by replacing the gradient of the optimal solution with a subgradient for which the first-order optimality condition holds.

Lemma 4. *Let \mathbf{X}^* be an optimal solution of rank r^* to Problem (1). Let $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ be a subgradient at \mathbf{X}^* such that $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$. Then, $\lambda_{n-r^*}(\mathbf{G}^*) = \lambda_n(\mathbf{G}^*)$ if and only if for any arbitrarily small $\zeta > 0$ it holds that $\text{rank}(\Pi_{(1+\zeta)\mathcal{S}_n}[\mathbf{X}^* - \eta \mathbf{G}^*]) > r^*$, where $\eta > 0$, $(1 + \zeta)\mathcal{S}_n = \{(1 + \zeta)\mathbf{X} \mid \mathbf{X} \in \mathcal{S}_n\}$, and $\Pi_{(1+\zeta)\mathcal{S}_n}[\cdot]$ denotes the Euclidean projection onto the set $(1 + \zeta)\mathcal{S}_n$.*

2.1 The challenge of applying low-rank projected subgradient steps

We now demonstrate the difficulty of replacing the full-rank SVD computations required in projected subgradient steps over the spectrahedron, with their low-rank SVD counterparts when attempting to solve Problem (1). We prove that a projected subgradient step from a point arbitrarily close to a low-rank optimal solution — even one that satisfies strict complementarity (Definition 1), may result in a higher rank matrix. The problem on which we demonstrate this phenomena is a well known convex formulation of the *sparse PCA* problem [12].

Lemma 5 (failure of low-rank subgradient descent on sparse PCA). *Consider the problem*

$$\min_{\mathbf{X} \in \mathcal{S}_n} \{g(\mathbf{X}) := -\langle \mathbf{z}\mathbf{z}^\top + \mathbf{z}_\perp \mathbf{z}_\perp^\top, \mathbf{X} \rangle + \frac{1}{2k} \|\mathbf{X}\|_1\},$$

where $\mathbf{z} = (1/\sqrt{k}, \dots, 1/\sqrt{k}, 0, \dots, 0)^\top$ is supported on the first k entries, $\mathbf{z}_\perp = (0, \dots, 0, 1/\sqrt{n-k}, \dots, 1/\sqrt{n-k})^\top$ is supported on the last $n - k$ entries, and $k \leq n/4$. Then, $\mathbf{z}\mathbf{z}^\top$ is a rank-one optimal solution for which strict complementarity holds. However, for any $\eta < \frac{2}{3}$ and any $\mathbf{v} \in \mathbb{R}^n$ such that $\|\mathbf{v}\| = 1$, $\text{support}(\mathbf{v}) \subseteq \text{support}(\mathbf{z})$, and $\langle \mathbf{z}, \mathbf{v} \rangle^2 = 1 - \frac{1}{2} \|\mathbf{v}\mathbf{v}^\top - \mathbf{z}\mathbf{z}^\top\|_F^2 \geq 1 - \frac{1}{2k^2}$, it holds that

$$\text{rank}(\Pi_{\mathcal{S}_n}[\mathbf{v}\mathbf{v}^\top - \eta \mathbf{G}_{\mathbf{v}\mathbf{v}^\top}]) > 1,$$

where $\mathbf{G}_{\mathbf{v}\mathbf{v}^\top} = -\mathbf{z}\mathbf{z}^\top - \mathbf{z}_\perp \mathbf{z}_\perp^\top + \frac{1}{2k} \text{sign}(\mathbf{v}\mathbf{v}^\top) \in \partial g(\mathbf{v}\mathbf{v}^\top)$.

Note that the subgradient of the ℓ_1 -norm which we choose for the projected subgradient step simply corresponds to the sign function, which is arguably the most natural choice.

Proof. $\mathbf{z}\mathbf{z}^\top$ is a rank-one optimal solution for this problem since for the subgradient $k\mathbf{z}\mathbf{z}^\top + 2k\mathbf{z}_\perp \mathbf{z}_\perp^\top \in \partial(\|\mathbf{z}\mathbf{z}^\top\|_1)$ the first-order optimality condition holds. Indeed, for all $\mathbf{X} \in \mathcal{S}_n$

$$\langle \mathbf{X} - \mathbf{z}\mathbf{z}^\top, -\mathbf{z}\mathbf{z}^\top - \mathbf{z}_\perp \mathbf{z}_\perp^\top + \frac{1}{2}\mathbf{z}\mathbf{z}^\top + \mathbf{z}_\perp \mathbf{z}_\perp^\top \rangle = \langle \mathbf{X} - \mathbf{z}\mathbf{z}^\top, -\frac{1}{2}\mathbf{z}\mathbf{z}^\top \rangle \geq 0. \quad (4)$$

For the subgradient $-\frac{1}{2}\mathbf{z}\mathbf{z}^\top \in \partial g(\mathbf{z}\mathbf{z}^\top)$ there is a gap $\lambda_{n-1}(-\frac{1}{2}\mathbf{z}\mathbf{z}^\top) - \lambda_n(-\frac{1}{2}\mathbf{z}\mathbf{z}^\top) = \frac{1}{2} > 0$, and as we showed in (4) the first order optimality condition holds for $-\frac{1}{2}\mathbf{z}\mathbf{z}^\top$. Thus, by Lemma 2 the optimal solution $\mathbf{z}\mathbf{z}^\top$ satisfies standard strict complementarity.

We will show that the projection onto the spectrahedron of a subgradient step from $\mathbf{v}\mathbf{v}^\top$ with respect to the natural subgradient of the ℓ_1 -norm $\text{sign}(\mathbf{v}\mathbf{v}^\top) \in \partial(\|\mathbf{v}\mathbf{v}^\top\|_1)$ returns a rank-2 solution.

It holds that

$$1 - \frac{1}{2k^2} \leq \langle \mathbf{z}\mathbf{z}^\top, \mathbf{v}\mathbf{v}^\top \rangle = \frac{1}{2} \left(\|\mathbf{z}\mathbf{z}^\top\|_F^2 + \|\mathbf{v}\mathbf{v}^\top\|_F^2 - \|\mathbf{v}\mathbf{v}^\top - \mathbf{z}\mathbf{z}^\top\|_F^2 \right) = 1 - \frac{1}{2} \|\mathbf{v}\mathbf{v}^\top - \mathbf{z}\mathbf{z}^\top\|_F^2,$$

and equivalently

$$\sum_{i=1}^k \sum_{j=1}^k \left(\frac{1}{k} - (\mathbf{v}\mathbf{v}^\top)_{ij} \right)^2 = \|\mathbf{v}\mathbf{v}^\top - \mathbf{z}\mathbf{z}^\top\|_F^2 \leq \frac{1}{k^2}.$$

Therefore, for every $i, j \in \{1, \dots, k\}$ it holds that

$$\left| (\mathbf{v}\mathbf{v}^\top)_{ij} - \frac{1}{k} \right| \leq \frac{1}{k},$$

which implies that $0 \leq (\mathbf{v}\mathbf{v}^\top)_{ij} \leq \frac{2}{k}$. Therefore, $k\mathbf{z}\mathbf{z}^\top = \text{sign}(\mathbf{v}\mathbf{v}^\top) \in \partial(\|\mathbf{v}\mathbf{v}^\top\|_1)$.

Taking a projected subgradient step from $\mathbf{v}\mathbf{v}^\top$ with respect to the subgradient $-\mathbf{z}\mathbf{z}^\top - \mathbf{z}_\perp \mathbf{z}_\perp^\top + \frac{1}{2} \mathbf{z}\mathbf{z}^\top \in \partial g(\mathbf{v}\mathbf{v}^\top)$ has the form

$$\Pi_{S_n} \left[\mathbf{v}\mathbf{v}^\top - \eta \left(-\mathbf{z}\mathbf{z}^\top - \mathbf{z}_\perp \mathbf{z}_\perp^\top + \frac{1}{2} \mathbf{z}\mathbf{z}^\top \right) \right] = \Pi_{S_n} \left[\mathbf{v}\mathbf{v}^\top + \frac{\eta}{2} \mathbf{z}\mathbf{z}^\top + \eta \mathbf{z}_\perp \mathbf{z}_\perp^\top \right].$$

Since $\text{support}(\mathbf{v}) \subseteq \text{support}(\mathbf{z})$ it holds that $(\mathbf{v}\mathbf{v}^\top + \frac{\eta}{2} \mathbf{z}\mathbf{z}^\top) \perp \mathbf{z}_\perp \mathbf{z}_\perp^\top$. $\mathbf{v}\mathbf{v}^\top + \frac{\eta}{2} \mathbf{z}\mathbf{z}^\top$ is a rank-2 matrix and so we can denote the eigen-decomposition of $\mathbf{v}\mathbf{v}^\top + \frac{\eta}{2} \mathbf{z}\mathbf{z}^\top + \eta \mathbf{z}_\perp \mathbf{z}_\perp^\top$ as $\mathbf{v}\mathbf{v}^\top + \frac{\eta}{2} \mathbf{z}\mathbf{z}^\top + \eta \mathbf{z}_\perp \mathbf{z}_\perp^\top = \lambda_1 \mathbf{v}_1 \mathbf{v}_1^\top + \lambda_2 \mathbf{v}_2 \mathbf{v}_2^\top + \eta \mathbf{z}_\perp \mathbf{z}_\perp^\top$, where $\lambda_1 \geq \lambda_2$. Thus, invoking (2) to calculate the projection we need to find the scalar $\lambda \in \mathbb{R}$ for which the following holds.

$$\max\{\lambda_1 - \lambda, 0\} + \max\{\lambda_2 - \lambda, 0\} + \max\{\eta - \lambda, 0\} + \sum_{i=4}^n \max\{0 - \lambda, 0\} = 1.$$

λ_1 is the largest eigenvalue of $\mathbf{v}\mathbf{v}^\top + \frac{\eta}{2} \mathbf{z}\mathbf{z}^\top + \eta \mathbf{z}_\perp \mathbf{z}_\perp^\top$ since, under our assumption that $\eta < 2/3$, we have that

$$\lambda_1 \geq \frac{1}{2}(\lambda_1 + \lambda_2) = \frac{1}{2} \text{Tr}(\mathbf{v}\mathbf{v}^\top + \frac{\eta}{2} \mathbf{z}\mathbf{z}^\top) = \frac{1}{2} + \frac{\eta}{4} > \eta.$$

Therefore, $\lambda < \lambda_1 \leq \lambda_1 + \lambda_2 = \text{Tr}(\mathbf{v}\mathbf{v}^\top + \frac{\eta}{2} \mathbf{z}\mathbf{z}^\top) = 1 + \frac{\eta}{2}$.

In addition, $\max\{\lambda_2, \eta\} \geq \eta$. Therefore,

$$\lambda_1 - \max\{\lambda_2, \eta\} \leq \lambda_1 + \lambda_2 - \eta = 1 + \frac{\eta}{2} - \eta < 1,$$

and so we must have that $\lambda < \max\{\lambda_2, \eta\} \leq \lambda_1$.

This implies that both $\max\{\lambda_1 - \lambda, 0\} > 0$ and $\max\{\max\{\lambda_2, \eta\} - \lambda, 0\} > 0$. Thus, using (2) we conclude that

$$\text{rank} \left(\Pi_{S_n} \left[\mathbf{v}\mathbf{v}^\top - \eta \left(-\mathbf{z}\mathbf{z}^\top - \mathbf{z}_\perp \mathbf{z}_\perp^\top + \frac{1}{2} \mathbf{z}\mathbf{z}^\top \right) \right] \right) \geq 2.$$

□

3 From Nonsmooth to Saddle-Point Problems

To circumvent the difficulty demonstrated in Lemma 5 in incorporating low-rank SVDs into standard subgradient methods for solving Problem (1), we propose tackling the non-smooth problem with saddle-point methods.

We assume the nonsmooth Problem (1) can be written as a maximum of smooth functions, i.e., $g(\mathbf{X}) = \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}, \mathbf{y})$, where $\mathcal{K} \subset \mathbb{Y}$ is some compact and convex subset of the finite linear space over the reals \mathbb{Y} onto which it is efficient to compute Euclidean projections. We assume $f(\cdot, \mathbf{y})$ is convex for all $\mathbf{y} \in \mathcal{K}$ and $f(\mathbf{X}, \cdot)$ is concave for all $\mathbf{X} \in \mathcal{S}_n$. That is, we rewrite Problem (1) as the following equivalent saddle-point problem:

$$\min_{\mathbf{X} \in \mathcal{S}_n} \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}, \mathbf{y}). \quad (5)$$

Finding an optimal solution to problem (5) is equivalent to finding a saddle-point $(\mathbf{X}^*, \mathbf{y}^*) \in \mathcal{S}_n \times \mathcal{K}$ such that for all $\mathbf{X} \in \mathcal{S}_n$ and $\mathbf{y} \in \mathcal{K}$,

$$f(\mathbf{X}^*, \mathbf{y}) \leq f(\mathbf{X}^*, \mathbf{y}^*) \leq f(\mathbf{X}, \mathbf{y}^*).$$

We make a standard assumption that $f(\cdot, \cdot)$ is smooth with respect to all the components. That is, we assume there exist $\beta_X, \beta_y, \beta_{Xy}, \beta_{yX} \geq 0$ such that for any $\mathbf{X}, \tilde{\mathbf{X}} \in \mathcal{S}_n$ and $\mathbf{y}, \tilde{\mathbf{y}} \in \mathcal{K}$ the following four inequalities hold:

$$\begin{aligned} \|\nabla_{\mathbf{X}} f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{X}} f(\tilde{\mathbf{X}}, \mathbf{y})\|_F &\leq \beta_X \|\mathbf{X} - \tilde{\mathbf{X}}\|_F, \\ \|\nabla_{\mathbf{y}} f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{y}} f(\mathbf{X}, \tilde{\mathbf{y}})\|_2 &\leq \beta_y \|\mathbf{y} - \tilde{\mathbf{y}}\|_2, \\ \|\nabla_{\mathbf{X}} f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{X}} f(\mathbf{X}, \tilde{\mathbf{y}})\|_F &\leq \beta_{Xy} \|\mathbf{y} - \tilde{\mathbf{y}}\|_2, \\ \|\nabla_{\mathbf{y}} f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{y}} f(\tilde{\mathbf{X}}, \mathbf{y})\|_2 &\leq \beta_{yX} \|\mathbf{X} - \tilde{\mathbf{X}}\|_F, \end{aligned}$$

where $\nabla_{\mathbf{X}} f = \frac{\partial f}{\partial \mathbf{X}}$ and $\nabla_{\mathbf{y}} f = \frac{\partial f}{\partial \mathbf{y}}$.

