

000 001 002 003 004 005 006 007 008 009 010 ADMM FOR NONCONVEX OPTIMIZATION UNDER MINIMAL CONTINUITY ASSUMPTION

005 **Anonymous authors**

006 Paper under double-blind review

ABSTRACT

011 This paper introduces a novel approach to solving multi-block nonconvex com-
 012 posite optimization problems through a proximal linearized Alternating Direction
 013 Method of Multipliers (ADMM). This method incorporates an Increasing Pen-
 014 alization and Decreasing Smoothing (IPDS) strategy. Distinguishing itself from
 015 existing ADMM-style algorithms, our approach (denoted IPDS-ADMM) impos-
 016 es a less stringent condition, specifically requiring continuity in just one block
 017 of the objective function. IPDS-ADMM requires that the penalty increases and
 018 the smoothing parameter decreases, both at a controlled pace. When the asso-
 019 ciated linear operator is bijective, IPDS-ADMM uses an over-relaxation stepsize
 020 for faster convergence; however, when the linear operator is surjective, IPDS-
 021 ADMM uses an under-relaxation stepsize for global convergence. We devise a
 022 novel potential function to facilitate our convergence analysis and prove an oracle
 023 complexity $\mathcal{O}(\epsilon^{-3})$ to achieve an ϵ -approximate critical point. To the best of our
 024 knowledge, this is the first complexity result for using ADMM to solve this class
 025 of nonsmooth nonconvex problems. Finally, some experiments on the sparse PCA
 026 problem are conducted to demonstrate the effectiveness of our approach.

027 1 INTRODUCTION

029 We consider the following multi-block nonconvex nonsmooth composite optimization problem:

$$031 \quad \min_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n} \sum_{i=1}^n [f_i(\mathbf{x}_i) + h_i(\mathbf{x}_i)], \text{ s.t. } [\sum_{i=1}^n \mathbf{A}_i \mathbf{x}_i] = \mathbf{b}, \quad (1)$$

034 where $\mathbf{b} \in \mathbb{R}^{m \times 1}$, $\mathbf{A}_i \in \mathbb{R}^{m \times d_i}$, $\mathbf{x}_i \in \mathbb{R}^{d_i \times 1}$, and $i \in [n] \triangleq \{1, 2, \dots, n\}$. We assume
 035 $f_i(\cdot) : \mathbb{R}^{d_i \times 1} \mapsto (-\infty, \infty)$ is differentiable and potentially nonconvex for all $i \in [n]$. The function
 036 $h_i(\cdot) : \mathbb{R}^{d_i \times 1} \mapsto (-\infty, \infty]$ is assumed to be closed, proper, lower semi-continuous, and poten-
 037 tially nonsmooth. While $h_n(\cdot)$ is convex, we do not require convexity for $h_i(\cdot)$ for $i \in [n-1]$.
 038 Additionally, we assume the nonconvex proximal operator of $h_i(\cdot)$ is easy to compute for all $i \in [n]$.

039 Problem (1) has a wide range of applications in machine learning. The function $f_i(\cdot)$ plays a cru-
 040 cial role in handling empirical loss, including neural network activation functions (Liu et al., 2022;
 041 Zeng et al., 2021; Wang et al., 2019a; Huang et al., 2019). Incorporating multiple nonsmooth reg-
 042 ularization terms $h_i(\cdot)$ enables diverse prior information integration, including structured sparsity,
 043 low-rank, binary, orthogonality, and non-negativity constraints, enhancing regularization model accu-
 044 racy. These capabilities extend to various applications such as sparse PCA, overlapping group
 045 Lasso, graph-guided fused Lasso, and phase retrieval.

046 ▶ **ADMM Literature.** The Alternating Direction Method of Multipliers (ADMM) is a versatile
 047 optimization tool suitable for solving composite constrained problems as in Problem (1), which
 048 pose challenges for other standard optimization methods, such as the accelerated proximal gradient
 049 method (Nesterov, 2003) and the augmented Lagrangian method (Zeng et al., 2022; Lu & Zhang,
 050 2012; Zhu et al., 2023; Lin et al., 2022). The standard ADMM was initially introduced in (Gabay
 051 & Mercier, 1976), and its complexity analysis for the convex settings was first conducted in (He &
 052 Yuan, 2012; Monteiro & Svaiter, 2013). Since then, numerous papers have explored the iteration
 053 complexity of ADMM in diverse settings. These settings include acceleration through multi-step
 054 updates (Pock & Sabach, 2016; Li et al., 2016; Ouyang et al., 2015; Shen et al., 2017; Tran Dinh,

Table 1: Comparison of existing nonconvex ADMM approaches. CVX: convex. NC: nonconvex. LCONT: Lipschitz continuous. WC: weakly convex. RWC: restricted weakly convex. \mathbb{I} : \mathbf{A}_n is identity. SU: \mathbf{A}_n is surjective with $\lambda_{\min}(\mathbf{A}_n \mathbf{A}_n^\top) > 0$. IN: \mathbf{A}_n is injective with $\lambda_{\min}(\mathbf{A}_n^\top \mathbf{A}_n) > 0$. BI: \mathbf{A}_n is bijective (both surjective and injective). IM: $\text{Im}([\mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_{n-1}]) \subseteq \text{Im}(\mathbf{A}_n)$ with Im being the image of the matrix.

Reference	Optimization Problems and Main Assumptions			Complexity	Parameter σ
	Blocks	Functions $f_i(\cdot)$ and $h_i(\cdot)^a$	Matrices \mathbf{A}_i		
(He & Yuan, 2012)	$n = 2$	CVX: $f_1, h_1, \forall i \in [2]$	feasible	$\mathcal{O}(\epsilon^{-2})^b$	$\sigma = 1$
(Li & Pong, 2015)	$n = 2$	NC: $h_1, f_2; f_1 = h_2 = 0$	SU	$\mathcal{O}(\epsilon^{-2})$	$\sigma = 1$
(Yang et al., 2017) ^c	$n = 3$	CVX: $h_1, f_3; \text{NC: } h_2, f_1 = f_2 = h_3 = 0$	\mathbb{I}	$\mathcal{O}(\epsilon^{-2})$	$\sigma \in [1, 2)$
(Yashtini, 2022)	$n = 2$	NC: $f_{[1,2]}, h_{[1,2]}; h_2 = 0$	BI	$\mathcal{O}(\epsilon^{-2})$	$\sigma \in (0, 1)$
(Yashtini, 2021)	$n \geq 2$	WC: $f_{[1,n-1]}, h_{[1,n]} = 0$	BI, IM	$\mathcal{O}(\epsilon^{-2})$	$\sigma \in (0, 1)$
(Wang et al., 2019b)	$n \geq 2$	RWC: $h_{[1,n-1]}, h_n = 0$	IN, IM	$\mathcal{O}(\epsilon^{-2})$	$\sigma = 1$
(Bȯt & Nguyen, 2020)	$n = 2$	NC: $h_{[1,n]}, f_{[1,n]}; f_1 = h_2 = 0$	\mathbb{I}	$\mathcal{O}(\epsilon^{-2})$	$\sigma \in [1, 2)$
(Bȯt et al., 2019)	$n = 2$	NC: $h_{[1,n]}, f_{[1,n]}; f_1 = h_2 = 0$	SU	$\mathcal{O}(\epsilon^{-2})$	$\sigma \in (0, 1)$
(Huang et al., 2019)	$n \geq 2$	CVX: $h_{[1,n]}; h_n = 0$	BI	$\mathcal{O}(\epsilon^{-2})$	$\sigma = 1$
(Li et al., 2022) ^d	$n = 2$	NC: $f_1, h_1; \text{CVX: } h_2; f_2 = 0; \text{LCONT: } h_2$	\mathbb{I}	$\mathcal{O}(\epsilon^{-4})$	$\sigma = 1$
Ours	$n \geq 2$	NC: $h_{[1,n-1]}, f_{[1,n]}; \text{CVX: } h_n; \text{LCONT: } h_n, f_n$	BI	$\mathcal{O}(\epsilon^{-3})$	$\sigma \in [1, 2)$
Ours	$n \geq 2$	NC: $h_{[1,n-1]}, f_{[1,n]}; \text{CVX: } h_n; \text{LCONT: } h_n, f_n$	SU	$\mathcal{O}(\epsilon^{-3})$	$\sigma \in (0, 1)$

Note *a*: The notation $h_n = 0$ indicates that, for the n -th block, the non-smooth part is absent and the objective function is smooth.

Note *b*: The iteration complexity relies on the variational inequality of the convex problem.

Note *c*: We adapt their application model into our optimization framework in Equation (1) with $(L, S, Z) = (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$, as their model additionally requires the linear operator for the other two blocks to be injective.

Note *d*: This paper focuses manifold optimization problem with a fixed large penalty and a fixed small stepsize.

2018), asynchronous updates (Zhang & Kwok, 2014), Jacobi updates (Deng et al., 2017), non-Euclidean proximal updates (Gonçalves et al., 2017b), and extensions to handle more specific or general functions such as strongly convex functions (Nishihara et al., 2015; Lin et al., 2015a; Ouyang et al., 2015), nonlinear constrained functions (Lin et al., 2022), and multi-block composite functions (Lin et al., 2015b; Xu et al., 2017).

► **Nonconvex ADMM.** Compared to the classical Subgradient Methods (Li et al., 2021; Davis & Drusvyatskiy, 2019) and Smoothing Proximal Gradient Methods (SPGM) (Böhm & Wright, 2021), designed for general nonconvex optimization, ADMM-type methods potentially offer faster convergence, enhanced parallelization, and greater numerical stability. However, the convergence analysis of the nonconvex ADMM is challenging due to the absence of Fejér monotonicity in iterations. In the past decade, significant research has focused on exploring various nonconvex ADMM variants (Li & Pong, 2015; Hong et al., 2016; Yang et al., 2017). (Li & Pong, 2015) establishes the convergence of a class of nonconvex problems when a specific potential function associated with the augmented Lagrangian satisfies the Kurdyka-Łojasiewicz inequality. (Yang et al., 2017) analyzes ADMM variants for solving low-rank and sparse optimization problems. (Hong et al., 2016) investigates ADMM variants for nonconvex consensus and sharing problems. Some researchers have examined ADMM variants under weaker conditions, such as restricted weak convexity (Wang et al., 2019b), restricted strong convexity (Barber & Sidky, 2024), and the Hoffman error bound (Zhang & Luo, 2020). However, existing methods all assume the smoothness of at least one block. In contrast, our approach imposes the fewest conditions on the objective function by employing an Increasing Penalization and Decreasing Smoothing (IPDS) strategy.

► **Over-Relaxed and Under-Relaxed ADMM.** Prior studies have analyzed ADMM using either under-relaxation stepsizes $\sigma \in (0, 1)$, or over-relaxation stepsizes $\sigma \in [1, 2)$, for updating the dual variable. This contrasts with earlier approaches that employed fixed values such as 1 or the golden ratio $(\sqrt{5} + 1)/2$. In nonconvex settings, most existing works require that the associated matrix of the problem be bijective (Gonçalves et al., 2017a; Yang et al., 2017; Yashtini, 2022; 2021; Bȯt & Nguyen, 2020). However, the work of (Bȯt et al., 2019) demonstrates that ADMM can still be applied when the associated matrix is surjective, provided that an under-relaxation stepsize is employed. Inspired by these findings, our work shows that when the associated linear operator is bijective, IPDS-ADMM uses an over-relaxation stepsize for faster convergence. In contrast, when the linear operator is surjective, we employ under-relaxation stepsizes to achieve global convergence.

► **Other Works on Accelerating ADMM.** Significant research interest has focused on accelerating ADMM for nonconvex problems. The work by (Hien et al., 2022) explores the use of an inertial force, an approach further investigated in studies by (Pock & Sabach, 2016; Le et al., 2020; Bȯt

et al., 2023; Phan & Gillis, 2023), to enhance the performance of nonconvex ADMM. Additionally, studies by (Huang et al., 2019; Bian et al., 2021; Liu et al., 2020) have employed variance-reduced stochastic gradient descent to decrease the incremental first-order oracle complexity in addressing composite problems characterized by finite-sum structures.

► **Existing Challenges.** We consider the linearly-constrained optimization problem in Problem (1), which involves $(n - 1)$ potentially nonsmooth, nonconvex, and non-Lipschitz composite functions $h_i(\cdot)$ for $i \in [n - 1]$, and one convex, non-smooth composite function $h_n(\cdot)$. Existing ADMM-type methods are unable to solve this problem as they require at least one of the composite functions to be smooth (i.e., $h_n(\cdot) = 0$). In the special case where $\mathbf{A}_n = \mathbf{I}$ and $h_n(\cdot)$ is the indicator function of orthogonality constraints, the Riemannian ADMM (RADMM) algorithm (Li et al., 2022) can solve Problem (1). However, its iteration complexity is suboptimal compared to our method, and it is unable to handle linearly-constrained problems, particularly when \mathbf{A}_n is subjective. We make a comparison of existing nonconvex ADMM approaches in Table 1.

► **Our Contributions.** Our main contributions are summarized as follows. (i) We introduce IPDS-ADMM to solve the nonconvex nonsmooth optimization problem as in Problem (1). This approach imposes the least stringent condition, specifically requiring continuity in only one block of the objective function. It employs an Increasing Penalization and Decreasing Smoothing (IPDS) strategy to ensure convergence (See Section 3). (ii) IPDS-ADMM achieves global convergence when the associated matrix is either bijective or surjective. We establish that IPDS-ADMM converges to an ϵ -critical point with a time complexity of $\mathcal{O}(1/\epsilon^3)$ (See Section 4). (iii) We have conducted experiments on the sparse PCA problem to demonstrate the effectiveness of our approach. (See Section 5).

► **Assumptions.** Through this paper, we impose the following assumptions on Problem (1).

Assumption 1.1. Each function $f_i(\cdot)$ is L_i -smooth for all $i \in [n]$ such that $\|\nabla f_i(\mathbf{x}_i) - \nabla f_i(\hat{\mathbf{x}}_i)\| \leq L_i \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|$ holds for all $\mathbf{x}_i \in \mathbb{R}^{d_i \times 1}$ and $\hat{\mathbf{x}}_i \in \mathbb{R}^{d_i \times 1}$. This implies that $|f_i(\mathbf{x}_i) - f_i(\hat{\mathbf{x}}_i) - \langle \nabla f_i(\hat{\mathbf{x}}_i), \mathbf{x}_i - \hat{\mathbf{x}}_i \rangle| \leq \frac{L_i}{2} \|\mathbf{x}_i - \hat{\mathbf{x}}_i\|_2^2$ (cf. Lemma 1.2.3 in (Nesterov, 2003)).

Assumption 1.2. The functions $f_n(\cdot)$ and $h_n(\cdot)$ are Lipschitz continuous with some constants C_f and C_h , satisfying $\|\nabla f_n(\mathbf{x}_n)\| \leq C_f$ and $\|\partial h_n(\mathbf{x}_n)\| \leq C_h$ for all \mathbf{x}_n .

Assumption 1.3. We define $\bar{\lambda} \triangleq \lambda_{\max}(\mathbf{A}_n \mathbf{A}_n^\top)$, $\underline{\lambda} \triangleq \lambda_{\min}(\mathbf{A}_n \mathbf{A}_n^\top)$, $\lambda' = \lambda_{\min}(\mathbf{A}_n^\top \mathbf{A}_n)$. Either of these two conditions holds for matrix \mathbf{A}_n :

a) Condition $\mathbb{B}\mathbb{I}$: \mathbf{A}_n is bijective (i.e., $\underline{\lambda} = \lambda' > 0$), and it holds that $\kappa \triangleq \bar{\lambda}/\underline{\lambda} < 2$.

b) Condition $\mathbb{S}\mathbb{U}$: \mathbf{A}_n is surjective (i.e., $\underline{\lambda} > 0$, and λ' could be zero).

Assumption 1.4. Given any constant $\bar{\beta} \geq 0$, we let $\underline{\Theta}' \triangleq \inf_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n} \sum_{i=1}^n [f_i(\mathbf{x}_i) + h_i(\mathbf{x}_i)] + \frac{\bar{\beta}}{2} \|\sum_{i=1}^n \mathbf{A}_i \mathbf{x}_i - \mathbf{b}\|_2^2$. We assert that $\underline{\Theta}' > -\infty$.

Assumption 1.5. Let $i \in [n]$. The proximal operator $\text{Prox}_i(\mathbf{x}'_i; \mu) \triangleq \min_{\mathbf{x}_i} \frac{\mu}{2} \|\mathbf{x}_i - \mathbf{x}'_i\|_2^2 + h_i(\mathbf{x}_i)$ can be computed efficiently and exactly for any given $\mathbf{x}'_i \in \mathbb{R}^{d_i \times 1}$ and $\mu > 0$.

Assumption 1.6. If $\sum_{i=1}^n [f_i(\mathbf{x}_i) + h_i(\mathbf{x}_i)] < +\infty$, it follows that $\|\mathbf{x}_i\| < +\infty$ for all $i \in [n]$.

Assumption 1.7. Let $i \in [n]$. Assume the vector $\mathbf{x}'_i \in \mathbb{R}^{d_i \times 1}$ is bounded. Then, for any $\mu \in (0, \infty)$, the set $\text{Prox}_i(\mathbf{x}'_i; \mu)$ is also bounded.

Remarks. (i) Assumption 1.1 is commonly used in the convergence analysis of nonconvex algorithms. (ii) Assumption 1.2 imposes a continuity assumption only for the last block, allowing other blocks of the function $h_i(\mathbf{x}_i)_{i=1}^{n-1}$ to be nonsmooth and non-Lipschitz, such as indicator functions of constraint sets. It ensures bounded (sub-)gradients for $f_n(\cdot)$ and $h_n(\cdot)$, a relatively mild requirement that has found use in nonsmooth optimization (Li et al., 2022; 2021; Huang et al., 2019; Böhm & Wright, 2021). (iii) Assumption 1.3 demands a condition on the linear matrix \mathbf{A}_i for the last block ($i = n$), while leaving \mathbf{A}_i unrestricted for $i \in [n - 1]$. (iv) Assumption 1.4 ensures the well-defined nature of the penalty function associated with the problem, as has been used in (Gonçalves et al., 2017a). Furthermore, Assumption 1.4 can be satisfied if $\sum_{i=1}^n [f_i(\mathbf{x}_i) + h_i(\mathbf{x}_i)] > -\infty$. (v) Assumption 1.5 is frequently employed in nonconvex ADMM frameworks (Li & Pong, 2015; Bo̧t et al., 2019). Common examples of functions $h_i(\mathbf{x}_i)$ arising in practical applications include those discussed in (Gong et al., 2013), ℓ_0 regularization, $\ell_{1/2}$ regularization (Zeng et al., 2014), and indicator functions of cardinality constraints, matrices with orthogonality constraints (Lai & Osher,

162 2014), and matrices with rank constraints, among others. **(vi)** Assumptions 1.6 and 1.7 are used to
 163 guarantee the boundedness of the solution.

164
► Notations. We define $[n] \triangleq \{1, 2, \dots, n\}$ and $\mathbf{x} \triangleq \mathbf{x}_{[n]} \triangleq \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$. For any $j \geq i$, we
 165 denote $\mathbf{x}_{[i,j]} \triangleq \{\mathbf{x}_i, \mathbf{x}_{i+1}, \dots, \mathbf{x}_j\}$. We define $\lambda_{\min}(\mathbf{M})$ and $\lambda_{\max}(\mathbf{M})$ as the smallest and largest
 166 eigenvalue of the given matrix \mathbf{M} , respectively. We denote $\|\mathbf{A}_i\|$ as the spectral norm of the matrix
 167 \mathbf{A}_i . We denote $\mathbf{Ax} \triangleq \sum_{j=1}^n \mathbf{A}_j \mathbf{x}_j$, and $\|\mathbf{x}^+ - \mathbf{x}\|_2^2 = \sum_{i=1}^n \|\mathbf{x}_i^+ - \mathbf{x}_i\|_2^2$. Further notations and
 168 technical preliminaries are provided in Appendix A.
 169

170 2 MOTIVATING APPLICATIONS

171 Many machine learning and data science models can be formulated as Problem (1). Below, we
 172 present two examples, with additional applications provided in Appendix B.

173
► Sparse PCA. Sparse PCA (Chen et al., 2016; Lu & Zhang, 2012) Sparse PCA focuses on identifying a subset of informative variables with sparse loadings to enhance interpretability and
 174 reduce model complexity. It is formulated as: $\min_{\mathbf{V} \in \mathbb{R}^{d \times r}} \frac{1}{2m} \|\mathbf{D} - \mathbf{DV}\mathbf{V}^\top\|_F^2 + \dot{\rho} \|\mathbf{V}\|_1$, s.t. $\mathbf{V} \in$
 175 $\mathcal{M} \triangleq \{\mathbf{V} | \mathbf{V}^\top \mathbf{V} = \mathbf{I}\}$, where $\mathbf{D} \in \mathbb{R}^{m \times d}$ is the data matrix, and $\dot{\rho} \geq 0$. Introducing an additional variable \mathbf{Y} , this problem can be formulated as: $\min_{\mathbf{V}, \mathbf{Y}} \frac{1}{2m} \|\mathbf{D} - \mathbf{DV}\mathbf{V}^\top\|_F^2 + \dot{\rho} \|\mathbf{V}\|_1 +$
 176 $\iota_{\mathcal{M}}(\mathbf{Y})$, s. t. $-\mathbf{Y} + \mathbf{V} = \mathbf{0}$. It corresponds to Problem (1) with $\mathbf{x}_1 = \text{vec}(\mathbf{Y})$, $\mathbf{x}_2 = \text{vec}(\mathbf{V})$,
 177 $f_1(\mathbf{x}_1) = 0$, $h_1(\mathbf{x}_1) = \iota_{\mathcal{M}}(\mathbf{Y})$, $f_2(\mathbf{x}_2) = \frac{1}{2m} \|\mathbf{D} - \mathbf{DV}\mathbf{V}^\top\|_F^2$, $h_2(\mathbf{x}_2) = \dot{\rho} \|\mathbf{V}\|_1$, $\mathbf{A}_1 = -\mathbf{I}$,
 178 $\mathbf{A}_2 = \mathbf{I}$, $\mathbf{b} = \mathbf{0}$, and Condition \mathbb{B} .

