

ADMM FOR NONSMOOTH COMPOSITE OPTIMIZATION UNDER ORTHOGONALITY CONSTRAINTS

Ganzhao Yuan

Peng Cheng Laboratory, China
yuangzh@pcl.ac.cn

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 polynomial convergence rate or super-exponential 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 inexact 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:

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

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

(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: $\mathbf{X}^{t+1} \approx \arg \min_{\mathbf{X}} F_2(\mathbf{X}, \mathbf{Y}^t) + \langle \mathbf{Z}^t, \mathbf{X} - \mathbf{Y}^t \rangle + \frac{\beta}{2} \|\mathbf{X} - \mathbf{Y}\|_F^2$, $\mathbf{Y}^{t+1} \in \arg \min_{\mathbf{Y}} F_2(\mathbf{X}^{t+1}, \mathbf{Y}) + \langle \mathbf{Z}^t, \mathbf{X}^{t+1} - \mathbf{Y} \rangle + \frac{\beta}{2} \|\mathbf{X}^{t+1} - \mathbf{Y}\|_F^2$, and $\mathbf{Z}^{t+1} = \mathbf{Z}^t + \beta(\mathbf{X}^{t+1} - \mathbf{Y}^{t+1})$, where β 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	$\sigma \in [1, 2], \alpha > 0$	$\mathcal{O}(\epsilon^{-3})$	✓ Theorem 5.9
OADM-RR[ours]	weakly convex $h(\cdot)$	convex	$\sigma \in [1, 2], \text{MBB}$	$\mathcal{O}(\epsilon^{-3})$	✓ Theorem 5.9

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 polynomial convergence rate or super-exponential 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 $\mathcal{A}^\top(\cdot)$ to denote the adjoint operator of $\mathcal{A}(\cdot)$ with $\langle \mathcal{A}(\mathbf{X}), \mathbf{z} \rangle = \langle \mathbf{X}, \mathcal{A}^\top(\mathbf{z}) \rangle$ for all $\mathbf{X} \in \mathbb{R}^{n \times r}$ and $\mathbf{z} \in \mathbb{R}^m$. We define $\bar{\mathcal{A}} \triangleq \max_{\mathbf{V}} \|\mathcal{A}(\mathbf{V})\|_F / \|\mathbf{V}\|_F$. We use $\iota_{\mathcal{M}}(\mathbf{X})$ to denote the indicator function of orthogonality constraints. 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. ((Böhm & Wright, 2021)) 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. 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 $\mathbb{P}_\mu(\mathbf{y}) \triangleq \arg \min_{\tilde{\mathbf{y}}} h(\tilde{\mathbf{y}}) + \frac{1}{2\mu} \|\tilde{\mathbf{y}} - \mathbf{y}\|_2^2$. 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 (μ^{-1}) -smooth and (μ^{-1}) -weakly convex.

Lemma 2.3. (Proof in Appendix B.2) 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.3) 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.4) 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 $\check{\mathbf{y}} = \arg \min_{\mathbf{y}} h(\mathbf{y}) + \frac{\beta}{2(1 + \mu\beta)} \|\mathbf{y} - \mathbf{b}\|_2^2 = \mathbb{P}(\mathbf{b}; \mu + 1/\beta)$. (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 ϵ -stationary point of Problem (1) as follows.

Definition 2.7. (First-Order Optimality Conditions, (Chen et al., 2020; Li et al., 2022)) The solution $(\tilde{\mathbf{X}}, \tilde{\mathbf{y}}, \tilde{\mathbf{z}})$ with $\tilde{\mathbf{X}} \in \mathcal{M}$ is called an ϵ -stationary point of Problem (1) if: $\text{Crit}(\tilde{\mathbf{X}}, \tilde{\mathbf{y}}, \tilde{\mathbf{z}}) \leq \epsilon$, 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 & Malick, 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 \mathbf{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 $\{\dot{k}, \ddot{k}\}$ such that $\|\text{Retr}_{\mathbf{X}}(\Delta) - \mathbf{X}\|_F \leq \dot{k} \|\Delta\|_F$, and $\|\text{Retr}_{\mathbf{X}}(\Delta) - \mathbf{X} - \Delta\|_F \leq \frac{1}{2} \ddot{k} \|\Delta\|_F^2$.

Furthermore, we present the following three insightful lemmas.

Lemma 2.11. (Proof in Appendix B.5) 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.6) For any $\rho > 0$, $\mathbf{G} \in \mathbb{R}^{n \times r}$, 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) $\max(1, 2\rho) \cdot \langle \mathbf{G}, \mathbb{G}_\rho \rangle \geq \|\mathbb{G}_\rho\|_F^2 \geq \min(1, \rho^2) \|\mathbb{G}_{1/2}\|_F^2$. (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.7) 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: $\text{dist}(\mathbf{0}, \partial I_{\mathcal{M}}(\mathbf{X}) + \nabla f(\mathbf{X})) \leq \|\nabla f(\mathbf{X}) - \mathbf{X} \nabla f(\mathbf{X})^\top \mathbf{X}\|_F$.

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 ρ 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:

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

where $\mu \rightarrow 0$, and $h_\mu(\mathbf{y})$ is the Moreau Envelope of $h(\mathbf{y})$. Importantly, $h_\mu(\mathbf{y})$ is (μ^{-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:

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

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

In simple terms, OADMM updates are performed by minimizing the augmented Lagrangian function $\mathcal{L}(\mathbf{X}, \mathbf{y}, \mathbf{z}, \beta)$ 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. (i) 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). (ii) To accelerate primal convergence in OADMM-EP, we incorporate a Nesterov extrapolation strategy with parameter $\alpha \in (0, 1)$. (iii) 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 curvature¹. The parameters $\{\gamma, \delta\}$ represent the decay rate and sufficient decrease parameter, commonly used in line search procedures (Chen et al., 2020). (iv) The \mathbf{X} -subproblem is solved as: $\mathbf{X}^{t+1} = \arg \min_{\mathbf{X} \in \mathcal{M}} \|\mathbf{X} - \mathbf{X}'\|_F^2 = \dot{\mathbf{U}} \dot{\mathbf{V}}^T$, where $\mathbf{X}' = \mathbf{X}_c^t - \mathbf{G}^t/(\theta \ell(\beta^t))$, and $\dot{\mathbf{U}} \text{diag}(\dot{\mathbf{x}}) \dot{\mathbf{V}}^T = \mathbf{X}'$ is the using singular value decomposition of \mathbf{X}' . (v) For practical implementation, we recommend the following default parameters: $p = 1/3$, $\theta = 1.01$, $\sigma = 1.1$, $\rho = 1$, $\gamma = 1/2$, $\delta = 10^{-3}$, $\xi = 1$, $\alpha = \frac{\theta-1}{(\theta+1)(\xi+2)} - 10^{-12}$.

4 ORACLE COMPLEXITY

This section details the oracle complexity of Algorithm 1.

¹Following (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 $p, \xi \in (0, 1), \theta \in (1, \infty), \sigma \in [1, 2)$.

Choose $\tau \in [\frac{4}{2-\sigma}, \infty)$. Choose β^0 with $\beta^0 \geq 2\tau W_h$.

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 $\beta^t = \beta^0(1 + \xi t^p)$.

S2) Update the primal variable \mathbf{X} :

if OADMM-EP **then**

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

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

end

if OADMM-RR **then**

Set $\mathbf{G}^t \in \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t) - \partial g(\mathbf{X}^t)$. 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 $\mathbf{X}^{t+1} = \text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t)$.

Here, $\eta^t \triangleq \frac{b^t \gamma^j}{\beta^t}$, and $j \in \{0, 1, \dots\}$ is the smallest integer such that:

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

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 $\mu^t = \tau/\beta^t$,

$\mathbf{b} \triangleq \mathbf{y}^t - \frac{1}{\beta^t} \nabla_{\mathbf{y}} \mathcal{S}(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t, \beta^t)$. It can be solved using Lemma 2.5 as:

$\mathbf{y}^{t+1} = \frac{\check{\mathbf{y}}^{t+1} + \mu^t \beta^t \mathbf{b}}{1 + \mu^t \beta^t}$, where $\check{\mathbf{y}}^{t+1} = \mathbb{P}(\mathbf{b}; \mu^t + 1/\beta^t)$.

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

Notations. We define $\varepsilon_z = \frac{1}{4}, \varepsilon_\beta > 0, \varepsilon_y \triangleq \frac{2-\sigma}{8}, \ddot{\sigma} \triangleq \frac{12\sigma^2}{p(2-\sigma)^2} C_h^2, c \triangleq \varepsilon_\beta + \tau C_h^2 + \frac{2}{\sigma} \ddot{\sigma}$. We define a sequence associated with the potential function (or Lyapunov function) for all $t \geq 1$, as follows:

$$\Theta^t \triangleq L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) + c/\beta^t + \mathbb{P}^t + \mathbb{D}^t, \quad (6)$$

where $\mathbb{P}^t \triangleq \frac{\alpha(\theta+1)\ell(\beta^t)}{2} \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F^2$, $\mathbb{D}^t \triangleq 2\beta^{t-1} \frac{\sigma-1}{2-\sigma} \|\sigma(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)\|_2^2$. We define $\mathcal{B}_t \triangleq \sqrt{(\frac{1}{\beta^{t-1}} - \frac{1}{\beta^t}) \frac{1}{\beta^{t-1}}}$, $\mathcal{Y}_t \triangleq \|\mathbf{y}^t - \mathbf{y}^{t-1}\|$, and $\mathcal{Z}_t \triangleq \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|$. For OADMM-EP, we define $\mathcal{X}_t \triangleq \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$, and for OADMM-RR, $\mathcal{X}_t \triangleq \|\frac{1}{\beta^t} \mathbb{G}_{1/2}^{t-1}\|_F$.

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)** $\forall t \geq 1, \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \leq \frac{\sigma-1}{2-\sigma} (\|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2 - \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2) + \dot{\sigma}(\beta^t)^2 \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \ddot{\sigma}(\frac{\beta^0}{\beta^t} - \frac{\beta^0}{\beta^{t+1}})$, where $\dot{\sigma} \triangleq \frac{2\sigma^2}{(2-\sigma)^2} \frac{1}{\tau^2}$.

Remark 4.2. For OADMM-RR, we set $\alpha = 0$, resulting in $\mathbb{P}^t = 0$ for all t . 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 $\varepsilon_z \beta^t \mathcal{Z}_{t+1}^2 + \varepsilon_y \beta^t \mathcal{Y}_{t+1}^2 + \varepsilon_\beta \beta^t \mathcal{B}_{t+1}^2 + L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^{t+1}, \beta^{t+1}) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) + c/\beta^{t+1} - c/\beta^t + \mathbb{D}^{t+1} - \mathbb{D}^t \leq \mathfrak{X}$, where $\mathfrak{X} \triangleq L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t, \beta^t) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t)$.

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 $\mathfrak{X} \leq -\varepsilon_x \beta^t \mathcal{X}_{t+1}^2 + \mathbb{P}^t - \mathbb{P}^{t+1}$.

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 define $\mathcal{X}_t \triangleq \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$. We have:

- (a) $\beta^t \{\varepsilon_\beta \mathcal{B}_{t+1}^2 + \varepsilon_z \mathcal{Z}_{t+1}^2 + \varepsilon_y \mathcal{Y}_{t+1}^2 + \varepsilon_x \mathcal{X}_{t+1}^2\} \leq \Theta^t - \Theta^{t+1}$.
- (b) $\frac{1}{T} \sum_{t=1}^T \beta^t [\mathcal{B}_{t+1} + \mathcal{Z}_{t+1} + \mathcal{Y}_{t+1} + \mathcal{X}_{t+1}] \leq \mathcal{O}(T^{(p-1)/2})$.

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: $\frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{X}^{t+1}, \check{\mathbf{y}}^{t+1}, \mathbf{z}^{t+1}) \leq \mathcal{O}(T^{(p-1)/2}) + \mathcal{O}(T^{-p}) = \mathcal{O}(T^{-1/3})$. In other words, there exists $\bar{t} \leq T$ such that: $\text{Crit}(\mathbf{X}^{\bar{t}+1}, \check{\mathbf{y}}^{\bar{t}+1}, \mathbf{z}^{\bar{t}+1}) \leq \epsilon$, provided that $T \geq \mathcal{O}(1/\epsilon^3)$.

Remark 4.9. (i) We notice that $\frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{X}^{t+1}, \check{\mathbf{y}}^{t+1}, \mathbf{z}^{t+1}) \leq \mathcal{O}(T^{(p-1)/2}) + \mathcal{O}(T^{-p})$. Minimizing the worse-case complexity of the right-hand side of this inequality with respect to p yields: $\arg \min_{p \in (0,1)} \max((p-1)/2, -p) = 1/3$. Thus, setting $p = 1/3$ achieves the optimal trade-off between the two terms, leading to the best complexity bounds. (ii) 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: $\mathfrak{X} \leq -\frac{\varepsilon_x}{\beta^t} \|\mathbb{G}_{1/2}^t\|_F^2$. Here, \bar{g} is a constant that $\|\mathbb{G}^t\|_F \leq \bar{g}$, $\{\bar{k}, \bar{b}\}$ are defined in Lemma 2.10, and $\{\rho, \gamma, \delta, \bar{b}, \bar{b}\}$ are defined in Algorithm 1.

Remark 4.11. By Lemma 4.10(a), since $\bar{\gamma}$ is a universal constant and γ^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 define $\mathcal{X}_t \triangleq \|\frac{1}{\beta^t} \mathbb{G}_{1/2}^{t-1}\|_F$. We have:

- (a) $\beta^t \{\varepsilon_\beta \mathcal{B}_{t+1}^2 + \varepsilon_z \mathcal{Z}_{t+1}^2 + \varepsilon_y \mathcal{Y}_{t+1}^2 + \varepsilon_x \mathcal{X}_{t+1}^2\} \leq \Theta^t - \Theta^{t+1}$.
- (b) $\frac{1}{T} \sum_{t=1}^T \beta^t [\mathcal{B}_{t+1} + \mathcal{Z}_{t+1} + \mathcal{Y}_{t+1} + \mathcal{X}_{t+1}] \leq \mathcal{O}(T^{(p-1)/2})$.

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: $\frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{X}^{t+1}, \check{\mathbf{y}}^{t+1}, \mathbf{z}^{t+1}) \leq \mathcal{O}(T^{(p-1)/2}) + \mathcal{O}(T^{-p}) = \mathcal{O}(T^{-1/3})$. In other words, there exists $\bar{t} \leq T$ such that: $\text{Crit}(\mathbf{X}^{\bar{t}+1}, \check{\mathbf{y}}^{\bar{t}+1}, \mathbf{z}^{\bar{t}+1}) \leq \epsilon$, 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).

For simplicity, we only consider the case where $\alpha = 0$, $\sigma = 1$, and $g(\mathbf{X}) = 0$. We only focus on $p = 1/3$, as it gives the best oracle complexity.

We define the Lyapunov function as: $\Theta(\mathbf{X}, \mathbf{y}, \mathbf{z}, \beta) \triangleq L(\mathbf{X}, \mathbf{y}, \mathbf{z}, \beta) + c/\beta$. We define $\mathcal{W} \triangleq \{\mathbf{X}, \mathbf{y}, \mathbf{z}, \beta\}$, $\mathcal{W}^t \triangleq \{\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t\}$. Consequently, $\Theta^t = \Theta(\mathcal{W}^t)$. We denote \mathcal{W}^∞ as a limiting point of Algorithm 1. We make the following additional assumption.

Assumption 5.1. *The function $\Theta(\mathcal{W})$ is a KL function w.r.t. \mathcal{W} .*

Proposition 5.2. (Kurdyka-Łojasiewicz Inequality (Attouch et al., 2010)). *Consider a semi-algebraic function $\Theta(\mathcal{W})$ with $\mathcal{W} \in \text{dom}(\Theta)$. There exist constants $\tilde{\eta} \in (0, +\infty)$, $\tilde{\sigma} \in [0, 1)$, a neighborhood Υ of \mathcal{W}^∞ , 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 $\mathcal{W} \in \Upsilon$ satisfying $\Theta(\mathcal{W}) - \Theta(\mathcal{W}^\infty) \in (0, \tilde{\eta})$, it holds that: $\varphi'(\Theta(\mathcal{W}) - \Theta(\mathcal{W}^\infty)) \cdot \text{dist}(\mathbf{0}, \partial\Theta(\mathcal{W})) \geq 1$.*

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) *For both OADMM-EP and OADMM-RR, there exists a constant $K > 0$ such that: $\text{dist}(\mathbf{0}, \partial\Theta(\mathcal{W}^t)) \leq K\beta^t(\mathcal{B}_t + \mathcal{X}_t + \mathcal{Y}_t + \mathcal{Z}_t)$.*

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 $d^t \triangleq \sum_{i=t}^\infty e^{i+1}$, where $e^t \triangleq \mathcal{B}_t + \mathcal{X}_t + \mathcal{Y}_t + \mathcal{Z}_t$. We define $\varphi^t \triangleq \varphi(\Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^\infty))$, where $\varphi(\cdot)$ is the desingularization function defined in Assumption 5.1. We have the following results for both OADMM-EP and OADMM-RR.*

- (a) $(e^{t+1})^2 \leq (\varphi^t - \varphi^{t+1})K'e^t$, where $K' = \frac{4K}{\min(\varepsilon_x, \varepsilon_y, \varepsilon_z, \varepsilon_\beta)}$, and K is defined in Lemma 5.4.
- (b) *It holds that $\forall t \geq 1$, $d^t \leq e^t + 2K'\varphi^t$. The sequence $\{\mathcal{W}^t\}_{t=1}^\infty$ has the finite length property that $d^t \leq e^1 + 2K'\varphi^1 < +\infty$.*

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) *For both OADMM-EP and OADMM-RR, we have:*

- (a) *There exists a constant ϖ such that $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F \leq \varpi d^t$.*
- (b) *We have $d^t \leq d^{t-1} - d^t + \nu[\beta^t(d^{t-1} - d^t)]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}$, where ν is some universal constant.*

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 $p = 1/3$. We have:

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

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 polynomial convergence rate when $\tilde{\sigma} \in (0, 1)$, and converges with a super-exponential rate when $\tilde{\sigma} \in (\frac{1}{4}, \frac{1}{2}]$ for the gap $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F$. (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:

$$\min_{\mathbf{X} \in \mathbb{R}^{n \times r}} \frac{1}{2\tilde{m}} \|\mathbf{X}\mathbf{X}^\top \mathbf{D} - \mathbf{D}\|_F^2 + \dot{\rho}(\|\mathbf{X}\|_1 - \|\mathbf{X}\|_{[k]}), \text{ s.t. } \mathbf{X}^\top \mathbf{X} = \mathbf{I}_r, \quad (7)$$

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 ℓ_1 norm of the k largest (in magnitude) elements of the matrix \mathbf{X} . Here, we consider the DC ℓ_1 -largest- k function (Gotoh et al., 2018) to induce sparsity in the solution. One advantage of this model is that when $\dot{\rho}$ is sufficient large, we have $\|\mathbf{X}\|_1 \approx \|\mathbf{X}\|_{[k]}$, leading to a k -sparsity solution \mathbf{X} . Problem (7) coincides with the optimization model in Problem (1), where $f(\mathbf{X}) = \frac{1}{2\tilde{m}} \|\mathbf{X}\mathbf{X}^\top \mathbf{D} - \mathbf{D}\|_F^2$, $f(\mathbf{X}) = \dot{\rho}\|\mathbf{X}\|_{[k]}$, and $h(\mathcal{A}(\mathbf{X})) = \dot{\rho}\|\mathbf{X}\|_1$.