We denote by β the full Lipschitz parameter of the gradient, that is for any $\mathbf{X}, \tilde{\mathbf{X}} \in \mathcal{S}_n$ and $\mathbf{y}, \tilde{\mathbf{y}} \in \mathcal{K}$,

$$\|(\nabla_{\mathbf{X}} f(\mathbf{X}, \mathbf{y}), -\nabla_{\mathbf{y}} f(\mathbf{X}, \mathbf{y})) - (\nabla_{\mathbf{X}} f(\tilde{\mathbf{X}}, \tilde{\mathbf{y}}), -\nabla_{\mathbf{y}} f(\tilde{\mathbf{X}}, \tilde{\mathbf{y}}))\| \leq \beta \|(\mathbf{X}, \mathbf{y}) - (\tilde{\mathbf{X}}, \tilde{\mathbf{y}})\|,$$

where $\|\cdot\|$ denotes the Euclidean norm over the product space $\mathcal{S}_n \times \mathbb{Y}$.

To establish the relationship between β and $\beta_X, \beta_y, \beta_{Xy}, \beta_{yX}$, we can see that for all $\mathbf{X}, \tilde{\mathbf{X}} \in \mathcal{S}_n$ and all $\mathbf{y}, \tilde{\mathbf{y}} \in \mathcal{K}$

$$\begin{aligned} &\|(\nabla_{\mathbf{X}} f(\mathbf{X}, \mathbf{y}), -\nabla_{\mathbf{y}} f(\mathbf{X}, \mathbf{y})) - (\nabla_{\mathbf{X}} f(\tilde{\mathbf{X}}, \tilde{\mathbf{y}}), -\nabla_{\mathbf{y}} f(\tilde{\mathbf{X}}, \tilde{\mathbf{y}}))\|^2 \\ &= \|\nabla_{\mathbf{X}} f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{X}} f(\tilde{\mathbf{X}}, \mathbf{y})\|_F^2 + \|\nabla_{\mathbf{y}} f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{y}} f(\tilde{\mathbf{X}}, \mathbf{y})\|_2^2 \\ &\leq 2\|\nabla_{\mathbf{X}} f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{X}} f(\tilde{\mathbf{X}}, \mathbf{y})\|_F^2 + 2\|\nabla_{\mathbf{X}} f(\tilde{\mathbf{X}}, \mathbf{y}) - \nabla_{\mathbf{X}} f(\tilde{\mathbf{X}}, \tilde{\mathbf{y}})\|_F^2 \\ &\quad + 2\|\nabla_{\mathbf{y}} f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{y}} f(\tilde{\mathbf{X}}, \mathbf{y})\|_2^2 + 2\|\nabla_{\mathbf{y}} f(\tilde{\mathbf{X}}, \mathbf{y}) - \nabla_{\mathbf{y}} f(\tilde{\mathbf{X}}, \tilde{\mathbf{y}})\|_2^2 \\ &\leq 2(\beta_X^2 + \beta_{yX}^2) \|\mathbf{X} - \tilde{\mathbf{X}}\|_F^2 + 2(\beta_y^2 + \beta_{Xy}^2) \|\mathbf{y} - \tilde{\mathbf{y}}\|_2^2 \\ &\leq 2 \max\{\beta_X^2 + \beta_{yX}^2, \beta_y^2 + \beta_{Xy}^2\} \|(\mathbf{X}, \mathbf{y}) - (\tilde{\mathbf{X}}, \tilde{\mathbf{y}})\|^2. \end{aligned}$$

Therefore, $\beta = \sqrt{2} \max \left\{ \sqrt{\beta_X^2 + \beta_{yX}^2}, \sqrt{\beta_y^2 + \beta_{Xy}^2} \right\}$.

The following lemma highlights a connection between the gradient of a saddle-point of (5) and subgradients of an optimal solution to (1) for which the first order optimality condition holds. One of the connections we will be interested in, is that GSC for Problem (1) implies GSC (with the same parameters) for Problem (5). However, to prove this specific connection we require an additional structural assumption on the objective function $g(\cdot)$. We note that this assumption holds for all applications mentioned in this paper.

Assumption 2. $g(\mathbf{X})$ is of the form $g(\mathbf{X}) = h(\mathbf{X}) + \max_{\mathbf{y} \in \mathcal{K}} \mathbf{y}^\top (\mathcal{A}(\mathbf{X}) - \mathbf{b})$, where $h(\cdot)$ is smooth and convex, and \mathcal{A} is a linear map.

Lemma 6. If $(\mathbf{X}^*, \mathbf{y}^*)$ is a saddle-point of Problem (5) then \mathbf{X}^* is an optimal solution to Problem (1), $\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*) \in \partial g(\mathbf{X}^*)$, and for all $\mathbf{X} \in \mathcal{S}_n$ it holds that $\langle \mathbf{X} - \mathbf{X}^*, \nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*) \rangle \geq 0$. Conversely, under Assumption 2, if \mathbf{X}^* is an optimal solution to Problem (1), and $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ which satisfies $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$, then there exists $\mathbf{y}^* \in \arg \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}^*, \mathbf{y})$ such that $(\mathbf{X}^*, \mathbf{y}^*)$ is a saddle-point of Problem (5), and $\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*) = \mathbf{G}^*$.

The proof is given in Appendix B. The connection between the gradient of an optimal solution to the saddle-point problem and a subgradient of a corresponding optimal solution in the equivalent nonsmooth problem established in Lemma 6, naturally leads to the formulation of the following generalized strict complementarity assumption for saddle-point problems.

Assumption 3 (generalized strict complementarity for saddle-points). We say a saddle-point $(\mathbf{X}^*, \mathbf{y}^*) \in \mathcal{S}_n \times \mathcal{K}$ of Problem (5) with $\text{rank}(\mathbf{X}^*) = r^*$ satisfies the generalized strict complementarity assumption with parameters $r \geq r^*, \delta > 0$, if $\lambda_{n-r}(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) \geq \delta$.

Remark 1. Note that under Assumption 2, due to Lemma 6, GSC with parameters r, δ for some optimal solution \mathbf{X}^* to Problem (1) implies GSC with parameters r, δ to a corresponding saddle-point $(\mathbf{X}^*, \mathbf{y}^*)$ of Problem (5). Nevertheless, Assumption 2 is not necessary for proving our convergence results for Problem (5), which are directly stated in terms of Assumption 3.

4 Projected Extragradient Method with Low-Rank Projections

In this section we formally state and prove our main result: the projected extragradient method for the saddle-point Problem (5), when initialized in the proximity of a saddle-point which satisfies GSC (Assumption 3), converges with its original $O(1/t)$ rate while requiring only two low-rank SVD computations per iteration.

Algorithm 1 Projected extragradient method for saddle-point problems (see also [24, 31])

Input: sequence of step-sizes $\{\eta_t\}_{t \geq 1}$
Initialization: $(\mathbf{X}_1, \mathbf{y}_1) \in \mathcal{S}_n \times \mathcal{K}$
for $t = 1, 2, \dots$ **do**
 $\mathbf{Z}_{t+1} = \Pi_{\mathcal{S}_n}[\mathbf{X}_t - \eta_t \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t)]$
 $\mathbf{w}_{t+1} = \Pi_{\mathcal{K}}[\mathbf{y}_t + \eta_t \nabla_{\mathbf{y}} f(\mathbf{X}_t, \mathbf{y}_t)]$
 $\mathbf{X}_{t+1} = \Pi_{\mathcal{S}_n}[\mathbf{X}_t - \eta_t \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1})]$
 $\mathbf{y}_{t+1} = \Pi_{\mathcal{K}}[\mathbf{y}_t + \eta_t \nabla_{\mathbf{y}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1})]$
end for

First, in the following lemma we state the standard convergence result of the projected extragradient method, which is a well known result⁴. For completeness we include the proof in Appendix C.

Lemma 7. *Let $\{(\mathbf{X}_t, \mathbf{y}_t)\}_{t \geq 1}$ and $\{(\mathbf{Z}_t, \mathbf{w}_t)\}_{t \geq 2}$ be the sequences generated by Algorithm 1 with a fixed step-size $\eta_t = \eta \leq \min \left\{ \frac{1}{\beta_X + \beta_{Xy}}, \frac{1}{\beta_y + \beta_{yX}}, \frac{1}{\beta_X + \beta_{yX}}, \frac{1}{\beta_y + \beta_{Xy}} \right\}$ then*

$$\max_{\mathbf{y} \in \mathcal{K}} f \left(\frac{1}{T} \sum_{t=1}^T \mathbf{Z}_{t+1}, \mathbf{y} \right) - \min_{\mathbf{X} \in \mathcal{S}_n} f \left(\mathbf{X}, \frac{1}{T} \sum_{t=1}^T \mathbf{w}_{t+1} \right) \leq \frac{D^2}{2\eta T},$$

where $D := \sup_{(\mathbf{X}, \mathbf{y}), (\tilde{\mathbf{X}}, \tilde{\mathbf{y}}) \in \mathcal{S}_n \times \mathcal{K}} \|(\mathbf{X}, \mathbf{y}) - (\tilde{\mathbf{X}}, \tilde{\mathbf{y}})\|$.

We can now state our main theorem.

Theorem 1 (main theorem). *Fix an optimal solution $(\mathbf{X}^*, \mathbf{y}^*) \in \mathcal{S}_n \times \mathcal{K}$ to Problem (5). Let \tilde{r} denote the multiplicity of $\lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*))$ and for any $r \geq \tilde{r}$ define $\delta(r) = \lambda_{n-r}(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*) - \lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*))$. Let $\{(\mathbf{X}_t, \mathbf{y}_t)\}_{t \geq 1}$ and $\{(\mathbf{Z}_t, \mathbf{w}_t)\}_{t \geq 2}$ be the sequences of iterates generated by Algorithm 1 with a fixed step-size*

$$\eta = \min \left\{ \frac{1}{2\sqrt{\beta_X^2 + \beta_{yX}^2}}, \frac{1}{2\sqrt{\beta_y^2 + \beta_{Xy}^2}}, \frac{1}{\beta_X + \beta_{Xy}}, \frac{1}{\beta_y + \beta_{yX}} \right\}.$$

Assume the initialization $(\mathbf{X}_1, \mathbf{y}_1)$ satisfies $\|(\mathbf{X}_1, \mathbf{y}_1) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq R_0(r)$, where

$$R_0(r) := \frac{\eta}{(1 + \sqrt{2})(1 + (2 + \sqrt{2})\eta \max\{\beta_X, \beta_{Xy}\})} \max \left\{ \frac{\sqrt{\tilde{r}}\delta(r - \tilde{r} + 1)}{2}, \frac{\delta(r)}{(1 + 1/\sqrt{\tilde{r}})} \right\}.$$

Then, for all $t \geq 1$, the projections $\Pi_{\mathcal{S}_n}[\mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t)]$ and $\Pi_{\mathcal{S}_n}[\mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1})]$ can be replaced with their rank- r truncated counterparts (see (3)) without changing the sequences $\{(\mathbf{X}_t, \mathbf{y}_t)\}_{t \geq 1}$ and $\{(\mathbf{Z}_t, \mathbf{w}_t)\}_{t \geq 2}$, and for any $T \geq 0$ it holds that

$$\begin{aligned} & \max_{\mathbf{y} \in \mathcal{K}} f \left(\frac{1}{T} \sum_{t=1}^T \mathbf{Z}_{t+1}, \mathbf{y} \right) - \min_{\mathbf{X} \in \mathcal{S}_n} f \left(\mathbf{X}, \frac{1}{T} \sum_{t=1}^T \mathbf{w}_{t+1} \right) \\ & \leq \frac{D^2 \max \left\{ \sqrt{\beta_X^2 + \beta_{yX}^2}, \sqrt{\beta_y^2 + \beta_{Xy}^2}, \frac{1}{2}(\beta_X + \beta_{Xy}), \frac{1}{2}(\beta_y + \beta_{yX}) \right\}}{T}, \end{aligned}$$

where $D := \sup_{(\mathbf{X}, \mathbf{y}), (\mathbf{Z}, \mathbf{w}) \in \mathcal{S}_n \times \mathcal{K}} \|(\mathbf{X}, \mathbf{y}) - (\mathbf{Z}, \mathbf{w})\|$.

⁴[7, 31] prove this result with respect to the ergodic series. A small adjustment of the proof proves the same with respect to the minimum and maximum iterates.

Remark 2. Note that Theorem 1 implies that if standard strict complementarity holds for Problem (5), that is Assumption 3 holds with $r = r^* = \text{rank}(\mathbf{X}^*)$ and some $\delta > 0$, then only rank- r^* SVDs are required so that Algorithm 1 converges with the guaranteed convergence rate of $O(1/t)$, when initialized with a “warm-start”. Furthermore, by using SVDs of rank $r > r^*$, with moderately higher values of r , we can increase the radius of the ball in which Algorithm 1 needs to be initialized quite significantly.

To prove Theorem 1 we first prove two technical lemmas. We begin by proving that the iterates of Algorithm 1 always remain inside a ball of a certain radius around an optimal solution.