179
► Structured Sparse Phase Retrieval. Sparse phase retrieval (Duchi & Ruan, 2018) aims to
 180 recover a sparse signal from the magnitudes of linear measurements. By incorporating additional linear constraints, recovery accuracy can be further improved. The problem is formulated as:
 181 $\min_{\mathbf{v}} \|(\mathbf{G}\mathbf{v}) \odot (\mathbf{G}\mathbf{v}) - \mathbf{z}\|_2^2 + \dot{\rho} \|\mathbf{v}\|_1$, s. t. $\mathbf{D}\mathbf{v} \geq \mathbf{0}$, where $\dot{\rho} \geq 0$, $\mathbf{G} \in \mathbb{R}^{m \times d}$, $\mathbf{z} \in \mathbb{R}^m$, $\mathbf{D} \in \mathbb{R}^{r \times d}$,
 182 with \mathbf{D} being surjective that $\mathbf{D}\mathbf{D}^\top \succ \mathbf{0}$. Introducing a new variable \mathbf{y} , this problem can be formulated as: $\min_{\mathbf{v}, \mathbf{y}} \|(\mathbf{G}\mathbf{v}) \odot (\mathbf{G}\mathbf{v}) - \mathbf{z}\|_2^2 + \dot{\rho} \|\mathbf{v}\|_1 + \iota_{\geq 0}(\mathbf{y})$, s.t. $\mathbf{y} - \mathbf{D}\mathbf{v} = \mathbf{0}$. This corresponds to
 183 Problem (1) with $\mathbf{x}_1 = \mathbf{y}$, $\mathbf{x}_2 = \mathbf{v}$, $f_1(\mathbf{x}_1) = 0$, $h_1(\mathbf{x}_1) = \iota_{\geq 0}(\mathbf{y})$, $f_2(\mathbf{x}_2) = \frac{1}{2} \|(\mathbf{G}\mathbf{v}) \odot (\mathbf{G}\mathbf{v}) - \mathbf{z}\|_2^2$,
 184 $h_2(\mathbf{x}_2) = \dot{\rho} \|\mathbf{v}\|_1$, $\mathbf{A}_1 = \mathbf{I}$, $\mathbf{A}_2 = -\mathbf{D}$, $\mathbf{b} = \mathbf{0}$, and Condition \mathbb{S} .

192 3 THE PROPOSED IPDS-ADMM ALGORITHM

193 This section describes the proposed IPDS-ADMM algorithm for solving Problem (1), featuring with
 194 using a new Increasing Penalization and Decreasing Smoothing (IPDS) strategy.

195 3.1 INCREASING PENALTY UPDATE STRATEGY

196 We employ an increasing penalty update strategy that is crucial to our algorithm. A natural choice
 197 for this penalty update rule is to use functions from the ℓ_p family. Throughout this paper, we consider
 198 the following penalty update rule $\{\beta^t\}_{t=0}^\infty$ for any given parameters $\xi, \delta, p \in (0, 1)$:

$$203 \quad \beta^t = \beta^0(1 + \xi t^p), \quad \beta^0 \geq L_n / (\delta \bar{\lambda}). \quad (2)$$

204 Here, L_n and $\bar{\lambda}$ are defined in Assumption 1.1 and Assumption 1.3, respectively.

205 We obtain the following useful lemma regarding the penalty update rule.

206
Lemma 3.1. (Proof in Appendix C.1) Given $\xi, \delta, p \in (0, 1)$, assume Formulation (2) is used to
 207 choose $\{\beta^t\}_{t=0}^\infty$. We have: **(a)** $\beta^t \leq \beta^{t+1} \leq (1 + \xi) \beta^t$, **(b)** $L_n \leq \delta \beta^t \bar{\lambda}$.

208
Remarks (i) The increasing penalty update strategy is closely coupled with the decreasing smoothing
 209 strategy and the diminishing stepsize approach in the literature. These strategies are frequently
 210 employed in subgradient methods (Li et al., 2021), smoothing gradient methods (Böhm & Wright,
 211 2021; Sun & Sun, 2023; Lei Yang, 2021), penalty decomposition methods (Lu & Zhang, 2013), and
 212 stochastic optimization algorithms like ADAM (Kingma & Ba, 2015; Chen et al., 2022), but are less
 213 commonly utilized in ADMM frameworks. We examine this approach within ADMM but limit our
 214 discussion to specific form and condition as in Formulation (2). **(ii)** The condition $\beta^0 \geq L_n / (\delta \bar{\lambda})$
 215

in Formulation (2) essentially mandates that the initial penalty value be sufficiently large. This condition can be automatically satisfied since an increasing penalty update is used. **(iii)** The result $\beta^{t+1} \leq (1 + \xi)\beta^t$ in Lemma 3.1 implies that the penalty parameter grows, but not excessively fast, with a constant ξ to prevent rapid escalation.

3.2 DECREASING MEREAU ENVELOPE SMOOTHING APPROACH

IPDS-ADMM is built upon the Moreau envelope smoothing technique (Li et al., 2022; Zeng et al., 2022; Sun & Sun, 2023; Böhm & Wright, 2021). Initially, we provide the following useful definition.

Definition 3.2. *The Moreau envelope of a proper convex and Lipschitz continuous function $h(\mathbf{u}) : \mathbb{R}^{d \times 1} \mapsto \mathbb{R}$ with parameter $\mu \in (0, \infty)$ is defined as $h(\mathbf{u}; \mu) \triangleq \min_{\mathbf{v} \in \mathbb{R}^{d \times 1}} h(\mathbf{v}) + \frac{1}{2\mu} \|\mathbf{v} - \mathbf{u}\|_2^2$.*

We offer some useful properties of Moreau envelop functions.

Lemma 3.3. *((Beck, 2017) Chapter 6) Suppose the function $h(\mathbf{u})$ is C_h -Lipschitz continuous and convex w.r.t. \mathbf{u} . We have: (a) The function $h(\mathbf{u}; \mu)$ is C_h -Lipschitz continuous w.r.t. \mathbf{u} . (b) The function $h(\mathbf{u}; \mu)$ is $(1/\mu)$ -smooth w.r.t. \mathbf{u} , and its gradient can be computed as: $\nabla h(\mathbf{u}; \mu) = \frac{1}{\mu}(\mathbf{u} - \mathbb{P}_h(\mathbf{u}; \mu))$, where $\mathbb{P}_h(\mathbf{u}; \mu) = \arg \min_{\mathbf{v}} h(\mathbf{v}) + \frac{1}{2\mu} \|\mathbf{v} - \mathbf{u}\|_2^2$. (c) $0 \leq h(\mathbf{u}) - h(\mathbf{u}; \mu) \leq \frac{1}{2} \mu C_h^2$.*

Lemma 3.4. *(Proof in Appendix C.2) Assuming $0 < \mu_2 < \mu_1$ and fixing $\mathbf{u} \in \mathbb{R}^{d \times 1}$, we have: $0 \leq \frac{h(\mathbf{u}; \mu_2) - h(\mathbf{u}; \mu_1)}{\mu_1 - \mu_2} \leq \frac{1}{2} C_h^2$.*

Lemma 3.5. *(Proof in Appendix C.3) Assuming $0 < \mu_2 < \mu_1$ and fixing $\mathbf{u} \in \mathbb{R}^{d \times 1}$, we have: $\|\nabla h(\mathbf{u}; \mu_1) - \nabla h(\mathbf{u}; \mu_2)\| \leq (\frac{\mu_1}{\mu_2} - 1) \cdot C_h$.*

Lemma 3.6. *(Proof in Appendix C.4) Given constants $\{\mathbf{c}, \mu, \rho\}$, we consider the convex problem in problem $\bar{\mathbf{x}}_n = \arg \min_{\mathbf{x}_n} h_n(\mathbf{x}_n; \mu) + \frac{\rho}{2} \|\mathbf{x}_n - \mathbf{c}\|_2^2$. We have: (a) $\bar{\mathbf{x}}_n = \frac{\mu}{1+\mu\rho} (\frac{1}{\mu} \check{\mathbf{x}}_n + \rho \mathbf{c})$, where $\check{\mathbf{x}}_n = \arg \min_{\mathbf{x}_n} h_n(\check{\mathbf{x}}_n) + \frac{1}{2} \cdot \frac{\rho}{1+\mu\rho} \|\check{\mathbf{x}}_n - \mathbf{c}\|_F^2 = \text{Prox}_n(\mathbf{c}; \mu + 1/\rho)$. (b) $\rho(\mathbf{c} - \bar{\mathbf{x}}_n) \in \partial h(\check{\mathbf{x}}_n)$. (c) $\|\mathbf{x}_n - \check{\mathbf{x}}_n\| \leq \mu C_h$.*

Remark 3.7. *(i) We highlight that Lemmas 3.4 and 3.5 are novel contributions of this paper and are instrumental for analyzing the proposed IPDS-ADMM algorithm. (ii) Lemma 3.6 is crucial for establishing the iteration complexity of Algorithm 1 to a critical point. The results of Lemma 3.6 are analogous to those of Lemma 1 in (Li et al., 2022).*

3.3 THE PROPOSED IPDS-ADMM ALGORITHM

This subsection provides the proposed IPDS-ADMM algorithm. Initially, we consider the following alternative optimization problem:

$$\min_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n} h_n(\mathbf{x}_n; \mu) + [\sum_{i=1}^{n-1} h_i(\mathbf{x}_i)] + [\sum_{i=1}^n f_i(\mathbf{x}_i)], \text{ s.t. } [\sum_{i=1}^n \mathbf{A}_i \mathbf{x}_i] = \mathbf{b}, \quad (3)$$

where $\mu \rightarrow 0$, and $h_n(\mathbf{x}_n; \mu) \triangleq \min_{\mathbf{v} \in \mathbb{R}^{d_n \times 1}} h(\mathbf{v}) + \frac{1}{2\mu} \|\mathbf{v} - \mathbf{x}_n\|_2^2$ is the Moreau envelope of $h_n(\mathbf{x}_n)$ with parameter μ . Lemma 3.3 confirms that $h_n(\mathbf{x}_n, \mu)$ is a $(1/\mu)$ -smooth function assuming $h_n(\cdot)$ is convex. We present the augmented Lagrangian function for Problem (3), as follows:

$$\mathcal{L}(\mathbf{x}, \mathbf{z}; \beta, \mu) \triangleq h_n(\mathbf{x}_n; \mu) + \{\sum_{i=1}^{n-1} h_i(\mathbf{x}_i)\} + G(\mathbf{x}, \mathbf{z}; \beta), \quad (4)$$

where $G(\mathbf{x}, \mathbf{z}; \beta)$ is differentiable and defined as:

$$G(\mathbf{x}, \mathbf{z}; \beta) \triangleq \sum_{i=1}^n f_i(\mathbf{x}_i) + \langle [\sum_{i=1}^n \mathbf{A}_i \mathbf{x}_i] - \mathbf{b}, \mathbf{z} \rangle + \frac{\beta}{2} \|[\sum_{i=1}^n \mathbf{A}_i \mathbf{x}_i] - \mathbf{b}\|_2^2.$$

Here, $\mu \in (0, \infty)$, $\beta \in (0, \infty)$, and $\mathbf{z} \in \mathbb{R}^{m \times 1}$ are the smoothing parameter, the penalty parameter, and the dual variable, respectively. We employ an increasing penalty and decreasing smoothing update scheme throughout all iterations $t = \{0, 1, \dots, \infty\}$ with $\beta^t \rightarrow +\infty$ and $\mu^t \propto \frac{1}{\beta^t} \rightarrow 0$. Notably, the function $G(\mathbf{x}^t, \mathbf{z}^t; \beta^t)$ is L_i^t -smooth w.r.t. \mathbf{x}_i for all $i \in [m]$, where $\mathsf{L}_i^t = L_i + \beta^t \|\mathbf{A}_i\|_2^2$. For notation simplicity, for all $i \in [n]$, we denote $\ddot{\mathbf{g}}_i^t \triangleq \nabla_{\mathbf{x}_i} G(\mathbf{x}_{[1, i-1]}^{t+1}, \mathbf{x}_i^t, \mathbf{x}_{[i+1, n]}^t, \mathbf{z}^t; \beta^t)$ as the gradient of $G(\mathbf{x}, \mathbf{z}^t; \beta^t)$ w.r.t. \mathbf{x}_i at the point \mathbf{x}_i^t .

In each iteration, we select suitable parameters $\{\beta^t, \mu^t\}$ and sequentially update the variables $(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n, \mathbf{z})$. We employ the proximal linearized method to cyclically update the variables $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$. Specifically, we update each variable \mathbf{x}_i by solving the following subproblem for all $i \in [n]$: $\mathbf{x}_i^{t+1} \approx \arg \min_{\mathbf{x}_i \in \mathbb{R}^{d_i \times 1}} \mathcal{L}(\mathbf{x}_{[1, i-1]}^t, \mathbf{x}_i^t, \mathbf{x}_{[i+1, n]}^t, \mathbf{z}^t; \beta^t, \mu^t)$. To address

the \mathbf{x}_i -subproblem, we employ a proximal linearized minimization strategy for all $i \in [n - 1]$:
 $\mathbf{x}_i^{t+1} \in \arg \min_{\mathbf{x}_i} h_i(\mathbf{x}_i) + \frac{\theta_1 L_i^t}{2} \|\mathbf{x}_i - \mathbf{x}_i^t\|_2^2 + \langle \mathbf{x}_i - \mathbf{x}_i^t, \ddot{\mathbf{g}}_i^t, \mathbf{z}^t; \beta^t \rangle$. However, for the final block of the problem, we consider a subtly different proximal linearized minimization strategy:
 $\mathbf{x}_n^{t+1} = \arg \min_{\mathbf{x}_n} h_n(\mathbf{x}_n; \mu^t) + \frac{\theta_2 L_n^t}{2} \|\mathbf{x}_n - \mathbf{x}_n^t\|_2^2 + \langle \mathbf{x}_n - \mathbf{x}_n^t, \ddot{\mathbf{g}}_n^t \rangle$. Importantly, we assign θ_1 to blocks $[1, n - 1]$ and θ_2 to block n . Our algorithm updates the dual variable \mathbf{z}^t using either an under-relaxation stepsize $\sigma \in (0, 1)$ or an over-relaxation stepsize $\sigma \in [1, 2)$.

Algorithm 1: IPDS-ADMM: The Proposed Proximal Linearized ADMM for Problem (1).

Choose suitable parameters $\{p, \xi, \delta\}$ and $\{\sigma, \theta_1, \theta_2\}$ using Formula (5) or Formula (6).

Initialize $\{\mathbf{x}^0, \mathbf{z}^0\}$. Choose $\beta^0 \geq L_n / (\delta \lambda)$.

for t from 0 to T **do**

S1) IPDS Strategy: Set $\beta^t = \beta^0(1 + \xi t^p)$, $\mu^t = 1 / (\bar{\lambda} \delta \beta^t)$.

We define $\ddot{\mathbf{g}}_i^t \triangleq \nabla_{\mathbf{x}_i} G(\mathbf{x}_{[1, i-1]}^{t+1}, \mathbf{x}_i^t, \mathbf{x}_{[i+1, n]}^t, \mathbf{z}^t; \beta^t)$.

S2) $\mathbf{x}_1^{t+1} \in \arg \min_{\mathbf{x}_1} h_1(\mathbf{x}_1) + \langle \mathbf{x}_1 - \mathbf{x}_1^t, \ddot{\mathbf{g}}_1^t \rangle + \frac{\theta_1 L_1^t}{2} \|\mathbf{x}_1 - \mathbf{x}_1^t\|_2^2$

S3) $\mathbf{x}_2^{t+1} \in \arg \min_{\mathbf{x}_2} h_2(\mathbf{x}_2) + \langle \mathbf{x}_2 - \mathbf{x}_2^t, \ddot{\mathbf{g}}_2^t \rangle + \frac{\theta_1 L_2^t}{2} \|\mathbf{x}_2 - \mathbf{x}_2^t\|_2^2$

...

S4) $\mathbf{x}_{n-1}^{t+1} \in \arg \min_{\mathbf{x}_{n-1}} h_{n-1}(\mathbf{x}_{n-1}) + \langle \mathbf{x}_{n-1} - \mathbf{x}_{n-1}^t, \ddot{\mathbf{g}}_{n-1}^t \rangle + \frac{\theta_1 L_{n-1}^t}{2} \|\mathbf{x}_{n-1} - \mathbf{x}_{n-1}^t\|_2^2$

S5) $\mathbf{x}_n^{t+1} \in \arg \min_{\mathbf{x}_n} h_n(\mathbf{x}_n; \mu) + \langle \mathbf{x}_n - \mathbf{x}_n^t, \ddot{\mathbf{g}}_n^t \rangle + \frac{\theta_2 L_n^t}{2} \|\mathbf{x}_n - \mathbf{x}_n^t\|_2^2$. It can be solved using Lemma 3.6 as $\mathbf{x}_n^{t+1} = \frac{1}{1+\mu\rho}(\check{\mathbf{x}}_n^{t+1} + \mu\rho\mathbf{c})$, where $\check{\mathbf{x}}_n^{t+1} = \text{Prox}_n(\mathbf{c}; \mu + 1/\rho)$, $\mu = \mu^t$,

$\rho \triangleq \theta_2 L_n^t$, and $\mathbf{c} \triangleq \mathbf{x}_n^t - \ddot{\mathbf{g}}_n^t / \rho$.

S6) $\mathbf{z}^{t+1} = \mathbf{z}^t + \sigma \beta^t ([\sum_{j=1}^n \mathbf{A}_j \mathbf{x}_j^{t+1}] - \mathbf{b})$

end

We present IPDS-ADMM in Algorithm 1, and have the following remarks.

Remark 3.8. (i) Algorithm 1 can be viewed as a generalized cyclic coordinate descent method applied to the augmented Lagrangian function in Equation (4). (ii) The Moreau envelope smoothing technique has been used in the design of augmented Lagrangian methods (Zeng et al., 2022) and ADMMs (Li et al., 2022), and minimax optimization (Zhang et al., 2020). However, these algorithms typically utilize constant penalties, whereas we adopt an Increasing Penalization and Decreasing Smoothing (IPDS) strategy to improve the iteration complexity of RADMM (Li et al., 2022), reducing it from $\mathcal{O}(1/\epsilon^4)$ to $\mathcal{O}(1/\epsilon^3)$. (iii) Algorithm 1 is a fully splitting algorithm, where each step reduces to computing a proximal operator. For the first $(n - 1)$ blocks, we have: $\mathbf{x}_i^{t+1} \in \text{Prox}_i(\mathbf{x}_i^t - \ddot{\mathbf{g}}_i^t / \dot{\rho}; 1 / \dot{\rho})$, where $\dot{\rho} = \theta_1 L_i^t$. For the last block, Lemma 3.6 can be applied to compute the proximal operator of the smoothed function $h_n(\mathbf{x}_n; \mu)$ using the proximal operator of the original function $h_n(\mathbf{x}_n)$. (iv) The point $\check{\mathbf{x}}_n^{t+1}$ in Step S5 of Algorihtm 1 plays a crucial role. As will be seen later in Theorem 4.18, the point $(\mathbf{x}_1^t, \mathbf{x}_2^t, \dots, \mathbf{x}_{n-1}^t, \check{\mathbf{x}}_n^t, \mathbf{z}^t)$, rather than the point $(\mathbf{x}_1^t, \mathbf{x}_2^t, \dots, \mathbf{x}_{n-1}^t, \mathbf{x}_n^t, \mathbf{z}^t)$, will serve as an approximate critical point of Problem (1) in our complexity results. (v) RADMM (Li et al., 2022) uses a fixed large penalty parameter $\mathcal{O}(1/\epsilon)$ and a fixed small smoothing parameter $\mathcal{O}(\epsilon)$ to achieve an ϵ -approximate critical point. However, this leads to overly conservative step sizes for the primal and dual updates, potentially hindering the algorithm's practical performance. (vi) We apply the smoothing strategy only to the last block to bound the dual variables via the primal ones. This leverages the Lipschitz continuity of the smoothed function to estimate $\frac{1}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2$ and construct a suitable potential function. (vii) Some may worry that using an increasing penalty could cause the parameter to become unbounded. However, by setting $\xi \ll 1$, we ensure $\beta^t \leq \beta^{t+1} \leq (1 + \xi)\beta^t$, meaning the penalty grows very slowly in practice.

3.4 CHOOSING SUITABLE PARAMETERS $\{p, \xi, \delta\}$ AND $\{\sigma, \theta_1, \theta_2\}$

Selecting appropriate parameters $\{p, \xi, \delta\}$ and $\{\sigma, \theta_1, \theta_2\}$ is essential to ensuring the global convergence of Algorithm 1. In our theoretical analysis and empirical experiments, we suggest the

following choices for $\{p, \xi, \delta\}$ and $\{\sigma, \theta_1, \theta_2\}$:

$$\mathbb{B}\mathbb{I} : p = \frac{1}{3}, \xi \in (0, \infty), \delta \in (0, \frac{1}{3}(\frac{2}{\kappa} - 1)), \sigma \in [1, 2], \theta_1 = 1.01, \theta_2 = \frac{1/\kappa-\delta}{1+\delta} + \frac{1}{2\chi_0(1+\delta)^2}. \quad (5)$$

$$\mathbb{S}\mathbb{U} : p = \frac{1}{3}, \xi = \delta = \sigma = \frac{0.01}{\kappa}, \theta_1 = 1.01, \theta_2 = 1.5. \quad (6)$$

Here, $\chi_0 \triangleq 6\omega\sigma_1\kappa$, and $\omega \triangleq 1 + \frac{\xi}{2\sigma} + \sigma\xi$. Notably, the parameter θ_2 in (5) depends on (ξ, δ, σ) .

Remark 3.9. (i) From (5), we find that $\frac{1/\kappa-\delta}{1+\delta} \geq \{1/\kappa - \frac{2}{3\kappa} + \frac{1}{3}\}/\{1 + \frac{2}{3\kappa} - \frac{1}{3}\} = 1/2$, leading to $\theta_2 > 1/2$. (ii) From (6), we observe that the parameters $\{\xi, \delta, \sigma\}$ is inversely proportional to the condition number κ . Such settings are partly consistent with those in (Bōt et al., 2019) (refer to Lemma 5 in (Bōt et al., 2019)). (iii) Introducing the relaxation parameter $\sigma \in (0, 2)$ enables handling cases where the matrix is surjective. Specifically, when the matrix is bijective, we can use an over-relaxation step size for faster convergence, whereas for surjective matrices, the algorithm requires conservative step sizes to ensure global convergence.