► **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 $\dot{\rho} \in \{10, 50, 100, 500, 1000\}$. For OADMM, default parameters are used, with $\beta^0 = 10\dot{\rho}$ and corresponding values $\xi = 0.5$ for each $\dot{\rho}$. 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 $\dot{\rho}$. 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.

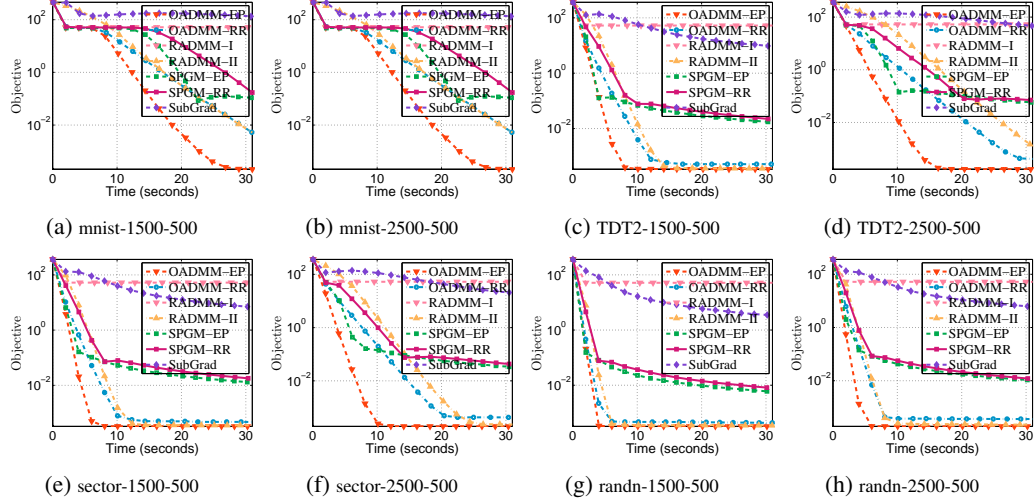


Figure 1: The convergence curve of the compared methods with $\dot{\rho} = 50$.

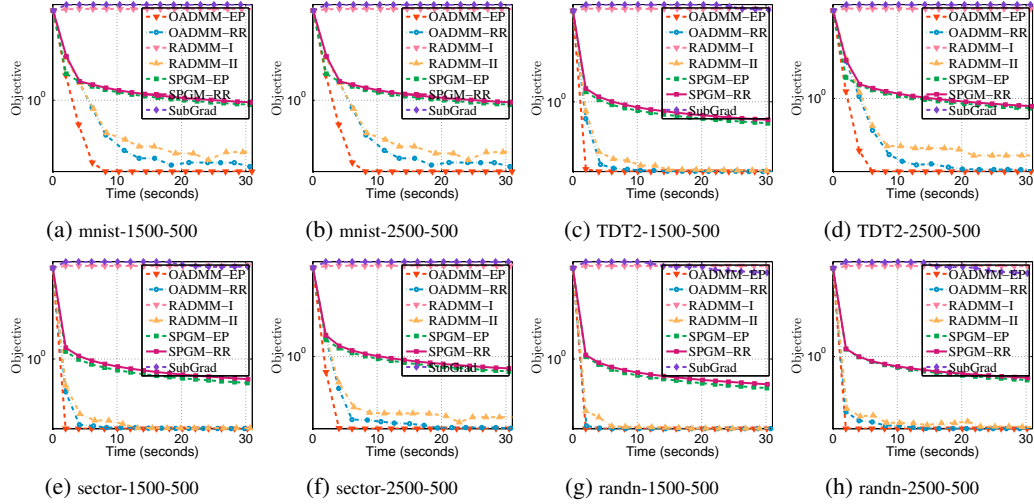


Figure 2: The convergence curve of the compared methods with $\dot{\rho} = 500$.

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 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

- Traian E Abrudan, Jan Eriksson, and Visa Koivunen. Steepest descent algorithms for optimization under unitary matrix constraint. *IEEE Transactions on Signal Processing*, 56(3):1134–1147, 2008.
- P-A Absil and Jérôme Malick. Projection-like retractions on matrix manifolds. *SIAM Journal on Optimization*, 22(1):135–158, 2012.
- P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2008a.
- Pierre-Antoine Absil, Robert E. Mahony, and Rodolphe Sepulchre. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2008b.
- Hedy Attouch, Jérôme Bolte, Patrick Redont, and Antoine Soubeyran. Proximal alternating minimization and projection methods for nonconvex problems: An approach based on the kurdyka-lojasiewicz inequality. *Mathematics of Operations Research*, 35(2):438–457, 2010.
- Nitin Bansal, Xiaohan Chen, and Zhangyang Wang. Can we gain more from orthogonality regularizations in training deep networks? *Advances in Neural Information Processing Systems*, 31, 2018.
- Amir Beck and Israel Rosset. A dynamic smoothing technique for a class of nonsmooth optimization problems on manifolds. *SIAM Journal on Optimization*, 33(3):1473–1493, 2023.
- Radu Ioan Boț and Dang-Khoa Nguyen. The proximal alternating direction method of multipliers in the nonconvex setting: convergence analysis and rates. *Mathematics of Operations Research*, 45(2):682–712, 2020.
- Radu Ioan Boț, Erno Robert Csetnek, and Dang-Khoa Nguyen. A proximal minimization algorithm for structured nonconvex and nonsmooth problems. *SIAM Journal on Optimization*, 29(2):1300–1328, 2019. doi: 10.1137/18M1190689.
- Axel Böhm and Stephen J. Wright. Variable smoothing for weakly convex composite functions. *Journal of Optimization Theory and Applications*, 188(3):628–649, 2021.
- Jérôme Bolte, Shoham Sabach, and Marc Teboulle. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1-2):459–494, 2014.
- Nicolas Boumal, Pierre-Antoine Absil, and Coralía Cartis. Global rates of convergence for nonconvex optimization on manifolds. *IMA Journal of Numerical Analysis*, 39(1):1–33, 2019.
- Shixiang Chen, Shiqian Ma, Anthony Man-Cho So, and Tong Zhang. Proximal gradient method for nonsmooth optimization over the stiefel manifold. *SIAM Journal on Optimization*, 30(1): 210–239, 2020.
- Weiqiang Chen, Hui Ji, and Yanfei You. An augmented lagrangian method for ℓ_1 -regularized optimization problems with orthogonality constraints. *SIAM Journal on Scientific Computing*, 38(4): B570–B592, 2016.
- Xiaojun Chen. Smoothing methods for nonsmooth, nonconvex minimization. *Mathematical programming*, 134(1):71–99, 2012.
- Minhyung Cho and Jaehyung Lee. Riemannian approach to batch normalization. *Advances in Neural Information Processing Systems*, 30, 2017.
- Michael Cogswell, Faruk Ahmed, Ross B. Girshick, Larry Zitnick, and Dhruv Batra. Reducing overfitting in deep networks by decorrelating representations. In Yoshua Bengio and Yann LeCun (eds.), *International Conference on Learning Representations (ICLR)*, 2016.
- Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019.

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

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

Jiawei Zhang, Peijun Xiao, Ruoyu Sun, and Zhi-Quan Luo. A single-loop smoothed gradient descent-ascent algorithm for nonconvex-concave min-max problems. In *Advances in Neural Information Processing Systems*, 2020a.

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

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 $\mathbf{X} = \sum_{ij} |\mathbf{X}_{ij}|$
- $\|\mathbf{X}\|_{[k]}$: ℓ_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} .
- $\iota_{\mathcal{M}}(\mathbf{x})$: Indicator function of \mathcal{M} with $\iota_{\mathcal{M}}(\mathbf{x}) = 0$ if $\mathbf{x} \in \mathcal{M}$ and otherwise $+\infty$.

We employ the following parameters in Algorithm 1.

- θ : proximal parameter
- τ : correlation coefficient between μ^t and β^t , such that $\mu^t \beta^t = \tau$
- σ : over-relaxation parameter with $\sigma \in [1, 2)$
- α : Nesterov extrapolation parameter with $\alpha \in [0, 1)$
- ρ : search descent parameter with $\rho \in (0, \infty)$
- γ : decay rate parameter in the line search procedure with $\gamma \in (0, 1)$
- δ : sufficient decrease parameter in the line search procedure with $\delta \in (0, \infty)$
- p : exponent parameter used in the penalty update rule with $p \in (0, 1)$
- ξ : 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 $\hat{\partial}F(\mathbf{X}) \triangleq \{\xi \in \mathbb{R}^{n \times r} : \lim_{\mathbf{Z} \rightarrow \mathbf{X}} \inf_{\mathbf{Z} \neq \mathbf{X}} \frac{F(\mathbf{Z}) - F(\mathbf{X}) - \langle \xi, \mathbf{Z} - \mathbf{X} \rangle}{\|\mathbf{Z} - \mathbf{X}\|_F} \geq 0\}$, while the limiting subdifferential of $F(\mathbf{X})$ at $\mathbf{X} \in \text{dom}(F)$ is denoted as $\partial F(\mathbf{X}) \triangleq \{\xi \in \mathbb{R}^{n \times r} : \exists \mathbf{X}^t \rightarrow \mathbf{X}, F(\mathbf{X}^t) \rightarrow F(\mathbf{X}), \xi^t \in \hat{\partial}F(\mathbf{X}^t) \rightarrow \xi, \forall t\}$. 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., $\partial F(\mathbf{X}) = \hat{\partial}F(\mathbf{X}) = \{\xi \in \mathbb{R}^{n \times r} : F(\mathbf{Z}) \geq F(\mathbf{X}) + \langle \xi, \mathbf{Z} - \mathbf{X} \rangle, \forall \mathbf{Z} \in \mathbb{R}^{n \times r}\}$. (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 $\mathcal{N}_{\mathcal{M}}(\mathbf{X}) = \partial \iota_{\mathcal{M}}(\mathbf{X}) = \{\mathbf{Z} \in \mathbb{R}^{n \times r} : \langle \mathbf{Z}, \mathbf{X} \rangle \geq \langle \mathbf{Z}, \mathbf{Y} \rangle, \forall \mathbf{Y} \in \mathcal{M}\}$. Moreover, the tangent and normal space to \mathcal{M} at $\mathbf{X} \in \mathcal{M}$ are respectively denoted as $\mathbf{T}_{\mathbf{X}}\mathcal{M}$ and $\mathbf{N}_{\mathbf{X}}\mathcal{M}$. We have: $\mathbf{T}_{\mathbf{X}}\mathcal{M} = \{\mathbf{Y} \in \mathbb{R}^{n \times r} | \mathcal{A}_{\mathbf{X}}(\mathbf{Y}) = \mathbf{0}\}$ and $\mathbf{N}_{\mathbf{X}}\mathcal{M} = \{2\mathbf{X}\Lambda | \Lambda = \Lambda^T, \Lambda \in \mathbb{R}^{r \times r}\}$, 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

We present a collection of useful lemmas, each of which is independent of context and methodology.

Lemma A.1. Assume $\mathbf{a}, \mathbf{b} \in \mathbb{R}^n$, and $\alpha \geq 0$. 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: $-\|\mathbf{a} - \alpha\mathbf{b}\|_2^2 = -\|\mathbf{a}\|_2^2 - \|\alpha\mathbf{b}\|_2^2 + 2\alpha\langle \mathbf{a}, \mathbf{b} \rangle \leq -\|\mathbf{a}\|_2^2 - \|\alpha\mathbf{b}\|_2^2 + 2\alpha \cdot (\frac{1}{2}\|\mathbf{a}\|_2^2 + \frac{1}{2}\|\mathbf{b}\|_2^2) = (\alpha - 1)\|\mathbf{a}\|_2^2 - (\alpha^2 - \alpha)\|\mathbf{b}\|_2^2. \quad \square$

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

Proof. We have: $\|\mathbf{a}^+\|_2^2 = \|\varrho\mathbf{a} + \mathbf{b}\|_2^2 = \|\varrho\mathbf{a} + (1 - \varrho) \cdot \frac{\mathbf{b}}{1-\varrho}\|_2^2 \leq \varrho\|\mathbf{a}\|_2^2 + (1 - \varrho) \cdot \|\frac{\mathbf{b}}{1-\varrho}\|_2^2 = \varrho\|\mathbf{a}\|_2^2 + \frac{1}{1-\varrho}\|\mathbf{b}\|_2^2$, where the inequality holds due to the convexity of $\|\cdot\|_2^2$. Subtracting $\varrho\|\mathbf{a}^+\|_2^2$ from both sides gives: $(1 - \varrho)\|\mathbf{a}^+\|_2^2 \leq \varrho(\|\mathbf{a}\|_2^2 - \|\mathbf{a}^+\|_2^2) + \frac{1}{1-\varrho}\|\mathbf{b}\|_2^2$. Dividing through by $(1 - \varrho)$ yields the resulting inequality. □

Lemma A.3. Assume $p \in (0, 1)$. We have:

- (a) $p(t + 1)^{p-1} \leq (t + 1)^p - t^p \leq pt^{p-1}$ for all integer $t \geq 1$.
- (b) $p(t + 1)^{p-1} \leq (t + 1)^p - t^p$ for all integer $t \geq 0$.

Proof. The proof of this lemma follows from the concavity of $h(x) = x^p$ for all $x \geq 0$ and $p \in (0, 1)$, and is omitted for brevity. □

Lemma A.4. Assume $p \in (0, 1)$. For all $t \geq 1$, we have $\frac{1}{1-p}(t+1)^{1-p} - \frac{1}{1-p} - (1-p)t^{1-p} \geq 0$.

Proof. We define $f(t) \triangleq \frac{1}{1-p}(t+1)^{1-p} - \frac{1}{1-p} - (1-p)t^{1-p}$, where $p \in (0, 1)$ and $t \geq 1$.

Part (a). We now show that $(1-p)^{1/p} \leq \frac{1}{\exp(1)}$. 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 $h(p) \triangleq (1-p)^{1/p}$ is a decreasing function on $p \in (0, 1)$, we have $h(p) \leq \lim_{p \rightarrow 0+} (1-p)^{1/p} = \frac{1}{\exp(1)}$.

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

Part (c). We now show that $f(t)$ is an increasing function. We have: $\nabla f(t) = (t+1)^{-p} - (1-p)^2 t^{-p} = (t+1)^{-p} \cdot (1 - (1-p)^2 (\frac{t+1}{t})^p) \stackrel{\textcircled{1}}{\geq} (t+1)^{-p} \cdot (1 - (1-p)^2 2^p) \stackrel{\textcircled{2}}{\geq} (t+1)^{-p} \cdot (1 - (\frac{2}{\exp(1)^2})^p) \stackrel{\textcircled{3}}{\geq} 0$, 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 $\frac{2}{\exp(1)^2} \approx 0.2707 < 1$.

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

Lemma A.5. 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 now 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.4.

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 $\nabla h(x) = x^{-p}$. □

Lemma A.6. Assume that $\mathbf{a}^t \leq \varrho \mathbf{a}^{t-1} + c$, where $\varrho \in [0, 1)$, $c \geq 0$, and $\{\mathbf{a}^i\}_{i=0}^\infty$ is a non-negative sequence. We have: $\mathbf{a}^t \leq \mathbf{a}^0 + \frac{c}{1-\varrho}$ for all $t \geq 0$.

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

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

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

Lemma A.7. For all $p \in (0, 1)$, it holds that $(p-3)2^p + (p+3) \leq 0$.

Proof. We define $f(p) = (p-3)2^p + (p+3) \leq 0$.

We define $g(p) \triangleq f'(p) = 2^p \cdot [1 + (p-3)\log(2)] + 1$.

We define $h(p) \triangleq f''(p) = 2^p \log(2) \cdot [2 + (p-3)\log(2)]$.

Part (a). Clearly, the equation $h(p) = 0$ has one unique solution $p^* \triangleq 3 - \frac{2}{\log(2)} \approx 0.1146$. Notably, $h(p) \leq 0$ for all $p \in (0, p^*]$, and $h(p) \geq 0$ for all $p \in [p^*, 1)$.

Part (b). We now show that the equation $g(p) = 0$ has a unique solution. Noting that $g(p)$ is decreasing on $(0, p^*)$ and increasing on $(p^*, 1)$, and observing that $g(p) < 0$ for all $(0, p^*)$ and $g(1) > 0$, we conclude, by the Mean Value Theorem, that $g(p) = 0$ has a unique solution.

Part (c). Given that $g(p) = 0$ has a unique solution, we conclude that $f(p)$ has exactly one critical point on $(0, 1)$. Since $f(0) = f(1) = 0$, $f'(0) < 0$, and $f'(1) > 0$, it follows that $f(p) \leq 0$ for all $p \in (0, 1)$. □

Lemma A.8. Let $\beta^t = \beta^0(1 + \xi t^p)$, where $t \geq 1$, $\beta^0 > 0$, $\xi, p \in (0, 1)$. For all integer $t \geq 1$, we have: $(\frac{1}{\beta^{t-1}} - \frac{1}{\beta^t}) / (\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}) \leq \frac{3}{p}$.

Proof. We derive:

$$\begin{aligned} \Psi &\triangleq (\frac{1}{\beta^{t-1}} - \frac{1}{\beta^t}) / (\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}) \\ &\stackrel{\textcircled{1}}{=} (\frac{1}{1+\xi(t-1)^p} - \frac{1}{1+\xi t^p}) / (\frac{1}{1+\xi t^p} - \frac{1}{1+\xi(t+1)^p}) \\ &= (\frac{\xi t^p - \xi(t-1)^p}{[1+\xi(t-1)^p][1+\xi t^p]}) / (\frac{\xi(t+1)^p - \xi t^p}{[1+\xi(t+1)^p][1+\xi t^p]}) \\ &= \frac{t^p - (t-1)^p}{(t+1)^p - t^p} \cdot \frac{1+\xi(t+1)^p}{1+\xi(t-1)^p}, \end{aligned} \tag{8}$$

where step ① uses the definition of β^t .