Lemma 8. Let $\{(\mathbf{X}_t, \mathbf{y}_t)\}_{t \geq 1}$ and $\{(\mathbf{Z}_t, \mathbf{w}_t)\}_{t \geq 2}$ be the sequences generated by Algorithm 1 with a step-size $\eta_t \leq \frac{1}{\beta}$, and let $(\mathbf{X}^*, \mathbf{y}^*)$ be some optimal solution to Problem (5). Then for all $t \geq 1$ it holds that

$$\begin{aligned} \|(\mathbf{X}_{t+1}, \mathbf{y}_{t+1}) - (\mathbf{X}^*, \mathbf{y}^*)\| &\leq \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\|, \\ \|(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) - (\mathbf{X}^*, \mathbf{y}^*)\| &\leq \left(1 + \frac{1}{\sqrt{1 - \eta_t^2 \beta^2}}\right) \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\|. \end{aligned}$$

Proof. A known inequality of the EG algorithm (see for example Lemma 12.1.10 in [15]) is

$$\begin{aligned} \|(\mathbf{X}_{t+1}, \mathbf{y}_{t+1}) - (\mathbf{X}^*, \mathbf{y}^*)\|^2 &\leq \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\|^2 - (1 - \eta_t^2 \beta^2) \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{Z}_{t+1}, \mathbf{w}_{t+1})\|^2. \end{aligned} \quad (6)$$

Since $\eta_t^2 \beta^2 \leq 1$ it follows that

$$\|(\mathbf{X}_{t+1}, \mathbf{y}_{t+1}) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\|.$$

In addition, using (6)

$$\|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{Z}_{t+1}, \mathbf{w}_{t+1})\| \leq \sqrt{(1 - \eta_t^2 \beta^2)^{-1}} \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\|.$$

Therefore,

$$\begin{aligned} \|(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) - (\mathbf{X}^*, \mathbf{y}^*)\| &\leq \|(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) - (\mathbf{X}_t, \mathbf{y}_t)\| + \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\| \\ &\leq \sqrt{(1 - \eta_t^2 \beta^2)^{-1}} \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\| \\ &\quad + \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\| \\ &= \left(1 + \frac{1}{\sqrt{1 - \eta_t^2 \beta^2}}\right) \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\|. \end{aligned}$$

□

We now prove that when close enough to a low-rank saddle-point of Problem (5), under an assumption of an eigen-gap in the gradient of the saddle-point, both projections onto the spectrahedron that are necessary in each iteration of Algorithm 1, result in low-rank matrices.

Lemma 9. Let $(\mathbf{X}^*, \mathbf{y}^*)$ be an optimal solution to Problem (5). Let \tilde{r} denote the multiplicity of $\lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*))$ and for any $r \geq \tilde{r}$ denote $\delta(r) := \lambda_{n-r}(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*))$. Then, for any $\eta \geq 0$ and $(\mathbf{X}, \mathbf{y}) \in \mathcal{S}_n \times \mathcal{K}$, if

$$\begin{aligned} & \|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\| \\ & \leq \frac{\eta}{1 + \sqrt{2}\eta \max\{\beta_X, \beta_{Xy}\} \left(1 + \frac{1}{\sqrt{1-\eta^2\beta^2}}\right)} \max \left\{ \frac{\sqrt{\tilde{r}}\delta(r - \tilde{r} + 1)}{2}, \frac{\delta(r)}{(1 + 1/\sqrt{\tilde{r}})} \right\} \end{aligned}$$

then $\text{rank}(\Pi_{\mathcal{S}_n}[\mathbf{X} - \eta \nabla_{\mathbf{X}} f(\mathbf{X}, \mathbf{y})]) \leq r$ and $\text{rank}(\Pi_{\mathcal{S}_n}[\mathbf{X} - \eta \nabla_{\mathbf{X}} f(\mathbf{Z}_+, \mathbf{w}_+)]) \leq r$ where $\mathbf{Z}_+ = \Pi_{\mathcal{S}_n}[\mathbf{X} - \eta \nabla_{\mathbf{X}} f(\mathbf{X}, \mathbf{y})]$ and $\mathbf{w}_+ = \Pi_{\mathcal{K}}[\mathbf{y} - \eta \nabla_{\mathbf{y}} f(\mathbf{X}, \mathbf{y})]$.

Proof. Denote $\mathbf{P}^* = \mathbf{X}^* - \eta \nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)$. By Lemma 6, $\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)$ is a subgradient of the corresponding nonsmooth objective $g(\mathbf{X}) = \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}, \mathbf{y})$ at the point \mathbf{X}^* . Moreover, this subgradient also satisfies the first-order optimality condition. Hence, invoking Lemma 3 with this subgradient we have that

$$\begin{aligned} \forall i \leq \text{rank}(\mathbf{X}^*) : \lambda_i(\mathbf{P}^*) &= \lambda_i(\mathbf{X}^*) - \eta \lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)); \\ \forall i > \text{rank}(\mathbf{X}^*) : \lambda_i(\mathbf{P}^*) &= -\eta \lambda_{n-i+1}(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)). \end{aligned} \quad (7)$$

Therefore, using (7) and the fact that $\lambda_{n-i+1}(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) = \lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*))$ for all $i \leq \tilde{r}$ we have,

$$\begin{aligned} \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{P}^*) &= \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{X}^* - \eta \nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) = \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{X}^*) - \eta \sum_{i=1}^{\tilde{r}} \lambda_{n-i+1}(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) \\ &= \sum_{i=1}^{\text{rank}(\mathbf{X}^*)} \lambda_i(\mathbf{X}^*) - \eta \sum_{i=1}^{\tilde{r}} \lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) = 1 - \eta \tilde{r} \lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)). \end{aligned} \quad (8)$$

Let $\mathbf{P} \in \mathbb{S}^n$. From the structure of the Euclidean projection onto the spectrahedron (see Eq. (2)), it follows that a sufficient condition so that $\text{rank}(\Pi_{\mathcal{S}_n}[\mathbf{P}]) \leq r$ is that $\sum_{i=1}^r \lambda_i(\mathbf{P}) - r \lambda_{r+1}(\mathbf{P}) \geq 1$. We will bound the LHS of this inequality.

First, it holds that

$$\begin{aligned} \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{P}) &\stackrel{(a)}{\geq} \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{P}^*) - \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{P}^* - \mathbf{P}) \geq \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{P}^*) - \sqrt{\tilde{r} \sum_{i=1}^{\tilde{r}} \lambda_i^2(\mathbf{P} - \mathbf{P}^*)} \\ &\geq \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{P}^*) - \sqrt{\tilde{r} \sum_{i=1}^n \lambda_i^2(\mathbf{P} - \mathbf{P}^*)} \geq \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{P}^*) - \sqrt{\tilde{r}} \|\mathbf{P} - \mathbf{P}^*\|_F \\ &\stackrel{(b)}{\geq} 1 - \eta \tilde{r} \lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) - \sqrt{\tilde{r}} \|\mathbf{P} - \mathbf{P}^*\|_F, \end{aligned} \quad (9)$$

where (a) holds from Ky Fan's inequality for eigenvalues and (b) holds from (8).

Now, for any $r \geq \tilde{r}$ using Weyl's inequality and (7)

$$\begin{aligned} \lambda_{r+1}(\mathbf{P}) &\leq \lambda_{r+1}(\mathbf{P}^*) + \lambda_1(\mathbf{P} - \mathbf{P}^*) \leq \lambda_{r+1}(\mathbf{P}^*) + \|\mathbf{P} - \mathbf{P}^*\|_F \\ &= -\eta \lambda_{n-r}(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) + \|\mathbf{P} - \mathbf{P}^*\|_F. \end{aligned} \quad (10)$$

Thus, combining (9) and (10) we obtain

$$\sum_{i=1}^r \lambda_i(\mathbf{P}) - r\lambda_{r+1}(\mathbf{P}) \geq \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{P}) - \tilde{r}\lambda_{r+1}(\mathbf{P}) \geq 1 + \eta\tilde{r}\delta(r) - (\tilde{r} + \sqrt{\tilde{r}})\|\mathbf{P} - \mathbf{P}^*\|_F. \quad (11)$$

Alternatively, if $r \geq 2\tilde{r} - 1$ then using the general Weyl inequality and (7) we obtain

$$\begin{aligned} \lambda_{r+1}(\mathbf{P}) &\leq \lambda_{r-\tilde{r}+2}(\mathbf{P}^*) + \lambda_{\tilde{r}}(\mathbf{P} - \mathbf{P}^*) = \lambda_{r-\tilde{r}+2}(\mathbf{P}^*) + \sqrt{\lambda_{\tilde{r}}^2(\mathbf{P} - \mathbf{P}^*)} \\ &\leq \lambda_{r-\tilde{r}+2}(\mathbf{P}^*) + \frac{1}{\sqrt{\tilde{r}}}\|\mathbf{P} - \mathbf{P}^*\|_F = -\eta\lambda_{n-r+\tilde{r}-1}(\nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*)) + \frac{1}{\sqrt{\tilde{r}}}\|\mathbf{P} - \mathbf{P}^*\|_F. \end{aligned} \quad (12)$$

Thus, combining (9) and (12) we obtain

$$\sum_{i=1}^r \lambda_i(\mathbf{P}) - r\lambda_{r+1}(\mathbf{P}) \geq \sum_{i=1}^{\tilde{r}} \lambda_i(\mathbf{P}) - \tilde{r}\lambda_{r+1}(\mathbf{P}) \geq 1 + \eta\tilde{r}\delta(r - \tilde{r} + 1) - 2\sqrt{\tilde{r}}\|\mathbf{P} - \mathbf{P}^*\|_F. \quad (13)$$

Now we are left with bounding $\|\mathbf{P} - \mathbf{P}^*\|_F$. Note that by the smoothness of f , for any $(\mathbf{X}, \mathbf{y}) \in \mathcal{S}_n \times \mathcal{K}$ it holds that

$$\begin{aligned} &\|\nabla_{\mathbf{X}}f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*)\|_F \\ &\leq \|\nabla_{\mathbf{X}}f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y})\|_F + \|\nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}) - \nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*)\|_F \\ &\leq \beta_X\|\mathbf{X} - \mathbf{X}^*\|_F + \beta_{Xy}\|\mathbf{y} - \mathbf{y}^*\|_2. \end{aligned} \quad (14)$$

Taking $\mathbf{P} = \mathbf{X} - \eta\nabla_{\mathbf{X}}f(\mathbf{X}, \mathbf{y})$ we get

$$\begin{aligned} \|\mathbf{P} - \mathbf{P}^*\|_F &= \|\mathbf{X} - \eta\nabla_{\mathbf{X}}f(\mathbf{X}, \mathbf{y}) - \mathbf{X}^* + \eta\nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*)\|_F \\ &\leq \|\mathbf{X} - \mathbf{X}^*\|_F + \eta\|\nabla_{\mathbf{X}}f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*)\|_F \\ &\leq \|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\| + \eta\|\nabla_{\mathbf{X}}f(\mathbf{X}, \mathbf{y}) - \nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*)\|_F \\ &\leq \|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\| + \eta\beta_X\|\mathbf{X} - \mathbf{X}^*\|_F + \eta\beta_{Xy}\|\mathbf{y} - \mathbf{y}^*\|_2, \end{aligned}$$

where the last inequality holds from (14).

For any $a, b \geq 0$ it holds that

$$\begin{aligned} a\|\mathbf{X} - \mathbf{X}^*\|_F + b\|\mathbf{y} - \mathbf{y}^*\|_2 &\leq \max\{a, b\}(\|\mathbf{X} - \mathbf{X}^*\|_F + \|\mathbf{y} - \mathbf{y}^*\|_2) \\ &\leq \sqrt{2}\max\{a, b\}\|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\|. \end{aligned}$$

Thus, by taking $a = \eta\beta_X$ and $b = \eta\beta_{Xy}$ we obtain

$$\|\mathbf{P} - \mathbf{P}^*\|_F \leq \left(1 + \sqrt{2}\eta\max\{\beta_X, \beta_{Xy}\}\right)\|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\|. \quad (15)$$

Therefore, plugging (15) into (11) we obtain that the condition $\sum_{i=1}^r \lambda_i(\mathbf{P}) - r\lambda_{r+1}(\mathbf{P}) \geq 1$ holds if

$$\|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq \frac{\eta\delta(r)}{(1 + 1/\sqrt{\tilde{r}})(1 + \sqrt{2}\eta\max\{\beta_X, \beta_{Xy}\})}.$$

Alternatively, plugging (15) into (13) we obtain that if $r \geq 2\tilde{r} - 1$ then the condition $\sum_{i=1}^r \lambda_i(\mathbf{P}) - r\lambda_{r+1}(\mathbf{P}) \geq 1$ holds if

$$\|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq \frac{\eta\sqrt{\tilde{r}}\delta(r - \tilde{r} + 1)}{2(1 + \sqrt{2}\eta \max\{\beta_X, \beta_{Xy}\})}.$$

Note that $\delta(r - \tilde{r} + 1) > 0$ only if $r \geq 2\tilde{r} - 1$. Therefore, we can combine the last two inequalities to conclude that for any $r \geq \tilde{r}$ if

$$\|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq \frac{\eta}{1 + \sqrt{2}\eta \max\{\beta_X, \beta_{Xy}\}} \max \left\{ \frac{\sqrt{\tilde{r}}\delta(r - \tilde{r} + 1)}{2}, \frac{\delta(r)}{(1 + 1/\sqrt{\tilde{r}})} \right\} \quad (16)$$

then $\text{rank}(\Pi_{\mathcal{S}_n}[\mathbf{X} - \nabla_{\mathbf{X}}f(\mathbf{X}, \mathbf{y})]) \leq r$.

Similarly, taking $\mathbf{P} = \mathbf{X} - \eta\nabla_{\mathbf{X}}f(\mathbf{Z}_+, \mathbf{w}_+)$ we get

$$\begin{aligned} \|\mathbf{P} - \mathbf{P}^*\|_F &= \|\mathbf{X} - \eta\nabla_{\mathbf{X}}f(\mathbf{Z}_+, \mathbf{w}_+) - \mathbf{X}^* + \eta\nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*)\|_F \\ &\leq \|\mathbf{X} - \mathbf{X}^*\|_F + \eta\|\nabla_{\mathbf{X}}f(\mathbf{Z}_+, \mathbf{w}_+) - \nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*)\|_F \\ &\leq \|\mathbf{X} - \mathbf{X}^*\|_F + \eta\beta_X\|\mathbf{Z}_+ - \mathbf{X}^*\|_F + \eta\beta_{Xy}\|\mathbf{w}_+ - \mathbf{y}^*\|_2, \end{aligned}$$

where the last inequality holds from (14).

For any $a, b, c \geq 0$ it holds that

$$\begin{aligned} &a\|\mathbf{X} - \mathbf{X}^*\|_F + b\|\mathbf{Z}_+ - \mathbf{X}^*\|_F + c\|\mathbf{w}_+ - \mathbf{y}^*\|_2 \\ &\leq a\|\mathbf{X} - \mathbf{X}^*\|_F + \max\{b, c\}(\|\mathbf{Z}_+ - \mathbf{X}^*\|_F + \|\mathbf{w}_+ - \mathbf{y}^*\|_2) \\ &\leq a\|\mathbf{X} - \mathbf{X}^*\|_F + \sqrt{2}\max\{b, c\}\|(\mathbf{Z}_+, \mathbf{w}_+) - (\mathbf{X}^*, \mathbf{y}^*)\| \\ &\leq a\|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\| + \sqrt{2}\max\{b, c\}\|(\mathbf{Z}_+, \mathbf{w}_+) - (\mathbf{X}^*, \mathbf{y}^*)\| \\ &\leq \left(a + \sqrt{2}\max\{b, c\} \left(1 + \frac{1}{\sqrt{1 - \eta^2\beta^2}} \right) \right) \|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\|, \end{aligned}$$

where the second to last inequality holds from Lemma 8.