4 GLOBAL CONVERGENCE

This section establishes the global convergence of Algorithm 1.

We begin with a high-level overview of the proof strategy. First, using the Lagrangian function, we derive sufficient decrease conditions for the four parameter sets: primal variables, dual variables, the penalty parameter, and the smoothing parameter. Next, using the first-order optimality conditions and dual update rules, we bound the difference in dual variables using primal by the difference in primal variables. Lastly, we show that the tail error term related to the smoothing parameter is constant, establishing the summability of the sequence linked to a potential function.

We provide the following three useful lemmas.

Lemma 4.1. (Proof in Appendix D.1, A Sufficient Decrease Property) Fix $\varepsilon_3 \triangleq \xi$ and $\varepsilon_1 \triangleq \frac{1}{2}\theta_1 - \frac{1}{2}$. Let $\varepsilon_2 \in \mathbb{R}$. For all $t \geq 1$, we have:

$$\mathcal{E}^{t+1} + \Theta_L^{t+1} - \Theta_L^t \leq (\frac{1}{2} - \theta_2 + \varepsilon_2) \cdot \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2, \quad (7)$$

where $\mathcal{E}^{t+1} \triangleq [\varepsilon_1 \sum_{i=1}^{n-1} \mathsf{L}_i^t \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2] + \varepsilon_2 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \frac{\varepsilon_3}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2$.

Furthermore, $\Theta_L^t \triangleq \mathcal{L}(\mathbf{x}^t, \mathbf{z}^t; \beta^t, \mu^t) + \frac{1}{2}C_h\mu^t$, $\mathsf{L}_i^t = L_i + \beta^t \|\mathbf{A}_i\|_2^2$, and $\omega \triangleq 1 + \frac{\xi}{2\sigma} + \sigma\xi$.

Lemma 4.2. (Proof in Appendix D.2, First-Order Optimality Condition) Assume $\sigma \in (0, 2)$. For all $t \geq 1$ and $i \in [n-1]$, we have the following results.

(a) Let $\mathbf{w}_i^{t+1} \in \partial h_i(\mathbf{x}_i^{t+1}) + \nabla f_i(\mathbf{x}_i^t)$, and $\mathbf{u}_i^{t+1} \triangleq \theta_1 \mathsf{L}_i^t (\mathbf{x}_i^{t+1} - \mathbf{x}_i^t) - \beta^t \mathbf{A}_i^\top [\sum_{j=i}^n \mathbf{A}_j (\mathbf{x}_j^{t+1} - \mathbf{x}_j^t)]$. It holds that: $\mathbf{0} = \sigma \mathbf{A}_i^\top \mathbf{z}^t + \mathbf{A}_i^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) + \sigma \mathbf{w}_i^{t+1} + \sigma \mathbf{u}_i^{t+1}$.

(b) Let $\mathbf{w}_n^{t+1} \triangleq \nabla h_n(\mathbf{x}_n^{t+1}, \mu^t) + \nabla f_n(\mathbf{x}_n^t)$, and $\mathbf{u}_n^{t+1} \triangleq \mathbf{Q}^t (\mathbf{x}_n^{t+1} - \mathbf{x}_n^t)$, where $\mathbf{Q}^t \triangleq \theta_2 \mathsf{L}_n^t \mathbf{I} - \beta^t \mathbf{A}_n^\top \mathbf{A}_n$. It holds that: $\mathbf{0} = \sigma \mathbf{A}_n^\top \mathbf{z}^t + \mathbf{A}_n^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) + \sigma \mathbf{w}_n^{t+1} + \sigma \mathbf{u}_n^{t+1}$.

(c) We have the following two different identities:

$$\mathbb{B}\mathbb{I} : \begin{cases} \mathbf{a}^{t+1} = (1 - \sigma)\mathbf{a}^t + \sigma\mathbf{c}^t, \\ \text{where } \mathbf{a}^{t+1} \triangleq \mathbf{A}_n^\top (\mathbf{z}^{t+1} - \mathbf{z}^t), \text{ and } \mathbf{c}^t \triangleq \mathbf{u}_n^t - \mathbf{u}_n^{t+1} + \mathbf{w}_n^t - \mathbf{w}_n^{t+1}. \end{cases} \quad (8)$$

$$\mathbb{S}\mathbb{U} : \begin{cases} \mathbf{a}^{t+1} = (1 - \sigma)\mathbf{a}^t + \sigma\mathbf{c}^t, \\ \text{where } \mathbf{a}^{t+1} \triangleq \mathbf{A}_n^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) + \sigma \mathbf{u}_n^{t+1}, \text{ and } \mathbf{c}^t \triangleq \sigma \mathbf{u}_n^t + \mathbf{w}_n^t - \mathbf{w}_n^{t+1}. \end{cases} \quad (9)$$

Lemma 4.3. (Proof in Appendix D.3) For all $t \geq 0$, we have: (a) $\mathsf{L}_n^t \leq \beta^t \bar{\lambda}(1 + \delta)$; (b) $\|\mathbf{Q}^t\| \leq \beta^t \bar{\lambda} q$, where $q \triangleq \theta_2(1 + \delta) - \underline{\lambda}'/\bar{\lambda}$; (c) $\|\mathbf{u}_n^{t+1}\| \leq q\bar{\lambda}\beta^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|$.

We provide convergence analysis of Algorithm 1 under two conditions: Condition $\mathbb{B}\mathbb{I}$ using Formulation (8), and Condition $\mathbb{S}\mathbb{U}$ using Formulation (9).

We first define the following parameters for different Conditions $\mathbb{B}\mathbb{I}$ and $\mathbb{S}\mathbb{U}$:

$$\mathbb{B}\mathbb{I} : \left\{ K_a \triangleq \frac{\omega\sigma_2}{\underline{\lambda}}, K_u \triangleq \frac{3\omega\sigma_1}{\underline{\lambda}}, \Theta_a^t \triangleq \frac{K_a}{\beta^t} \|\mathbf{a}^t\|_2^2, \Theta_u^t = \frac{K_u}{\beta^t} (L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \|\mathbf{u}_n^t\|)^2 \right\}. \quad (10)$$

$$\mathbb{S}\mathbb{U} : \left\{ K_a \triangleq \frac{2\omega\sigma_2}{\underline{\lambda}}, K_u \triangleq \frac{6\omega\sigma_1}{\underline{\lambda}}, \Theta_a^t \triangleq \frac{K_a}{\beta^t} \|\mathbf{a}^t\|_2^2, \Theta_u^t = \frac{K_u}{\beta^t} (L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \sigma \|\mathbf{u}_n^t\|)^2 \right\}. \quad (11)$$

378 Here, $\sigma \in (0, 2)$, and $\{\sigma_1, \sigma_2\}$ are defined as: $\sigma_1 \triangleq \frac{\sigma}{(1-|\sigma|)^2}$, $\sigma_2 \triangleq \frac{|1-\sigma|}{\sigma(1-|1-\sigma|)}$. Using the
 379 parameters $\{K_a, K_u\}$, we construct a sequence associated with the potential (or Lyapunov) function
 380 as follows: $\Theta^t = \Theta_L^t + \Theta_a^t + \Theta_u^t$.
 381

382 4.1 ANALYSIS FOR CONDITION $\mathbb{B}\mathbb{I}$

384 We provide a convergence analysis of Algorithm 1 under Condition $\mathbb{B}\mathbb{I}$, where \mathbf{A}_n is a bijective
 385 matrix. We assume an over-relaxation stepsize is used with $\sigma \in [1, 2)$.

386 The subsequent lemma uses Equation (8) to establish an upper bound for the term $\frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2$.
 387

388 **Lemma 4.4.** (*Proof in Appendix D.4, Bounding Dual Using Primal*) *We define ω as in Lemma 4.1.
 389 For all $t \geq 1$, we have:*

$$390 \quad \frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \leq \Theta_{au}^t - \Theta_{au}^{t+1} + \chi_1 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \Gamma_\mu^t, \quad (12)$$

392 where $\chi_1 \triangleq \chi_0(\delta + \theta_2 + \theta_2\delta - 1/\kappa)^2$, $\chi_0 \triangleq 6\omega\sigma_1\kappa$, $\Theta_{au}^t \triangleq \Theta_a^t + \Theta_u^t$, and $\{K_a, K_u\}$ are defined
 393 in Equation (10), and $\Gamma_\mu^t \triangleq C_h^2 \frac{K_u}{\beta^t} \cdot (\frac{\mu^{t-1}}{\mu^t} - 1)^2$.
 394

395 Assume Equation (5) is used to choose $\{p, \xi, \delta, \sigma, \theta_1, \theta_2\}$. We have the following two lemmas.

396 **Lemma 4.5.** (*Proof in Appendix D.5*) *We have: $\varepsilon_1 \triangleq \frac{1}{2}\theta_1 - \frac{1}{2} > 0$, and $\varepsilon_2 \triangleq \theta_2 - \frac{1}{2} - \chi_1 \geq \frac{1}{8\chi_0} > 0$.
 397 Here, $\{\chi_1, \chi_0\}$ are defined in Lemma 4.4.*

398 **Lemma 4.6.** (*Proof in Appendix D.6, Decrease on a Potential Function*) *For all $t \geq 1$, we have
 399 $\mathcal{E}^{t+1} \leq \Theta^t - \Theta^{t+1} + \Gamma_\mu^t$*

402 4.2 ANALYSIS FOR CONDITION $\mathbb{S}\mathbb{U}$

403 We provide a convergence analysis of Algorithm 1 under Condition $\mathbb{S}\mathbb{U}$, where \mathbf{A}_n is a surjective
 404 matrix. We assume an under-relaxation stepsize is used with $\sigma \in (0, 1)$.

405 The following lemma utilizes Equation (9) to establish an upper bound for the term $\frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2$.
 406

407 **Lemma 4.7.** (*Proof in Appendix D.7, Bounding Dual Using Primal*) *We define ω as in Lemma 4.1.
 408 For all $t \geq 1$, we have:*

$$409 \quad \frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \leq \Theta_{au}^t - \Theta_{au}^{t+1} + \chi_2 \cdot \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \Gamma_\mu^t, \quad (13)$$

410 where $\chi_2 \triangleq \frac{2\omega\kappa}{\sigma} \cdot \{\sigma^2 q^2 + 3\delta^2 + 3(\delta + \sigma q)^2\}$, $q \triangleq \theta_2 + \theta_2\delta$, $\Theta_{au}^t \triangleq \Theta_a^t + \Theta_u^t$, and $\{K_a, K_u\}$ are
 411 defined in Equation (11), and $\Gamma_\mu^t \triangleq C_h^2 \frac{K_u}{\beta^t} \cdot (\frac{\mu^{t-1}}{\mu^t} - 1)^2$.
 412

413 Assume Equation (6) is used to choose $\{p, \xi, \delta, \sigma, \theta_1, \theta_2\}$. We have the following two lemmas.

414 **Lemma 4.8.** (*Proof in Appendix D.8*) *We have: $\varepsilon_1 \triangleq \frac{1}{2}\theta_1 - \frac{1}{2} > 0$, and $\varepsilon_2 \triangleq \theta_2 - \frac{1}{2} - \chi_2 \geq 0.02 > 0$.*

415 **Lemma 4.9.** (*Proof in Appendix D.9, Decrease on a Potential Function*). *For all $t \geq 1$, we have:
 416 $\mathcal{E}^{t+1} \leq \Theta^t - \Theta^{t+1} + \Gamma_\mu^t$*

421 4.3 CONTINUING ANALYSIS FOR CONDITIONS $\mathbb{B}\mathbb{I}$ AND $\mathbb{S}\mathbb{U}$

422 The following lemma demonstrates that Θ^t is consistently lower bounded.

423 **Lemma 4.10.** (*Proof in Appendix D.10*) *For all $t \geq 1$, there exists a constant $\underline{\Theta}$ suc that $\Theta^t \geq \underline{\Theta}$.*

424 The following lemma shows that $\sum_{t=1}^{\infty} \Gamma_\mu^t$ is always upper bounded.

425 **Lemma 4.11.** (*Proof in Appendix D.11*) *We define Γ_μ^t as in Lemma 4.4 and Lemma 4.7. There
 426 exists a universal positive constant C_μ such that $\sum_{t=1}^{\infty} \Gamma_\mu^t \leq C_\mu$.*
 427

428 We present the following theorem concerning a summable property of the sequence $\{\mathcal{E}^{t+1}\}_{t=1}^{\infty}$.
 429

430 **Theorem 4.12.** (*Proof in Appendix D.12*) *Letting $K_e \triangleq \Theta^1 - \underline{\Theta} + C_\mu$, we have: $\sum_{t=1}^{\infty} \mathcal{E}^{t+1} \leq K_e$.*

432 The following lemmas are useful to provide upper bounds for the dual and primal variables.
 433

434 **Lemma 4.13.** (*Proof in Appendix D.13*) There exist constants $\{K_z, \tilde{K}_z\}$ such that $\forall t \geq 1, \frac{1}{\beta^t} \|\mathbf{z}^t\|_2^2 \leq K_z$, and $\sum_{t=1}^{\infty} \frac{1}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \leq \tilde{K}_z$.
 435

436 **Lemma 4.14.** (*Proof in Appendix D.14*) We have $\|\mathbf{x}_i^{t+1}\| < +\infty$ for all $i \in [n]$.
 437

438 Finally, we have the following theorem regrading to the global convergence of IPDS-ADMM.
 439

440 **Theorem 4.15.** (*Proof in Appendix D.15*) We define $K_c \triangleq K_e / \min\{\epsilon_3, \min(\epsilon_1, \epsilon_2) \mathbf{A}\}$, where $\mathbf{A} \triangleq 441 \min_{i=1}^n \|\mathbf{A}_i\|_2^2$. We have the following results: (a) $\sum_{t=1}^T \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + \|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2 \leq K_c \beta^T$.
 442 (b) There exists an index \bar{t} with $\bar{t} \leq T$ such that $\|\mathbf{z}^{\bar{t}+1} - \mathbf{z}^{\bar{t}}\|_2^2 + \|\beta^{\bar{t}}(\mathbf{x}^{\bar{t}+1} - \mathbf{x}^{\bar{t}})\|_2^2 \leq \frac{K_c \beta^T}{T}$.
 443

444 **Remark 4.16.** (i) With the choice $\beta^T = \mathcal{O}(T^p)$ with $p \in (0, 1)$, we observe $\ddot{e}^{\bar{t}} \triangleq \|\mathbf{z}^{\bar{t}+1} - \mathbf{z}^{\bar{t}}\|_2^2 + 445 \|\beta^{\bar{t}}(\mathbf{x}^{\bar{t}+1} - \mathbf{x}^{\bar{t}})\|_2^2 = \mathcal{O}(T^{p-1})$, indicating convergence of $\ddot{e}^{\bar{t}}$ towards 0. (ii) In light of Theorem
 446 4.15, a reasonable stopping criterion for Algorithm 1 is $\|\mathbf{z}^{\bar{t}+1} - \mathbf{z}^{\bar{t}}\| + \|\beta^{\bar{t}}(\mathbf{x}^{\bar{t}+1} - \mathbf{x}^{\bar{t}})\| \leq \epsilon$, where
 447 $\epsilon \geq 0$ is a user-defined parameter.
 448

449 4.4 ITERATION COMPLEXITY

450 We now establish the iteration complexity of Algorithm 1. We first restate the following standard
 451 definition of approximated critical points.
 452

453 **Definition 4.17.** (ϵ -Critical Point) A solution $(\check{\mathbf{x}}, \check{\mathbf{z}})$ is an ϵ -critical point if it holds that: $\text{Crit}(\check{\mathbf{x}}, \check{\mathbf{z}}) \leq 454 \epsilon^2$, where $\text{Crit}(\check{\mathbf{x}}, \check{\mathbf{z}}) \triangleq \|\mathbf{A}\check{\mathbf{x}} - \mathbf{b}\|_2^2 + \sum_{i=1}^n \text{dist}^2(\mathbf{0}, \nabla f_i(\check{\mathbf{x}}_i) + \partial h_i(\check{\mathbf{x}}_i) + \mathbf{A}_i^\top \check{\mathbf{z}})$, and $\text{dist}^2(\Omega, \Omega') \triangleq 455 \inf_{\mathbf{w} \in \Omega, \mathbf{w}' \in \Omega'} \|\mathbf{w} - \mathbf{w}'\|_2^2$ is the squared distance between two sets.
 456

457 We obtain the following iteration complexity results.
 458

459 **Theorem 4.18.** (*Proof in Appendix D.16*) We define $\mathbf{q}^t \triangleq \{\mathbf{x}_1^t, \mathbf{x}_2^t, \dots, \mathbf{x}_{n-1}^t, \check{\mathbf{x}}_n^t\}$. Let the sequence
 460 $\{\mathbf{q}^t, \mathbf{z}^t\}_{t=0}^T$ be generated by Algorithm 1. If $p \in (0, \frac{1}{2})$, we have: $\frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{q}^{t+1}, \mathbf{z}^{t+1}) \leq 461 \mathcal{O}(T^{p-1}) + \mathcal{O}(T^{-1}) + \mathcal{O}(T^{-2p})$. In particular, with the choice $p = 1/3$, we have
 462 $\frac{1}{T} \sum_{t=1}^T \text{Crit}(\mathbf{q}^{t+1}, \mathbf{z}^{t+1}) \leq \mathcal{O}(T^{-2/3})$. In other words, there exists $\bar{t} \leq T$ such that:
 463 $\text{Crit}(\mathbf{q}^{\bar{t}+1}, \mathbf{z}^{\bar{t}+1}) \leq \epsilon^2$, provided that $T \geq \mathcal{O}(1/\epsilon^3)$.
 464

465 **Remark 4.19.** To the best of our knowledge, this represents the first complexity result for using
 466 ADMM to solve this class of nonsmooth and nonconvex problems. Remarkably, we observe that it
 467 aligns with the iteration bound found in smoothing proximal gradient methods (Böhm & Wright,
 2021).
 468

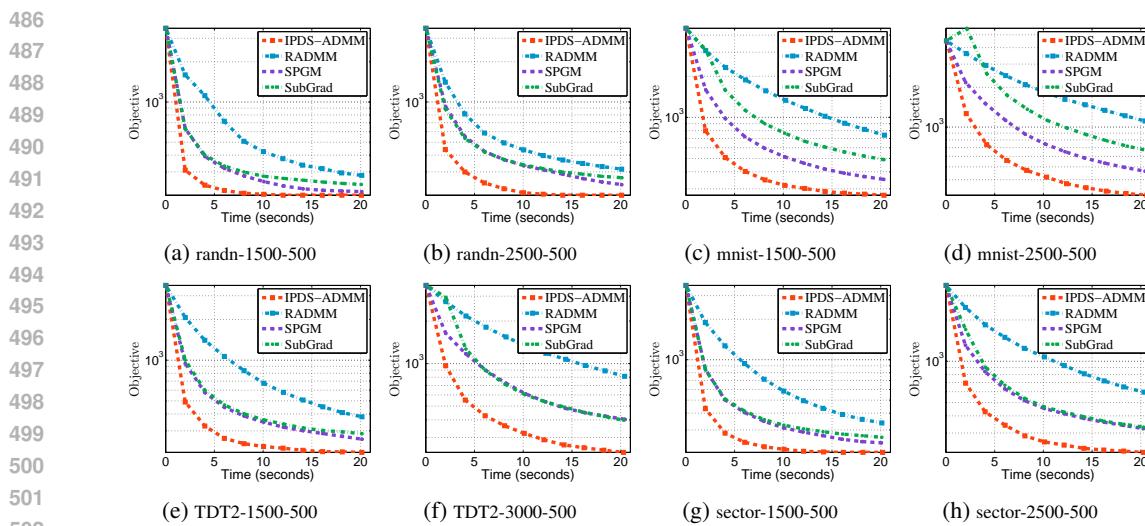
469 4.5 ON THE BOUNDEDNESS AND CONVERGENCE OF THE MULTIPLIERS

470 Questions may arise regarding whether the multipliers \mathbf{z}^t in Algorithm 1 are bounded, given that
 471 $\|\mathbf{z}^t\|_2^2 \leq K_z \beta^t$, as stated in Lemma 4.13. We argue that the boundedness of the multipliers is not
 472 an issue. We propose the following variable substitution: $\frac{\mathbf{z}^t}{\sqrt{\beta^t}} \triangleq \hat{\mathbf{z}}^t$ for all t . Consequently, we
 473 can implement the following update rule to replace the dual variable update rule of Algorithm 1:
 474 $\hat{\mathbf{z}}^{t+1} = \hat{\mathbf{z}}^t \frac{\sqrt{\beta^t}}{\sqrt{\beta^{t+1}}} + \frac{\beta^t}{\sqrt{\beta^{t+1}}} \cdot \sigma(\mathbf{A}\mathbf{x}^{t+1} - \mathbf{b})$. Additionally, \mathbf{z}^t should be replaced with $\sqrt{\beta^t} \cdot \hat{\mathbf{z}}^t$ in
 475 the remaining steps of Algorithm 1. Importantly, such a substitution does not essentially alter the
 476 algorithm or our analysis throughout this paper.
 477

478 We have the following results for the new multipliers $\hat{\mathbf{z}}^t$:
 479

480 **Lemma 4.20.** (*Proof in Appendix D.17*) We have: (a) $\forall t \geq 0$, $\|\hat{\mathbf{z}}^t\|_2^2 \leq K_z$; (b) $\sum_{t=1}^{\infty} \|\hat{\mathbf{z}}^{t+1} - 481 \hat{\mathbf{z}}^t\|_2^2 \leq 2\tilde{K}_z + K_z$. Here, $\{\tilde{K}_z, K_z\}$ are bounded constants defined in Lemma 4.13.
 482

483 **Remark 4.21.** Thanks to the variable substitution, the new multiplier $\|\hat{\mathbf{z}}^t\|$ is bounded and convergent with $(\min_{t=1}^T \|\hat{\mathbf{z}}^{t+1} - \hat{\mathbf{z}}^t\|_2^2) \leq \frac{1}{T} \sum_{t=1}^T \|\hat{\mathbf{z}}^{t+1} - \hat{\mathbf{z}}^t\|_2^2 \leq \mathcal{O}(1/T)$.
 484

Figure 1: Convergence curves of methods for sparse PCA with $\dot{\rho} = 10$ and $\beta^0 = 50\dot{\rho}$.