We now focus on Inequality (8). Case (i). When $t = 1$, we have:

$$\Psi = \frac{1}{2^p-1} \cdot \frac{1+\xi 2^p}{1} = \frac{2^p+1}{2^p-1} \stackrel{\textcircled{1}}{\leq} \frac{3}{p},$$

where step ① uses the fact that $\frac{2^p+1}{2^p-1} \leq \frac{3}{p}$ for all $p \in (0, 1)$, which is implied by Lemma A.7.

Case (ii). When $t \geq 2$, we have:

$$\begin{aligned} \Psi &= \frac{t^p - (t-1)^p}{(t+1)^p - t^p} \cdot \frac{1+\xi(t+1)^p}{1+\xi(t-1)^p} \\ &\stackrel{\textcircled{1}}{\leq} \frac{p(t-1)^{p-1}}{p(t+1)^{p-1}} \cdot \frac{1+\xi(t+1)^p}{1+\xi(t-1)^p} \\ &\stackrel{\textcircled{2}}{\leq} \frac{p(t-1)^{p-1}}{p(t+1)^{p-1}} \cdot \frac{(t+1)^p}{(t-1)^p} \\ &= \frac{t+1}{t-1} \leq \frac{3}{p}, \end{aligned}$$

where step ① uses Lemma A.3; step ② uses $\frac{a+b}{c+d} \leq \max(\frac{a}{c}, \frac{b}{d})$ for all $a, b, c, d > 0$. □

Lemma A.9. 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: $(\frac{\beta^t}{\beta^{t-1}} - 1)^2 \leq (\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}) \cdot \frac{6\beta^0}{p}$.

Proof. First, we derive:

$$\begin{aligned} \Psi &= (\frac{\beta^t}{\beta^{t-1}} - 1)^2 / (\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}) \\ &= \beta^t (\frac{\beta^t}{\beta^{t-1}} - 1) \cdot \{(\frac{1}{\beta^{t-1}} - \frac{1}{\beta^t}) / (\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}})\} \\ &\stackrel{\textcircled{1}}{\leq} \frac{3}{p} \beta^t (\frac{\beta^t}{\beta^{t-1}} - 1) \\ &= \frac{3}{p} \xi \beta^0 (1 + \xi t^p) (\frac{t^p - (t-1)^p}{1 + \xi(t-1)^p}), \end{aligned} \tag{9}$$

where step ① uses Lemma A.8. We now focus on Inequality (9). Case (i). When $t = 1$, we have:

$$\Psi \leq \frac{3}{p}\xi\beta^0(1+\xi) \stackrel{\textcircled{1}}{\leq} \frac{6}{p}\beta^0,$$

where step ① uses $\xi \in (0, 1)$.

Case (ii). When $t \geq 2$, we have:

$$\begin{aligned} \Psi &\leq \frac{3}{p}\xi\beta^0 \frac{1+\xi t^p}{1+\xi(t-1)^p} \cdot (t^p - (t-1)^p) \\ &\stackrel{\textcircled{1}}{\leq} \frac{3}{p}\xi\beta^0 \frac{1+\xi t^p}{1+\xi(t-1)^p} \cdot (p(t-1)^{p-1}) \\ &\stackrel{\textcircled{2}}{\leq} \frac{3}{p}\xi\beta^0 \frac{t^p}{(t-1)^p} \cdot (p(t-1)^{p-1}) \\ &= 3\xi\beta^0 \frac{t^p}{t-1} \\ &\stackrel{\textcircled{3}}{\leq} 3\xi\beta^0 \frac{t}{t-1} \\ &\stackrel{\textcircled{4}}{\leq} 6\beta^0/p, \end{aligned}$$

where step ① uses Lemma A.3; step ② uses $\frac{a+b}{c+d} \leq \max(\frac{a}{c}, \frac{b}{d})$ for all $a, b, c, d > 0$; step ③ uses $p \in (0, 1)$; step ④ uses $p, \xi \in (0, 1)$, and $\frac{t}{t-1} \leq 2$ for all $t \geq 2$. □

Lemma A.10. Assume $\beta^t = \beta^0(1 + \xi t^p)$, where $p \in [1/4, 1)$, and $t \geq 1$ is an integer. We have: $(\frac{1}{\beta^t})^3 \leq c\sqrt{(\frac{1}{\beta^{t-1}} - \frac{1}{\beta^t})\frac{1}{\beta^{t-1}}}$, where $c = \frac{2}{(\xi\beta^0)^2}$.

Proof. For all $t \geq 1$, we have:

$$\begin{aligned} (\frac{1}{\beta^t})^3 / \sqrt{(\frac{1}{\beta^{t-1}} - \frac{1}{\beta^t})\frac{1}{\beta^{t-1}}} &\stackrel{\textcircled{1}}{\leq} (\frac{1}{\beta^t})^2 / \sqrt{\frac{\beta^{t-1}-1}{\beta^{t-1}} - \frac{\beta^{t-1}}{\beta^t}} \\ &\stackrel{\textcircled{2}}{=} \frac{1}{(\beta^0)^2} \cdot \frac{1}{(1+\xi t^p)^2} \cdot \sqrt{\frac{1+\xi t^p}{\xi t^p - \xi(t-1)^p}} \\ &= \frac{1}{(\beta^0)^2} \cdot \frac{1}{(1+\xi t^p)^{3/2}} \cdot \sqrt{\frac{1}{\xi t^p - \xi(t-1)^p}} \\ &\stackrel{\textcircled{3}}{\leq} \frac{1}{(\beta^0)^2} \cdot \frac{1}{(\xi t^p)^{3/2}} \cdot \sqrt{\frac{1}{\xi p t^{p-1}}} \\ &= \frac{1}{(\xi\beta^0)^2} \frac{1}{\sqrt{p}} \cdot \sqrt{t^{-3p}} \cdot \sqrt{t^{1-p}} \\ &= \frac{1}{(\xi\beta^0)^2} \frac{1}{\sqrt{p}} \cdot \sqrt{t^{1-4p}} \\ &\stackrel{\textcircled{4}}{\leq} \frac{2}{(\xi\beta^0)^2}, \end{aligned}$$

where step ① uses $\beta^t \leq \beta^{t+1}$; step ② uses $\beta^t = \beta^0(1 + \xi t^p)$; step ③ uses Lemma A.3 that $t^p - \xi(t-1)^p \geq p t^{p-1}$ for all integer $t \geq 1$; step ④ uses $p \geq 1/4$. □

Lemma A.11. 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) $\|\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t\| \leq \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \bar{A}\|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$.

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

Lemma A.12. Let $\mathbf{P}, \tilde{\mathbf{P}} \in \mathbb{R}^{n \times r}$, and $\mathbf{X}, \tilde{\mathbf{X}} \in \mathcal{M}$. We have:

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

Proof. First, we obtain:

$$\begin{aligned} & \|\mathbf{X}\mathbf{P}^\top\mathbf{X} - \tilde{\mathbf{X}}\tilde{\mathbf{P}}^\top\tilde{\mathbf{X}}\|_F \\ &= \|(\mathbf{X} - \tilde{\mathbf{X}})\mathbf{P}^\top\mathbf{X} + \tilde{\mathbf{X}}\mathbf{P}^\top(\mathbf{X} - \tilde{\mathbf{X}}) + \tilde{\mathbf{X}}(\mathbf{P} - \tilde{\mathbf{P}})^\top\tilde{\mathbf{X}}\|_F \\ &\stackrel{\textcircled{1}}{\leq} \|\mathbf{X} - \tilde{\mathbf{X}}\|_F\|\mathbf{P}^\top\mathbf{X}\| + \|\tilde{\mathbf{X}}\mathbf{P}^\top\|\|\mathbf{X} - \tilde{\mathbf{X}}\|_F + \|\tilde{\mathbf{X}}(\mathbf{P} - \tilde{\mathbf{P}})^\top\tilde{\mathbf{X}}\|_F \\ &\stackrel{\textcircled{2}}{\leq} 2\sqrt{r}\|\mathbf{P}\|\|\mathbf{X} - \tilde{\mathbf{X}}\|_F + \|\mathbf{P} - \tilde{\mathbf{P}}\|_F, \end{aligned} \tag{10}$$

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

Second, we have:

$$\begin{aligned} & \|\mathbf{X}\mathbf{X}^\top\mathbf{P} - \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top\tilde{\mathbf{P}}\|_F \\ &= \|(\mathbf{X} - \tilde{\mathbf{X}})\mathbf{X}^\top\mathbf{P} + \tilde{\mathbf{X}}(\mathbf{X} - \tilde{\mathbf{X}})^\top\mathbf{P} + \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top(\mathbf{P} - \tilde{\mathbf{P}})\|_F \\ &\stackrel{\textcircled{1}}{\leq} \|\mathbf{X} - \tilde{\mathbf{X}}\|_F\|\mathbf{X}^\top\mathbf{P}\| + \|\tilde{\mathbf{X}}\| \cdot \|\mathbf{X} - \tilde{\mathbf{X}}\|_F \cdot \|\mathbf{P}\| + \|\tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top\| \cdot \|\mathbf{P} - \tilde{\mathbf{P}}\|_F \\ &\stackrel{\textcircled{2}}{\leq} 2\sqrt{r}\|\mathbf{P}\|\|\mathbf{X} - \tilde{\mathbf{X}}\|_F + \|\mathbf{P} - \tilde{\mathbf{P}}\|_F, \end{aligned} \tag{11}$$

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

Finally, we derive:

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

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

Lemma A.13. Assume $(e^{t+1})^2 \leq (p^t - p^{t+1}) \cdot (e^t + w^t)$, where $\{p^t\}_{t=1}^\infty$ is a nonnegative decreasing sequence. We have: $\sum_{t=i}^\infty e^{t+1} \leq df$.

Proof. \square

Lemma A.14. Assume $(e^{t+1})^2 \leq e^t(p^t - p^{t+1})$ 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 $1 \leq i < T$. We let $\alpha > 0$ with $1 - \sqrt{\frac{\alpha}{2}} > 0$.

We obtain the following results:

$$\begin{aligned}
e^{t+1} &\stackrel{\textcircled{1}}{\leq} \sqrt{e^t w_t} \\
&\stackrel{\textcircled{2}}{\leq} \sqrt{\frac{\alpha}{2}(e^t)^2 + (w_t)^2/(2\alpha)}, \\
&\stackrel{\textcircled{3}}{\leq} \sqrt{\frac{\alpha}{2}} \cdot e^t + w_t \sqrt{1/(2\alpha)},
\end{aligned} \tag{12}$$

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

Telescoping Inequality (12) over t from i to T , we obtain:

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

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

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

step ① uses the fact that $(1 - \sqrt{\frac{\alpha}{2}})^{-1} \cdot \sqrt{\frac{\alpha}{2}} = 1$ and $(1 - \sqrt{\frac{\alpha}{2}})^{-1} \cdot \sqrt{\frac{1}{2\alpha}} = 2$ when $\alpha = 1/2$; step ② uses $p^{T+1} \geq 0$. Letting $T \rightarrow \infty$, we conclude this lemma. \square

Lemma A.15. Let $\{d^t\}_{t=1}^\infty$ be a nonnegative sequence satisfying $d^t \leq d^{t-1} - d^t + c \cdot t^{pu} \cdot (d^{t-1} - d^t)^u$, where $u > 0$, $c > 0$, $p \in (0, 1)$. Then we have: $d^t \leq \mathcal{O}(t^{-\zeta})$, where $\zeta = (1-p)u$.

Proof. We derive the following results:

$$\begin{aligned}
d^t &\leq d^{t-1} - d^t + c t^{pu} (d^{t-1} - d^t)^u \\
&\stackrel{\textcircled{1}}{\leq} 2 \max(d^{t-1} - d^t, c t^{pu} (d^{t-1} - d^t)^u),
\end{aligned} \tag{13}$$

where step ① uses $a + b \leq 2 \max(a, b)$ for all $a, b \geq 0$.

We now consider two cases for Inequality (13). Case (i). When $d^{t-1} - d^t \geq c t^{pu} (d^{t-1} - d^t)^u$, we have:

$$d^t \leq 2(d^{t-1} - d^t) \leq \frac{2}{3} d^{t-1} \leq \left(\frac{2}{3}\right)^t d^0.$$

Case (ii). When $d^{t-1} - d^t \leq c t^{pu} (d^{t-1} - d^t)^u$, we have: $d^t \leq 2 c t^{pu} (d^{t-1} - d^t)^u$, leading to:

$$(2c)^{-1/u} \frac{(d^t)^{1/u}}{t^p} \leq d^{t-1} - d^t. \tag{14}$$

Summing over Inequality (14) over t from 1 to T yields:

$$\begin{aligned}
0 &\leq \sum_{t=1}^T \{d^{t-1} - d^t - (2c)^{-1/u} \frac{(d^t)^{1/u}}{t^p}\} \\
&\stackrel{\textcircled{1}}{\leq} d^0 - (2c)^{-1/u} \sum_{t=1}^T \left\{ \frac{(d^t)^{1/u}}{t^p} \right\} \\
&\stackrel{\textcircled{2}}{\leq} d^0 - (2c)^{-1/u} (d^T)^{1/u} \cdot \sum_{t=1}^T \frac{1}{t^p} \\
&\stackrel{\textcircled{3}}{\leq} d^0 - (2c)^{-1/u} (d^T)^{1/u} \cdot (1-p) T^{1-p},
\end{aligned}$$

where step ① uses $\sum_{t=1}^T d^{t-1} - d^t = d^0 - d^T \leq d^0$; step ② uses the fact that $\{d^t\}_{t=1}^T$ is decreasing; step ③ uses Lemma A.5 that $\sum_{t=1}^T \frac{1}{t^p} \geq (1-p)T^{(1-p)}$. This further leads to:

$$\begin{aligned} & (2c)^{-1/u} (d^T)^{1/u} \cdot (1-p)T^{1-p} \leq d^0 \\ \Rightarrow & d^T \leq (d^0)^u (2c)(1-p)^{-u} T^{u(p-1)} = \mathcal{O}(T^{(p-1)u}). \end{aligned}$$

□

Lemma A.16. Assume that $d^t \leq c \cdot t^{pu} \cdot (d^{t-1} - d^t)^u$, where $c > 0$, $u, p \in (0, 1)$. Then we have: $d^t \leq \mathcal{O}(t^{-\zeta})$, where $\zeta = \frac{1-p}{1/u-1}$.

Proof. We define $r \triangleq \frac{1}{u} - 1 > 0$, and $g(s) \triangleq s^{-r-1}$.

From the inequality $d^t \leq c \cdot t^{pu} \cdot (d^{t-1} - d^t)^u$, we obtain:

$$c^{1/u} t^p (d^{t-1} - d^t) \geq (d^t)^{1/u} \stackrel{\text{①}}{=} (d^t)^{r+1} \stackrel{\text{②}}{=} \frac{1}{g(d^t)}, \quad (15)$$

where step ① uses $r+1 = \frac{1}{u}$; step ② uses the definition of $g(s)$.

We let $\kappa > 1$ be any constant and consider two cases for $g(d^t)/g(d^{t-1})$.

Case (1). $g(d^t) \leq \kappa g(d^{t-1})$. We define $f(s) \triangleq -\frac{1}{r} \cdot s^{-r}$. We derive:

$$\begin{aligned} 1 & \stackrel{\text{①}}{\leq} c^{1/u} t^p \cdot (d^{t-1} - d^t) \cdot g(d^t) \\ & \stackrel{\text{②}}{\leq} c^{1/u} t^p \cdot (d^{t-1} - d^t) \cdot \kappa g(d^{t-1}) \\ & \stackrel{\text{③}}{\leq} c^{1/u} t^p \cdot \kappa \int_{d^t}^{d^{t-1}} g(s) ds \\ & \stackrel{\text{④}}{=} c^{1/u} t^p \cdot \kappa \cdot (f(d^{t-1}) - f(d^t)) \\ & \stackrel{\text{⑤}}{=} c^{1/u} t^p \cdot \kappa \cdot \frac{1}{r} \cdot ([d^t]^{-r} - [d^{t-1}]^{-r}), \end{aligned}$$

where step ① uses Inequality (15); step ② uses $g(d^t) \leq \kappa g(d^{t-1})$; step ③ uses the fact that $g(s)$ is a nonnegative and increasing function that $(a-b)g(a) \leq \int_b^a g(s) ds$ for all $a, b \in [0, \infty)$; step ④ uses the fact that $\nabla f(s) = g(s)$; step ⑤ uses the definition of $f(\cdot)$. This leads to:

$$[d^t]^{-r} - [d^{t-1}]^{-r} \geq \frac{r}{c^{1/u} \kappa t^p}. \quad (16)$$

Case (2). $g(d^t) > \kappa g(d^{t-1})$. We have:

$$\begin{aligned} g(d^t) > \kappa g(d^{t-1}) & \stackrel{\text{①}}{\Rightarrow} [d^t]^{-(r+1)} > \kappa \cdot [d^{t-1}]^{-(r+1)} \\ & \stackrel{\text{②}}{\Rightarrow} ([d^t]^{-(r+1)})^{\frac{r}{r+1}} > \kappa^{\frac{r}{r+1}} \cdot ([d^{t-1}]^{-(r+1)})^{\frac{r}{r+1}} \\ & \Rightarrow [d^t]^{-r} > \kappa^{\frac{r}{r+1}} \cdot [d^{t-1}]^{-r}, \end{aligned} \quad (17)$$

where step ① uses the definition of $g(\cdot)$; step ② uses the fact that if $a > b > 0$, then $a^{\dot{r}} > b^{\dot{r}}$ for any exponent $\dot{r} \triangleq \frac{r}{r+1} \in (0, 1)$. We further derive:

$$\begin{aligned} [d^t]^{-r} - [d^{t-1}]^{-r} & \stackrel{\text{①}}{\geq} (\kappa^{\frac{r}{r+1}} - 1) \cdot [d^{t-1}]^{-r} \\ & \stackrel{\text{②}}{\geq} (\kappa^{\frac{r}{r+1}} - 1) \cdot [d^0]^{-r}, \end{aligned} \quad (18)$$

where step ① uses Inequality (17); step ② uses $r > 0$ and $d^{t-1} \leq d^0$ for all t .