Thus, by taking $a = 1$, $b = \eta\beta_X$, and $c = \eta\beta_{Xy}$ we obtain

$$\|\mathbf{P} - \mathbf{P}^*\|_F \leq \left(1 + \sqrt{2}\eta \max\{\beta_X, \beta_{Xy}\} \left(1 + \frac{1}{\sqrt{1 - \eta^2\beta^2}} \right) \right) \|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\|. \quad (17)$$

Therefore, plugging (17) into (11) we obtain that the condition $\sum_{i=1}^r \lambda_i(\mathbf{P}) - r\lambda_{r+1}(\mathbf{P}) \geq 1$ holds if

$$\|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq \frac{\eta\delta(r)}{(1 + 1/\sqrt{r^*}) \left(1 + \sqrt{2}\eta \max\{\beta_X, \beta_{Xy}\} \left(1 + \frac{1}{\sqrt{1 - \eta^2\beta^2}} \right) \right)}.$$

Alternatively, plugging (17) into (13) we obtain that if $r \geq 2\tilde{r} - 1$ then the condition $\sum_{i=1}^r \lambda_i(\mathbf{P}) - r\lambda_{r+1}(\mathbf{P}) \geq 1$ holds if

$$\|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq \frac{\eta\sqrt{\tilde{r}}\delta(r - \tilde{r} + 1)}{1 + \sqrt{2}\eta \max\{\beta_X, \beta_{Xy}\} \left(1 + \frac{1}{\sqrt{1 - \eta^2\beta^2}} \right)}.$$

Note that $\delta(r - \tilde{r} + 1) > 0$ only if $r \geq 2\tilde{r} - 1$. Therefore, we can combine the last two inequalities to conclude that for any $r \geq \tilde{r}$ if

$$\begin{aligned} & \|(\mathbf{X}, \mathbf{y}) - (\mathbf{X}^*, \mathbf{y}^*)\| \\ & \leq \frac{\eta}{1 + \sqrt{2}\eta \max\{\beta_X, \beta_{Xy}\} \left(1 + \frac{1}{\sqrt{1 - \eta^2 \beta^2}}\right)} \max \left\{ \frac{\sqrt{\tilde{r}}\delta(r - \tilde{r} + 1)}{2}, \frac{\delta(r)}{(1 + 1/\sqrt{\tilde{r}})} \right\} \quad (18) \end{aligned}$$

then $\text{rank}(\Pi_{\mathcal{S}_n}[\mathbf{X} - \nabla_{\mathbf{X}} f(\mathbf{Z}_+, \mathbf{w}_+)]) \leq r$.

Taking the minimum between (16) and (18) gives us the bound on the radius in the lemma. □

Now we can prove Theorem 1.

Proof of Theorem 1. We will prove by induction that for all $t \geq 1$ it holds that $\|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq R_0(r)$ and $\|(\mathbf{Z}_t, \mathbf{w}_t) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq (1 + \sqrt{2}) R_0(r)$, thus implying through Lemma 9 that all projections $\Pi_{\mathcal{S}_n}[\mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t)]$ and $\Pi_{\mathcal{S}_n}[\mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1})]$ can be replaced with their rank- r truncated counterparts given in (3), without any change to the result.

The initialization $\|(\mathbf{X}_1, \mathbf{y}_1) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq R_0(r)$ holds trivially. Now, by Lemma 8, using recursion, we have that for all $t \geq 1$,

$$\|(\mathbf{X}_{t+1}, \mathbf{y}_{t+1}) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq \dots \leq \|(\mathbf{X}_1, \mathbf{y}_1) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq R_0(r),$$

and for $\beta = \sqrt{2} \max \left\{ \sqrt{\beta_X^2 + \beta_{yX}^2}, \sqrt{\beta_y^2 + \beta_{Xy}^2} \right\}$ we have that,

$$\begin{aligned} \|(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) - (\mathbf{X}^*, \mathbf{y}^*)\| & \leq \left(1 + \frac{1}{\sqrt{1 - \eta_t^2 \beta^2}}\right) \|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}^*, \mathbf{y}^*)\| \\ & \leq \left(1 + \frac{1}{\sqrt{1 - \eta_t^2 \beta^2}}\right) \|(\mathbf{X}_1, \mathbf{y}_1) - (\mathbf{X}^*, \mathbf{y}^*)\| \\ & \leq (1 + \sqrt{2}) \|(\mathbf{X}_1, \mathbf{y}_1) - (\mathbf{X}^*, \mathbf{y}^*)\| \leq (1 + \sqrt{2}) R_0(r). \end{aligned}$$

Therefore, under the assumptions of the theorem, Algorithm 1 can be run using only rank- r truncated projections, while maintaining its original convergence rate stated in Lemma 7. □

Remark 3. A downside of considering the saddle-point formulation (5) when attempting to solve Problem (1) that arises from Theorem 1, is that not only do we need a “warm-start” initialization for the original primal matrix variable \mathbf{X} , in the saddle-point formulation we need a “warm-start” for the primal-dual pair (\mathbf{X}, \mathbf{y}) . Nevertheless, as we demonstrate extensively in Section 5, it seems that very simple initialization schemes work very well in practice.

4.1 Back to nonsmooth problems

Corollary 1. Fix an optimal solution $\mathbf{X}^* \in \mathcal{S}_n$ to Problem (1) and assume Assumption 2 holds. Let $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ which satisfies that $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$. Let \tilde{r} denote the multiplicity of $\lambda_n(\mathbf{G}^*)$ and for any $r \geq \tilde{r}$ define $\delta(r) := \lambda_{n-r}(\mathbf{G}^*) - \lambda_n(\mathbf{G}^*)$. Define f as in Problem (5) and let $\{(\mathbf{X}_t, \mathbf{y}_t)\}_{t \geq 1}$ and $\{(\mathbf{Z}_t, \mathbf{w}_t)\}_{t \geq 2}$ be the sequences of iterates generated by Algorithm 1 where η and $R_0(r)$ are as defined in Theorem 1. Then, for all $t \geq 1$ the projections $\Pi_{\mathcal{S}_n}[\mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t)]$ and $\Pi_{\mathcal{S}_n}[\mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1})]$ can be replaced with rank- r truncated projections (3) without changing the sequences $\{(\mathbf{X}_t, \mathbf{y}_t)\}_{t \geq 1}$ and $\{(\mathbf{Z}_t, \mathbf{w}_t)\}_{t \geq 2}$, and for any $T \geq 0$ it holds that

$$\begin{aligned} & g\left(\frac{1}{T} \sum_{t=1}^T \mathbf{Z}_{t+1}\right) - g(\mathbf{X}^*) \\ & \leq \frac{D^2 \max \left\{ \sqrt{\beta_X^2 + \beta_{yX}^2}, \sqrt{\beta_y^2 + \beta_{Xy}^2}, \frac{1}{2}(\beta_X + \beta_{Xy}), \frac{1}{2}(\beta_y + \beta_{yX}) \right\}}{T}, \end{aligned}$$

where $D := \sup_{(\mathbf{X}, \mathbf{y}), (\mathbf{Z}, \mathbf{w}) \in \mathcal{S}_n \times \mathcal{K}} \|(\mathbf{X}, \mathbf{y}) - (\mathbf{Z}, \mathbf{w})\|$.

Proof. Since Assumption 2 holds, invoking Lemma 6 we obtain that there exists a point $\mathbf{y}^* \in \arg \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}^*, \mathbf{y})$ such that $(\mathbf{X}^*, \mathbf{y}^*)$ is a saddle-point of Problem (5), and $\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*) = \mathbf{G}^*$. Therefore, the assumptions of Theorem 1 hold and so by Theorem 1 we get that

$$\begin{aligned} & \max_{\mathbf{y} \in \mathcal{K}} f\left(\frac{1}{T} \sum_{t=1}^T \mathbf{Z}_{t+1}, \mathbf{y}\right) - \min_{\mathbf{X} \in \mathcal{S}_n} f\left(\mathbf{X}, \frac{1}{T} \sum_{t=1}^T \mathbf{w}_{t+1}\right) \\ & \leq \frac{D^2 \max \left\{ \sqrt{\beta_X^2 + \beta_{yX}^2}, \sqrt{\beta_y^2 + \beta_{Xy}^2}, \frac{1}{2}(\beta_X + \beta_{Xy}), \frac{1}{2}(\beta_y + \beta_{yX}) \right\}}{T}. \end{aligned} \quad (19)$$

From the definition of g it holds that

$$g\left(\frac{1}{T} \sum_{t=1}^T \mathbf{Z}_{t+1}\right) = \max_{\mathbf{y} \in \mathcal{K}} f\left(\frac{1}{T} \sum_{t=1}^T \mathbf{Z}_{t+1}, \mathbf{y}\right) \quad (20)$$

and

$$\min_{\mathbf{X} \in \mathcal{S}_n} f\left(\mathbf{X}, \frac{1}{T} \sum_{t=1}^T \mathbf{w}_{t+1}\right) \leq f\left(\mathbf{X}^*, \frac{1}{T} \sum_{t=1}^T \mathbf{w}_{t+1}\right) \leq \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}^*, \mathbf{y}) = g(\mathbf{X}^*). \quad (21)$$

Plugging (20) and (21) into the RHS of (19) we obtain the required result. \square

4.2 Efficiently-computable certificates for correctness of low-rank projections

Since Theorem 1 only applies in some neighborhood of an optimal solution, it is of interest to have a procedure for verifying if the rank- r truncated projection of a given point indeed

equals the exact Euclidean projection. In particular, from a practical point of view, it does not matter whether the conditions of Theorem 1 hold. In practice, as long as the truncated projection (see (3)) equals the exact projection (see (2)), we are guaranteed that Algorithm 1 converges correctly with rate $O(1/t)$, without needing to verify any other condition. Luckily, the expression in (2) which characterizes the structure of the Euclidean projection onto the spectrahedron, yields exactly such a verification procedure. As already noted in [17], for any $\mathbf{X} \in \mathbb{S}^n$, we have $\widehat{\Pi}_{\mathbf{S}_n}^r[\mathbf{X}] = \Pi_{\mathbf{S}_n}[\mathbf{X}]$ if and only if the condition

$$\sum_{i=1}^r \lambda_i(\mathbf{X}) \geq 1 + r \cdot \lambda_{r+1}(\mathbf{X})$$

holds. Note that verifying this condition simply requires increasing the rank of the SVD computation by one, i.e., computing a rank- $(r+1)$ SVD of the matrix to project rather than a rank- r SVD.

5 Empirical Evidence

The goal of this section is to bring empirical evidence in support of our theoretical approach. We consider various tasks that take the form of minimizing a composite objective, i.e., the sum of a smooth convex function and a nonsmooth convex function, where the nonsmoothness comes from either an ℓ_1 -norm or ℓ_2 -norm regularizer / penalty term, over a τ -scaled spectrahedron. In all cases the nonsmooth objective can be written as a saddle-point with function $f(\mathbf{X}, \mathbf{y})$ which is linear in \mathbf{y} and in particular satisfies Assumption 2.

The tasks considered include 1. sparse PCA, 2. robust PCA, 3. low-rank and sparse recovery, 4. phase synchronization, and 5. linearly-constrained low-rank estimation, under variety of parameters.

For all tasks considered we generate random instances, and examine the sequences of iterates generated by Algorithm 1 $\{(\mathbf{X}_t, \mathbf{y}_t)\}_{t \geq 1}$, $\{(\mathbf{Z}_t, \mathbf{w}_t)\}_{t \geq 2}$, when initialized with simple initialization procedures. Out of both sequences generated, we choose our candidate for the optimal solution to be the iterate for which the dual-gap, which is a certificate for optimality, is smallest. See Appendix D.

In all tasks considered the goal is to recover a ground-truth low-rank matrix $\mathbf{M}_0 \in \mathbb{S}^n$ from some noisy observation of it $\mathbf{M} = \mathbf{M}_0 + \mathbf{N}$, where $\mathbf{N} \in \mathbb{S}^n$ is a noise matrix. We measure the signal-to-noise ratio (SNR) as $\|\mathbf{M}_0\|_F^2 / \|\mathbf{N}\|_F^2$. In all experiments we measure the relative initialization error by $\left\| \frac{\text{Tr}(\mathbf{M}_0)}{\tau} \mathbf{X}_1 - \mathbf{M}_0 \right\|_F^2 / \|\mathbf{M}_0\|_F^2$, and similarly we measure the relative recovery error by $\left\| \frac{\text{Tr}(\mathbf{M}_0)}{\tau} \mathbf{X}^* - \mathbf{M}_0 \right\|_F^2 / \|\mathbf{M}_0\|_F^2$. Note that in some of the experiments we take $\tau < \text{Tr}(\mathbf{M}_0)$ to prevent the method from overfitting the noise. In addition, we measure the (standard) strict complementarity parameter which corresponds to the eigen-gap $\text{gap}(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) := \lambda_{n-r}(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*))$, $r = \text{rank}(\mathbf{M}_0)$.

In all experiments we use SVDs of rank $r = \text{rank}(\mathbf{M}_0)$ to compute the projections in Algorithm 1 according to the truncated projection given in (3). To certify the correctness

of these low-rank projections (that is, that they equal the exact Euclidean projection) we confirm that the inequality

$$\sum_{i=1}^r \lambda_i(\mathbf{P}_j) \geq \tau + r \cdot \lambda_{r+1}(\mathbf{P}_j)$$

always holds for $\mathbf{P}_1 = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t)$ and $\mathbf{P}_2 = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{W}_{t+1})$ (see also Section 4.2). Indeed, we can now already state our main observation from the experiments:

In all tasks considered and for all random instances generated, throughout all iterations of Algorithm 1, when initialized with a simple “warm-start” strategy and when computing only rank- r truncated projections, $r = \text{rank}(\mathbf{M}_0)$, the truncated projections of $\mathbf{P}_1 = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{Y}_t)$ and $\mathbf{P}_2 = \mathbf{X}_t - \eta \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{W}_{t+1})$ equal their exact full-rank counterparts. That is, Algorithm 1, using only rank- r SVDs, computed exactly the same sequences of iterates it would have computed if using full-rank SVDs.

