

# ADMM FOR NONSMOOTH COMPOSITE OPTIMIZATION UNDER ORTHOGONALITY CONSTRAINTS

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We consider a class of structured, nonconvex, nonsmooth optimization problems under orthogonality constraints, where the objectives combine a smooth function, a nonsmooth concave function, and a nonsmooth weakly convex function. This class of problems finds diverse applications in statistical learning and data science. Existing methods for addressing these problems often fail to exploit the specific structure of orthogonality constraints, struggle with nonsmooth functions, or result in suboptimal oracle complexity. We propose OADMM, an Alternating Direction Method of Multipliers (ADMM) designed to solve this class of problems using efficient proximal linearized strategies. Two specific variants of OADMM are explored: one based on Euclidean Projection (OADMM-EP) and the other on Riemannian Retraction (OADMM-RR). Under mild assumptions, we prove that OADMM converges to a critical point of the problem with an ergodic convergence rate of  $\mathcal{O}(1/\epsilon^3)$ . Additionally, we establish a super-exponential convergence rate or polynomial convergence rate for OADMM, depending on the specific setting, under the Kurdyka-Lojasiewicz (KL) inequality. To the best of our knowledge, this is the first non-ergodic convergence result for this class of nonconvex nonsmooth optimization problems. Numerical experiments demonstrate that the proposed algorithm achieves state-of-the-art performance.

**Keywords:** Orthogonality Constraints; Nonconvex Optimization; Nonsmooth Composite Optimization; ADMM; Convergence Analysis

## 1 INTRODUCTION

This paper focuses on the following nonsmooth composite optimization problem under orthogonality constraints (‘ $\triangleq$ ’ means define):

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times r}} F(\mathbf{X}) \triangleq f(\mathbf{X}) - g(\mathbf{X}) + h(\mathcal{A}(\mathbf{X})), \text{ s.t. } \mathbf{X}^\top \mathbf{X} = \mathbf{I}_r. \quad (1)$$

Here,  $n \geq r$ ,  $\mathcal{A}(\mathbf{X}) \in \mathbb{R}^m$  is a linear mapping of  $\mathbf{X}$ , and  $\mathbf{I}_r$  is a  $r \times r$  identity matrix. For conciseness, the orthogonality constraints  $\mathbf{X}^\top \mathbf{X} = \mathbf{I}_r$  in Problem (1) is rewritten as  $\mathbf{X} \in \mathcal{M} \in \mathbb{R}^{n \times r}$ , with  $\mathcal{M}$  representing the Stiefel manifold in the literature (Edelman et al., 1998; Absil et al., 2008b).

We impose the following assumptions on Problem (1) throughout this paper. (A-i)  $f(\mathbf{X})$  is  $L_f$ -smooth, satisfying  $\|\nabla f(\mathbf{X}) - \nabla f(\mathbf{X}')\|_F \leq L_f \|\mathbf{X} - \mathbf{X}'\|_F$  holds for all  $\mathbf{X}, \mathbf{X}' \in \mathbb{R}^{n \times r}$ . This implies:  $|f(\mathbf{X}) - f(\mathbf{X}') - \langle \nabla f(\mathbf{X}'), \mathbf{X} - \mathbf{X}' \rangle| \leq \frac{L_f}{2} \|\mathbf{X} - \mathbf{X}'\|_F^2$  (cf. Lemma 1.2.3 in (Nesterov, 2003)). We also assume that  $f(\mathbf{X})$  demonstrates  $C_f$ -Lipschitz continuity, with  $\|\nabla f(\mathbf{X})\|_F \leq C_f$  for all  $\mathbf{X} \in \mathcal{M}$ . The convexity of  $f(\mathbf{X})$  is not assumed. (A-ii) The function  $g(\cdot)$  is convex, proper, and  $C_g$ -Lipschitz continuous, though it is not necessarily smooth. (A-iii) The function  $h(\cdot)$  is proper, lower semicontinuous,  $C_h$ -Lipschitz continuous, and potentially nonsmooth. Also, it is weakly convexity with constant  $W_h \geq 0$ , which implies that the function  $h(\mathbf{y}) + \frac{W_h}{2} \|\mathbf{y}\|_2^2$  is convex for all  $\mathbf{y} \in \mathbb{R}^m$ . (A-iv) The proximal operator,  $\mathbb{P}_\mu(\mathbf{y}') \triangleq \arg \min_{\mathbf{y}} \frac{1}{2\mu} \|\mathbf{y} - \mathbf{y}'\|_2^2 + h(\mathbf{y})$ , can be computed efficiently and exactly for any given  $\mu > 0$  and  $\mathbf{y}' \in \mathbb{R}^m$ .

Problem (1) represents an optimization framework that plays a crucial role in a variety of statistical learning and data science models. These models include sparse Principal Component Analysis (PCA) (Journée et al., 2010; Lu & Zhang, 2012), deep neural networks (Cho & Lee, 2017; Xie et al.,

2017; Bansal et al., 2018; Cogswell et al., 2016; Huang & Gao, 2023), orthogonal nonnegative matrix factorization (Jiang et al., 2022), range-based independent component analysis (Selvan et al., 2015), and dictionary learning (Zhai et al., 2020).

## 1.1 RELATED WORK

► **Optimization under Orthogonality Constraints.** Solving Problem (1) is challenging due to the computationally expensive and non-convex orthogonality constraints. Existing methods can be divided into three classes. *(i)* Geodesic-like methods (Edelman et al., 1998; Abrudan et al., 2008; Absil et al., 2008b; Jiang & Dai, 2015). These methods involve calculating geodesics by solving ordinary differential equations, which can introduce significant computational complexity. To mitigate this, geodesic-like methods iteratively compute the geodesic logarithm using simple linear algebra calculations. Efficient constraint-preserving update schemes have been integrated with the Barzilai-Borwein (BB) stepsize strategy (Wen & Yin, 2013; Jiang & Dai, 2015) for minimizing smooth functions under orthogonality constraints. *(ii)* Projection and retractions methods (Absil et al., 2008b; Golub & Van Loan, 2013). These methods maintain orthogonality constraints through projection or retraction. They reduce the objective value by using its current Euclidean gradient direction or Riemannian tangent direction, followed by an orthogonal projection operation. This projection can be computed using polar decomposition or singular value decomposition, or approximated with QR factorization. *(iii)* Multiplier correction methods (Gao et al., 2018; 2019; Xiao et al., 2022). Leveraging the insight that the Lagrangian multiplier associated with the orthogonality constraint is symmetric and has an explicit closed-form expression at the first-order optimality condition, these methods tackle an alternative unconstrained nonlinear objective minimization problem, rather than the original smooth function under orthogonality constraints.

► **Optimization with Nonsmooth Objectives.** Another challenge in addressing Problem (1) stems from the nonsmooth nature of the objective function. Existing methods for tackling this challenge fall into three main categories. *(i)* Subgradient methods (Ferreira & Oliveira, 1998; Hwang et al., 2015; Li et al., 2021). Subgradient methods, analogous to gradient descent methods, can incorporate various geodesic-like and projection-like techniques. However, they often exhibit slower convergence rates compared to other approaches. *(ii)* Proximal gradient methods (Chen et al., 2020). These methods use a semi-smooth Newton approach to solve a strongly convex minimization problem over the tangent space, finding a descent direction while preserving the orthogonality constraint through a retraction operation. *(iii)* Operator splitting methods (Lai & Osher, 2014; Chen et al., 2016; Zhang et al., 2020b). These methods introduce linear constraints to break down the original problem into simpler subproblems that can be solved separately and exactly. Among these, ADMM is a promising solution for Problem (1) due to its capability to handle nonsmooth objectives and nonconvex constraints separately and alternately. Several ADMM-like algorithms have been proposed for solving nonconvex problems (Boş & Nguyen, 2020; Boş et al., 2019; Wang et al., 2019; Li & Pong, 2015; He & Yuan, 2012; Yuan, 2024; Zhang et al., 2020b), but these methods fail to exploit the specific structure of orthogonality constraints or cannot be adapted to solve Problem (1). *(iv)* Other methods. OBCD (Yuan, 2023) has been proposed to solve a specific class of our problems, while the exact augmented Lagrangian method ManIAL was introduced in (Deng et al., 2024).

► **Detailed Discussions on Operator Splitting Methods.** We list some popular variants of operator splitting methods for tackling Problem (1). Initially, two natural splitting strategies are used in the literature:

$$\min_{\mathbf{X}, \mathbf{y}} F_1(\mathbf{X}, \mathbf{y}) \triangleq f(\mathbf{X}) - g(\mathbf{X}) + h(\mathbf{y}) + \mathcal{I}_{\mathcal{M}}(\mathbf{X}), \text{ s.t. } \mathcal{A}(\mathbf{X}) = \mathbf{y} \quad (2)$$

$$\min_{\mathbf{X}, \mathbf{Y}} F_2(\mathbf{X}, \mathbf{Y}) \triangleq f(\mathbf{X}) - g(\mathbf{X}) + h(\mathcal{A}(\mathbf{X})) + \mathcal{I}_{\mathcal{M}}(\mathbf{Y}), \text{ s.t. } \mathbf{X} = \mathbf{Y}. \quad (3)$$

*(a)* Smoothing Proximal Gradient Methods (SPGM, (Beck & Rosset, 2023; Böhm & Wright, 2021)) incorporate a penalty (or smoothing) parameter  $\mu \rightarrow 0$  to penalize the squared error in the constraints, resulting in the subsequent minimization problem (Beck & Rosset, 2023; Böhm & Wright, 2021; Chen, 2012):  $\min_{\mathbf{X}, \mathbf{y}} F_1(\mathbf{X}, \mathbf{y}) + \frac{1}{2\mu} \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|_2^2$ . During each iteration, SPGM employs proximal gradient strategies to alternatively minimize *w.r.t.*  $\mathbf{X}$  and  $\mathbf{y}$ . *(b)* Splitting Orthogonality Constraints Methods (SOCM, (Lai & Osher, 2014)) use the following iteration scheme:  $\mathbf{X}^{t+1} \approx \arg \min_{\mathbf{X}} F_2(\mathbf{X}, \mathbf{Y}^t) + \langle \mathbf{Z}^t, \mathbf{X} - \mathbf{Y}^t \rangle + \frac{\beta}{2} \|\mathbf{X} - \mathbf{Y}^t\|_F^2$ ,  $\mathbf{Y}^{t+1} \in \arg \min_{\mathbf{Y}} F_2(\mathbf{X}^{t+1}, \mathbf{Y}) + \langle \mathbf{Z}^t, \mathbf{X}^{t+1} - \mathbf{Y} \rangle + \frac{\beta}{2} \|\mathbf{X}^{t+1} - \mathbf{Y}\|_F^2$ , and  $\mathbf{Z}^{t+1} = \mathbf{Z}^t + \beta(\mathbf{X}^{t+1} - \mathbf{Y}^{t+1})$ , where  $\beta$  is a fixed penalty constant, and  $\mathbf{Z}^t$  is the multiplier associated with the constraint  $\mathbf{X} = \mathbf{Y}$  at

Table 1: Comparison of existing methods for solving Problem (1).

Reference	$h(\mathcal{A}(\mathbf{X}))$	$g(\mathbf{X})$	Notable Features	Complexity	Conv. Rate
SOCM (Lai & Osher, 2014)	convex $h(\cdot)$	empty	$\sigma = 1, \alpha = 0$	unknown	unknown
MADMM (Kovnatsky et al., 2016)	convex $h(\cdot)$	empty	$\sigma = 1, \alpha = 0$	unknown	unknown
RSG (Li et al., 2021)	weakly convex $h(\cdot)$	empty	–	$\mathcal{O}(\epsilon^{-4})$	unknown
ManPG (Chen et al., 2020)	$h(\mathcal{A}(\mathbf{X})) = \ \mathbf{X}\ _1$	empty	hard subproblem	$\mathcal{O}(\epsilon^{-2})$	unknown
OBCD (Yuan, 2023)	separable $h(\cdot)$	empty	hard subproblem	$\mathcal{O}(\epsilon^{-2})$	unknown
RADMM (Li et al., 2022)	convex $h(\cdot)$	empty	$\sigma = 1, \alpha = 0$	$\mathcal{O}(\epsilon^{-4})$	unknown
ManIAL (Deng et al., 2024)	convex $h(\cdot)$	empty	inexact subproblem	$\mathcal{O}(\epsilon^{-3})$	unknown
SPGM (Beck & Rosset, 2023)	convex $h(\cdot)$	empty	–	$\mathcal{O}(\epsilon^{-3})$	unknown
OADM-EP[ours]	weakly convex $h(\cdot)$	convex	$\sigma \in [1, 2], \alpha > 0$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(1/\exp(T^u)), \dot{u} \in (0, \frac{2}{3}]$ *
OADM-RR[ours]	weakly convex $h(\cdot)$	convex	$\sigma \in [1, 2], \text{MBB}$	$\mathcal{O}(\epsilon^{-3})$	or $\mathcal{O}(1/T^{\ddot{u}}), \ddot{u} \in (0, +\infty)$ †

Note \*: This is known as super-exponential convergence, please refer to Theorem 5.9(a) for more details.

Note †: This is known as polynomial convergence, please refer to Theorem 5.9(b) for more details.

iteration  $t$ . (c) Similarly, Manifold ADMM (MADMM, (Kovnatsky et al., 2016)) iterates as follows:  $\mathbf{X}^{t+1} \approx \arg \min_{\mathbf{X}} F_1(\mathbf{X}, \mathbf{y}^t) + \langle \mathbf{z}^t, \mathcal{A}(\mathbf{X}) - \mathbf{y}^t \rangle + \frac{\beta}{2} \|\mathcal{A}(\mathbf{X}) - \mathbf{y}^t\|_F^2$ ,  $\mathbf{y}^{t+1} \in \arg \min_{\mathbf{y}} F_1(\mathbf{X}^{t+1}, \mathbf{y}) + \langle \mathbf{z}^t, \mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y} \rangle + \frac{\beta}{2} \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}\|_F^2$ , and  $\mathbf{z}^{t+1} = \mathbf{z}^t + \beta(\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1})$ , where  $\mathbf{z}^t$  is the multiplier associated with the constraint  $\mathcal{A}(\mathbf{X}) - \mathbf{y} = \mathbf{0}$  at iteration  $t$ . (d) Like MADMM, Riemannian ADMM (RADMM, (Li et al., 2022)) operates using the first splitting strategy in Equation (2). In contrast, it employs a Riemannian retraction strategy to solve the  $\mathbf{X}$ -subproblem and a Moreau envelope smoothing strategy to solve the  $\mathbf{y}$ -subproblem.

**Contributions.** We compare existing methods for solving Problem (1) in Table 1, and our main contributions are summarized as follows. (i) We introduce OADM, a specialized ADMM designed for structured nonsmooth composite optimization problems under orthogonality constraints in Problem (1). Two specific variants of OADM are explored: one based on Euclidean Projection (OADM-EP) and the other on Riemannian Retraction (OADM-RR). Notably, while many existing works primarily address cases where  $g(\mathbf{X}) = 0$  and  $h(\cdot)$  is convex, our approach considers a more general setting where  $h(\cdot)$  is weakly convex and  $g(\mathbf{X})$  is convex. (ii) OADM could demonstrate fast convergence by incorporating Nesterov’s extrapolation (Nesterov, 2003) into OADM-EP and a Monotone Barzilai-Borwein (MBB) stepsize strategy (Wen & Yin, 2013) into OADM-RR to potentially accelerate primal convergence. Both variants also employ an over-relaxation strategy to enhance dual convergence (Gonçalves et al., 2017; Yang et al., 2017; Li et al., 2016). (iii) By introducing a novel Lyapunov function, we establish the convergence of OADM to critical points of Problem (1) within an oracle complexity of  $\mathcal{O}(1/\epsilon^3)$ , matching the best-known results to date (Beck & Rosset, 2023; Böhm & Wright, 2021). This is achieved through a decreasing step size for updating primal and dual variables. In contrast, RADMM employs a small constant step size for such updates, resulting in a sub-optimal oracle complexity of  $\mathcal{O}(\epsilon^{-4})$  (Li et al., 2022). (iv) We establish a super-exponential convergence rate or polynomial convergence rate for OADM, depending on the specific setting, under the Kurdyka-Lojasiewicz (KL) inequality, providing *the first non-ergodic convergence result* for this class of non-convex nonsmooth optimization problems.

## 2 TECHNICAL PRELIMINARIES

This section provides some technical preliminaries on Moreau envelopes for weakly convex functions and manifold optimization.

**Notations.** We define  $[n] \triangleq \{1, 2, \dots, n\}$ . We use  $\mathcal{A}^T(\cdot)$  to denote the adjoint operator of  $\mathcal{A}(\cdot)$  with  $\langle \mathcal{A}(\mathbf{X}), \mathbf{z} \rangle = \langle \mathbf{X}, \mathcal{A}^T(\mathbf{z}) \rangle$  for all  $\mathbf{X} \in \mathbb{R}^{n \times r}$  and  $\mathbf{z} \in \mathbb{R}^m$ . We define  $\bar{\mathbf{A}} \triangleq \max_{\mathbf{V}} \|\mathcal{A}(\mathbf{V})\|_F / \|\mathbf{V}\|_F$ . We use  $\mathcal{I}_{\mathcal{M}}(\mathbf{X})$  to denote the indicator function of orthogonality constants. Further notations, technical preliminaries, and relevant lemmas are detailed in Appendix Section A.

### 2.1 MOREAU ENVELOPES FOR WEAKLY CONVEX FUNCTIONS

We provide the following useful definition.

**Definition 2.1.** For a proper convex, and Lipschitz continuous function  $h(\mathbf{y}) : \mathbb{R}^m \mapsto \mathbb{R}$ , the Moreau envelope of  $h(\mathbf{y})$  with the parameter  $\mu > 0$  is given by  $h_\mu(\mathbf{y}) \triangleq \min_{\check{\mathbf{y}}} h(\check{\mathbf{y}}) + \frac{1}{2\mu} \|\check{\mathbf{y}} - \mathbf{y}\|_2^2$ .

We show some useful properties of Moreau envelope for weakly convex functions.

**Lemma 2.2.** *Let  $h : \mathbb{R}^m \mapsto \mathbb{R}$  to be a proper,  $W_h$ -weakly convex, and lower semicontinuous function. Assume  $\mu \in (0, W_h^{-1})$ . We have the following results (Böhm & Wright, 2021). (a) The function  $h_\mu(\cdot)$  is  $C_h$ -Lipschitz continuous. (b) The function  $h_\mu(\cdot)$  is continuously differentiable with gradient  $\nabla h_\mu(\mathbf{y}) = \frac{1}{\mu}(\mathbf{y} - \mathbb{P}_\mu(\mathbf{y}))$  for all  $\mathbf{y}$ , where  $\mathbb{P}_\mu(\mathbf{y}) \triangleq \arg \min_{\check{\mathbf{y}}} h(\check{\mathbf{y}}) + \frac{1}{2\mu} \|\check{\mathbf{y}} - \mathbf{y}\|_2^2$ . This gradient is  $\max(\mu^{-1}, \frac{W_h}{1-\mu W_h})$ -Lipschitz continuous. In particular, when  $\mu \in (0, \frac{1}{2W_h}]$ , the condition  $\mu^{-1} \geq \frac{W_h}{1-\mu W_h}$  ensures that  $h_\mu(\mathbf{y})$  is  $(\mu^{-1})$ -smooth and  $(\mu^{-1})$ -weakly convex.*

**Lemma 2.3.** *(Proof in Appendix B.1) Assume  $0 < \mu_2 < \mu_1 < \frac{1}{W_h}$ , and fixing  $\mathbf{y} \in \mathbb{R}^m$ . We have:  $0 \leq h_{\mu_2}(\mathbf{y}) - h_{\mu_1}(\mathbf{y}) \leq \min\{\frac{\mu_1}{2\mu_2}, 1\} \cdot (\mu_1 - \mu_2) C_h^2$ .*

**Lemma 2.4.** *(Proof in Appendix B.2) Assume  $0 < \mu_2 < \mu_1 \leq \frac{1}{2W_h}$ , and fixing  $\mathbf{y} \in \mathbb{R}^m$ . We have:  $\|\nabla h_{\mu_1}(\mathbf{y}) - \nabla h_{\mu_2}(\mathbf{y})\| \leq (\frac{\mu_1}{\mu_2} - 1) C_h$ .*

**Lemma 2.5.** *(Proof in Appendix B.3) Assume that  $h(\mathbf{y})$  is  $W_h$ -weakly convex,  $\mu \in (0, \frac{1}{2W_h}]$ ,  $\beta > \mu^{-1}$ . Consider the following strongly convex optimization problem:  $\bar{\mathbf{y}} = \arg \min_{\mathbf{y}} h_\mu(\mathbf{y}) + \frac{\beta}{2} \|\mathbf{y} - \mathbf{b}\|_2^2$ , which is equivalent to:  $(\bar{\mathbf{y}}, \check{\mathbf{y}}) = \arg \min_{\mathbf{y}, \mathbf{y}'}$   $h(\mathbf{y}') + \frac{1}{2\mu} \|\mathbf{y}' - \mathbf{y}\|_2^2 + \frac{\beta}{2} \|\mathbf{y} - \mathbf{b}\|_2^2$ . We have: (a)  $\bar{\mathbf{y}} = \frac{(\check{\mathbf{y}} + \mu\beta\mathbf{b})}{1+\mu\beta}$ , where  $\check{\mathbf{y}} = \arg \min_{\mathbf{y}} h(\mathbf{y}) + \frac{\beta}{2(1+\mu\beta)} \|\mathbf{y} - \mathbf{b}\|_2^2 = \mathbb{P}_{[\mu+1/\beta]}(\mathbf{b})$ . (b)  $\beta(\mathbf{b} - \bar{\mathbf{y}}) \in \partial h(\check{\mathbf{y}})$ . (c)  $\|\bar{\mathbf{y}} - \check{\mathbf{y}}\| \leq \mu C_h$ .*

**Remark 2.6.** (i) Lemmas 2.3 and 2.4 presented in this paper are novel. (ii) The upper bound in Lemma 2.3 is slightly better than the bound established in Lemma 4.1 of (Böhm & Wright, 2021). (iii) Lemma 2.5 is very critical in our algorithm development and theoretical analysis.

## 2.2 MANIFOLD OPTIMIZATION

We define the  $\epsilon$ -stationary point of Problem (1) as follows.

**Definition 2.7.** (First-Order Optimality Conditions, (Chen et al., 2020; Li et al., 2022; Beck & Rosset, 2023)) The solution  $(\check{\mathbf{X}}, \check{\mathbf{y}}, \check{\mathbf{z}})$  with  $\check{\mathbf{X}} \in \mathcal{M}$  is called an  $\epsilon$ -stationary point of Problem (1) if:  $\text{Crit}(\check{\mathbf{X}}, \check{\mathbf{y}}, \check{\mathbf{z}}) \leq \epsilon$ , where  $\text{Crit}(\mathbf{X}, \mathbf{y}, \mathbf{z}) \triangleq \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\| + \|\partial h(\mathbf{y}) - \mathbf{z}\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\nabla f(\mathbf{X}) - \partial g(\mathbf{X}) + \mathcal{A}^\top(\mathbf{z}))\|_F$ . Here, according to (Absil et al., 2008a), for all  $\mathbf{X} \in \mathcal{M}$  and  $\Delta \in \mathbb{R}^{n \times r}$ , we have:  $\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\Delta) = \Delta - \frac{1}{2}\mathbf{X}(\Delta^\top \mathbf{X} + \mathbf{X}^\top \Delta)$ .

The proposed algorithm is an iterative procedure. After shifting the current iterate  $\mathbf{X} \in \mathcal{M}$  in the search direction, it may no longer reside on  $\mathcal{M}$ . Therefore, we must retract the point onto  $\mathcal{M}$  to form the next iterate. The following definition is useful in this context.

**Definition 2.8.** A retraction on  $\mathcal{M}$  is a smooth map (Absil et al., 2008a):  $\text{Retr}_{\mathbf{X}}(\Delta) \in \mathcal{M}$  with  $\mathbf{X} \in \mathcal{M}$  and  $\Delta \in \mathbb{R}^{n \times r}$  satisfying  $\text{Retr}_{\mathbf{X}}(\mathbf{0}) = \mathbf{X}$ , and  $\lim_{\mathbf{T}_{\mathbf{X}}\mathcal{M} \ni \Delta \rightarrow \mathbf{0}} \frac{\|\text{Retr}_{\mathbf{X}}(\Delta) - \mathbf{X} - \Delta\|_F}{\|\Delta\|_F} = 0$  for any  $\mathbf{X} \in \mathcal{M}$ .

**Remark 2.9.** Several retractions on the Stiefel manifold have been explored in literature (Absil & Mallick, 2012; Absil et al., 2008b). We present two examples below. (i) Polar Decomposition-Based Retraction:  $\text{Retr}_{\mathbf{X}}(\Delta) = (\mathbf{X} + \Delta)(\mathbf{I}_r + \Delta^\top \Delta)^{-1/2}$ . (ii) QR-Decomposition-Based Retraction:  $\text{Retr}_{\mathbf{X}}(\Delta) = \text{qf}(\mathbf{X} + \Delta)$ , where  $\text{qf}(\mathbf{X})$  is the Q-factor in the thin QR-decomposition of  $\mathbf{X}$ .

The following lemma concerning the retraction operator is useful for our subsequent analysis.

**Lemma 2.10.** ((Boumal et al., 2019)) Let  $\mathbf{X} \in \mathcal{M}$  and  $\Delta \in \mathbf{T}_{\mathbf{X}}\mathcal{M}$ . There exists positive constants  $\{\check{k}, \check{k}\}$  such that  $\|\text{Retr}_{\mathbf{X}}(\Delta) - \mathbf{X}\|_F \leq \check{k}\|\Delta\|_F$ , and  $\|\text{Retr}_{\mathbf{X}}(\Delta) - \mathbf{X} - \Delta\|_F \leq \frac{1}{2}\check{k}\|\Delta\|_F^2$ .

Furthermore, we present the following three insightful lemmas.

**Lemma 2.11.** (Proof in Appendix B.4) Let  $\mathbf{X} \in \mathcal{M}$  and  $\Delta \in \mathbb{R}^{n \times r}$ , we have  $\|\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\Delta)\|_F \leq \|\Delta\|_F$ .

**Lemma 2.12.** (Proof in Appendix B.5) Let  $\rho > 0$ ,  $\mathbf{G} \in \mathbb{R}^{n \times r}$ , and  $\mathbf{X} \in \mathcal{M}$ . We define  $\mathbb{G}_\rho \triangleq \mathbf{G} - \rho \mathbf{X} \mathbf{G}^\top \mathbf{X} - (1-\rho) \mathbf{X} \mathbf{X}^\top \mathbf{G}$ . It follows that: (a)  $\max(1, 2\rho) \cdot \langle \mathbf{G}, \mathbb{G}_\rho \rangle \geq \|\mathbb{G}_\rho\|_F^2 \geq \min(1, \rho^2) \|\mathbb{G}_1\|_F^2$ . (b)  $\min(1, 2\rho) \|\mathbb{G}_{1/2}\|_F \leq \|\mathbb{G}_\rho\|_F \leq \max(1, 2\rho) \|\mathbb{G}_{1/2}\|_F$ .

**Lemma 2.13.** (Proof in Appendix B.6) Consider the following optimization problem:  $\min_{\mathbf{X} \in \mathcal{M}} f(\mathbf{X})$ , where  $f(\mathbf{X})$  is differentiable. For all  $\mathbf{X} \in \mathcal{M}$ , we have:  $\text{dist}(\mathbf{0}, \partial I_{\mathcal{M}}(\mathbf{X}) + \nabla f(\mathbf{X})) \leq \|\nabla f(\mathbf{X}) - \mathbf{X} \nabla f(\mathbf{X})^{\top} \mathbf{X}\|_{\text{F}}$ .

**Remark 2.14.** The matrix  $\mathbb{G}_{\rho} \in \mathbb{R}^{n \times r}$  in Lemma 2.12 is closely related to the search descent direction of the proposed OADMM-RR algorithm. While one can set  $\rho$  to typical values such as 1 or 1/2, we consider the setting  $\rho \in (0, \infty)$  to enhance the versatility of OADMM-RR, aligning with (Liu et al., 2016; Jiang & Dai, 2015).

### 3 THE PROPOSED OADMM ALGORITHM

This section provides the proposed OADMM algorithm for solving Problem (1), featuring two variants, one is based on Euclidean Projection (OADMM-EP) and the other on Riemannian Retraction (OADMM-RR).

Using the Moreau envelope smoothing technique, we consider the following optimization problem:

$$\min_{\mathbf{X}, \mathbf{y}} f(\mathbf{X}) - g(\mathbf{X}) + h_{\mu}(\mathbf{y}) + \mathcal{I}_{\mathcal{M}}(\mathbf{X}), \text{ s.t. } \mathcal{A}(\mathbf{X}) = \mathbf{y}, \quad (4)$$

where  $\mu \rightarrow 0$ , and  $h_{\mu}(\mathbf{y})$  is the Moreau Envelope of  $h(\mathbf{y})$ . Importantly,  $h_{\mu}(\mathbf{y})$  is  $(\mu^{-1})$ -smooth when  $\mu \leq \frac{1}{2W_h}$ , according to Lemma 2.2. It is worth noting that similar smoothing techniques have been used in the design of augmented Lagrangian methods (Zeng et al., 2022), and minimax optimization (Zhang et al., 2020a), and ADMMs (Li et al., 2022). We define the augmented Lagrangian function of Problem (4) as follows:

$$\mathcal{L}(\mathbf{X}, \mathbf{y}; \mathbf{z}; \beta, \mu) = \underbrace{f(\mathbf{X}) + \langle \mathbf{z}, \mathcal{A}(\mathbf{X}) - \mathbf{y} \rangle + \frac{\beta}{2} \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|_2^2}_{\triangleq \mathcal{S}(\mathbf{X}, \mathbf{y}; \mathbf{z}; \beta)} - g(\mathbf{X}) + h_{\mu}(\mathbf{y}) + \mathcal{I}_{\mathcal{M}}(\mathbf{X}). \quad (5)$$

Here,  $\mathbf{z}$  is the dual variable for the equality constraint,  $\mu$  is the smoothing parameter linked to the function  $h(\mathbf{y})$ ,  $\beta$  is the penalty parameter associated with the equality constraint, and  $\mathcal{I}_{\mathcal{M}}(\mathbf{X})$  is the indicator function of the set  $\mathcal{M}$ .

In simple terms, OADMM updates are performed by minimizing the augmented Lagrangian function  $\mathcal{L}(\mathbf{X}, \mathbf{y}, \mathbf{z}; \beta, \mu)$  over the primal variables  $\{\mathbf{X}^t, \mathbf{y}^t\}$  at each iteration, while keeping all other primal and dual variables fixed. The dual variables are updated using gradient ascent on the dual problem.

For updating the primal variable  $\mathbf{X}$ , we use different strategies, resulting in distinct variants of OADMM. We first observe that the function  $\mathcal{S}(\mathbf{X}, \mathbf{y}^t; \mathbf{z}^t; \beta^t)$  is  $\ell(\beta^t)$ -smooth w.r.t.  $\mathbf{X}$ , where  $\ell(\beta^t) \triangleq \beta^t \bar{\Lambda}^2 + L_f$ . In OADMM-EP, we adopt a proximal linearized method based on Euclidean projection (Lai & Osher, 2014), while in OADMM-RR, we apply line-search methods on the Stiefel manifold (Liu et al., 2016).

We detail iteration steps of OADMM in Algorithm 1, and have the following remarks.

- (a) To achieve possible faster dual convergence, we apply an over-relaxation step size with  $\sigma \in (1, 2)$  for updating the dual variable  $\mathbf{z}$ , as suggested by previous studies (Gonçalves et al., 2017; Yang et al., 2017; Li et al., 2016; 2023).
- (b) To accelerate primal convergence in OADMM-EP, we incorporate a Nesterov extrapolation strategy with parameter  $\alpha \in (0, 1)$ .
- (c) To enhance primal convergence in OADMM-RR, we use a Monotone Barzilai-Borwein (MBB) strategy (Wen & Yin, 2013) with a dynamically adjusted parameter  $b^t$  to capture the problem's curvature<sup>1</sup>. The parameters  $\{\gamma, \delta\}$  represent the decay rate and sufficient decrease parameter, commonly used in line search procedures (Chen et al., 2020).
- (d) The  $\mathbf{X}$ -subproblem is solved as:  $\mathbf{X}^{t+1} = \arg \min_{\mathbf{X} \in \mathcal{M}} \|\mathbf{X} - \mathbf{X}'\|_{\text{F}}^2 = \dot{\mathbf{U}} \dot{\mathbf{V}}^{\top}$ , where  $\mathbf{X}' = \mathbf{X}_c^t - \mathbf{G}^t / (\theta \ell(\beta^t))$ , and  $\dot{\mathbf{U}} \text{diag}(\dot{\mathbf{x}}) \dot{\mathbf{V}}^{\top} = \mathbf{X}'$  is the using singular value decomposition of  $\mathbf{X}'$ .
- (e) The  $\mathbf{y}$ -subproblem can be solved using the result from Lemma 2.5.
- (f) For practical implementation, we recommend the following default parameters:  $p = 1/3$ ,  $\theta = 1.01$ ,  $\sigma = 1.1$ ,  $\rho = 1$ ,  $\gamma = 1/2$ ,  $\delta = 10^{-3}$ ,  $\xi = 1$ ,  $\alpha = \frac{\theta-1}{(\theta+1)(\xi+2)} - 10^{-12}$ .