In view of Inequalities (16) and (18), we have:

$$\begin{aligned} [d^t]^{-r} - [d^{t-1}]^{-r} & \geq \min\left(\frac{\kappa^{-1}r}{c^{1/u} t^p}, (\kappa^{\frac{r}{r+1}} - 1) \cdot [d^0]^{-r}\right) \\ & = \mathcal{O}\left(\frac{1}{t^p}\right). \end{aligned} \quad (19)$$

We now focus on Inequality (19). Telescoping Inequality (19) over $t = \{1, 2, \dots, T\}$, we have:

$$[d^T]^{-r} - [d^0]^{-r} \geq \mathcal{O}(\sum_{t=1}^T \frac{1}{t^p}) \stackrel{\textcircled{1}}{=} \mathcal{O}((1-p)T^{1-p}) = \mathcal{O}(T^{1-p}),$$

where step ① use Lemma A.5. This leads to:

$$d^T = ([d^T]^{-r})^{-1/r} \leq \mathcal{O}(T^{1-p})^{-1/r} = \mathcal{O}(T^{-\zeta}).$$

□

Lemma A.17. Assume that $d^t/d^{t-1} \leq \frac{ct^q}{ct^q+1}$, where $c \geq 0$ and $q \in (0, 1)$. Then we have: $d^t \leq \mathcal{O}(1/\exp(t^\zeta))$, where $\zeta = 1 - q$.

Proof. We define $\gamma^t \triangleq \frac{1}{ct^q+1} \in (0, 1)$.

First, we derive the following results:

$$\begin{aligned} \sum_{t=1}^T \gamma^t &= \sum_{t=1}^T \frac{1}{ct^q+1} \\ &\stackrel{\textcircled{1}}{\geq} \frac{1}{1+c} \sum_{t=1}^T \frac{1}{t^q} \\ &\stackrel{\textcircled{2}}{\geq} \frac{1}{1+c} (1-q)T^{(1-q)} = \mathcal{O}(T^{1-q}), \end{aligned} \quad (20)$$

where step ① uses $ct^q + 1 \leq (1+c)t^q$ since $t^q \geq 1$ if $t \geq 1$; step ② uses Lemma A.5 that $(1-p)T^{(1-p)} \leq \sum_{t=1}^T \frac{1}{t^p}$ for all $p \in (0, 1)$.

Second, noticing that $\frac{d^t}{d^{t-1}} = \frac{ct^q}{ct^q+1} = 1 - \gamma^t$, we have:

$$\frac{d^T}{d^0} \leq (1 - \gamma^1)(1 - \gamma^2)(1 - \gamma^3) \dots (1 - \gamma^T). \quad (21)$$

This further leads to:

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

where step ① uses Inequality (21); step ② uses $\log(ab) = \log(a) + \log(b)$ for all $a, b > 0$; 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 $a, b > 0$; step ⑤ uses Inequality (20).

□

B PROOFS FOR SECTION 2

B.1 PROOF OF LEMMA 2.3

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

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 $\mathbb{P}_{\mu_1}(\mathbf{y}) = \arg \min_{\mathbf{v}} h(\mathbf{v}) + \frac{1}{2\mu_1} \|\mathbf{v} - \mathbf{y}\|_2^2$.

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 $\mathbb{P}_{\mu_2}(\mathbf{y}) = \arg \min_{\mathbf{v}} h(\mathbf{v}) + \frac{1}{2\mu_2} \|\mathbf{v} - \mathbf{y}\|_2^2$.

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

$$\mathbf{y} - \mathbb{P}_{\mu_1}(\mathbf{y}) \in \mu_1 \partial h(\mathbb{P}_{\mu_1}(\mathbf{y})), \quad (22)$$

$$\mathbf{y} - \mathbb{P}_{\mu_2}(\mathbf{y}) \in \mu_2 \partial h(\mathbb{P}_{\mu_2}(\mathbf{y})). \quad (23)$$

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:

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

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 ③ uses the optimality of $\mathbb{P}_{\mu_1}(\mathbf{y})$ and $\mathbb{P}_{\mu_2}(\mathbf{y})$ in Equations (24) and (25); step ④ uses $W_h \leq \frac{1}{\mu_1}$; step ⑤ uses $1 \geq \frac{\mu_2}{\mu_1}$.

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 $\mathbf{s}_1 \in \partial h(\mathbb{P}_{\mu_1}(\mathbf{y}))$ and $\mathbf{s}_2 \in \partial h(\mathbb{P}_{\mu_2}(\mathbf{y}))$, we have:

$$\begin{aligned}
& h_{\mu_2}(\mathbf{y}) - h_{\mu_1}(\mathbf{y}) \\
& \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})) \\
& \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 \\
& \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 \\
& = -\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 \\
& \stackrel{\textcircled{4}}{\leq} \min\left\{-\frac{\mu_1}{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. \\
& \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\} \\
& = \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\} \\
& \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\} \\
& \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 \\
& = \min\left\{\frac{\mu_1}{2\mu_2}, 1\right\} \cdot (\mu_1 - \mu_2) \cdot C_h^2,
\end{aligned}$$

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 $h(\cdot)$; step ③ uses the optimality of $\mathbb{P}_{\mu_2}(\mathbf{y})$ and $\mathbb{P}_{\mu_1}(\mathbf{y})$ in Equations (24) and (25); step ④ uses $W_h \leq \frac{1}{\mu_1}$ 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 $\langle \mathbf{s}_1, \mathbf{s}_2 \rangle \leq \|\mathbf{s}_1\| \cdot \|\mathbf{s}_2\| \leq C_h^2$.

□

B.2 PROOF OF LEMMA 2.3

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

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 $\mathbb{P}_{\mu_1}(\mathbf{y}) = \arg \min_{\mathbf{v}} h(\mathbf{v}) + \frac{1}{2\mu_1} \|\mathbf{v} - \mathbf{y}\|_2^2$.

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 $\mathbb{P}_{\mu_2}(\mathbf{y}) = \arg \min_{\mathbf{v}} h(\mathbf{v}) + \frac{1}{2\mu_2} \|\mathbf{v} - \mathbf{y}\|_2^2$.

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

$$\mathbf{y} - \mathbb{P}_{\mu_1}(\mathbf{y}) \in \mu_1 \partial h(\mathbb{P}_{\mu_1}(\mathbf{y})) \quad (24)$$

$$\mathbf{y} - \mathbb{P}_{\mu_2}(\mathbf{y}) \in \mu_2 \partial h(\mathbb{P}_{\mu_2}(\mathbf{y})). \quad (25)$$

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:

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

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 ③ uses the optimality of $\mathbb{P}_{\mu_1}(\mathbf{y})$ and $\mathbb{P}_{\mu_2}(\mathbf{y})$ in Equations (24) and (25); step ④ uses $W_h \leq \frac{1}{\mu_1}$; step ⑤ uses $1 \geq \frac{\mu_2}{\mu_1}$.

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 $\mathbf{s}_1 \in \partial h(\mathbb{P}_{\mu_1}(\mathbf{y}))$ and $\mathbf{s}_2 \in \partial h(\mathbb{P}_{\mu_2}(\mathbf{y}))$, we have:

$$\begin{aligned} & h_{\mu_2}(\mathbf{y}) - h_{\mu_1}(\mathbf{y}) \\ \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})) \\ \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 \\ \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 \\ = & -\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 \\ \stackrel{\textcircled{4}}{\leq} & \min\left\{-\frac{\mu_1}{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. \\ & \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\} \\ = & \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\} \\ \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\} \\ \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 \\ = & \min\left\{\frac{\mu_1}{2\mu_2}, 1\right\} \cdot (\mu_1 - \mu_2) \cdot C_h^2, \end{aligned}$$

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 $h(\cdot)$; step ③ uses the optimality of $\mathbb{P}_{\mu_2}(\mathbf{y})$ and $\mathbb{P}_{\mu_1}(\mathbf{y})$ in Equations (24) and (25); step ④ uses $W_h \leq \frac{1}{\mu_1}$ 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 $\langle \mathbf{s}_1, \mathbf{s}_2 \rangle \leq \|\mathbf{s}_1\| \cdot \|\mathbf{s}_2\| \leq C_h^2$.

□

B.3 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: $\|\nabla h_{\mu_1}(\mathbf{y}) - \nabla h_{\mu_2}(\mathbf{y})\|_2 \leq |1/\mu_1 - 1/\mu_2| C_h = (\mu_1/\mu_2 - 1)C_h$. \square

B.4 PROOF OF LEMMA 2.5

Proof. We consider the following optimization problem:

$$\bar{\mathbf{y}} = \arg \min_{\mathbf{y}} h_\mu(\mathbf{y}) + \frac{\beta}{2} \|\mathbf{y} - \mathbf{b}\|_2^2. \quad (26)$$

Given $h_\mu(\mathbf{y})$ being (μ^{-1}) -weakly convex and $\beta > \mu^{-1}$, Problem (26) 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 $(\bar{\mathbf{y}}, \check{\mathbf{y}})$:

$$\frac{1}{\mu}(\bar{\mathbf{y}} - \check{\mathbf{y}}) = \beta(\mathbf{b} - \bar{\mathbf{y}}) \quad (27)$$

$$\frac{1}{\mu}(\bar{\mathbf{y}} - \check{\mathbf{y}}) \in \partial h(\check{\mathbf{y}}). \quad (28)$$

Part (a). We have the following results:

$$\begin{aligned} \mathbf{0} &\stackrel{\textcircled{1}}{\in} \partial h(\check{\mathbf{y}}) + \frac{1}{\mu}(\check{\mathbf{y}} - \bar{\mathbf{y}}) \\ &\stackrel{\textcircled{2}}{=} \partial h(\check{\mathbf{y}}) + \frac{1}{\mu}(\check{\mathbf{y}} - \frac{1}{1/\mu + \beta}(\frac{1}{\mu}\check{\mathbf{y}} + \beta\mathbf{b})) \\ &= \partial h(\check{\mathbf{y}}) + \frac{\beta}{1+\mu\beta}(\check{\mathbf{y}} - \mathbf{b}), \end{aligned} \quad (29)$$

where step $\textcircled{1}$ uses Equality (28); step $\textcircled{2}$ uses Equality (27) that $\bar{\mathbf{y}} = \frac{1}{1/\mu + \beta}(\frac{1}{\mu}\check{\mathbf{y}} + \beta\mathbf{b})$. The inclusion in (29) implies that:

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

Part (b). Combining Equalities (27) and (28), we have: $\beta(\mathbf{b} - \bar{\mathbf{y}}) \in \partial h(\check{\mathbf{y}})$.

Part (c). In view of Equation (28), 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.5 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:

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

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 $\|\mathbf{X}\mathbf{P}\|_F^2 = \text{tr}(\mathbf{P}\mathbf{X}^\top \mathbf{X}\mathbf{P}) = \|\mathbf{P}\|_F^2$ 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 $(\mathbf{U} + \mathbf{U}^\top)$. \square

B.6 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:

$$\begin{aligned}
\langle \mathbf{G}, \mathbb{G}_\rho \rangle &= \langle \mathbf{G}, \mathbf{G} - \rho \mathbf{X} \mathbf{G}^\top \mathbf{X} - (1 - \rho) \mathbf{X} \mathbf{X}^\top \mathbf{G} \rangle \\
&= \langle \mathbf{G}, \mathbf{G} \rangle - \rho \text{tr}(\mathbf{G}^\top \mathbf{X} \mathbf{G}^\top \mathbf{X}) - (1 - \rho) \text{tr}(\mathbf{G}^\top \mathbf{X} \mathbf{X}^\top \mathbf{G}) \\
&\stackrel{\textcircled{1}}{=} \langle \mathbf{G}, \mathbf{G} \rangle - \rho \text{tr}(\mathbf{U} \mathbf{U}) - (1 - \rho) \text{tr}(\mathbf{U} \mathbf{U}^\top),
\end{aligned} \tag{30}$$

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

Second, we derive the following equalities:

$$\begin{aligned}
\|\mathbb{G}_\rho\|_F^2 &= \langle \rho \mathbf{X} \mathbf{G}^\top \mathbf{X} + (1 - \rho) \mathbf{X} \mathbf{X}^\top \mathbf{G} - \mathbf{G}, \rho \mathbf{X} \mathbf{G}^\top \mathbf{X} + (1 - \rho) \mathbf{X} \mathbf{X}^\top \mathbf{G} - \mathbf{G} \rangle \\
&\stackrel{\textcircled{1}}{=} \rho^2 \text{tr}(\mathbf{U}^\top \mathbf{U}) + \rho(1 - \rho) \text{tr}(\mathbf{U}^\top \mathbf{U}^\top) - \rho \text{tr}(\mathbf{U}^\top \mathbf{U}^\top) \\
&\quad + (1 - \rho)\rho \text{tr}(\mathbf{U} \mathbf{U}) + (1 - \rho)^2 \text{tr}(\mathbf{U} \mathbf{U}^\top) - (1 - \rho) \text{tr}(\mathbf{U} \mathbf{U}^\top) \\
&\quad - \rho \text{tr}(\mathbf{U} \mathbf{U}) - (1 - \rho) \text{tr}(\mathbf{U} \mathbf{U}^\top) + \langle \mathbf{G}, \mathbf{G} \rangle \\
&\stackrel{\textcircled{2}}{=} (2\rho^2 - 1) \cdot \text{tr}(\mathbf{U}^\top \mathbf{U}) - 2\rho^2 \cdot \text{tr}(\mathbf{U} \mathbf{U}) + \langle \mathbf{G}, \mathbf{G} \rangle,
\end{aligned} \tag{31}$$

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:

$$\text{tr}(\mathbf{G}^\top \mathbf{G}) - \text{tr}(\mathbf{U}^\top \mathbf{U}) \stackrel{\textcircled{1}}{=} \langle \mathbf{G} \mathbf{G}^\top, \mathbf{I}_n - \mathbf{X} \mathbf{X}^\top \rangle \stackrel{\textcircled{2}}{\geq} 0, \tag{32}$$

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 $\max(1, 2\rho) \langle \mathbf{G}, \mathbb{G}_\rho \rangle \geq \|\mathbb{G}_\rho\|_F^2$. We discuss two cases. Case (i): $\rho \in (0, \frac{1}{2}]$. We have:

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

where step ① uses Inequalities (30) and (31); step ② uses $2\rho^2 - \rho \leq 0$ for all $\rho \in (0, \frac{1}{2}]$, and $\text{tr}(\mathbf{U} \mathbf{U}) \leq \text{tr}(\mathbf{U} \mathbf{U}^\top)$ for all $\mathbf{U} \in \mathbb{R}^{r \times r}$.

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

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

where step ① uses Inequalities (30) and (31); step ② uses $2\rho - 1 \geq 0$ for all $\rho \in [\frac{1}{2}, \infty)$, and Inequality (32). Therefore, we conclude that: $\max(1, 2\rho)\langle \mathbf{G}, \mathbb{G}_\rho \rangle \geq \|\mathbb{G}_\rho\|_F^2$.

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

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

where step ① uses Inequalities (30) and (31); step ② uses $1 - \rho^2 \geq 0$, and Inequality (32).

Case (ii): $\rho \in (1, \infty)$. We have:

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

where step ① uses Inequality (31); step ② uses $4\rho^2 - 1 \leq 0$ for all $\rho \in (0, \frac{1}{2}]$, and 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: $\min(1, \rho^2)\|\mathbb{G}_1\|_F^2 \leq \|\mathbb{G}_\rho\|_F^2$.

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

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

where step ① uses Inequality (31); step ② uses $4\rho^2 - 1 \leq 0$ for all $\rho \in (0, \frac{1}{2}]$, and Inequality (32).

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

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

where step ① uses Inequalities (30) and (31); step ② uses $2\rho^2 - \frac{1}{2} \geq 0$ for all $\rho \in (\frac{1}{2}, \infty)$, and 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 $\|\mathbb{G}_\rho\|_F \geq \min(1, 2\rho)\|\mathbb{G}_{1/2}\|_F$.

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

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

where step ① uses Inequality (31); step ② uses $2\rho^2 - \frac{1}{2} \leq 0$ for all $\rho \in (0, \frac{1}{2}]$, and the fact that $\text{tr}(\mathbf{U}\mathbf{U}) - \text{tr}(\mathbf{U}\mathbf{U}^\top) \leq 0$ for all $\mathbf{U} \in \mathbb{R}^{r \times r}$.

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

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

where step ① uses Inequalities (30) and (31); step ② uses $4\rho^2 - 1 \geq 0$ for all $\rho \in (\frac{1}{2}, \infty)$, and Inequality (32). Therefore, we conclude that: $\|\mathbb{G}_\rho\|_F \leq \max(1, 2\rho)\|\mathbb{G}_{1/2}\|_F$.

□

B.7 PROOF OF LEMMA 2.13

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

$$(\mathbf{0} \in \partial \iota_{\mathcal{M}}(\mathbf{X}) + \nabla f(\mathbf{X})) \Leftrightarrow (\mathbf{0} \in \text{Proj}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}(\nabla f(\mathbf{X}))). \quad (33)$$

Therefore, we derive the following results:

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

where step ① uses Formulation (33); step 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 ③ uses the norm inequality $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\| \|\mathbf{B}\|_F$, 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.

□

C PROOFS FOR SECTION 4

C.1 PROOF OF LEMMA 4.1

Proof. We define $\mu^t = \tau/\beta^t$.

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

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

$$\begin{aligned} \mathbf{0} &= \nabla h_{\mu^t}(\mathbf{y}^{t+1}) + \beta^t(\mathbf{y}^{t+1} - \mathbf{y}^t) + \nabla_{\mathbf{y}} \mathcal{S}(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \\ &\stackrel{\text{①}}{=} \nabla h_{\mu^t}(\mathbf{y}^{t+1}) + \beta^t(\mathbf{y}^{t+1} - \mathbf{y}^t) - \mathbf{z}^t + \beta^t(\mathbf{y}^t - \mathcal{A}(\mathbf{X}^{t+1})) \\ &= \nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \mathbf{z}^t + \beta^t(\mathbf{y}^{t+1} - \mathcal{A}(\mathbf{X}^{t+1})) \\ &\stackrel{\text{②}}{=} \nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \mathbf{z}^t + \frac{1}{\sigma}(\mathbf{z}^t - \mathbf{z}^{t+1}), \end{aligned} \quad (34)$$

where step 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 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:

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

where step 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 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 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 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:

$$\begin{aligned} &\|\nabla h_{\mu^{t-1}}(\mathbf{y}^t) - \nabla h_{\mu^t}(\mathbf{y}^{t+1})\| \\ &\stackrel{\text{①}}{\leq} \|\nabla h_{\mu^{t-1}}(\mathbf{y}^t) - \nabla h_{\mu^t}(\mathbf{y}^t)\| + \|\nabla h_{\mu^t}(\mathbf{y}^t) - \nabla h_{\mu^t}(\mathbf{y}^{t+1})\| \\ &\stackrel{\text{②}}{\leq} \|\nabla h_{\mu^t}(\mathbf{y}^t) - \nabla h_{\mu^{t-1}}(\mathbf{y}^t)\| + \frac{1}{\mu^t} \|\mathbf{y}^{t+1} - \mathbf{y}^t\| \\ &\stackrel{\text{③}}{\leq} C_h \left(\frac{\mu^{t-1}}{\mu^t} - 1 \right) + \frac{1}{\mu^t} \|\mathbf{y}^{t+1} - \mathbf{y}^t\| \\ &\stackrel{\text{④}}{=} C_h \left(\frac{\beta^t}{\beta^{t-1}} - 1 \right) + \frac{\beta^t}{\tau} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|, \end{aligned} \quad (35)$$

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; step ④ uses $\mu^t = \frac{\tau}{\beta^t}$.