Aside from the above observation, in the sequel we demonstrate that all models considered indeed satisfy that: 1. the returned solution, denoted $(\mathbf{X}^*, \mathbf{y}^*)$, is of the same rank as the ground-truth matrix and satisfies the strict complementarity condition with non-negligible parameter (measured by the eigengap $\lambda_{n-r}(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*))$), 2. the recovery error of the returned solution indeed improves significantly over the error of the initialization point.

5.1 Sparse PCA

We consider the sparse PCA problem in a well known convex formulation taken from [12] and its equivalent saddle-point formulation:

$$\min_{\substack{\text{Tr}(\mathbf{X})=1, \\ \mathbf{X} \succeq 0}} \langle \mathbf{X}, -\mathbf{M} \rangle + \lambda \|\mathbf{X}\|_1 = \min_{\substack{\text{Tr}(\mathbf{X})=1, \\ \mathbf{X} \succeq 0}} \max_{\|\mathbf{Y}\|_\infty \leq 1} \{ \langle \mathbf{X}, -\mathbf{M} \rangle + \lambda \langle \mathbf{X}, \mathbf{Y} \rangle \},$$

where $\mathbf{M} = \mathbf{z}\mathbf{z}^\top + \frac{c}{2}(\mathbf{N} + \mathbf{N}^\top)$ is a noisy observation of a rank-one matrix $\mathbf{z}\mathbf{z}^\top$, with \mathbf{z} being a sparse unit vector. Each entry \mathbf{z}_i is chosen to be 0 with probability 0.9 and $U\{1, \dots, 10\}$ with probability 0.1, and then we normalize \mathbf{z} to be of unit norm.

We test the results obtained when adding different magnitudes of Gaussian or uniform noise. We set the signal-to-noise ratio (SNR) to be a constant. Thus, we set the noise level to $c = \frac{2}{\text{SNR} \cdot \|\mathbf{N} + \mathbf{N}^\top\|_F}$ for our choice of SNR.

We initialize the \mathbf{X} variable with the rank-one approximation of \mathbf{M} . That is, we take $\mathbf{X}_1 = \mathbf{u}_1 \mathbf{u}_1^\top$, where \mathbf{u}_1 is the top eigenvector of \mathbf{M} . For the \mathbf{Y} variable we initialize it with $\mathbf{Y}_1 = \text{sign}(\mathbf{X}_1)$ which is a subgradient of $\|\mathbf{X}_1\|_1$.

We set the step-size to $\eta = 1/(2\lambda)$ and we set the number of iterations to $T = 1000$ and for any set of parameters we average the measurements over 10 i.i.d. runs.

dimension (n)	100	200	400	600
$\downarrow \mathbf{N} \sim U[0, 1], \text{SNR} = 1 \downarrow$				
λ	0.008	0.004	0.002	0.0013
initialization error	0.5997	0.6009	0.5990	0.6002
recovery error	0.0054	0.0040	0.0035	0.0043
dual gap	4.1×10^{-5}	7.9×10^{-5}	4.9×10^{-5}	3.4×10^{-6}
$\lambda_{n-1}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	0.8840	0.8898	0.8938	0.8777
$\downarrow \mathbf{N} \sim U[0, 1], \text{SNR} = 0.05 \downarrow$				
λ	0.04	0.02	0.01	0.0067
initialization error	1.7456	1.7494	1.7566	1.7625
recovery error	0.0425	0.0244	0.0149	0.0100
dual gap	2.0×10^{-9}	5.8×10^{-6}	4.5×10^{-4}	0.0018
$\lambda_{n-1}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	0.7092	0.7854	0.8340	0.8622
$\downarrow \mathbf{N} \sim \mathcal{N}(0.5, \mathbf{I}_n), \text{SNR} = 1 \downarrow$				
λ	0.006	0.003	0.0015	0.001
initialization error	0.1584	0.1464	0.1443	0.1411
recovery error	0.0059	0.0033	0.0019	0.0015
dual gap	8.6×10^{-4}	0.0031	0.0053	0.0060
$\lambda_{n-1}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	0.8406	0.8869	0.9178	0.9331
$\downarrow \mathbf{N} \sim \mathcal{N}(0.5, \mathbf{I}_n), \text{SNR} = 0.05 \downarrow$				
λ	0.04	0.02	0.01	0.005
initialization error	1.6701	1.6620	1.6542	1.6610
recovery error	0.0502	0.0234	0.0137	0.0109
dual gap	1.9×10^{-5}	0.0041	0.0534	0.0409
$\lambda_{n-1}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	0.2200	0.4076	0.5460	0.6788

Table 1: Numerical results for the sparse PCA problem.

5.2 Low-rank and sparse matrix recovery

We consider the problem of recovering a simultaneously low-rank and sparse covariance matrix [36], which can be written as the following saddle-point optimization problem:

$$\min_{\substack{\text{Tr}(\mathbf{X})=1, \\ \mathbf{X} \succeq 0}} \frac{1}{2} \|\mathbf{X} - \mathbf{M}\|_F^2 + \lambda \|\mathbf{X}\|_1 = \min_{\substack{\text{Tr}(\mathbf{X})=\tau, \\ \mathbf{X} \succeq 0}} \max_{\|\mathbf{Y}\|_\infty \leq 1} \frac{1}{2} \|\mathbf{X} - \mathbf{M}\|_F^2 + \lambda \langle \mathbf{X}, \mathbf{Y} \rangle,$$

where $\mathbf{M} = \mathbf{Z}_0 \mathbf{Z}_0^\top + \frac{\epsilon}{2}(\mathbf{N} + \mathbf{N}^\top)$ is a noisy observation of some low-rank and sparse covariance matrix $\mathbf{Z}_0 \mathbf{Z}_0^\top$. We choose $\mathbf{Z}_0 \in \mathbb{R}^{n \times r}$ to be a sparse matrix where each entry $\mathbf{Z}_{0,i,j}$ is chosen to be 0 with probability 0.9 and $U\{1, \dots, 10\}$ with probability 0.1, and then we normalize \mathbf{Z}_0 to be of unit Frobenius norm. We choose $\mathbf{N} \sim \mathcal{N}(0.5, \mathbf{I}_n)$.

We test the model with $\text{rank}(\mathbf{Z}_0 \mathbf{Z}_0^\top) = 1, 5, 10$. We set the signal-to-noise ratio (SNR) to be a constant and set the noise level to $c = \frac{2\|\mathbf{Z}_0 \mathbf{Z}_0^\top\|_F}{\text{SNR} \cdot \|\mathbf{N} + \mathbf{N}^\top\|_F}$ for our choice of SNR.

We initialize the \mathbf{X} variable with the rank- r approximation of \mathbf{M} . That is, we take $\mathbf{X}_1 = \mathbf{U}_r \text{diag}(\Pi_{\Delta_{\tau,r}}[\text{diag}(-\Lambda_r)]) \mathbf{U}_r^\top$, where $\mathbf{U}_r \Lambda_r \mathbf{U}_r^\top$ is the rank- r eigen-decomposition of \mathbf{M} and $\Delta_{\tau,r} = \{\mathbf{z} \in \mathbb{R}^r \mid \mathbf{z} \geq 0, \sum_{i=1}^r \mathbf{z}_i = \tau\}$ is the simplex of radius τ in \mathbb{R}^r . For the \mathbf{Y} variable we initialize it with $\mathbf{Y}_1 = \text{sign}(\mathbf{X}_1)$ which is a subgradient of $\|\mathbf{X}_1\|_1$.

We set the step-size to $\eta = 1$, $\tau = 0.7 \cdot \text{Tr}(\mathbf{Z}_0 \mathbf{Z}_0^\top)$, and the number of iterations in each experiment to $T = 2000$. For each value of r and n we average the measurements over over 10 i.i.d. runs.

dimension (n)	100	200	400	600
$\downarrow r = \text{rank}(\mathbf{Z}_0 \mathbf{Z}_0^\top) = 1, \text{SNR} = 0.48 \downarrow$				
λ	0.0012	0.0035	0.0016	0.001
initialization error	0.4562	0.4471	0.4507	0.4450
recovery error	0.0364	0.0193	0.0160	0.0168
dual gap	0.0083	0.0086	0.0020	4.2×10^{-4}
$\lambda_{n-r}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	0.0628	0.1439	0.1258	0.1069
$\downarrow r = \text{rank}(\mathbf{Z}_0 \mathbf{Z}_0^\top) = 5, \text{SNR} = 2.4 \downarrow$				
λ	0.0012	0.0006	0.0003	0.0002
initialization error	0.2132	0.2103	0.1983	0.1907
recovery error	0.0641	0.0478	0.0349	0.0274
dual gap	9.0×10^{-4}	4.3×10^{-4}	1.4×10^{-4}	7.3×10^{-5}
$\lambda_{n-r}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	0.0148	0.0200	0.0257	0.0277
$\downarrow r = \text{rank}(\mathbf{Z}_0 \mathbf{Z}_0^\top) = 10, \text{SNR} = 4.8 \downarrow$				
λ	0.0007	0.0004	0.0002	0.0001
initialization error	0.1855	0.1661	0.1527	0.1473
recovery error	0.0702	0.0403	0.0268	0.0356
dual gap	4.9×10^{-4}	6.6×10^{-4}	4.2×10^{-4}	3.4×10^{-5}
$\lambda_{n-r}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	0.0072	0.0142	0.0187	0.0160

Table 2: Numerical results for the low-rank and sparse matrix recovery problem.

5.3 Robust PCA

We consider the robust PCA problem [29] in the following formulation:

$$\min_{\substack{\text{Tr}(\mathbf{X})=\tau, \\ \mathbf{X} \succeq 0}} \|\mathbf{X} - \mathbf{M}\|_1 = \min_{\substack{\text{Tr}(\mathbf{X})=\tau, \\ \mathbf{X} \succeq 0}} \max_{\|\mathbf{Y}\|_\infty \leq 1} \langle \mathbf{X} - \mathbf{M}, \mathbf{Y} \rangle,$$

where $\mathbf{M} = r\mathbf{Z}_0\mathbf{Z}_0^\top + \frac{1}{2}(\mathbf{N} + \mathbf{N}^\top)$ is a sparsely-corrupted observation of some rank- r matrix $\mathbf{Z}_0\mathbf{Z}_0^\top$. We choose $\mathbf{Z}_0 \in \mathbb{R}^{n \times r}$ to be a random unit Frobenius norm matrix. For $\mathbf{N} \in \mathbb{R}^{n \times n}$, we choose each entry to be 0 with probability $1 - 1/\sqrt{n}$ and otherwise 1 or -1 with equal probability.

We initialize the \mathbf{X} variable with the projection $\mathbf{X}_1 = \Pi_{\{\text{Tr}(\mathbf{X})=\tau, \mathbf{X} \succeq 0\}}[\mathbf{M}]$, and the \mathbf{Y} variable with $\mathbf{Y}_1 = \text{sign}(\mathbf{X}_1 - \mathbf{M})$.

We test the model with $\text{rank}(\mathbf{Z}_0\mathbf{Z}_0^\top) = 1, 5, 10$. For $\text{rank}(\mathbf{Z}_0\mathbf{Z}_0^\top) = 1$ we set the step-size to $\eta = n/10$ and for $\text{rank}(\mathbf{Z}_0\mathbf{Z}_0^\top) = 5, 10$ we set it to $\eta = 1$. We set the trace bound to $\tau = 0.95 \cdot \text{Tr}(r\mathbf{Z}_0\mathbf{Z}_0^\top)$. For every set of parameters we average the measurements over 10 i.i.d. runs.

dimension (n)	100	200	400	600
$\downarrow r = \text{rank}(\mathbf{Z}_0 \mathbf{Z}_0^\top) = 1, T = 3000 \downarrow$				
SNR	0.0021	7.2×10^{-4}	2.5×10^{-4}	1.3×10^{-4}
initialization error	1.3511	1.3430	1.2889	1.2606
recovery error	0.0084	0.0107	0.0109	0.0107
dual gap	0.0016	0.0029	0.0044	0.0069
$\lambda_{n-r}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	15.5944	41.2139	85.8117	140.5349
$\downarrow r = \text{rank}(\mathbf{Z}_0 \mathbf{Z}_0^\top) = 5, T = 20,000 \downarrow$				
SNR	0.0110	0.0038	0.0013	6.9×10^{-4}
initialization error	1.5501	1.5527	1.5221	1.4833
recovery error	0.0092	0.0092	0.0087	0.0075
dual gap	0.0084	0.0390	0.1866	0.4721
$\lambda_{n-r}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	7.6734	26.2132	66.1113	108.7215
$\downarrow r = \text{rank}(\mathbf{Z}_0 \mathbf{Z}_0^\top) = 10, T = 30,000 \downarrow$				
SNR	0.0229	0.0077	0.0026	0.0014
initialization error	1.5729	1.6485	1.6317	1.5949
recovery error	0.0079	0.0081	0.0073	0.0065
dual gap	0.0139	0.0338	0.1533	0.3561
$\lambda_{n-r}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	1.7945	16.9890	48.9799	82.2727

Table 3: Numerical results for the robust PCA problem.

5.4 Phase synchronization

We consider the phase synchronization problem (see for instance [42]) which can be written as:

$$\max_{\substack{\mathbf{z} \in \mathbb{C}^n, \\ |z_j|=1 \ \forall j \in [n]}} \mathbf{z}^* \mathbf{M} \mathbf{z}, \quad (22)$$

where $\mathbf{M} = \mathbf{z}_0 \mathbf{z}_0^* + c \mathbf{N}$ is a noisy observation of some rank-one matrix $\mathbf{z}_0 \mathbf{z}_0^*$ such that $\mathbf{z}_0 \in \mathbb{C}^n$ and $\mathbf{z}_{0j} = e^{i\theta_j}$ where $\theta_j \in [0, 2\pi]$. We follow the statistical model in [42] where the noise matrix $\mathbf{N} \in \mathbb{C}^{n \times n}$ is chosen such that every entry is

$$\mathbf{N}_{jk} = \begin{cases} \mathcal{N}(0, 1) + i\mathcal{N}(0, 1) & j < k \\ \overline{\mathbf{N}_{kj}} & j > k \\ 0 & j = k \end{cases}.$$

It is known that for a large n and $c = \mathcal{O}\left(\sqrt{\frac{n}{\log n}}\right)$, with high probability the SDP relaxation of (22) is able to recover the original signal (see [42]).