<sup>1</sup>Following (Wen & Yin, 2013), one can set  $b^t = \langle \mathbf{S}^t, \mathbf{S}^t \rangle / \langle \mathbf{S}^t, \mathbf{Z}^t \rangle$  or  $b^t = \langle \mathbf{S}^t, \mathbf{Z}^t \rangle / \langle \mathbf{Z}^t, \mathbf{Z}^t \rangle$ , where  $\mathbf{S}^t = \mathbf{X}^t - \mathbf{X}^{t-1}$  and  $\mathbf{Z}^t = \mathbb{G}_1^{t-1} - \mathbb{G}_1^t$ , with  $\mathbb{G}_1^t$  being the Riemannian gradient.

**Algorithm 1: OADMM: The Proposed ADMM for Solving Problem (1).****Initialization:**

Choose  $\{\mathbf{X}^0, \mathbf{y}^0, \mathbf{z}^0\}$ . Choose  $p \in (0, 1)$ ,  $\xi \in (0, \infty)$ ,  $\theta \in (1, \infty)$ ,  $\sigma \in [1, 2)$ .

Choose  $\chi \in (1 + 4\omega\ddot{\sigma}, \infty)$ , where  $\omega \triangleq \frac{1}{\sigma} + \frac{\xi}{2\sigma^2} + \frac{\varepsilon_z}{\sigma^2}$ ,  $\ddot{\sigma} \triangleq (\sigma/(2-\sigma))^2$ ,  $\varepsilon_z = \xi$ .

Choose  $\beta^0$  sufficiently large such that  $\beta^0 \geq 2\chi W_h$ .

For OADMM-EP, choose  $\alpha \in [0, \frac{\theta-1}{(\theta+1)(\xi+2)})$ .

For OADMM-RR, choose  $\alpha = 0$ ,  $\rho \in (0, \infty)$ ,  $\gamma \in (0, 1)$ ,  $\delta \in (0, \frac{1}{\max(1, 2\rho)})$ .

**for**  $t$  **from** 0 **to**  $T$  **do**

S1) Set  $\beta^t = \beta^0(1 + \xi t^p)$ ,  $\mu^t = \chi/\beta^t$ .

S2) Update the primal variable  $\mathbf{X}$ : **if** OADMM-EP **then**

Set  $\mathbf{X}_c^t = \mathbf{X}^t + \alpha(\mathbf{X}^t - \mathbf{X}^{t-1})$ ,  $\mathbf{G}^t \in \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}_c^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) - \partial g(\mathbf{X}^t)$ .

$\mathbf{X}^{t+1} \in \arg \min_{\mathbf{X} \in \mathcal{M}} \langle \mathbf{X} - \mathbf{X}^t, \mathbf{G}^t \rangle + \frac{\theta \ell(\beta^t)}{2} \|\mathbf{X} - \mathbf{X}_c^t\|_F^2$ , where  $\ell(\beta^t) \triangleq \beta^t \bar{\mathbf{A}}^2 + L_f$ .

**end**

**if** OADMM-RR **then**

Set  $\mathbf{G}^t \in \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) - \partial g(\mathbf{X}^t)$ ,  $\dot{\mathcal{L}}(\mathbf{X}) \triangleq L(\mathbf{X}, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^t)$ . Set

$\mathbb{G}_\rho^t \triangleq \mathbf{G}^t - \rho \mathbf{X}^t [\mathbf{G}^t]^\top \mathbf{X}^t - (1 - \rho) \mathbf{X}^t [\mathbf{X}^t]^\top \mathbf{G}^t$ . Set  $b^t \in (\underline{b}, \bar{b})$  as the BB step size,

where  $\underline{b}, \bar{b} \in (0, \infty)$ . Set  $\mathbf{X}^{t+1} = \text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t)$ , where  $\eta^t \triangleq \frac{b^t \gamma^j}{\beta^t}$ , and

$j \in \{0, 1, 2, \dots\}$  is the smallest integer that:

$\dot{\mathcal{L}}(\text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t)) - \dot{\mathcal{L}}(\mathbf{X}^t) \leq -\delta \eta^t \|\mathbb{G}_\rho^t\|_F^2$ .

**end**

S3) Update the primal variable  $\mathbf{y}$ :  $\mathbf{y}^{t+1} = \arg \min_{\mathbf{y}} h_{\mu^t}(\mathbf{y}) + \frac{\beta^t}{2} \|\mathbf{y} - \mathbf{b}\|_2^2$ , where

$\mathbf{b} \triangleq \mathbf{y}^t - \frac{1}{\beta^t} \nabla_{\mathbf{y}} \mathcal{S}(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t; \beta^t)$ . It can be solved as:  $\mathbf{y}^{t+1} = \frac{\check{\mathbf{y}}^{t+1} + \mu^t \beta^t \mathbf{b}}{1 + \mu^t \beta^t}$ , where

$\check{\mathbf{y}}^{t+1} = \mathbb{P}_{[\mu^t + 1/\beta^t]}(\mathbf{b})$ .

S4) Update the dual variable  $\mathbf{z}$ :  $\mathbf{z}^{t+1} = \mathbf{z}^t + \sigma \beta^t (\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1})$

**end**

## 4 ORACLE COMPLEXITY

This section details the oracle complexity of Algorithm 1.

We define  $\varepsilon_z = \xi$ ,  $\varepsilon_y \triangleq \frac{1}{2}(1 - \frac{1+4\omega\ddot{\sigma}}{\chi})$ ,  $\dot{\sigma} \triangleq (\sigma - 1)/(2 - \sigma)$ ,  $\ddot{\sigma} \triangleq (\sigma/(2 - \sigma))^2$ ,  $\omega \triangleq \frac{1}{\sigma} + \frac{\xi}{2\sigma^2} + \frac{\varepsilon_z}{\sigma^2}$ .

We define the potential function (or Lyapunov function) for all  $t \geq 1$ , as follows:

$$\begin{aligned} \Theta^t &\triangleq \Theta(\mathbf{X}^t, \mathbf{X}^{t-1}, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \beta^{t-1}, \mu^{t-1}, t) \\ &\triangleq L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \mu^{t-1}) + \mu^{t-1} C_h^2 + \mathbb{T}^t + \mathbb{Z}^t + \mathbb{X}^t, \end{aligned} \quad (6)$$

where  $\mathbb{T}^t \triangleq \frac{4\omega\ddot{\sigma}}{\beta^0} C_h^2 \frac{1}{t}$ ,  $\mathbb{Z}^t \triangleq \omega \dot{\sigma} \sigma^2 \beta^{t-1} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|_2^2$ , and  $\mathbb{X}^t \triangleq \frac{\alpha(\theta+1)\ell(\beta^t)}{2} \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F^2$ .

Additionally, we define:

$$e^t \triangleq \begin{cases} \|\mathbf{y}^t - \mathbf{y}^{t-1}\| + \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F, & \text{OADMM-EP;} \\ \|\mathbf{y}^t - \mathbf{y}^{t-1}\| + \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \frac{1}{\beta^t} \|\mathbb{G}_{1/2}^{t-1}\|_F, & \text{OADMM-RR.} \end{cases} \quad (7)$$

We have the following useful lemma, derived using the first-order optimality condition of  $\mathbf{y}^{t+1}$ .

**Lemma 4.1.** (Proof in Section C.1, Bounding Dual using Primal) We have: (a)  $\forall t \geq 0$ ,  $\mathbf{z}^t - \frac{1}{\sigma}(\mathbf{z}^t - \mathbf{z}^{t+1}) = \nabla h_{\mu^t}(\mathbf{y}^{t+1}) \in \partial h(\check{\mathbf{y}}^{t+1})$ . (b)  $\forall t \geq 1$ ,  $\|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \leq \dot{\sigma}(\|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2 - \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2) + 2\ddot{\sigma}(\beta^t/\chi)^2 \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + 2\ddot{\sigma} C_h^2 (\frac{2}{t} - \frac{2}{t+1})$ .

**Remark 4.2.** Here, for OADMM-RR, we set  $\alpha = 0$ , resulting in  $\mathbb{X}^t = 0$  for all  $t$ . (i) With the choice  $\sigma = 1$ , we have:  $\nabla h_{\mu^{t-1}}(\mathbf{y}^t) = \mathbf{z}^t$ , and  $\|\mathbf{z}^{t+1} - \mathbf{z}^t\| \leq \|\nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \nabla h_{\mu^{t-1}}(\mathbf{y}^t)\|$ .

**Lemma 4.3.** (Proof in Appendix C.2) (a) It holds that  $\beta^{t+1} \leq \beta^t(1 + \xi)$ . (b) There exists constant  $\{\underline{\ell}, \bar{\ell}\}$  such that  $\beta^t \underline{\ell} \leq \ell(\beta^t) \leq \beta^t \bar{\ell}$ .

The subsequent lemma demonstrates that the sequence  $\{\Theta^t\}_{t=1}^\infty$  is always lower bounded.

**Lemma 4.4.** (Proof in Section C.3) For all  $t \geq 1$ , there exists constants  $\{\bar{X}, \bar{z}, \bar{y}, \underline{\Theta}\}$  such that  $\|\mathbf{X}^t\|_F \leq \bar{X}$ ,  $\|\mathbf{z}^t\| \leq \bar{z}$ ,  $\|\mathbf{y}^t\| \leq \bar{y}$ , and  $\Theta^t \geq \underline{\Theta}$ .

The following lemma is useful for our subsequent analysis, applicable to both OADMM-EP and OADMM-RR.

**Lemma 4.5.** (Proof in Appendix C.4, Sufficient Decrease for Variables  $\{\mathbf{y}, \mathbf{z}, \beta, \mu\}$ ) We have  $L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^{t+1}; \beta^{t+1}, \mu^t) - L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \mu^{t-1}) + (\mu^t - \mu^{t-1})C_h^2 + \mathbb{T}^{t+1} - \mathbb{T}^t + \mathbb{Z}^{t+1} - \mathbb{Z}^t + \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 \leq -\varepsilon_y \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2$ .

In the remaining content of this section, we provide separate analyses for OADMM-EP and OADMM-RR.

#### 4.1 ANALYSIS FOR OADMM-EP

Using the optimality condition of  $\mathbf{X}^{t+1}$ , we derive the following lemma.

**Lemma 4.6.** (Proof in Appendix C.5, Sufficient Decrease for Variable  $\mathbf{X}$ ) We define  $\varepsilon_x \triangleq \frac{1}{2}\varepsilon'_x \ell$ , where  $\varepsilon'_x \triangleq \theta - 1 - \alpha(2 + \xi)(1 + \theta) > 0$ . We have  $L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \mu^{t-1}) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \mu^{t-1}) \leq -\varepsilon_x \beta^t \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 + \mathbb{X}^t - \mathbb{X}^{t+1}$ .

Combining the results from Lemmas 4.5, and 4.6, we arrive at the following lemma.

**Lemma 4.7.** (Proof in Appendix C.6) We have: (a)  $\beta^t \{\varepsilon_z \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + \varepsilon_y \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \varepsilon_x \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2\} \leq \Theta^t - \Theta^{t+1}$ . (b)  $\frac{1}{T} \sum_{t=1}^T \beta^t e^{t+1} \leq \mathcal{O}(T^{(p-1)/2})$ .

Finally, we have the following theorem regarding the oracle complexity of OADMM-EP.

**Theorem 4.8.** (Proof in Appendix C.7) Let  $p = 1/3$ . We have:  $\frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{X}^{t+1}, \check{\mathbf{y}}^{t+1}, \mathbf{z}^{t+1}) \leq \mathcal{O}(T^{-1/3})$ . In other words, there exists  $\bar{t} \leq T$  such that:  $\frac{1}{\bar{t}} \sum_{t=1}^{\bar{t}} \text{Crit}(\mathbf{X}^{t+1}, \check{\mathbf{y}}^{t+1}, \mathbf{z}^{t+1}) \leq \epsilon$ , provided that  $T \geq \mathcal{O}(1/\epsilon^3)$ .

**Remark 4.9.** The oracle complexity of OADMM-EP matches the best-known complexities currently available to date (Beck & Rosset, 2023; Böhm & Wright, 2021).

#### 4.2 ANALYSIS FOR OADMM-RR

Using the properties of the line search procedure for updating the variable  $\mathbf{X}^{t+1}$ , we deduce the following lemma.

**Lemma 4.10.** (Proof in Appendix C.8, Sufficient Decrease for Variable  $\mathbf{X}$ ) We define  $\varepsilon_x \triangleq \delta \bar{\gamma} \gamma \bar{b} \min(1, 2\rho)^2 > 0$ , where  $\bar{\gamma} \triangleq 2(1/\max(1, 2\rho) - \delta)/(\bar{\ell} \bar{k} \bar{b} + \bar{g} \bar{k} \bar{b}/\beta^0) > 0$ . We have: (a) For any  $t \geq 0$ , if  $j$  is large enough such that  $\gamma^j \in (0, \bar{\gamma})$ , then the condition of the line search procedure is satisfied. (b) It follows that:  $L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \mu^t) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \mu^t) \leq -\frac{\varepsilon_x}{\beta^t} \|\mathbb{G}_{1/2}^t\|_F^2$ .

Here,  $\bar{g}$  is a constant that  $\|\mathbb{G}^t\|_F \leq \bar{g}$ ,  $\{\bar{k}, \check{k}\}$  are defined in Lemma 2.10, and  $\{\rho, \gamma, \delta, \bar{b}, \check{b}\}$  are defined in Algorithm 1.

**Remark 4.11.** By Lemma 4.10(a), since  $\bar{\gamma}$  is a universal constant and  $\gamma^j$  decreases exponentially, the line search procedure of OADMM-RR will terminate in  $\log(\bar{\gamma})/\log(\gamma) + 1 = \mathcal{O}(1)$  time.

Combining the results from Lemmas 4.5, and 4.10, we obtain the following lemma.

**Lemma 4.12.** (Proof in Appendix C.9) We have: (a)  $\beta^t \{\varepsilon_z \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + \varepsilon_y \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \varepsilon_x \|\frac{1}{\beta^t} \mathbb{G}_{1/2}^t\|_F^2\} \leq \Theta^t - \Theta^{t+1}$ . (b)  $\frac{1}{T} \sum_{t=1}^T \beta^t e^{t+1} \leq \mathcal{O}(T^{(p-1)/2})$ .

Finally, we derive the following theorem on the oracle complexity of OADMM-RR.

**Theorem 4.13.** (Proof in Appendix C.10) Let  $p = 1/3$ . We have:  $\frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{X}^{t+1}, \check{\mathbf{y}}^{t+1}, \mathbf{z}^{t+1}) \leq \mathcal{O}(T^{-1/3})$ . In other words, there exists  $\bar{t} \leq T$  such that:  $\frac{1}{\bar{t}} \sum_{t=1}^{\bar{t}} \text{Crit}(\mathbf{X}^{t+1}, \check{\mathbf{y}}^{t+1}, \mathbf{z}^{t+1}) \leq \epsilon$ , provided that  $T \geq \mathcal{O}(1/\epsilon^3)$ .

**Remark 4.14.** Theorem 4.13 mirrors Theorem 4.8, and OADMM-RR shares the same oracle complexity as OADMM-EP.

## 5 CONVERGENCE RATE

This section provides convergence rate of OADMM-EP and OADMM-RR. Our analyses are based on a non-convex analysis tool called KL inequality (Attouch et al., 2010; Bolte et al., 2014; Li & Lin, 2015; Li et al., 2023).

We define the Lyapunov function as:  $\Theta(\mathbf{X}, \mathbf{X}^-, \mathbf{y}, \mathbf{z}; \beta, \beta^-, \mu^-, t) \triangleq L(\mathbf{X}, \mathbf{y}, \mathbf{z}; \beta, \mu^-) + \omega \tilde{\sigma} \sigma^2 \beta^- \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|_2^2 + \frac{\alpha(\theta+1)\ell(\beta)}{2} \|\mathbf{X} - \mathbf{X}^-\|_F^2 + \frac{4\omega\tilde{\sigma}}{\beta^0} C_h^2 \frac{1}{t} + C_h^2 \mu^-$ , where we let  $\alpha = 0$  for OADMM-RR. We define  $\mathbb{w} \triangleq \{\mathbf{X}, \mathbf{X}^-, \mathbf{y}, \mathbf{z}\}$ ,  $\mathbb{w}^t \triangleq \{\mathbf{X}^t, \mathbf{X}^{t-1}, \mathbf{y}^t, \mathbf{z}^t\}$ ,  $\mathbb{u} \triangleq \{\beta, \beta^-, \mu^-, t\}$ , and  $\mathbb{u}^t \triangleq \{\beta^t, \beta^{t-1}, \mu^{t-1}, t\}$ . Thus, we have  $\Theta^t = \Theta(\mathbb{w}^t; \mathbb{u}^t)$ . We denote  $\mathbb{w}^\infty$  as a limiting point of Algorithm 1.

We make the following additional assumptions.

**Assumption 5.1.** (Kurdyka-Łojasiewicz Inequality (Attouch et al., 2010)). Consider a semi-algebraic function  $\Theta(\mathbb{w}^t; \mathbb{u}^t)$  w.r.t.  $\mathbb{w}^t$  for all  $t$ , where  $\mathbb{w}^t$  is in the effective domain of  $\Theta(\mathbb{w}^t; \mathbb{u}^t)$ . There exist  $\tilde{\eta} \in (0, +\infty)$ ,  $\tilde{\sigma} \in [0, 1)$ , a neighborhood  $\Upsilon$  of  $\mathbb{w}^\infty$ , and a continuous and concave desingularization function  $\varphi(s) \triangleq \tilde{c}s^{1-\tilde{\sigma}}$  with  $\tilde{c} > 0$  and  $s \in [0, \tilde{\eta})$  such that, for all  $\mathbb{w}^t \in \Upsilon$  satisfying  $\Theta(\mathbb{w}^t; \mathbb{u}^t) - \Theta(\mathbb{w}^\infty; \mathbb{u}^\infty) \in (0, \tilde{\eta})$ , it holds that:  $\varphi'(\Theta(\mathbb{w}^t; \mathbb{u}^t) - \Theta(\mathbb{w}^\infty; \mathbb{u}^\infty)) \cdot \text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbb{u}^t)) \geq 1$ . Here,  $\text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbb{u}^t)) \triangleq \{\text{dist}^2(\mathbf{0}, \partial_{\mathbf{X}}\Theta(\mathbb{w}^t; \mathbb{u}^t)) + \text{dist}^2(\mathbf{0}, \partial_{\mathbf{X}^-}\Theta(\mathbb{w}^t; \mathbb{u}^t)) + \text{dist}^2(\mathbf{0}, \partial_{\mathbf{y}}\Theta(\mathbb{w}^t; \mathbb{u}^t)) + \text{dist}^2(\mathbf{0}, \partial_{\mathbf{z}}\Theta(\mathbb{w}^t; \mathbb{u}^t))\}^{1/2}$ .

**Assumption 5.2.** The function  $g(\mathbf{X})$  is  $L_g$ -smooth such that  $\|\nabla g(\mathbf{X}) - \nabla g(\mathbf{X}')\|_F \leq L_g \|\mathbf{X} - \mathbf{X}'\|_F$  holds for all  $\mathbf{X} \in \mathcal{M}$  and  $\mathbf{X}' \in \mathcal{M}$ .

**Remark 5.3.** Semi-algebraic functions, including real polynomial functions, finite combinations, and indicator functions of semi-algebraic sets, commonly exhibit the KL property and find extensive use in applications (Attouch et al., 2010).

We present the following lemma regarding subgradient bounds for each iteration.

**Lemma 5.4.** (Proof in Section D.1, Subgradient Bounds) (a) For OADMM-EP, there exists a constant  $K > 0$  such that:  $\text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbb{u}^t)) \leq \beta^t K (e^t + e^{t-1})$ . (b) For OADMM-RR, there exists a constant  $K > 0$  such that:  $\text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbb{u}^t)) \leq \beta^t K e^t$ .

**Remark 5.5.** Lemma 5.4 significantly differs from prior work that used a constant penalty due to the crucial role played by the increasing penalty.

The following theorem establishes a finite length property of OADMM.

**Theorem 5.6.** (Proof in Section D.2, A Finite Length Property) We define  $d^t \triangleq \sum_{i=t}^\infty e^{i+1}$ . We define  $\varphi^t \triangleq \varphi(\Theta(\mathbb{w}^t; \mathbb{u}^t) - \Theta(\mathbb{w}^\infty; \mathbb{u}^\infty))$ , where  $\varphi(\cdot)$  is the desingularization function defined in Assumption 5.1. (a) We have the following recursive inequality for both OADMM-EP and OADMM-RR:  $(e^{t+1})^2 \leq (e^t + e^{t-1}) \cdot \dot{K}(\varphi^t - \varphi^{t+1})$ , where  $\dot{K} = \frac{3K}{\min(\varepsilon_z, \varepsilon_y, \varepsilon_x)}$ , and  $K$  is defined in Lemma 5.4. (b) It holds that  $\forall t \geq 1$ ,  $d^t \leq e^t + e^{t-1} + 4\dot{K}\varphi^t$ . The sequence  $\{\mathbb{w}^t\}_{t=1}^\infty$  has the finite length property that  $d^1 \leq e^1 + e^0 + 4\dot{K}\varphi^1 < +\infty$ .

**Remark 5.7.** The finite length property in Theorem 5.6 represents much stronger convergence results compared to those outlined in Theorems 4.8 and 4.13.

We prove a lemma demonstrating that the convergence of  $d^t \triangleq \sum_{i=t}^\infty e^{i+1}$  is sufficient to establish the convergence of  $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F$ .

**Lemma 5.8.** (Proof in Section D.3) We define  $d^t \triangleq \sum_{i=t}^\infty e^{i+1}$ . For both OADMM-EP and OADMM-RR, we have: (a) There exists a constant  $\check{c}$  such that  $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F \leq \check{c} \cdot d^t$ . (b) We have  $d^t \leq d^{t-2} - d^t + \ddot{K}[\beta^t(d^{t-2} - d^t)]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}$ , where  $\ddot{K} \triangleq 4\dot{K}\check{c} \cdot [\check{c}(1-\tilde{\sigma})K]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}$ .

Finally, we establish the convergence rate of OADMM with exploiting the KL exponent  $\tilde{\sigma}$ .

**Theorem 5.9.** (Proof in Section D.4, Convergence Rate) We fix  $p = 1/3$ . There exists  $t'$  such that for all  $t \geq t'$ , we have:

- (a) If  $\tilde{\sigma} \in (\frac{1}{4}, \frac{1}{2}]$ , then we have  $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F \leq \mathcal{O}(1/\exp(t^{1-u}))$ , where  $u = \frac{p(1-\tilde{\sigma})}{\tilde{\sigma}} \in [\frac{1}{3}, 1)$ .
- (b) If  $\tilde{\sigma} \in (\frac{1}{2}, 1)$ , then we have:  $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F \leq \mathcal{O}(1/(t^{(1-p)/\tau}))$ , where  $\tau = \frac{\tilde{\sigma}}{1-\tilde{\sigma}} - 1 \in (0, \infty)$ .

**Remark 5.10.** (i) To the best of our knowledge, Theorem 5.9 represents the first non-ergodic convergence rate for solving this class of nonconvex and nonsmooth problem in Problem (1). It is worth noting that the work of (Li et al., 2023) establishes a non-ergodic convergence rate for subgradient methods with diminishing stepsizes by further exploring the KL exponent. (ii) Under the KL inequality assumption, with the desingularizing function chosen in the form of  $\varphi(s) \triangleq \tilde{c}s^{1-\tilde{\sigma}}$  with  $\tilde{\sigma} \in (0, 1)$ , OADMM converges with a super-exponential rate when  $\tilde{\sigma} \in (\frac{1}{4}, \frac{1}{2}]$ , and converges with a polynomial convergence rate when  $\tilde{\sigma} \in (\frac{1}{2}, 1)$  for the gap  $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F$ . Notably, super-exponential convergence is faster than polynomial convergence. (iii) Our result generalizes the classical findings of (Attouch et al., 2010; Bolte et al., 2014), which characterize the convergence rate of proximal gradient methods for a specific class of nonconvex composite optimization problems.

## 6 APPLICATIONS AND NUMERICAL EXPERIMENTS

In this section, we assess the effectiveness of the proposed algorithm OADMM on the sparse PCA problem by comparing it against existing non-convex, non-smooth optimization algorithms.

► **Application to Sparse PCA.** Sparse PCA is a method to produce modified principal components with sparse loadings, which helps reduce model complexity and increase model interpretation (Chen et al., 2016). It can be formulated as:

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times r}} \frac{1}{2\tilde{m}} \|\mathbf{X}\mathbf{X}^T \mathbf{D} - \mathbf{D}\|_F^2 + \hat{\rho}(\|\mathbf{X}\|_1 - \|\mathbf{X}\|_{[k]}), \text{ s.t. } \mathbf{X}^T \mathbf{X} = \mathbf{I}_r,$$

where  $\mathbf{D} \in \mathbb{R}^{n \times \tilde{m}}$  is the data matrix,  $\tilde{m}$  is the number of data points, and  $\|\mathbf{X}\|_{[k]}$  is the  $\ell_1$  norm of the  $k$  largest (in magnitude) elements of the matrix  $\mathbf{X}$ . Here, we consider the DC  $\ell_1$ -largest- $k$  function (Gotoh et al., 2018) to induce sparsity in the solution. One advantage of this model is that when  $\hat{\rho}$  is sufficient large, we have  $\|\mathbf{X}\|_1 \approx \|\mathbf{X}\|_{[k]}$ , leading to a  $k$ -sparsity solution  $\mathbf{X}$ .

► **Compared Methods.** We compare OADMM-EP and OADMM-RR against four state-of-the-art optimization algorithms: (i) RADMM: ADMM using Riemannian retraction with fixed and small stepsizes (Li et al., 2022), tested with two different penalty parameters  $\forall t, \beta^t \in \{100, 10000\}$ , leading to two variants: RADMM-I and RADMM-II. (ii) SPGM-EP: Smoothing Proximal Gradient Method using Euclidean projection (Böhm & Wright, 2021). (iii) SPGM-EP: SPGM utilizing Riemannian retraction (Beck & Rosset, 2023). (iv) Sub-Grad: Subgradient methods with Euclidean projection (Davis & Drusvyatskiy, 2019; Li et al., 2021).

► **Experiment Settings.** All methods are implemented in MATLAB on an Intel 2.6 GHz CPU with 64 GB RAM. For all retraction-based methods, we use only polar decomposition-based retraction. We evaluate different regularization parameters  $\hat{\rho} \in \{10, 50, 100, 500, 1000\}$ . For OADMM, default parameters are used, with  $\beta^0 = 10\hat{\rho}$  and corresponding values  $\xi = \{1, 2, 5, 8, 10\}$  for each  $\hat{\rho}$ . For simplicity, we omit the Barzilai-Borwein strategy and instead use a fixed constant  $b^t = 1$  for all iterations. All algorithms start with a common initial solution  $\mathbf{x}^0$ , generated from a standard normal distribution. Our code for reproducing the experiments is available in the **supplemental material**.

► **Experiment Results.** We report the objective values for different methods with varying parameters  $\hat{\rho}$ . The experimental results presented in Figures 1 and 2 reveal the following insights: (i) Sub-Grad essentially fails to solve this problem, as the subgradient is inaccurately estimated when the solution is sparse. (ii) SPGM-EP and SPGM-RR, which rely on a variable smoothing strategy, exhibit slower performance than the multiplier-based variable splitting method. This observation aligns with the commonly accepted notion that primal-dual methods are generally more robust and faster than primal-only methods. (iii) The proposed OADMM-EP and OADMM-RR demonstrate similar results and generally achieve lower objective function values than the other methods.

## 7 CONCLUSIONS

This paper introduces OADMM, an Alternating Direction Method of Multipliers (ADMM) tailored for solving structured nonsmooth composite optimization problems under orthogonality constraints. OADMM integrates either a Nesterov extrapolation strategy or a Monotone Barzilai-Borwein (MBB) stepsize strategy to potentially accelerate primal convergence, complemented by an over-relaxation stepsize strategy for rapid dual convergence. We adjust the penalty and smoothing parameters at

486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

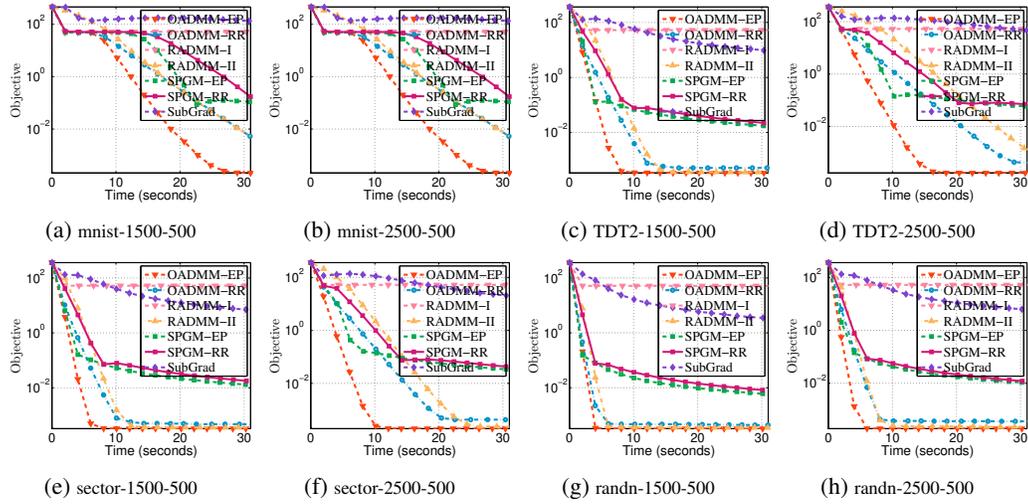


Figure 1: The convergence curve of the compared methods with  $\rho = 50$ .

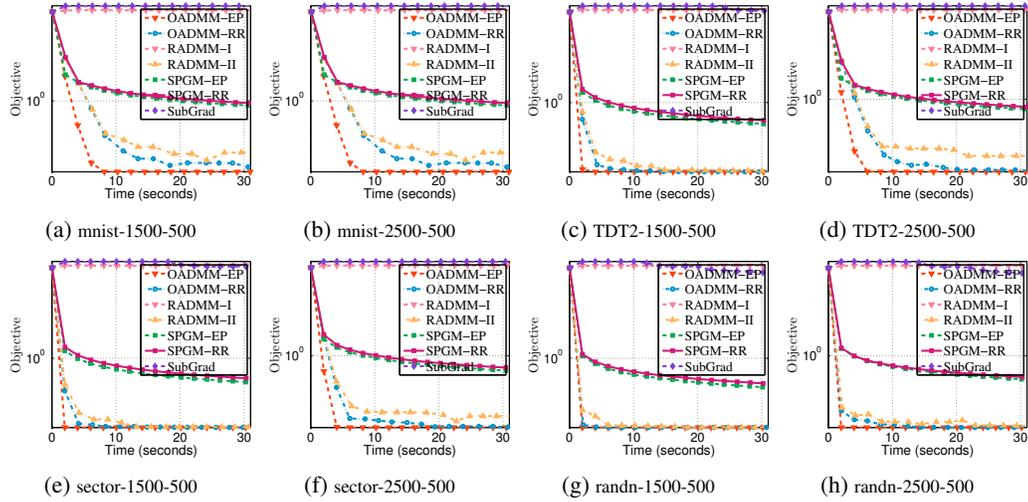


Figure 2: The convergence curve of the compared methods with  $\rho = 500$ .

a controlled rate. Additionally, we develop a novel Lyapunov function to rigorously analyze the oracle complexity of OADMM and establish the first non-ergodic convergence rate for this method. Finally, numerical experiments show that our OADMM achieves state-of-the-art performance.