5 EXPERIMENTS

This section assesses the performance of IPDS-ADMM in solving the sparse PCA problem, as shown in Section 2.

► **Compared Methods.** We compare IPDS-ADMM against three state-of-the-art general-purpose algorithms that solve Problem (1) (i) the Subgradient method (SubGrad) (Li et al., 2021; Davis & Drusvyatskiy, 2019), (ii) the Smoothing Proximal Gradient Method (SPGM) (Böhm & Wright, 2021), (iii) the Riemannian ADMM with fixed and large penalty (RADMM) (Li et al., 2022).

► **Experimental Settings.** All methods are implemented in MATLAB on an Intel 2.6 GHz CPU with 64 GB RAM. We incorporate a set of 8 datasets into our experiments, comprising both randomly generated and publicly available real-world data. Appendix Section E describes how to generate the data used in the experiments. For IPDS-ADMM, we set $(\beta^0, p, \xi, \delta, \sigma, \theta) = (50\dot{\rho}, 1/3, 0.5, 1/4, 1.618, 1.01)$. The penalty parameter for RADMM is set to a reasonably large constant $\beta = 100\dot{\rho}$. We fix $\dot{r} = 20$ and compare objective values for all methods after running T' seconds, where T' is reasonably large to ensure the proposed method converges. We provide our code in the supplemental material.

► **Experiment Results.** The experimental results depicted in Figure 1 offer the following insights: (i) Sub-Grad tends to be less efficient in comparison to other methods. (ii) SPGM, utilizing a variable smoothing strategy, generally demonstrates slower performance than the multiplier-based variable splitting method. This observation corroborates the widely accepted notion that primal-dual methods are typically more robust and quicker than primal-only methods. (iii) The proposed IPDS-ADMM generally attains the lowest objective function values among all methods examined.

6 CONCLUSIONS

In this paper, we introduce IPDS-ADMM, a proximal linearized ADMM that uses an Increasing Penalization and Decreasing Smoothing (IPDS) strategy for solving general multi-block nonconvex composite optimization problems. IPDS-ADMM operates under a relatively relaxed condition, requiring continuity in just one block of the objective function. It incorporates relaxed strategies for dual variable updates when the associated linear operator is either bijective or surjective. We increase the penalty parameter and decrease the smoothing parameter at a controlled pace, and introduce a Lyapunov function for convergence analysis. We also derive the iteration complexity of IPDS-ADMM. Finally, we conduct experiments to demonstrate the effectiveness of our approaches.

540 REFERENCES
541

- 542 Rina Foygel Barber and Emil Y Sidky. Convergence for nonconvex admm, with applications to ct
543 imaging. *Journal of Machine Learning Research*, 25(38):1–46, 2024.
- 544 Amir Beck. *First-order methods in optimization*. SIAM, 2017.
545
- 546 Dimitri Bertsekas. *Convex optimization algorithms*. Athena Scientific, 2015.
547
- 548 Fengmiao Bian, Jingwei Liang, and Xiaoqun Zhang. A stochastic alternating direction method of
549 multipliers for non-smooth and non-convex optimization. *Inverse Problems*, 37(7):075009, 2021.
- 550 Radu Ioan Boț and Dang-Khoa Nguyen. The proximal alternating direction method of multipliers
551 in the nonconvex setting: convergence analysis and rates. *Mathematics of Operations Research*,
552 45(2):682–712, 2020.
- 553 Radu Ioan Boț, Erno Robert Csetnek, and Dang-Khoa Nguyen. A proximal minimization algorithm
554 for structured nonconvex and nonsmooth problems. *SIAM Journal on Optimization*, 29(2):1300–
555 1328, 2019. doi: 10.1137/18M1190689.
556
- 557 Axel Böhm and Stephen J. Wright. Variable smoothing for weakly convex composite functions.
558 *Journal of Optimization Theory and Applications*, 188(3):628–649, 2021.
559
- 560 Radu Ioan Boț, Minh N Dao, and Guoyin Li. Inertial proximal block coordinate method for a class of
561 nonsmooth sum-of-ratios optimization problems. *SIAM Journal on Optimization*, 33(2):361–393,
562 2023.
- 563 E. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis? *Journal of the ACM*,
564 58(3), May 2011.
565
- 566 Congliang Chen, Li Shen, Fangyu Zou, and Wei Liu. Towards practical adam: Non-convexity,
567 convergence theory, and mini-batch acceleration. *Journal of Machine Learning Research*, 23
568 (229):1–47, 2022. URL <http://jmlr.org/papers/v23/20-1438.html>.
- 569 Weiqiang Chen, Hui Ji, and Yanfei You. An augmented lagrangian method for ℓ_1 -regularized opti-
570 mization problems with orthogonality constraints. *SIAM Journal on Scientific Computing*, 38(4):
571 B570–B592, 2016.
572
- 573 Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex
574 functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019.
575
- 576 Wei Deng, Ming-Jun Lai, Zhimin Peng, and Wotao Yin. Parallel multi-block admm with o (1/k)
577 convergence. *Journal of Scientific Computing*, 71:712–736, 2017.
578
- 579 John C Duchi and Feng Ruan. Solving (most) of a set of quadratic equalities: composite optimization
580 for robust phase retrieval. *Information and Inference: A Journal of the IMA*, 8(3):471–529, 2018.
581
- 582 Daniel Gabay and Bertrand Mercier. A dual algorithm for the solution of nonlinear variational
583 problems via finite element approximation. *Computers & mathematics with applications*, 2(1):
17–40, 1976.
584
- 585 Max LN Gonçalves, Jefferson G Melo, and Renato DC Monteiro. Convergence rate bounds for
586 a proximal admm with over-relaxation stepsize parameter for solving nonconvex linearly con-
587 strained problems. *arXiv preprint arXiv:1702.01850*, 2017a.
588
- 589 Max LN Gonçalves, Jefferson G Melo, and Renato DC Monteiro. Improved pointwise iteration-
590 complexity of a regularized admm and of a regularized non-euclidean hpe framework. *SIAM
Journal on Optimization*, 27(1):379–407, 2017b.
591
- 592 Pinghua Gong, Changshui Zhang, Zhaosong Lu, Jianhua Huang, and Jieping Ye. A general iterative
593 shrinkage and thresholding algorithm for non-convex regularized optimization problems. In *In-
ternational Conference on Machine Learning, ICML 2013, Atlanta, GA, USA, 16-21 June 2013*,
volume 28, pp. 37–45, 2013.

- 594 Bingsheng He and Xiaoming Yuan. On the $\mathcal{O}(1/n)$ convergence rate of the douglas-rachford alter-
 595 nating direction method. *SIAM Journal on Numerical Analysis*, 50(2):700–709, 2012.
 596
- 597 Le Thi Khanh Hien, Duy Nhat Phan, and Nicolas Gillis. Inertial alternating direction method of mul-
 598 tipliers for non-convex non-smooth optimization. *Computational Optimization and Applications*,
 599 83(1):247–285, 2022.
- 600 Mingyi Hong, Zhi-Quan Luo, and Meisam Razaviyayn. Convergence analysis of alternating direc-
 601 tion method of multipliers for a family of nonconvex problems. *SIAM Journal on Optimization*,
 602 26(1):337–364, 2016.
- 603 Feihu Huang, Songcan Chen, and Heng Huang. Faster stochastic alternating direction method of
 604 multipliers for nonconvex optimization. In *International Conference on Machine Learning (ICM-
 605 L)*, volume 97, pp. 2839–2848, 2019.
- 606 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio
 607 and Yann LeCun (eds.), *International Conference on Learning Representations (ICLR)*, 2015.
- 609 Rongjie Lai and Stanley Osher. A splitting method for orthogonality constrained problems. *Journal
 610 of Scientific Computing*, 58(2):431–449, 2014.
- 612 Hien Le, Nicolas Gillis, and Panagiotis Patrinos. Inertial block proximal methods for non-convex
 613 non-smooth optimization. In *International Conference on Machine Learning*, pp. 5671–5681.
 614 PMLR, 2020.
- 615 Shuhuang Xiang Lei Yang, Xiaojun Chen. Sparse solutions of a class of constrained optimization
 616 problems. *Mathematics of Operations Research*, 2021.
- 617 Guoyin Li and Ting Kei Pong. Global convergence of splitting methods for nonconvex composite
 618 optimization. *SIAM Journal on Optimization*, 25(4):2434–2460, 2015.
- 620 Jiaxiang Li, Shiqian Ma, and Tejes Srivastava. A riemannian admm. *arXiv preprint arX-
 621 iv:2211.02163*, 2022.
- 623 Min Li, Defeng Sun, and Kim-Chuan Toh. A majorized admm with indefinite proximal terms for
 624 linearly constrained convex composite optimization. *SIAM Journal on Optimization*, 26(2):922–
 625 950, 2016.
- 626 Xiao Li, Shixiang Chen, Zengde Deng, Qing Qu, Zhihui Zhu, and Anthony Man-Cho So. Weakly
 627 convex optimization over stiefel manifold using riemannian subgradient-type methods. *SIAM
 628 Journal on Optimization*, 31(3):1605–1634, 2021.
- 630 Qihang Lin, Runchao Ma, and Yangyang Xu. Complexity of an inexact proximal-point penalty
 631 method for constrained smooth non-convex optimization. *Computational optimization and appli-
 632 cations*, 82(1):175–224, 2022.
- 633 Tian-Yi Lin, Shi-Qian Ma, and Shu-Zhong Zhang. On the sublinear convergence rate of multi-block
 634 admm. *Journal of the Operations Research Society of China*, 3:251–274, 2015a.
- 636 Tianyi Lin, Shiqian Ma, and Shuzhong Zhang. On the global linear convergence of the admm with
 637 multiblock variables. *SIAM Journal on Optimization*, 25(3):1478–1497, 2015b.
- 638 Dekai Liu, Song Li, and Yi Shen. One-bit compressive sensing with projected subgradient method
 639 under sparsity constraints. *IEEE Transactions on Information Theory*, 65(10):6650–6663, 2019.
- 640 Wei Liu, Xin Liu, and Xiaojun Chen. Linearly constrained nonsmooth optimization for training
 641 autoencoders. *SIAM Journal on Optimization*, 32(3):1931–1957, 2022.
- 643 Yuanyuan Liu, Fanhua Shang, Hongying Liu, Lin Kong, Licheng Jiao, and Zhouchen Lin. Acceler-
 644 ated variance reduction stochastic admm for large-scale machine learning. *IEEE Transactions on
 645 Pattern Analysis and Machine Intelligence*, 43(12):4242–4255, 2020.
- 646 Zhaosong Lu and Yong Zhang. An augmented lagrangian approach for sparse principal component
 647 analysis. *Mathematical Programming*, 135:149–193, 2012.

- 648 Zhaosong Lu and Yong Zhang. Sparse approximation via penalty decomposition methods. *SIAM
649 Journal on Optimization*, 23(4):2448–2478, 2013.
650
- 651 Renato DC Monteiro and Benar F Svaiter. Iteration-complexity of block-decomposition algorithms
652 and the alternating direction method of multipliers. *SIAM Journal on Optimization*, 23(1):475–
653 507, 2013.
- 654 Boris S. Mordukhovich. Variational analysis and generalized differentiation i: Basic theory. *Berlin
655 Springer*, 330, 2006.
656
- 657 Y. E. Nesterov. *Introductory lectures on convex optimization: a basic course*, volume 87 of *Applied
658 Optimization*. Kluwer Academic Publishers, 2003.
- 659 Robert Nishihara, Laurent Lessard, Ben Recht, Andrew Packard, and Michael Jordan. A general
660 analysis of the convergence of admm. In *International Conference on Machine Learning*, pp.
661 343–352. PMLR, 2015.
662
- 663 Yuyuan Ouyang, Yunmei Chen, Guanghui Lan, and Eduardo Pasiliao Jr. An accelerated linearized
664 alternating direction method of multipliers. *SIAM Journal on Imaging Sciences*, 8(1):644–681,
665 2015.
- 666 Duy Nhat Phan and Nicolas Gillis. An inertial block majorization minimization framework for
667 nonsmooth nonconvex optimization. *Journal of Machine Learning Research*, 24:1–41, 2023.
668
- 669 Thomas Pock and Shoham Sabach. Inertial proximal alternating linearized minimization (ipalm)
670 for nonconvex and nonsmooth problems. *SIAM Journal on Imaging Sciences*, 9(4):1756–1787,
671 2016.
- 672 R. Tyrrell Rockafellar and Roger J-B. Wets. Variational analysis. *Springer Science & Business
673 Media*, 317, 2009.
674
- 675 Li Shen, Wei Liu, Ganzhao Yuan, and Shiqian Ma. Gsos: Gauss-seidel operator splitting algo-
676 rithm for multi-term nonsmooth convex composite optimization. In *International Conference on
677 Machine Learning*, pp. 3125–3134. PMLR, 2017.
- 678 Kaizhao Sun and Xu Andy Sun. Algorithms for difference-of-convex programs based on difference-
679 of-moreau-envelopes smoothing. *INFORMS Journal on Optimization*, 5(4):321–339, 2023.
680
- 681 Quoc Tran Dinh. Non-ergodic alternating proximal augmented lagrangian algorithms with optimal
682 rates. *Advances in Neural Information Processing Systems*, 31, 2018.
- 683 Manolis C. Tsakiris and René Vidal. Dual principal component pursuit. *J. Mach. Learn. Res.*, 19:
684 18:1–18:50, 2018. URL <https://jmlr.org/papers/v19/17-436.html>.
685
- 686 Junxiang Wang, Fuxun Yu, Xiang Chen, and Liang Zhao. ADMM for efficient deep learning with
687 global convergence. In *ACM International Conference on Knowledge Discovery & Data Mining
(SIGKDD)*, pp. 111–119, 2019a.
- 688 Yu Wang, Wotao Yin, and Jinshan Zeng. Global convergence of admm in nonconvex nonsmooth
689 optimization. *Journal of Scientific Computing*, 78(1):29–63, 2019b.
690
- 691 Yi Xu, Mingrui Liu, Qihang Lin, and Tianbao Yang. Admm without a fixed penalty parameter:
692 Faster convergence with new adaptive penalization. In *Advances in Neural Information Process-
693 ing Systems*, volume 30. Curran Associates, Inc., 2017.
694
- 695 Lei Yang, Ting Kei Pong, and Xiaojun Chen. Alternating direction method of multipliers for a class
696 of nonconvex and nonsmooth problems with applications to background/foreground extraction.
697 *SIAM Journal on Imaging Sciences*, 10(1):74–110, 2017.
698
- 699 Maryam Yashtini. Multi-block nonconvex nonsmooth proximal admm: Convergence and rates un-
700 der kurdyka–lojasiewicz property. *Journal of Optimization Theory and Applications*, 190(3):
701 966–998, 2021. doi: 10.1007/s10957-021-01919-7. URL <https://doi.org/10.1007/s10957-021-01919-7>.

702 Maryam Yashtini. Convergence and rate analysis of a proximal linearized ADMM for nonconvex
 703 nonsmooth optimization. *Journal of Global Optimization*, 84(4):913–939, 2022.

704
 705 Jinshan Zeng, Shaobo Lin, Yao Wang, and Zongben Xu. $l_{1/2}$ regularization: Convergence of iterative
 706 half thresholding algorithm. *IEEE Trans. Signal Process.*, 62(9):2317–2329, 2014.

707
 708 Jinshan Zeng, Shao-Bo Lin, Yuan Yao, and Ding-Xuan Zhou. On ADMM in deep learning: Con-
 709 vergence and saturation-avoidance. *Journal of Machine Learning Research*, 22:199:1–199:67,
 710 2021.

711
 712 Jinshan Zeng, Wotao Yin, and Ding-Xuan Zhou. Moreau envelope augmented lagrangian method
 713 for nonconvex optimization with linear constraints. *Journal of Scientific Computing*, 91(2):61,
 714 2022.

715
 716 Jiawei Zhang and Zhi-Quan Luo. A proximal alternating direction method of multiplier for linearly
 717 constrained nonconvex minimization. *SIAM Journal on Optimization*, 30(3):2272–2302, 2020.

718
 719 Jiawei Zhang, Peijun Xiao, Ruoyu Sun, and Zhi-Quan Luo. A single-loop smoothed gradient
 720 descent-ascent algorithm for nonconvex-concave min-max problems. In Hugo Larochelle,
 721 Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in
 722 Neural Information Processing Systems*, 2020.

723
 724 Ruiliang Zhang and James Kwok. Asynchronous distributed admm for consensus optimization. In
 725 *International Conference on Machine Learning*, pp. 1701–1709. PMLR, 2014.

726
 727 Daoli Zhu, Lei Zhao, and Shuzhong Zhang. A first-order primal-dual method for nonconvex con-
 728 strained optimization based on the augmented lagrangian. *Mathematics of Operations Research*,
 729 2023.

730
 731
 732
 733
 734
 735
 736
 737
 738
 739
 740
 741
 742
 743
 744
 745
 746
 747
 748
 749
 750
 751
 752
 753
 754
 755

756 Appendix

758 The organization of the appendix is as follows:

759 Appendix A covers notations, technical preliminaries, and relevant lemmas.

760 Appendix B provides additional motivating applications.

761 Appendix C contains proofs related to Section 3.

762 Appendix D offers proofs related to Section 4.

763 Appendix E includes additional experiments details and results.

767 A NOTATIONS, TECHNICAL PRELIMINARIES, AND RELEVANT LEMMAS

768 A.1 NOTATIONS

769 We use the following notations in this paper.

- 773 • $[n]: \{1, 2, \dots, n\}.$
- 774 • $\mathbf{x}: \mathbf{x} \triangleq \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\} = \mathbf{x}_{[n]}.$
- 775 • $\mathbf{x}_{[i,j]}: \mathbf{x}_{[i,j]} \triangleq \{\mathbf{x}_i, \mathbf{x}_{i+1}, \mathbf{x}_{i+2}, \dots, \mathbf{x}_j\}$, where $j \geq i$.
- 776 • $L_i^t: L_i^t = L_i + \beta^t \|\mathbf{A}_i\|_2^2$. Note that the function $G(\mathbf{x}, \mathbf{z}^t; \beta^t)$ is L_i^t -smooth w.r.t. \mathbf{x} .
- 777 • $\sigma_1: \sigma_1 \triangleq \frac{\sigma}{(1-|1-\sigma|)^2} \in \mathbb{R}$, where $\sigma \in (0, 2)$. Refer to Lemma A.2.
- 778 • $\sigma_2: \sigma_2 \triangleq \frac{|1-\sigma|}{\sigma(1-|1-\sigma|)} \in \mathbb{R}$, where $\sigma \in (0, 2)$. Refer to Lemma A.2.
- 779 • $\|\mathbf{x}\|$: Euclidean norm: $\|\mathbf{x}\| = \|\mathbf{x}\|_2 = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}.$
- 780 • $\langle \mathbf{a}, \mathbf{b} \rangle$: Euclidean inner product, i.e., $\langle \mathbf{a}, \mathbf{b} \rangle = \sum_i \mathbf{a}_i \mathbf{b}_i$.
- 781 • \mathbf{A}^\top : the transpose of the matrix \mathbf{A} .
- 782 • \mathbf{x}_i : the i -th block of the vector $\mathbf{x} \in \mathbb{R}^{(d_1+d_2+\dots+d_n) \times 1}$ with $\mathbf{x}_i \in \mathbb{R}^{d_i \times 1}$.
- 783 • $\bar{\lambda}$: the largest eigenvalue of the matrix $\mathbf{A}_n \mathbf{A}_n^\top$.
- 784 • $\underline{\lambda}$: the smallest eigenvalue of the matrix $\mathbf{A}_n \mathbf{A}_n^\top$.
- 785 • $\underline{\lambda}'$: the smallest eigenvalue of the matrix $\mathbf{A}_n^\top \mathbf{A}_n$.
- 786 • $\|\mathbf{A}\|$: the spectral norm of the matrix \mathbf{A} .
- 787 • \mathbf{I}_r : $\mathbf{I}_r \in \mathbb{R}^{r \times r}$, Identity matrix; the subscript is omitted sometimes.
- 788 • $\iota_\Omega(\mathbf{x})$: Indicator function of a set Ω with $\iota_\Omega(\mathbf{x}) = 0$ if $\mathbf{x} \in \Omega$ and otherwise $+\infty$.
- 789 • $\text{vec}(\mathbf{V})$: Vector formed by stacking the column vectors of \mathbf{V} with $\text{vec}(\mathbf{V}) \in \mathbb{R}^{d' \times r'}$.
- 790 • $\text{mat}(\mathbf{x})$: Convert $\mathbf{x} \in \mathbb{R}^{(d' \cdot r') \times 1}$ into a matrix with $\text{mat}(\text{vec}(\mathbf{V})) = \mathbf{V}$ with $\text{mat}(\mathbf{x}) \in \mathbb{R}^{d' \times r'}$.
- 791 • $\text{dist}^2(\Omega, \Omega')$: squared distance between two sets with $\text{dist}^2(\Omega, \Omega') \triangleq \inf_{\mathbf{w} \in \Omega, \mathbf{w}' \in \Omega'} \|\mathbf{w} - \mathbf{w}'\|_2^2$.