We solve a penalized version of the SDP relaxation of (22) which can be written as the following saddle-point optimization problem:

$$\min_{\substack{\text{Tr}(\mathbf{X})=n, \\ \mathbf{X} \succeq 0}} \langle \mathbf{X}, -\mathbf{M} \rangle + \lambda \|\text{diag}(\mathbf{X}) - \vec{\mathbf{1}}\|_2 = \min_{\substack{\text{Tr}(\mathbf{X})=n, \\ \mathbf{X} \succeq 0}} \max_{\|\mathbf{y}\|_2 \leq 1} \langle \mathbf{X}, -\mathbf{M} \rangle + \lambda \langle \text{diag}(\mathbf{X}) - \vec{\mathbf{1}}, \mathbf{y} \rangle,$$

where $\vec{\mathbf{1}}$ is the all-ones vector.

While the phase synchronization problem is formulated over the complex numbers, extending our results to handle this model is straightforward.

We initialize the \mathbf{X} variable with the rank-one approximation of \mathbf{M} . That is, we take $\mathbf{X}_1 = n\mathbf{u}_1\mathbf{u}_1^*$, where \mathbf{u}_1 is the top eigenvector of \mathbf{M} . For the \mathbf{y} variable we initialize it with $\mathbf{y}_1 = (\text{diag}(\mathbf{X}_1) - \vec{\mathbf{1}})/\|\text{diag}(\mathbf{X}_1) - \vec{\mathbf{1}}\|_2$.

We set the noise level to $c = 0.18\sqrt{n}$. We set the number of iterations in each experiment to $T = 10,000$ and for each choice of n we average the measurements over 10 i.i.d. runs.

dimension (n)	100	200	400	600
SNR	0.1553	0.0775	0.0387	0.0258
λ	200	600	1600	2800
η	1/400	1/800	1/1800	1/1800
initialization error	0.1270	0.1255	0.1284	0.1323
recovery error	0.0698	0.0659	0.0683	0.0719
dual gap	7.8×10^{-8}	3.9×10^{-5}	0.1553	0.5112
$\lambda_{n-1}(\nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*))$	39.8591	78.9982	150.3524	217.06
$\ \text{diag}(\mathbf{X}^*) - \vec{\mathbf{1}}\ _2$	3.2×10^{-10}	2.1×10^{-8}	5.1×10^{-7}	3.7×10^{-7}

Table 4: Numerical results for the phase synchronization problem.

5.5 Linearly constrained low-rank matrix estimation

Consider the following penalized formulation:

$$\min_{\substack{\text{Tr}(\mathbf{X})=1, \\ \mathbf{X} \succeq 0}} \langle \mathbf{X}, -\mathbf{M} \rangle + \lambda \|\mathcal{A}(\mathbf{X}) - \mathbf{b}\|_2 = \min_{\substack{\text{Tr}(\mathbf{X})=1, \\ \mathbf{X} \succeq 0}} \max_{\|\mathbf{y}\|_2 \leq 1} \langle \mathbf{X}, -\mathbf{M} \rangle + \lambda \langle \mathcal{A}(\mathbf{X}) - \mathbf{b}, \mathbf{y} \rangle,$$

where $\mathbf{M} = \mathbf{z}_0\mathbf{z}_0^\top + \frac{c}{2}(\mathbf{N} + \mathbf{N}^\top)$ is the noisy observation of some rank-one matrix $\mathbf{z}_0\mathbf{z}_0^\top$ such that $\|\mathbf{z}_0\|_2 = 1$ and the noise matrix is chosen $\mathbf{N} \sim \mathcal{N}(0, \mathbf{I}_n)$. We take $\mathcal{A}(\mathbf{X}) = (\langle \mathbf{A}_1, \mathbf{X} \rangle, \dots, \langle \mathbf{A}_m, \mathbf{X} \rangle)^\top$ with matrices $\mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{S}^n$ of the form $\mathbf{A}_i = \mathbf{v}_i\mathbf{v}_i^\top$ such that $\mathbf{v}_i \sim \mathcal{N}(0, 1)$. We take $\mathbf{b} \in \mathbb{R}^m$ such that $b_i = \langle \mathbf{A}_i, \mathbf{z}_0\mathbf{z}_0^\top \rangle$.

We initialize the \mathbf{X} variable with the rank-one approximation of \mathbf{M} . That is, we take $\mathbf{X}_1 = \mathbf{u}_1\mathbf{u}_1^\top$, where \mathbf{u}_1 is the top eigenvector of \mathbf{M} . The \mathbf{y} variable is initialized with $\mathbf{y}_1 = (\mathcal{A}(\mathbf{X}_1) - \mathbf{b})/\|\mathcal{A}(\mathbf{X}_1) - \mathbf{b}\|_2$.

We set the number of constraints to $m = n$, the penalty parameter to $\lambda = 2$, and the step-size to $\eta = 1/(2\lambda)$. We set the number of iterations in each experiment to $T = 2000$ and for each value of n we average the measurements over 10 i.i.d. runs.

dimension (n)	100	200	400	600
SNR	0.15	0.075	0.04	0.027
initialization error	0.1219	0.1324	0.1242	0.1228
recovery error	0.0437	0.0617	0.0685	0.0735
dual gap	5.3×10^{-11}	5.0×10^{-12}	8.5×10^{-12}	2.3×10^{-11}
$\lambda_{n-1}(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*)) - \lambda_n(\nabla_{\mathbf{x}} f(\mathbf{X}^*, \mathbf{y}^*))$	0.2941	0.3409	0.4690	0.5069
$\ \mathcal{A}(\mathbf{X}^*) - \mathbf{b}\ _2$	0.0080	0.0082	0.0079	0.0073

Table 5: Numerical results for the linearly constrained low-rank matrix estimation problem.

6 Discussion

This work expands upon a line of research that aims to harness the ability of convex relaxations to produce low-rank and high-quality solutions to important low-rank matrix optimization problems, while insisting on methods that, at least locally, store and manipulate only low-rank matrices. Focusing on the challenging case of nonsmooth objective functions and following our evidence for the difficulties of obtaining such a result for subgradient methods (Lemma 5), we consider tackling nonsmooth objectives via saddle-point formulations. We prove that indeed under a generalized strict complementarity condition, a state-of-the-art method for convex-concave saddle-point problems converges locally while storing and manipulating only low-rank matrices. Extensive experiments over several tasks demonstrate that our conceptual approach of utilizing low-rank projections for more efficient optimization is not only of theoretical merit, but indeed seems to work well in practice.

Acknowledgements

This research was supported by the ISRAEL SCIENCE FOUNDATION (grant No. 1108/18).

References

- [1] Emmanuel Abbe. Community detection and stochastic block models: recent developments. *The Journal of Machine Learning Research*, 18(1):6446–6531, 2017.
- [2] Praneeth Netrapalli, UN Niranjan, Sujay Sanghavi, Animashree Anandkumar and Prateek Jain. Non-convex robust pca. *Advances in Neural Information Processing Systems*, page 1107–1115, 2014.
- [3] Amir Beck. First-order methods in optimization. *SIAM*, 2017.
- [4] Amir Beck and Marc Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, 2(1):183–202, 2009.
- [5] Dimitri P Bertsekas. Nonlinear programming. *Athena Scientific*, 1999.

- [6] Afonso S. Bandeira, Nicolas Boumal, and Amit Singer. Tightness of the maximum likelihood semidefinite relaxation for angular synchronization. *Mathematical Programming*, 163(1-2):145–167, 2017.
- [7] Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends in Machine Learning*, 8(3-4):231–357, 2015.
- [8] Vasileios Charisopoulos, Yudong Chen, Damek Davis, Mateo Díaz, Lijun Ding, and Dmitriy Drusvyatskiy. Low-rank matrix recovery with composite optimization: good conditioning and rapid convergence. *Foundations of Computational Mathematics*, pages 1–89, 2021.
- [9] PXinyang Yi, Dohyung Park, Yudong Chen and Constantine Caramanis. Fast algorithms for robust pca via gradient descent. *Advances in Neural Information Processing Systems*, page 4152–4160, 2016.
- [10] Yudong Chen, Huan Xu, Constantine Caramanis, and Sujay Sanghavi. Robust matrix completion and corrupted columns. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pages 873–880. Citeseer, 2011.
- [11] Frank H. Clarke. *Optimization and Nonsmooth Analysis*. SIAM, 1990.
- [12] Alexandre d’Aspremont, Laurent El Ghaoui, Michael I Jordan, and Gert RG Lanckriet. A direct formulation for sparse pca using semidefinite programming. *SIAM review*, 49(3):434–448, 2007.
- [13] Lijun Ding, Jicong Fan, and Madeleine Udell. *kfw*: A frank-wolfe style algorithm with stronger subproblem oracles, 2020.
- [14] Dmitriy Drusvyatskiy and Adrian S Lewis. Error bounds, quadratic growth, and linear convergence of proximal methods. *Mathematics of Operations Research*, 43(3):919–948, 2018.
- [15] Francisco Facchinei and Jong-Shi Pang. *Finite-Dimensional Variational Inequalities and Complementarity Problems*, volume II. Springer-Verlag New York, 2003.
- [16] Dan Garber. Linear convergence of frank-wolfe for rank-one matrix recovery without strong convexity. *arXiv preprint arXiv:1912.01467*, 2019.
- [17] Dan Garber. On the convergence of projected-gradient methods with low-rank projections for smooth convex minimization over trace-norm balls and related problems. *SIAM Journal on Optimization*, 2019.
- [18] Dan Garber. On the convergence of stochastic gradient descent with low-rank projections for convex low-rank matrix problems. *Conference on Learning Theory, COLT*, 125:1666–1681, 2020.
- [19] Dan Garber and Atara Kaplan. Fast stochastic algorithms for low-rank and nonsmooth matrix problems. *The 22nd International Conference on Artificial Intelligence and Statistics, AISTATS*, 89:286–294, 2019.
- [20] Dan Garber and Atara Kaplan. On the efficient implementation of the matrix exponentiated gradient algorithm for low-rank matrix optimization. *arXiv preprint arXiv:2012.10469*, 2020.

- [21] Hui Zou, Trevor Hastie, and Robert Tibshirani. Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, 15:2006, 2004.
- [22] Atara Kaplan and Dan Garber. Low-rank extragradient method for nonsmooth and low-rank matrix optimization problems. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 26332–26344. Curran Associates, Inc., 2021.
- [23] Olga Klopp, Karim Lounici, and Alexandre B Tsybakov. Robust matrix completion. *Probability Theory and Related Fields*, 169(1):523–564, 2017.
- [24] G.M. Korpelevich. The extragradient method for finding saddle points and other problems. *Matecon*, 12:747–756, 1976.
- [25] Srinadh Bhojanapalli, Anastasios Kyrillidis, and Sujay Sanghavi. Dropping convexity for faster semi-definite optimization. In *29th Annual Conference on Learning Theory*, volume 49 of *Proceedings of Machine Learning Research*, pages 530–582, 2016.
- [26] Xiao Li, Zhihui Zhu, Anthony Man-Cho So, and Rene Vidal. Nonconvex robust low-rank matrix recovery. *SIAM Journal on Optimization*, 30(1):660–686, 2020.
- [27] Francesco Locatello, Alp Yurtsevert, Olivier Fercoq, and Volkan Cevher. Stochastic frank-wolfe for composite convex minimization. *Advances in Neural Information Processing Systems*, 32, 2019.
- [28] Canyi Lu, Jinhui Tang, Shuicheng Yan, and Zhouchen Lin. Generalized nonconvex nonsmooth low-rank minimization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4130–4137, 2014.
- [29] Emmanuel J. Candès, Xiaodong Li, Yi Ma and John Wright. Robust principal component analysis? *Journal of the ACM*, 58, 2009.
- [30] Adel Javanmard, Andrea Montanari, and Federico Ricci-Tersenghi. Phase transitions in semidefinite relaxations. *Proceedings of the National Academy of Sciences*, 113(16):E2218–E2223, 2016.
- [31] Arkadi Nemirovski. Prox-method with rate of convergence $o(1/t)$ for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM J. on Optimization*, 15(1):229–251, 2005.
- [32] Gergely Odor, Yen-Huan Li, Alp Yurtsever, Ya-Ping Hsieh, Quoc Tran-Dinh, Marwa El Halabi, and Volkan Cevher. Frank-wolfe works for non-lipschitz continuous gradient objectives: scalable poisson phase retrieval. In *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6230–6234. Ieee, 2016.
- [33] Stephen Boyd, Neal Parikh, Eric Chu, Borja Peleato, and Jonathan Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning*, 3(1):1–122, 2011.
- [34] John Wright, Arvind Ganesh, Shankar Rao, Yigang Peng and Yi Ma. Robust principal component analysis: Exact recovery of corrupted low-rank matrices via convex optimization. *Advances in Neural Information Processing Systems 22*, pages 2080–2088, 2009.

- [35] Julien Mairal, Francis Bach, Jean Ponce, and Guillermo Sapiro. Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11:19–60, 2010.
- [36] Emile Richard, Pierre-Andr e Savalle and Nicolas Vayatis. Estimation of simultaneously sparse and low rank matrices. *Proceedings of the 29th International Conference on Machine Learning*, 2012.
- [37] Kiran K Thekumparampil, Prateek Jain, Praneeth Netrapalli, and Sewoong Oh. Projection efficient subgradient method and optimal nonsmooth frank-wolfe method. In *Advances in Neural Information Processing Systems*, volume 33, pages 12211–12224. Curran Associates, Inc., 2020.
- [38] Nicolas Boumal, Vladislav Voroninski, and Afonso S. Bandeira. Deterministic guarantees for burer-monteiro factorizations of smooth semidefinite programs. *Communications on Pure and Applied Mathematics*, 73(3):581 – 608, 2020.
- [39] Cun Mu, Yuqian Zhang, John Wright and Donald Goldfarb. Scalable robust matrix recovery: Frank-wolfe meets proximal methods. *SIAM Journal on Scientific Computing*, 38(5):A3291–A3317, 2016.
- [40] Lijun Ding, Yingjie Fei, Qiantong Xu and Chengrun Yang. Spectral frank-wolfe algorithm: Strict complementarity and linear convergence. *ICML*, 2020.
- [41] Quanming Yao, James T Kwok, Taifeng Wang, and Tie-Yan Liu. Large-scale low-rank matrix learning with nonconvex regularizers. *IEEE transactions on pattern analysis and machine intelligence*, 41(11):2628–2643, 2018.
- [42] Yiqiao Zhong and Nicolas Boumal. Near-optimal bounds for phase synchronization. *SIAM Journal on Optimization*, 28(2):989–1016, 2018.
- [43] Zirui Zhou and Anthony Man-Cho So. A unified approach to error bounds for structured convex optimization problems. *Mathematical Programming*, 165(2):689–728, 2017.