## REFERENCES

- 540  
541  
542 Traian E Abrudan, Jan Eriksson, and Visa Koivunen. Steepest descent algorithms for optimiza-  
543 tion under unitary matrix constraint. *IEEE Transactions on Signal Processing*, 56(3):1134–1147,  
544 2008.
- 545 P-A Absil and Jérôme Malick. Projection-like retractions on matrix manifolds. *SIAM Journal on*  
546 *Optimization*, 22(1):135–158, 2012.
- 547 P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*.  
548 Princeton University Press, 2008a.
- 550 Pierre-Antoine Absil, Robert E. Mahony, and Rodolphe Sepulchre. *Optimization Algorithms on*  
551 *Matrix Manifolds*. Princeton University Press, 2008b.
- 553 Hedy Attouch, Jérôme Bolte, Patrick Redont, and Antoine Soubeyran. Proximal alternating mini-  
554 mization and projection methods for nonconvex problems: An approach based on the kurdyka-  
555 lojasiewicz inequality. *Mathematics of Operations Research*, 35(2):438–457, 2010.
- 556 Nitin Bansal, Xiaohan Chen, and Zhangyang Wang. Can we gain more from orthogonality regu-  
557 larizations in training deep networks? *Advances in Neural Information Processing Systems*, 31,  
558 2018.
- 560 Amir Beck and Israel Rosset. A dynamic smoothing technique for a class of nonsmooth optimization  
561 problems on manifolds. *SIAM Journal on Optimization*, 33(3):1473–1493, 2023.
- 562 Radu Ioan Boț and Dang-Khoa Nguyen. The proximal alternating direction method of multipliers  
563 in the nonconvex setting: convergence analysis and rates. *Mathematics of Operations Research*,  
564 45(2):682–712, 2020.
- 566 Radu Ioan Boț, Erno Robert Csetnek, and Dang-Khoa Nguyen. A proximal minimization algorithm  
567 for structured nonconvex and nonsmooth problems. *SIAM Journal on Optimization*, 29(2):1300–  
568 1328, 2019. doi: 10.1137/18M1190689.
- 569 Axel Böhm and Stephen J. Wright. Variable smoothing for weakly convex composite functions.  
570 *Journal of Optimization Theory and Applications*, 188(3):628–649, 2021.
- 572 Jérôme Bolte, Shoham Sabach, and Marc Teboulle. Proximal alternating linearized minimization  
573 for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1-2):459–494, 2014.
- 574 Nicolas Boumal, Pierre-Antoine Absil, and Coralia Cartis. Global rates of convergence for noncon-  
575 vex optimization on manifolds. *IMA Journal of Numerical Analysis*, 39(1):1–33, 2019.
- 577 Shixiang Chen, Shiqian Ma, Anthony Man-Cho So, and Tong Zhang. Proximal gradient method  
578 for nonsmooth optimization over the stiefel manifold. *SIAM Journal on Optimization*, 30(1):  
579 210–239, 2020.
- 580 Weiqiang Chen, Hui Ji, and Yanfei You. An augmented lagrangian method for  $\ell_1$ -regularized opti-  
581 mization problems with orthogonality constraints. *SIAM Journal on Scientific Computing*, 38(4):  
582 B570–B592, 2016.
- 584 Xiaojun Chen. Smoothing methods for nonsmooth, nonconvex minimization. *Mathematical pro-*  
585 *gramming*, 134(1):71–99, 2012.
- 586 Minhung Cho and Jaehyung Lee. Riemannian approach to batch normalization. *Advances in*  
587 *Neural Information Processing Systems*, 30, 2017.
- 589 Michael Cogswell, Faruk Ahmed, Ross B. Girshick, Larry Zitnick, and Dhruv Batra. Reducing  
590 overfitting in deep networks by decorrelating representations. In Yoshua Bengio and Yann LeCun  
591 (eds.), *International Conference on Learning Representations (ICLR)*, 2016.
- 592 Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex  
593 functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019.

- 594 Kangkang Deng, Jiang Hu, and Zaiwen Wen. Oracle complexity of augmented lagrangian methods  
595 for nonsmooth manifold optimization. *arXiv preprint arXiv:2404.05121*, 2024.  
596
- 597 Dmitriy Drusvyatskiy and Courtney Paquette. Efficiency of minimizing compositions of convex  
598 functions and smooth maps. *Mathematical Programming*, 178:503–558, 2019.
- 599 Alan Edelman, Tomás A. Arias, and Steven Thomas Smith. The geometry of algorithms with orthog-  
600 onality constraints. *SIAM Journal on Matrix Analysis and Applications*, 20(2):303–353, 1998.  
601
- 602 OP Ferreira and PR1622188 Oliveira. Subgradient algorithm on riemannian manifolds. *Journal of*  
603 *Optimization Theory and Applications*, 97:93–104, 1998.
- 604 Bin Gao, Xin Liu, Xiaojun Chen, and Ya-xiang Yuan. A new first-order algorithmic framework  
605 for optimization problems with orthogonality constraints. *SIAM Journal on Optimization*, 28(1):  
606 302–332, 2018.  
607
- 608 Bin Gao, Xin Liu, and Ya-xiang Yuan. Parallelizable algorithms for optimization problems with  
609 orthogonality constraints. *SIAM Journal on Scientific Computing*, 41(3):A1949–A1983, 2019.
- 610 Gene H Golub and Charles F Van Loan. *Matrix computations*. 2013.  
611
- 612 Max LN Gonçalves, Jefferson G Melo, and Renato DC Monteiro. Convergence rate bounds for  
613 a proximal admm with over-relaxation stepsize parameter for solving nonconvex linearly con-  
614 strained problems. *arXiv preprint arXiv:1702.01850*, 2017.
- 615 Jun-ya Gotoh, Akiko Takeda, and Katsuya Tono. Dc formulations and algorithms for sparse opti-  
616 mization problems. *Mathematical Programming*, 169(1):141–176, 2018.  
617
- 618 Bingsheng He and Xiaoming Yuan. On the  $o(1/n)$  convergence rate of the douglas–rachford alter-  
619 nating direction method. *SIAM Journal on Numerical Analysis*, 50(2):700–709, 2012.
- 620 Feihu Huang and Shangqian Gao. Gradient descent ascent for minimax problems on riemannian  
621 manifolds. *IEEE Trans. Pattern Anal. Mach. Intell.*, 45(7):8466–8476, 2023.  
622
- 623 Seong Jae Hwang, Maxwell D. Collins, Sathya N. Ravi, Vamsi K. Ithapu, Nagesh Adluru, Sterling C.  
624 Johnson, and Vikas Singh. A projection free method for generalized eigenvalue problem with a  
625 nonsmooth regularizer. In *IEEE International Conference on Computer Vision (ICCV)*, pp. 1841–  
626 1849, 2015.
- 627 Bo Jiang and Yu-Hong Dai. A framework of constraint preserving update schemes for optimization  
628 on stiefel manifold. *Mathematical Programming*, 153(2):535–575, 2015.  
629
- 630 Bo Jiang, Xiang Meng, Zaiwen Wen, and Xiaojun Chen. An exact penalty approach for optimization  
631 with nonnegative orthogonality constraints. *Mathematical Programming*, pp. 1–43, 2022.
- 632 Michel Journée, Yurii Nesterov, Peter Richtárik, and Rodolphe Sepulchre. Generalized power  
633 method for sparse principal component analysis. *Journal of Machine Learning Research*, 11  
634 (2):517–553, 2010.  
635
- 636 Artiom Kovnatsky, Klaus Glashoff, and Michael M Bronstein. Madmm: a generic algorithm for non-  
637 smooth optimization on manifolds. In *The European Conference on Computer Vision (ECCV)*,  
638 pp. 680–696. Springer, 2016.
- 639 Rongjie Lai and Stanley Osher. A splitting method for orthogonality constrained problems. *Journal*  
640 *of Scientific Computing*, 58(2):431–449, 2014.  
641
- 642 Guoyin Li and Ting Kei Pong. Global convergence of splitting methods for nonconvex composite  
643 optimization. *SIAM Journal on Optimization*, 25(4):2434–2460, 2015.
- 644 Huan Li and Zhouchen Lin. Accelerated proximal gradient methods for nonconvex programming.  
645 *Advances in neural information processing systems*, 28, 2015.  
646
- 647 Jiaxiang Li, Shiqian Ma, and Tejes Srivastava. A riemannian admm. *arXiv preprint*  
*arXiv:2211.02163*, 2022.

- 648 Min Li, Defeng Sun, and Kim-Chuan Toh. A majorized admm with indefinite proximal terms for  
649 linearly constrained convex composite optimization. *SIAM Journal on Optimization*, 26(2):922–  
650 950, 2016.
- 651 Xiao Li, Shixiang Chen, Zengde Deng, Qing Qu, Zhihui Zhu, and Anthony Man-Cho So. Weakly  
652 convex optimization over stiefel manifold using riemannian subgradient-type methods. *SIAM*  
653 *Journal on Optimization*, 31(3):1605–1634, 2021.
- 654 Xiao Li, Andre Milzarek, and Junwen Qiu. Convergence of random reshuffling under the kurdyka–  
655 lojasiewicz inequality. *SIAM Journal on Optimization*, 33(2):1092–1120, 2023.
- 656 Huikang Liu, Weijie Wu, and Anthony Man-Cho So. Quadratic optimization with orthogonality  
657 constraints: Explicit lojasiewicz exponent and linear convergence of line-search methods. In  
658 *International Conference on Machine Learning (ICML)*, pp. 1158–1167, 2016.
- 659 Zhaosong Lu and Yong Zhang. An augmented lagrangian approach for sparse principal component  
660 analysis. *Mathematical Programming*, 135(1-2):149–193, 2012.
- 661 Boris S. Mordukhovich. Variational analysis and generalized differentiation i: Basic theory. *Berlin*  
662 *Springer*, 330, 2006.
- 663 Y. E. Nesterov. *Introductory lectures on convex optimization: a basic course*, volume 87 of *Applied*  
664 *Optimization*. Kluwer Academic Publishers, 2003.
- 665 R. Tyrrell Rockafellar and Roger J-B. Wets. Variational analysis. *Springer Science & Business*  
666 *Media*, 317, 2009.
- 667 S Easter Selvan, S Thomas George, and R Balakrishnan. Range-based ica using a nonsmooth quasi-  
668 newton optimizer for electroencephalographic source localization in focal epilepsy. *Neural com-*  
669 *putation*, 27(3):628–671, 2015.
- 670 Yu Wang, Wotao Yin, and Jinshan Zeng. Global convergence of admm in nonconvex nonsmooth  
671 optimization. *Journal of Scientific Computing*, 78(1):29–63, 2019.
- 672 Zaiwen Wen and Wotao Yin. A feasible method for optimization with orthogonality constraints.  
673 *Mathematical Programming*, 142(1):397–434, 2013.
- 674 Nachuan Xiao, Xin Liu, and Ya-Xiang Yuan. A class of smooth exact penalty function methods for  
675 optimization problems with orthogonality constraints. *Optimization Methods and Software*, 37  
676 (4):1205–1241, 2022.
- 677 Di Xie, Jiang Xiong, and Shiliang Pu. All you need is beyond a good init: Exploring better solution  
678 for training extremely deep convolutional neural networks with orthonormality and modulation.  
679 In *IEEE Conference on Computer Vision and Pattern Recognition*, pp. 6176–6185, 2017.
- 680 Lei Yang, Ting Kei Pong, and Xiaojun Chen. Alternating direction method of multipliers for a class  
681 of nonconvex and nonsmooth problems with applications to background/foreground extraction.  
682 *SIAM Journal on Imaging Sciences*, 10(1):74–110, 2017.
- 683 Ganzhao Yuan. A block coordinate descent method for nonsmooth composite optimization under  
684 orthogonality constraints. *arXiv preprint arXiv:2304.03641*, 2023.
- 685 Ganzhao Yuan. Admm for nonconvex optimization under minimal continuity assumption. *arXiv*  
686 *preprint*, 2024.
- 687 Jinshan Zeng, Wotao Yin, and Ding-Xuan Zhou. Moreau envelope augmented lagrangian method  
688 for nonconvex optimization with linear constraints. *Journal of Scientific Computing*, 91(2):61,  
689 2022.
- 690 Yuexiang Zhai, Zitong Yang, Zhenyu Liao, John Wright, and Yi Ma. Complete dictionary learning  
691 via l4-norm maximization over the orthogonal group. *Journal of Machine Learning Research*, 21  
692 (165):1–68, 2020.

702 Jiawei Zhang, Peijun Xiao, Ruoyu Sun, and Zhi-Quan Luo. A single-loop smoothed gradi-  
703 ent descent-ascent algorithm for nonconvex-concave min-max problems. In Hugo Larochelle,  
704 Marc’ Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances*  
705 *in Neural Information Processing Systems*, 2020a.

706  
707 Junyu Zhang, Shiqian Ma, and Shuzhong Zhang. Primal-dual optimization algorithms over rieman-  
708 nian manifolds: an iteration complexity analysis. *Mathematical Programming*, 184(1):445–490,  
709 2020b.

710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

# Appendix

The appendix is organized as follows.

Appendix A provides notations, technical preliminaries, and relevant lemmas.

Appendix B contains the proofs for Section 2.

Appendix C includes the proofs for Section 4.

Appendix D encompasses the proofs for Section 5.

Appendix E presents additional experiments details and results.

## A NOTATIONS, TECHNICAL PRELIMINARIES, AND RELEVANT LEMMAS

### A.1 NOTATIONS

In this paper, lowercase boldface letters signify vectors, while uppercase letters denote real-valued matrices. The following notations are utilized throughout this paper.

- $[n]$ :  $\{1, 2, \dots, n\}$
- $\|\mathbf{x}\|$ : Euclidean norm:  $\|\mathbf{x}\| = \|\mathbf{x}\|_2 = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$
- $\mathbf{X}^\top$ : the transpose of the matrix  $\mathbf{X}$
- $\mathbf{0}_{n,r}$ : A zero matrix of size  $n \times r$ ; the subscript is omitted sometimes
- $\mathbf{I}_r$ :  $\mathbf{I}_r \in \mathbb{R}^{r \times r}$ , Identity matrix
- $\mathcal{M}$ : Orthogonality constraint set (a.k.a., Stiefel manifold:  $\mathcal{M} = \{\mathbf{X} \in \mathbb{R}^{n \times r} \mid \mathbf{X}^\top \mathbf{X} = \mathbf{I}_r\}$ ).
- $\mathbf{X} \succeq \mathbf{0}$  (or  $\succ \mathbf{0}$ ): the Matrix  $\mathbf{X}$  is symmetric positive semidefinite (or definite)
- $\text{tr}(\mathbf{A})$ : Sum of the elements on the main diagonal  $\mathbf{A}$ :  $\text{tr}(\mathbf{A}) = \sum_i \mathbf{A}_{i,i}$
- $\|\mathbf{X}\|$ : Operator/Spectral norm: the largest singular value of  $\mathbf{X}$
- $\|\mathbf{X}\|_F$ : Frobenius norm:  $(\sum_{ij} \mathbf{X}_{ij}^2)^{1/2}$
- $\|\mathbf{X}\|_1$ : Absolute sum of the elements in  $\mathbf{X}$  with  $\mathbf{X} = \sum_{ij} |\mathbf{X}_{ij}|$
- $\|\mathbf{X}\|_{[k]}$ :  $\ell_1$  norm the the  $k$  largest (in magnitude) elements of the matrix  $\mathbf{X}$
- $\partial g(\mathbf{X})$ : (limiting) Euclidean subdifferential of  $g(\mathbf{X})$  at  $\mathbf{X}$
- $\text{Proj}_{\Xi}(\mathbf{X}')$ : Orthogonal projection of  $\mathbf{X}'$  with  $\text{Proj}_{\Xi}(\mathbf{X}') = \arg \arg \min_{\mathbf{X} \in \Xi} \|\mathbf{X}' - \mathbf{X}\|_F^2$
- $\text{dist}(\Xi, \Xi')$ : the distance between two sets with  $\text{dist}(\Xi, \Xi') \triangleq \inf_{\mathbf{X} \in \Xi, \mathbf{X}' \in \Xi'} \|\mathbf{X} - \mathbf{X}'\|_F$
- $\|\partial g(\mathbf{X})\|_F$ :  $\|\partial g(\mathbf{X})\|_F = \inf_{\mathbf{Y} \in \partial g(\mathbf{X})} \|\mathbf{Y}\|_F = \text{dist}(\mathbf{0}, \partial g(\mathbf{X}))$ .
- $\ell(\beta^t)$ : the smoothness parameter of the function  $\mathcal{S}(\mathbf{X}, \mathbf{y}^t; \mathbf{z}^t; \beta^t)$  w.r.t.  $\mathbf{X}$ .
- $\mathcal{I}_{\mathcal{M}}(\mathbf{x})$ : Indicator function of  $\mathcal{M}$  with  $\mathcal{I}_{\mathcal{M}}(\mathbf{x}) = 0$  if  $\mathbf{x} \in \mathcal{M}$  and otherwise  $+\infty$ .

We employ the following parameters in Algorithm 1.

- $\theta$ : proximal parameter
- $\chi$ : correlation coefficient between  $\mu^t$  and  $\beta^t$ , such that  $\mu^t \beta^t = \chi$
- $\sigma$ : over-relaxation parameter with  $\sigma \in [1, 2)$
- $\alpha$ : Nesterov extrapolation parameter with  $\alpha \in [0, 1)$
- $\rho$ : search descent parameter with  $\rho \in (0, \infty)$
- $\gamma$ : decay rate parameter in the line search procedure with  $\gamma \in (0, 1)$
- $\delta$ : sufficient decrease parameter in the line search procedure with  $\delta \in (0, \infty)$
- $p$ : exponent parameter used in the penalty update rule with  $p \in (0, 1)$
- $\xi$ : growth factor parameter used in the penalty update rule with  $\xi \in (0, \infty)$

## A.2 TECHNICAL PRELIMINARIES

**Non-convex Non-smooth Optimization.** Given the potential non-convexity and non-smoothness of the function  $F(\cdot)$ , we introduce tools from non-smooth analysis (Mordukhovich, 2006; Rockafellar & Wets., 2009). The domain of any extended real-valued function  $F : \mathbb{R}^{n \times r} \rightarrow (-\infty, +\infty]$  is defined as  $\text{dom}(F) \triangleq \{\mathbf{X} \in \mathbb{R}^{n \times r} : |F(\mathbf{X})| < +\infty\}$ . At  $\mathbf{X} \in \text{dom}(F)$ , the Fréchet subdifferential of  $F$  is defined as  $\hat{\partial}F(\mathbf{X}) \triangleq \{\boldsymbol{\xi} \in \mathbb{R}^{n \times r} : \lim_{\mathbf{Z} \rightarrow \mathbf{X}} \inf_{\mathbf{Z} \neq \mathbf{X}} \frac{F(\mathbf{Z}) - F(\mathbf{X}) - \langle \boldsymbol{\xi}, \mathbf{Z} - \mathbf{X} \rangle}{\|\mathbf{Z} - \mathbf{X}\|_F} \geq 0\}$ , while the limiting subdifferential of  $F(\mathbf{X})$  at  $\mathbf{X} \in \text{dom}(F)$  is denoted as  $\partial F(\mathbf{X}) \triangleq \{\boldsymbol{\xi} \in \mathbb{R}^{n \times r} : \exists \mathbf{X}^t \rightarrow \mathbf{X}, F(\mathbf{X}^t) \rightarrow F(\mathbf{X}), \boldsymbol{\xi}^t \in \hat{\partial}F(\mathbf{X}^t) \rightarrow \boldsymbol{\xi}, \forall t\}$ . The gradient of  $F(\cdot)$  at  $\mathbf{X}$  in the Euclidean space is denoted as  $\nabla F(\mathbf{X})$ . The following relations hold among  $\hat{\partial}F(\mathbf{X})$ ,  $\partial F(\mathbf{X})$ , and  $\nabla F(\mathbf{X})$ : (i)  $\hat{\partial}F(\mathbf{X}) \subseteq \partial F(\mathbf{X})$ . (ii) If the function  $F(\cdot)$  is convex,  $\partial F(\mathbf{X})$  and  $\hat{\partial}F(\mathbf{X})$  represent the classical subdifferential for convex functions, i.e.,  $\partial F(\mathbf{X}) = \hat{\partial}F(\mathbf{X}) = \{\boldsymbol{\xi} \in \mathbb{R}^{n \times r} : F(\mathbf{Z}) \geq F(\mathbf{X}) + \langle \boldsymbol{\xi}, \mathbf{Z} - \mathbf{X} \rangle, \forall \mathbf{Z} \in \mathbb{R}^{n \times r}\}$ . (iii) If the function  $F(\cdot)$  is differentiable, then  $\hat{\partial}F(\mathbf{X}) = \partial F(\mathbf{X}) = \{\nabla F(\mathbf{X})\}$ .

**Optimization with Orthogonality Constraints.** We introduce some prior knowledge of optimization involving orthogonality constraints (Absil et al., 2008b). The nearest orthogonality matrix to any arbitrary matrix  $\mathbf{Y} \in \mathbb{R}^{n \times r}$  is determined as  $\mathbb{P}_{\mathcal{M}}(\mathbf{Y}) = \check{\mathbf{U}}\check{\mathbf{V}}^T$ , where  $\mathbf{Y} = \check{\mathbf{U}}\text{Diag}(\mathbf{s})\check{\mathbf{V}}^T$  represents the singular value decomposition of  $\mathbf{Y}$ . We use  $\mathcal{N}_{\mathcal{M}}(\mathbf{X})$  to denote the limiting normal cone to  $\mathcal{M}$  at  $\mathbf{X}$ , thus defined as  $\mathcal{N}_{\mathcal{M}}(\mathbf{X}) = \partial \mathcal{I}_{\mathcal{M}}(\mathbf{X}) = \{\mathbf{Z} \in \mathbb{R}^{n \times r} : \langle \mathbf{Z}, \mathbf{X} \rangle \geq \langle \mathbf{Z}, \mathbf{Y} \rangle, \forall \mathbf{Y} \in \mathcal{M}\}$ . Moreover, the tangent and normal space to  $\mathcal{M}$  at  $\mathbf{X} \in \mathcal{M}$  are respectively denoted as  $T_{\mathbf{X}}\mathcal{M}$  and  $N_{\mathbf{X}}\mathcal{M}$ . We have:  $T_{\mathbf{X}}\mathcal{M} = \{\mathbf{Y} \in \mathbb{R}^{n \times r} | \mathcal{A}_{\mathbf{X}}(\mathbf{Y}) = \mathbf{0}\}$  and  $N_{\mathbf{X}}\mathcal{M} = 2\mathbf{X}\boldsymbol{\Lambda} | \boldsymbol{\Lambda} = \boldsymbol{\Lambda}^T, \boldsymbol{\Lambda} \in \mathbb{R}^{r \times r}$ , where  $\mathcal{A}_{\mathbf{X}}(\mathbf{Y}) \triangleq \mathbf{X}^T\mathbf{Y} + \mathbf{Y}^T\mathbf{X}$  for  $\mathbf{Y} \in \mathbb{R}^{n \times r}$  and  $\mathbf{X} \in \mathcal{M}$ .

**Weakly Convex Functions.** The function  $h(\mathbf{y})$  is weakly convex if there exists a constant  $W_h \geq 0$  such that  $h(\mathbf{y}) + \frac{1}{2}W_h\|\mathbf{y}\|_2^2$  is convex; the smallest such  $W_h$  is termed the modulus of weak convexity. Weakly convex functions encompass a diverse range, including convex functions, differentiable functions with Lipschitz continuous gradient, and compositions of convex, Lipschitz-continuous functions with  $C^1$ -smooth mappings having Lipschitz continuous Jacobians (Drusvyatskiy & Paquette, 2019).

## A.3 RELEVANT LEMMAS

**Lemma A.1.** Let  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$ , and  $\alpha \geq 0$  be any constant. We have:  $-\|\mathbf{a} - \alpha\mathbf{b}\|_2^2 \leq (\alpha - 1)\|\mathbf{a}\|_2^2 - (\alpha^2 - \alpha)\|\mathbf{b}\|_2^2$ .

*Proof.* We have:  $-\|\mathbf{a} - \alpha\mathbf{b}\|_2^2 = -\|\mathbf{a}\|_2^2 - \|\alpha\mathbf{b}\|_2^2 + 2\alpha\langle \mathbf{a}, \mathbf{b} \rangle \leq -\|\mathbf{a}\|_2^2 - \|\alpha\mathbf{b}\|_2^2 + 2\alpha \cdot (\frac{1}{2}\|\mathbf{a}\|_2^2 + \frac{1}{2}\|\mathbf{b}\|_2^2) = (\alpha - 1)\|\mathbf{a}\|_2^2 - (\alpha^2 - \alpha)\|\mathbf{b}\|_2^2$ .  $\square$

**Lemma A.2.** Assume  $t \geq 1$  and  $p \in (0, 1)$ . We have:  $\frac{t^p - (t-1)^p}{1 + (t-1)^p} \leq \frac{1}{t}$ .

*Proof.* We let  $t \geq 1$  and  $p \in (0, 1)$ .

First, we define  $f(t) \triangleq t^p - 2(t-1)^p - 1$ . We have  $\nabla f(t) = pt^{p-1} - 2p(t-1)^{p-1} = p(t-1)^{p-1}\{(\frac{t}{t-1})^{p-1} - 2\} \leq p(t-1)^{p-1}\{(\frac{t}{t-1}) - 2\} \leq p(t-1)^{p-1}\{\frac{1}{2} - 2\} \leq 0$ . This implies that  $f(t)$  is decreasing. Noting that  $f(1) = 0$ , we conclude that

$$f(t) \triangleq t^p - 2(t-1)^p - 1 \leq 0. \quad (8)$$

Second, we have:

$$g(t) \triangleq t^p - \frac{t^{p+1}}{t+1} - (t-1)^p \stackrel{\textcircled{1}}{\leq} t^p - \frac{t^p}{t} - (t-1)^p = t^p((1 - \frac{1}{t}) - (1 - \frac{1}{t})^p) \stackrel{\textcircled{2}}{\leq} 0, \quad (9)$$

where step  $\textcircled{1}$  uses  $\frac{t^{p+1}}{t+1} \geq \frac{t^p}{t}$  as  $t \geq t^p$ ; step  $\textcircled{2}$  uses  $a \leq a^p$  for all  $a \in (0, 1)$  and  $p \in (0, 1)$ .

Finally, we derive the following results:

$$\frac{t^p - (t-1)^p}{1 + (t-1)^p} \cdot t \stackrel{\textcircled{1}}{\leq} \frac{t}{t+1} \cdot \frac{1+t^p}{1+(t-1)^p} \stackrel{\textcircled{2}}{\leq} \frac{1+t^p}{2+2(t-1)^p} \stackrel{\textcircled{3}}{\leq} 1,$$

where step ① uses Inequality (9); step ② uses  $\frac{t}{t+1} \leq \frac{1}{2}$ ; step ③ uses Inequality (8). □

**Lemma A.3.** Let  $\beta^t = \beta^0(1 + \xi t^p)$ , where  $t \geq 0$ ,  $\beta^0 > 0$ ,  $\xi, p \in (0, 1)$ . For all  $t \geq 1$ , we have:  $(\frac{\beta^t}{\beta^{t-1}} - 1)^2 \leq \frac{2}{t} - \frac{2}{t+1}$ .

*Proof.* We derive:  $(\frac{\beta^t}{\beta^{t-1}} - 1)^2 \stackrel{\text{①}}{=} (\frac{1+\xi t^p}{1+\xi(t-1)^p} - 1)^2 = (\frac{\xi t^p - \xi(t-1)^p}{1+\xi(t-1)^p})^2 \stackrel{\text{②}}{\leq} (\frac{t^p - (t-1)^p}{1+(t-1)^p})^2 \stackrel{\text{③}}{\leq} (\frac{1}{t})^2 \stackrel{\text{④}}{\leq} \frac{2}{t} - \frac{2}{t+1}$ , where step ① uses  $\beta^t = \beta^0(1 + \xi t^p)$ ; step ② uses  $\frac{\xi}{1+\xi a} < \frac{1}{1+a}$  for all  $a \geq 0$  when  $\xi \in (0, 1)$ ; step ③ uses Lemma A.2; step ④ uses the fact that  $\frac{1}{t^2} \leq \frac{2}{t} - \frac{2}{t+1}$  for all  $t \geq 1$ . □

**Lemma A.4.** Assume  $\mathbf{a}^+ = \rho \mathbf{a} + \mathbf{b}$ , where  $\mathbf{a}, \mathbf{b}, \mathbf{a}^+ \in \mathbb{R}^m$ , and  $\rho \in [0, 1)$ . We have:  $\|\mathbf{a}^+\|_2^2 \leq \frac{\rho}{1-\rho}(\|\mathbf{a}\|_2^2 - \|\mathbf{a}^+\|_2^2) + \frac{1}{(1-\rho)^2}\|\mathbf{b}\|_2^2$ .

*Proof.* We have:  $\|\mathbf{a}^+\|_2^2 = \|\rho \mathbf{a} + \mathbf{b}\|_2^2 = \|\rho \mathbf{a} + (1-\rho) \cdot \frac{\mathbf{b}}{1-\rho}\|_2^2 \leq \rho \|\mathbf{a}\|_2^2 + (1-\rho) \cdot \|\frac{\mathbf{b}}{1-\rho}\|_2^2 = \rho \|\mathbf{a}\|_2^2 + \frac{1}{1-\rho} \|\mathbf{b}\|_2^2$ , where the inequality holds due to the convexity of  $\|\cdot\|_2^2$ . □

**Lemma A.5.** Assume that  $\mathbf{a}^t \leq \rho \mathbf{a}^{t-1} + c$ , where  $\rho \in [0, 1)$ ,  $c \geq 0$ , and  $\{\mathbf{a}^i\}_{i=0}^\infty$  is a non-negative sequence. We have:  $\mathbf{a}^t \leq \mathbf{a}^0 + \frac{c}{1-\rho}$  for all  $t \geq 0$ .

*Proof.* Using basic induction, we have the following results:

$$\begin{aligned} t = 1, \quad \mathbf{a}^1 &\leq \rho \mathbf{a}^0 + c \\ t = 2, \quad \mathbf{a}^2 &\leq \rho \mathbf{a}^1 + c \leq \rho(\rho \mathbf{a}^0 + c) + c = \rho^2 \mathbf{a}^0 + c(1 + \rho) \\ t = 3, \quad \mathbf{a}^3 &\leq \rho \mathbf{a}^2 + c \leq \rho(\rho^2 \mathbf{a}^0 + (c + \rho c)) + c = \rho^3 \mathbf{a}^0 + c(1 + \rho + \rho^2) \\ &\dots \\ t = n, \quad \mathbf{a}^n &\leq \rho \mathbf{a}^{n-1} + c \leq \rho^n \mathbf{a}^0 + c \cdot (1 + \rho + \dots + \rho^{n-1}). \end{aligned}$$

Therefore, we obtain:  $\mathbf{a}^n \leq \rho^n \mathbf{a}^0 + c \cdot (1 + \rho + \dots + \rho^{n-1}) \stackrel{\text{①}}{\leq} a_0 + \frac{c}{1-\rho}$ , where step ① uses  $\rho^n \leq \rho < 1$ , and the summation formula of geometric sequences that  $1 + \rho^1 + \rho^2 + \dots + \rho^{t-1} = \frac{1-\rho^t}{1-\rho} < \frac{1}{1-\rho}$ . □

**Lemma A.6.** Assume  $\mathbf{X}_c^t = \mathbf{X}^t + \alpha(\mathbf{X}^t - \mathbf{X}^{t-1})$ , where  $\alpha \in [0, 1)$ , and  $\mathbf{X}^t, \mathbf{X}^{t-1} \in \mathcal{M}$ . We have:

- (a)  $\|\mathbf{X}^t - \mathbf{X}_c^t\|_F \leq \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$ .
- (b)  $\|\mathbf{X}^{t+1} - \mathbf{X}_c^t\|_F \leq \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F + \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$ .
- (c)  $\|\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t\| \leq \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \bar{\mathbf{A}}\|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$ .

*Proof. Part (a).* We have:  $\|\mathbf{X}^t - \mathbf{X}_c^t\|_F \stackrel{\text{①}}{=} \alpha \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F \stackrel{\text{②}}{\leq} \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$ , where step ① uses  $\mathbf{X}_c^t = \mathbf{X}^t + \alpha(\mathbf{X}^t - \mathbf{X}^{t-1})$ ; step ② uses  $\alpha \in [0, 1)$ .

*Part (b).* We have:  $\|\mathbf{X}^{t+1} - \mathbf{X}_c^t\|_F \stackrel{\text{①}}{=} \|\mathbf{X}^{t+1} - \mathbf{X}^t - \alpha(\mathbf{X}^t - \mathbf{X}^{t-1})\|_F \stackrel{\text{②}}{\leq} \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F + \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$ , where step ① uses  $\mathbf{X}_c^t = \mathbf{X}^t + \alpha(\mathbf{X}^t - \mathbf{X}^{t-1})$ ; step ② uses the triangle inequality and  $\alpha \in [0, 1)$ .