792 A.2 TECHNICAL PRELIMINARIES

800 We present some tools in non-smooth analysis including Fréchet subdifferential, limiting (Fréchet)
 801 subdifferential, and directional derivative (Mordukhovich, 2006; Rockafellar & Wets., 2009; Bert-
 802 sekas, 2015). For any extended real-valued (not necessarily convex) function $F : \mathbb{R}^n \rightarrow$
 803 $(-\infty, +\infty]$, its domain is defined by $\text{dom}(F) \triangleq \{\mathbf{x} \in \mathbb{R}^n : |F(\mathbf{x})| < +\infty\}$. The Fréchet
 804 subdifferential of F at $\mathbf{x} \in \text{dom}(F)$, denoted as $\hat{\partial}F(\mathbf{x})$, is defined as $\hat{\partial}F(\mathbf{x}) \triangleq \{\mathbf{v} \in \mathbb{R}^n : \lim_{\mathbf{z} \rightarrow \mathbf{x}} \inf_{\mathbf{z} \neq \mathbf{x}} \frac{F(\mathbf{z}) - F(\mathbf{x}) - \langle \mathbf{v}, \mathbf{z} - \mathbf{x} \rangle}{\|\mathbf{z} - \mathbf{x}\|} \geq 0\}$. The limiting subdifferential of $F(\mathbf{x})$ at $\mathbf{x} \in \text{dom}(F)$
 805 is defined as: $\partial F(\mathbf{x}) \triangleq \{\mathbf{v} \in \mathbb{R}^n : \exists \mathbf{x}^k \rightarrow \mathbf{x}, F(\mathbf{x}^k) \rightarrow F(\mathbf{x}), \mathbf{v}^k \in \hat{\partial}F(\mathbf{x}^k) \rightarrow \mathbf{v}, \forall k\}$. Note that
 806 $\hat{\partial}F(\mathbf{x}) \subseteq \partial F(\mathbf{x})$. If $F(\cdot)$ is differentiable at \mathbf{x} , then $\hat{\partial}F(\mathbf{x}) = \partial F(\mathbf{x}) = \{\nabla F(\mathbf{x})\}$ with $\nabla F(\mathbf{x})$
 807 being the gradient of $F(\cdot)$ at \mathbf{x} . When $F(\cdot)$ is convex, $\hat{\partial}F(\mathbf{x})$ and $\partial F(\mathbf{x})$ reduce to the classical sub-
 808 differential for convex functions, i.e., $\hat{\partial}F(\mathbf{x}) = \partial F(\mathbf{x}) = \{\mathbf{v} \in \mathbb{R}^n : F(\mathbf{z}) - F(\mathbf{x}) - \langle \mathbf{v}, \mathbf{z} - \mathbf{x} \rangle \geq$
 809

810 $0, \forall \mathbf{z} \in \mathbb{R}^n\}$. The directional derivative of $F(\cdot)$ at \mathbf{x} in the direction \mathbf{v} is defined (if it exists) by
 811 $F'(\mathbf{x}; \mathbf{v}) \triangleq \lim_{t \rightarrow 0^+} \frac{1}{t}(F(\mathbf{x} + t\mathbf{v}) - F(\mathbf{x}))$.
 812

813 A.3 RELEVANT LEMMAS

815 We present several useful lemmas, each independent of context and specific methodology.

816 **Lemma A.1.** (Pythagoras Relation) For any vectors $\mathbf{a} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^n$, $\mathbf{c} \in \mathbb{R}^n$, we have:

$$\begin{aligned} \frac{1}{2}\|\mathbf{a} - \mathbf{b}\|_2^2 - \frac{1}{2}\|\mathbf{c} - \mathbf{b}\|_2^2 &= \frac{1}{2}\|\mathbf{a} - \mathbf{c}\|_2^2 + \langle \mathbf{b} - \mathbf{c}, \mathbf{c} - \mathbf{a} \rangle. \\ \frac{1}{2}\|\mathbf{b}\|_2^2 - \frac{1}{2}\|\mathbf{c} - \mathbf{b}\|_2^2 &= \frac{1}{2}\|\mathbf{c}\|_2^2 + \langle \mathbf{b} - \mathbf{c}, \mathbf{c} \rangle. \end{aligned}$$

821 **Lemma A.2.** Assume $\sigma \in (0, 2)$. Let $\mathbf{b}^+ = \sigma\mathbf{a} + (1 - \sigma)\mathbf{b}$, where $\mathbf{b}^+ \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^n$, and $\mathbf{a} \in \mathbb{R}^n$.
 822 We have:

$$\frac{1}{\sigma}\|\mathbf{b}^+\|_2^2 \leq \sigma_1\|\mathbf{a}\|_2^2 + \sigma_2(\|\mathbf{b}\|_2^2 - \|\mathbf{b}^+\|_2^2),$$

825 where $\sigma_1 \triangleq \frac{\sigma}{(1 - |1 - \sigma|)^2}$, and $\sigma_2 \triangleq \frac{|1 - \sigma|}{\sigma(1 - |1 - \sigma|)}$.

827 *Proof.* (a) When $\sigma = 1$, we have $\sigma_1 = 1$, $\sigma_2 = 0$, and $\mathbf{b}^+ = \mathbf{a}$. The conclusion of this lemma
 828 clearly holds.

829 (b) We now focus on the case when $\sigma \neq 1$. Noticing $|1 - \sigma| \neq 0$ and $1 - |1 - \sigma| \neq 0$, we rewrite
 830 $\mathbf{b}^+ = (1 - \sigma)\mathbf{b} + \sigma\mathbf{a}$ into the following equivalent equality

$$\mathbf{b}^+ = (1 - |1 - \sigma|) \cdot \frac{\sigma\mathbf{a}}{1 - |1 - \sigma|} + |1 - \sigma| \cdot \frac{(1 - \sigma)\mathbf{b}}{|1 - \sigma|}.$$

834 Using the fact that the function $\|\cdot\|_2^2$ is convex and $|1 - \sigma| \in (0, 1)$, we derive the following results:

$$\begin{aligned} \|\mathbf{b}^+\|_2^2 &\leq (1 - |1 - \sigma|) \cdot \left\| \frac{\sigma\mathbf{a}}{1 - |1 - \sigma|} \right\|_2^2 + |1 - \sigma| \cdot \left\| \frac{(1 - \sigma)\mathbf{b}}{|1 - \sigma|} \right\|_2^2 \\ &\leq \frac{\sigma^2}{1 - |1 - \sigma|} \cdot \|\mathbf{a}\|_2^2 + |1 - \sigma| \cdot \|\mathbf{b}\|_2^2. \end{aligned}$$

839 Subtracting $(|1 - \sigma| \cdot \|\mathbf{b}^+\|_2^2)$ from both sides of the above inequality, we have:

$$(1 - |1 - \sigma|)\|\mathbf{b}^+\|_2^2 \leq \frac{\sigma^2}{1 - |1 - \sigma|} \cdot \|\mathbf{a}\|_2^2 + |1 - \sigma|(\|\mathbf{b}\|_2^2 - \|\mathbf{b}^+\|_2^2).$$

843 Dividing both sides by $\sigma(1 - |1 - \sigma|)$, we have:

$$\frac{1}{\sigma}\|\mathbf{b}^+\|_2^2 \leq \frac{\sigma}{(1 - |1 - \sigma|)^2} \|\mathbf{a}\|_2^2 + \frac{|1 - \sigma|}{\sigma(1 - |1 - \sigma|)} (\|\mathbf{b}\|_2^2 - \|\mathbf{b}^+\|_2^2).$$

846 Using the definition of σ_1 and σ_2 , we finish the proof of this lemma.

847 \square

849 **Lemma A.3.** We let $t \geq 1$, and $q \in (0, 1)$. We have: $\frac{1}{q}(t + 1)^q - \frac{1}{q} \geq \frac{1}{2}t^q$.

851 *Proof.* We let $h(t) \triangleq (t + 1)^q - 1 - \frac{q}{2}t^q$.

853 Initially, we prove that $f(q) \triangleq 2^q - \frac{q}{2} - 1 \geq 0$ for all $q \geq 0$. Given $\nabla f(q) = 2^q \log(2) - \frac{1}{2} \geq 2^0 \log(2) - \frac{1}{2} = 0.1931 > 0$, the function $f(q)$ is increasing for all $q \geq 0$. Combining with the fact
 854 that $f(0) = 0$, we have: $f(q) \geq 0$ for all $q \geq 0$.

856 We derive the following inequalities:

$$\nabla h(t) = qt^{q-1} \cdot \left\{ \left(\frac{t+1}{t} \right)^{q-1} - \frac{q}{2} \right\} \stackrel{\textcircled{1}}{\geq} qt^{p-1} \cdot \left\{ 2^{q-1} - \frac{q}{2} \right\} \stackrel{\textcircled{2}}{\geq} qt^{q-1} \cdot \left\{ \frac{q/2+1}{2} - \frac{q}{2} \right\} \stackrel{\textcircled{3}}{\geq} 0,$$

860 where step ① uses $\frac{t+1}{t} \leq 2$ and $q - 1 \leq 0$; step ② uses $2^q \geq \frac{q}{2} + 1$ for all $q \geq 0$; step ③ uses
 861 $1 - q \geq 0$. Therefore, $h(t)$ is an increasing function.

862 Finally, noticing that $h(1) = 2^q - 1 - \frac{q}{2} \geq 0$, we conclude that $h(t) \geq 0$ for all $t \geq 1$.

863 \square

864 **Lemma A.4.** We let $p \in (0, 1)$ and $t \geq 1$. We have: $(t+1)^p - t^p \leq pt^{p-1}$.
 865

866 *Proof.* We notice that $h(t) \triangleq t^p$ is concave for all $t \geq 1$ and $p \in (0, 1)$ since $\nabla h(t) = pt^{p-1}$ and
 867 $\nabla^2 h(t) = p(p-1)t^{p-2} < 0$. It follows that: $\forall x, y \geq 1, h(y) - h(x) \leq \langle y-x, \nabla h(x) \rangle$. Letting
 868 $x = t$ and $y = t+1$, for all $t \geq 1$ and $p \in (0, 1)$, we have: $(t+1)^p - t^p \leq pt^{p-1}$.
 869

□

870 **Lemma A.5.** We let $p \in (0, 1)$. We have: $\sum_{t=1}^{\infty} \left(\frac{(t+1)^p - t^p}{t^p} \right)^2 \leq 2$.
 871

872 *Proof.* We have:

$$\sum_{t=1}^{\infty} \left(\frac{(t+1)^p - t^p}{t^p} \right)^2 \stackrel{\textcircled{1}}{\leq} \sum_{t=1}^{\infty} \frac{1}{t^{2p}} t^{2p-2} = \sum_{t=1}^{\infty} t^{-2} \stackrel{\textcircled{2}}{\leq} 2,$$

873 where step ① uses Lemma A.4 and $p \leq 1$; step ② uses $\sum_{t=1}^{\infty} \frac{1}{t^2} \leq \sum_{t=1}^{\infty} \frac{1}{t^2} = \frac{\pi^2}{6} < 2$.
 874

□

875 **Lemma A.6.** We let $p \in (0, 1)$. We have: $\frac{1}{2}T^{1-p} \leq \sum_{t=1}^T t^{-p} \leq \frac{T^{(1-p)}}{1-p}$.
 876

877 *Proof.* We define $h(x) = x^{-p}$ and $g(x) = \frac{1}{1-p}x^{1-p}$. Clearly, we have: $\nabla g(x) = h(x)$.
 878

879 By employing the integral test for convergence ¹, we obtain:

$$\int_1^{T+1} h(x)dx \leq \sum_{t=1}^T h(t) \leq h(1) + \int_1^T h(x)dx. \quad (14)$$

880 **(a)** We have: $\sum_{t=1}^T t^{-p} \stackrel{\textcircled{1}}{\geq} \int_1^{T+1} x^{-p}dx \stackrel{\textcircled{2}}{=} g(T+1) - g(1) = \frac{1}{1-p}(T+1)^{1-p} - \frac{1}{1-p} \stackrel{\textcircled{3}}{\geq} \frac{1}{2}T^{1-p}$,
 881 where step ① uses the first inequality in (14); step ② uses $\nabla g(x) = h(x) = x^{-p}$; step ③ uses
 882 Lemma A.3 with $q = 1 - p$ and $t = T$.
 883

□

884 **(b)** We have: $\sum_{t=1}^T t^{-p} \stackrel{\textcircled{1}}{\leq} h(1) + \int_1^T x^{-p}dx \stackrel{\textcircled{2}}{=} 1 + g(T) - g(1) = 1 + \frac{1}{1-p}(T)^{1-p} - \frac{1}{1-p} =$
 885 $\frac{T^{(1-p)} - p}{1-p} < \frac{T^{(1-p)}}{1-p}$, where step ① uses the second inequality in (14); step ② uses $h(1) = 1$, and
 886 $\nabla g(x) = h(x) = x^{-p}$.
 887

888 *Proof.* Given $\sigma \in (0, 2)$, we define $\sigma_* \triangleq |1 - \sigma| \in [0, 1]$.
 889

890 We derive the following results:
 891

$$\begin{aligned} t = 1, \quad e^2 &\leq \sigma_* e^1 + \sigma \Psi^1 \\ t = 2, \quad e^3 &\leq \sigma_* e^2 + \sigma \Psi^2 \leq \sigma_*^2 e^1 + \sigma_* \sigma \Psi^1 + \sigma \Psi^2 \\ t = 3, \quad e^4 &\leq \sigma_* e^3 + \sigma \Psi^3 \leq \sigma_*^3 e^1 + \sigma_*^2 \sigma \Psi^1 + \sigma_* \sigma \Psi^2 + \sigma \Psi^3 \\ &\dots \\ t = T, \quad e^{T+1} &\leq \sigma_* e^T + \sigma \Psi^T \leq \sigma_*^T e^1 + \sigma \sum_{i=1}^T \sigma_*^{T-i} \Psi^i. \end{aligned}$$

892 Therefore, we have:
 893

$$\begin{aligned} e^{T+1} &\leq \sigma_*^T e^1 + \sigma \sum_{i=1}^T \sigma_*^{T-i} \Psi^i \\ &\stackrel{\textcircled{1}}{\leq} e^1 + \sigma \{ \max_{i=1}^T \Psi^i \} \{ \sum_{i=1}^T \sigma_*^{T-i} \} \\ &\stackrel{\textcircled{2}}{\leq} e^1 + \sigma \{ \max_{i=1}^T \Psi^i \} \frac{1}{1-\sigma_*}, \end{aligned}$$

894 where step ① uses $\sigma_*^T \leq 1$; step ② uses the fact that:
 895

$$\sum_{i=1}^T \sigma_*^{T-i} = \sigma_*^{T-1} + \dots + \sigma_*^1 + \sigma_*^0 = \frac{1-\sigma_*^T}{1-\sigma_*} \leq \frac{1}{1-\sigma_*}.$$

□

¹https://en.wikipedia.org/wiki/Integral_test_for_convergence

918 **B ADDITIONAL MOTIVATING APPLICATIONS**
919

920 ► **Robust Sparse Regression.** Robust sparse regression (Liu et al., 2019) utilizes the ℓ_1 -norm of
921 the residuals to ensure robustness against outliers while enforcing sparsity via ℓ_0 -norm constraints to
922 identify key variables. The problem is formulated as: $\min_{\mathbf{v}} \|\mathbf{G}\mathbf{v} - \mathbf{z}\|_1$, s.t. $\mathbf{v} \in \Omega \triangleq \{\mathbf{v} \mid \|\mathbf{v}\|_0 \leq$
923 $\dot{s}\}$, where $\dot{s} \geq 0$ is an integer, $\mathbf{G} \in \mathbb{R}^{m \times d}$, and $\mathbf{z} \in \mathbb{R}^m$. By introducing a new variable \mathbf{y} , this
924 problem can be formulated as: $\min_{\mathbf{v}, \mathbf{y}} \iota_\Omega(\mathbf{v}) + \|\mathbf{y}\|_1$, s.t. $-\mathbf{G}\mathbf{v} + \mathbf{y} = -\mathbf{z}$. It corresponds to
925 Problem (1) with $\mathbf{x}_1 = \mathbf{v}$, $\mathbf{x}_2 = \mathbf{y}$, $f_1(\mathbf{x}_1) = f_2(\mathbf{x}_2) = 0$, $h_1(\mathbf{x}_1) = \iota_\Omega(\mathbf{v})$, $h_2(\mathbf{x}_2) = \|\mathbf{y}\|_1$, and
926 $\mathbf{A}_1 = -\mathbf{G}$, $\mathbf{A}_2 = \mathbf{I}$, $\mathbf{b} = -\mathbf{z}$, and Condition $\mathbb{B}\mathbb{I}$.

927 ► **Dual Principal Component Pursuit.** Dual principal component pursuit (Tsakiris & Vidal,
928 2018) is used primarily in subspace clustering and outlier detection, aiming to robustly represent
929 data structures across different subspaces in the presence of noise and outliers. The problem
930 is formulated as: $\min_{\mathbf{V}} \|\mathbf{G}\mathbf{V}\|_{2,1}$, s.t. $\mathbf{V} \in \Omega \triangleq \{\mathbf{V} \mid \mathbf{V}^\top \mathbf{V} = \mathbf{I}\}$, where $\mathbf{G} \in \mathbb{R}^{m \times d}$, and
931 $\|\mathbf{Y}\|_{2,1} \triangleq \sum_i \|\mathbf{Y}(i, :) \|$. By introducing a new variable \mathbf{Y} , this problem can be formulated as:
932 $\min_{\mathbf{V}, \mathbf{Y}} \iota_\Omega(\mathbf{V}) + \|\mathbf{Y}\|_{2,1}$, s.t. $-\mathbf{G}\mathbf{V} + \mathbf{Y} = \mathbf{0}$. It corresponds to Problem (1) with $\mathbf{x}_1 = \text{vec}(\mathbf{V})$,
933 $\mathbf{x}_2 = \text{vec}(\mathbf{Y})$, $f_1(\mathbf{x}_1) = f_2(\mathbf{x}_1) = 0$, $h_1(\mathbf{x}_1) = \iota_\Omega(\mathbf{V})$, $h_2(\mathbf{x}_2) = \|\mathbf{Y}\|_{2,1}$, and $\mathbf{A}_1 = -\mathbf{G}$,
934 $\mathbf{A}_2 = \mathbf{I}$, $\mathbf{b} = \mathbf{0}$, and Condition $\mathbb{B}\mathbb{I}$.

935 ► **Robust Low-Rank Approximation.** Robust low-rank approximation (Candès et al., 2011) uses
936 the ℓ_1 -norm of the residuals to ensure robustness against outliers while imposing a low-rank
937 constraint on the solution matrix. The problem is formulated as: $\min_{\mathbf{V}} \|\mathbf{G}(\mathbf{V}) - \mathbf{z}\|_1$, s.t. $\mathbf{V} \triangleq$
938 $\{\mathbf{V} \mid \text{rank}(\mathbf{V}) \leq \dot{s}\}$, where $\dot{s} \geq 0$ is an integer, $\mathbf{G}(\cdot) : \mathbb{R}^{d \times r} \mapsto \mathbb{R}^m$, and $\mathbf{z} \in \mathbb{R}^m$. By introducing a
939 new variable \mathbf{y} , this problem can be formulated as: $\min_{\mathbf{V}, \mathbf{y}} \iota_\Omega(\mathbf{V}) + \|\mathbf{y}\|_1$, s.t. $-\mathbf{G}(\mathbf{V}) + \mathbf{y} = -\mathbf{z}$. It corresponds to Problem (1) with $\mathbf{x}_1 = \text{vec}(\mathbf{V})$, $\mathbf{x}_2 = \mathbf{y}$, $f_1(\mathbf{x}_1) = f_2(\mathbf{x}_1) = 0$, $h_1(\mathbf{x}_1) = \iota_\Omega(\mathbf{V})$,
940 $h_2(\mathbf{x}_2) = \|\mathbf{y}\|_1$, $\mathbf{A}_1 \mathbf{x}_1 = -\mathbf{G}(\mathbf{V})$, $\mathbf{A}_2 = \mathbf{I}$, $\mathbf{b} = -\mathbf{z}$, and Condition $\mathbb{B}\mathbb{I}$.

943 **C PROOFS FOR SECTION 3**

944 **C.1 PROOF OF LEMMA 3.1**

945 *Proof.* Consider the update rule $\beta^t = \beta^0 + \beta^0 \xi t^p$, where $p \in (0, 1)$.

946 (a) We have:

947
$$\beta^{t+1} - \beta^t - \xi \beta^t \stackrel{\textcircled{1}}{=} \beta^0 \xi ((t+1)^p - t^p) - \xi \beta^0 \stackrel{\textcircled{2}}{\leq} \beta^0 \xi - \beta^0 \xi = 0,$$

948 where step ① uses the update rule $\beta^t = \beta^0 + \beta^0 \xi t^p$; step ② uses the fact that the function $h(t) \triangleq$
949 $(t+1)^p - t^p$ is monotonically decreasing w.r.t. t that: $h(t) \leq h(0) = 1$.

950 (b) We derive: $L_n \leq \beta^0 \delta \bar{\lambda} \stackrel{\textcircled{1}}{\leq} \beta^t \delta \bar{\lambda}$, where step ① uses $\beta^t \geq \beta^0$. □

951 **C.2 PROOF OF LEMMA 3.4**

952 *Proof.* We let \mathbf{u} be a fixed constant vector. We assume $0 < \mu_2 < \mu_1$.

953 We define: $h(\mathbf{u}; \mu) \triangleq \min_{\mathbf{v} \in \mathbb{R}^{d \times 1}} h(\mathbf{v}) + \frac{1}{2\mu} \|\mathbf{v} - \mathbf{u}\|_2^2$.

954 We define $\mathbb{P}_h(\mathbf{u}; \mu) \triangleq \arg \min_{\mathbf{v} \in \mathbb{R}^{d \times 1}} h(\mathbf{v}) + \frac{1}{2\mu} \|\mathbf{v} - \mathbf{u}\|_2^2$.