Second, we have from Equality (34):

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

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.2 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:

$$\begin{aligned} & \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 - \frac{\varrho}{1-\varrho}(\|\mathbf{z}^{t-1} - \mathbf{z}^t\|_2^2 - \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2) \\ & \leq \frac{\sigma^2}{(1-\varrho)^2} \|\nabla h_{\mu^t}(\mathbf{y}^{t+1}) - \nabla h_{\mu^{t-1}}(\mathbf{y}^t)\|_2^2 \\ & \stackrel{\textcircled{1}}{\leq} \frac{2\sigma^2}{(1-\varrho)^2} \cdot \frac{(\beta^t)^2}{\tau^2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \frac{2\sigma^2}{(1-\varrho)^2} \cdot C_h^2 \left(\frac{\beta^t}{\beta^{t-1}} - 1\right)^2 \\ & \stackrel{\textcircled{2}}{\leq} \underbrace{\frac{2\sigma^2}{(1-\varrho)^2} \frac{1}{\tau^2}}_{\triangleq \hat{\sigma}} \cdot (\beta^t)^2 \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \underbrace{\frac{2\sigma^2}{(1-\varrho)^2} \cdot C_h^2 \cdot \frac{6}{p}}_{\triangleq \hat{\sigma}} \cdot \left(\frac{\beta^0}{\beta^t} - \frac{\beta^0}{\beta^{t+1}}\right), \end{aligned}$$

where step ① uses Inequality (35), and the inequality $(a + b)^2 \leq 2a^2 + 2b^2$ for all $a, b \in \mathbb{R}$; step ② uses Lemma A.9 that $(\frac{\beta^t}{\beta^{t-1}} - 1)^2 \leq \frac{6\beta^0}{p} (\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}})$ 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 $\beta^t = \beta^0(1 + \xi t^p)$; 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.6 with $\mathbf{a}_t = \|\mathbf{z}^{t+1}\|$, $c = \sigma C_h$, and $\varrho = \sigma - 1 \in [0, 1)$, we have:

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

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

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

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 $\|\mathbf{z}^t\| \leq \bar{z}$.

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

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

where step ① uses the definition of $L(\mathbf{X}, \mathbf{y}, \mathbf{z}, \beta)$ and the positivity of $\{\mu^t, \mathbb{T}^t, \mathbb{D}^t, \mathbb{P}^t\}$; 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 $\|\mathbf{z}^t\| \leq \bar{z}$. □

C.4 PROOF OF LEMMA 4.5

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

We define $\mu^t \triangleq \tau/\beta^t$. We define $\mathbb{D}^t \triangleq \frac{\sigma-1}{2-\sigma} \frac{2}{\beta^{t-1}} \|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2$.

Part (a). We focus on the sufficient decrease for variable $\{\beta\}$. We have:

$$\begin{aligned}
&\varepsilon_\beta \beta^t \mathcal{B}_{t+1}^2 + L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^{t+1}, \beta^{t+1}) - L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^{t+1}, \beta^t) \\
&\stackrel{\textcircled{1}}{=} \varepsilon_\beta \beta^t \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}} \right) \cdot \frac{1}{\beta^t} + h_{\mu^{t+1}}(\mathbf{y}^{t+1}) - h_{\mu^t}(\mathbf{y}^{t+1}) + \frac{\beta^{t+1} - \beta^t}{2} \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 \\
&\stackrel{\textcircled{2}}{\leq} \varepsilon_\beta \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}} \right) + \tau C_h^2 \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}} \right) + \frac{\beta^{t+1} - \beta^t}{2(\sigma\beta^t)^2} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \\
&\stackrel{\textcircled{3}}{\leq} \varepsilon_\beta \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}} \right) + \tau C_h^2 \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}} \right) + \frac{(1+\xi)\beta^t - \beta^t}{\sigma^2 \beta^t} \frac{1}{2\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \\
&\stackrel{\textcircled{4}}{\leq} (\varepsilon_\beta + \tau C_h^2) \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}} \right) + \frac{\xi}{\sigma^2} \frac{1}{2\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2,
\end{aligned} \tag{36}$$

where step ① uses the definition of $L(\cdot, \cdot, \cdot, \cdot)$; step ② uses Lemma 2.3, and the fact that $\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1} = \frac{1}{\sigma\beta^t}(\mathbf{z}^{t+1} - \mathbf{z}^t)$; step ③ uses $\beta^{t+1} \leq (1+\xi)\beta^t$; step ④ uses $\beta^t(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}) \leq 1$.

Part (b). We focus on the sufficient decrease for variable $\{\mathbf{z}\}$. We have:

$$\begin{aligned}
&\varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^{t+1}, \beta^t) - L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^t, \beta^t) \\
&\stackrel{\textcircled{1}}{=} \varepsilon_z \beta^t \|\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}\|_2^2 + \langle \mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}, \mathbf{z}^{t+1} - \mathbf{z}^t \rangle \\
&\stackrel{\textcircled{2}}{=} \left(\frac{\varepsilon_z}{\sigma^2} + \frac{1}{\sigma} \right) \frac{1}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2,
\end{aligned} \tag{37}$$

where step ① uses the definition of $L(\mathbf{X}, \mathbf{y}, \mathbf{z}, \beta)$; step ② uses $\mathbf{z}^{t+1} - \mathbf{z}^t = \sigma\beta^t(\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1})$.

Part (c). We focus on the sufficient decrease for variable $\{\mathbf{y}\}$. We have:

$$\begin{aligned}
& L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^t, \beta^t) - L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \\
&= h_{\mu^t}(\mathbf{y}^{t+1}) - h_{\mu^t}(\mathbf{y}^t) + \langle \mathbf{y}^t - \mathbf{y}^{t+1}, \mathbf{z}^t \rangle + \frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathcal{A}(\mathbf{X}^{t+1})\|_2^2 - \frac{\beta^t}{2} \|\mathbf{y}^t - \mathcal{A}(\mathbf{X}^{t+1})\|_2^2 \\
&\stackrel{\textcircled{1}}{=} h_{\mu^t}(\mathbf{y}^{t+1}) - h_{\mu^t}(\mathbf{y}^t) + \langle \mathbf{y}^t - \mathbf{y}^{t+1}, \mathbf{z}^t + \beta^t(\mathcal{A}(\mathbf{X}^{t+1}) - \mathbf{y}^{t+1}) \rangle - \frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 \\
&\stackrel{\textcircled{2}}{=} h_{\mu^t}(\mathbf{y}^{t+1}) - h_{\mu^t}(\mathbf{y}^t) - \frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \langle \mathbf{y}^t - \mathbf{y}^{t+1}, \mathbf{z}^t + \frac{1}{\sigma}(\mathbf{z}^{t+1} - \mathbf{z}^t) \rangle \\
&\stackrel{\textcircled{3}}{=} h_{\mu^t}(\mathbf{y}^{t+1}) - h_{\mu^t}(\mathbf{y}^t) - \frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \langle \mathbf{y}^t - \mathbf{y}^{t+1}, \nabla h_{\mu^t}(\mathbf{y}^{t+1}) \rangle \\
&\stackrel{\textcircled{4}}{\leq} \frac{1}{2\mu^t} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 - \frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 \\
&\stackrel{\textcircled{5}}{=} \left(\frac{1}{\tau} - 1\right) \frac{\beta^t}{2} \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2, \tag{38}
\end{aligned}$$

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 $h_{\mu^t}(\mathbf{y})$ is $(1/\mu^t)$ -weakly convex w.r.t \mathbf{y} ; step ⑤ uses $\mu^t\beta^t = \tau$.

Part (d). We focus on the sufficient decrease for variable $\{\mathbf{X}\}$. We have:

$$L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t, \beta^t) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) = \mathfrak{X}. \tag{39}$$

Adding Inequalities (36), (37), (38), and (39) together, we have:

$$\begin{aligned}
& \varepsilon_\beta \beta^t \mathcal{B}_{t+1}^2 + \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 \\
&+ L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^{t+1}, \beta^{t+1}) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \\
&\leq (\varepsilon_\beta + \tau C_h^2) \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}\right) + \left(\frac{1}{2\tau} - \frac{1}{2}\right) \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \left(\frac{\xi}{2\sigma^2} + \frac{\varepsilon_z}{\sigma^2} + \frac{1}{\sigma}\right) \cdot \frac{1}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \\
&\stackrel{\textcircled{1}}{\leq} (\varepsilon_\beta + \tau C_h^2) \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}\right) + \left(\frac{1}{2\tau} - \frac{1}{2}\right) \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \left(\frac{1}{2\sigma} + \frac{1}{4\sigma} + \frac{1}{\sigma}\right) \cdot \frac{1}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \\
&\stackrel{\textcircled{2}}{\leq} (\varepsilon_\beta + \tau C_h^2) \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}\right) + \left(\frac{1}{2\tau} - \frac{1}{2}\right) \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \frac{2}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \\
&\stackrel{\textcircled{3}}{\leq} (\varepsilon_\beta + \tau C_h^2) \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}\right) + \left(\frac{1}{2\tau} - \frac{1}{2}\right) \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 \\
&\quad + \frac{2}{\beta^t} \left\{ \frac{\sigma-1}{2-\sigma} (\|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2 - \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2) + \dot{\sigma}(\beta^t)^2 \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 + \ddot{\sigma} \left(\frac{\beta^0}{\beta^t} - \frac{\beta^0}{\beta^{t+1}}\right) \right\} \\
&\stackrel{\textcircled{4}}{\leq} \underbrace{(\varepsilon_\beta + \tau C_h^2 + \frac{2}{\sigma} \ddot{\sigma}) \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}\right)}_{\triangleq_c} + \underbrace{\left(\frac{\sigma-1}{2-\sigma} \left(\frac{2}{\beta^{t-1}} \|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2 - \frac{2}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2\right)\right)}_{=\mathbb{D}^t - \mathbb{D}^{t+1}} \\
&\quad + \left(\frac{1}{2\tau} + 2\dot{\sigma} - \frac{1}{2}\right) \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2,
\end{aligned}$$

where step ① uses $\varepsilon_z = \frac{1}{4}$, $\xi \leq 1$, $\sigma \geq 1$; step ② uses $\sigma \geq 1$; step ③ uses the upper bound for $\|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2$ as shown in Lemma 4.1(b); step ④ uses $\frac{1}{\beta^t} \leq \frac{1}{\beta^{t-1}}$. This further leads to:

$$\begin{aligned}
& \varepsilon_\beta \beta^t \mathcal{B}_{t+1}^2 + \varepsilon_y \beta^t \mathcal{Y}_{t+1}^2 + \varepsilon_z \beta^t \mathcal{Z}_{t+1}^2 + \mathbb{D}^{t+1} - \mathbb{D}^t - \frac{c}{\beta^t} + \frac{c}{\beta^{t+1}} \\
&+ L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^{t+1}, \beta^{t+1}) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \\
&\leq \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 \cdot \left(\frac{1}{2\tau} + 2\dot{\sigma} - \frac{1}{2}\right) \\
&\stackrel{\textcircled{1}}{=} \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 \cdot \left(\frac{1}{2\tau} + 2\frac{2\sigma^2}{(2-\sigma)^2} \frac{1}{\tau^2} - \frac{1}{2}\right) \\
&\stackrel{\textcircled{2}}{\leq} \beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 \cdot \left(\frac{2-\sigma}{8} + \frac{4\sigma}{16} - \frac{1}{2}\right) \\
&= -\beta^t \|\mathbf{y}^{t+1} - \mathbf{y}^t\|_2^2 \cdot \underbrace{\left(\frac{2-\sigma}{8}\right)}_{\triangleq_{\varepsilon_y}},
\end{aligned}$$

where step ① uses the definition of $\dot{\sigma} \triangleq \frac{2\sigma^2}{(2-\sigma)^2} \frac{1}{\tau^2}$, as shown in Lemma 4.1; step ② uses $\tau \geq \frac{4}{2-\sigma}$. \square

C.5 PROOF OF LEMMA 4.6

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

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

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

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

$$\langle \mathbf{X}^{t+1} - \mathbf{X}^t, \mathbf{G}^t \rangle + \frac{\theta\ell(\beta^t)}{2} \|\mathbf{X}^{t+1} - \mathbf{X}_c^t\|_F^2 \leq \langle \mathbf{X}^t - \mathbf{X}^t, \mathbf{G}^t \rangle + \frac{\theta\ell(\beta^t)}{2} \|\mathbf{X}^t - \mathbf{X}_c^t\|_F^2. \quad (40)$$

Second, we have:

$$\begin{aligned} & L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t, \beta^t) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \\ &= \mathcal{S}(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t, \beta^t) - \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) + g(\mathbf{X}^t) - g(\mathbf{X}^{t+1}) \\ &\stackrel{\textcircled{1}}{\leq} \frac{\ell(\beta^t)}{2} \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 + \langle \mathbf{X}^{t+1} - \mathbf{X}^t, \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \rangle + \langle \mathbf{X}^t - \mathbf{X}^{t+1}, \partial g(\mathbf{X}^t) \rangle, \end{aligned} \quad (41)$$

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 $g(\mathbf{X})$.

Third, we derive:

$$\begin{aligned} & \langle \mathbf{X}^{t+1} - \mathbf{X}^t, \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) - \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}_c^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \rangle \\ &\stackrel{\textcircled{1}}{\leq} \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F \cdot \|\nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) - \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}_c^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t)\|_F \\ &\stackrel{\textcircled{2}}{\leq} \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F \cdot \ell(\beta^t) \|\mathbf{X}^t - \mathbf{X}_c^t\|_F \\ &\stackrel{\textcircled{3}}{\leq} \alpha\ell(\beta^t) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F \cdot \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F \\ &\stackrel{\textcircled{4}}{\leq} \frac{\alpha\ell(\beta^t)}{2} \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 + \frac{\alpha\ell(\beta^t)}{2} \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F^2, \end{aligned} \quad (42)$$

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 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 $b \in \mathbb{R}$.

Summing Inequalities (40), (42), and (41), we obtain:

$$\begin{aligned} & L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t, \beta^t) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \\ &\leq \frac{\ell(\beta^t)}{2} \{ (1 + \alpha) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 + \alpha \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F + \theta \|\mathbf{X}^t - \mathbf{X}_c^t\|_F^2 - \theta \|\mathbf{X}^{t+1} - \mathbf{X}_c^t\|_F^2 \} \\ &\stackrel{\textcircled{1}}{=} \frac{\ell(\beta^t)}{2} \{ (1 + \alpha) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 + (\alpha + \theta\alpha^2) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F^2 - \theta \|\mathbf{X}^{t+1} - \mathbf{X}^t - \alpha(\mathbf{X}^t - \mathbf{X}^{t-1})\|_F^2 \} \\ &\stackrel{\textcircled{2}}{\leq} \frac{\ell(\beta^t)}{2} \{ (1 + \alpha) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 + (\alpha + \theta\alpha^2) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F^2 \\ &\quad + \theta(\alpha - 1) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 - \theta\alpha(\alpha - 1) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F^2 \} \\ &= \underbrace{\frac{1}{2}(\alpha + \theta\alpha)\ell(\beta^t)\|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F^2}_{\triangleq \mathbb{P}^t} + \frac{\ell(\beta^t)}{2} \cdot \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 \cdot \{1 + \alpha + \theta\alpha - \theta\} \\ &= \mathbb{P}^t - \mathbb{P}^{t+1} + \frac{1}{2} \cdot \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 \cdot \{ \ell(\beta^t)(1 + \alpha + \theta\alpha - \theta) + \ell(\beta^{t+1})(\alpha + \theta\alpha) \} \\ &\stackrel{\textcircled{3}}{\leq} \mathbb{P}^t - \mathbb{P}^{t+1} + \frac{1}{2} \cdot \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 \cdot \ell(\beta^t) \underbrace{\{ (1 + \alpha + \theta\alpha - \theta) + (1 + \xi)(\alpha + \theta\alpha) \}}_{\triangleq -\varepsilon'_x} \\ &\stackrel{\textcircled{4}}{\leq} \mathbb{P}^t - \mathbb{P}^{t+1} - \frac{1}{2} \cdot \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 \cdot \varepsilon'_x \cdot \beta^t \bar{\ell} \\ &\stackrel{\textcircled{5}}{=} \mathbb{P}^t - \mathbb{P}^{t+1} - \varepsilon_x \cdot \beta^t \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2, \end{aligned}$$

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 $\mathbf{a} = \mathbf{X}^{t+1} - \mathbf{X}^t$, 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 $\beta^{t+1} \leq (1 + \xi)\beta^t$; step ④ uses Lemma 4.3 that $\beta^t \bar{\ell} \leq \ell(\beta^t) \leq \beta^t \bar{\ell}$; step ⑤ uses $\varepsilon_x \triangleq \frac{1}{2}\varepsilon'_x\bar{\ell} > 0$.

□

C.6 PROOF OF LEMMA 4.7

Proof. We define: $\Theta^t \triangleq L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) + c/\beta^t + \mathbb{P}^t + \mathbb{D}^t$.