Appendices

A Proof of Lemma 2

We first restate the lemma and then prove it.

Lemma 10. *Let $\mathbf{X}^* \in \mathcal{S}_n$ be a rank- r^* optimal solution to Problem (1). \mathbf{X}^* satisfies the (standard) strict complementarity assumption with parameter $\delta > 0$ if and only if there exists a subgradient $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ such that $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$ and $\lambda_{n-r^*}(\mathbf{G}^*) - \lambda_n(\mathbf{G}^*) \geq \delta$.*

Proof. By Slater’s condition strong duality holds for Problem (1). Therefore, the KKT conditions for Problem (1) hold for the optimal solution \mathbf{X}^* and some optimal dual solu-

tion (\mathbf{Z}^*, s^*) . The Lagrangian of Problem (1) can be written as

$$\mathcal{L}(\mathbf{X}, \mathbf{Z}, s) = g(\mathbf{X}) + s(1 - \text{Tr}(\mathbf{X})) - \langle \mathbf{Z}, \mathbf{X} \rangle.$$

Thus, using the generalized KKT conditions for nonsmooth optimization problems (see Theorem 6.1.1 in [11]), this implies that for the primal and dual optimal solutions

$$\begin{aligned} \mathbf{0} &\in \partial g(\mathbf{X}^*) - \mathbf{Z}^* - s^* \mathbf{I}, \\ \langle \mathbf{X}^*, \mathbf{Z}^* \rangle &= 0, \\ \text{Tr}(\mathbf{X}^*) &= 1, \\ \mathbf{X}^*, \mathbf{Z}^* &\succeq 0. \end{aligned}$$

The generalized first order optimality condition for unconstrained minimization implies that there exists some $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ for which $\mathbf{0} = \mathbf{G}^* - \mathbf{Z}^* - s^* \mathbf{I}$. It remains to be shown that $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$.

The cone of positive semidefinite matrices is self-dual, that is $\mathbf{Z}^* \succeq 0$ if and only if $\langle \mathbf{X}, \mathbf{Z}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$. Therefore, $\mathbf{Z}^* \succeq 0$ if and only if for all $\mathbf{X} \in \mathcal{S}_n$ it holds that

$$\begin{aligned} 0 \leq \langle \mathbf{X}, \mathbf{Z}^* \rangle &= \langle \mathbf{X}, \mathbf{Z}^* \rangle - \langle \mathbf{X}^*, \mathbf{Z}^* \rangle + s^* \langle \mathbf{X} - \mathbf{X}^*, \mathbf{I} \rangle = \langle \mathbf{X} - \mathbf{X}^*, \mathbf{Z}^* + s^* \mathbf{I} \rangle \\ &= \langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \end{aligned}$$

as desired. The first equality holds using the complementarity condition and the property that $\text{Tr}(\mathbf{X}) = \text{Tr}(\mathbf{X}^*) = 1$.

Using the equality $\mathbf{G}^* = \mathbf{Z}^* + s^* \mathbf{I}$ it holds that

$$\begin{aligned} \lambda_{n-r^*}(\mathbf{Z}^*) &= \lambda_{n-r^*}(\mathbf{Z}^*) + s^* - \lambda_n(\mathbf{Z}^*) - s^* = \lambda_{n-r^*}(\mathbf{Z}^* + s^* \mathbf{I}) - \lambda_n(\mathbf{Z}^* + s^* \mathbf{I}) \\ &= \lambda_{n-r^*}(\mathbf{G}^*) - \lambda_n(\mathbf{G}^*). \end{aligned}$$

Thus, \mathbf{X}^* satisfies the strict complementarity assumption with parameter $\delta > 0$, i.e., $\lambda_{n-r^*}(\mathbf{Z}^*) \geq \delta$, if and only if $\lambda_{n-r^*}(\mathbf{G}^*) - \lambda_n(\mathbf{G}^*) \geq \delta$. \square

B Proof of Lemma 6

We first restate the lemma and then prove it.

Lemma 11. *If $(\mathbf{X}^*, \mathbf{y}^*)$ is a saddle-point of Problem (5) then \mathbf{X}^* is an optimal solution to Problem (1), $\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*) \in \partial g(\mathbf{X}^*)$, and for all $\mathbf{X} \in \mathcal{S}_n$ it holds that $\langle \mathbf{X} - \mathbf{X}^*, \nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*) \rangle \geq 0$. Conversely, under Assumption 2, if \mathbf{X}^* is an optimal solution to Problem (1), and $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ which satisfies $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$, then there exists $\mathbf{y}^* \in \arg \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}^*, \mathbf{y})$ such that $(\mathbf{X}^*, \mathbf{y}^*)$ is a saddle-point of Problem (5), and $\nabla_{\mathbf{X}} f(\mathbf{X}^*, \mathbf{y}^*) = \mathbf{G}^*$.*

Proof. For the first direction of the lemma, we first observe that for any $\mathbf{X}_1, \mathbf{X}_2 \in \mathcal{S}_n$ and $\tilde{\mathbf{y}}_1 \in \arg \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}_1, \mathbf{y})$, $\tilde{\mathbf{y}}_2 \in \arg \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}_2, \mathbf{y})$, using the gradient inequality for $f(\cdot, \tilde{\mathbf{y}}_2)$, it holds that

$$\begin{aligned} g(\mathbf{X}_1) &= f(\mathbf{X}_1, \tilde{\mathbf{y}}_1) \geq f(\mathbf{X}_1, \tilde{\mathbf{y}}_2) \geq f(\mathbf{X}_2, \tilde{\mathbf{y}}_2) + \langle \nabla_{\mathbf{X}} f(\mathbf{X}_2, \tilde{\mathbf{y}}_2), \mathbf{X}_1 - \mathbf{X}_2 \rangle \\ &= g(\mathbf{X}_2) + \langle \nabla_{\mathbf{X}} f(\mathbf{X}_2, \tilde{\mathbf{y}}_2), \mathbf{X}_1 - \mathbf{X}_2 \rangle. \end{aligned}$$

Thus, $\nabla_{\mathbf{X}}f(\mathbf{X}_2, \tilde{\mathbf{y}}_2)$ is a subgradient of $g(\cdot)$ at \mathbf{X}_2 , i.e., $\nabla_{\mathbf{X}}f(\mathbf{X}_2, \tilde{\mathbf{y}}_2) \in \partial g(\mathbf{X}_2)$.

In particular, for a saddle-point $(\mathbf{X}^*, \mathbf{y}^*) \in \mathcal{S}_n \times \mathcal{K}$ it holds that $\mathbf{y}^* \in \arg \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}^*, \mathbf{y})$, and therefore, it follows that $\nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*) \in \partial g(\mathbf{X}^*)$. In addition, for all $\mathbf{X} \in \mathcal{S}_n$ and $\tilde{\mathbf{y}} \in \arg \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}, \mathbf{y})$ we have

$$g(\mathbf{X}^*) = f(\mathbf{X}^*, \mathbf{y}^*) \leq f(\mathbf{X}, \mathbf{y}^*) \leq f(\mathbf{X}, \tilde{\mathbf{y}}) = g(\mathbf{X}),$$

which implies that \mathbf{X}^* is an optimal solution to $\min_{\mathbf{X} \in \mathcal{S}_n} g(\mathbf{X})$.

Finally, we need to show that the subgradient $\nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*) \in \partial g(\mathbf{X}^*)$ indeed satisfies the first-order optimality condition for $g(\cdot)$ at \mathbf{X}^* . To see this, we observe that since \mathbf{X}^* is an optimal solution to $\min_{\mathbf{X} \in \mathcal{S}_n} f(\mathbf{X}, \mathbf{y}^*)$, it follows from the first-order optimality condition for the problem $\min_{\mathbf{X} \in \mathcal{S}_n} f(\mathbf{X}, \mathbf{y}^*)$, that for all $\mathbf{W} \in \mathcal{S}_n$

$$\langle \mathbf{W} - \mathbf{X}^*, \nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*) \rangle \geq 0,$$

as needed.

For the second direction, let $\mathbf{X}^* \in \arg \min_{\mathbf{X} \in \mathcal{S}_n} g(\mathbf{X})$ and let $\mathbf{G}^* \in \partial g(\mathbf{X}^*)$ such that $\langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle \geq 0$ for all $\mathbf{X} \in \mathcal{S}_n$. By Assumption 2 and using Danskin's theorem (see for instance [5]), the subdifferential set of $g(\mathbf{X}^*) = h(\mathbf{X}^*) + \max_{\mathbf{y} \in \mathcal{K}} \mathbf{y}^\top (\mathcal{A}(\mathbf{X}^*) - \mathbf{b})$ can be written as

$$\begin{aligned} \partial g(\mathbf{X}^*) &= \nabla h(\mathbf{X}^*) + \text{conv} \left\{ \mathcal{A}^\top(\mathbf{y}) \mid \mathbf{y} \in \arg \max_{\mathbf{y} \in \mathcal{K}} \mathbf{y}^\top (\mathcal{A}(\mathbf{X}^*) - \mathbf{b}) \right\} \\ &= \nabla h(\mathbf{X}^*) + \mathcal{A}^\top \left(\text{conv} \left\{ \mathbf{y} \mid \mathbf{y} \in \arg \max_{\mathbf{y} \in \mathcal{K}} \mathbf{y}^\top (\mathcal{A}(\mathbf{X}^*) - \mathbf{b}) \right\} \right) \\ &= \nabla h(\mathbf{X}^*) + \mathcal{A}^\top \left(\left\{ \mathbf{y} \mid \mathbf{y} \in \arg \max_{\mathbf{y} \in \mathcal{K}} \mathbf{y}^\top (\mathcal{A}(\mathbf{X}^*) - \mathbf{b}) \right\} \right) \\ &= \nabla h(\mathbf{X}^*) + \mathcal{A}^\top \left(\left\{ \mathbf{y} \mid \mathbf{y} \in \arg \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}^*, \mathbf{y}) \right\} \right), \end{aligned}$$

where $\text{conv}\{\cdot\}$ denotes the convex hull operation and the third equality follows from the convexity of \mathcal{K} .

Thus, there exists some $\mathbf{y}^* \in \arg \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}^*, \mathbf{y})$ such that $\mathbf{G}^* = \nabla h(\mathbf{X}^*) + \mathcal{A}^\top(\mathbf{y}^*) = \nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*)$.

Since $\mathbf{y}^* \in \arg \max_{\mathbf{y} \in \mathcal{K}} f(\mathbf{X}^*, \mathbf{y})$, it follows that for all $\mathbf{y} \in \mathcal{K}$, $f(\mathbf{X}^*, \mathbf{y}^*) \geq f(\mathbf{X}^*, \mathbf{y})$. In addition, using the fact that \mathbf{G}^* satisfies the first-order optimality condition, and using gradient inequality w.r.t. $f(\cdot, \mathbf{y}^*)$, we have that for all $\mathbf{X} \in \mathcal{S}_n$,

$$0 \leq \langle \mathbf{X} - \mathbf{X}^*, \mathbf{G}^* \rangle = \langle \mathbf{X} - \mathbf{X}^*, \nabla_{\mathbf{X}}f(\mathbf{X}^*, \mathbf{y}^*) \rangle \leq f(\mathbf{X}, \mathbf{y}^*) - f(\mathbf{X}^*, \mathbf{y}^*).$$

Thus, it follows that $f(\mathbf{X}, \mathbf{y}^*) \geq f(\mathbf{X}^*, \mathbf{y}^*)$. Therefore, $(\mathbf{X}^*, \mathbf{y}^*)$ is indeed a saddle-point of f .

□

C Proof of Lemma 7

We first restate the lemma and then prove it.

Lemma 12. *Let $\{(\mathbf{X}_t, \mathbf{y}_t)\}_{t \geq 1}$ and $\{(\mathbf{Z}_t, \mathbf{w}_t)\}_{t \geq 2}$ be the sequences generated by Algorithm 1 with a fixed step-size $\eta_t = \eta \leq \min \left\{ \frac{1}{\beta_X + \beta_{Xy}}, \frac{1}{\beta_y + \beta_{yX}}, \frac{1}{\beta_X + \beta_{yX}}, \frac{1}{\beta_y + \beta_{Xy}} \right\}$ then*

$$\max_{\mathbf{y} \in \mathcal{K}} f \left(\frac{1}{T} \sum_{t=1}^T \mathbf{Z}_{t+1}, \mathbf{y} \right) - \min_{\mathbf{X} \in \mathcal{S}_n} f \left(\mathbf{X}, \frac{1}{T} \sum_{t=1}^T \mathbf{w}_{t+1} \right) \leq \frac{D^2}{2\eta T},$$

where $D := \sup_{(\mathbf{X}, \mathbf{y}), (\tilde{\mathbf{X}}, \tilde{\mathbf{y}}) \in \mathcal{S}_n \times \mathcal{K}} \|(\mathbf{X}, \mathbf{y}) - (\tilde{\mathbf{X}}, \tilde{\mathbf{y}})\|$.