955 Initially, by the optimality of $\mathbb{P}_h(\mathbf{u}; \mu_1)$ and $\mathbb{P}_h(\mathbf{u}; \mu_2)$, we obtain:

956
$$\mathbf{u} - \mathbb{P}_h(\mathbf{u}; \mu_1) \in \mu_1 \partial h(\mathbb{P}_h(\mathbf{u}; \mu_1)), \tag{15}$$

957
$$\mathbf{u} - \mathbb{P}_h(\mathbf{u}; \mu_2) \in \mu_2 \partial h(\mathbb{P}_h(\mathbf{u}; \mu_2)). \tag{16}$$

958 For notation simplicity, we define:

959
$$\mathbf{p}_1 \triangleq \mathbb{P}_h(\mathbf{u}; \mu_1), \mathbf{g}_1 \in \partial h(\mathbb{P}_h(\mathbf{u}; \mu_1))$$

960
$$\mathbf{p}_2 \triangleq \mathbb{P}_h(\mathbf{u}; \mu_2), \mathbf{g}_2 \in \partial h(\mathbb{P}_h(\mathbf{u}; \mu_2)).$$

972 Equations (15) and (16) can be rewritten as:
 973

$$\mathbf{u} - \mathbf{p}_1 = \mu_1 \mathbf{g}_1, \quad (17)$$

$$\mathbf{u} - \mathbf{p}_2 = \mu_2 \mathbf{g}_2. \quad (18)$$

977 (a) We now prove that $0 \leq \frac{h(\mathbf{u}; \mu_2) - h(\mathbf{u}; \mu_1)}{\mu_1 - \mu_2}$. We have:
 978

$$\begin{aligned} h(\mathbf{u}; \mu_1) - h(\mathbf{u}; \mu_2) &\stackrel{\textcircled{1}}{=} \frac{1}{2\mu_1} \|\mathbf{u} - \mathbf{p}_1\|_2^2 - \frac{1}{2\mu_2} \|\mathbf{u} - \mathbf{p}_2\|_2^2 + h(\mathbf{p}_1) - h(\mathbf{p}_2) \\ &\stackrel{\textcircled{2}}{\leq} \frac{1}{2\mu_1} \|\mathbf{u} - \mathbf{p}_1\|_2^2 - \frac{1}{2\mu_2} \|\mathbf{u} - \mathbf{p}_2\|_2^2 + \langle \mathbf{p}_1 - \mathbf{p}_2, \mathbf{g}_1 \rangle \\ &\stackrel{\textcircled{3}}{=} \frac{\mu_1}{2} \|\mathbf{g}_1\|_2^2 - \frac{\mu_2}{2} \|\mathbf{g}_2\|_2^2 + \langle \mu_2 \mathbf{g}_2 - \mu_1 \mathbf{g}_1, \mathbf{g}_1 \rangle \\ &= -\frac{\mu_1}{2} \|\mathbf{g}_1\|_2^2 - \frac{\mu_2}{2} \|\mathbf{g}_2\|_2^2 + \mu_2 \langle \mathbf{g}_2, \mathbf{g}_1 \rangle \\ &\stackrel{\textcircled{4}}{\leq} -\frac{\mu_2}{2} \|\mathbf{g}_1\|_2^2 - \frac{\mu_2}{2} \|\mathbf{g}_2\|_2^2 + \mu_2 \langle \mathbf{g}_2, \mathbf{g}_1 \rangle \\ &= -\frac{\mu_2}{2} \|\mathbf{g}_2 - \mathbf{g}_1\|_2^2 \leq 0, \end{aligned}$$

988 where step ① uses the definition of $h(\mathbf{u}; \mu)$; step ② uses the convexity of $h(\cdot)$; step ③ uses the
 989 optimality of $\mathbf{p}_1 \triangleq \mathbb{P}_h(\mathbf{u}; \mu_1)$ and $\mathbf{p}_2 \triangleq \mathbb{P}_h(\mathbf{u}; \mu_2)$ as in (17) and (18); step ④ uses $\mu_2 < \mu_1$.
 990

991 (b) We now prove that $\frac{h(\mathbf{u}; \mu_2) - h(\mathbf{u}; \mu_1)}{\mu_1 - \mu_2} \leq \frac{1}{2} C_g^2$. We have:
 992

$$\begin{aligned} h(\mathbf{u}; \mu_2) - h(\mathbf{u}; \mu_1) &\stackrel{\textcircled{1}}{=} \frac{1}{2\mu_2} \|\mathbf{u} - \mathbf{p}_2\|_2^2 - \frac{1}{2\mu_1} \|\mathbf{u} - \mathbf{p}_1\|_2^2 + h(\mathbf{p}_2) - h(\mathbf{p}_1) \\ &\stackrel{\textcircled{2}}{\leq} \frac{1}{2\mu_2} \|\mathbf{u} - \mathbf{p}_2\|_2^2 - \frac{1}{2\mu_1} \|\mathbf{u} - \mathbf{p}_1\|_2^2 + \langle \mathbf{p}_2 - \mathbf{p}_1, \mathbf{g}_2 \rangle \\ &\stackrel{\textcircled{3}}{=} \frac{\mu_2}{2} \|\mathbf{g}_2\|_2^2 - \frac{\mu_1}{2} \|\mathbf{g}_1\|_2^2 + \langle \mu_1 \mathbf{g}_1 - \mu_2 \mathbf{g}_2, \mathbf{g}_2 \rangle \\ &= -\frac{\mu_2}{2} \|\mathbf{g}_2\|_2^2 - \frac{\mu_1}{2} \|\mathbf{g}_1\|_2^2 + \mu_1 \langle \mathbf{g}_2, \mathbf{g}_1 \rangle \\ &\stackrel{\textcircled{4}}{\leq} \frac{\mu_1}{2} \|\mathbf{g}_2\|_2^2 - \frac{\mu_2}{2} \|\mathbf{g}_2\|_2^2 \\ &\stackrel{\textcircled{5}}{\leq} \frac{\mu_1 - \mu_2}{2} \cdot C_g^2, \end{aligned}$$

1003 where step ① uses the definition of $h(\mathbf{u}; \mu)$; step ② uses the convexity of $h(\cdot)$; step ③ uses the
 1004 optimality of $\mathbf{p}_1 \triangleq \mathbb{P}_h(\mathbf{u}; \mu_1)$ and $\mathbf{p}_2 \triangleq \mathbb{P}_h(\mathbf{u}; \mu_2)$ as in (17) and (18); step ④ uses the inequality
 1005 that: $-\frac{1}{2} \|\mathbf{g}_1\|_2^2 + \langle \mathbf{g}_1, \mathbf{g}_2 \rangle \leq \frac{1}{2} \|\mathbf{g}_2\|_2^2$ for all $\mathbf{g}_1 \in \mathbb{R}^{d \times 1}$ and $\mathbf{g}_2 \in \mathbb{R}^{d \times 1}$; step ⑤ uses $\|\mathbf{g}_2\| \leq C_g$. \square
 1006

1007 C.3 PROOF OF LEMMA 3.5

1008 *Proof.* We let \mathbf{u} be a fixed constant vector. We assume $0 < \mu_2 < \mu_1$.
 1009

1010 We define: $h(\mathbf{u}; \mu) \triangleq \min_{\mathbf{v} \in \mathbb{R}^{d \times 1}} h(\mathbf{v}) + \frac{1}{2\mu} \|\mathbf{v} - \mathbf{u}\|_2^2$.
 1011

1012 We define: $\mathbb{P}_h(\mathbf{u}; \mu) \triangleq \arg \min_{\mathbf{v} \in \mathbb{R}^{d \times 1}} h(\mathbf{v}) + \frac{1}{2\mu} \|\mathbf{v} - \mathbf{u}\|_2^2$.
 1013

1014 Using Claim (b) of Lemma 3.3, we establish that $h(\mathbf{u}; \mu)$ is smooth w.r.t. \mathbf{u} , and its gradient can be
 1015 computed as:
 1016

$$\nabla h(\mathbf{u}; \mu) = \mu^{-1} (\mathbf{u} - \mathbb{P}_h(\mathbf{u}; \mu)).$$

1017 We examine the following mapping $\mathcal{H}(v) \triangleq v(\mathbf{u} - \mathbb{P}_h(\mathbf{u}; \frac{1}{v}))$ with $\mathcal{H}(v) : \mathbb{R} \mapsto \mathbb{R}^n$. We derive:
 1018

$$\begin{aligned} \lim_{\delta \rightarrow 0} \frac{\mathcal{H}(v+\delta) - \mathcal{H}(v)}{\delta} &= \lim_{\delta \rightarrow 0} \frac{(v+\delta)(\mathbf{u} - \mathbb{P}_h(\mathbf{u}; \frac{1}{v+\delta})) - v(\mathbf{u} - \mathbb{P}_h(\mathbf{u}; \frac{1}{v}))}{\delta} \\ &= \lim_{\delta \rightarrow 0} \frac{\delta \mathbf{u} - (v+\delta)\mathbb{P}_h(\mathbf{u}; \frac{1}{v}) + v\mathbb{P}_h(\mathbf{u}; \frac{1}{v})}{\delta} = \mathbf{u} - \mathbb{P}_h(\mathbf{u}; \frac{1}{v}). \end{aligned}$$

1023 Therefore, the first-order derivative of the mapping $\mathcal{H}(v)$ w.r.t. v always exists and can be computed
 1024 as $\nabla_v \mathcal{H}(v) = \mathbf{u} - \mathbb{P}_h(\mathbf{u}; \frac{1}{v})$, leading to:
 1025

$$\forall v, v' > 0, \frac{\|\mathcal{H}(v) - \mathcal{H}(v')\|}{|v - v'|} \leq \|\mathbf{u} - \mathbb{P}_h(\mathbf{u}; \frac{1}{v})\|.$$

1026 Letting $v = 1/\mu_1$ and $v' = 1/\mu_2$, we derive:
1027

$$1028 \frac{\|\nabla h(\mathbf{u}; \mu_1) - \nabla h(\mathbf{u}; \mu_2)\|}{|1/\mu_1 - 1/\mu_2|} \leq \|\mathbf{u} - \mathbb{P}_h(\mathbf{u}; \mu_1)\| \stackrel{\textcircled{1}}{=} \mu_1 \|\partial h(\mathbb{P}_h(\mathbf{u}; \mu_1))\| \stackrel{\textcircled{2}}{\leq} \mu_1 C_h,$$

1030 where step ① uses the optimality of $\mathbb{P}_h(\mathbf{u}; \mu)$ that $\mathbf{0} \in \partial h(\mathbb{P}_h(\mathbf{u}; \mu)) + \frac{1}{\mu}(\mathbb{P}_h(\mathbf{u}; \mu) - \mathbf{u})$ for all μ ;
1031 step ② uses the Lipschitz continuity of $h(\cdot)$. We further obtain:
1032

$$1033 \|\nabla h(\mathbf{u}; \mu_1) - \nabla h(\mathbf{u}; \mu_2)\| \leq |1/\mu_1 - 1/\mu_2| \mu_1 C_h = (\mu_1/\mu_2 - 1) \cdot C_h.$$

1034
1035
1036

C.4 PROOF OF LEMMA 3.6

1038 *Proof.* The proof of this lemma is similar to that of Lemma 1 in (Li et al., 2022). For completeness,
1039 we include the proof here.
1040

1041 We consider the following strongly convex problems:

$$1042 \bar{\mathbf{x}}_n = \arg \min_{\mathbf{x}_n} h_n(\mathbf{x}_n; \mu) + \frac{\rho}{2} \|\mathbf{x}_n - \mathbf{c}\|_2^2 \\ 1043 \Leftrightarrow (\bar{\mathbf{x}}_n, \check{\mathbf{x}}_n) = \arg \min_{\mathbf{x}_n, \check{\mathbf{x}}_n} h_n(\check{\mathbf{x}}_n) + \frac{1}{2\mu} \|\mathbf{x}_n - \check{\mathbf{x}}_n\|_2^2 + \frac{\rho}{2} \|\mathbf{x}_n - \mathbf{c}\|_2^2.$$

1044 We have the following first-order optimality conditions:
1045

$$1046 \mathbf{0} = \frac{1}{\mu}(\bar{\mathbf{x}}_n - \check{\mathbf{x}}_n) + \rho(\bar{\mathbf{x}}_n - \mathbf{c}) \quad (19)$$

$$1047 \mathbf{0} \in \partial h_n(\check{\mathbf{x}}_n) + \frac{1}{\mu}(\check{\mathbf{x}}_n - \bar{\mathbf{x}}_n). \quad (20)$$

1048 (a) Using (19), we obtain: $\bar{\mathbf{x}}_n = \frac{1}{1/\mu+\rho}(\frac{1}{\mu}\check{\mathbf{x}}_n + \rho\mathbf{c})$. Plugging this equation into (20) yields:
1049

$$1050 \mathbf{0} \in \partial h_n(\check{\mathbf{x}}_n) + \frac{1}{\mu}(\check{\mathbf{x}}_n - \frac{1}{1/\mu+\rho}(\frac{1}{\mu}\check{\mathbf{x}}_n + \rho\mathbf{c})) \\ 1051 = \partial h_n(\check{\mathbf{x}}_n) + \frac{\rho}{1+\mu\rho}(\check{\mathbf{x}}_n - \mathbf{c}).$$

1052 The inclusion above implies that:
1053

$$1054 \check{\mathbf{x}}_n = \arg \min_{\check{\mathbf{x}}_n} h_n(\check{\mathbf{x}}_n) + \frac{1}{2} \cdot \frac{\rho}{1+\mu\rho} \|\check{\mathbf{x}}_n - \mathbf{c}\|_2^2.$$

1055 (b) We derive:
1056

$$1057 -\rho(\bar{\mathbf{x}}_n - \mathbf{c}) \stackrel{\textcircled{1}}{=} \frac{1}{\mu}(\bar{\mathbf{x}}_n - \check{\mathbf{x}}_n) \stackrel{\textcircled{2}}{\in} \partial h_n(\check{\mathbf{x}}_n),$$

1058 where step ① uses (19); step ② uses (20).
1059

1060 (c) Using (20), we have: $\check{\mathbf{x}}_n - \bar{\mathbf{x}}_n = -\mu \partial h_n(\check{\mathbf{x}}_n)$. This leads to $\|\check{\mathbf{x}}_n - \bar{\mathbf{x}}_n\| \leq \mu C_h$.
1061

1062 \square

D PROOFS FOR SECTION 4

D.1 PROOF OF LEMMA 4.1

1063 *Proof.* (a) We now focus on sufficient decrease for variables $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n-1}\}$. We define $\Phi_i^t = G(\mathbf{x}_{[1,i-1]}^{t+1}, \mathbf{x}_i^{t+1}, \mathbf{x}_{[i+1,n]}^t, \mathbf{z}^t; \beta^t) - G(\mathbf{x}_{[1,i-1]}^{t+1}, \mathbf{x}_i^t, \mathbf{x}_{[i+1,n]}^t, \mathbf{z}^t; \beta^t) + h_i(\mathbf{x}_i^{t+1}) - h_i(\mathbf{x}_i^t)$, where $i \in [n-1]$.
1064

1065 Noticing the function $G(\mathbf{x}_{[1,i-1]}^{t+1}, \mathbf{x}_i, \mathbf{x}_{[i+1,n]}^t, \mathbf{z}^t; \beta^t)$ is L_i^t -smooth w.r.t. \mathbf{x}_i for the t -th iteration, we
1066 have:
1067

$$1068 G(\mathbf{x}_{[1,i-1]}^{t+1}, \mathbf{x}_i^{t+1}, \mathbf{x}_{[i+1,n]}^t, \mathbf{z}^t; \beta^t) - G(\mathbf{x}_{[1,i-1]}^{t+1}, \mathbf{x}_i^t, \mathbf{x}_{[i+1,n]}^t, \mathbf{z}^t; \beta^t) \\ 1069 \leq \langle \mathbf{x}_i^{t+1} - \mathbf{x}_i^t, \nabla_{\mathbf{x}_i} G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_i^t, \mathbf{x}_{[i+1,n]}^t, \mathbf{z}^t; \beta^t) \rangle + \frac{L_i^t}{2} \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2. \quad (21)$$

Given \mathbf{x}_i^{t+1} is the minimizer of the following optimization problem:

$$\mathbf{x}_i^{t+1} \in \arg \min_{\mathbf{x}_i} h_i(\mathbf{x}_i) + \langle \mathbf{x}_i - \mathbf{x}_i^t, \nabla_{\mathbf{x}_i} G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_i^t, \mathbf{x}_{[i+1,n]}^t, \mathbf{z}^t; \beta^t) \rangle + \frac{\theta_1 \mathsf{L}_i^t}{2} \|\mathbf{x}_i - \mathbf{x}_i^t\|_2^2.$$

The optimality of \mathbf{x}_i^{t+1} leads to:

$$h_i(\mathbf{x}_i^{t+1}) - h_i(\mathbf{x}_i^t) + \langle \mathbf{x}_i^{t+1} - \mathbf{x}_i^t, \nabla_{\mathbf{x}_i} G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_i^t, \mathbf{x}_{[i+1,n]}^t, \mathbf{z}^t; \beta^t) \rangle \leq -\frac{\theta_1 \mathsf{L}_i^t}{2} \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2. \quad (22)$$

Combining equations (21) and (22), we derive the following expressions:

$$\Phi_i^t \leq (\frac{1}{2} - \frac{\theta_1}{2}) \cdot \mathsf{L}_i^t \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2.$$

Telescoping the above inequality over i from 1 to $(n-1)$ leads to:

$$\sum_{i=1}^{n-1} \Phi_i^t \leq \sum_{i=1}^{n-1} \{(\frac{1}{2} - \frac{\theta_1}{2}) \cdot \mathsf{L}_i^t \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2\}.$$

Therefore, we obtain:

$$\mathcal{L}(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_n^t, \mathbf{z}^t; \beta^t, \mu^t) - \mathcal{L}(\mathbf{x}^t, \mathbf{z}^t; \beta^t, \mu^t) \leq \sum_{i=1}^{n-1} \{(\frac{1}{2} - \frac{\theta_1}{2}) \cdot \mathsf{L}_i^t \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2\}. \quad (23)$$

(b) We now focus on sufficient decrease for variable $\{\mathbf{x}_n\}$. Noticing the function $G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_n, \mathbf{z}^t; \beta^t)$ is L_n^t -smooth w.r.t. \mathbf{x}_n for the t -th iteration, we have:

$$\begin{aligned} & G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_n^{t+1}, \mathbf{z}^t; \beta^t) - G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_n^t, \mathbf{z}^t; \beta^t) \\ & \leq \langle \mathbf{x}_n^{t+1} - \mathbf{x}_n^t, \nabla_{\mathbf{x}_n} G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_n^t, \mathbf{z}^t; \beta^t) \rangle + \frac{\mathsf{L}_n^t}{2} \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2. \end{aligned} \quad (24)$$

Since $h_n(\mathbf{x}_n; \mu^t)$ is convex, we have:

$$\begin{aligned} & h_n(\mathbf{x}_n^{t+1}; \mu^t) - h_n(\mathbf{x}_n^t; \mu^t) \\ & \leq \langle \mathbf{x}_n^{t+1} - \mathbf{x}_n^t, \nabla h_n(\mathbf{x}_n^{t+1}; \mu^t) \rangle \\ & \stackrel{\textcircled{1}}{=} \langle \mathbf{x}_n^{t+1} - \mathbf{x}_n^t, -\nabla_{\mathbf{x}_n} G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_n^t, \mathbf{z}^t; \beta^t) \rangle - \theta_2 \mathsf{L}_n^t (\mathbf{x}_n^{t+1} - \mathbf{x}_n^t), \end{aligned} \quad (25)$$

where step ① uses the first-order optimality condition of \mathbf{x}_n^{t+1} that:

$$\mathbf{0} = \nabla h_n(\mathbf{x}_n^{t+1}; \mu^t) + \nabla_{\mathbf{x}_n} G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_n^t, \mathbf{z}^t; \beta^t) + \theta_2 \mathsf{L}_n^t (\mathbf{x}_n^{t+1} - \mathbf{x}_n^t).$$

Adding Inequalities (24) and (25) together, we have:

$$\begin{aligned} & h_n(\mathbf{x}_n^{t+1}; \mu^t) - h_n(\mathbf{x}_n^t; \mu^t) + G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_n^{t+1}, \mathbf{z}^t; \beta^t) - G(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_n^t, \mathbf{z}^t; \beta^t) \\ & \leq \frac{\mathsf{L}_n^t}{2} \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 - \theta_2 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 \\ & = (\frac{1}{2} - \theta_2) \cdot \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2. \end{aligned}$$

This results in the following inequality:

$$\mathcal{L}(\mathbf{x}^{t+1}, \mathbf{z}^t; \beta^t, \mu^t) - \mathcal{L}(\mathbf{x}_{[1,n-1]}^{t+1}, \mathbf{x}_n^t, \mathbf{z}^t; \beta^t, \mu^t) \leq (\frac{1}{2} - \theta_2) \cdot \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2. \quad (26)$$

(c) We now focus on sufficient decrease for variable $\{\mathbf{z}\}$. We have:

$$\begin{aligned} & \mathcal{L}(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}; \beta^t, \mu^t) - \mathcal{L}(\mathbf{x}^{t+1}, \mathbf{z}^t; \beta^t, \mu^t) \\ & = \langle \mathbf{A}\mathbf{x}^{t+1} - \mathbf{b}, \mathbf{z}^{t+1} - \mathbf{z}^t \rangle \\ & \stackrel{\textcircled{1}}{=} \langle \frac{1}{\sigma \beta^t} (\mathbf{z}^{t+1} - \mathbf{z}^t), \mathbf{z}^{t+1} - \mathbf{z}^t \rangle \\ & = \frac{1}{\sigma \beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2, \end{aligned} \quad (27)$$

where step ① uses $\mathbf{z}^{t+1} = \mathbf{z}^t + \sigma \beta^t (\mathbf{A}\mathbf{x}^{t+1} - \mathbf{b})$ with $\mathbf{A}\mathbf{x}^{t+1} \triangleq \sum_{j=1}^n \mathbf{A}_j \mathbf{x}_j^{t+1}$.