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

$$\begin{aligned} & L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^{t+1}, \beta^{t+1}) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \\ & \leq c/\beta^t - c/\beta^{t+1} + \mathbb{D}^t - \mathbb{D}^{t+1} - \varepsilon_y \beta^t \mathcal{Y}_{t+1}^2 - \varepsilon_z \beta^t \mathcal{Z}_{t+1}^2 - \varepsilon_\beta \beta^t \mathcal{B}_{t+1}^2 + \mathfrak{X}. \end{aligned} \quad (43)$$

Using Lemma 4.6, we have:

$$\mathfrak{X} \leq \mathbb{P}^t - \mathbb{P}^{t+1} - \varepsilon_x \beta^t \mathcal{X}_{t+1}^2.$$

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

$$\begin{aligned} \Theta^t - \Theta^{t+1} & \geq \varepsilon_\beta \beta^t \mathcal{B}_{t+1}^2 + \varepsilon_y \beta^t \mathcal{Y}_{t+1}^2 + \varepsilon_x \beta^t \mathcal{X}_{t+1}^2 + \varepsilon_z \beta^t \mathcal{Z}_{t+1}^2 \\ & \geq \min(\varepsilon_y, \varepsilon_x, \varepsilon_z, \varepsilon_\beta) \cdot \beta^t \cdot (\mathcal{B}_{t+1}^2 + \mathcal{X}_{t+1}^2 + \mathcal{Y}_{t+1}^2 + \mathcal{Z}_{t+1}^2). \end{aligned}$$

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

$$\begin{aligned} \sum_{t=1}^T \beta^t (\mathcal{B}_{t+1}^2 + \mathcal{X}_{t+1}^2 + \mathcal{Y}_{t+1}^2 + \mathcal{Z}_{t+1}^2) & \leq \frac{1}{\min(\varepsilon_y, \varepsilon_x, \varepsilon_z)} \cdot \sum_{t=1}^T (\Theta^t - \Theta^{t+1}) \\ & = \frac{1}{\min(\varepsilon_y, \varepsilon_x, \varepsilon_z)} \cdot (\Theta^1 - \Theta^{T+1}) \\ & \stackrel{\textcircled{1}}{\leq} \frac{1}{\min(\varepsilon_y, \varepsilon_x, \varepsilon_z)} \cdot (\Theta^1 - \underline{\Theta}), \end{aligned} \quad (44)$$

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

$$\begin{aligned} & \sum_{t=1}^T \beta^t (\mathcal{B}_{t+1}^2 + \mathcal{X}_{t+1}^2 + \mathcal{Y}_{t+1}^2 + \mathcal{Z}_{t+1}^2) \\ & = \sum_{t=1}^T \frac{1}{\beta^t} (\beta^t)^2 (\mathcal{B}_{t+1}^2 + \mathcal{X}_{t+1}^2 + \mathcal{Y}_{t+1}^2 + \mathcal{Z}_{t+1}^2) \\ & \geq \frac{1}{\beta^T} \sum_{t=1}^T (\beta^t)^2 (\mathcal{B}_{t+1}^2 + \mathcal{X}_{t+1}^2 + \mathcal{Y}_{t+1}^2 + \mathcal{Z}_{t+1}^2) \\ & \stackrel{\textcircled{1}}{\geq} \frac{1}{4T\beta^T} \left(\sum_{t=1}^T \beta^t [(\mathcal{B}_{t+1} + \mathcal{X}_{t+1} + \mathcal{Y}_{t+1} + \mathcal{Z}_{t+1})]^2 \right), \end{aligned} \quad (45)$$

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 (44) and (45), we have:

$$\sum_{t=1}^T \beta^t [(\mathcal{B}_{t+1} + \mathcal{X}_{t+1} + \mathcal{Y}_{t+1} + \mathcal{Z}_{t+1})] \leq \left\{ \frac{\Theta^1 - \underline{\Theta}}{\min(\varepsilon_y, \varepsilon_x, \varepsilon_z, \varepsilon_\beta)} \cdot 4T\beta^T \right\}^{1/2} = \mathcal{O}(T^{(1+p)/2}).$$

□

C.7 PROOF OF THEOREM 4.8

Proof. We define $e^t \triangleq \mathcal{B}_t + \mathcal{Z}_t + \mathcal{Y}_t + \mathcal{X}_t$.

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

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:

$$\begin{aligned} & \|\ddot{\mathbf{G}} - \dot{\mathbf{G}}\|_\mathbb{F} \\ & \stackrel{\textcircled{1}}{=} \|\nabla f(\mathbf{X}^t) - \nabla f(\mathbf{X}_c^t) - \beta^t \mathcal{A}^\top(\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t) - \theta \ell(\beta^t)(\mathbf{X}^{t+1} - \mathbf{X}_c^t)\|_\mathbb{F} \\ & \stackrel{\textcircled{2}}{\leq} L_f \|\mathbf{X}^t - \mathbf{X}_c^t\|_\mathbb{F} + \beta^t \bar{\mathcal{A}} \|\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t\| + \theta \ell(\beta^t) \|\mathbf{X}^{t+1} - \mathbf{X}_c^t\|_\mathbb{F} \\ & \stackrel{\textcircled{3}}{\leq} L_f \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_\mathbb{F} + \beta^t \bar{\mathcal{A}} \{\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \bar{\mathcal{A}} \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_\mathbb{F}\} \\ & \quad + \theta \ell(\beta^t) (\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_\mathbb{F} + \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_\mathbb{F}) \\ & \stackrel{\textcircled{4}}{\leq} (L_f + \beta^t \bar{\mathcal{A}}^2 + \theta \ell(\beta^t)) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_\mathbb{F} + \beta^t \bar{\mathcal{A}} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \theta \ell(\beta^t) \|\mathbf{X}^{t+1} - \mathbf{X}^t\|_\mathbb{F} \\ & \stackrel{\textcircled{5}}{=} \mathcal{O}(\beta^{t-1} e^t) + \mathcal{O}(\beta^t e^{t+1}), \end{aligned} \quad (46)$$

where step ① uses the definitions of $\{\ddot{\mathbf{G}}, \dot{\mathbf{G}}\}$; step ② uses the triangle inequality; step 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 $\|\mathcal{A}(\mathbf{X}_c^t) - \mathbf{y}^t\| \leq \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|_F + \bar{A}\|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$, as shown in Lemma A.11.

We derive the following inequalities:

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

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

Finally, we derive:

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{X}^t, \check{\mathbf{y}}^t, \mathbf{z}^t) \\
& \stackrel{\text{①}}{=} \frac{1}{T} \sum_{t=1}^T \{\|\mathcal{A}(\mathbf{X}^t) - \check{\mathbf{y}}^t\| + \|\partial h(\check{\mathbf{y}}^t) - \mathbf{z}^t\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\dot{\mathbf{G}})\|_F\} \\
& \stackrel{\text{②}}{\leq} \frac{1}{T} \sum_{t=1}^T \{\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \|\check{\mathbf{y}}^t - \mathbf{y}^t\| + \|(1 - \frac{1}{\sigma})(\mathbf{z}^t - \mathbf{z}^{t-1})\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\dot{\mathbf{G}})\|_F\} \\
& \stackrel{\text{③}}{=} \frac{1}{T} \sum_{t=1}^T \{\mathcal{O}(\beta^{t-1}e^t) + \mathcal{O}(\beta^te^{t+1})\} + \frac{1}{T} \sum_{t=1}^T \|\check{\mathbf{y}}^t - \mathbf{y}^t\| \\
& \stackrel{\text{④}}{=} \frac{1}{T} \sum_{t=1}^T \{\mathcal{O}(\beta^te^{t+1}) + \mathcal{O}(\beta^{t-1}e^t)\} + \frac{1}{T} \mathcal{O}(\sum_{t=1}^T \frac{1}{t^p}) \\
& \stackrel{\text{⑤}}{=} \mathcal{O}(T^{(p-1)/2}) + \mathcal{O}(T^{1-p-1}) \\
& \stackrel{\text{⑥}}{=} \mathcal{O}(T^{-1/3}),
\end{aligned}$$

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 Lemma 2.5(c) that $\|\check{\mathbf{y}}^t - \mathbf{y}^t\| \leq \mu^t C_h = \mathcal{O}(\frac{1}{t^p})$; step ⑤ uses Lemma A.5 that $\sum_{t=1}^T \frac{1}{t^p} \leq \mathcal{O}(T^{1-p})$, and Lemma 4.7(b) that $\frac{1}{T} \sum_{t=1}^T \beta^te^{t+1} \leq \mathcal{O}(T^{(p-1)/2})$; step ⑥ uses the choice $p = 1/3$ and Lemma 4.7(b). \square

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

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

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

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

We derive the following inequalities:

$$\begin{aligned}
& L(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t, \beta^t) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) = \dot{L}(\mathbf{X}^{t+1}) - \dot{L}(\mathbf{X}^t) \\
& \stackrel{\text{①}}{=} \{\mathcal{S}^t(\mathbf{X}^{t+1}, \mathbf{y}^t, \mathbf{z}^t, \beta^t) - g(\mathbf{X}^{t+1})\} - \{\mathcal{S}^t(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) - g(\mathbf{X}^t)\} \\
& \stackrel{\text{②}}{\leq} \frac{1}{2}\ell(\beta^t)\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F^2 + \langle \mathbf{G}^t, \mathbf{X}^{t+1} - \mathbf{X}^t \rangle \\
& \stackrel{\text{③}}{=} \frac{1}{2}\ell(\beta^t)\|\text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t) - \mathbf{X}^t\|_F^2 + \langle \mathbf{G}^t, \text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t) - \mathbf{X}^t + \eta^t \mathbb{G}_\rho^t \rangle - \eta^t \langle \mathbf{G}^t, \mathbb{G}_\rho^t \rangle \\
& \stackrel{\text{④}}{\leq} \frac{1}{2}\ell(\beta^t)\|\text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t) - \mathbf{X}^t\|_F^2 + \bar{g}\|\text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t) - \mathbf{X}^t + \eta^t \mathbb{G}_\rho^t\|_F - \frac{\eta^t}{\max(1, 2\rho)}\|\mathbb{G}_\rho^t\|_F^2 \\
& \stackrel{\text{⑤}}{\leq} \frac{1}{2}\ell(\beta^t)\dot{k}\|\eta^t \mathbb{G}_\rho^t\|_F^2 + \frac{1}{2}\bar{g}\ddot{k}\|\eta^t \mathbb{G}_\rho^t\|_F^2 - \frac{\eta^t}{\max(1, 2\rho)}\|\mathbb{G}_\rho^t\|_F^2 \\
& \stackrel{\text{⑥}}{=} \eta^t\|\mathbb{G}_\rho^t\|_F^2 \cdot \left\{ \frac{1}{2}\ell(\beta^t)\dot{k}\frac{b^t\gamma^j}{\beta^t} + \frac{1}{2}\bar{g}\ddot{k}\frac{b^t\gamma^j}{\beta^t} - \frac{1}{\max(1, 2\rho)} \right\} \\
& \stackrel{\text{⑦}}{\leq} \eta^t\|\mathbb{G}_\rho^t\|_F^2 \cdot \left\{ \left(\frac{\bar{b}}{2}\dot{k}\bar{\ell} + \frac{\bar{b}}{2\beta^0}\ddot{k}\bar{g}\right)\gamma^j - \frac{1}{\max(1, 2\rho)} \right\} \\
& \stackrel{\text{⑧}}{\leq} \eta^t\|\mathbb{G}_\rho^t\|_F^2 \cdot \{-\delta\},
\end{aligned} \tag{47}$$

where step ① uses the definitions of $L(\mathbf{X}, \mathbf{y}, \mathbf{z}, \beta)$; 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, $\|\mathbf{G}^t\|_F \leq \bar{g}$, and Lemma 2.12(a) that $\langle \mathbf{G}^t, \mathbb{G}_\rho^t \rangle \geq \frac{1}{\max(1, 2\rho)}\|\mathbb{G}_\rho^t\|_F^2$; 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 γ^j is sufficiently small such that:

$$\gamma^j \leq \frac{2(\frac{1}{\max(1, 2\rho)} - \delta)}{\bar{\ell}\bar{k}\bar{b} + \bar{g}\ddot{k}\bar{b}/\beta^0} \triangleq \bar{\gamma}. \tag{48}$$

Given Inequality (47) coincides with the condition of the line search procedure, we complete the proof.

Part (b). We derive the following inequalities:

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

where step ① uses Inequality (47); step ② uses Lemma 2.12(b) that $\|\mathbb{G}_\rho\|_F \geq \min(1, 2\rho)\|\mathbb{G}_{1/2}\|_F$; 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.

□

C.9 PROOF OF LEMMA 4.12

Proof. We define: $\Theta^t \triangleq L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) + c/\beta^t + \mathbb{D}^t + 0 \times \mathbb{P}^t$,

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

$$\begin{aligned} & L(\mathbf{X}^{t+1}, \mathbf{y}^{t+1}, \mathbf{z}^{t+1}, \beta^{t+1}) - L(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \\ & \leq c/\beta^t - c/\beta^{t+1} + \mathfrak{X} + \mathbb{D}^t - \mathbb{D}^{t+1} - \varepsilon_y \beta^t \mathcal{Y}_{t+1}^2 - \varepsilon_z \beta^t \mathcal{Z}_{t+1}^2. \end{aligned} \quad (49)$$

Using Lemma 4.10, we have:

$$\mathfrak{X} \leq 0 \times \mathbb{P}^t - 0 \times \mathbb{P}^{t+1} - \varepsilon_x \beta^t \mathcal{X}_{t+1}^2.$$

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

$$\begin{aligned} \Theta^t - \Theta^{t+1} & \geq \varepsilon_\beta \beta^t \mathcal{B}_{t+1}^2 + \varepsilon_y \beta^t \mathcal{Y}_{t+1}^2 + \varepsilon_x \beta^t \mathcal{X}_{t+1}^2 + \varepsilon_z \beta^t \mathcal{Z}_{t+1}^2 \\ & \geq \min(\varepsilon_y, \varepsilon_x, \varepsilon_z, \varepsilon_\beta) \cdot \beta^t \cdot (\mathcal{B}_{t+1}^2 + \mathcal{X}_{t+1}^2 + \mathcal{Y}_{t+1}^2 + \mathcal{Z}_{t+1}^2). \end{aligned}$$

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

C.10 PROOF OF THEOREM 4.13

Proof. We define $e^t \triangleq \mathcal{B}_t + \mathcal{Z}_t + \mathcal{Y}_t + \mathcal{X}_t$.

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

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

First, we obtain:

$$\begin{aligned} \mathbb{G}_{1/2}^t & \stackrel{\textcircled{1}}{=} \mathbf{G} - \frac{1}{2} \mathbf{X}^t \mathbf{G}^\top \mathbf{X}^t - \frac{1}{2} \mathbf{X}^t [\mathbf{X}^t]^\top \mathbf{G} \\ & \stackrel{\textcircled{2}}{=} (\dot{\mathbf{G}} - \frac{1}{2} \mathbf{X}^t \dot{\mathbf{G}}^\top \mathbf{X}^t - \frac{1}{2} \mathbf{X}^t [\mathbf{X}^t]^\top \dot{\mathbf{G}}) + (\ddot{\mathbf{G}} - \frac{1}{2} \mathbf{X}^t \ddot{\mathbf{G}}^\top \mathbf{X}^t - \frac{1}{2} \mathbf{X}^t [\mathbf{X}^t]^\top \ddot{\mathbf{G}}) \\ & \stackrel{\textcircled{3}}{=} \text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\dot{\mathbf{G}}) + \text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\ddot{\mathbf{G}}) \end{aligned}$$

where step 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 ② uses $\mathbf{G} \in \dot{\mathbf{G}} + \ddot{\mathbf{G}}$; step 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:

$$\begin{aligned} \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\dot{\mathbf{G}})\|_F & = \|\mathbb{G}_{1/2}^t - \text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\ddot{\mathbf{G}})\|_F \\ & \stackrel{\textcircled{1}}{\leq} \|\mathbb{G}_{1/2}^t\|_F + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t}\mathcal{M}}(\ddot{\mathbf{G}})\|_F \\ & \stackrel{\textcircled{2}}{\leq} \|\mathbb{G}_{1/2}^t\|_F + \|\ddot{\mathbf{G}}\|_F \\ & \leq \|\mathbb{G}_{1/2}^t\|_F + \beta^t \bar{\mathcal{A}} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| \\ & \leq \beta^t e^{t+1} + \mathcal{O}(\beta^{t-1} e^t), \end{aligned}$$

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

Finally, we derive:

$$\begin{aligned}
& \frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{X}^t, \check{\mathbf{y}}^t, \mathbf{z}^t) \\
& \stackrel{\textcircled{1}}{=} \frac{1}{T} \sum_{t=1}^T \{ \|\mathcal{A}(\mathbf{X}^t) - \check{\mathbf{y}}^t\| + \|\partial h(\check{\mathbf{y}}^t) - \mathbf{z}^t\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t} \mathcal{M}}(\dot{\mathbf{G}})\|_{\mathbb{F}} \} \\
& \stackrel{\textcircled{2}}{\leq} \frac{1}{T} \sum_{t=1}^T \{ \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\| + \|\check{\mathbf{y}}^t - \mathbf{y}^t\| + \|(1 - \frac{1}{\sigma})(\mathbf{z}^t - \mathbf{z}^{t-1})\| + \|\text{Proj}_{\mathbf{T}_{\mathbf{X}^t} \mathcal{M}}(\dot{\mathbf{G}})\|_{\mathbb{F}} \} \\
& \stackrel{\textcircled{3}}{=} \frac{1}{T} \sum_{t=1}^T \{ \mathcal{O}(\beta^t e^{t+1}) + \mathcal{O}(\beta^{t-1} e^t) \} + \frac{1}{T} \sum_{t=1}^T \|\check{\mathbf{y}}^t - \mathbf{y}^t\| \\
& \stackrel{\textcircled{4}}{=} \frac{1}{T} \sum_{t=1}^T \{ \mathcal{O}(\beta^t e^{t+1}) + \mathcal{O}(\beta^{t-1} e^t) \} + \frac{1}{T} \mathcal{O}(\sum_{t=1}^T \frac{1}{t^p}) \\
& \stackrel{\textcircled{5}}{=} \mathcal{O}(T^{(p-1)/2}) + \mathcal{O}(T^{1-p-1}) \\
& \stackrel{\textcircled{6}}{=} \mathcal{O}(T^{-1/3}),
\end{aligned}$$

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 Lemma 2.5(c) that $\|\check{\mathbf{y}}^t - \mathbf{y}^t\| \leq \mu^t C_h = \mathcal{O}(\frac{1}{t^p})$; step ⑤ uses Lemma A.5 that $\sum_{t=1}^T \frac{1}{t^p} \leq \mathcal{O}(T^{1-p})$, and Lemma 4.7(b) that $\frac{1}{T} \sum_{t=1}^T \beta^t e^{t+1} \leq \mathcal{O}(T^{(p-1)/2})$; step ⑥ uses the choice $p = 1/3$ and Lemma 4.7(b). \square

D PROOFS FOR SECTION 5

D.1 PROOF OF LEMMA 5.4

We begin by presenting the following four useful lemmas.