*Part (c).* We have:  $\|\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t\| \stackrel{\text{①}}{\leq} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \|\mathcal{A}(\mathbf{X}^t) - \mathcal{A}(\mathbf{X}_c^t)\| \leq \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \bar{\mathbf{A}}\|\mathbf{X}^t - \mathbf{X}_c^t\| \stackrel{\text{②}}{\leq} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \bar{\mathbf{A}}\|\mathbf{X}^t - \mathbf{X}^{t-1}\|$ , where step ① uses the triangle inequality; step ② uses Claim (a) of this lemma. □

**Lemma A.7.** Let  $\mathbf{P}, \tilde{\mathbf{P}} \in \mathbb{R}^{n \times r}$ , and  $\mathbf{X}, \tilde{\mathbf{X}} \in \mathcal{M}$ . We have:

$$\|\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\mathbf{P}) - \text{Proj}_{\mathbf{T}_{\tilde{\mathbf{X}}}\mathcal{M}}(\tilde{\mathbf{P}})\|_{\text{F}} \leq 2\|\mathbf{P} - \tilde{\mathbf{P}}\|_{\text{F}} + 2\sqrt{r}\|\mathbf{P}\|\|\mathbf{X} - \tilde{\mathbf{X}}\|_{\text{F}}.$$

*Proof.* First, we obtain:

$$\begin{aligned} & \|\mathbf{X}\mathbf{P}^{\text{T}}\mathbf{X} - \tilde{\mathbf{X}}\tilde{\mathbf{P}}^{\text{T}}\tilde{\mathbf{X}}\|_{\text{F}} \\ &= \|(\mathbf{X} - \tilde{\mathbf{X}})\mathbf{P}^{\text{T}}\mathbf{X} + \tilde{\mathbf{X}}\mathbf{P}^{\text{T}}(\mathbf{X} - \tilde{\mathbf{X}}) + \tilde{\mathbf{X}}(\mathbf{P} - \tilde{\mathbf{P}})^{\text{T}}\tilde{\mathbf{X}}\|_{\text{F}} \\ &\stackrel{\textcircled{1}}{\leq} \|\mathbf{X} - \tilde{\mathbf{X}}\|_{\text{F}}\|\mathbf{P}^{\text{T}}\mathbf{X}\| + \|\tilde{\mathbf{X}}\mathbf{P}^{\text{T}}\|\|\mathbf{X} - \tilde{\mathbf{X}}\|_{\text{F}} + \|\tilde{\mathbf{X}}(\mathbf{P} - \tilde{\mathbf{P}})^{\text{T}}\tilde{\mathbf{X}}\|_{\text{F}} \\ &\stackrel{\textcircled{2}}{\leq} 2\sqrt{r}\|\mathbf{P}\|\|\mathbf{X} - \tilde{\mathbf{X}}\|_{\text{F}} + \|\mathbf{P} - \tilde{\mathbf{P}}\|_{\text{F}}, \end{aligned} \quad (10)$$

where step  $\textcircled{1}$  uses the triangle inequality; step  $\textcircled{2}$  uses  $\|\mathbf{AB}\|_{\text{F}} \leq \|\mathbf{A}\| \cdot \|\mathbf{B}\|_{\text{F}}$ , and  $\|\tilde{\mathbf{X}}\| \leq 1$ .

Second, we have:

$$\begin{aligned} & \|\mathbf{X}\mathbf{X}^{\text{T}}\mathbf{P} - \tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\text{T}}\tilde{\mathbf{P}}\|_{\text{F}} \\ &= \|(\mathbf{X} - \tilde{\mathbf{X}})\mathbf{X}^{\text{T}}\mathbf{P} + \tilde{\mathbf{X}}(\mathbf{X} - \tilde{\mathbf{X}})^{\text{T}}\mathbf{P} + \tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\text{T}}(\mathbf{P} - \tilde{\mathbf{P}})\|_{\text{F}} \\ &\stackrel{\textcircled{1}}{\leq} \|\mathbf{X} - \tilde{\mathbf{X}}\|_{\text{F}}\|\mathbf{X}^{\text{T}}\mathbf{P}\| + \|\tilde{\mathbf{X}}\| \cdot \|\mathbf{X} - \tilde{\mathbf{X}}\|_{\text{F}} \cdot \|\mathbf{P}\| + \|\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\text{T}}\| \cdot \|\mathbf{P} - \tilde{\mathbf{P}}\|_{\text{F}} \\ &\stackrel{\textcircled{2}}{\leq} 2\sqrt{r}\|\mathbf{P}\|\|\mathbf{X} - \tilde{\mathbf{X}}\|_{\text{F}} + \|\mathbf{P} - \tilde{\mathbf{P}}\|_{\text{F}}, \end{aligned} \quad (11)$$

where step  $\textcircled{1}$  uses the triangle inequality; step  $\textcircled{2}$  uses  $\|\mathbf{AB}\|_{\text{F}} \leq \|\mathbf{A}\| \cdot \|\mathbf{B}\|_{\text{F}}$ , and  $\|\tilde{\mathbf{X}}\| \leq 1$ .

Finally, we derive:

$$\begin{aligned} & \|\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\mathbf{P}) - \text{Proj}_{\mathbf{T}_{\tilde{\mathbf{X}}}\mathcal{M}}(\tilde{\mathbf{P}})\|_{\text{F}} \\ &\stackrel{\textcircled{1}}{=} \|\mathbf{P} - \frac{1}{2}\mathbf{X}\mathbf{P}^{\text{T}}\mathbf{X} - \frac{1}{2}\mathbf{X}\mathbf{X}^{\text{T}}\mathbf{P}\|_{\text{F}} - \|\tilde{\mathbf{P}} - \frac{1}{2}\tilde{\mathbf{X}}\tilde{\mathbf{P}}^{\text{T}}\tilde{\mathbf{X}} - \frac{1}{2}\tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\text{T}}\tilde{\mathbf{P}}\|_{\text{F}} \\ &\stackrel{\textcircled{2}}{\leq} \|\mathbf{P} - \tilde{\mathbf{P}}\|_{\text{F}} + \frac{1}{2}\|\mathbf{X}\mathbf{P}^{\text{T}}\mathbf{X} - \tilde{\mathbf{X}}\tilde{\mathbf{P}}^{\text{T}}\tilde{\mathbf{X}}\|_{\text{F}} + \frac{1}{2}\|\mathbf{X}\mathbf{X}^{\text{T}}\mathbf{P} - \tilde{\mathbf{X}}\tilde{\mathbf{X}}^{\text{T}}\tilde{\mathbf{P}}\|_{\text{F}} \\ &\stackrel{\textcircled{3}}{\leq} \|\mathbf{P} - \tilde{\mathbf{P}}\|_{\text{F}} + 2\sqrt{r}\|\mathbf{P}\|\|\mathbf{X} - \tilde{\mathbf{X}}\|_{\text{F}} + \|\mathbf{P} - \tilde{\mathbf{P}}\|_{\text{F}} \end{aligned}$$

where step  $\textcircled{1}$  uses  $\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\mathbf{\Delta}) = \mathbf{\Delta} - \frac{1}{2}\mathbf{X}(\mathbf{\Delta}^{\text{T}}\mathbf{X} + \mathbf{X}^{\text{T}}\mathbf{\Delta})$  for all  $\mathbf{\Delta} \in \mathbb{R}^{n \times r}$  (Absil et al., 2008a); step  $\textcircled{2}$  uses the triangle inequality; step  $\textcircled{3}$  uses Inequalities (10) and (11).  $\square$

**Lemma A.8.** We let  $p \in (0, 1)$ . We define  $g(t) \triangleq \frac{1}{1-p}(t+1)^{(1-p)} - \frac{1}{1-p} - (1-p)t^{(1-p)}$ . We have  $g(t) \geq 0$  for all  $t \geq 1$ .

*Proof.* We assume  $p \in (0, 1)$ .

First, we show that  $h(p) \triangleq (1-p)^{1/p} \leq \frac{1}{\exp(1)}$ . Recall that it holds:  $\lim_{p \rightarrow 0^+} (1+p)^{1/p} = \exp(1)$  and  $\lim_{p \rightarrow 0^+} (1-p)^{1/p} = 1/\exp(1)$ . Given the function  $h(p)$  is a decreasing function on  $p \in (0, 1)$ , we have  $h(p) \leq \lim_{p \rightarrow 0^+} (1-p)^{1/p} = \frac{1}{\exp(1)}$ .

Second, we show that  $f(q) = 2^q - 1 - q^2 \geq 0$  for all  $q \in (0, 1)$ . We have  $\nabla f(q) = \log(2)2^q - 2q$ , and  $\nabla^2 f(q) = 2^q(\log(2))^2 - 2 \leq 2(\log(2))^2 - 2 \leq 0$ , implying that the function  $f(q)$  is concave on  $q \in (0, 1)$ . Noticing  $f(0) = f(1) = 0$ , we conclude that  $f(q) \geq 0$ .

Third, we show that  $g(t)$  is an increasing function. We have:  $\nabla g(t) = (t+1)^{-p} - (1-p)^2 t^{-p} = (t+1)^{-p} \cdot (1 - (1-p)^2 (\frac{t+1}{t})^p) \stackrel{\textcircled{1}}{\geq} (t+1)^{-p} \cdot (1 - (1-p)^2 2^p) \stackrel{\textcircled{2}}{\geq} (t+1)^{-p} \cdot (1 - (\frac{2}{\exp(1)^2})^p) \stackrel{\textcircled{3}}{\geq} 0$ , where step  $\textcircled{1}$  uses  $\frac{t+1}{t} \leq 2$  for all  $t \geq 1$ ; step  $\textcircled{2}$  uses  $1-p \leq (\frac{1}{\exp(1)})^p$  for all  $p \in (0, 1)$ ; step  $\textcircled{3}$  uses  $\frac{2}{\exp(1)^2} \approx 0.2707 < 1$ .

Finally, we have:  $\forall t \geq 1, g(t) \stackrel{\textcircled{1}}{\geq} g(1) = (1-p)^{-1} \cdot \{2^{(1-p)} - 1 - (1-p)^2\} \stackrel{\textcircled{2}}{\geq} 0$ , where step  $\textcircled{1}$  uses the fact that  $g(t)$  is an increasing function; step  $\textcircled{2}$  uses  $2^q - 1 - q^2 \geq 0$  for all  $q = 1-p \in (0, 1)$ .  $\square$

**Lemma A.9.** Assume  $p \in (0, 1)$ . We have:  $(1-p)T^{(1-p)} \leq \sum_{t=1}^T \frac{1}{t^p} \leq \frac{T^{(1-p)}}{1-p}$ .

*Proof.* We define  $g(t) \triangleq \frac{1}{t^p}$  and  $h(t) \triangleq \frac{1}{1-p}t^{(1-p)}$ .

Using the integral test for convergence, we obtain:  $\int_1^{T+1} g(x)dx \leq \sum_{t=1}^T g(t) \leq g(1) + \int_1^T g(x)dx$ .

**Part (a).** We first consider the lower bound. We obtain:  $\sum_{t=1}^T t^{-p} \geq \sum_{t=1}^T \int_t^{t+1} x^{-p} dx = \int_1^{T+1} x^{-p} dx \stackrel{\textcircled{1}}{\geq} h(T+1) - h(1) = \frac{1}{1-p}(T+1)^{1-p} - \frac{1}{1-p} \stackrel{\textcircled{2}}{\geq} (1-p)T^{1-p}$ , where step  $\textcircled{1}$  uses  $\nabla h(x) = x^{-p}$ ; step  $\textcircled{2}$  uses Lemma A.8.

**Part (b).** We now consider the upper bound. We have:  $\sum_{t=1}^T t^{-p} \leq h(1) + \int_1^T x^{-p} dx \stackrel{\textcircled{1}}{=} 1 + h(T) - h(1) = 1 + \frac{1}{1-p}(T)^{1-p} - \frac{1}{1-p} = \frac{T^{(1-p)} - p}{1-p} < \frac{T^{(1-p)}}{1-p}$ , where step  $\textcircled{1}$  uses  $\nabla h(x) = x^{-p}$ .  $\square$

**Lemma A.10.** Assume  $(e^{t+1})^2 \leq (e^t + e^{t-1})(p^t - p^{t+1})$  and  $p^t \geq p^{t+1}$ , where  $\{e^t, p^t\}_{t=0}^\infty$  are two nonnegative sequences. For all  $i \geq 1$ , we have:  $\sum_{t=i}^\infty e^{t+1} \leq e^i + e^{i-1} + 4p^i$ .

*Proof.* We define  $w_t \triangleq p^t - p^{t+1}$ . We let  $1 \leq i < T$ .

First, for any  $i \geq 1$ , we have:

$$\sum_{t=i}^T w_t = \sum_{t=i}^T (p^t - p^{t+1}) = p^i - p^{T+1} \stackrel{\textcircled{1}}{\leq} p^i, \quad (12)$$

where step  $\textcircled{1}$  uses  $p^i \geq 0$  for all  $i$ .

Second, we obtain:

$$\begin{aligned} e^{t+1} &\stackrel{\textcircled{1}}{\leq} \sqrt{(e^t + e^{t-1})w_t} \\ &\stackrel{\textcircled{2}}{\leq} \sqrt{\frac{\alpha}{2}(e^t + e^{t-1})^2 + (w_t)^2/(2\alpha)}, \forall \alpha > 0 \\ &\stackrel{\textcircled{3}}{\leq} \sqrt{\frac{\alpha}{2}} \cdot (e^t + e^{t-1}) + w_t \sqrt{1/(2\alpha)}, \forall \alpha > 0. \end{aligned} \quad (13)$$

Here, step  $\textcircled{1}$  uses  $(e^{t+1})^2 \leq (e^t + e^{t-1})(p^t - p^{t+1})$  and  $w_t \triangleq p^t - p^{t+1}$ ; step  $\textcircled{2}$  uses the fact that  $ab \leq \frac{\alpha}{2}a^2 + \frac{1}{2\alpha}b^2$  for all  $\alpha > 0$ ; step  $\textcircled{3}$  uses the fact that  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for all  $a, b \geq 0$ .

Assume the parameter  $\alpha$  is sufficiently small that  $1 - 2\sqrt{\frac{\alpha}{2}} > 0$ . Telescoping Inequality (13) over  $t$  from  $i$  to  $T$ , we obtain:

$$\begin{aligned} &\sum_{t=i}^T w_t \sqrt{1/(2\alpha)} \\ &\geq \left\{ \sum_{t=i}^T e^{t+1} \right\} - \sqrt{\frac{\alpha}{2}} \left\{ \sum_{t=i}^T e^t \right\} - \sqrt{\frac{\alpha}{2}} \left\{ \sum_{t=i}^T e^{t-1} \right\} \\ &= \left\{ e^T + e^{T+1} + \sum_{t=i}^{T-2} e^{t+1} \right\} - \sqrt{\frac{\alpha}{2}} \left\{ e^i + e^T + \sum_{t=i}^{T-2} e^{t+1} \right\} \\ &\quad - \sqrt{\frac{\alpha}{2}} \left\{ e^{i-1} + e^i + \sum_{t=i}^{T-2} e^{t+1} \right\} \\ &= e^T + e^{T+1} - \sqrt{\frac{\alpha}{2}}(e^i + e^T + e^{i-1} + e^i) + (1 - 2\sqrt{\frac{\alpha}{2}}) \sum_{t=i}^{T-2} e^{t+1} \\ &\stackrel{\textcircled{1}}{\geq} e^T(1 - \sqrt{\frac{\alpha}{2}}) - \sqrt{\frac{\alpha}{2}}(e^i + e^{i-1} + e^i) + (1 - 2\sqrt{\frac{\alpha}{2}}) \sum_{t=i}^{T-2} e^{t+1} \\ &\stackrel{\textcircled{2}}{\geq} -2\sqrt{\frac{\alpha}{2}}(e^i + e^{i-1}) + (1 - 2\sqrt{\frac{\alpha}{2}}) \sum_{t=i}^{T-2} e^{t+1}, \end{aligned}$$

where step  $\textcircled{1}$  uses  $e^{T+1} \geq 0$ ; step  $\textcircled{2}$  uses  $1 - \sqrt{\frac{\alpha}{2}} > 1 - 2\sqrt{\frac{\alpha}{2}} > 0$ . This leads to:

$$\begin{aligned} \sum_{t=i}^{T-2} e^{t+1} &\leq (1 - 2\sqrt{\frac{\alpha}{2}})^{-1} \cdot \left\{ 2\sqrt{\frac{\alpha}{2}}(e^i + e^{i-1}) + \sqrt{\frac{1}{2\alpha}} \sum_{t=i}^T w_t \right\} \\ &\stackrel{\textcircled{1}}{=} (e^i + e^{i-1}) + 4 \sum_{t=i}^T w_t \\ &\stackrel{\textcircled{2}}{=} (e^i + e^{i-1}) + 4p^i, \end{aligned}$$

step ① uses the fact that  $(1 - 2\sqrt{\frac{\alpha}{2}})^{-1} \cdot 2\sqrt{\frac{\alpha}{2}} = 1$  and  $(1 - 2\sqrt{\frac{\alpha}{2}})^{-1} \cdot \sqrt{\frac{1}{2\alpha}} = 4$  when  $\alpha = 1/8$ ; step ② uses Inequalities (12). Letting  $T \rightarrow \infty$ , we conclude this lemma.  $\square$

**Lemma A.11.** Assume  $\sum_{t=1}^T (1/\tilde{\beta}^t) \geq \mathcal{O}(T^a)$ , where  $a \geq 0$  is a constant, and  $\{\tilde{\beta}^t\}_{t=1}^T$  is a nonnegative increasing sequence. If  $T$  is an even number, we have:  $\sum_{t=1}^{T/2} (1/\tilde{\beta}^{2t}) \geq \mathcal{O}(T^a)$ .

*Proof.* We have:  $\sum_{t=1}^{T/2} \frac{1}{\tilde{\beta}^{2t}} = \frac{1}{2} \sum_{t=1}^{T/2} (\frac{1}{\tilde{\beta}^{2t}} + \frac{1}{\tilde{\beta}^{2t}}) \stackrel{\textcircled{1}}{\geq} \frac{1}{2} \sum_{t=1}^{T/2} (\frac{1}{\tilde{\beta}^{2t}} + \frac{1}{\tilde{\beta}^{2t+1}}) = \frac{1}{\tilde{\beta}^{2T+1}} - \frac{1}{\tilde{\beta}^1} + \sum_{t=1}^T \frac{1}{\tilde{\beta}^t} = \mathcal{O}(\sum_{t=1}^T \frac{1}{\tilde{\beta}^t}) \geq \mathcal{O}(T^a)$ , where step ① uses the fact that  $\{\tilde{\beta}^t\}_{t=1}^T$  is increasing.  $\square$

**Lemma A.12.** Assume that  $\frac{d^t}{d^{t-2}} \leq \frac{\dot{\beta}^t+1}{\dot{\beta}^t+2}$ , and  $\sum_{i=0}^T (1/\dot{\beta}^i) \geq \mathcal{O}(T^a)$ , where  $a \geq 0$  is a positive constant,  $\{d^t\}_{t=0}^\infty$  and  $\{\dot{\beta}^t\}_{t=0}^\infty$  are two nonnegative sequences. Assume that  $\{\dot{\beta}^t\}_{t=0}^\infty$  is increasing. We have:  $d^T \leq \mathcal{O}(1/\exp(T^a))$ .

*Proof.* We define  $\gamma^t \triangleq \frac{1}{\dot{\beta}^t+2} \in (0, 1)$ .

Given  $\frac{d^t}{d^{t-2}} \leq \frac{\dot{\beta}^t+1}{\dot{\beta}^t+2}$ , we have  $\frac{d^t}{d^{t-2}} \leq 1 - \gamma^t$ , leading to:

$$d^{2t} \leq d^0 (1 - \gamma^2)(1 - \gamma^4)(1 - \gamma^6) \dots (1 - \gamma^{2t}). \quad (14)$$

**Part (a).** When  $T$  is an even number, we have:

$$\begin{aligned} d^T &= \exp(\log(d^T)) \\ &\stackrel{\textcircled{1}}{\leq} \exp(\log(d^0 \cdot \prod_{t=1}^{T/2} (1 - \gamma^{2t}))) \\ &\stackrel{\textcircled{2}}{=} \exp(\log(d^0) + \sum_{t=1}^{T/2} \log(1 - \gamma^{2t})) \\ &\stackrel{\textcircled{3}}{\leq} \exp(\log(d^0) + \sum_{t=1}^{T/2} (-\gamma^{2t})) \\ &\stackrel{\textcircled{4}}{\leq} \exp(\log(d^0)) \times \{\exp(\sum_{t=1}^{T/2} (\gamma^{2t}))\}^{-1} \\ &\stackrel{\textcircled{5}}{\leq} d^0 \times \{\exp(\mathcal{O}(T^a))\}^{-1} = \mathcal{O}(1/\exp(T^a)), \end{aligned}$$

where step ① uses Inequality (14); step ② uses  $\log(ab) = \log(a) + \log(b)$  for all  $a > 0$  and  $b > 0$ ; step ③ uses  $\log(1 - x) \leq -x$  for all  $x \in (0, 1)$ , and  $1 - \gamma^t \in (0, 1)$  for all  $t$ ; step ④ uses  $\exp(a + b) = \exp(a)\exp(b)$  for all  $a > 0$  and  $b > 0$ ; step ⑤ uses Lemma A.11 with  $\tilde{\beta}^t = 1/\gamma^t = \dot{\beta}^t + 2$ .

**Part (b).** When  $T$  is an odd number, analogous strategies result in the same complexity outcome.  $\square$

**Lemma A.13.** Assume that  $[d^t]^{\tau+1} \leq \dot{\beta}^t (d^{t-2} - d^t)$ , and  $\sum_{i=1}^T (1/\dot{\beta}^i) \geq \mathcal{O}(T^a)$ , where  $\tau, a > 0$  are positive constants,  $\{d^t\}_{t=0}^\infty$  and  $\{\dot{\beta}^t\}_{t=0}^\infty$  are two nonnegative sequences. Assume that  $\{\dot{\beta}^t\}_{t=0}^\infty$  is increasing. We have:  $d^T \leq \mathcal{O}(1/(T^a/\tau))$ .

*Proof.* We let  $\kappa > 1$  be any constant. We define  $h(s) = s^{-\tau-1}$ , where  $\tau > 0$ .

We consider two cases for  $h(d^t)/h(d^{t-2})$ .



1134 We define  $h_{\mu_1}(\mathbf{y}) \triangleq \min_{\mathbf{v}} h(\mathbf{v}) + \frac{1}{2\mu_1} \|\mathbf{v} - \mathbf{y}\|_2^2$ , and  $\mathbb{P}_{\mu_1}(\mathbf{y}) = \arg \min_{\mathbf{v}} h(\mathbf{v}) + \frac{1}{2\mu_1} \|\mathbf{v} - \mathbf{y}\|_2^2$ .

1136 We define  $h_{\mu_2}(\mathbf{y}) \triangleq \min_{\mathbf{v}} h(\mathbf{v}) + \frac{1}{2\mu_2} \|\mathbf{v} - \mathbf{y}\|_2^2$ , and  $\mathbb{P}_{\mu_2}(\mathbf{y}) = \arg \min_{\mathbf{v}} h(\mathbf{v}) + \frac{1}{2\mu_2} \|\mathbf{v} - \mathbf{y}\|_2^2$ .

1137 By the optimality of  $\mathbb{P}_{\mu_1}(\mathbf{y})$  and  $\mathbb{P}_{\mu_2}(\mathbf{y})$ , we obtain:

$$1139 \quad \mathbf{y} - \mathbb{P}_{\mu_1}(\mathbf{y}) \in \mu_1 \partial h(\mathbb{P}_{\mu_1}(\mathbf{y})) \quad (19)$$

$$1140 \quad \mathbf{y} - \mathbb{P}_{\mu_2}(\mathbf{y}) \in \mu_2 \partial h(\mathbb{P}_{\mu_2}(\mathbf{y})). \quad (20)$$

1142 **Part (a).** We now prove that  $0 \leq h_{\mu_2}(\mathbf{y}) - h_{\mu_1}(\mathbf{y})$ . For any  $\mathbf{s}_1 \in \partial h(\mathbb{P}_{\mu_1}(\mathbf{y}))$  and  $\mathbf{s}_2 \in \partial h(\mathbb{P}_{\mu_2}(\mathbf{y}))$ , we have:

$$\begin{aligned} 1145 & h_{\mu_1}(\mathbf{y}) - h_{\mu_2}(\mathbf{y}) \\ 1146 & \stackrel{\textcircled{1}}{=} \frac{1}{2\mu_1} \|\mathbf{y} - \mathbb{P}_{\mu_1}(\mathbf{y})\|_2^2 - \frac{1}{2\mu_2} \|\mathbf{y} - \mathbb{P}_{\mu_2}(\mathbf{y})\|_2^2 + h(\mathbb{P}_{\mu_1}(\mathbf{y})) - h(\mathbb{P}_{\mu_2}(\mathbf{y})) \\ 1147 & \stackrel{\textcircled{2}}{\leq} \frac{1}{2\mu_1} \|\mathbf{y} - \mathbb{P}_{\mu_1}(\mathbf{y})\|_2^2 - \frac{1}{2\mu_2} \|\mathbf{y} - \mathbb{P}_{\mu_2}(\mathbf{y})\|_2^2 + \langle \mathbb{P}_{\mu_1}(\mathbf{y}) - \mathbb{P}_{\mu_2}(\mathbf{y}), \mathbf{s}_1 \rangle + \frac{W_h}{2} \|\mathbb{P}_{\mu_2}(\mathbf{y}) - \mathbb{P}_{\mu_1}(\mathbf{y})\|_2^2 \\ 1148 & \stackrel{\textcircled{3}}{=} \frac{1}{2\mu_1} \|\mu_1 \mathbf{s}_1\|_2^2 - \frac{1}{2\mu_2} \|\mu_2 \mathbf{s}_2\|_2^2 + \langle \mu_2 \mathbf{s}_2 - \mu_1 \mathbf{s}_1, \mathbf{s}_1 \rangle + \frac{W_h}{2} \|\mu_1 \mathbf{s}_1 - \mu_2 \mathbf{s}_2\|_2^2 \\ 1149 & \stackrel{\textcircled{4}}{\leq} \frac{1}{2\mu_1} \|\mu_1 \mathbf{s}_1\|_2^2 - \frac{1}{2\mu_2} \|\mu_2 \mathbf{s}_2\|_2^2 + \langle \mu_2 \mathbf{s}_2 - \mu_1 \mathbf{s}_1, \mathbf{s}_1 \rangle + \frac{1}{2\mu_1} \|\mu_1 \mathbf{s}_1 - \mu_2 \mathbf{s}_2\|_2^2 \\ 1150 & = -\frac{\mu_2}{2} \|\mathbf{s}_2\|_2^2 \cdot \left(1 - \frac{\mu_2}{\mu_1}\right) \\ 1151 & \stackrel{\textcircled{5}}{\leq} 0, \end{aligned}$$

1152 where step ① uses the definition of  $h_{\mu_1}(\mathbf{y})$  and  $h_{\mu_2}(\mathbf{y})$ ; step ② uses weakly convexity of  $h(\cdot)$ ; step  
1153 ③ uses the optimality of  $\mathbb{P}_{\mu_1}(\mathbf{y})$  and  $\mathbb{P}_{\mu_2}(\mathbf{y})$  in Equations (19) and (20); step ④ uses  $W_h \leq \frac{1}{\mu_1}$ ; step  
1154 ⑤ uses  $1 \geq \frac{\mu_2}{\mu_1}$ .

1155 **Part (b).** We now prove that  $h_{\mu_2}(\mathbf{y}) - h_{\mu_1}(\mathbf{y}) \leq \min\{\frac{\mu_1}{2\mu_2}, 1\} \cdot (\mu_1 - \mu_2) C_h^2$ . For any  $\mathbf{s}_1 \in$   
1156  $\partial h(\mathbb{P}_{\mu_1}(\mathbf{y}))$  and  $\mathbf{s}_2 \in \partial h(\mathbb{P}_{\mu_2}(\mathbf{y}))$ , we have:

$$\begin{aligned} 1164 & h_{\mu_2}(\mathbf{y}) - h_{\mu_1}(\mathbf{y}) \\ 1165 & \stackrel{\textcircled{1}}{=} \frac{1}{2\mu_2} \|\mathbf{y} - \mathbb{P}_{\mu_2}(\mathbf{y})\|_2^2 - \frac{1}{2\mu_1} \|\mathbf{y} - \mathbb{P}_{\mu_1}(\mathbf{y})\|_2^2 + h(\mathbb{P}_{\mu_2}(\mathbf{y})) - h(\mathbb{P}_{\mu_1}(\mathbf{y})) \\ 1166 & \stackrel{\textcircled{2}}{\leq} \frac{1}{2\mu_2} \|\mathbf{y} - \mathbb{P}_{\mu_2}(\mathbf{y})\|_2^2 - \frac{1}{2\mu_1} \|\mathbf{y} - \mathbb{P}_{\mu_1}(\mathbf{y})\|_2^2 + \langle \mathbb{P}_{\mu_2}(\mathbf{y}) - \mathbb{P}_{\mu_1}(\mathbf{y}), \mathbf{s}_1 \rangle + \frac{W_h}{2} \|\mathbb{P}_{\mu_2}(\mathbf{y}) - \mathbb{P}_{\mu_1}(\mathbf{y})\|_2^2 \\ 1167 & \stackrel{\textcircled{3}}{=} \frac{\mu_2}{2} \|\mathbf{s}_1\|_2^2 - \frac{\mu_1}{2} \|\mathbf{s}_2\|_2^2 + \langle \mu_1 \mathbf{s}_2 - \mu_2 \mathbf{s}_1, \mathbf{s}_1 \rangle + \frac{W_h}{2} \|\mu_1 \mathbf{s}_2 - \mu_2 \mathbf{s}_1\|_2^2 \\ 1168 & = -\frac{\mu_2}{2} \|\mathbf{s}_1\|_2^2 - \frac{\mu_1}{2} \|\mathbf{s}_2\|_2^2 + \mu_1 \langle \mathbf{s}_1, \mathbf{s}_2 \rangle + \frac{W_h}{2} \|\mu_1 \mathbf{s}_2 - \mu_2 \mathbf{s}_1\|_2^2 \\ 1169 & \stackrel{\textcircled{4}}{\leq} \min\left\{-\frac{\mu_2}{2} \|\mathbf{s}_2\|_2^2 + \mu_1 \langle \mathbf{s}_1, \mathbf{s}_2 \rangle + \frac{1}{2\mu_2} \|\mu_1 \mathbf{s}_2 - \mu_2 \mathbf{s}_1\|_2^2 - \frac{\mu_2}{2} \|\mathbf{s}_1\|_2^2,\right. \\ 1170 & \quad \left. -\frac{\mu_1}{2} \|\mathbf{s}_2\|_2^2 + \mu_1 \langle \mathbf{s}_1, \mathbf{s}_2 \rangle + \frac{1}{2\mu_1} \|\mu_1 \mathbf{s}_2 - \mu_2 \mathbf{s}_1\|_2^2 - \frac{\mu_2}{2} \|\mathbf{s}_1\|_2^2\right\} \\ 1171 & = \min\left\{(-\mu_2 + \mu_1) \cdot \frac{\mu_1}{2\mu_2} \|\mathbf{s}_2\|_2^2, (\mu_1 - \mu_2) \langle \mathbf{s}_1, \mathbf{s}_2 \rangle - \frac{\mu_2}{2} \|\mathbf{s}_1\|_2^2 + \frac{\mu_2^2}{2\mu_1} \|\mathbf{s}_1\|_2^2\right\} \\ 1172 & \stackrel{\textcircled{5}}{\leq} \min\left\{\frac{\mu_1}{2\mu_2} \|\mathbf{s}_2\|_2^2 \cdot (\mu_1 - \mu_2), (\mu_1 - \mu_2) \langle \mathbf{s}_1, \mathbf{s}_2 \rangle\right\} \\ 1173 & \stackrel{\textcircled{6}}{\leq} \min\left\{\frac{\mu_1}{2\mu_2} \cdot (\mu_1 - \mu_2), (\mu_1 - \mu_2)\right\} \cdot C_h^2 \\ 1174 & = \min\left\{\frac{\mu_1}{2\mu_2}, 1\right\} \cdot (\mu_1 - \mu_2) \cdot C_h^2, \end{aligned}$$

1175 where step ① uses the definition of  $h_{\mu_1}(\mathbf{y})$  and  $h_{\mu_2}(\mathbf{y})$ ; step ② uses the weakly convexity of  $h(\cdot)$ ;  
1176 step ③ uses the optimality of  $\mathbb{P}_{\mu_2}(\mathbf{y})$  and  $\mathbb{P}_{\mu_1}(\mathbf{y})$  in Equations (19) and (20); step ④ uses  $W_h \leq \frac{1}{\mu_1}$   
1177 and  $W_h \leq \frac{1}{\mu_2}$ ; step ⑤ uses  $\mu_2 \leq \mu_1$ ; step ⑥ uses  $\|\mathbf{s}_1\| \leq C_h$ ,  $\|\mathbf{s}_2\| \leq C_h$ , and  $\langle \mathbf{s}_1, \mathbf{s}_2 \rangle \leq$   
1178  $\|\mathbf{s}_1\| \cdot \|\mathbf{s}_2\| \leq C_h^2$ .

1187  $\square$

## B.2 PROOF OF LEMMA 2.4

*Proof.* Assume  $0 < \mu_2 < \mu_1 \leq \frac{1}{2W_h}$ , and fixing  $\mathbf{y} \in \mathbb{R}^m$ .