1134 (d) We now focus on sufficient decrease for variable $\{\beta\}$. We have:
 1135

$$\begin{aligned}
 & \mathcal{L}(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}; \beta^{t+1}, \mu^t) - \mathcal{L}(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}; \beta^t, \mu^t) \\
 &= (\frac{\beta^{t+1}}{2} - \frac{\beta^t}{2}) \|\mathbf{A}\mathbf{x}^{t+1} - \mathbf{b}\|_2^2 \\
 &\stackrel{\textcircled{1}}{=} (\frac{\beta^{t+1}}{2} - \frac{\beta^t}{2}) \|\frac{1}{\sigma\beta^t}(\mathbf{z}^{t+1} - \mathbf{z}^t)\|_2^2 \\
 &\stackrel{\textcircled{2}}{\leq} (\frac{(1+\xi)\beta^t}{2} - \frac{\beta^t}{2}) \|\frac{1}{\sigma\beta^t}(\mathbf{z}^{t+1} - \mathbf{z}^t)\|_2^2 \\
 &= \frac{\xi}{2\sigma} \cdot \frac{1}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2,
 \end{aligned} \tag{28}$$

1144 where step ① uses $\mathbf{z}^{t+1} = \mathbf{z}^t + \sigma\beta^t(\mathbf{A}\mathbf{x}^{t+1} - \mathbf{b})$; step ② uses Lemma 3.1 that $\beta^{t+1} \leq \beta^t(1 + \xi)$.
 1145

1146 (e) We now focus on sufficient decrease for variable $\{\mu\}$. We have:
 1147

$$\begin{aligned}
 & \mathcal{L}(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}; \beta^{t+1}, \mu^{t+1}) - \mathcal{L}(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}; \beta^{t+1}, \mu^t) \\
 &= h_n(\mathbf{x}_n^{t+1}; \mu^{t+1}) - h_n(\mathbf{x}_n^{t+1}; \mu^t) \\
 &\stackrel{\textcircled{1}}{\leq} \frac{1}{2}C_h(\mu^t - \mu^{t+1}),
 \end{aligned} \tag{29}$$

1152 where step ① uses Lemma 3.4.
 1153

Combining Inequalities (23), (26), (27), (28), and (29), we have:
 1154

$$\begin{aligned}
 & \mathcal{L}(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}; \beta^{t+1}, \mu^{t+1}) - \mathcal{L}(\mathbf{x}^t, \mathbf{z}^t; \beta^t, \mu^t) \\
 &\leq [\sum_{i=1}^{n-1} \{(\frac{1}{2} - \frac{\theta_1}{2}) \cdot \mathsf{L}_i^t \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2\}] + (\frac{1}{2} - \theta_2) \cdot \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 \\
 &\quad (1 + \frac{\xi}{2\sigma}) \cdot \frac{1}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + \frac{1}{2}C_h(\mu^t - \mu^{t+1})
 \end{aligned} \tag{30}$$

1159 We define $\Theta_L^t \triangleq \mathcal{L}(\mathbf{x}^t, \mathbf{z}^t; \beta^t, \mu^t) + \frac{1}{2}C_h\mu^t$, $\varepsilon_3 \triangleq \xi$, $\varepsilon_1 \triangleq \frac{1}{2}\theta_1 - \frac{1}{2}$, and $\mathcal{E}^{t+1} \triangleq \frac{\varepsilon_3}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + \varepsilon_2 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \varepsilon_1 \sum_{i=1}^{n-1} \mathsf{L}_i^t \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2$. We have:
 1160

$$\begin{aligned}
 & \mathcal{E}^{t+1} + \Theta_L^{t+1} - \Theta_L^t \\
 &\leq (\frac{1}{2} - \theta_2 + \varepsilon_2) \cdot \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + (1 + \frac{\xi}{2\sigma} + \sigma\xi) \cdot \frac{1}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2.
 \end{aligned}$$

1166 \square

1167 D.2 PROOF OF LEMMA 4.2

1169 *Proof.* For any $i \in [n]$, we define $\mathbf{u}_i^{t+1} \triangleq \theta_i \mathsf{L}_i^t [\mathbf{x}_i^{t+1} - \mathbf{x}_i^t] - \beta^t \mathbf{A}_i^\top [\sum_{j=i}^n \mathbf{A}_j (\mathbf{x}_j^{t+1} - \mathbf{x}_j^t)]$, and let
 1170 $\mathbf{w}_i^{t+1} \in \partial h_i(\mathbf{x}_i^{t+1}) + \nabla f_i(\mathbf{x}_i^t)$.
 1171

1172 We notice that \mathbf{x}_i^{t+1} is the minimizer of the following problem:
 1173

$$\mathbf{x}_i^{t+1} \in \arg \min_{\mathbf{x}_i} \frac{\theta \mathsf{L}_i^t}{2} \|\mathbf{x}_i - \mathbf{x}_i^t\|_2^2 + h_i(\mathbf{x}_i) + \langle \mathbf{x}_i - \mathbf{x}_i^t, \nabla_{\mathbf{x}_i} G(\mathbf{x}_{[1,i-1]}^{t+1}, \mathbf{x}_{[i,n]}^t, \mathbf{z}^t; \beta^t) \rangle.$$

1176 Using the necessary first-order optimality condition of the solution \mathbf{x}_i^{t+1} , we have:
 1177

$$\nabla_{\mathbf{x}_i} G(\mathbf{x}_{[1,i-1]}^{t+1}, \mathbf{x}_{[i,n]}^t, \mathbf{z}^t; \beta^t) \in -\partial h_i(\mathbf{x}_i^{t+1}) - \theta \mathsf{L}_i^t (\mathbf{x}_i^{t+1} - \mathbf{x}_i^t). \tag{31}$$

1180 Using the definition of the function $G(\mathbf{x}, \mathbf{z}; \beta) \triangleq \langle [\sum_{j=1}^n \mathbf{A}_j \mathbf{x}_j] - \mathbf{b}, \mathbf{z} \rangle + \frac{\beta}{2} \|[\sum_{j=1}^n \mathbf{A}_j \mathbf{x}_j] - \mathbf{b}\|_2^2 + \sum_{j=1}^n f_j(\mathbf{x}_j)$, we have:
 1181

$$\begin{aligned}
 & \nabla_{\mathbf{x}_i} G(\mathbf{x}_{[1,i-1]}^{t+1}, \mathbf{x}_{[i,n]}^t, \mathbf{z}^t; \beta^t) \\
 &= \nabla f_i(\mathbf{x}_i^t) + \mathbf{A}_i^\top \mathbf{z}^t + \beta^t \mathbf{A}_i^\top \{[\sum_{j=1}^{i-1} \mathbf{A}_j \mathbf{x}_j^{t+1}] + [\sum_{j=i}^n \mathbf{A}_j \mathbf{x}_j^t] - \mathbf{b}\} \\
 &= \nabla f_i(\mathbf{x}_i^t) + \mathbf{A}_i^\top \mathbf{z}^t + \beta^t \mathbf{A}_i^\top \{\mathbf{A}\mathbf{x}^{t+1} - \mathbf{b} + [\sum_{j=i}^n \mathbf{A}_j (\mathbf{x}_j^t - \mathbf{x}_j^{t+1})]\} \\
 &\stackrel{\textcircled{1}}{=} \nabla f_i(\mathbf{x}_i^t) + \mathbf{A}_i^\top \mathbf{z}^t + \frac{1}{\sigma} \mathbf{A}_i^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) + \beta^t \mathbf{A}_i^\top \{\sum_{j=i}^n \mathbf{A}_j (\mathbf{x}_j^t - \mathbf{x}_j^{t+1})\},
 \end{aligned} \tag{32}$$

where step ① uses the update rule of \mathbf{z}^{t+1} that $\mathbf{z}^{t+1} - \mathbf{z}^t = \sigma\beta^t(\sum_{i=1}^n \mathbf{A}_i \mathbf{x}_i^{t+1} - \mathbf{b})$. Combining the Equalities (31) and (32), we obtain the following result:

$$\begin{aligned} \mathbf{0} &\in \partial h_i(\mathbf{x}_i^{t+1}) + \boldsymbol{\theta}_i \mathbf{L}_i^t[\mathbf{x}_i^{t+1} - \mathbf{x}_i^t] + \nabla f_i(\mathbf{x}_i^t) \\ &\quad + \mathbf{A}_i^\top \mathbf{z}^t + \beta^t \mathbf{A}_i^\top [\sum_{j=i}^n \mathbf{A}_j(\mathbf{x}_j^t - \mathbf{x}_j^{t+1})] + \frac{1}{\sigma} \mathbf{A}_i^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) \end{aligned}$$

Using the definition of \mathbf{w}_i^{t+1} and \mathbf{u}_i^{t+1} for all $i \in [n]$, we have: $\mathbf{0} = \mathbf{w}_i^{t+1} + \mathbf{u}_i^{t+1} + \mathbf{A}_i^\top \mathbf{z}^t + \frac{1}{\sigma} \mathbf{A}_i^\top (\mathbf{z}^{t+1} - \mathbf{z}^t)$. Multiplying both sides by $\sigma \in (0, 2)$, for all $t \geq 0$, we have:

$$\mathbf{0} = \sigma \mathbf{w}_i^{t+1} + \sigma \mathbf{A}_i^\top \mathbf{z}^t + \mathbf{A}_i^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) + \sigma \mathbf{u}_i^{t+1}. \quad (33)$$

Given that t can take on any integer value, for all $t \geq 1$, we derive:

$$\mathbf{0} = \sigma \mathbf{w}_i^t + \sigma \mathbf{A}_i^\top \mathbf{z}^{t-1} + \mathbf{A}_i^\top (\mathbf{z}^t - \mathbf{z}^{t-1}) + \sigma \mathbf{u}_i^t. \quad (34)$$

Combining Equality (33) and Equality (34), for all $t \geq 1$, we have:

$$\mathbf{A}_i^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) = (1 - \sigma) \mathbf{A}_i^\top (\mathbf{z}^t - \mathbf{z}^{t-1}) - \sigma (\mathbf{w}_i^{t+1} - \mathbf{w}_i^t) - \sigma (\mathbf{u}_i^{t+1} - \mathbf{u}_i^t) \quad (35)$$

In view of (35), we let $i = n$ and arrive at the following two distinct identities:

$$\begin{aligned} \mathbb{B}\mathbb{I} : \underbrace{\mathbf{A}_n^\top (\mathbf{z}^{t+1} - \mathbf{z}^t)}_{\triangleq \mathbf{a}^{t+1}} &= (1 - \sigma) \underbrace{(\mathbf{A}_n^\top (\mathbf{z}^t - \mathbf{z}^{t-1}))}_{\triangleq \mathbf{a}^t} + \sigma \underbrace{(\mathbf{u}_n^t - \mathbf{u}_n^{t+1} + \mathbf{w}_n^t - \mathbf{w}_n^{t+1})}_{\triangleq \mathbf{c}^t}. \\ \mathbb{S}\mathbb{U} : \underbrace{\mathbf{A}_n^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) + \sigma \mathbf{u}_n^{t+1}}_{\triangleq \mathbf{a}^{t+1}} &= (1 - \sigma) \underbrace{(\mathbf{A}_n^\top (\mathbf{z}^t - \mathbf{z}^{t-1}) + \sigma \mathbf{u}_n^t)}_{\triangleq \mathbf{a}^t} + \sigma \underbrace{(\sigma \mathbf{u}_n^t + \mathbf{w}_n^t - \mathbf{w}_n^{t+1})}_{\triangleq \mathbf{c}^t}. \end{aligned}$$

□

D.3 PROOF OF LEMMA 4.3

Proof. We denote $\mathbf{Q}^t \triangleq \theta_2 \mathbf{L}_n^t \mathbf{I} - \beta^t \mathbf{A}_n^\top \mathbf{A}_n \in \mathbb{R}^{\mathbf{d}_i \times \mathbf{d}_i}$.

We assume $\mathbf{A}_n^\top \mathbf{A}_n$ has the singular value decomposition $\mathbf{A}_n^\top \mathbf{A}_n = \tilde{\mathbf{U}}^\top \text{diag}(\boldsymbol{\lambda}) \tilde{\mathbf{U}}$, where $\tilde{\mathbf{U}} \in \mathbb{R}^{\mathbf{d}_i \times \mathbf{d}_i}$, $\boldsymbol{\lambda} \in \mathbb{R}^{\mathbf{d}_i \times 1}$, and $\tilde{\mathbf{U}}^\top \tilde{\mathbf{U}} = \tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top = \mathbf{I}_{\mathbf{d}_i}$. Here, $\text{diag}(\boldsymbol{\lambda})$ denotes a diagonal matrix with $\boldsymbol{\lambda}$ as the main diagonal entries.

(a) We derive:

$$\mathbf{L}_n^t \triangleq L_n + \beta^t \bar{\lambda} \stackrel{\textcircled{1}}{\leq} \beta^t \bar{\lambda}(\delta + 1), \quad (36)$$

where step ① uses Lemma 3.1 that $L_n \leq \delta \beta^t \bar{\lambda}$.

(b) We have:

$$\|\mathbf{Q}^t\| \stackrel{\textcircled{1}}{=} \|\theta_2 \mathbf{L}_n^t - \beta^t \boldsymbol{\lambda}\|_\infty \stackrel{\textcircled{2}}{=} \theta_2 \mathbf{L}_n^t - \min(\beta^t \boldsymbol{\lambda}) \stackrel{\textcircled{3}}{\leq} \bar{\lambda} \beta^t \cdot \underbrace{(\theta_2(1 + \delta) - \lambda'/\bar{\lambda})}_{\triangleq q},$$

where step ① uses $\|\theta_2 \mathbf{L}_n^t \mathbf{I} - \beta^t \mathbf{A}_n^\top \mathbf{A}_n\| = \|\tilde{\mathbf{U}}^\top \text{diag}(\theta_2 \mathbf{L}_n^t - \beta^t \boldsymbol{\lambda}) \tilde{\mathbf{U}}\| = \|\theta_2 \mathbf{L}_n^t - \beta^t \boldsymbol{\lambda}\|_\infty$; step ② uses the fact that $\|\rho - \mathbf{x}\|_\infty = \max(\rho - \mathbf{x}) = \rho - \min(\mathbf{x})$ whenever $\rho \geq \max(\mathbf{x})$ for all ρ and \mathbf{x} ; step ③ uses Inequality (36).

(c) Given $\mathbf{u}_n^{t+1} \triangleq \mathbf{Q}^t(\mathbf{x}_n^{t+1} - \mathbf{x}_n^t)$ as presented in Lemma 4.2, we have: $\|\mathbf{u}_n^{t+1}\| \leq \|\mathbf{Q}^t\| \cdot \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\| \leq q \bar{\lambda} \beta^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|$.

□

D.4 PROOF OF LEMMA 4.4

Proof. For any $\sigma \in [1, 2)$, we define $\sigma_1 \triangleq \frac{\sigma}{(1 - |\sigma|)^2}$, and $\sigma_2 \triangleq \frac{|1 - \sigma|}{\sigma(1 - |\sigma|)}$.

We define $\mathbf{w}_n^{t+1} = \nabla h_n(\mathbf{x}_n^{t+1}; \mu^t) + \nabla f_n(\mathbf{x}_n^t)$.

We define $\mathbf{a}^{t+1} \triangleq \mathbf{A}_n^\top(\mathbf{z}^{t+1} - \mathbf{z}^t)$, and $\mathbf{c}^t \triangleq \mathbf{u}_n^t - \mathbf{u}_n^{t+1} + \mathbf{w}_n^t - \mathbf{w}_n^{t+1}$.

We define $\Theta_a^t \triangleq \frac{K_a}{\beta^t} \|\mathbf{a}^t\|_2^2$, where $K_a = \frac{\omega\sigma_2}{\lambda}$.

We define $\Theta_u^t \triangleq \frac{K_u}{\beta^t} (L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \|\mathbf{u}_n^t\|)^2$, where $K_u = \frac{3\omega\sigma_1}{\lambda}$.

We define $\Gamma_\mu^t \triangleq \frac{C_h^2 K_u}{\beta^t} \cdot (\frac{\mu^{t-1}}{\mu^t} - 1)^2$.

First, we bound the term $\|\mathbf{c}^t\|$. For all $t \geq 1$, we have:

$$\begin{aligned} \|\mathbf{c}^t\| &= \|\mathbf{w}_n^t - \mathbf{w}_n^{t+1} + \mathbf{u}_n^t - \mathbf{u}_n^{t+1}\| \\ &\stackrel{(1)}{\leq} \|\nabla h_n(\mathbf{x}_n^{t+1}; \mu^t) - \nabla h_n(\mathbf{x}_n^t; \mu^{t-1})\| + \|\nabla f_n(\mathbf{x}_n^t) - \nabla f_n(\mathbf{x}_n^{t-1})\| + \|\mathbf{u}_n^t - \mathbf{u}_n^{t+1}\| \\ &\stackrel{(2)}{\leq} \|\nabla h_n(\mathbf{x}_n^{t+1}; \mu^t) - \nabla h_n(\mathbf{x}_n^t; \mu^{t-1})\| + L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \|\mathbf{u}_n^t - \mathbf{u}_n^{t+1}\| \\ &= \|\nabla h_n(\mathbf{x}_n^{t+1}; \mu^t) - \nabla h_n(\mathbf{x}_n^t; \mu^t) + \nabla h_n(\mathbf{x}_n^t; \mu^t) - \nabla h_n(\mathbf{x}_n^t; \mu^{t-1})\| \\ &\quad + L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \|\mathbf{u}_n^t - \mathbf{u}_n^{t+1}\| \\ &\stackrel{(3)}{\leq} \frac{1}{\mu^t} \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\| + (\frac{\mu^{t-1}}{\mu^t} - 1) C_h + L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \|\mathbf{u}_n^t\| + \|\mathbf{u}_n^{t+1}\|, \end{aligned} \quad (37)$$

where step ① uses the triangle inequality; step ② uses the fact that $f_n(\mathbf{x})$ is L_n -smooth; step ③ uses Lemma 3.5 and Lemma 3.3.

Second, we bound the term $\frac{\omega\sigma_1}{\lambda\beta^t} \|\mathbf{c}^t\|_2^2$. For all $t \geq 1$, we have:

$$\begin{aligned} \frac{\omega\sigma_1}{\lambda\beta^t} \|\mathbf{c}^t\|_2^2 &\stackrel{(1)}{\leq} \frac{3\omega\sigma_1}{\lambda\beta^t} \left(\frac{1}{\mu^t} \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\| + \|\mathbf{u}_n^{t+1}\| \right)^2 + \underbrace{\frac{3\omega\sigma_1}{\lambda\beta^t} C_h^2 \left(\frac{\mu^{t-1}}{\mu^t} - 1 \right)^2}_{\triangleq \Gamma_\mu^t} + \underbrace{\frac{3\omega\sigma_1}{\lambda\beta^t} (L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \|\mathbf{u}_n^t\|)^2}_{\triangleq \Theta_u^t} \\ &\stackrel{(2)}{=} \frac{3\omega\sigma_1}{\lambda\beta^t} \left\{ \left(\frac{1}{\mu^t} \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\| + \|\mathbf{u}_n^{t+1}\| \right)^2 + (L_n \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\| + \|\mathbf{u}_n^{t+1}\|)^2 \right\} + \Gamma_\mu^t + \Theta_u^t - \Theta_u^{t+1} \\ &\stackrel{(3)}{\leq} \frac{3\omega\sigma_1}{\lambda\beta^t} \cdot 2((\delta + q)\bar{\lambda}\beta^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|)^2 + \Gamma_\mu^t + \Theta_u^t - \Theta_u^{t+1} \\ &= \underbrace{6\omega\sigma_1\kappa(\delta + q)^2 \cdot \bar{\lambda}\beta^t \cdot \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2}_{\triangleq \chi_1} + \Gamma_\mu^t + \Theta_u^t - \Theta_u^{t+1} \\ &\stackrel{(4)}{\leq} \chi_1 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \Gamma_\mu^t + \Theta_u^t - \Theta_u^{t+1}, \end{aligned} \quad (38)$$

where step ① uses Inequality 41 and the fact that $(a + b + c)^2 \leq 3a^2 + 3b^2 + 3c^2$ for all $a \in \mathbb{R}$, $b \in \mathbb{R}$, and $c \in \mathbb{R}$; step ② uses the definitions of $\{K_u, \Theta_u^t, \Gamma_\mu^t\}$; step ③ uses Lemma 4.3 that: $\frac{1}{\mu^t} \leq \delta\bar{\lambda}\beta^t$, $L_n \leq \delta\bar{\lambda}\beta^t$, and $\|\mathbf{u}_n^{t+1}\| \leq \|\mathbf{Q}^t\| \cdot \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\| \leq q\bar{\lambda}\beta^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|$; step ④ uses $\beta^t \bar{\lambda} \leq \mathsf{L}_n^t \triangleq \beta^t \bar{\lambda} + L_n$.