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

$$(\mathbf{d}_{\mathbf{X}}, \mathbf{d}_{\mathbf{y}}, \mathbf{d}_{\mathbf{z}}, \mathbf{d}_{\beta}) \in \partial \Theta(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t), \quad (50)$$

where $\mathbf{d}_{\mathbf{X}} \triangleq \mathbb{A}^t + \beta^t \mathbf{A}^T(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)$, $\mathbf{d}_{\mathbf{y}} \triangleq \nabla h_{[\tau/\beta^t]}(\mathbf{y}^t) - \mathbf{z}^t + \beta^t(\mathbf{y}^t - \mathcal{A}(\mathbf{X}^t))$, $\mathbf{d}_{\mathbf{z}} \triangleq \mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t$, $\mathbf{d}_{\beta} \triangleq \frac{1}{2} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|_2^2 + \partial_{\beta}(h_{[\tau/\beta^t]}(\mathbf{y}^t)) - \frac{c}{(\beta^t)^2}$, $\mathbb{A}^t \triangleq \partial_{\mathbf{t}} \mathcal{M}(\mathbf{X}^t) + \nabla f(\mathbf{X}^t) + \mathcal{A}^T(\mathbf{z}^t)$.

Proof. We define the Lyapunov function as: $\Theta(\mathbf{X}, \mathbf{y}, \mathbf{z}, \beta) \triangleq L(\mathbf{X}, \mathbf{y}, \mathbf{z}, \beta) + c/\beta$. Using this definition, we can promptly derive the conclusion of the lemma. \square

Lemma D.2. *For OADMM-EP, we define $\{\mathbf{d}_{\mathbf{X}}, \mathbf{d}_{\mathbf{y}}, \mathbf{d}_{\mathbf{z}}, \mathbf{d}_{\beta}\}$ as in Lemma D.1. There exists a constant K such that:*

$$\frac{1}{\beta^t} \{ \|\mathbf{d}_{\mathbf{X}}\|_{\mathbb{F}} + \|\mathbf{d}_{\mathbf{y}}\| + \|\mathbf{d}_{\mathbf{z}}\| + |\mathbf{d}_{\beta}| \} \leq K \{ \mathcal{B}_t + \mathcal{X}^t + \mathcal{Y}^t + \mathcal{Z}^t \}. \quad (51)$$

Here, $\mathcal{B}_t \triangleq \sqrt{(\frac{1}{\beta^{t-1}} - \frac{1}{\beta^t}) \frac{1}{\beta^{t-1}}}$, $\mathcal{X}_t \triangleq \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}}$, $\mathcal{Y}_t \triangleq \|\mathbf{y}^t - \mathbf{y}^{t-1}\|$, and $\mathcal{Z}_t \triangleq \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|$.

Proof. First, we obtain:

$$\begin{aligned}
& \frac{1}{\beta^t} \|\mathbb{A}^t\|_{\mathbb{F}} = \|\partial I_{\mathcal{M}}(\mathbf{X}^t) + \nabla f(\mathbf{X}^t) + \mathcal{A}^T(\mathbf{z}^t)\|_{\mathbb{F}} \\
& \stackrel{\textcircled{1}}{=} \frac{1}{\beta^t} \|\nabla f(\mathbf{X}^t) - \nabla f(\mathbf{X}^{t-1}) - \theta \ell(\beta^{t-1})(\mathbf{X}^t - \mathbf{X}^{t-1}) \\
& \quad + \mathbf{A}^T(\mathbf{z}^t - \mathbf{z}^{t-1}) - \beta^{t-1} \mathbf{A}^T(\mathcal{A}(\mathbf{X}^{t-1}) - \mathbf{y}^{t-1})\|_{\mathbb{F}} \\
& \stackrel{\textcircled{2}}{=} \frac{1}{\beta^t} \|\nabla f(\mathbf{X}^t) - \nabla f(\mathbf{X}^{t-1}) - \theta \ell(\beta^{t-1})(\mathbf{X}^t - \mathbf{X}^{t-1})\|_{\mathbb{F}} \\
& \quad + \frac{1}{\beta^t} \|\sigma \beta^{t-1} \mathbf{A}^T(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t) - \beta^{t-1} \mathbf{A}^T(\mathcal{A}(\mathbf{X}^{t-1}) - \mathbf{y}^{t-1})\|_{\mathbb{F}} \\
& \stackrel{\textcircled{3}}{\leq} \frac{1}{\beta^t} (L_f + \theta \ell(\beta^{t-1})) \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}} \\
& \quad + \beta^{t-1} \|(\sigma - 1) \mathbf{A}^T(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t) + \mathbf{A}^T(\mathcal{A}(\mathbf{X}^{t-1}) - \mathbf{y}^{t-1}) + \mathbf{A}^T(\mathbf{y}^{t-1} - \mathbf{y}^t)\| \\
& = \mathcal{O}(\|\mathbf{X}^t - \mathbf{X}^{t-1}\|_{\mathbb{F}}) + \mathcal{O}(\|\mathbf{y}^t - \mathbf{y}^{t-1}\|) + \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|), \quad (52)
\end{aligned}$$

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

$$\begin{aligned} & \partial I_{\mathcal{M}}(\mathbf{X}^{t+1}) \\ & \ni -\theta \ell(\beta^t)(\mathbf{X}^{t+1} - \mathbf{X}^t) - \nabla_{\mathbf{X}} \mathcal{S}(\mathbf{X}^t, \mathbf{y}^t, \mathbf{z}^t, \beta^t) \\ & = -\theta \ell(\beta^t)(\mathbf{X}^{t+1} - \mathbf{X}^t) - \nabla f(\mathbf{X}^t) - \mathcal{A}^\top[\mathbf{z}^t + \beta^t(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)]; \end{aligned} \quad (53)$$

step 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.11(c); step ③ uses the upper bound of $\|\mathbf{X}^t - \mathbf{X}_c^{t-1}\|_F$, and $\mathbf{z}^t - \mathbf{z}^{t-1} = \sigma \beta^{t-1}(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)$.

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

$$\begin{aligned} \frac{1}{\beta^t} \|\mathbf{d}_{\mathbf{X}}\|_F & \stackrel{\text{①}}{=} \frac{1}{\beta^t} \|\mathbf{A}^t + \beta^t \mathcal{A}^\top(\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t)\|_F \\ & \stackrel{\text{③}}{\leq} \mathcal{O}(\|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F) + \mathcal{O}(\|\mathbf{y}^t - \mathbf{y}^{t-1}\|) + \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|), \end{aligned}$$

where step ① uses the definition of $\mathbf{d}_{\mathbf{X}}$ 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 (52).

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

$$\begin{aligned} \frac{1}{\beta^t} \|\mathbf{d}_{\mathbf{y}}\| & \stackrel{\text{①}}{=} \frac{1}{\beta^t} \|\nabla h_{\tau/\beta^t}(\mathbf{y}^t) - \mathbf{z}^t + \beta^t(\mathbf{y}^t - \mathcal{A}(\mathbf{X}^t))\| \\ & \stackrel{\text{②}}{=} \frac{1}{\beta^t} \|\nabla h_{\tau/\beta^t}(\mathbf{y}^t) - \nabla h_{\tau/\beta^{t-1}}(\mathbf{y}^t) + \beta^t(\mathbf{y}^t - \mathcal{A}(\mathbf{X}^t))\| \\ & \stackrel{\text{③}}{\leq} \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|) + C_h \frac{1}{\beta^t} \left(\frac{\beta^t}{\beta^{t-1}} - 1 \right) \\ & = \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|) + C_h \sqrt{\frac{1}{\beta^{t-1}} - \frac{1}{\beta^t}} \sqrt{\frac{1}{\beta^{t-1}}} \cdot \sqrt{\frac{\beta^{t-1}}{\beta^{t-1}} - \frac{\beta^{t-1}}{\beta^t}} \\ & \stackrel{\text{④}}{\leq} \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|) + C_h \underbrace{\sqrt{\frac{1}{\beta^{t-1}} - \frac{1}{\beta^t}} \sqrt{\frac{1}{\beta^{t-1}}}}_{\triangleq \mathcal{B}_t} \cdot \sqrt{1 - 0}, \end{aligned} \quad (54)$$

where step ① uses the definition of $\mathbf{d}_{\mathbf{y}}$ 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 $\beta^{t-1} = \mathcal{O}(\beta^t)$; step ④ uses $\mathcal{B}_t \triangleq \sqrt{(\frac{1}{\beta^{t-1}} - \frac{1}{\beta^t}) \frac{1}{\beta^{t-1}}}$.

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

$$\begin{aligned} \frac{1}{\beta^t} |\mathbf{d}_{\beta}| & \stackrel{\text{①}}{=} \frac{1}{\beta^t} \cdot \left\{ \frac{1}{2} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|_2^2 + \partial_{\beta}(h_{[\tau/\beta^t]}(\mathbf{y}^t)) - \frac{c}{(\beta^t)^2} \right\} \\ & = \frac{1}{\beta^t} \left\{ \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|) + \partial_{[\tau/\beta^t]}(h_{[\tau/\beta^t]}(\mathbf{y}^t)) \cdot \frac{\tau}{(\beta^t)^2} - \frac{c}{(\beta^t)^2} \right\} \\ & \stackrel{\text{②}}{\leq} \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|) + (\tau C_h^2 + c) \frac{1}{(\beta^t)^3} \\ & \stackrel{\text{③}}{\leq} \mathcal{O}(\|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|) + (\tau C_h^2 + c) \frac{2}{(\xi \beta^0)^2} \cdot \mathcal{B}^t, \end{aligned}$$

where step ① uses the definition of \mathbf{d}_{β} in Lemma D.1; step ② uses Lemma 2.3 that the function $h_{\tau/\beta^t}(\mathbf{y})$ is C_h^2 -Lipschitz continuous w.r.t. $\mu^t \triangleq \tau/\beta^t$; step ③ uses Lemma A.10.

Part (e). We bound the term $\frac{1}{\beta^t} \|\mathbf{d}_{\mathbf{z}}\|_F$. We have: $\frac{1}{\beta^t} \|\mathbf{d}_{\mathbf{z}}\| \leq \frac{1}{\beta^0} \|\mathcal{A}(\mathbf{X}^t) - \mathbf{y}^t\|$.

Part (f). Combining the upper bounds for the terms $\{\frac{1}{\beta^t} \|\mathbf{d}_{\mathbf{X}}\|_F, \frac{1}{\beta^t} \|\mathbf{d}_{\mathbf{y}}\|, \frac{1}{\beta^t} |\mathbf{d}_{\beta}|, \frac{1}{\beta^t} \|\mathbf{d}_{\mathbf{z}}\|\}$, we finish the proof of this lemma. \square

Lemma D.3. For OADMM-RR, we define $\{\mathbf{d}_{\mathbf{X}}, \mathbf{d}_{\mathbf{y}}, \mathbf{d}_{\mathbf{z}}, \mathbf{d}_{\beta}\}$ as in Lemma D.1. There exists a constant K such that :

$$\frac{1}{\beta^t} \{\|\mathbf{d}_{\mathbf{X}}\|_F + \|\mathbf{d}_{\mathbf{y}}\| + \|\mathbf{d}_{\mathbf{z}}\| + |\mathbf{d}_{\beta}|\} \leq K \{\mathcal{B}_t + \mathcal{X}^t + \mathcal{Y}^t + \mathcal{Z}^t\},$$

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

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

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

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

$$\begin{aligned} \frac{1}{\beta^t} \|\mathbf{d}_\mathbf{X}\|_\mathbf{F} &= \frac{1}{\beta^t} \|\partial I_\mathcal{M}(\mathbf{X}^t) + \nabla \dot{\mathcal{L}}(\mathbf{X}^t)\|_\mathbf{F} \\ &\stackrel{\textcircled{1}}{\leq} \frac{1}{\beta^t} \|\nabla \dot{\mathcal{L}}(\mathbf{X}^t) - \mathbf{X}^t [\nabla \dot{\mathcal{L}}(\mathbf{X}^t)]^\top \mathbf{X}^t\|_\mathbf{F} \\ &\stackrel{\textcircled{2}}{=} \frac{1}{\beta^t} \|\mathbf{G}^t - \mathbf{X}^t [\mathbf{G}^t]^\top \mathbf{X}^t\|_\mathbf{F} = \frac{1}{\beta^t} \|\mathbb{G}_1^t\|_\mathbf{F} \\ &\stackrel{\textcircled{3}}{\leq} \frac{1}{\beta^t} \max(1, 1/\rho) \cdot \|\mathbb{G}_{1/2}\|_\mathbf{F} = \mathcal{O}(\mathcal{X}_t), \end{aligned} \quad (55)$$

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 $\|\mathbb{G}_1\|_\mathbf{F} \leq \max(1, 1/\rho) \|\mathbb{G}_\rho\|_\mathbf{F}$, as shown in Lemma 2.12(b).

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

Part (c). Combining the upper bounds for the terms $\{\frac{1}{\beta^t} \|\mathbf{d}_\mathbf{X}\|_\mathbf{F}, \frac{1}{\beta^t} \|\mathbf{d}_\mathbf{y}\|, \frac{1}{\beta^t} |\mathbf{d}_\beta|, \frac{1}{\beta^t} \|\mathbf{d}_\mathbf{z}\|_\mathbf{F}\}$, we finish the proof of this lemma. \square

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

Lemma D.4. (Subgradient Bounds) For both OADMM-EP and OADMM-RR, there exists a constant $K > 0$ such that:

$$\text{dist}(\mathbf{0}, \partial\Theta(\mathcal{W}^t)) \leq \beta^t K e^t.$$

Here, $\text{dist}^2(\mathbf{0}, \partial\Theta(\mathcal{W})) \triangleq \text{dist}^2(\mathbf{0}, \partial_\beta\Theta(\mathcal{W})) + \text{dist}^2(\mathbf{0}, \partial_\mathbf{X}\Theta(\mathcal{W})) + \text{dist}^2(\mathbf{0}, \partial_\mathbf{y}\Theta(\mathcal{W})) + \text{dist}^2(\mathbf{0}, \partial_\mathbf{z}\Theta(\mathcal{W}))$.

Proof. For both OADMM-EP and OADMM-RR, we have:

$$\begin{aligned} \text{dist}(\mathbf{0}, \partial\Theta(\mathcal{W}^t)) &= \sqrt{\|\mathbf{d}_\mathbf{X}\|_\mathbf{F}^2 + \|\mathbf{d}_\mathbf{y}\|^2 + \|\mathbf{d}_\mathbf{z}\|_\mathbf{F}^2 + |\mathbf{d}_\beta|^2} \\ &\stackrel{\textcircled{1}}{\leq} \|\mathbf{d}_\mathbf{X}\|_\mathbf{F} + \|\mathbf{d}_\mathbf{y}\| + |\mathbf{d}_\beta| + \|\mathbf{d}_\mathbf{z}\|_\mathbf{F} \\ &\stackrel{\textcircled{2}}{\leq} K\beta^t\{\mathcal{X}^t + \mathcal{Y}^t + \mathcal{B}^t + \mathcal{Z}^t\}, \end{aligned}$$

where step ① uses $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ for all $a, b \geq 0$; step ② uses Lemmas D.2 and D.3. \square

D.2 PROOF OF THEOREM 5.6

Proof. We define $K' \triangleq 4K / \min(\varepsilon_x, \varepsilon_y, \varepsilon_z, \varepsilon_\beta)$.

Firstly, using Assumption 5.1, we have:

$$\varphi'(\Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^\infty)) \cdot \text{dist}(\mathbf{0}, \partial\Theta(\mathcal{W}^t)) \geq 1. \quad (56)$$

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 $a = \Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^\infty)$ and $b = \Theta(\mathcal{W}^{t+1}) - \Theta(\mathcal{W}^\infty)$, we have:

$$\begin{aligned} &(\Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^{t+1})) \cdot \varphi'(\Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^\infty)) \\ &\leq \underbrace{\varphi(\Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^\infty))}_{\triangleq \varphi^t} - \underbrace{\varphi(\Theta(\mathcal{W}^{t+1}) - \Theta(\mathcal{W}^\infty))}_{\triangleq \varphi^{t+1}}. \end{aligned} \quad (57)$$

Third, we define $\mathcal{X}_t \triangleq \|\mathbf{X}^t - \mathbf{X}^{t-1}\|_F$, and derive the following inequalities for OADMM-EP:

$$\begin{aligned}
& \min(\varepsilon_z, \varepsilon_y, \varepsilon_x, \varepsilon_\beta) \beta^t \{\mathcal{B}_{t+1}^2 + \mathcal{Z}_{t+1}^2 + \mathcal{Y}_{t+1}^2 + \mathcal{X}_{t+1}^2\} \\
& \stackrel{\textcircled{1}}{\leq} \Theta^t - \Theta^{t+1} = \Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^{t+1}) \\
& \stackrel{\textcircled{2}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot \frac{1}{\varphi'(\Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^\infty))} \\
& \stackrel{\textcircled{3}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot \text{dist}(\mathbf{0}, \partial\Theta(\mathcal{W}^t)) \\
& \stackrel{\textcircled{4}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot K \beta^t (\mathcal{B}_t + \mathcal{Z}_t + \mathcal{Y}_t + \mathcal{X}_t), \tag{58}
\end{aligned}$$

where step ① uses Lemma 4.7; step ③ uses Inequality (57); step ④ uses Inequality (56); step ⑤ uses Lemma 5.4. We further derive the following inequalities:

$$\begin{aligned}
& (\mathcal{B}_{t+1} + \mathcal{Z}_{t+1} + \mathcal{Y}_{t+1} + \mathcal{X}_{t+1})^2 \\
& \stackrel{\textcircled{1}}{\leq} 4 \cdot \{\mathcal{B}_{t+1}^2 + \mathcal{Z}_{t+1}^2 + \mathcal{Y}_{t+1}^2 + \mathcal{X}_{t+1}^2\} \\
& \stackrel{\textcircled{2}}{\leq} \{4K / \min(\varepsilon_z, \varepsilon_y, \varepsilon_x, \varepsilon_\beta)\} \cdot (\varphi^t - \varphi^{t+1}) \cdot (\mathcal{B}_t + \mathcal{Z}_t + \mathcal{Y}_t + \mathcal{X}_t), \tag{59}
\end{aligned}$$

where step ① uses the norm inequality that $(a + b + c + d)^2 \leq 4(a^2 + b^2 + c^2 + d^2)$ for any $a, b, c, d \in \mathbb{R}$; step ② uses Inequality (58).