Using the result in Lemma 2.2, we establish that the gradient of  $h_\mu(\mathbf{y})$  w.r.t  $\mathbf{y}$  can be computed as:

$$\nabla h_\mu(\mathbf{y}) = \mu^{-1}(\mathbf{y} - \mathbb{P}_\mu(\mathbf{y})).$$

The gradient of the mapping  $\nabla h_\mu(\mathbf{y})$  w.r.t. the variable  $1/\mu$  can be computed as:  $\nabla_{1/\mu}(\nabla h_\mu(\mathbf{y})) = \mathbf{y} - \mathbb{P}_\mu(\mathbf{y})$ . We further obtain:

$$\|\nabla_{1/\mu}(\nabla h_\mu(\mathbf{y}))\| = \|\mathbf{y} - \mathbb{P}_\mu(\mathbf{y})\| \stackrel{\textcircled{1}}{=} \mu \|\partial h(\mathbb{P}_\mu(\mathbf{y}))\| \leq \mu C_h.$$

Here, step  $\textcircled{1}$  uses the optimality of  $\mathbb{P}_\mu(\mathbf{y})$  that:  $\mathbf{0} \in \partial h(\mathbb{P}_\mu(\mathbf{y})) + \frac{1}{\mu}(\mathbb{P}_\mu(\mathbf{y}) - \mathbf{y})$ . Therefore, for all  $\mu \in (0, \frac{1}{2W_h}]$ , we have:

$$\frac{\|\nabla h_\mu(\mathbf{y}) - \nabla h_{\mu'}(\mathbf{y})\|_2}{|1/\mu - 1/\mu'|} \leq \mu C_h.$$

Letting  $\mu = \mu_1$  and  $\mu' = \mu_2$ , we have:  $\|\nabla h_{\mu_1}(\mathbf{y}) - \nabla h_{\mu_2}(\mathbf{y})\|_2 \leq |1 - \mu_1/\mu_2| C_h = (\mu_1/\mu_2 - 1) C_h$ .  $\square$

## B.3 PROOF OF LEMMA 2.5

*Proof.* We consider the following optimization problem:

$$\bar{\mathbf{y}} = \arg \min_{\mathbf{y}} h_\mu(\mathbf{y}) + \frac{\beta}{2} \|\mathbf{y} - \mathbf{b}\|_2^2. \quad (21)$$

Given  $h_\mu(\mathbf{y})$  being  $(\mu^{-1})$ -weakly convex and  $\beta > \mu^{-1}$ , Problem (21) becomes strongly convex and has a unique optimal solution, which leads to the following equivalent problem:

$$(\bar{\mathbf{y}}, \check{\mathbf{y}}) = \arg \min_{\mathbf{y}, \mathbf{y}'} h(\mathbf{y}') + \frac{1}{2\mu} \|\mathbf{y} - \mathbf{y}'\|_2^2 + \frac{\beta}{2} \|\mathbf{y} - \mathbf{b}\|_2^2,$$

We have the following first-order optimality conditions for  $(\bar{\mathbf{y}}, \check{\mathbf{y}})$ :

$$\frac{1}{\mu}(\bar{\mathbf{y}} - \check{\mathbf{y}}) = \beta(\mathbf{b} - \bar{\mathbf{y}}) \quad (22)$$

$$\frac{1}{\mu}(\bar{\mathbf{y}} - \check{\mathbf{y}}) \in \partial h(\check{\mathbf{y}}). \quad (23)$$

**Part (a).** We have the following results:

$$\begin{aligned} \mathbf{0} &\stackrel{\textcircled{1}}{\in} \partial h(\check{\mathbf{y}}) + \frac{1}{\mu}(\check{\mathbf{y}} - \bar{\mathbf{y}}) \\ &\stackrel{\textcircled{2}}{=} \partial h(\check{\mathbf{y}}) + \frac{1}{\mu}(\check{\mathbf{y}} - \frac{1}{1/\mu + \beta}(\frac{1}{\mu}\check{\mathbf{y}} + \beta\mathbf{b})) \\ &= \partial h(\check{\mathbf{y}}) + \frac{\beta}{1 + \mu\beta}(\check{\mathbf{y}} - \mathbf{b}), \end{aligned} \quad (24)$$

where step  $\textcircled{1}$  uses Equality (23); step  $\textcircled{2}$  uses Equality (22) that  $\bar{\mathbf{y}} = \frac{1}{1/\mu + \beta}(\frac{1}{\mu}\check{\mathbf{y}} + \beta\mathbf{b})$ . The inclusion in (24) implies that:

$$\check{\mathbf{y}} = \arg \min_{\check{\mathbf{y}}} h(\check{\mathbf{y}}) + \frac{1}{2} \cdot \frac{\beta}{1 + \mu\beta} \|\check{\mathbf{y}} - \mathbf{b}\|_2^2.$$

**Part (b).** Combining Equalities (22) and (23), we have:  $\beta(\mathbf{b} - \bar{\mathbf{y}}) \in \partial h(\check{\mathbf{y}})$ .

**Part (c).** In view of Equation (23), we have:  $\bar{\mathbf{y}} - \check{\mathbf{y}} = \mu \partial h(\check{\mathbf{y}})$ , leading to:  $\|\check{\mathbf{y}} - \bar{\mathbf{y}}\| \leq \mu C_h$ .  $\square$

## B.4 PROOFS FOR LEMMA 2.11

*Proof.* We let  $\Delta \in \mathbb{R}^{n \times r}$  and  $\mathbf{X} \in \mathcal{M}$ . We define  $\mathbf{U} \triangleq \Delta^\top \mathbf{X} \in \mathbb{R}^{r \times r}$ .

We derive the following results:

$$\begin{aligned}
& \|\text{Proj}_{\mathbf{T}_X \mathcal{M}}(\Delta)\|_{\mathbb{F}}^2 - \|\Delta\|_{\mathbb{F}}^2 \\
& \stackrel{\textcircled{1}}{=} \|\Delta - \frac{1}{2} \mathbf{X}(\Delta^\top \mathbf{X} + \mathbf{X}^\top \Delta)\|_{\mathbb{F}}^2 - \|\Delta\|_{\mathbb{F}}^2 \\
& = \frac{1}{4} \|\mathbf{X}(\Delta^\top \mathbf{X} + \mathbf{X}^\top \Delta)\|_{\mathbb{F}}^2 - \langle \Delta, \mathbf{X}(\Delta^\top \mathbf{X} + \mathbf{X}^\top \Delta) \rangle \\
& \stackrel{\textcircled{2}}{=} \frac{1}{4} \|\Delta^\top \mathbf{X} + \mathbf{X}^\top \Delta\|_{\mathbb{F}}^2 - \langle \Delta, \mathbf{X}(\Delta^\top \mathbf{X} + \mathbf{X}^\top \Delta) \rangle \\
& \stackrel{\textcircled{3}}{=} \frac{1}{4} \|\mathbf{U} + \mathbf{U}^\top\|_{\mathbb{F}}^2 - \langle \mathbf{U} + \mathbf{U}^\top, \mathbf{U} \rangle \\
& \stackrel{\textcircled{4}}{=} \frac{1}{4} \|\mathbf{U} + \mathbf{U}^\top\|_{\mathbb{F}}^2 - \langle \mathbf{U} + \mathbf{U}^\top, \mathbf{U} + \mathbf{U}^\top \rangle \cdot \frac{1}{2} \\
& = -\frac{1}{4} \|\mathbf{U} + \mathbf{U}^\top\|_{\mathbb{F}}^2 \leq 0,
\end{aligned}$$

where step ① uses  $\text{Proj}_{\mathbf{T}_X \mathcal{M}}(\Delta) = \Delta - \frac{1}{2} \mathbf{X}(\Delta^\top \mathbf{X} + \mathbf{X}^\top \Delta)$  for all  $\Delta \in \mathbb{R}^{n \times r}$  (Absil et al., 2008a); step ② uses the fact that  $\|\mathbf{X}\mathbf{P}\|_{\mathbb{F}}^2 = \text{tr}(\mathbf{P}\mathbf{X}^\top \mathbf{X}\mathbf{P}^\top) = \|\mathbf{P}\|_{\mathbb{F}}^2$  for all  $\mathbf{X} \in \mathcal{M}$ ; step ③ uses the definition of  $\mathbf{U} \triangleq \Delta^\top \mathbf{X}$ ; step ④ uses the symmetric properties of the matrix  $(\mathbf{U} + \mathbf{U}^\top)$ .

□

## B.5 PROOF OF LEMMA 2.12

*Proof.* We let  $\rho > 0$ ,  $\mathbf{G} \in \mathbb{R}^{n \times r}$ , and  $\mathbf{X} \in \mathcal{M}$ .

We define  $\mathbf{U} \triangleq \mathbf{G}^\top \mathbf{X}$ , and  $\mathbb{G}_\rho \triangleq \mathbf{G} - \rho \mathbf{X}\mathbf{G}^\top \mathbf{X} - (1 - \rho) \mathbf{X}\mathbf{X}^\top \mathbf{G}$ .

First, we have the following equalities:

$$\begin{aligned}
\langle \mathbf{G}, \mathbb{G}_\rho \rangle &= \langle \mathbf{G}, \mathbf{G} - \rho \mathbf{X}\mathbf{G}^\top \mathbf{X} - (1 - \rho) \mathbf{X}\mathbf{X}^\top \mathbf{G} \rangle \\
&= \langle \mathbf{G}, \mathbf{G} \rangle - \rho \text{tr}(\mathbf{G}^\top \mathbf{X}\mathbf{G}^\top \mathbf{X}) - (1 - \rho) \text{tr}(\mathbf{G}^\top \mathbf{X}\mathbf{X}^\top \mathbf{G}) \\
&\stackrel{\textcircled{1}}{=} \langle \mathbf{G}, \mathbf{G} \rangle - \rho \text{tr}(\mathbf{U}\mathbf{U}) - (1 - \rho) \text{tr}(\mathbf{U}\mathbf{U}^\top), \tag{25}
\end{aligned}$$

where step ① uses  $\mathbf{U} \triangleq \mathbf{G}^\top \mathbf{X}$ .

Second, we derive the following equalities:

$$\begin{aligned}
\|\mathbb{G}_\rho\|_{\mathbb{F}}^2 &= \langle \rho \mathbf{X}\mathbf{G}^\top \mathbf{X} + (1 - \rho) \mathbf{X}\mathbf{X}^\top \mathbf{G} - \mathbf{G}, \rho \mathbf{X}\mathbf{G}^\top \mathbf{X} + (1 - \rho) \mathbf{X}\mathbf{X}^\top \mathbf{G} - \mathbf{G} \rangle \\
&\stackrel{\textcircled{1}}{=} \rho^2 \text{tr}(\mathbf{U}^\top \mathbf{U}) + \rho(1 - \rho) \text{tr}(\mathbf{U}^\top \mathbf{U}^\top) - \rho \text{tr}(\mathbf{U}^\top \mathbf{U}^\top) \\
&\quad + (1 - \rho)\rho \text{tr}(\mathbf{U}\mathbf{U}) + (1 - \rho)^2 \text{tr}(\mathbf{U}\mathbf{U}^\top) - (1 - \rho) \text{tr}(\mathbf{U}\mathbf{U}^\top) \\
&\quad - \rho \text{tr}(\mathbf{U}\mathbf{U}) - (1 - \rho) \text{tr}(\mathbf{U}\mathbf{U}^\top) + \langle \mathbf{G}, \mathbf{G} \rangle \\
&\stackrel{\textcircled{2}}{=} (2\rho^2 - 1) \cdot \text{tr}(\mathbf{U}^\top \mathbf{U}) - 2\rho^2 \cdot \text{tr}(\mathbf{U}\mathbf{U}) + \langle \mathbf{G}, \mathbf{G} \rangle, \tag{26}
\end{aligned}$$

where step ① uses  $\mathbf{U} \triangleq \mathbf{G}^\top \mathbf{X}$  and  $\mathbf{X}^\top \mathbf{X} = \mathbf{I}_r$ ; step ② uses  $\text{tr}(\mathbf{U}^\top \mathbf{U}^\top) = \text{tr}(\mathbf{U}\mathbf{U})$ .

Third, we have:

$$\text{tr}(\mathbf{G}^\top \mathbf{G}) - \text{tr}(\mathbf{U}^\top \mathbf{U}) \stackrel{\textcircled{1}}{=} \langle \mathbf{G}\mathbf{G}^\top, \mathbf{I}_n - \mathbf{X}\mathbf{X}^\top \rangle \stackrel{\textcircled{2}}{\geq} 0, \tag{27}$$

where step ① uses  $\mathbf{U} \triangleq \mathbf{G}^\top \mathbf{X}$ ; step ② uses the fact that the matrix  $(\mathbf{I}_n - \mathbf{X}\mathbf{X}^\top)$  only contains eigenvalues that are 0 or 1.

**Part (a-i).** We now prove that  $\max(1, 2\rho) \langle \mathbf{G}, \mathbb{G}_\rho \rangle \geq \|\mathbb{G}_\rho\|_{\mathbb{F}}^2$ . We discuss two cases. Case (i):  $\rho \in (0, \frac{1}{2}]$ . We have:

$$\|\mathbb{G}_\rho\|_{\mathbb{F}}^2 - \langle \mathbf{G}, \mathbb{G}_\rho \rangle \stackrel{\textcircled{1}}{=} (2\rho^2 - \rho) \cdot (\text{tr}(\mathbf{U}\mathbf{U}^\top) - \text{tr}(\mathbf{U}\mathbf{U})) \stackrel{\textcircled{2}}{\leq} 0,$$

1296 where step ① uses Inequalities (25) and (26); step ② uses  $2\rho^2 - \rho \leq 0$  for all  $\rho \in (0, \frac{1}{2}]$ , and  
 1297  $\text{tr}(\mathbf{U}\mathbf{U}) \leq \text{tr}(\mathbf{U}\mathbf{U}^\top)$  for all  $\mathbf{U} \in \mathbb{R}^{r \times r}$ .  
 1298 Case (ii):  $\rho \in [\frac{1}{2}, \infty)$ . We have:

$$1299 \quad \|\mathbb{G}_\rho\|_{\mathbb{F}}^2 - 2\rho\langle \mathbf{G}, \mathbb{G}_\rho \rangle \stackrel{\textcircled{1}}{=} (2\rho - 1)(\text{tr}(\mathbf{U}\mathbf{U}^\top) - \langle \mathbf{G}, \mathbf{G} \rangle) \stackrel{\textcircled{2}}{\leq} 0,$$

1302 where step ① uses Inequalities (25) and (26); step ② uses  $2\rho - 1 \geq 0$  for all  $\rho \in [\frac{1}{2}, \infty)$ , and  
 1303 Inequality(27). Therefore, we conclude that:  $\max(1, 2\rho)\langle \mathbf{G}, \mathbb{G}_\rho \rangle \geq \|\mathbb{G}_\rho\|_{\mathbb{F}}^2$ .

1304 **Part (a-ii).** We now prove that  $\|\mathbb{G}_\rho\|_{\mathbb{F}}^2 \geq \min(1, \rho^2)\|\mathbb{G}_1\|_{\mathbb{F}}^2$ . We consider two cases. Case (i):  
 1305  $\rho \in (0, 1]$ . We have:

$$1307 \quad \rho^2\|\mathbb{G}_1\|_{\mathbb{F}}^2 - \|\mathbb{G}_\rho\|_{\mathbb{F}}^2 \stackrel{\textcircled{1}}{=} (1 - \rho^2)(\text{tr}(\mathbf{U}^\top\mathbf{U}) - \langle \mathbf{G}, \mathbf{G} \rangle) \stackrel{\textcircled{2}}{\leq} 0,$$

1309 where step ① uses Inequalities (25) and (26); step ② uses  $1 - \rho^2 \geq 0$ , and Inequality (27).  
 1310 Case (ii):  $\rho \in (1, \infty)$ . We have:

$$1312 \quad \|\mathbb{G}_1\|_{\mathbb{F}}^2 - \|\mathbb{G}_\rho\|_{\mathbb{F}}^2 \stackrel{\textcircled{1}}{=} (2 - 2\rho^2)(\text{tr}(\mathbf{U}^\top\mathbf{U}) - \text{tr}(\mathbf{U}\mathbf{U})) \leq 0,$$

1314 where step ① uses Inequality (26); step ② uses  $4\rho^2 - 1 \leq 0$  for all  $\rho \in (0, \frac{1}{2}]$ , and the fact that  
 1315  $\text{tr}(\mathbf{U}\mathbf{U}) - \text{tr}(\mathbf{U}\mathbf{U}^\top) \leq 0$  for all  $\mathbf{U} \in \mathbb{R}^{r \times r}$ . Therefore, we conclude that:  $\min(1, \rho^2)\|\mathbb{G}_1\|_{\mathbb{F}}^2 \leq$   
 1316  $\|\mathbb{G}_\rho\|_{\mathbb{F}}^2$ .

1317 **Part (b-i).** We now prove that  $\|\mathbb{G}_\rho\|_{\mathbb{F}} \geq \min(1, 2\rho)\|\mathbb{G}_{1/2}\|_{\mathbb{F}}$ . We consider two cases. Case (i):  
 1318  $\rho \in (0, \frac{1}{2}]$ . We have:

$$1320 \quad (2\rho)^2\|\mathbb{G}_{1/2}\|_{\mathbb{F}}^2 - \|\mathbb{G}_\rho\|_{\mathbb{F}}^2 \stackrel{\textcircled{1}}{=} (4\rho^2 - 1) \cdot (\text{tr}(\mathbf{G}^\top\mathbf{G}) - \text{tr}(\mathbf{U}^\top\mathbf{U})) \stackrel{\textcircled{2}}{\leq} 0,$$

1322 where step ① uses Inequality (26); step ② uses  $4\rho^2 - 1 \leq 0$  for all  $\rho \in (0, \frac{1}{2}]$ , and Inequality (27).  
 1323 Case (ii):  $\rho \in (\frac{1}{2}, \infty)$ . We have:

$$1325 \quad \|\mathbb{G}_{1/2}\|_{\mathbb{F}}^2 - \|\mathbb{G}_\rho\|_{\mathbb{F}}^2 \stackrel{\textcircled{1}}{=} (2\rho^2 - \frac{1}{2}) \cdot (\text{tr}(\mathbf{U}\mathbf{U}) - \text{tr}(\mathbf{U}^\top\mathbf{U})) \stackrel{\textcircled{2}}{\leq} 0,$$

1327 where step ① uses Inequalities (25) and (26); step ② uses  $2\rho^2 - \frac{1}{2} \geq 0$  for all  $\rho \in (\frac{1}{2}, \infty)$ , and  
 1328 the fact that  $\text{tr}(\mathbf{U}\mathbf{U}) - \text{tr}(\mathbf{U}\mathbf{U}^\top) \leq 0$  for all  $\mathbf{U} \in \mathbb{R}^{r \times r}$ . Therefore, we conclude that  $\|\mathbb{G}_\rho\|_{\mathbb{F}} \geq$   
 1329  $\min(1, 2\rho)\|\mathbb{G}_{1/2}\|_{\mathbb{F}}$ .

1331 **Part (b-ii).** We now prove that  $\|\mathbb{G}_\rho\|_{\mathbb{F}} \leq \max(1, 2\rho)\|\mathbb{G}_{1/2}\|_{\mathbb{F}}$ . We consider two cases. Case (i):  
 1332  $\rho \in (0, \frac{1}{2}]$ . We have:

$$1334 \quad \|\mathbb{G}_{1/2}\|_{\mathbb{F}}^2 - \|\mathbb{G}_\rho\|_{\mathbb{F}}^2 \stackrel{\textcircled{1}}{=} (2\rho^2 - \frac{1}{2}) \cdot (\text{tr}(\mathbf{U}\mathbf{U}) - \text{tr}(\mathbf{U}^\top\mathbf{U})) \stackrel{\textcircled{2}}{\geq} 0,$$

1336 where step ① uses Inequality (26); step ② uses  $2\rho^2 - \frac{1}{2} \leq 0$  for all  $\rho \in (0, \frac{1}{2}]$ , and the fact that  
 1337  $\text{tr}(\mathbf{U}\mathbf{U}) - \text{tr}(\mathbf{U}\mathbf{U}^\top) \leq 0$  for all  $\mathbf{U} \in \mathbb{R}^{r \times r}$ .

1338 Case (ii):  $\rho \in (\frac{1}{2}, \infty)$ . We have:

$$1339 \quad (2\rho)^2\|\mathbb{G}_{1/2}\|_{\mathbb{F}}^2 - \|\mathbb{G}_\rho\|_{\mathbb{F}}^2 \stackrel{\textcircled{1}}{=} (4\rho^2 - 1) \cdot (\text{tr}(\mathbf{G}^\top\mathbf{G}) - \text{tr}(\mathbf{U}^\top\mathbf{U})) \stackrel{\textcircled{2}}{\geq} 0,$$

1342 where step ① uses Inequalities (25) and (26); step ② uses  $4\rho^2 - 1 \geq 0$  for all  $\rho \in (\frac{1}{2}, \infty)$ , and  
 1343 Inequality (27). Therefore, we conclude that:  $\|\mathbb{G}_\rho\|_{\mathbb{F}} \geq \min(1, 2\rho)\|\mathbb{G}_{1/2}\|_{\mathbb{F}}$ .

□

## 1346 B.6 PROOF OF LEMMA 2.13

1347 *Proof.* Recall that the following first-order optimality conditions are equivalent for all  $\mathbf{X} \in \mathbb{R}^{n \times r}$ :

$$1348 \quad (\mathbf{0} \in \partial\mathcal{I}_{\mathcal{M}}(\mathbf{X}) + \nabla f(\mathbf{X})) \Leftrightarrow (\mathbf{0} \in \text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\nabla f(\mathbf{X}))). \quad (28)$$

Therefore, we derive the following results:

$$\begin{aligned}
\text{dist}(\mathbf{0}, \partial \mathcal{I}_{\mathcal{M}}(\mathbf{X}) + \nabla f(\mathbf{X})) &= \inf_{\mathbf{R} \in \nabla f(\mathbf{X}) + \partial \mathcal{I}_{\mathcal{M}}(\mathbf{X})} \|\mathbf{R}\|_{\text{F}} \\
&\stackrel{\textcircled{1}}{=} \inf_{\mathbf{R} \in \text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\nabla f(\mathbf{X}))} \|\mathbf{R}\|_{\text{F}} \\
&= \|\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\nabla f(\mathbf{X}))\|_{\text{F}} \\
&\stackrel{\textcircled{2}}{=} \|\nabla f(\mathbf{X}) - \frac{1}{2}\mathbf{X}(\mathbf{X}^{\top}\nabla f(\mathbf{X}) + \nabla f(\mathbf{X})^{\top}\mathbf{X})\|_{\text{F}} \\
&= \|(\mathbf{I} - \frac{1}{2}\mathbf{X}\mathbf{X}^{\top})(\nabla f(\mathbf{X}) - \mathbf{X}\nabla f(\mathbf{X})^{\top}\mathbf{X})\|_{\text{F}} \\
&\stackrel{\textcircled{3}}{\leq} \|\nabla f(\mathbf{X}) - \mathbf{X}\nabla f(\mathbf{X})^{\top}\mathbf{X}\|_{\text{F}},
\end{aligned}$$

where step  $\textcircled{1}$  uses Formulation (28); step  $\textcircled{2}$  uses  $\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\Delta) = \Delta - \frac{1}{2}\mathbf{X}(\Delta^{\top}\mathbf{X} + \mathbf{X}^{\top}\Delta)$  for all  $\Delta \in \mathbb{R}^{n \times r}$  (Absil et al., 2008a); step  $\textcircled{3}$  uses the norm inequality  $\|\mathbf{A}\mathbf{B}\|_{\text{F}} \leq \|\mathbf{A}\|_{\text{F}}\|\mathbf{B}\|_{\text{F}}$ , and fact that the matrix  $\mathbf{I} - \frac{1}{2}\mathbf{X}\mathbf{X}^{\top}$  only contains eigenvalues that are  $\frac{1}{2}$  or 1.  $\square$

## C PROOFS FOR SECTION 4

### C.1 PROOF OF LEMMA 4.1

*Proof.* We define  $L(\mathbf{X}, \mathbf{y}; \mathbf{z}; \beta, \mu) \triangleq f(\mathbf{X}) - g(\mathbf{X}) + h_{\mu}(\mathbf{y}) + \langle \mathbf{z}, \mathcal{A}(\mathbf{X}) - \mathbf{y} \rangle + \frac{\beta}{2}\|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|_2^2$ .

We define  $\dot{\sigma} \triangleq (\sigma - 1)/(2 - \sigma)$ , and  $\ddot{\sigma} \triangleq (\sigma/(2 - \sigma))^2$ .

**Part (a-i).** Using the first-order optimality condition of  $\mathbf{y}^{t+1} \in \arg \min_{\mathbf{y}} L(\mathbf{X}^{t+1}, \mathbf{y}, \mathbf{z}^t; \beta^t, \mu^t)$  in Algorithm 1, for all  $t \geq 0$ , we have:

$$\begin{aligned}
\mathbf{0} &= \nabla h_{\mu^t}(\mathbf{y}^{t+1}) + \beta^t(\mathbf{y}^{t+1} - \mathbf{y}^t) + \nabla_{\mathbf{y}}\mathcal{S}(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t) \\
&\stackrel{\textcircled{1}}{=} \nabla h_{\mu^t}(\mathbf{y}^{t+1}) + \beta^t(\mathbf{y}^{t+1} - \mathbf{y}^t) - \mathbf{z}^t + \beta^t(\mathbf{y}^t - \mathcal{A}(\mathbf{X}^{t+1})) \\
&= \nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \mathbf{z}^t + \beta^t(\mathbf{y}^{t+1} - \mathcal{A}(\mathbf{X}^{t+1})) \\
&\stackrel{\textcircled{2}}{=} \nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \mathbf{z}^t + \frac{1}{\sigma}(\mathbf{z}^t - \mathbf{z}^{t+1}), \tag{29}
\end{aligned}$$

where step  $\textcircled{1}$  uses  $\nabla_{\mathbf{y}}\mathcal{S}(\mathbf{X}^{t+1}, \mathbf{y}; \mathbf{z}^t; \beta^t) = -\mathbf{z}^t + \beta^t(\mathbf{y} - \mathcal{A}(\mathbf{X}^{t+1}))$ ; step  $\textcircled{2}$  uses  $\mathbf{z}^{t+1} = \mathbf{z}^t + \sigma\beta^t(\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1})$ .

**Part (a-ii).** We obtain:

$$\begin{aligned}
\partial h(\check{\mathbf{y}}^{t+1}) - \mathbf{z}^t &\stackrel{\textcircled{1}}{\ni} \beta^t(\mathbf{b} - \mathbf{y}^{t+1}) - \mathbf{z}^t \\
&\stackrel{\textcircled{2}}{=} \beta^t\mathbf{y}^t - \nabla_{\mathbf{y}}\mathcal{S}^t(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t) - \beta^t\mathbf{y}^{t+1} - \mathbf{z}^t \\
&\stackrel{\textcircled{3}}{=} \beta^t\mathbf{y}^t - \beta^t(\mathbf{y}^t - \mathcal{A}(\mathbf{X}^{t+1})) - \beta^t\mathbf{y}^{t+1} \\
&= \beta^t(\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}) \\
&\stackrel{\textcircled{4}}{=} \frac{1}{\sigma}(\mathbf{z}^{t+1} - \mathbf{z}^t),
\end{aligned}$$

where step  $\textcircled{1}$  uses the result in Lemma 2.5 that  $\beta^t(\mathbf{b} - \mathbf{y}^{t+1}) \in \partial h(\check{\mathbf{y}}^{t+1})$ ; step  $\textcircled{2}$  uses  $\mathbf{b} \triangleq \mathbf{y}^t - \nabla_{\mathbf{y}}\mathcal{S}^t(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t)/\beta^t$ , as shown in Algorithm 1; step  $\textcircled{3}$  uses  $\nabla_{\mathbf{y}}\mathcal{S}^t(\mathbf{X}^{t+1}, \mathbf{y}; \mathbf{z}^t; \beta^t) = -\mathbf{z}^t + \beta^t(\mathbf{y} - \mathcal{A}(\mathbf{X}^{t+1}))$ ; step  $\textcircled{4}$  uses  $\mathbf{z}^{t+1} - \mathbf{z}^t = \sigma\beta^t(\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1})$ .

**Part (b).** First, we derive:

$$\begin{aligned}
&\|\nabla h_{\mu^{t-1}}(\mathbf{y}^t) - \nabla h_{\mu^t}(\mathbf{y}^{t+1})\| \\
&\stackrel{\textcircled{1}}{\leq} \|\nabla h_{\mu^{t-1}}(\mathbf{y}^t) - \nabla h_{\mu^t}(\mathbf{y}^t)\| + \|\nabla h_{\mu^t}(\mathbf{y}^t) - \nabla h_{\mu^t}(\mathbf{y}^{t+1})\| \\
&\stackrel{\textcircled{2}}{\leq} \|\nabla h_{\mu^t}(\mathbf{y}^t) - \nabla h_{\mu^{t-1}}(\mathbf{y}^t)\| + \frac{1}{\mu^t}\|\mathbf{y}^{t+1} - \mathbf{y}^t\| \\
&\stackrel{\textcircled{3}}{\leq} C_h\left(\frac{\mu^{t-1}}{\mu^t} - 1\right) + \frac{\beta^t}{\lambda}\|\mathbf{y}^{t+1} - \mathbf{y}^t\|, \tag{30}
\end{aligned}$$

where step ① uses  $\|\mathbf{a} - \mathbf{b}\| \leq \|\mathbf{a} - \mathbf{c}\| + \|\mathbf{c} - \mathbf{b}\|$ ; step ② uses the fact that the function  $h_{\mu^t}(\mathbf{y})$  is  $\frac{1}{\mu^t}$ -smooth w.r.t.  $\mathbf{y}$  that:  $\|\nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \nabla h_{\mu^t}(\mathbf{y}^t)\| \leq \frac{1}{\mu^t} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|$ ; step ③ uses the fact that  $\|\nabla h_{\mu^t}(\mathbf{y}^t) - \nabla h_{\mu^{t-1}}(\mathbf{y}^t)\| \leq (\mu^{t-1}/\mu^t - 1)C_h$  which holds due to Lemma 2.4, and the equality  $\mu^t \beta^t = \chi$ .

Second, we have from Equality (29):

$$\begin{aligned} \forall t \geq 0, \mathbf{0} &\in \sigma \nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \sigma \mathbf{z}^t + (\mathbf{z}^t - \mathbf{z}^{t+1}), \\ \forall t \geq 1, \mathbf{0} &\in \sigma \nabla h_{\mu^{t-1}}(\mathbf{y}^t) - \sigma \mathbf{z}^{t-1} + (\mathbf{z}^{t-1} - \mathbf{z}^t). \end{aligned}$$

Combining these two equalities yields:

$$\forall t \geq 1, \mathbf{z}^{t+1} - \mathbf{z}^t = (\sigma - 1)(\mathbf{z}^{t-1} - \mathbf{z}^t) + \sigma(\nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \nabla h_{\mu^{t-1}}(\mathbf{y}^t)).$$

Applying Lemma A.4 with  $\mathbf{a}^+ = \mathbf{z}^{t+1} - \mathbf{z}^t$ ,  $\mathbf{a} = \mathbf{z}^{t-1} - \mathbf{z}^t$ ,  $\mathbf{b} = \sigma\{\nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \nabla h_{\mu^{t-1}}(\mathbf{y}^t)\}$ , and  $\varrho = \sigma - 1 \in [0, 1)$ , we have:

$$\begin{aligned} &\|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \\ &\leq \frac{\varrho}{1-\varrho} (\|\mathbf{z}^{t-1} - \mathbf{z}^t\|_2^2 - \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2) + \frac{1}{(1-\varrho)^2} \|\sigma(\nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \nabla h_{\mu^{t-1}}(\mathbf{y}^t))\|_2^2 \\ &\stackrel{\textcircled{1}}{=} \dot{\sigma} (\|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2 - \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2) + \ddot{\sigma} \|\nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \nabla h_{\mu^{t-1}}(\mathbf{y}^t)\|_2^2 \\ &\stackrel{\textcircled{2}}{\leq} \dot{\sigma} (\|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2 - \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2) + 2\ddot{\sigma} \left\{ \frac{(\beta^t)^2}{\chi^2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + C_h^2 (\mu^{t-1}/\mu^t - 1)^2 \right\} \\ &\stackrel{\textcircled{3}}{\leq} \dot{\sigma} (\|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2 - \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2) + 2\ddot{\sigma} \frac{(\beta^t)^2}{\chi^2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + 2\ddot{\sigma} C_h^2 \left( \frac{2}{t} - \frac{2}{t+1} \right), \end{aligned}$$

where step ① uses the definitions of  $\{\dot{\sigma}, \ddot{\sigma}\}$ ; step ② uses Inequality (30), and the inequality  $(a+b)^2 \leq 2a^2 + 2b^2$  for all  $a, b \in \mathbb{R}$ ; step ③ uses Lemma A.3 that  $(\frac{\beta^t}{\beta^{t-1}} - 1)^2 \leq \frac{2}{t} - \frac{2}{t+1}$  for all  $t \geq 1$ ;

□

## C.2 PROOF OF LEMMA 4.3

*Proof.* **Part (a).** We have:

$$\beta^{t+1} - \beta^t \cdot (1 + \xi) \stackrel{\textcircled{1}}{=} \beta^0 \xi (t+1)^p - \beta^0 \xi t^p - \beta^t \xi \stackrel{\textcircled{2}}{\leq} \beta^0 \xi - \beta^t \xi \stackrel{\textcircled{3}}{\leq} 0,$$

where step ① uses  $\beta^t = \beta^0(1 + \xi^t)$ ; step ② uses  $(t+1)^p - t^p \leq 1$  for all  $p \in (0, 1)$ ; step ③ uses  $\beta^0 \leq \beta^t$  and  $\xi > 0$ .