Proof. The projection theorem states that projecting some point \mathbf{s} onto some closed and convex set \mathcal{C} satisfies that for all $\mathbf{z} \in \mathcal{C}$ it holds that $\langle \Pi_{\mathcal{C}}[\mathbf{s}] - \mathbf{s}, \Pi_{\mathcal{C}}[\mathbf{s}] - \mathbf{z} \rangle \leq 0$. In particular, for any $\mathbf{X} \in \mathcal{S}_n$, using the updates for \mathbf{X}_{t+1} and \mathbf{Z}_{t+1} , the two following inequalities hold:

$$\eta_t \langle \mathbf{Z}_{t+1} - \mathbf{X}, \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t) \rangle \leq \langle \mathbf{X}_t - \mathbf{Z}_{t+1}, \mathbf{Z}_{t+1} - \mathbf{X} \rangle \quad (23)$$

$$\eta_t \langle \mathbf{X}_{t+1} - \mathbf{X}, \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) \rangle \leq \langle \mathbf{X}_t - \mathbf{X}_{t+1}, \mathbf{X}_{t+1} - \mathbf{X} \rangle. \quad (24)$$

By the gradient inequality, for any $\mathbf{X} \in \mathcal{S}_n$

$$\begin{aligned} & f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) - f(\mathbf{X}, \mathbf{w}_{t+1}) \\ & \leq \langle \mathbf{Z}_{t+1} - \mathbf{X}, \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) \rangle \\ & = \langle \mathbf{X}_{t+1} - \mathbf{X}, \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) \rangle + \langle \mathbf{Z}_{t+1} - \mathbf{X}_{t+1}, \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t) \rangle \\ & \quad + \langle \mathbf{Z}_{t+1} - \mathbf{X}_{t+1}, \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) - \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t) \rangle. \end{aligned} \quad (25)$$

We will bound these three terms separately.

For the first term, using (24) and the Pythagoras identity

$$\begin{aligned} & \langle \mathbf{X}_{t+1} - \mathbf{X}, \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) \rangle \\ & \leq \frac{1}{\eta_t} \langle \mathbf{X}_t - \mathbf{X}_{t+1}, \mathbf{X}_{t+1} - \mathbf{X} \rangle \\ & = -\frac{1}{2\eta_t} \|\mathbf{X}_t - \mathbf{X}_{t+1}\|_F^2 + \frac{1}{2\eta_t} \|\mathbf{X}_t - \mathbf{X}\|_F^2 - \frac{1}{2\eta_t} \|\mathbf{X}_{t+1} - \mathbf{X}\|_F^2. \end{aligned} \quad (26)$$

For the second term, using (23) with $\mathbf{X} = \mathbf{X}_{t+1}$ and the Pythagoras identity

$$\begin{aligned} & \langle \mathbf{Z}_{t+1} - \mathbf{X}_{t+1}, \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t) \rangle \\ & \leq \frac{1}{\eta_t} \langle \mathbf{X}_t - \mathbf{Z}_{t+1}, \mathbf{Z}_{t+1} - \mathbf{X}_{t+1} \rangle \\ & = -\frac{1}{2\eta_t} \|\mathbf{X}_t - \mathbf{Z}_{t+1}\|_F^2 + \frac{1}{2\eta_t} \|\mathbf{X}_t - \mathbf{X}_{t+1}\|_F^2 - \frac{1}{2\eta_t} \|\mathbf{Z}_{t+1} - \mathbf{X}_{t+1}\|_F^2. \end{aligned} \quad (27)$$

For the third term, using the Cauchy–Schwarz inequality, the β_X and β_{Xy} smoothness,

and the inequality $2ab \leq a^2 + b^2$ we obtain

$$\begin{aligned}
& \langle \mathbf{Z}_{t+1} - \mathbf{X}_{t+1}, \nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) - \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t) \rangle \\
& \leq \|\mathbf{Z}_{t+1} - \mathbf{X}_{t+1}\|_F \cdot \|\nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) - \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t)\|_F \\
& \leq \|\mathbf{Z}_{t+1} - \mathbf{X}_{t+1}\|_F \cdot \|\nabla_{\mathbf{X}} f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) - \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{w}_{t+1})\|_F \\
& \quad + \|\mathbf{Z}_{t+1} - \mathbf{X}_{t+1}\|_F \cdot \|\nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{w}_{t+1}) - \nabla_{\mathbf{X}} f(\mathbf{X}_t, \mathbf{y}_t)\|_F \\
& \leq (\beta_X \|\mathbf{Z}_{t+1} - \mathbf{X}_t\|_F + \beta_{Xy} \|\mathbf{w}_{t+1} - \mathbf{y}_t\|_2) \cdot \|\mathbf{Z}_{t+1} - \mathbf{X}_{t+1}\|_F \\
& \leq \frac{\beta_X}{2} (\|\mathbf{Z}_{t+1} - \mathbf{X}_t\|_F^2 + \|\mathbf{Z}_{t+1} - \mathbf{X}_{t+1}\|_F^2) + \frac{\beta_{Xy}}{2} (\|\mathbf{w}_{t+1} - \mathbf{y}_t\|_2^2 + \|\mathbf{Z}_{t+1} - \mathbf{X}_{t+1}\|_F^2).
\end{aligned} \tag{28}$$

Plugging (26), (27), and (28) into (25) we obtain

$$\begin{aligned}
& f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) - f(\mathbf{X}, \mathbf{w}_{t+1}) \\
& \leq \frac{1}{2\eta_t} (\|\mathbf{X}_t - \mathbf{X}\|_F^2 - \|\mathbf{X}_{t+1} - \mathbf{X}\|_F^2) + \left(\frac{\beta_X}{2} - \frac{1}{2\eta_t} \right) \|\mathbf{Z}_{t+1} - \mathbf{X}_t\|_F^2 \\
& \quad + \left(\frac{\beta_X + \beta_{Xy}}{2} - \frac{1}{2\eta_t} \right) \|\mathbf{Z}_{t+1} - \mathbf{X}_{t+1}\|_F^2 + \frac{\beta_{Xy}}{2} \|\mathbf{w}_{t+1} - \mathbf{y}_t\|_2^2.
\end{aligned}$$

Using similar arguments, for any $\mathbf{y} \in \mathcal{K}$

$$\begin{aligned}
& f(\mathbf{Z}_{t+1}, \mathbf{y}) - f(\mathbf{Z}_{t+1}, \mathbf{w}_{t+1}) \\
& \leq \frac{1}{2\eta_t} (\|\mathbf{y}_t - \mathbf{y}\|_2^2 - \|\mathbf{y}_{t+1} - \mathbf{y}\|_2^2) + \left(\frac{\beta_y}{2} - \frac{1}{2\eta_t} \right) \|\mathbf{w}_{t+1} - \mathbf{y}_t\|_2^2 \\
& \quad + \left(\frac{\beta_y + \beta_{yX}}{2} - \frac{1}{2\eta_t} \right) \|\mathbf{w}_{t+1} - \mathbf{y}_{t+1}\|_2^2 + \frac{\beta_{yX}}{2} \|\mathbf{Z}_{t+1} - \mathbf{X}_t\|_F^2.
\end{aligned}$$

Summing the last two inequalities, we obtain for $\eta_t \leq \min \left\{ \frac{1}{\beta_X + \beta_{Xy}}, \frac{1}{\beta_y + \beta_{yX}}, \frac{1}{\beta_X + \beta_{yX}}, \frac{1}{\beta_y + \beta_{Xy}} \right\}$

$$\begin{aligned}
f(\mathbf{Z}_{t+1}, \mathbf{y}) - f(\mathbf{X}, \mathbf{w}_{t+1}) & \leq \frac{1}{2\eta_t} (\|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}, \mathbf{y})\|^2 - \|(\mathbf{X}_{t+1}, \mathbf{y}_{t+1}) - (\mathbf{X}, \mathbf{y})\|^2) \\
& \quad + \left(\frac{\beta_X + \beta_{yX}}{2} - \frac{1}{2\eta_t} \right) \|\mathbf{Z}_{t+1} - \mathbf{X}_t\|_F^2 \\
& \quad + \left(\frac{\beta_y + \beta_{Xy}}{2} - \frac{1}{2\eta_t} \right) \|\mathbf{w}_{t+1} - \mathbf{y}_t\|_2^2 \\
& \quad + \left(\frac{\beta_X + \beta_{Xy}}{2} - \frac{1}{2\eta_t} \right) \|\mathbf{Z}_{t+1} - \mathbf{X}_{t+1}\|_F^2 \\
& \quad + \left(\frac{\beta_y + \beta_{yX}}{2} - \frac{1}{2\eta_t} \right) \|\mathbf{w}_{t+1} - \mathbf{y}_{t+1}\|_2^2 \\
& \leq \frac{1}{2\eta_t} (\|(\mathbf{X}_t, \mathbf{y}_t) - (\mathbf{X}, \mathbf{y})\|^2 - \|(\mathbf{X}_{t+1}, \mathbf{y}_{t+1}) - (\mathbf{X}, \mathbf{y})\|^2).
\end{aligned}$$

Averaging over $t = 1, \dots, T$ and taking a $\eta_t = \eta$

$$\frac{1}{T} \sum_{t=1}^T (f(\mathbf{Z}_{t+1}, \mathbf{y}) - f(\mathbf{X}, \mathbf{w}_{t+1})) \leq \frac{1}{2\eta T} \max_{(\mathbf{X}, \mathbf{y}) \in \mathcal{S}_n \times \mathcal{K}} \|(\mathbf{X}_1, \mathbf{y}_1) - (\mathbf{X}, \mathbf{y})\|^2 \leq \frac{D^2}{2\eta T}.$$

Taking the maximum over all $\mathbf{y} \in \mathcal{K}$ and minimum over all $\mathbf{X} \in \mathcal{S}_n$ and using the convexity of $f(\cdot, \mathbf{y})$ and concavity of $f(\mathbf{X}, \cdot)$,

$$\begin{aligned} & \max_{\mathbf{y} \in \mathcal{K}} f\left(\frac{1}{T} \sum_{t=1}^T \mathbf{Z}_{t+1}, \mathbf{y}\right) - \min_{\mathbf{X} \in \mathcal{S}_n} f\left(\mathbf{X}, \frac{1}{T} \sum_{t=1}^T \mathbf{w}_{t+1}\right) \\ & \leq \max_{\mathbf{y} \in \mathcal{K}} \frac{1}{T} \sum_{t=1}^T f(\mathbf{Z}_{t+1}, \mathbf{y}) - \min_{\mathbf{X} \in \mathcal{S}_n} \frac{1}{T} \sum_{t=1}^T f(\mathbf{X}, \mathbf{w}_{t+1}) \leq \frac{D^2}{2\eta T}. \end{aligned}$$

□

D Calculating the dual-gap in saddle-point problems

Set some point $(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) \in \{\text{Tr}(\mathbf{X}) = \tau, \mathbf{X} \succeq 0\} \times \mathcal{K}$. Using the concavity of $f(\hat{\mathbf{Z}}, \cdot)$ and convexity of $f(\cdot, \hat{\mathbf{w}})$, for all $\mathbf{y} \in \mathcal{K}$ and $\mathbf{X} \in \{\text{Tr}(\mathbf{X}) = \tau, \mathbf{X} \succeq 0\}$, it holds that

$$\begin{aligned} f(\hat{\mathbf{Z}}, \mathbf{y}) - f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) & \leq \langle \hat{\mathbf{w}} - \mathbf{y}, -\nabla_{\mathbf{y}} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) \rangle, \\ f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) - f(\mathbf{X}, \hat{\mathbf{w}}) & \leq \langle \hat{\mathbf{Z}} - \mathbf{X}, \nabla_{\mathbf{X}} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) \rangle. \end{aligned}$$

By taking the maximum of all $\mathbf{y} \in \mathcal{K}$ we obtain in particular that

$$\begin{aligned} f(\mathbf{X}^*, \mathbf{y}^*) - f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) & \leq f(\hat{\mathbf{Z}}, \mathbf{y}^*) - f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) \leq \max_{\mathbf{y} \in \mathcal{K}} f(\hat{\mathbf{Z}}, \mathbf{y}) - f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) \\ & \leq \max_{\mathbf{y} \in \mathcal{K}} \langle \hat{\mathbf{w}} - \mathbf{y}, -\nabla_{\mathbf{y}} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) \rangle, \end{aligned}$$

and taking the maximum of all $\mathbf{X} \in \{\text{Tr}(\mathbf{X}) = \tau, \mathbf{X} \succeq 0\}$

$$\begin{aligned} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) - f(\mathbf{X}^*, \mathbf{y}^*) & \leq f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) - f(\mathbf{X}^*, \hat{\mathbf{w}}) \leq f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) - \min_{\substack{\text{Tr}(\mathbf{X})=\tau, \\ \mathbf{X} \succeq 0}} f(\mathbf{X}, \hat{\mathbf{w}}) \\ & \leq \max_{\substack{\text{Tr}(\mathbf{X})=\tau, \\ \mathbf{X} \succeq 0}} \langle \hat{\mathbf{Z}} - \mathbf{X}, \nabla_{\mathbf{X}} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) \rangle. \end{aligned}$$

Summing these two inequalities, we obtain a bound on the dual-gap at $(\hat{\mathbf{Z}}, \hat{\mathbf{w}})$ which can be written as

$$\begin{aligned} g(\hat{\mathbf{Z}}) - g^* & \leq \max_{\mathbf{y} \in \mathcal{K}} f(\hat{\mathbf{Z}}, \mathbf{y}) - \min_{\mathbf{X} \in \mathcal{S}_n} f(\mathbf{X}, \hat{\mathbf{w}}) \\ & \leq \max_{\substack{\text{Tr}(\mathbf{X})=\tau, \\ \mathbf{X} \succeq 0}} \langle \hat{\mathbf{Z}} - \mathbf{X}, \nabla_{\mathbf{X}} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) \rangle - \min_{\mathbf{y} \in \mathcal{K}} \langle \hat{\mathbf{w}} - \mathbf{y}, \nabla_{\mathbf{y}} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) \rangle. \end{aligned}$$

It is easy to see that the maximizer of the first term in the RHS of the above is $\tau \mathbf{v}_n \mathbf{v}_n^\top$ where \mathbf{v}_n is the smallest eigenvector of $\nabla_{\mathbf{X}} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}})$, and the minimizer of the second term is $\mathbf{Y}_{i,j} = \text{sign}(\nabla_{\mathbf{Y}} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}})_{i,j})$ for $\mathcal{K} = \{\mathbf{Y} \in \mathbb{R}^{n \times n} \mid \|\mathbf{Y}\|_\infty \leq 1\}$ and $\nabla_{\mathbf{y}} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}}) / \|\nabla_{\mathbf{y}} f(\hat{\mathbf{Z}}, \hat{\mathbf{w}})\|_2$ for $\mathcal{K} = \{\mathbf{y} \in \mathbb{R}^n \mid \|\mathbf{y}\|_2 \leq 1\}$.