Fourth, we define $\mathcal{X}_t \triangleq \|\frac{1}{\beta} \mathbb{G}_{1/2}^t\|_F$, and derive the following inequalities for OADMM-RR:

$$\begin{aligned}
& \min(\varepsilon_z, \varepsilon_y, \varepsilon_x, \varepsilon_\beta) \beta^t \{\mathcal{B}_{t+1}^2 + \mathcal{Z}_{t+1}^2 + \mathcal{Y}_{t+1}^2 + \mathcal{X}_{t+1}^2\} \\
& \stackrel{\textcircled{1}}{\leq} \Theta^t - \Theta^{t+1} = \Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^{t+1}) \\
& \stackrel{\textcircled{2}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot \frac{1}{\varphi'(\Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^\infty))} \\
& \stackrel{\textcircled{3}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot \text{dist}(\mathbf{0}, \partial\Theta(\mathcal{W}^t)) \\
& \stackrel{\textcircled{4}}{\leq} (\varphi^t - \varphi^{t+1}) \cdot K \beta^t (\varphi^t - \varphi^{t+1}) \cdot (\mathcal{B}_t + \mathcal{Z}_t + \mathcal{Y}_t + \mathcal{X}_t), \tag{60}
\end{aligned}$$

where step ① uses Lemma 4.12; step ② uses Inequality (57); step ③ uses Inequality (56); step ④ uses Lemma 5.4. Furthermore, using the similar strategies as in deriving Inequality (59), we have:

$$\begin{aligned}
& (\mathcal{B}_{t+1} + \mathcal{Z}_{t+1} + \mathcal{Y}_{t+1} + \mathcal{X}_{t+1})^2 \\
& \leq \{4K / \min(\varepsilon_z, \varepsilon_y, \varepsilon_x, \varepsilon_\beta)\} \cdot (\varphi^t - \varphi^{t+1}) \cdot (\mathcal{B}_t + \mathcal{Z}_t + \mathcal{Y}_t + \mathcal{X}_t). \tag{61}
\end{aligned}$$

Part (a). We define $e^t \triangleq \mathcal{B}_t + \mathcal{X}_t + \mathcal{Y}_t + \mathcal{Z}_t$. Given Inequalities (59) and (61), we establish the following unified inequality applicable to both OADMM-EP and OADMM-RR:

$$(e^{t+1})^2 \leq e^t \cdot \underbrace{\{4K / \min(\varepsilon_z, \varepsilon_y, \varepsilon_x, \varepsilon_\beta)\}}_{\triangleq K'} \cdot (\varphi^t - \varphi^{t+1}). \tag{62}$$

Part (b). Considering Inequality (62) and applying Lemma A.14 with $p^t \triangleq K' \varphi^t$, we have:

$$\forall t, \sum_{i=t}^{\infty} e^{i+1} \leq e^t + 2K' \varphi^t.$$

Letting $t = 1$, we have: $\sum_{i=1}^{\infty} e^{i+1} \leq e^1 + 2K' \varphi^1$.

□

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 $t \geq 1$: $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F \stackrel{\textcircled{1}}{\leq} \sum_{i=t}^{\infty} \|\mathbf{X}^i - \mathbf{X}^{i+1}\|_F \leq \sum_{i=t}^{\infty} \left\{ \sqrt{\left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}}\right) \frac{1}{\beta^t}} + \|\mathbf{X}^{i+1} - \mathbf{X}^i\|_F + \|\mathbf{y}^{i+1} - \mathbf{y}^i\| + \|\mathcal{A}(\mathbf{X}^{i+1}) - \mathbf{y}^{i+1}\| \right\} = \sum_{i=t}^{\infty} e^{i+1} \triangleq d^t$, where step ① use the triangle inequality.

Part (a-ii). For OADMM-RR, we have: $\|\mathbf{X}^{t+1} - \mathbf{X}^t\|_F \stackrel{\textcircled{1}}{=} \|\text{Retr}_{\mathbf{X}^t}(-\eta^t \mathbb{G}_\rho^t) - \mathbf{X}^t\|_F \stackrel{\textcircled{2}}{\leq} \dot{k} \|\eta^t \mathbb{G}_\rho^t\|_F \stackrel{\textcircled{3}}{\leq} \dot{k} \eta^t \max(2\rho, 1) \|\mathbb{G}_{1/2}^t\|_F \stackrel{\textcircled{4}}{=} \dot{k} \max(2\rho, 1) \frac{b^t \gamma^j}{\beta^t} \|\mathbb{G}_{1/2}^t\|_F \stackrel{\textcircled{5}}{\leq} \dot{k} \max(2\rho, 1) \bar{b} \bar{\gamma} \cdot \|\frac{1}{\beta^t} \mathbb{G}_{1/2}^t\|_F = \mathcal{O}(\|\frac{1}{\beta^t} \mathbb{G}_{1/2}^t\|_F)$, where step ① uses the update rule of \mathbf{X}^{t+1} ; step ② uses Lemma 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 the fact that $\gamma^j \leq \bar{\gamma}$. Furthermore, we derive for all $t \geq 1$: $\|\mathbf{X}^t - \mathbf{X}^\infty\|_F \leq \sum_{i=t}^\infty \|\mathbf{X}^i - \mathbf{X}^{i+1}\|_F \leq \mathcal{O}(\sum_{i=t}^\infty \|\frac{1}{\beta^i} \mathbb{G}_{1/2}^i\|_F) \leq \mathcal{O}(\sum_{i=t}^\infty e^{i+1}) = \mathcal{O}(d^t)$.

Part (b). We define $\varphi^t \triangleq \varphi(s^t)$, where $s^t \triangleq \Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^\infty)$. We derive:

$$\begin{aligned}
d^t &\triangleq \sum_{i=t}^\infty e^{i+1} \\
&\stackrel{\textcircled{1}}{\leq} e^t + 2K' \varphi^t \\
&\stackrel{\textcircled{2}}{=} e^t + 2K' \tilde{c} \cdot [(s^t)^{\tilde{\sigma}}]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\
&\stackrel{\textcircled{3}}{=} e^t + 2K' \tilde{c} \cdot [\tilde{c}(1-\tilde{\sigma}) \cdot \frac{1}{\varphi'(s^t)}]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\
&\stackrel{\textcircled{4}}{\leq} e^t + 2K' \tilde{c} \cdot [\tilde{c}(1-\tilde{\sigma}) \cdot \text{dist}(\mathbf{0}, \partial\Theta(\mathcal{W}^t))]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\
&\stackrel{\textcircled{5}}{\leq} e^t + 2K' \tilde{c} \cdot [\tilde{c}(1-\tilde{\sigma}) \cdot \beta^t K e^t]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\
&\stackrel{\textcircled{6}}{=} d^{t-1} - d^t + 2K' \tilde{c} \cdot \{\tilde{c}(1-\tilde{\sigma}) \cdot \beta^t K (d^{t-1} - d^t)\}^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}} \\
&= d^{t-1} - d^t + \underbrace{2K' \tilde{c} \cdot [\tilde{c}(1-\tilde{\sigma}) K]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}}_{\triangleq \nu} \cdot [\beta^t (d^{t-1} - d^t)]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}},
\end{aligned}$$

where step ① uses $\sum_{i=t}^\infty e^{i+1} \leq e^t + e^{t-1} + 4K' \varphi^t$, as shown in Theorem 5.6(b); step ② uses the definitions that $\varphi^t \triangleq \varphi(s^t)$, $s^t \triangleq \Theta(\mathcal{W}^t) - \Theta(\mathcal{W}^\infty)$, and $\varphi(s) = \tilde{c} s^{1-\tilde{\sigma}}$; step ③ 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 that $1 \leq \text{dist}(\mathbf{0}, \partial\Theta(\mathcal{W}^t)) \cdot \varphi'(s^t)$; step ⑤ uses $\text{dist}(\mathbf{0}, \partial\Theta(\mathcal{W}^t)) \leq \beta^t K e^t$ for both OADMM-EP and OADMM-RR, as shown in Lemma 5.4; step ⑥ uses the fact that $e^t = d^{t-1} - d^t$. □

D.4 PROOF OF THEOREM 5.9

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

$$d^t \leq d^{t-1} - d^t + \nu \cdot [\beta^t (d^{t-1} - d^t)]^{\frac{1-\tilde{\sigma}}{\tilde{\sigma}}}. \quad (63)$$

By the definition of $e^t \triangleq d^{t-1} - d^t \triangleq \mathcal{B}_t + \mathcal{X}_t + \mathcal{Y}_t + \mathcal{Z}_t$, there exists a universal constant $\bar{e} > 0$ such that $e^t \leq \bar{e}$.

We consider three cases for Inequality (63).

Part (a). $\tilde{\sigma} \in (0, \frac{1}{2}]$. We let $u = \frac{1-\tilde{\sigma}}{\tilde{\sigma}} \in [1, \infty)$, and $\zeta = (1-p)u = \frac{2}{3} \cdot \frac{1-\tilde{\sigma}}{\tilde{\sigma}} \in [\frac{2}{3}, \infty]$.

From Inequality (63), we obtain:

$$\begin{aligned}
d^t &\leq d^{t-1} - d^t + \nu \cdot [\beta^t (d^{t-1} - d^t)]^u \\
&\stackrel{\textcircled{1}}{\leq} d^{t-1} - d^t + \nu' \cdot t^{pu} \cdot [d^{t-1} - d^t]^u \\
&\stackrel{\textcircled{2}}{\leq} \mathcal{O}(t^{-\zeta}),
\end{aligned}$$

where step ① uses $\beta^t \leq \nu' t^p$, where ν' is some certain constant; step ② uses Lemma A.15 with $c = \nu'$ and $u > 1$.

Part (b). $\tilde{\sigma} \in (\frac{1}{2}, 1)$. We let $u = \frac{1-\tilde{\sigma}}{\tilde{\sigma}} \in (0, 1)$, and $\zeta = \frac{1-p}{1-u-1} = \frac{2}{3} \cdot \frac{1-\tilde{\sigma}}{2\tilde{\sigma}-1} \in (0, \infty)$.

We have from Inequality (63):

$$\begin{aligned}
d^t &\leq \nu[\beta^t]^u \cdot (d^{t-1} - d^t)^u + d^{t-1} - d^t \\
&\stackrel{\textcircled{1}}{\leq} \nu[\beta^t]^u \cdot (d^{t-1} - d^t)^u + (d^{t-1} - d^t)^u \cdot \bar{e}^{1-u} \\
&\stackrel{\textcircled{2}}{\leq} \nu[\beta^t]^u \cdot (d^{t-1} - d^t)^u + (d^{t-1} - d^t)^u \cdot \bar{e}^{1-u} \cdot (\frac{\beta^t}{\beta^0})^u \\
&= [\nu + (\beta^0)^{-u}] \cdot [\beta^t]^u \cdot (d^{t-1} - d^t)^u \\
&\stackrel{\textcircled{3}}{=} \nu' t^{pu} \cdot (d^{t-1} - d^t)^u \\
&\stackrel{\textcircled{4}}{\leq} \mathcal{O}(t^{-\zeta}),
\end{aligned}$$

where step ① uses the fact that $\max_{x \in (0, \bar{e}]} \frac{x}{x^u} \leq \bar{e}^{1-u}$ if $u \in (0, 1)$; step ② uses $\beta^0 \leq \beta^t$ and $u \in (0, 1)$; step ③ uses the fact that $\beta^t \leq \nu' t^p$ for some certain constant ν' , and the choice that $\nu' = [\nu + (\beta^0)^{-u}] \nu''$; step ④ uses Lemma A.16 with $c = \nu'$.

Part (c). $\tilde{\sigma} \in (\frac{1}{4}, \frac{1}{2}]$. We let $u = \frac{1-\tilde{\sigma}}{\tilde{\sigma}} \in [1, 3)$, and $\zeta = 1 - pu = \frac{4\tilde{\sigma}-1}{3\tilde{\sigma}} \in (0, \frac{2}{3}]$.

We have from Inequality (63):

$$\begin{aligned}
d^t &\leq d^{t-1} - d^t + \nu \cdot [\beta^t]^u \cdot [d^{t-1} - d^t]^u \\
&\stackrel{\textcircled{1}}{\leq} (d^{t-1} - d^t) \cdot (1 + \bar{e}^{u-1} \nu \cdot [\beta^t]^u) \\
&\stackrel{\textcircled{2}}{\leq} (d^{t-1} - d^t) \cdot \nu' t^{pu} \\
&\leq d^{t-1} \cdot \frac{\nu' t^{pu}}{\nu' t^{pu} + 1} \\
&\stackrel{\textcircled{3}}{\leq} \mathcal{O}(1/\exp(t^\zeta)),
\end{aligned}$$

where step ① uses the fact that $\frac{x^u}{x} \leq \bar{e}^{u-1}$ for all $u \geq 1$, and $x = d^{t-1} - d^t = e^t \leq \bar{e}$; step ② uses $(1 + \bar{e}^{u-1} \nu \cdot [\beta^t]^u) \leq \nu' t^{pu}$ with ν' being some certain constant, which is implied by $\beta^t = \mathcal{O}(t^p)$; step ③ uses Lemma A.17 with $c = \nu'$ and $q = pu$.

Part (d). 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 $\mathbf{D} \in \mathbb{R}^{\hat{m} \times \hat{d}}$. They are denoted as follows: ‘mnist- \hat{m} - \hat{d} ’, ‘TDT2- \hat{m} - \hat{d} ’, ‘sector- \hat{m} - \hat{d} ’, and ‘randn- \hat{m} - \hat{d} ’, where $\text{randn}(m, n)$ generates a standard Gaussian random matrix of size $m \times n$. The construction of $\mathbf{D} \in \mathbb{R}^{\hat{m} \times \hat{d}}$ involves randomly selecting \hat{m} examples and \hat{d} dimensions from the original 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 $\mathbf{D} \leftarrow \mathbf{D} - \mathbf{1}\mathbf{1}^\top \mathbf{D}$.

► **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.

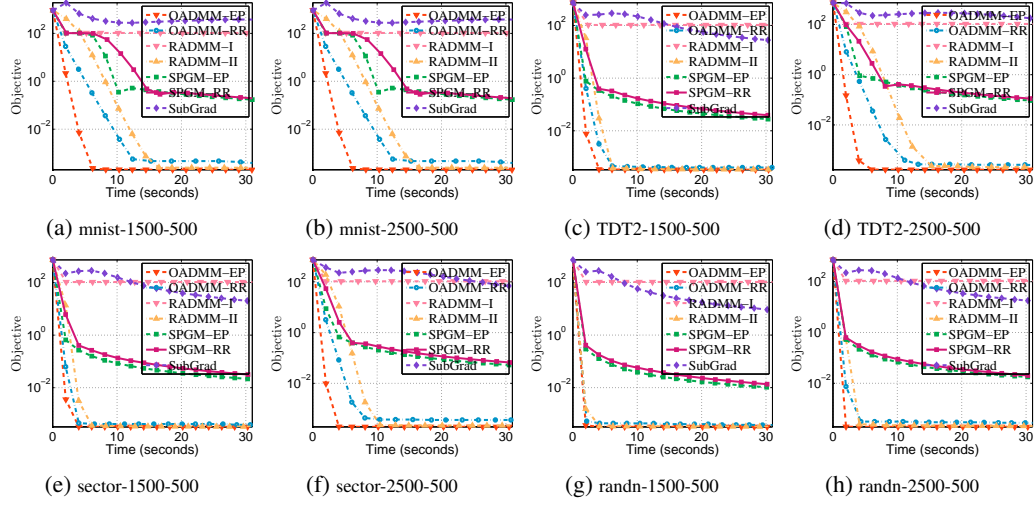


Figure 3: The convergence curve of the compared methods with $\dot{\rho} = 10$.

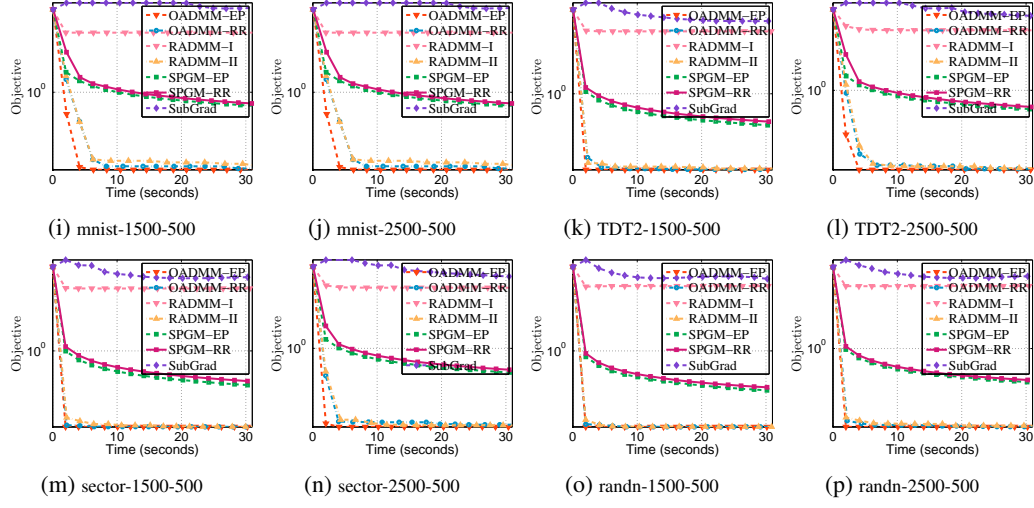


Figure 4: The convergence curve of the compared methods with $\dot{\rho} = 100$.

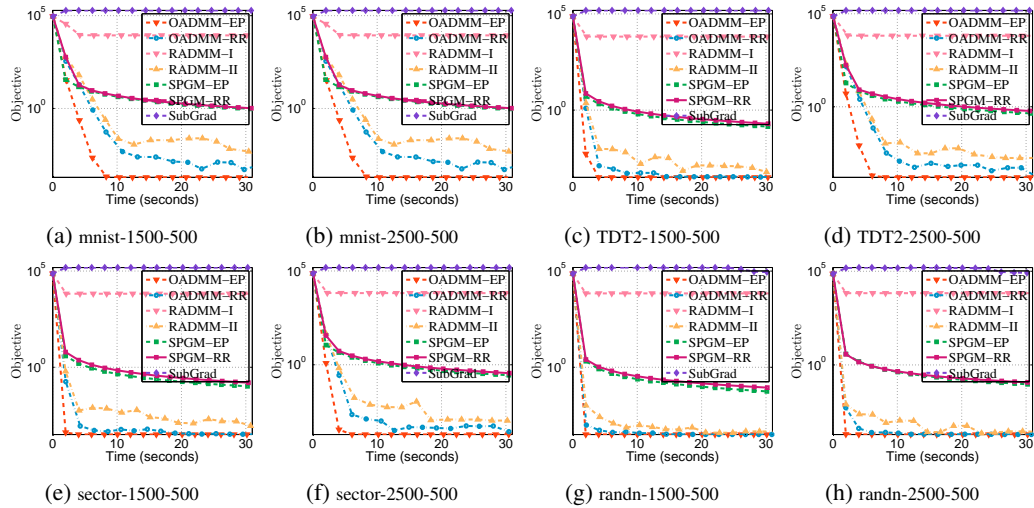


Figure 5: The convergence curve of the compared methods with $\dot{\rho} = 1000$.