**Part (b).** It holds with  $\underline{\ell} = \bar{A}^2$  and  $\bar{\ell} = \bar{A}^2 + L_f/\beta^0$ .

□

## C.3 PROOF OF LEMMA 4.4

*Proof.* We define  $\bar{X} \triangleq \sqrt{r}$ ,  $\bar{z} \triangleq \|\mathbf{z}^0\| + \frac{\sigma C_h}{2-\sigma}$ ,  $\bar{y} \triangleq \bar{A}\sqrt{r} + \frac{2\bar{z}}{\beta^0}$ , where  $\sigma \in [1, 2)$ .

We let  $\underline{\Theta} \triangleq F(\bar{\mathbf{X}}) - \mu^0 C_h^2 - C_h(\bar{A}\sqrt{r} + \bar{y}) - \frac{\bar{z}^2}{2\beta^0}$ , where  $\bar{\mathbf{X}}$  is the optimal solution of Problem (1).

**Part (a).** Given  $\mathbf{X}^{t+1} \in \mathcal{M}$ , we have:  $\|\mathbf{X}^t\|_F \leq \bar{X} \triangleq \sqrt{r}$ .

**Part (b).** We show that  $\|\mathbf{z}^t\| \leq \bar{z}$ . For all  $t \geq 0$ , we have:

$$\begin{aligned} \|\mathbf{z}^{t+1}\| &\stackrel{\textcircled{1}}{\leq} \|(\sigma - 1)\mathbf{z}^t\| + \|(\sigma - 1)\mathbf{z}^t + \mathbf{z}^{t+1}\| \\ &\stackrel{\textcircled{2}}{=} (\sigma - 1)\|\mathbf{z}^t\| + \|\sigma \partial h(\check{\mathbf{y}}^{t+1})\| \\ &\stackrel{\textcircled{3}}{=} (\sigma - 1)\|\mathbf{z}^t\| + \sigma C_h, \end{aligned}$$

step ① uses the triangle inequality; step ② uses  $\mathbf{z}^{t+1} + (\sigma - 1)\mathbf{z}^t \in \sigma \partial h(\check{\mathbf{y}}^{t+1})$ , as shown in Lemma 4.1(a); step ③ uses  $C_h$ -Lipschitz continuity of  $h(\mathbf{y})$ . Applying Lemma A.5 with  $\mathbf{a}_t = \|\mathbf{z}^{t+1}\|$ ,  $c = \sigma C_h$ , and  $\varrho = \sigma - 1 \in [0, 1)$ , we have:

$$\forall t \geq 0, \|\mathbf{z}^{t+1}\| \leq \|\mathbf{z}^0\| + \frac{c}{1-\varrho} = \|\mathbf{z}^0\| + \frac{\sigma C_h}{2-\sigma} \triangleq \bar{z}.$$

1458 **Part (c).** We show that  $\|\mathbf{y}^t\| \leq \bar{y}$ . For all  $t \geq 0$ , we have:

$$\begin{aligned}
1460 \quad \|\mathbf{y}^{t+1}\| &= \|\mathcal{A}(\mathbf{X}^{t+1}) - \frac{\mathbf{z}^{t+1} - \mathbf{z}^t}{\sigma\beta^t}\| \\
1461 &\stackrel{\textcircled{1}}{\leq} \|\mathcal{A}(\mathbf{X}^{t+1})\| + \frac{1}{\beta^0} \|\mathbf{z}^{t+1} - \mathbf{z}^t\| \\
1462 &\stackrel{\textcircled{2}}{\leq} \bar{A}\sqrt{r} + \frac{1}{\beta^0} \cdot 2\bar{z} \triangleq \bar{y},
\end{aligned}$$

1466 where step ① uses the triangle inequality,  $\sigma \geq 1$ , and  $\frac{1}{\beta^t} \leq \frac{1}{\beta^0}$ ; step ② uses  $\|\mathcal{A}(\mathbf{X})\|_F \leq \bar{A}\|\mathbf{X}\|_F \leq \bar{A}\sqrt{r}$ , and  $\|\mathbf{z}^t\| \leq \bar{z}$ .

1469 **Part (d).** We show that  $\Theta^t \geq \underline{\Theta}$ . For all  $t \geq 1$ , we have:

$$\begin{aligned}
1471 \quad \Theta^t &\triangleq L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \mu^{t-1}) + \mu^{t-1}C_h^2 + \mathbb{T}^t + \mathbb{Z}^t + \mathbb{X}^t \\
1472 &\stackrel{\textcircled{1}}{\geq} f(\mathbf{X}^t) - g(\mathbf{X}^t) + h_{\mu^{t-1}}(\mathbf{y}^t) + \langle \mathbf{z}^t, \mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t \rangle + \frac{\beta^t}{2} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|_2^2 \\
1473 &= f(\mathbf{X}^t) - g(\mathbf{X}^t) + h_{\mu^{t-1}}(\mathbf{y}^t) + \frac{\beta^t}{2} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t + \mathbf{z}^t/\beta^t\|_2^2 - \frac{\beta^t}{2} \|\mathbf{z}^t/\beta^t\|_2^2 \\
1474 &\stackrel{\textcircled{2}}{\geq} f(\mathbf{X}^t) - g(\mathbf{X}^t) + h_{\mu^{t-1}}(\mathcal{A}(\mathbf{X}^t)) - C_h \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| - \frac{1}{2\beta^t} \|\mathbf{z}^t\|_2^2 \\
1475 &\stackrel{\textcircled{3}}{\geq} f(\mathbf{X}^t) - g(\mathbf{X}^t) + h(\mathcal{A}(\mathbf{X}^t)) - \mu^{t-1}C_h^2 - C_h(\|\mathcal{A}(\mathbf{X}^t)\| + \|\mathbf{y}^t\|) - \frac{1}{2\beta^t} \|\mathbf{z}^t\|_2^2 \\
1476 &\stackrel{\textcircled{4}}{\geq} F(\bar{\mathbf{X}}) - \mu^0 C_h^2 - C_h(\bar{A}\sqrt{r} + \bar{y}) - \frac{\bar{z}^2}{2\beta^0} \triangleq \underline{\Theta},
\end{aligned}$$

1482 where step ① uses the definition of  $L(\mathbf{X}, \mathbf{y}; \mathbf{z}; \beta; \mu)$  and the positivity of  $\{\mu^t, \mathbb{T}^t, \mathbb{Z}^t, \mathbb{X}^t\}$ ; step ② uses the  $L_h$ -Lipschitz continuity of  $h_{\mu^{t-1}}(\mathbf{y})$ , ensuring  $h_{\mu^{t-1}}(\mathbf{y}^t) \geq h_{\mu^{t-1}}(\mathbf{y}) - C_h\|\mathbf{y}^t - \mathbf{y}\|$ , with the specific choice of  $\mathbf{y} = \mathcal{A}(\mathbf{X}^t)$ ; step ③ uses  $h(\mathbf{y}) - h_{\mu}(\mathbf{y}) \leq \mu C_h^2$ , which has been shown in Lemma 2.3; step ④ uses  $\mu^t \leq \mu^0$ ,  $\beta^t \geq \beta^0$ ,  $\|\mathcal{A}(\mathbf{X})\| \leq \bar{A}\|\mathbf{X}\|_F \leq \bar{A}\sqrt{r}$  for all  $\mathbf{X} \in \mathcal{M}$ ;  $\|\mathbf{y}^t\| \leq \bar{y}$ , and  $\|\mathbf{z}^t\| \leq \bar{z}$ .

□

#### 1490 C.4 PROOF OF LEMMA 4.5

1492 *Proof.* We define  $L(\mathbf{X}, \mathbf{y}; \mathbf{z}; \beta, \mu) \triangleq f(\mathbf{X}) - g(\mathbf{X}) + h_{\mu}(\mathbf{y}) + \langle \mathbf{z}, \mathcal{A}(\mathbf{X}) - \mathbf{y} \rangle + \frac{\beta}{2} \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|_2^2$ .

1494 We define  $\omega \triangleq \frac{1}{\sigma} + \frac{\xi}{2\sigma^2} + \frac{\varepsilon_z}{\sigma^2}$ .

1496 We define  $\mathbb{Z}^t \triangleq \omega \dot{\sigma}^2 \beta^{t-1} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|_2^2 = \frac{\omega \dot{\sigma}}{\beta^{t-1}} \|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2$ , where we use  $\mathbf{z}^{t+1} - \mathbf{z}^t = \beta^t \sigma (\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1})$ .

1498 **Part (a).** We focus on the sufficient decrease for variables  $\{\mu, \mathbf{y}\}$ . First, we have:

$$\begin{aligned}
1500 \quad \Xi &\triangleq \langle \mathbf{y}^t - \mathbf{y}^{t+1}, \mathbf{z}^t \rangle + \frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathcal{A}(\mathbf{X}^{t+1})\|_2^2 - \frac{\beta^t}{2} \|\mathbf{y}^t - \mathcal{A}(\mathbf{X}^{t+1})\|_2^2 \\
1501 &\stackrel{\textcircled{1}}{=} \langle \mathbf{y}^t - \mathbf{y}^{t+1}, \mathbf{z}^t + \beta^t (\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}) \rangle - \frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 \\
1502 &\stackrel{\textcircled{2}}{=} -\frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \langle \mathbf{y}^t - \mathbf{y}^{t+1}, \mathbf{z}^t + \frac{1}{\sigma} (\mathbf{z}^{t+1} - \mathbf{z}^t) \rangle \\
1503 &\stackrel{\textcircled{3}}{=} -\frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \langle \mathbf{y}^t - \mathbf{y}^{t+1}, \nabla h_{\mu^t}(\mathbf{y}^{t+1}) \rangle \\
1504 &\stackrel{\textcircled{4}}{\leq} \left\{ \frac{1}{\chi} - \beta^t \right\} \frac{1}{2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + h_{\mu^t}(\mathbf{y}^t) - h_{\mu^t}(\mathbf{y}^{t+1}),
\end{aligned} \tag{31}$$

1509 where step ① uses the Pythagoras Relation that  $\frac{1}{2} \|\mathbf{y}^+ - \mathbf{a}\|_2^2 - \frac{1}{2} \|\mathbf{y} - \mathbf{a}\|_2^2 = -\frac{1}{2} \|\mathbf{y}^+ - \mathbf{y}\|_2^2 + \langle \mathbf{y} - \mathbf{y}^+, \mathbf{a} - \mathbf{y}^+ \rangle$  for all  $\mathbf{y}, \mathbf{y}^+, \mathbf{a} \in \mathbb{R}^m$ ; step ② uses  $\mathbf{z}^{t+1} - \mathbf{z}^t = \sigma \beta^t (\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1})$ ; step ③ uses  $\nabla h_{\mu^t}(\mathbf{y}^{t+1}) = \mathbf{z}^t + \frac{1}{\sigma} (\mathbf{z}^{t+1} - \mathbf{z}^t)$ , as shown in Lemma 4.1(a); step ④ uses the fact that the function

1512  $h_{\mu^t}(\mathbf{y})$  is  $(1/\mu^t)$ -weakly convex w.r.t  $\mathbf{y}$ , and  $\mu^t \beta^t = \chi$ . Furthermore, we have:

$$\begin{aligned}
1513 & L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}; \mathbf{z}^t; \beta^t, \mu^t) - L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \mu^{t-1}) \\
1514 & \stackrel{\textcircled{1}}{=} h_{\mu^t}(\mathbf{y}^{t+1}) - h_{\mu^{t-1}}(\mathbf{y}^t) + \underbrace{\langle \mathbf{y}^t - \mathbf{y}^{t+1}, \mathbf{z}^t \rangle + \frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathcal{A}(\mathbf{X}^{t+1})\|_2^2 - \frac{\beta^t}{2} \|\mathbf{y}^t - \mathcal{A}(\mathbf{X}^{t+1})\|_2^2}_{=\Xi} \\
1515 & \stackrel{\textcircled{2}}{\leq} \frac{1/\chi - 1}{2} \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + h_{\mu^t}(\mathbf{y}^t) - h_{\mu^{t-1}}(\mathbf{y}^t) \\
1516 & \stackrel{\textcircled{3}}{=} \frac{1/\chi - 1}{2} \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + (\mu^{t-1} - \mu^t) C_h^2, \tag{32}
\end{aligned}$$

1521 where step ① uses the definition of  $L(\mathbf{X}, \mathbf{y}; \mathbf{z}; \beta, \mu)$ ; step ② uses Inequality (31); step ③ uses Lemma  
1522 2.3 that  $h_{\mu^t}(\mathbf{y}) - h_{\mu^{t-1}}(\mathbf{y}) \leq \min\{\frac{\mu^{t-1}}{2\mu^t}, 1\} \cdot (\mu^{t-1} - \mu^t) C_h^2 \leq (\mu^{t-1} - \mu^t) C_h^2$  for all  $\mathbf{y}$ .

1523 **Part (b).** We focus on the sufficient decrease for variables  $\{\mathbf{z}, \beta\}$ . We have:

$$\begin{aligned}
1524 & L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}; \mathbf{z}^{t+1}; \beta^{t+1}, \mu^t) - L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}; \mathbf{z}^t; \beta^t, \mu^t) + \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 \\
1525 & \stackrel{\textcircled{1}}{=} \langle \mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}, \mathbf{z}^{t+1} - \mathbf{z}^t \rangle + \frac{\beta^{t+1} - \beta^t}{2} \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 \\
1526 & \stackrel{\textcircled{2}}{=} \left\{ \frac{1}{\sigma} + \frac{\beta^{t+1} - \beta^t}{2\sigma^2 \beta^t} + \frac{\varepsilon_z}{\sigma^2} \right\} \cdot \frac{1}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \\
1527 & \stackrel{\textcircled{3}}{=} \underbrace{\left\{ \frac{1}{\sigma} + \frac{\xi}{2\sigma^2} + \frac{\varepsilon_z}{\sigma^2} \right\}}_{\triangleq \omega} \cdot \frac{1}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \\
1528 & \stackrel{\textcircled{4}}{\leq} \frac{\omega \check{\sigma}}{\beta^t} (\|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2 - \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2) + \frac{2\omega \check{\sigma}}{\chi^2} \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \frac{2\omega \check{\sigma}}{\beta^t} C_h^2 \left( \frac{2}{t} - \frac{2}{t+1} \right) \\
1529 & \stackrel{\textcircled{5}}{\leq} \underbrace{\frac{\omega \check{\sigma}}{\beta^{t-1}} \|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2 - \frac{\omega \check{\sigma}}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2}_{\triangleq \mathbb{Z}^t} + \frac{2\omega \check{\sigma}}{\chi^2} \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \underbrace{\frac{2\omega \check{\sigma}}{\beta^0} C_h^2 \left( \frac{2}{t} - \frac{2}{t+1} \right)}_{=\mathbb{T}^t - \mathbb{T}^{t+1}}, \tag{33}
\end{aligned}$$

1532 where step ① uses the definition of  $L(\mathbf{X}, \mathbf{y}; \mathbf{z}; \beta; \mu)$ ; step ② uses  $\mathbf{z}^{t+1} - \mathbf{z}^t = \sigma \beta^t (\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1})$ ;  
1533 step ③ uses  $\beta^{t+1} \leq (1 + \xi) \beta^t$ ; step ④ uses the upper bound for  $\|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2$  as shown in Lemma  
1534 4.1(b); step ⑤ uses  $\beta^t \geq \beta^{t-1} \geq \beta^0$ .

1535 Adding Inequalities (32) and (33) together, we have:

$$\begin{aligned}
1536 & L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}; \mathbf{z}^{t+1}; \beta^{t+1}, \mu^t) - L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \mu^{t-1}) + (\mu^t - \mu^{t-1}) C_h^2 \\
1537 & + \mathbb{T}^{t+1} - \mathbb{T}^t + \mathbb{Z}^{t+1} - \mathbb{Z}^t + \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 \\
1538 & \leq \frac{1}{2} \left\{ \frac{1}{\chi} - 1 + \frac{4\omega \check{\sigma}}{\chi^2} \right\} \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 \\
1539 & \stackrel{\textcircled{1}}{\leq} \frac{1}{2} \underbrace{\left\{ -1 + \frac{1+4\omega \check{\sigma}}{\chi} \right\}}_{\triangleq -\varepsilon_y} \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2,
\end{aligned}$$

1540 where step ① uses  $\chi \geq 1$ .

□

## 1555 C.5 PROOF OF LEMMA 4.6

1556 *Proof.* We define  $\mathcal{S}(\mathbf{X}, \mathbf{y}^t; \mathbf{z}^t; \beta^t) \triangleq f(\mathbf{X}) + \langle \mathbf{z}^t, \mathcal{A}(\mathbf{X}) - \mathbf{y}^t \rangle + \frac{\beta^t}{2} \|\mathcal{A}(\mathbf{X}) - \mathbf{y}^t\|_2^2$ .

1557 We let  $\mathbf{G}^t \in \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}_c^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) - \partial g(\mathbf{X}^t)$ .

1558 We define  $\mathbb{X}^t \triangleq \frac{1}{2} (\alpha + \theta \alpha) \ell(\beta^t) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F^2$ .

1559 We define  $\varepsilon'_x \triangleq (\theta - 1 - \alpha - \theta \alpha) - (1 + \xi)(\alpha + \theta \alpha) > 0$ , and  $\varepsilon_x \triangleq \frac{1}{2} \varepsilon'_x \ell > 0$ .

1560 First, using the optimality condition of  $\mathbf{X}^{t+1} \in \mathcal{M}$ , we have:

$$1561 \langle \mathbf{X}^{t+1} - \mathbf{X}^t, \mathbf{G}^t \rangle + \frac{\theta \ell(\beta^t)}{2} \|\mathbf{X}^{t+1} - \mathbf{X}_c^t\|_F^2 \leq \langle \mathbf{X}^t - \mathbf{X}^t, \mathbf{G}^t \rangle + \frac{\theta \ell(\beta^t)}{2} \|\mathbf{X}^t - \mathbf{X}_c^t\|_F^2. \tag{34}$$

1566 Second, we have:

$$\begin{aligned}
1567 & L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t; \mu^t, \beta^t) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t; \mu^t, \beta^t) \\
1568 & = \mathcal{S}(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t) - \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) + g(\mathbf{X}^t) - g(\mathbf{X}^{t+1}) \\
1569 & \stackrel{\textcircled{1}}{\leq} \frac{\ell(\beta^t)}{2} \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 + \langle \mathbf{X}^{t+1} - \mathbf{X}^t, \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) \rangle + \langle \mathbf{X}^t - \mathbf{X}^{t+1}, \partial g(\mathbf{X}^t) \rangle, \quad (35)
\end{aligned}$$

1572 where step ① uses the  $\ell(\beta^t)$ -smoothness of  $\mathcal{S}(\mathbf{X}, \mathbf{y}^t; \mathbf{z}^t; \beta^t)$  and convexity of  $g(\mathbf{X})$ ;

1574 Third, we derive:

$$\begin{aligned}
1575 & \langle \mathbf{X}^{t+1} - \mathbf{X}^t, \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) - \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}_c^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) \rangle \\
1576 & \stackrel{\textcircled{1}}{\leq} \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}} \cdot \|\nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) - \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}_c^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t)\|_{\mathbb{F}} \\
1577 & \stackrel{\textcircled{2}}{\leq} \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}} \cdot \ell(\beta^t) \|\mathbf{X}^t - \mathbf{X}_c^t\|_{\mathbb{F}} \\
1578 & \stackrel{\textcircled{3}}{\leq} \alpha \ell(\beta^t) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}} \cdot \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}} \\
1579 & \stackrel{\textcircled{4}}{\leq} \frac{\alpha \ell(\beta^t)}{2} \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 + \frac{\alpha \ell(\beta^t)}{2} \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}}^2, \quad (36)
\end{aligned}$$

1584 where step ① uses the norm inequality; step ② uses the  $\ell(\beta^t)$ -smoothness of  $\mathcal{S}(\mathbf{X}, \mathbf{y}^t; \mathbf{z}^t; \beta^t)$ ; step  
1585 ③ uses  $\mathbf{X}_c^t = \mathbf{X}^t + \alpha(\mathbf{X}^t - \mathbf{X}^{t-1})$ ; step ④ uses  $ab \leq \frac{1}{2}a^2 + \frac{1}{2}b^2$  for all  $a \in \mathbb{R}$  and  $b \in \mathbb{R}$ .

1586 Summing Inequalities (34),(36), and (35), we obtain:

$$\begin{aligned}
1587 & L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t; \mu^t, \beta^t) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t; \mu^t, \beta^t) \\
1588 & \leq \frac{\ell(\beta^t)}{2} \{(1 + \alpha) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 + \alpha \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}} + \theta \|\mathbf{X}^t - \mathbf{X}_c^t\|_{\mathbb{F}}^2 - \theta \|\mathbf{X}^{t+1} - \mathbf{X}_c^t\|_{\mathbb{F}}^2\} \\
1589 & \stackrel{\textcircled{1}}{=} \frac{\ell(\beta^t)}{2} \{(1 + \alpha) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 + (\alpha + \theta \alpha^2) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}}^2 - \theta \|\mathbf{X}^{t+1} - \mathbf{X}^t - \alpha(\mathbf{X}^t - \mathbf{X}^{t-1})\|_{\mathbb{F}}^2\} \\
1590 & \stackrel{\textcircled{2}}{\leq} \frac{\ell(\beta^t)}{2} \{(1 + \alpha) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 + (\alpha + \theta \alpha^2) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}}^2 \\
1591 & + \theta(\alpha - 1) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 - \theta \alpha(\alpha - 1) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}}^2\} \\
1592 & = \underbrace{\frac{1}{2}(\alpha + \theta \alpha) \ell(\beta^t) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}}^2}_{\triangleq \mathbb{X}^t} + \frac{\ell(\beta^t)}{2} \cdot \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 \cdot \{1 + \alpha + \theta \alpha - \theta\} \\
1593 & = \mathbb{X}^t - \mathbb{X}^{t+1} + \frac{1}{2} \cdot \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 \cdot \{\ell(\beta^t)(1 + \alpha + \theta \alpha - \theta) + \ell(\beta^{t+1})(\alpha + \theta \alpha)\} \\
1594 & \stackrel{\textcircled{3}}{\leq} \mathbb{X}^t - \mathbb{X}^{t+1} + \frac{1}{2} \cdot \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 \cdot \underbrace{\ell(\beta^t) \{(1 + \alpha + \theta \alpha - \theta) + (1 + \xi)(\alpha + \theta \alpha)\}}_{\triangleq -\varepsilon'_x} \\
1595 & \stackrel{\textcircled{4}}{\leq} \mathbb{X}^t - \mathbb{X}^{t+1} - \frac{1}{2} \cdot \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 \cdot \varepsilon'_x \cdot \beta^t \underline{\ell} \\
1596 & \stackrel{\textcircled{5}}{=} \mathbb{X}^t - \mathbb{X}^{t+1} - \varepsilon_x \cdot \beta^t \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2,
\end{aligned}$$

1597 where step ① uses  $\mathbf{X}_c^t = \mathbf{X}^t + \alpha(\mathbf{X}^t - \mathbf{X}^{t-1})$ ; step ② uses Lemma A.1 with  $\mathbf{a} = \mathbf{X}^{t+1} - \mathbf{X}^t$ ,  
1598 and  $\mathbf{b} = \mathbf{X}^t - \mathbf{X}^{t-1}$ ; step ③ uses the fact that  $\ell(\beta^{t+1}) \leq (1 + \xi)\ell(\beta^t)$ , which is implied by  
1599  $\beta^{t+1} \leq (1 + \xi)\beta^t$ ; step ④ uses Lemma 4.3 that  $\beta^t \underline{\ell} \leq \ell(\beta^t) \leq \beta^t \bar{\ell}$ ; step ⑤ uses  $\varepsilon_x \triangleq \frac{1}{2}\varepsilon'_x \underline{\ell} > 0$ .

1600

1601

## 1602 C.6 PROOF OF LEMMA 4.7

1603 *Proof.* We define:  $\Theta^t \triangleq L(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^{t-1}) + \mu^{t-1} C_h^2 + \mathbb{T}^t + \mathbb{Z}^t + \mathbb{X}^t$ ,

1604 We define  $\tilde{\varepsilon}_t \triangleq \|\mathbf{y}^t - \mathbf{y}^{t-1}\|^2 + \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|^2 + \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}}^2$ .

1605 **Part (a).** Using Lemma 4.5, we have:

$$\begin{aligned}
1606 & L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}; \mathbf{z}^{t+1}; \beta^{t+1}, \mu^t) - L(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^{t-1}) - (\mu^{t-1} - \mu^t) C_h^2 \\
1607 & \leq \mathbb{T}^t - \mathbb{T}^{t+1} + \mathbb{Z}^t - \mathbb{Z}^{t+1} - \varepsilon_y \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|^2 - \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2. \quad (37)
\end{aligned}$$

1608

1609

Using Lemma 4.6, we have:

$$L(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^{t-1}) - L(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^{t-1}) \leq \mathbb{X}^t - \mathbb{X}^{t+1} - \varepsilon_x \beta^t \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2.$$

Adding these two inequalities together and using the definition of  $\Theta^t$ , we have:

$$\begin{aligned} \Theta^t - \Theta^{t+1} &\geq \varepsilon_y \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \varepsilon_x \beta^t \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 + \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 \\ &\geq \min(\varepsilon_y, \varepsilon_x, \varepsilon_z) \cdot \beta^t \cdot \tilde{e}_{t+1}. \end{aligned}$$

**Part (b).** Telescoping this inequality over  $t$  from 1 to  $T$ , we have:

$$\begin{aligned} \sum_{t=1}^T \beta^t \tilde{e}_{t+1} &\leq \frac{1}{\min(\varepsilon_y, \varepsilon_x, \varepsilon_z)} \cdot \sum_{t=1}^T (\Theta^t - \Theta^{t+1}) \\ &= \frac{1}{\min(\varepsilon_y, \varepsilon_x, \varepsilon_z)} \cdot (\Theta^1 - \Theta^{T+1}) \\ &\stackrel{\textcircled{1}}{\leq} \frac{1}{\min(\varepsilon_y, \varepsilon_x, \varepsilon_z)} \cdot (\Theta^1 - \underline{\Theta}), \end{aligned} \quad (38)$$

where step  $\textcircled{1}$  uses  $\Theta^t \geq \underline{\Theta}$ . Furthermore, we have:

$$\sum_{t=1}^T \beta^t \tilde{e}_{t+1} = \sum_{t=1}^T \frac{1}{\beta^t} (\beta^t)^2 \tilde{e}_{t+1} \geq \frac{1}{\beta^T} \sum_{t=1}^T (\beta^t)^2 \tilde{e}_{t+1} \stackrel{\textcircled{1}}{\geq} \frac{1}{3T\beta^T} (\sum_{t=1}^T \beta^t e^{t+1})^2, \quad (39)$$

where step  $\textcircled{1}$  uses  $\sum_{i=1}^n \mathbf{x}_i^2 \geq \frac{1}{n} (\sum_{i=1}^n |\mathbf{x}_i|)^2$  for all  $\mathbf{x} \in \mathbb{R}^n$ . Combining Inequalities (38) and (39), we have:  $\sum_{t=1}^T \beta^t e^{t+1} \leq \left\{ \frac{\Theta^1 - \underline{\Theta}}{\min(\varepsilon_y, \varepsilon_x, \varepsilon_z)} \cdot 3T\beta^T \right\}^{1/2} = \mathcal{O}(T^{(1+p)/2})$ .

□

## C.7 PROOF OF THEOREM 4.8

*Proof.* We define  $\text{Crit}(\mathbf{X}, \mathbf{y}, \mathbf{z}) \triangleq \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\| + \|\partial h(\mathbf{y}) - \mathbf{z}\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\nabla f(\mathbf{X}) - \partial g(\mathbf{X}) + \mathcal{A}^\top(\mathbf{z}))\|_F$ .

We define  $\dot{\mathbf{G}} \triangleq \nabla f(\mathbf{X}^t) - \partial g(\mathbf{X}^t) + \mathcal{A}^\top(\mathbf{z}^t)$ .

We define  $\ddot{\mathbf{G}} \triangleq \nabla f(\mathbf{X}_c^t) - \partial g(\mathbf{X}^t) + \mathcal{A}^\top(\mathbf{z}^t + \beta^t \mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t) + \theta \ell(\beta^t)(\mathbf{X}^{t+1} - \mathbf{X}_c^t)$ .

We first derive the following inequalities:

$$\begin{aligned} &\|\ddot{\mathbf{G}} - \dot{\mathbf{G}}\|_F \\ &\stackrel{\textcircled{1}}{=} \|\nabla f(\mathbf{X}^t) - \nabla f(\mathbf{X}_c^t) - \beta^t \mathcal{A}^\top(\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t) - \theta \ell(\beta^t)(\mathbf{X}^{t+1} - \mathbf{X}_c^t)\|_F \\ &\stackrel{\textcircled{2}}{\leq} L_f \|\mathbf{X}^t - \mathbf{X}_c^t\|_F + \beta^t \bar{\mathbf{A}} \|\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t\| + \theta \ell(\beta^t) \|\mathbf{X}^{t+1} - \mathbf{X}_c^t\|_F \\ &\stackrel{\textcircled{3}}{\leq} L_f \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F + \beta^t \bar{\mathbf{A}} \{\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \bar{\mathbf{A}} \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F\} \\ &\quad + \theta \ell(\beta^t) (\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F + \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F) \\ &\stackrel{\textcircled{4}}{\leq} (L_f + \beta^t \bar{\mathbf{A}}^2 + \theta \ell(\beta^t)) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F + \beta^t \bar{\mathbf{A}} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \theta \ell(\beta^t) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F \\ &\stackrel{\textcircled{5}}{=} \mathcal{O}(\beta^{t-1} e^t) + \mathcal{O}(\beta^t e^{t+1}), \end{aligned} \quad (40)$$

where step  $\textcircled{1}$  uses the definitions of  $\{\ddot{\mathbf{G}}, \dot{\mathbf{G}}\}$ ; step  $\textcircled{2}$  uses the triangle inequality; step  $\textcircled{3}$  uses the fact that  $f(\mathbf{X})$  is  $L_f$ -smooth,  $\|\mathbf{X}^t - \mathbf{X}_c^t\|_F \leq \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$ ,  $\|\mathbf{X}^{t+1} - \mathbf{X}_c^t\|_F \leq \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F + \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$ , and  $\|\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t\| \leq \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \bar{\mathbf{A}} \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$ , as shown in Lemma A.6.

We derive the following inequalities:

$$\begin{aligned} &\|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\dot{\mathbf{G}})\|_F \\ &\stackrel{\textcircled{1}}{=} \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\dot{\mathbf{G}}) + \text{Proj}_{\mathbf{T}_{\mathbf{X}^{t+1}}\mathcal{M}}(\ddot{\mathbf{G}})\|_F \\ &\stackrel{\textcircled{2}}{\leq} 2\|\dot{\mathbf{G}} - \ddot{\mathbf{G}}\|_F + 2\sqrt{r}\|\dot{\mathbf{G}}\| \cdot \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F \\ &\stackrel{\textcircled{3}}{\leq} \mathcal{O}(\beta^{t-1} e^t) + \mathcal{O}(\beta^t e^{t+1}) + 2\sqrt{r}(C_f + C_g + \bar{\mathbf{A}}\bar{z})\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F \\ &= \mathcal{O}(\beta^{t-1} e^t) + \mathcal{O}(\beta^t e^{t+1}), \end{aligned}$$

where step ① uses the optimality of  $\mathbf{X}^{t+1}$  that:

$$\mathbf{0} = \text{Proj}_{\mathbf{T}_{\mathbf{x}^{t+1}}\mathcal{M}}(\ddot{\mathbf{G}});$$

step ② uses the result of Lemma A.7 by applying

$$\mathbf{X} = \mathbf{X}^t, \tilde{\mathbf{X}} = \mathbf{X}^{t+1}, \mathbf{P} = \dot{\mathbf{G}}, \text{ and } \tilde{\mathbf{P}} = \ddot{\mathbf{G}};$$

step ③ uses Inequality (40), and the fact that  $\|\dot{\mathbf{G}}\| = \|\nabla f(\mathbf{X}^t) - \partial g(\mathbf{X}^t) + \mathcal{A}^\top(\mathbf{z}^t)\| \leq \|\nabla f(\mathbf{X}^t) - \partial g(\mathbf{X}^t) + \mathcal{A}^\top(\mathbf{z}^t)\|_F \leq C_f + C_g + \bar{\mathbf{A}}\bar{\mathbf{z}}$ .

Finally, we derive:

$$\begin{aligned} & \frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{X}^t, \check{\mathbf{y}}^t, \mathbf{z}^t) \\ \stackrel{\textcircled{1}}{=} & \frac{1}{T} \sum_{t=1}^T \{\|\mathcal{A}(\mathbf{X}^t) - \check{\mathbf{y}}^t\| + \|\partial h(\check{\mathbf{y}}^t) - \mathbf{z}^t\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{x}^t}\mathcal{M}}(\dot{\mathbf{G}})\|_F\} \\ \stackrel{\textcircled{2}}{\leq} & \frac{1}{T} \sum_{t=1}^T \{\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \|(1 - \frac{1}{\sigma})(\mathbf{z}^t - \mathbf{z}^{t-1})\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{x}^t}\mathcal{M}}(\dot{\mathbf{G}})\|_F\} \\ \stackrel{\textcircled{3}}{=} & \frac{1}{T} \sum_{t=1}^T \{\mathcal{O}(\beta^{t-1}e^t) + \mathcal{O}(\beta^t e^{t+1})\} \\ \stackrel{\textcircled{4}}{=} & \mathcal{O}(T^{(p-1)/2}) = \mathcal{O}(T^{-1/3}), \end{aligned}$$

where step ① uses the definition of  $\text{Crit}(\mathbf{X}, \mathbf{y}, \mathbf{z})$ ; step ② uses  $\mathbf{z}^{t+1} - \partial h(\check{\mathbf{y}}^{t+1}) \ni (1 - \frac{1}{\sigma})(\mathbf{z}^{t+1} - \mathbf{z}^t)$ , as shown in Lemma 4.1; step ③ uses  $\|\mathbf{z}^t - \mathbf{z}^{t-1}\| = \|\sigma\beta^{t-1}(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)\| \leq 2\beta^t\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| = \mathcal{O}(\beta^{t-1}e^t)$ ; step ④ uses the choice  $p = 1/3$  and Lemma 4.7(b). □

## C.8 PROOF OF LEMMA 4.10

*Proof.* We define  $\mathcal{S}(\mathbf{X}, \mathbf{y}^t; \mathbf{z}^t; \beta^t) \triangleq f(\mathbf{X}) + \langle \mathbf{z}^t, \mathcal{A}(\mathbf{X}) - \mathbf{y}^t \rangle + \frac{\beta^t}{2} \|\mathcal{A}(\mathbf{X}) - \mathbf{y}^t\|_2^2$ .

We let  $\mathbf{G}^t \in \nabla_{\mathbf{X}}\mathcal{S}(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) - \partial g(\mathbf{X}^t)$ .

We define  $\eta^t \triangleq \frac{\beta^t \gamma^j}{\beta^t} \in (0, \infty)$ .

**Part (a).** Initially, we show that  $\|\mathbf{G}^t\|_F$  is always bounded for  $t$  with  $\mathbf{X} \in \mathcal{M}$ . We have:

$$\begin{aligned} \|\mathbf{G}^t\|_F &= \|\nabla f(\mathbf{X}^t) - \partial g(\mathbf{X}^t) + \mathcal{A}^\top[\mathbf{z}^t + \beta^t(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)]\|_F \\ &\stackrel{\textcircled{1}}{=} \|\nabla f(\mathbf{X}^t) - \partial g(\mathbf{X}^t) + \mathcal{A}^\top[\mathbf{z}^t + \frac{\beta^t}{\sigma\beta^{t-1}}(\mathbf{z}^t - \mathbf{z}^{t-1})]\|_F \\ &\stackrel{\textcircled{2}}{\leq} \|\nabla f(\mathbf{X}^t)\|_F + \|\partial g(\mathbf{X}^t)\|_F + \bar{\mathbf{A}} \cdot \{\|\mathbf{z}^t\| + \frac{\beta^t}{\sigma\beta^{t-1}}(\|\mathbf{z}^t\| + \|\mathbf{z}^{t-1}\|)\} \\ &\stackrel{\textcircled{3}}{\leq} C_f + C_g + \bar{\mathbf{A}} \cdot (\bar{\mathbf{z}} + 2(1 + \xi)\bar{\mathbf{z}}) \triangleq \bar{g}, \end{aligned}$$

where step ① uses  $\mathbf{z}^{t+1} = \mathbf{z}^t + \sigma\beta^t(\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1})$ ; step ② uses the triangle inequality; step ③ uses  $\|\nabla f(\mathbf{X}^t)\|_F \leq C_f$ ,  $\|\nabla g(\mathbf{X}^t)\|_F \leq C_g$ ,  $\|\nabla \mathcal{A}(\mathbf{X}^t)\|_F \leq \|\nabla \mathcal{A}(\mathbf{X}^t)\| \leq \bar{\mathbf{A}}$ ,  $\|\mathbf{z}^t\| \leq \bar{\mathbf{z}}$ ,  $\frac{1}{\sigma} \leq 1$ ,  $\beta^t \leq \beta^{t-1}(1 + \xi)$ ; step ④ uses  $\xi \leq 1$ .

We derive the following inequalities:

$$\begin{aligned}
& L(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^t) - L(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^t) = \dot{\mathcal{L}}(\mathbf{X}^{t+1}) - \dot{\mathcal{L}}(\mathbf{X}^t) \\
& \stackrel{\textcircled{1}}{=} \{\mathcal{S}^t(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t) - g(\mathbf{X}^{t+1})\} - \{\mathcal{S}^t(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) - g(\mathbf{X}^t)\} \\
& \stackrel{\textcircled{2}}{\leq} \frac{1}{2} \ell(\beta^t) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_{\mathbb{F}}^2 + \langle \mathbf{G}^t, \mathbf{X}^{t+1} - \mathbf{X}^t \rangle \\
& \stackrel{\textcircled{3}}{=} \frac{1}{2} \ell(\beta^t) \|\text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t) - \mathbf{X}^t\|_{\mathbb{F}}^2 + \langle \mathbf{G}^t, \text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t) - \mathbf{X}^t + \eta^t \mathbb{G}_\rho^t \rangle - \eta^t \langle \mathbf{G}^t, \mathbb{G}_\rho^t \rangle \\
& \stackrel{\textcircled{4}}{\leq} \frac{1}{2} \ell(\beta^t) \|\text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t) - \mathbf{X}^t\|_{\mathbb{F}}^2 + \bar{g} \|\text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t) - \mathbf{X}^t + \eta^t \mathbb{G}_\rho^t\|_{\mathbb{F}} - \frac{\eta^t}{\max(1, 2\rho)} \|\mathbb{G}_\rho^t\|_{\mathbb{F}}^2 \\
& \stackrel{\textcircled{5}}{\leq} \frac{1}{2} \ell(\beta^t) \dot{k} \|\eta^t \mathbb{G}_\rho^t\|_{\mathbb{F}}^2 + \frac{1}{2} \bar{g} \ddot{k} \|\eta^t \mathbb{G}_\rho^t\|_{\mathbb{F}}^2 - \frac{\eta^t}{\max(1, 2\rho)} \|\mathbb{G}_\rho^t\|_{\mathbb{F}}^2 \\
& \stackrel{\textcircled{6}}{=} \eta^t \|\mathbb{G}_\rho^t\|_{\mathbb{F}}^2 \cdot \left\{ \frac{1}{2} \ell(\beta^t) \dot{k} \frac{b^t \gamma^j}{\beta^t} + \frac{1}{2} \bar{g} \ddot{k} \frac{b^t \gamma^j}{\beta^t} - \frac{1}{\max(1, 2\rho)} \right\} \\
& \stackrel{\textcircled{7}}{\leq} \eta^t \|\mathbb{G}_\rho^t\|_{\mathbb{F}}^2 \cdot \left\{ \left( \frac{\bar{b}}{2} \dot{k} \bar{\ell} + \frac{\bar{b}}{2\beta^0} \ddot{k} \bar{g} \right) \gamma^j - \frac{1}{\max(1, 2\rho)} \right\} \\
& \stackrel{\textcircled{8}}{\leq} \eta^t \|\mathbb{G}_\rho^t\|_{\mathbb{F}}^2 \cdot \{-\delta\}, \tag{41}
\end{aligned}$$

where step ① uses the definitions of  $L(\mathbf{X}, \mathbf{y}; \mathbf{z}; \beta, \mu)$ ; step ② uses the fact that the function  $g(\mathbf{X})$  is convex and the function  $\mathcal{S}(\mathbf{X}, \mathbf{y}^t; \mathbf{z}^t; \beta^t)$  is  $\ell(\beta^t)$ -smooth w.r.t.  $\mathbf{X}$ ; step ③ uses  $\mathbf{X}^{t+1} = \text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t)$ ; step ④ uses the Cauchy-Schwarz Inequality,  $\|\mathbf{G}^t\|_{\mathbb{F}} \leq \bar{g}$ , and Lemma 2.12(a) that  $\langle \mathbf{G}^t, \mathbb{G}_\rho^t \rangle \geq \frac{1}{\max(1, 2\rho)} \|\mathbb{G}_\rho^t\|_{\mathbb{F}}^2$ ; step ⑤ uses Lemma 2.10 with  $\Delta \triangleq -\eta^t \mathbb{G}_\rho^t$  given that  $\mathbf{X}^t \in \mathcal{M}$  and  $\Delta \in \mathbf{T}_{\mathbf{X}^t} \mathcal{M}$ ; step ⑥ uses  $\eta^t \triangleq \frac{b^t \gamma^j}{\beta^t}$ ; step ⑦ uses  $\ell(\beta^t) \leq \beta^t \bar{\ell}$ ,  $\beta^0 \leq \beta^t$ , and  $b^t \leq \bar{b}$ ; step ⑧ uses the fact that  $\gamma^j$  is sufficiently small such that:

$$\gamma^j \leq \frac{2\left(\frac{1}{\max(1, 2\rho)} - \delta\right)}{\bar{\ell} \bar{k} \bar{b} + \bar{g} \ddot{k} \bar{b} / \beta^0} \triangleq \bar{\gamma}. \tag{42}$$

Given Inequality (41) coincides with the condition of the line search procedure, we complete the proof.

**Part (b).** We derive the following inequalities:

$$\begin{aligned}
& L(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^t) - L(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^t) \\
& \stackrel{\textcircled{1}}{\leq} -\|\mathbb{G}_\rho^t\|_{\mathbb{F}}^2 \delta \eta^t \\
& \stackrel{\textcircled{2}}{\leq} -\|\mathbb{G}_{1/2}^t\|_{\mathbb{F}}^2 \delta \eta^t \cdot \min(1, 2\rho)^2 \\
& \stackrel{\textcircled{3}}{=} -\frac{1}{\beta^t} \|\mathbb{G}_{1/2}^t\|_{\mathbb{F}}^2 \cdot \delta b^t \gamma^{j-1} \gamma \cdot \min(1, 2\rho)^2 \\
& \stackrel{\textcircled{4}}{\leq} -\frac{1}{\beta^t} \|\mathbb{G}_{1/2}^t\|_{\mathbb{F}}^2 \cdot \underbrace{\delta \bar{b} \bar{\gamma} \gamma \cdot \min(1, 2\rho)^2}_{\triangleq \varepsilon_x},
\end{aligned}$$

where step ① uses Inequality (41); step ② uses Lemma 2.12(b) that  $\|\mathbb{G}_\rho\|_{\mathbb{F}} \geq \min(1, 2\rho) \|\mathbb{G}_{1/2}\|_{\mathbb{F}}$ ; step ③ uses the definition  $\eta^t \triangleq \frac{b^t \gamma^j}{\beta^t}$ ; step ④ uses  $b^t \geq \bar{b}$ , and the following inequality:

$$\gamma^{j-1} \geq \bar{\gamma} \geq \gamma^j,$$

which can be implied by the stopping criteria of the line search procedure.  $\square$

## C.9 PROOF OF LEMMA 4.12

*Proof.* We define:  $\Theta^t \triangleq L(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^{t-1}) + \mu^{t-1} C_h^2 + \mathbb{T}^t + \mathbb{Z}^t + 0 \times \mathbb{X}^t$ ,

We define  $\tilde{\varepsilon}_t \triangleq \|\mathbf{y}^t - \mathbf{y}^{t-1}\|^2 + \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|^2 + \|\frac{1}{\beta^t} \mathbb{G}_{1/2}^t\|_{\mathbb{F}}^2$ .

1782 **Part (a).** Using Lemma 4.5, we have:

$$1783 \quad L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}; \mathbf{z}^{t+1}; \beta^{t+1}, \mu^t) - L(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^{t-1}) - (\mu^{t-1} - \mu^t)C_h^2 \\ 1784 \quad \leq \mathbb{T}^t - \mathbb{T}^{t+1} + \mathbb{Z}^t - \mathbb{Z}^{t+1} - \varepsilon_y \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 - \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2. \quad (43) \\ 1785$$

1786 Using Lemma 4.10, we have:

$$1787 \quad L(\mathbf{X}^{t+1}, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^{t-1}) - L(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^{t-1}) \leq 0 \times \mathbb{X}^t - 0 \times \mathbb{X}^{t+1} - \varepsilon_x \beta^t \|\frac{1}{\beta^t} \mathbb{G}_{1/2}^t\|_{\mathbb{F}}^2. \\ 1788$$

1789 Adding these two inequalities together and using the definition of  $\Theta^t$ , we have:

$$1790 \quad \Theta^t - \Theta^{t+1} \geq \varepsilon_y \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \varepsilon_x \beta^t \|\frac{1}{\beta^t} \mathbb{G}_{1/2}^t\|_{\mathbb{F}}^2 + \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 \\ 1791 \quad \geq \min(\varepsilon_y, \varepsilon_x, \varepsilon_z) \cdot \beta^t \cdot \tilde{e}_{t+1}. \\ 1792$$

1793 **Part (b).** Using the same strategy as in deriving Lemma 4.7(b), we finish the proof. □

### 1796 C.10 PROOF OF THEOREM 4.13

1797 *Proof.* We define  $\text{Crit}(\mathbf{X}, \mathbf{y}, \mathbf{z}) \triangleq \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\| + \|\partial h(\mathbf{y}) - \mathbf{z}\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\nabla f(\mathbf{X}) - \partial g(\mathbf{X}) + \mathcal{A}^\top(\mathbf{z}))\|_{\mathbb{F}}$ .

1800 We define  $\dot{\mathbf{G}} \triangleq \nabla f(\mathbf{X}^t) - \partial g(\mathbf{X}^t) + \mathcal{A}^\top(\mathbf{z}^t)$ , and  $\ddot{\mathbf{G}} \triangleq \beta^t \mathcal{A}^\top(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)$ .

1801 We let  $\mathbf{G} = \mathbf{G}^t \in \partial_{\mathbf{X}} L(\mathbf{X}^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^t)$ .

1802 First, we obtain:

$$1803 \quad \mathbb{G}_{1/2}^t \stackrel{\textcircled{1}}{=} \mathbf{G} - \frac{1}{2} \mathbf{X}^t \mathbf{G}^\top \mathbf{X}^t - \frac{1}{2} \mathbf{X}^t [\mathbf{X}^t]^\top \mathbf{G} \\ 1804 \quad \stackrel{\textcircled{2}}{=} (\dot{\mathbf{G}} - \frac{1}{2} \mathbf{X}^t \dot{\mathbf{G}}^\top \mathbf{X}^t - \frac{1}{2} \mathbf{X}^t [\mathbf{X}^t]^\top \dot{\mathbf{G}}) + (\ddot{\mathbf{G}} - \frac{1}{2} \mathbf{X}^t \ddot{\mathbf{G}}^\top \mathbf{X}^t - \frac{1}{2} \mathbf{X}^t [\mathbf{X}^t]^\top \ddot{\mathbf{G}}) \\ 1805 \quad \stackrel{\textcircled{3}}{=} \text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\dot{\mathbf{G}}) + \text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\ddot{\mathbf{G}}) \\ 1806$$

1807 where step  $\textcircled{1}$  uses the definition  $\mathbb{G}_\rho^t \triangleq \mathbf{G} - \rho \mathbf{X}^t \mathbf{G}^\top \mathbf{X}^t - (1 - \rho) \mathbf{X}^t [\mathbf{X}^t]^\top \mathbf{G}$ , as shown in Algorithm 1; step  $\textcircled{2}$  uses  $\mathbf{G} \in \dot{\mathbf{G}} + \ddot{\mathbf{G}}$ ; step  $\textcircled{3}$  uses the fact that  $\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\Delta) = \Delta - \frac{1}{2} \mathbf{X}(\Delta^\top \mathbf{X} + \mathbf{X}^\top \Delta)$  for all  $\Delta \in \mathbb{R}^{n \times r}$  (Absil et al., 2008a). This leads to:

$$1808 \quad \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\dot{\mathbf{G}})\|_{\mathbb{F}} = \|\mathbb{G}_{1/2}^t - \text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\ddot{\mathbf{G}})\|_{\mathbb{F}} \\ 1809 \quad \stackrel{\textcircled{1}}{\leq} \|\mathbb{G}_{1/2}^t\|_{\mathbb{F}} + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\ddot{\mathbf{G}})\|_{\mathbb{F}} \\ 1810 \quad \stackrel{\textcircled{2}}{\leq} \|\mathbb{G}_{1/2}^t\|_{\mathbb{F}} + \|\ddot{\mathbf{G}}\|_{\mathbb{F}} \\ 1811 \quad \leq \|\mathbb{G}_{1/2}^t\|_{\mathbb{F}} + \beta^t \bar{\mathbf{A}} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| \\ 1812 \quad \leq \beta^t e^{t+1} + \mathcal{O}(\beta^{t-1} e^t), \\ 1813$$

1814 where step  $\textcircled{1}$  uses the triangle inequality; step  $\textcircled{2}$  uses Lemma 2.11 that  $\|\text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\Delta)\|_{\mathbb{F}} \leq \|\Delta\|_{\mathbb{F}}$  for all  $\Delta \in \mathbb{R}^{n \times r}$ .

1815 Finally, we derive:

$$1816 \quad \frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{X}^t, \check{\mathbf{y}}^t, \mathbf{z}^t) \\ 1817 \quad \stackrel{\textcircled{1}}{=} \frac{1}{T} \sum_{t=1}^T \{\|\mathcal{A}(\mathbf{X}^t) - \check{\mathbf{y}}^t\| + \|\partial h(\check{\mathbf{y}}^t) - \mathbf{z}^t\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\dot{\mathbf{G}})\|_{\mathbb{F}}\} \\ 1818 \quad \stackrel{\textcircled{2}}{\leq} \frac{1}{T} \sum_{t=1}^T \{\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \|(1 - \frac{1}{\sigma})(\mathbf{z}^t - \mathbf{z}^{t-1})\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\dot{\mathbf{G}})\|_{\mathbb{F}}\} \\ 1819 \quad \stackrel{\textcircled{3}}{=} \frac{1}{T} \sum_{t=1}^T \{\mathcal{O}(\beta^t e^{t+1}) + \mathcal{O}(\beta^{t-1} e^t)\} \\ 1820 \quad \stackrel{\textcircled{4}}{=} \mathcal{O}(T^{(p-1)/2}) = \mathcal{O}(T^{-1/3}), \\ 1821$$

1822 where step  $\textcircled{1}$  uses the definition of  $\text{Crit}(\mathbf{X}, \mathbf{y}, \mathbf{z})$ ; step  $\textcircled{2}$  uses  $\mathbf{z}^{t+1} - \partial h(\check{\mathbf{y}}^{t+1}) \ni (1 - \frac{1}{\sigma})(\mathbf{z}^{t+1} - \mathbf{z}^t)$ , as shown in Lemma 4.1; step  $\textcircled{3}$  uses  $\|\mathbf{z}^t - \mathbf{z}^{t-1}\| = \|\sigma \beta^{t-1} (\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)\| \leq 2\beta^t \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| = \mathcal{O}(\beta^{t-1} e^t)$ ; step  $\textcircled{4}$  uses the choice  $p = 1/3$  and Lemma 4.7(b). □

1836 D PROOFS FOR SECTION 5

1837 D.1 PROOF OF LEMMA 5.4

1838 We begin by presenting the following four useful lemmas.

1839 **Lemma D.1.** *For both OADMM-EP and OADMM-RR, we have:*

$$1840 \quad (\mathbf{d}_X, \mathbf{d}_{X^-}, \mathbf{d}_Y, \mathbf{d}_Z) \in \partial\Theta(\mathbf{X}^t, \mathbf{X}^{t-1}, \mathbf{y}^t, \mathbf{z}^t; \beta^t, \beta^{t-1}, \mu^{t-1}, t), \quad (44)$$

1841 where  $\mathbf{d}_X \triangleq \mathbb{A}^t + \{\beta^t + 2\omega\ddot{\sigma}\sigma^2\beta^{t-1}\} \cdot \mathcal{A}^\top(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t) + \alpha(\theta+1)\ell(\beta^t)(\mathbf{X}^t - \mathbf{X}^{t-1})$ ,  $\mathbf{d}_{X^-} \triangleq \alpha(\theta+1)\ell(\beta^t)(\mathbf{X}^{t-1} - \mathbf{X}^t)$ ,  $\mathbf{d}_Y \triangleq \nabla h_{\mu^{t-1}}(\mathbf{y}^t) - \mathbf{z}^t + (\mathbf{y}^t - \mathcal{A}(\mathbf{X}^t)) \cdot (\beta^t + 2\omega\ddot{\sigma}\sigma^2\beta^{t-1})$ ,  $\mathbf{d}_Z \triangleq \mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t$ .  
 1842 Here,  $\mathbb{A}^t \triangleq \partial I_{\mathcal{M}}(\mathbf{X}^t) + \nabla f(\mathbf{X}^t) - \nabla g(\mathbf{X}^t) + \mathcal{A}^\top(\mathbf{z}^t)$ .

1843 *Proof.* We define the Lyapunov function as:  $\Theta(\mathbf{X}, \mathbf{X}^-, \mathbf{y}, \mathbf{z}; \beta, \beta^-, \mu^-, t) \triangleq L(\mathbf{X}, \mathbf{y}; \mathbf{z}; \beta, \mu^-) + \omega\ddot{\sigma}\sigma^2\beta^- \|\mathcal{A}(\mathbf{X}) - \mathbf{y}\|_2^2 + \frac{\alpha(\theta+1)\ell(\beta)}{2} \|\mathbf{X} - \mathbf{X}^-\|_F^2 + \frac{4\omega\ddot{\sigma}}{\beta^0} C_h^2 \frac{1}{t} + C_h^2 \mu^-$ .

1844 Using this definition, we can promptly derive the conclusion of the lemma.

1845  $\square$

1846 **Lemma D.2.** *For OADMM-EP, we define  $\{\mathbf{d}_X, \mathbf{d}_{X^-}, \mathbf{d}_Y, \mathbf{d}_Z\}$  as in Lemma D.1. There exists a constant  $K$  such that:*

$$1847 \quad \frac{1}{\beta^t} \{\|\mathbf{d}_X\|_F + \|\mathbf{d}_{X^-}\|_F + \|\mathbf{d}_Y\| + \|\mathbf{d}_Z\|\} \leq K \{\mathcal{X}^t + \mathcal{Z}^t + \mathcal{X}^{t-1} + \mathcal{Z}^{t-1}\}. \quad (45)$$

1848 Here,  $\mathcal{X}^t \triangleq \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$ , and  $\mathcal{Z}^t \triangleq \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|$ .

1849 *Proof.* First, we obtain:

$$1850 \quad \begin{aligned} & \frac{1}{\beta^t} \|\mathbb{A}^t\|_F = \|\partial I_{\mathcal{M}}(\mathbf{X}^t) + \nabla f(\mathbf{X}^t) - \nabla g(\mathbf{X}^t) + \mathcal{A}^\top(\mathbf{z}^t)\|_F \\ & \stackrel{\textcircled{1}}{=} \frac{1}{\beta^t} \|\nabla g(\mathbf{X}^{t-1}) - \nabla g(\mathbf{X}^t) + \nabla f(\mathbf{X}^t) - \nabla f(\mathbf{X}_c^{t-1}) - \theta\ell(\beta^{t-1})(\mathbf{X}^t - \mathbf{X}_c^{t-1}) \\ & \quad + \mathcal{A}^\top(\mathbf{z}^t - \mathbf{z}^{t-1}) - \beta^{t-1}\mathcal{A}^\top(\mathcal{A}(\mathbf{X}_c^{t-1}) - \mathbf{y}^{t-1})\|_F \\ & \stackrel{\textcircled{2}}{\leq} \frac{1}{\beta^t} L_g \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F + \frac{1}{\beta^t} (L_f + \theta\ell(\beta^{t-1})) \|\mathbf{X}^t - \mathbf{X}_c^{t-1}\|_F \\ & \quad + \frac{1}{\beta^t} \bar{\mathbb{A}} \|\mathbf{z}^t - \mathbf{z}^{t-1}\| + \frac{1}{\beta^t} \beta^{t-1} \bar{\mathbb{A}} \{\|\mathcal{A}(\mathbf{X}^{t-1}) - \mathbf{y}^{t-1}\| + \bar{\mathbb{A}} \|\mathbf{X}^{t-1} - \mathbf{X}^{t-2}\|_F\} \\ & = \mathcal{O}(\|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F) + \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|) \\ & \quad + \mathcal{O}(\|\mathbf{X}^{t-1} - \mathbf{X}^{t-2}\|_F) + \mathcal{O}(\|\mathcal{A}(\mathbf{X}^{t-1}) - \mathbf{y}^{t-1}\|), \end{aligned} \quad (46)$$

1851 where step  $\textcircled{1}$  uses the optimality of  $\mathbf{X}^{t+1}$  for OADMM-EP that:

$$1852 \quad \begin{aligned} & \partial I_{\mathcal{M}}(\mathbf{X}^{t+1}) - \nabla g(\mathbf{X}^t) \\ & \ni -\theta\ell(\beta^t)(\mathbf{X}^{t+1} - \mathbf{X}_c^t) - \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}_c^t, \mathbf{y}^t; \mathbf{z}^t; \beta^t) \\ & = -\theta\ell(\beta^t)(\mathbf{X}^{t+1} - \mathbf{X}_c^t) - \nabla f(\mathbf{X}_c^t) - \mathcal{A}^\top[\mathbf{z}^t + \beta^t(\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t)]; \end{aligned} \quad (47)$$

1853 step  $\textcircled{2}$  uses the triangle inequality, the  $L_f$ -Lipschitz continuity of  $\nabla f(\mathbf{X})$  for all  $\mathbf{X}$ ; the  $L_g$ -Lipschitz continuity of  $\nabla g(\mathbf{X})$ , and the upper bound of  $\|\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t\|$  as shown in Lemma A.6(c); step  $\textcircled{3}$  uses the upper bound of  $\|\mathbf{X}^t - \mathbf{X}_c^{t-1}\|_F$ , and  $\mathbf{z}^t - \mathbf{z}^{t-1} = \sigma\beta^{t-1}(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)$ .

1854 **Part (a).** We bound the term  $\frac{1}{\beta^t} \|\mathbf{d}_X\|_F$ . We have:

$$1855 \quad \begin{aligned} & \frac{1}{\beta^t} \|\mathbf{d}_X\|_F \\ & \stackrel{\textcircled{1}}{=} \frac{1}{\beta^t} \|\mathbb{A}^t + (\beta^t + 2\omega\ddot{\sigma}\sigma^2\beta^{t-1})\mathcal{A}^\top(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t) + \alpha(\theta+1)\ell(\beta^t)(\mathbf{X}^t - \mathbf{X}^{t-1})\|_F \\ & \stackrel{\textcircled{2}}{\leq} \frac{1}{\beta^t} \|\mathbb{A}^t\|_F + (1 + 2\omega\ddot{\sigma}\sigma^2)\bar{\mathbb{A}} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|_F + \alpha(\theta+1)\bar{\ell} \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F \\ & \stackrel{\textcircled{3}}{\leq} \mathcal{O}(\|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F) + \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|) + \mathcal{O}(\|\mathbf{X}^{t-1} - \mathbf{X}^{t-2}\|_F) + \mathcal{O}(\|\mathcal{A}(\mathbf{X}^{t-1}) - \mathbf{y}^{t-1}\|), \end{aligned}$$

where step ① uses the definition of  $\mathbf{d}_X$  in Lemma D.1; step ② uses the triangle inequality,  $\beta^{t-1} \leq \beta^t$ , and  $\ell(\beta^t) \leq \beta^t \bar{\ell}$ ; step ③ uses Inequality (46).

**Part (b).** We bound the term  $\frac{1}{\beta^t} \|\mathbf{d}_{X^-}\|_F$ . We have:

$$\frac{1}{\beta^t} \|\mathbf{d}_{X^-}\|_F \stackrel{\text{①}}{=} \frac{1}{\beta^t} \alpha(\theta + 1) \ell(\beta^t) \|\mathbf{X}^{t-1} - \mathbf{X}^t\|_F \stackrel{\text{②}}{=} \mathcal{O}(\|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F), \quad (48)$$

where step ① uses the definition of  $\mathbf{d}_{X^-}$  in Lemma D.1; step ② uses  $\ell(\beta^t) \leq \beta^t \bar{\ell}$ .

**Part (c).** We bound the term  $\frac{1}{\beta^t} \|\mathbf{d}_Y\|_F$ . We have:

$$\begin{aligned} \frac{1}{\beta^t} \|\mathbf{d}_Y\|_F &\stackrel{\text{①}}{=} \frac{1}{\beta^t} \|\nabla h_{\mu^{t-1}}(\mathbf{y}^t) - \mathbf{z}^t + (\mathbf{y}^t - \mathcal{A}(\mathbf{X}^t)) \cdot (\beta^t + 2\omega\ddot{\sigma}\sigma^2\beta^{t-1})\| \\ &\stackrel{\text{②}}{=} \frac{1}{\beta^t} \|(1 - \frac{1}{\sigma})(\mathbf{z}^{t-1} - \mathbf{z}^t) + (\mathbf{y}^t - \mathcal{A}(\mathbf{X}^t)) \cdot (\beta^t + 2\omega\ddot{\sigma}\sigma^2\beta^{t-1})\| \\ &\stackrel{\text{③}}{=} \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|), \end{aligned}$$

where step ① uses the definition of  $\mathbf{d}_Y$  in Lemma D.1; step ② uses the fact that  $\mathbf{z}^t - \frac{1}{\sigma}(\mathbf{z}^t - \mathbf{z}^{t+1}) = \nabla h_{\mu^t}(\mathbf{y}^{t+1})$ , as shown in Lemma 4.1; step ③ uses  $\frac{1}{\beta^t}(\mathbf{z}^{t+1} - \mathbf{z}^t) = \sigma(\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1})$ , and  $\beta^{t-1} = \mathcal{O}(\beta^t)$ .

**Part (d).** We bound the term  $\frac{1}{\beta^t} \|\mathbf{d}_Z\|_F$ . We have:  $\frac{1}{\beta^t} \|\mathbf{d}_Z\|_F \leq \frac{1}{\beta^t} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|$ .

**Part (e).** Combining the upper bounds for the terms  $\{\frac{1}{\beta^t} \|\mathbf{d}_X\|_F, \frac{1}{\beta^t} \|\mathbf{d}_{X^-}\|_F, \frac{1}{\beta^t} \|\mathbf{d}_Y\|_F, \frac{1}{\beta^t} \|\mathbf{d}_Z\|_F\}$ , we finish the proof of this lemma.  $\square$

**Lemma D.3.** For OADMM-RR, we define  $\{\mathbf{d}_X, \mathbf{d}_{X^-}, \mathbf{d}_Y, \mathbf{d}_Z\}$  as in Lemma D.1. There exists a constant  $K$  such that :

$$\frac{1}{\beta^t} \{\|\mathbf{d}_X\|_F + \|\mathbf{d}_{X^-}\|_F + \|\mathbf{d}_Y\|_F + \|\mathbf{d}_Z\|_F\} \leq K\{\mathcal{X}^t + \mathcal{Z}^t\},$$

Here,  $\mathcal{X}^t \triangleq \|\frac{1}{\beta^t} \mathbb{G}_{1/2}\|_F$ , and  $\mathcal{Z}^t \triangleq \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|$ .

*Proof.* We define  $\mathbf{G}^t \triangleq \nabla f(\mathbf{X}^t) - \nabla g(\mathbf{X}^t) + \mathbf{A}^\top(\mathbf{z}^t + \beta^t(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t))$ .

We define  $\dot{\mathcal{L}}(\mathbf{X}) \triangleq L(\mathbf{X}, \mathbf{y}^t; \mathbf{z}^t; \beta^t, \mu^t)$ , we have:  $\nabla \dot{\mathcal{L}}(\mathbf{X}^t) = \mathbf{G}^t$ .

First, given  $\mathbf{X}^t \in \mathcal{M}$ , we obtain:

$$\begin{aligned} \frac{1}{\beta^t} \|\partial I_{\mathcal{M}}(\mathbf{X}^t) + \nabla \dot{\mathcal{L}}(\mathbf{X}^t)\|_F &\stackrel{\text{①}}{\leq} \frac{1}{\beta^t} \|\nabla \dot{\mathcal{L}}(\mathbf{X}^t) - \mathbf{X}^t [\nabla \dot{\mathcal{L}}(\mathbf{X}^t)]^\top \mathbf{X}^t\|_F \\ &\stackrel{\text{②}}{=} \frac{1}{\beta^t} \|\mathbf{G}^t - \mathbf{X}^t [\mathbf{G}^t]^\top \mathbf{X}^t\|_F = \frac{1}{\beta^t} \|\mathbb{G}_1^t\|_F \\ &\stackrel{\text{③}}{\leq} \frac{1}{\beta^t} \max(1, 1/\rho) \cdot \|\mathbb{G}_{1/2}\|_F = \mathcal{O}(\mathcal{X}^t), \end{aligned} \quad (49)$$

where step ① uses Lemma 2.13; step ② uses the definitions of  $\{\mathbf{G}^t, \mathbf{D}_\rho^t\}$  as in Algorithm 1; step ③ uses  $\|\mathbb{G}_1\|_F \leq \max(1, 1/\rho) \|\mathbb{G}_\rho\|_F$ , as shown in Lemma 2.12(b).

**Part (a).** We bound the term  $\frac{1}{\beta^t} \|\mathbf{d}_X\|_F$ . We have:

$$\begin{aligned} &\frac{1}{\beta^t} \|\mathbf{d}_X\|_F \\ &\stackrel{\text{①}}{=} \frac{1}{\beta^t} \|\partial I_{\mathcal{M}}(\mathbf{X}^t) + \nabla \dot{\mathcal{L}}(\mathbf{X}^t) + 2\omega\ddot{\sigma}\sigma^2\beta^{t-1} \mathbf{A}^\top(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)\|_F \\ &\stackrel{\text{②}}{\leq} \frac{1}{\beta^t} \|\partial I_{\mathcal{M}}(\mathbf{X}^t) + \nabla \dot{\mathcal{L}}(\mathbf{X}^t)\|_F + 2\omega\ddot{\sigma}\sigma^2 \|\mathbf{A}^\top(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)\|_F \\ &\stackrel{\text{③}}{\leq} \mathcal{O}(\mathcal{X}^t) + \mathcal{O}(\mathcal{Z}^t), \end{aligned}$$

where step ① uses  $\mathbf{d}_X = \partial I_{\mathcal{M}}(\mathbf{X}^t) + \nabla f(\mathbf{X}^t) - \nabla g(\mathbf{X}^t) + \mathbf{A}^\top(\mathbf{z}^t) + \{\beta^t + 2\omega\ddot{\sigma}\sigma^2\beta^{t-1}\} \cdot \mathbf{A}^\top(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)$  with the choice  $\alpha = 0$  for OADMM-RR; step ② uses the triangle inequality and  $\beta^{t-1} \leq \beta^t$ ; step ③ uses Inequality (49).

**Part (b).** We bound the term  $\frac{1}{\beta^t} \|\mathbf{d}_{X^-}\|_F$ . Given  $\alpha = 0$ , we conclude that  $\frac{1}{\beta^t} \|\mathbf{d}_{X^-}\|_F = 0$ .

**Part (c).** We bound the terms  $\frac{1}{\beta^t} \|\mathbf{d}_y\|_F$  and  $\frac{1}{\beta^t} \|\mathbf{d}_z\|_F$ . Considering that the same strategies for updating  $\{\mathbf{y}^t, \mathbf{z}^t\}$  are employed, their bounds in OADMM-RR are identical to those in OADMM-ER.

**Part (d).** Combining the upper bounds for the terms  $\{\frac{1}{\beta^t} \|\mathbf{d}_x\|_F, \frac{1}{\beta^t} \|\mathbf{d}_{x^-}\|_F, \frac{1}{\beta^t} \|\mathbf{d}_y\|_F, \frac{1}{\beta^t} \|\mathbf{d}_z\|_F\}$ , we finish the proof of this lemma.  $\square$

Now, we proceed to prove the main result of this lemma.

**Lemma D.4.** (Subgradient Bounds) **(a)** For OADMM-EP, there exists a constant  $K > 0$  such that:  $\text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbf{u}^t)) \leq \beta^t K (e^t + e^{t-1})$ . **(b)** For OADMM-RR, there exists a constant  $K > 0$  such that:  $\text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbf{u}^t)) \leq \beta^t K e^t$ . Here,  $\text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbf{u}^t)) \triangleq \{\text{dist}^2(\mathbf{0}, \partial_{\mathbf{X}}\Theta(\mathbb{w}^t; \mathbf{u}^t)) + \text{dist}^2(\mathbf{0}, \partial_{\mathbf{X}^-}\Theta(\mathbb{w}^t; \mathbf{u}^t)) + \text{dist}^2(\mathbf{0}, \partial_{\mathbf{y}}\Theta(\mathbb{w}^t; \mathbf{u}^t)) + \text{dist}^2(\mathbf{0}, \partial_{\mathbf{z}}\Theta(\mathbb{w}^t; \mathbf{u}^t))\}^{1/2}$ .

*Proof.* For OADMM-EP, we have:

$$\begin{aligned} \text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbf{u}^t)) &= \sqrt{\|\mathbf{d}_x\|_F^2 + \|\mathbf{d}_{x^-}\|_F^2 + \|\mathbf{d}_y\|_F^2 + \|\mathbf{d}_z\|_F^2} \\ &\stackrel{\textcircled{1}}{\leq} \|\mathbf{d}_x\|_F + \|\mathbf{d}_{x^-}\|_F + \|\mathbf{d}_y\|_F + \|\mathbf{d}_z\|_F \\ &\stackrel{\textcircled{2}}{\leq} K\beta^t\{\mathcal{X}^t + \mathcal{Z}^t + \mathcal{X}^{t-1} + \mathcal{Z}^{t-1}\} \\ &\stackrel{\textcircled{3}}{\leq} K\beta^t(e^t + e^{t-1}), \end{aligned}$$

where step  $\textcircled{1}$  uses  $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$  for all  $a \geq 0$  and  $b \geq 0$ ; step  $\textcircled{2}$  uses Lemma D.2; step  $\textcircled{3}$  uses the definition of  $K$ .

For OADMM-RR, using Lemma D.3 and similar strategies, we have:  $\text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbf{u}^t)) \leq \beta^t K e^t$ .  $\square$

## D.2 PROOF OF THEOREM 5.6

*Proof.* We define  $\dot{K} \triangleq 3K / \min(\varepsilon_x, \varepsilon_y, \varepsilon_z)$ .

Firstly, using Assumption 5.1, we have:

$$\varphi'(\Theta(\mathbb{w}^t; \mathbf{u}^t) - \Theta(\mathbb{w}^\infty; \mathbf{u}^\infty)) \cdot \text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbf{u}^t)) \geq 1. \quad (50)$$

Secondly, given the desingularization function  $\varphi(\cdot)$  is concave, for any  $a, b \in \mathbb{R}$ , we have:  $\varphi(b) + (a-b)\varphi'(a) \leq \varphi(a)$ . Applying the inequality above with  $a = \Theta(\mathbb{w}^t; \mathbf{u}^t) - \Theta(\mathbb{w}^\infty; \mathbf{u}^\infty)$  and  $b = \Theta(\mathbb{w}^{t+1}; \mathbf{u}^{t+1}) - \Theta(\mathbb{w}^\infty; \mathbf{u}^\infty)$ , we have:

$$\begin{aligned} &(\Theta(\mathbb{w}^t; \mathbf{u}^t) - \Theta(\mathbb{w}^{t+1}; \mathbf{u}^{t+1})) \cdot \varphi'(\Theta(\mathbb{w}^t; \mathbf{u}^t) - \Theta(\mathbb{w}^\infty; \mathbf{u}^\infty)) \\ &\leq \underbrace{\varphi(\Theta(\mathbb{w}^t; \mathbf{u}^t) - \Theta(\mathbb{w}^\infty; \mathbf{u}^\infty))}_{\triangleq \varphi^t} - \underbrace{\varphi(\Theta(\mathbb{w}^{t+1}; \mathbf{u}^{t+1}) - \Theta(\mathbb{w}^\infty; \mathbf{u}^\infty))}_{\triangleq \varphi^{t+1}}. \end{aligned} \quad (51)$$

Third, we derive the following inequalities for OADMM-EP:

$$\begin{aligned} &\min(\varepsilon_z, \varepsilon_y, \varepsilon_x) \beta^t \{\|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2\} \\ &\stackrel{\textcircled{1}}{\leq} \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + \varepsilon_y \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \varepsilon_x \ell(\beta^t) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 \\ &\stackrel{\textcircled{2}}{\leq} \Theta^t - \Theta^{t+1} = \Theta(\mathbb{w}^t; \mathbf{u}^t) - \Theta(\mathbb{w}^{t+1}; \mathbf{u}^{t+1}) \\ &\stackrel{\textcircled{3}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot \frac{1}{\varphi'(\Theta(\mathbb{w}^t; \mathbf{u}^t) - \Theta(\mathbb{w}^\infty; \mathbf{u}^\infty))} \\ &\stackrel{\textcircled{4}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot \text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbf{u}^t)) \\ &\stackrel{\textcircled{5}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot K\beta^t(e^t + e^{t-1}), \end{aligned} \quad (52)$$

where step ① uses  $\ell(\beta^t) \geq \beta^t \ell$ ; step ② uses Lemma 4.7; step ③ uses Inequality (51); step ④ uses Inequality (50); step ⑤ uses Lemma 5.4. We further derive the following inequalities:

$$\begin{aligned}
(e^{t+1})^2 &\triangleq (\|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2 + \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2 + \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F)^2 \\
&\stackrel{\textcircled{1}}{\leq} 3 \cdot \{\|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2\} \\
&\stackrel{\textcircled{2}}{\leq} \{3K / \min(\varepsilon_z, \varepsilon_y, \varepsilon_x)\} \cdot (e^t + e^{t-1}) \cdot (\varphi^t - \varphi^{t+1}), \tag{53}
\end{aligned}$$

where step ① uses the norm inequality that  $(a + b + c)^2 \leq 3(a^2 + b^2 + c^2)$  for any  $a, b, c \in \mathbb{R}$ ; step ② uses Inequality (52).

Fourth, we derive the following inequalities for OADMM-RR:

$$\begin{aligned}
&\min(\varepsilon_z, \varepsilon_y, \varepsilon_x) \beta^t \{\|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \|\frac{1}{\beta} \mathbb{G}_{1/2}^t\|_F^2\} \\
&\leq \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + \varepsilon_y \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \frac{\varepsilon_x}{\beta} \|\mathbb{G}_{1/2}^t\|_F^2 \\
&\stackrel{\textcircled{1}}{\leq} \Theta^t - \Theta^{t+1} = \Theta(\mathbf{w}^t; \mathbf{u}^t) - \Theta(\mathbf{w}^{t+1}; \mathbf{u}^{t+1}) \\
&\stackrel{\textcircled{2}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot \frac{1}{\varphi'(\Theta(\mathbf{w}^t; \mathbf{u}^t) - \Theta(\mathbf{w}^\infty; \mathbf{u}^\infty))} \\
&\stackrel{\textcircled{3}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot \text{dist}(\mathbf{0}, \partial\Theta(\mathbf{w}^t; \mathbf{u}^t)) \\
&\stackrel{\textcircled{4}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot K \beta^t (e^t + e^{t-1}), \tag{54}
\end{aligned}$$

where step ① uses Lemma 4.12; step ② uses Inequality (51); step ③ uses Inequality (50); step ④ uses Lemma 5.4. We further derive the following inequalities:

$$\begin{aligned}
(e^{t+1})^2 &\triangleq (\|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2 + \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2 + \|\frac{1}{\beta} \mathbb{G}_{1/2}^t\|_F)^2 \\
&\stackrel{\textcircled{1}}{\leq} 3 \cdot \{\|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \|\frac{1}{\beta} \mathbb{G}_{1/2}^t\|_F^2\} \\
&\stackrel{\textcircled{2}}{\leq} \{3K / \min(\varepsilon_z, \varepsilon_y, \varepsilon_x)\} \cdot (\varphi^t - \varphi^{t+1}) \cdot (e^t + e^{t-1}), \tag{55}
\end{aligned}$$

where step ① uses the norm inequality that  $(a + b + c)^2 \leq 3(a^2 + b^2 + c^2)$  for any  $a, b, c \in \mathbb{R}$ ; step ② uses Inequality (54).

**Part (a).** Given Inequalities (53) and (55), we establish the following unified inequality applicable to both OADMM-EP and OADMM-RR:

$$(e^{t+1})^2 \leq (e^t + e^{t-1}) \cdot \underbrace{\{3K / \min(\varepsilon_z, \varepsilon_y, \varepsilon_x)\}}_{\triangleq \dot{K}} \cdot (\varphi^t - \varphi^{t+1}). \tag{56}$$

**Part (b).** Considering Inequality (56) and applying Lemma A.10 with  $p^t \triangleq \dot{K} \varphi^t$ , we have:

$$\forall t, \sum_{i=t}^{\infty} e^{i+1} \leq e^t + e^{t-1} + 4\dot{K} \varphi^t.$$

Letting  $t = 1$ , we have:  $\sum_{i=1}^{\infty} e^{i+1} \leq e^1 + e^0 + 4\dot{K} \varphi^1$ .

□

### D.3 PROOF OF LEMMA 5.8

*Proof.* We define  $d^t \triangleq \sum_{i=t}^{\infty} e^{i+1}$ .

**Part (a-i).** For OADMM-EP, we have for all  $t \geq 1$ :  $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F \stackrel{\textcircled{1}}{\leq} \sum_{i=t}^{\infty} \|\mathbf{X}^i - \mathbf{X}^{i+1}\|_F \leq \sum_{i=t}^{\infty} \{\|\mathbf{X}^{i+1} - \mathbf{X}^i\|_F + \|\mathbf{y}^{i+1} - \mathbf{y}^i\|_2 + \|\mathcal{A}(\mathbf{X}^{i+1}) - \mathbf{y}^{i+1}\|_2\} = \sum_{i=t}^{\infty} e^{i+1} \triangleq d^t$ , where step ① use the triangle inequality.

**Part (a-ii).** For OADMM-RR, we have:  $\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F \stackrel{\textcircled{1}}{=} \|\text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t) - \mathbf{X}^t\|_F \stackrel{\textcircled{2}}{\leq} \dot{k} \|\eta^t \mathbb{G}_\rho^t\|_F \stackrel{\textcircled{3}}{\leq} \dot{k} \eta^t \max(2\rho, 1) \|\mathbb{G}_{1/2}^t\|_F \stackrel{\textcircled{4}}{=} \dot{k} \max(2\rho, 1) \frac{b^t \gamma^j}{\beta^t} \|\mathbb{G}_{1/2}^t\|_F \stackrel{\textcircled{5}}{\leq} \dot{k} \max(2\rho, 1) \bar{b} \bar{\gamma} \cdot$

2052  $\|\frac{1}{\beta^t} \mathbb{G}_{1/2}^t\|_F = \mathcal{O}(\|\frac{1}{\beta^t} \mathbb{G}_{1/2}^t\|_F)$ , where step ① uses the update rule of  $\mathbf{X}^{t+1}$ ; step ② uses Lemma  
 2053 2.10; step ③ uses Lemma 2.12(c); step ④ uses the definition of  $\eta^t \triangleq \frac{b^t \gamma^j}{\beta^t}$ ; step ⑤ uses  $b^t \leq \bar{b}$ , and  
 2054 the fact that  $\gamma^j \leq \bar{\gamma}$ . Furthermore, we derive for all  $t \geq 1$ :  $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F \leq \sum_{i=t}^\infty \|\mathbf{X}^i - \mathbf{X}^{i+1}\|_F \leq$   
 2055  $\mathcal{O}(\sum_{i=t}^\infty \|\frac{1}{\beta^i} \mathbb{G}_{1/2}^i\|_F) \leq \mathcal{O}(\sum_{i=t}^\infty e^{i+1}) = \mathcal{O}(d^t)$ .

2056 **Part (b).** We define  $\varphi^t \triangleq \varphi(s^t)$ , where  $s^t \triangleq \Theta(\mathbb{w}^t; \mathbb{u}^t) - \Theta(\mathbb{w}^\infty; \mathbb{u}^\infty)$ . Using the definition of  $d^t$ ,  
 2057 we derive:

$$\begin{aligned}
 2060 \quad d^t &\triangleq \sum_{i=t}^\infty e^{i+1} \\
 2061 &\stackrel{\textcircled{1}}{\leq} e^t + e^{t-1} + 4\dot{K}\varphi^t \\
 2062 &\stackrel{\textcircled{2}}{=} e^t + e^{t-1} + 4\dot{K}\tilde{c} \cdot \{[s^t]^{\tilde{\sigma}}\}^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\
 2063 &\stackrel{\textcircled{3}}{=} e^t + e^{t-1} + 4\dot{K}\tilde{c} \cdot \{\tilde{c}(1-\tilde{\sigma}) \cdot \frac{1}{\varphi'(s^t)}\}^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\
 2064 &\stackrel{\textcircled{4}}{\leq} e^t + e^{t-1} + 4\dot{K}\tilde{c} \cdot \{\tilde{c}(1-\tilde{\sigma}) \cdot \text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbb{u}^t))\}^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\
 2065 &\stackrel{\textcircled{5}}{\leq} e^t + e^{t-1} + 4\dot{K}\tilde{c} \cdot \{\tilde{c}(1-\tilde{\sigma}) \cdot \beta^t K(e^t + e^{t-1})\}^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\
 2066 &\stackrel{\textcircled{6}}{=} d^{t-2} - d^t + 4\dot{K}\tilde{c} \cdot \{\tilde{c}(1-\tilde{\sigma}) \cdot \beta^t K(d^{t-2} - d^t)\}^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\
 2067 &= d^{t-2} - d^t + \underbrace{4\dot{K}\tilde{c} \cdot [\tilde{c}(1-\tilde{\sigma})K]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}}_{\triangleq \ddot{K}} \cdot \{(\beta^t(d^{t-2} - d^t))^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}\},
 \end{aligned}$$

2075 where step ① uses  $\sum_{i=t}^\infty e^{i+1} \leq e^t + e^{t-1} + 4\dot{K}\varphi^t$ , as shown in Theorem 5.6(b); step ② uses  
 2076 the definitions that  $\varphi^t \triangleq \varphi(s^t)$ ,  $s^t \triangleq \Theta(\mathbb{w}^t; \mathbb{u}^t) - \Theta(\mathbb{w}^\infty; \mathbb{u}^\infty)$ , and  $\varphi(s) = \tilde{c}s^{1-\tilde{\sigma}}$ ; step ③  
 2077 uses  $\varphi'(s) = \tilde{c}(1-\tilde{\sigma}) \cdot [s]^{-\tilde{\sigma}}$ , leading to  $[s^t]^{\tilde{\sigma}} = \tilde{c}(1-\tilde{\sigma}) \cdot \frac{1}{\varphi'(s^t)}$ ; step ④ uses Assumption 5.1  
 2078 that  $1 \leq \text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbb{u}^t)) \cdot \varphi'(s^t)$ ; step ⑤ uses  $\text{dist}(\mathbf{0}, \partial\Theta(\mathbb{w}^t; \mathbb{u}^t)) \leq K(e^t + e^{t-1})$  for both  
 2079 OADMM-EP and OADMM-RR, as shown in Lemma 5.4; step ⑥ uses the fact that  $e^t = d^{t-1} - d^t$ ,  
 2080 which implies:

$$2081 \quad e^t + e^{t-1} = (d^{t-1} - d^t) + (d^{t-2} - d^{t-1}) = d^{t-2} - d^t.$$

□

#### 2086 D.4 PROOF OF THEOREM 5.9

2087 *Proof.* Using Lemma 5.8(b), we have:

$$2088 \quad d^t \leq d^{t-2} - d^t + \ddot{K} \cdot \{(\beta^t(d^{t-2} - d^t))^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}\}. \quad (57)$$

2091 We consider two cases for Inequality (57).

2092 **Part (a).**  $\tilde{\sigma} \in (\frac{1}{4}, \frac{1}{2}]$ . We define  $u \triangleq \frac{p(1-\tilde{\sigma})}{\tilde{\sigma}} \in [\frac{1}{3}, 1)$ , where  $p = \frac{1}{3}$  is a fixed constant.

2093 We define  $\tilde{\beta}^t \triangleq \ddot{K}(\beta^t)^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}$ . We define  $t' \triangleq \{i \mid d^{i-2} - d^i \leq 1\}$ .

2094 For all  $t \geq t'$ , we have from Inequality (57):

$$\begin{aligned}
 2095 \quad d^t &\leq d^{t-2} - d^t + (d^{t-2} - d^t)^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \cdot \underbrace{\ddot{K}(\beta^t)^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}}_{\triangleq \tilde{\beta}^t} \\
 2096 &\stackrel{\textcircled{1}}{\leq} d^{t-2} - d^t + (d^{t-2} - d^t) \cdot \tilde{\beta}^t \\
 2097 &\leq d^{t-2} \cdot \frac{\tilde{\beta}^t + 1}{\tilde{\beta}^t + 2}, \quad (58)
 \end{aligned}$$

2098 where step ① uses the fact that  $[\Delta^{(1-\tilde{\sigma})/\tilde{\sigma}}]/\Delta = \Delta^{(1-2\tilde{\sigma})/\tilde{\sigma}} = \Delta^{(1/\tilde{\sigma}-2)} \leq \Delta^0 = 1$  for all  
 2099  $\Delta = d^{t-2} - d^t \in [0, 1]$  and  $\tilde{\sigma} \in (0, \frac{1}{2}]$ .

Furthermore, We derive:

$$\sum_{t=1}^T (\tilde{\beta}^t)^{-1} \stackrel{\textcircled{1}}{=} \mathcal{O} \left( \sum_{t=1}^T [t^p]^{-\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \right) \stackrel{\textcircled{2}}{=} \mathcal{O}(\sum_{t=1}^T t^{-u}) \stackrel{\textcircled{3}}{\geq} \mathcal{O}(T^{1-u}),$$

where step ① uses  $\tilde{\beta}^t \triangleq \ddot{K}(\beta^t)^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}$  and  $\beta^t \triangleq \beta^0(1 + \xi t^p) = \mathcal{O}(t^p)$ ; step ② uses the definition of  $u$ ; step ③ uses Lemma A.9 that:  $\sum_{t=1}^T t^{-u} \geq (1-u)T^{1-u} = \mathcal{O}(T^{1-u})$  for all  $u \in (0, 1)$ .

Applying Lemma Lemma A.12 with  $a = 1 - u$ , we have:

$$d^T \leq \mathcal{O} \left( \frac{1}{\exp(T^{1-u})} \right).$$

**Part (b).**  $\tilde{\sigma} \in (\frac{1}{2}, 1)$ . We define  $w \triangleq \frac{1-\tilde{\sigma}}{\tilde{\sigma}} \in (0, 1)$ , and  $\tau \triangleq 1/w - 1 \in (0, \infty)$ .

We define  $\tilde{\beta}^t = \dot{K}^{1/w} \beta^t$ , where  $\dot{K} \triangleq \ddot{K} + R^{1-w}(\beta^0)^{-w}$ , and  $R \triangleq d^0$ .

Notably, we have:  $d^{t-2} - d^t \leq d^0 \triangleq R$  for all  $t \geq 2$ .

For all  $t \geq 2$ , we have from Inequality (57):

$$\begin{aligned} d^t &\leq d^{t-2} - d^t + \ddot{K}(\beta^t)^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} (d^{t-2} - d^t)^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\ &\stackrel{\textcircled{1}}{=} \ddot{K} \{\beta^t (d^{t-2} - d^t)\}^w + d^{t-2} - d^t \\ &\stackrel{\textcircled{2}}{\leq} \dot{K} \{\beta^t (d^{t-2} - d^t)\}^w + (d^{t-2} - d^t)^w \cdot R^{1-w} \\ &\stackrel{\textcircled{3}}{\leq} \ddot{K} \{\beta^t (d^{t-2} - d^t)\}^w + (d^{t-2} - d^t)^w \cdot R^{1-w} \cdot (\frac{\beta^t}{\beta^0})^w \\ &= \{\beta^t (d^{t-2} - d^t)\}^w \cdot \underbrace{(\ddot{K} + R^{1-w} \cdot (\beta^0)^{-w})}_{\triangleq \dot{K}}, \end{aligned}$$

where step ① uses the the definition of  $w$ ; step ② uses the fact that  $\max_{x \in (0, R]} \frac{x}{x^w} \leq R^{1-w}$  if  $w \in (0, 1)$  and  $R > 0$ ; step ③ uses  $\beta^0 \leq \beta^t$  and  $w \in (0, 1)$ . We further obtain:

$$\underbrace{[d^t]^{1/w}}_{=[d^t]^{\tau+1}} \leq (d^{t-2} - d^t) \cdot \underbrace{\beta^t \dot{K}^{1/w}}_{\triangleq \tilde{\beta}^t}.$$

Additionally, we have:

$$\sum_{t=1}^T (1/\tilde{\beta}^t) \stackrel{\textcircled{1}}{=} \mathcal{O}(\sum_{t=1}^T (1/\beta^t)) \stackrel{\textcircled{2}}{=} \mathcal{O}(\sum_{t=1}^T t^{-p}) \stackrel{\textcircled{3}}{\geq} \mathcal{O}(T^{1-p}),$$

where step ① uses  $\tilde{\beta}^t = \dot{K}^{1/w} \beta^t$ ; step ② uses  $\beta^t \triangleq \beta^0(1 + \xi t^p) = \mathcal{O}(t^p)$ ; step ③ uses Lemma A.9 that:  $\sum_{t=1}^T t^{-u} \geq (1-p)T^{1-u} = \mathcal{O}(T^{1-p})$  for all  $p \in (0, 1)$ .

Applying Lemma A.13 with  $a = 1 - p$ , we have:

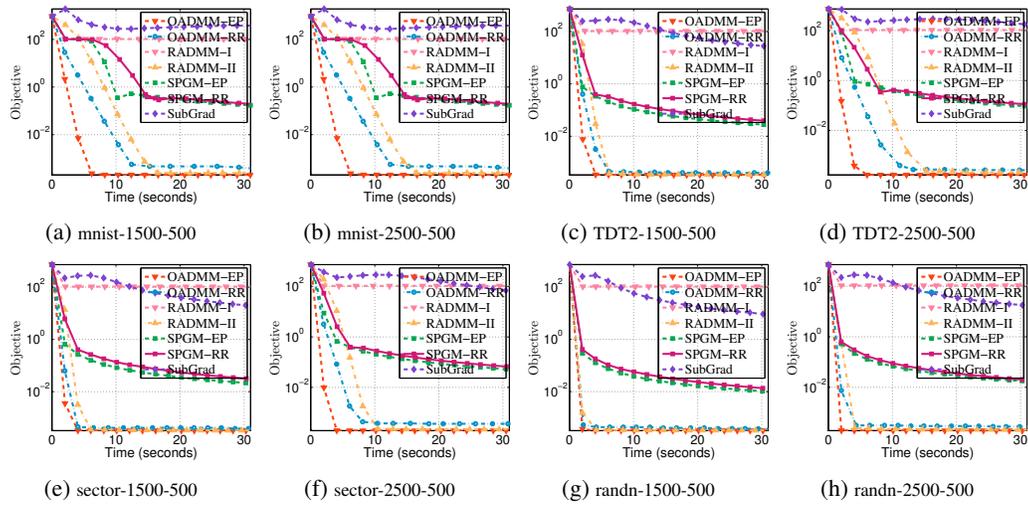
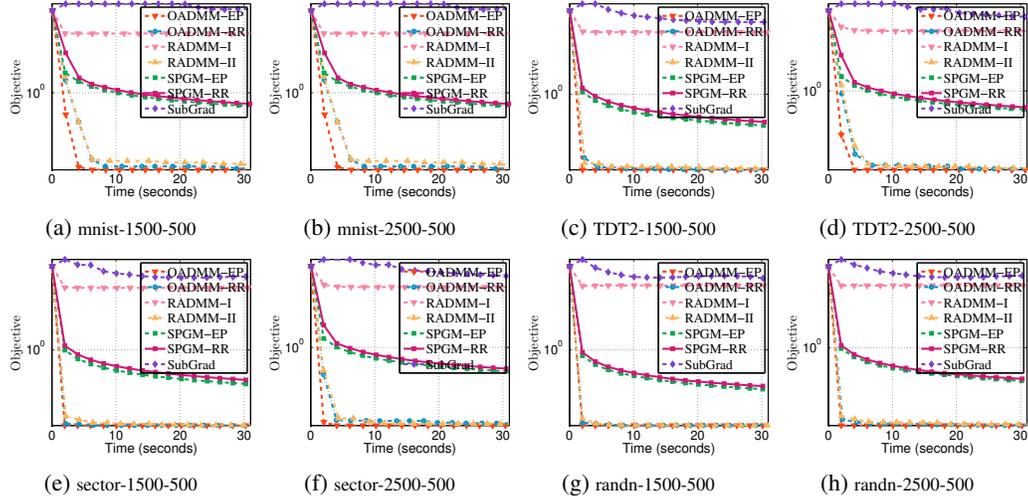
$$d^T \leq \mathcal{O}(1/(T^{(1-p)/\tau})).$$

**Part (c).** Finally, using the fact  $\|\mathbf{X}^T - \mathbf{X}^\infty\|_F \leq \mathcal{O}(d^T)$  as shown in Lemma D.3(b), we finish the proof of this theorem.

□

## E ADDITIONAL EXPERIMENTS DETAILS AND RESULTS

► **Datasets.** In our experiments, we utilize several datasets comprising both randomly generated and publicly available real-world data. These datasets are structured as data matrices  $\mathbf{D} \in \mathbb{R}^{\dot{m} \times \dot{d}}$ . They are denoted as follows: ‘mnist- $\dot{m}$ - $\dot{d}$ ’, ‘TDT2- $\dot{m}$ - $\dot{d}$ ’, ‘sector- $\dot{m}$ - $\dot{d}$ ’, and ‘randn- $\dot{m}$ - $\dot{d}$ ’, where  $\text{randn}(m, n)$  generates a standard Gaussian random matrix of size  $m \times n$ . The construction of  $\mathbf{D} \in \mathbb{R}^{\dot{m} \times \dot{d}}$  involves randomly selecting  $\dot{m}$  examples and  $\dot{d}$  dimensions from the original

Figure 3: The convergence curve of the compared methods with  $\rho = 10$ .Figure 4: The convergence curve of the compared methods with  $\rho = 100$ .

real-world dataset, sourced from <http://www.cad.zju.edu.cn/home/dengcai/Data/TextData.html> and <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>. Subsequently, we normalize each column of  $\mathbf{D}$  to possess a unit norm and center the data by subtracting the mean, denoted as  $\mathbf{D} \leftarrow \mathbf{D} - \mathbf{1}\mathbf{1}^T\mathbf{D}$ .

► **Additional experiment Results.** We present additional experimental results in Figures 3, 4, and 5. The figures demonstrate that the proposed OADMM method generally outperforms the other methods, with OADMM-EP surpassing OADMM-RR. These results reinforce our previous conclusions.

2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267

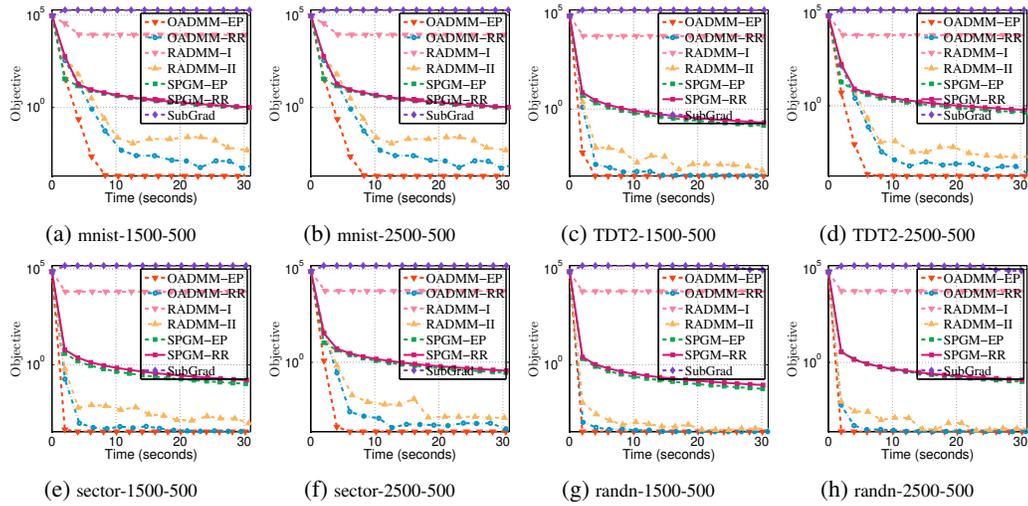


Figure 5: The convergence curve of the compared methods with  $\rho = 1000$ .