Finally, we derive the following inequalities for all $t \geq 1$:

$$\begin{aligned} \frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 &\stackrel{(1)}{\leq} \frac{\omega}{\lambda\sigma\beta^t} \|\mathbf{A}_n^\top(\mathbf{z}^{t+1} - \mathbf{z}^t)\|_2^2 = \frac{\omega}{\sigma\lambda\beta^t} \|\mathbf{a}^t\|_2^2 \\ &\stackrel{(2)}{\leq} \frac{\sigma_2\omega}{\lambda} \left(\frac{1}{\beta^t} \|\mathbf{a}^t\|_2^2 - \frac{1}{\beta^t} \|\mathbf{a}^{t+1}\|_2^2 \right) + \frac{\omega\sigma_1}{\lambda\beta^t} \|\mathbf{c}^t\|_2^2 \\ &\stackrel{(3)}{\leq} \underbrace{\frac{\sigma_2\omega}{\lambda} \cdot \frac{1}{\beta^t} \|\mathbf{a}^t\|_2^2}_{\triangleq \Theta_a^t} - \frac{\sigma_2\omega}{\lambda} \cdot \frac{1}{\beta^{t+1}} \|\mathbf{a}^{t+1}\|_2^2 + \frac{\omega\sigma_1}{\lambda} \cdot \frac{1}{\beta^t} \|\mathbf{c}^t\|_2^2 \\ &\stackrel{(4)}{\leq} \Theta_a^t - \Theta_a^{t+1} + \chi_1 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \Gamma_\mu^t + \Theta_u^t - \Theta_u^{t+1}, \end{aligned}$$

where step ① uses $\underline{\lambda} \|\mathbf{z}\|_2^2 \leq \|\mathbf{A}_n^\top \mathbf{z}\|_2^2$ for all \mathbf{z} ; step ② uses Lemma A.2 with $\mathbf{b} = \mathbf{a}^t$, $\mathbf{b}^+ = \mathbf{a}^{t+1}$, and $\mathbf{a} = \mathbf{c}^t$ that:

$$\frac{1}{\sigma\beta^t} \|\mathbf{a}^{t+1}\|_2^2 \leq \frac{\sigma_2}{\beta^t} (\|\mathbf{a}^t\|_2^2 - \|\mathbf{a}^{t+1}\|_2^2) + \frac{\sigma_1}{\beta^t} \|\mathbf{c}^t\|_2^2;$$

1296 step ③ uses $-\frac{1}{\beta^t} \leq -\frac{1}{\beta^{t+1}}$; step ④ uses Inequality (38).
 1297

□

1300 D.5 PROOF OF LEMMA 4.5

1301 *Proof.* (a) With the choice $\theta_1 = 1.01$, it clearly holds that $\varepsilon_1 \triangleq \frac{1}{2}\theta_1 - \frac{1}{2} > 0$.
 1302

1303 (b) We define $\chi_1 \triangleq \chi_0(\delta + \theta_2 + \theta_2\delta - 1/\kappa)^2$, where $\chi_0 \triangleq 6\omega\sigma_1\kappa$.
 1304

1305 With the choice $\theta_2 = \frac{1}{2\chi_0(1+\delta)^2} + \frac{1/\kappa-\delta}{1+\delta}$, we now prove that $\varepsilon_2 \triangleq \theta_2 - \frac{1}{2} - \chi_1 > 0$.
 1306

1307 We consider the following concave auxiliary function
 1308

$$f(\theta_2) \triangleq \theta_2 - \frac{1}{2} - \chi_0(\delta + \theta_2 + \delta\theta_2 - 1/\kappa)^2.$$

1309 Setting the gradient of $f(\theta_2)$ w.r.t. θ_2 yields: $1 - 2\chi_0(\delta + \theta_2 + \delta\theta_2 - 1/\kappa)(1 + \delta) = 0$. It follows
 1310 that the solution $\bar{\theta}_2 = \frac{1}{2(1+\delta)^2\chi_0} + \frac{1/\kappa-\delta}{\delta+1}$ is the maximizer of the concave auxiliary function. We
 1311 have:
 1312

$$\begin{aligned} f(\bar{\theta}_2) &\stackrel{\textcircled{1}}{=} \bar{\theta}_2 - \frac{1}{2} - \chi_0(\delta + \theta_2 + \delta\theta_2 - 1/\kappa)^2 \\ &= \frac{1}{4(1+\delta)^2\chi_0} + \frac{1/\kappa-\delta}{\delta+1} - \frac{1}{2} \\ &\stackrel{\textcircled{2}}{\geq} \frac{1}{4(1+\delta)^2\chi_0} + 0 \\ &\stackrel{\textcircled{3}}{\geq} \frac{1}{4(1+1/3)^2\chi_0} \\ &\stackrel{\textcircled{4}}{\geq} \frac{1}{8\chi_0}, \end{aligned}$$

1322 where step ① uses the definitions of $f(\theta_2)$ and $\bar{\theta}_2$; step ② uses the following derivations: $(\delta \leq \frac{2/\kappa-1}{3}) \Rightarrow (2/\kappa-1 \geq 3\delta) \Rightarrow (2/\kappa-2\delta \geq 1+\delta) \Rightarrow (\frac{1/\kappa-\delta}{1+\delta} \geq \frac{1}{2})$; step ③ uses the fact that $\delta \leq \frac{1}{3}$;
 1323 step ④ uses $4 \times (1+1/3)^2 < 8$.
 1324

□

1328 D.6 PROOF OF LEMMA 4.6

1329 *Proof.* We define $\mathcal{E}^{t+1} \triangleq [\varepsilon_1 \sum_{i=1}^{n-1} \mathsf{L}_i^t \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2] + \varepsilon_2 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \frac{\varepsilon_3}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2$.
 1330

1331 We define $\Theta^t \triangleq \Theta_L^t + \Theta_{au}^t$, where $\Theta_{au}^t \triangleq \Theta_a^t + \Theta_u^t$.
 1332

1333 Using the results from Lemma 4.1 and Lemma 4.4, we derive the following two respective inequalities:
 1334

$$\mathcal{E}^{t+1} + \Theta_L^{t+1} - \Theta_L^t \leq (\frac{1}{2} - \theta_2 + \varepsilon_2) \cdot \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \quad (39)$$

$$\frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \leq \Theta_{au}^t - \Theta_{au}^{t+1} + \chi_1 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \Gamma_\mu^t. \quad (40)$$

1339 Adding Inequalities (39) and (40) together, we have:
 1340

$$\mathcal{E}^{t+1} + \Theta^{t+1} - \Theta^t - \Gamma_\mu^t \leq \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 \cdot \{\frac{1}{2} - \theta_2 + \varepsilon_2 + \chi_1\} \stackrel{\textcircled{1}}{=} 0,$$

1343 where step ① uses the definition of $\varepsilon_2 \triangleq \theta_2 - \frac{1}{2} - \chi_1$ as in Lemma (4.5).
 1344

□

1346 D.7 PROOF OF LEMMA 4.7

1348 *Proof.* For any $\sigma \in (0, 1)$, we define $\sigma_1 \triangleq \frac{\sigma}{(1-|1-\sigma|)^2}$, and $\sigma_2 \triangleq \frac{|1-\sigma|}{\sigma(1-|1-\sigma|)}$.
 1349

We define $\mathbf{w}_n^{t+1} = \nabla h_n(\mathbf{x}_n^{t+1}; \mu^t) + \nabla f_n(\mathbf{x}_n^t)$.

We define $\mathbf{a}^{t+1} \triangleq \mathbf{A}_n^T(\mathbf{z}^{t+1} - \mathbf{z}^t) + \sigma \mathbf{u}_n^t$, and $\mathbf{c}^t \triangleq \sigma \mathbf{u}_n^t + \mathbf{w}_n^t - \mathbf{w}_n^{t+1}$.

We define $\Theta_a^t \triangleq \frac{K_a}{\beta^t} \|\mathbf{a}^t\|_2^2$, where $K_a \triangleq \frac{2\omega\sigma_2}{\lambda}$.

We define $\Theta_u^t \triangleq \frac{K_u}{\beta^t} (L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \sigma \|\mathbf{u}_n^t\|)^2$, where $K_u = \frac{6\omega\sigma_1}{\lambda}$.

We define $\Gamma_\mu^t \triangleq \frac{C_h^2 6\omega\sigma_1}{\lambda\beta^t} \cdot (\frac{\mu^{t-1}}{\mu^t} - 1)^2$.

First, we bound the term $\|\mathbf{c}^t\|$. For all $t \geq 1$, we have:

$$\begin{aligned} \|\mathbf{c}^t\| &= \|\mathbf{w}_n^t - \mathbf{w}_n^{t+1} + \sigma \mathbf{u}_n^t\| \\ &\stackrel{(1)}{\leq} \|\nabla h_n(\mathbf{x}_n^{t+1}; \mu^t) - \nabla h_n(\mathbf{x}_n^t; \mu^{t-1})\| + \|\nabla f_n(\mathbf{x}_n^t) - \nabla f_n(\mathbf{x}_n^{t-1})\| + \sigma \|\mathbf{u}_n^t\| \\ &\stackrel{(2)}{\leq} \|\nabla h_n(\mathbf{x}_n^{t+1}; \mu^t) - \nabla h_n(\mathbf{x}_n^t; \mu^{t-1})\| + L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \sigma \|\mathbf{u}_n^t\| \\ &= \|\nabla h_n(\mathbf{x}_n^{t+1}; \mu^t) - \nabla h_n(\mathbf{x}_n^t; \mu^t) + \nabla h_n(\mathbf{x}_n^t; \mu^t) - \nabla h_n(\mathbf{x}_n^t; \mu^{t-1})\| + L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \sigma \|\mathbf{u}_n^t\| \\ &\stackrel{(3)}{\leq} \frac{1}{\mu^t} \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\| + (\frac{\mu^{t-1}}{\mu^t} - 1) C_h + L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \sigma \|\mathbf{u}_n^t\|, \end{aligned} \quad (41)$$

where step ① uses the triangle inequality; step ② uses the fact that $f_n(\mathbf{x})$ is L_n -smooth; step ③ uses Lemma 3.3 and Lemma 3.5.

Second, we bound the term $\frac{2\omega\sigma}{\lambda\beta^t} \|\mathbf{u}_n^t\|_2^2 + \frac{2\omega}{\sigma\lambda\beta^t} \|\mathbf{c}^t\|_2^2$. For all $t \geq 1$, we have:

$$\begin{aligned} &\frac{2\omega\sigma}{\lambda\beta^t} \|\mathbf{u}_n^{t+1}\|_2^2 + \frac{2\omega\sigma_1}{\lambda\beta^t} \|\mathbf{c}^t\|_2^2 \\ &\stackrel{(1)}{\leq} \frac{2\omega\sigma}{\lambda\beta^t} \|\mathbf{u}_n^{t+1}\|_2^2 + \frac{6\omega\sigma_1}{\lambda\beta^t} (\frac{1}{\mu^t} \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|)^2 + \underbrace{\frac{6\omega\sigma_1}{\lambda\beta^t} (\frac{\mu^{t-1}}{\mu^t} - 1)^2 C_h^2}_{\Gamma_\mu^t} + \underbrace{\frac{6\omega\sigma_1}{\lambda\beta^t} (L_n \|\mathbf{x}_n^t - \mathbf{x}_n^{t-1}\| + \sigma \|\mathbf{u}_n^t\|)^2}_{\triangleq \Theta_u^t} \\ &\stackrel{(2)}{=} \frac{2\omega\sigma_1}{\beta^t \lambda} \cdot \{ \frac{\sigma}{\sigma_1} \|\mathbf{u}_n^{t+1}\|_2^2 + 3(\frac{1}{\mu^t} \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|)^2 + 3(L_n \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\| + \sigma \|\mathbf{u}_n^{t+1}\|)^2 \} + \Gamma_\mu^t + \Theta_u^t - \Theta_u^{t+1} \\ &\stackrel{(3)}{\leq} \frac{2\omega\sigma_1}{\beta^t \lambda} \cdot \bar{\lambda}^2 (\beta^t)^2 \cdot \{ \frac{\sigma}{\sigma_1} q^2 + 3\delta^2 + 3(\delta + \sigma q)^2 \} \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \Gamma_\mu^t + \Theta_u^t - \Theta_u^{t+1} \\ &\stackrel{(4)}{\leq} \underbrace{\frac{2\omega\kappa}{\sigma} \cdot \{ \sigma^2 q^2 + 3\delta^2 + 3(\delta + \sigma q)^2 \} \cdot \bar{\lambda} \beta^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2}_{\triangleq \chi_2} + \Gamma_\mu^t + \Theta_u^t - \Theta_u^{t+1} \\ &\stackrel{(5)}{\leq} \chi_2 \cdot \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \Gamma_\mu^t + \Theta_u^t - \Theta_u^{t+1}, \end{aligned} \quad (42)$$

where step ① uses Inequality 41 and the fact that $(a+b+c)^2 \leq 3a^2 + 3b^2 + 3c^2$ for all $a \in \mathbb{R}$, $b \in \mathbb{R}$, and $c \in \mathbb{R}$; step ② uses the definitions of $\{K_u, \Theta_u^t, \Gamma_\mu^t\}$; step ③ uses $\|\mathbf{u}_n^{t+1}\| \leq \|\mathbf{Q}^t\| \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\| \leq \beta^t \bar{\lambda} q \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|$ and $L_n \leq \bar{\lambda} \beta^t \delta$, as has been shown respectively in Lemma 4.3 and Lemma 3.1, as well as the fact that $\frac{1}{\mu^t} = \beta^t \bar{\lambda} \delta$; step ④ uses $\kappa = \bar{\lambda}/\lambda$, and the fact that $\sigma_1 = \frac{1}{\sigma}$ when $\sigma \in (0, 1)$; step ⑤ uses $\beta^t \bar{\lambda} \leq \mathsf{L}_n^t \triangleq \beta^t \bar{\lambda} + L_n$.

Finally, for all $t \geq 1$, we derive:

$$\begin{aligned} &\frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \\ &\stackrel{(1)}{\leq} \frac{\omega}{\sigma\beta^t \cdot \lambda} \|\mathbf{A}_n^T(\mathbf{z}^{t+1} - \mathbf{z}^t)\|_2^2 \\ &\stackrel{(2)}{=} \frac{\omega}{\lambda} \cdot \frac{1}{\sigma\beta^t} \|\mathbf{a}^{t+1} - \sigma \mathbf{u}_n^{t+1}\|_2^2 \\ &\stackrel{(3)}{\leq} \frac{2\omega}{\lambda} \cdot \{ \frac{1}{\sigma\beta^t} \|\mathbf{a}^{t+1}\|_2^2 + \frac{\sigma}{\beta^t} \|\mathbf{u}_n^{t+1}\|_2^2 \} \\ &\stackrel{(4)}{\leq} \frac{2\omega}{\lambda} \cdot \{ \frac{\sigma_2}{\beta^t} \|\mathbf{a}^t\|_2^2 - \frac{\sigma_2}{\beta^{t+1}} \|\mathbf{a}^{t+1}\|_2^2 + \frac{\sigma_1}{\beta^t} \|\mathbf{c}^t\|_2^2 \} + \frac{2\omega\sigma}{\beta^t \lambda} \|\mathbf{u}_n^{t+1}\|_2^2 \\ &\stackrel{(5)}{\leq} \underbrace{\frac{2\omega}{\lambda} \frac{\sigma_2}{\beta^t} \|\mathbf{a}^t\|_2^2}_{\triangleq \Theta_a^t} - \frac{2\omega}{\lambda} \frac{\sigma_2}{\beta^{t+1}} \|\mathbf{a}^{t+1}\|_2^2 + \frac{2\omega}{\lambda} \frac{\sigma_1}{\beta^t} \|\mathbf{c}^t\|_2^2 + \frac{2\omega\sigma}{\beta^t \lambda} \|\mathbf{u}_n^{t+1}\|_2^2 \\ &\stackrel{(6)}{\leq} \Theta_a^t - \Theta_a^{t+1} + \chi_2 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \Gamma_\mu^t + \Theta_u^t - \Theta_u^{t+1}, \end{aligned}$$

where step ① uses the fact that $\lambda\|\mathbf{x}\|_2^2 \leq \|\mathbf{A}_n^\top \mathbf{x}\|_2^2$ for all \mathbf{x} ; step ② uses the definition of \mathbf{a}^{t+1} ; step ③ uses the inequality $\|\mathbf{a} + \mathbf{b}\|_2^2 \leq 2\|\mathbf{a}\|_2^2 + 2\|\mathbf{b}\|_2^2$ for all \mathbf{a} and \mathbf{b} ; step ④ uses Lemma A.2 with $\mathbf{b} = \mathbf{a}^t$, $\mathbf{b}^+ = \mathbf{a}^{t+1}$, and $\mathbf{a} = \mathbf{c}^t$ that

$$\frac{1}{\sigma\beta^t}\|\mathbf{a}^{t+1}\|_2^2 \leq \frac{\sigma_1}{\beta^t}\|\mathbf{c}^t\|_2^2 + \frac{\sigma_2}{\beta^t}(\|\mathbf{a}^t\|_2^2 - \|\mathbf{a}^{t+1}\|_2^2);$$

step ⑤ uses $-\frac{1}{\beta^t} \leq -\frac{1}{\beta^{t+1}}$ and $\sigma_1 = \frac{1}{\sigma}$ when $\sigma \in (0, 1)$; step ⑥ uses Inequality (42). \square

D.8 PROOF OF LEMMA 4.8

Proof. We assume $\xi = \delta = \sigma = \frac{c}{\kappa}$, where $c \in (0, 1)$.

We have:

$$\omega \triangleq 1 + \frac{\xi}{\sigma} = 2 \quad (43)$$

$$q \triangleq \theta_2 + \theta_2\delta \stackrel{\textcircled{1}}{\leq} \theta_2 + \theta_2 c. \quad (44)$$

where step ① uses $\delta = c/\kappa \leq c$ since $\kappa \geq 1$. We further obtain:

$$\begin{aligned} \varepsilon_2 &\triangleq \theta_2 - \frac{1}{2} - \frac{6\omega\kappa}{\sigma}\left\{\frac{1}{3}\sigma^2q^2 + (\delta + \sigma q)^2 + \delta^2\right\} \\ &\stackrel{\textcircled{1}}{\geq} \theta_2 - \frac{1}{2} - \frac{12}{c}\left\{\frac{1}{3}c^2q^2 + (c + cq)^2 + c^2\right\} \\ &= \theta_2 - \frac{1}{2} - 12c\left\{\frac{1}{3}q^2 + (1 + q)^2 + 1\right\} \\ &\stackrel{\textcircled{2}}{\geq} \theta_2 - \frac{1}{2} - 12c\left\{\frac{(\theta_2 + \theta_2 c)^2}{3} + (1 + \theta_2 + \theta_2 c)^2 + 1\right\} \\ &\stackrel{\textcircled{3}}{>} 0.02, \end{aligned}$$

where step ① uses (43), $\sigma \leq c$, $\delta \leq c$; step ② uses (44); step ③ uses the choice $c = 0.01$ and $\theta_2 = 1.5$. \square

D.9 PROOF OF LEMMA 4.9

Proof. We define $\mathcal{E}^{t+1} \triangleq [\varepsilon_1 \sum_{i=1}^{n-1} \mathsf{L}_i^t \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2] + \varepsilon_2 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \frac{\varepsilon_3}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2$.

We define $\Theta^t \triangleq \Theta_L^t + \Theta_{au}^t$, where $\Theta_{au}^t \triangleq \Theta_a^t + \Theta_u^t$.

Using the results from Lemma 4.1 and Lemma 4.7, we derive the following two respective inequalities:

$$\mathcal{E}^{t+1} + \Theta_L^{t+1} - \Theta_L^t \leq \left(\frac{1}{2} - \theta_2 + \varepsilon_2\right) \cdot \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2,$$

$$\frac{\omega}{\sigma\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + \Theta_{au}^{t+1} - \Theta_{au}^t \leq \chi_2 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \Gamma_\mu^t.$$

Adding the two inequalities above together leads to:

$$\mathcal{E}^{t+1} + \Theta^{t+1} - \Theta^t - \Gamma_\mu^t \leq \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 \cdot \left\{\frac{1}{2} - \theta_2 + \varepsilon_2 + \chi_2\right\} \stackrel{\textcircled{1}}{=} 0,$$

where step ① uses the definition of $\varepsilon_2 \triangleq \theta_2 - \frac{1}{2} - \chi_2$ as in Lemma (4.8). \square

D.10 PROOF OF LEMMA 4.10

Proof. The proof of this lemma closely resembles that of Theorem 6 in (Bōt et al., 2019).

We denote $\underline{\Theta} \triangleq \underline{\Theta}' - \mu^0 C_h^2$, where $\underline{\Theta}'$ is defined in Assumption 1.4

Initially, for all $t \geq 1$, we have:

$$\begin{aligned}
\Theta^t &\stackrel{\textcircled{1}}{=} \mathcal{L}(\mathbf{x}^t, \mathbf{z}^t; \beta^t, \mu^t) + \frac{1}{2}C_h\mu^t + \Theta_a^t + \Theta_u^t \\
&\stackrel{\textcircled{2}}{\geq} \mathcal{L}(\mathbf{x}^t, \mathbf{z}^t; \beta^t, \mu^t) \\
&\stackrel{\textcircled{3}}{=} h_n(\mathbf{x}_n^t; \mu^t) + \{\sum_{i=1}^{n-1} h_i(\mathbf{x}_i^t)\} + \sum_{i=1}^n f_i(\mathbf{x}_i^t) + \langle \mathbf{A}\mathbf{x}^t - \mathbf{b}, \mathbf{z} \rangle + \frac{\beta^t}{2} \|\mathbf{A}\mathbf{x}^t - \mathbf{b}\|_2^2 \\
&\stackrel{\textcircled{4}}{\geq} -\mu^0 C_h^2 + \{\sum_{i=1}^n h_i(\mathbf{x}_i^t)\} + \{\sum_{i=1}^n f_i(\mathbf{x}_i^t)\} + \langle \mathbf{A}\mathbf{x}^t - \mathbf{b}, \mathbf{z} \rangle + \frac{\beta^t}{2} \|\mathbf{A}\mathbf{x}^t - \mathbf{b}\|_2^2 \\
&\stackrel{\textcircled{5}}{\geq} \langle \mathbf{A}\mathbf{x}^t - \mathbf{b}, \mathbf{z}^t \rangle - \mu^0 C_h^2 + \underline{\Theta}' \\
&\stackrel{\textcircled{6}}{\geq} \langle \mathbf{A}\mathbf{x}^t - \mathbf{b}, \mathbf{z}^t \rangle + \underline{\Theta} \tag{45}
\end{aligned}$$

where step ① uses the definition of Θ^t ; step ② uses the nonnegativity of the terms $\{\frac{1}{2}C_h\mu^t, \Theta_a^t, \Theta_u^t\}$; step ③ uses the definition of $\mathcal{L}(\mathbf{x}^t, \mathbf{z}^t; \beta^t, \mu^t)$ in Equation (4); step ④ uses $0 \leq h_n(\mathbf{u}) - h_n(\mathbf{u}; \mu) \leq \mu C_h^2$ as shown in Lemma 3.3, and the fact that $\mu^t \leq \mu^0$; step ⑤ uses Assumption 1.4; step ⑥ uses $\underline{\Theta} \triangleq \underline{\Theta}' - \mu^0 C_h^2$.

We now conclude the proof of this lemma through contradiction. Suppose that there exists $t_0 \geq 1$ such that $\Theta^{t_0} < \underline{\Theta}$. We derive the following inequalities:

$$\begin{aligned}
\sum_{t=1}^T (\Theta^t - \underline{\Theta}) &= [\sum_{t=1}^{t_0-1} (\Theta^t - \underline{\Theta})] + [\sum_{t=t_0}^T (\Theta^t - \underline{\Theta})] \\
&\leq [\sum_{t=1}^{t_0-1} (\Theta^t - \underline{\Theta})] + (T+1-t_0) \cdot \max_{t=t_0}^T (\Theta^t - \underline{\Theta}) \\
&\stackrel{\textcircled{1}}{\leq} [\sum_{t=1}^{t_0-1} (\Theta^t - \underline{\Theta})] + (T+1-t_0) \cdot (\Theta^{t_0} - \underline{\Theta}), \tag{46}
\end{aligned}$$

where step ① uses $\Theta^t \leq \Theta^{t_0}$ for all $t \geq t_0$. We closely examine Inequality (46). As t_0 is finite, the sum $\sum_{t=1}^{t_0-1} (\Theta^t - \underline{\Theta})$ is upper bounded. Considering the negativity of the term $(\Theta^{t_0} - \underline{\Theta})$, we deduce from Inequality (46):

$$\lim_{T \rightarrow \infty} \sum_{t=1}^T (\Theta^t - \underline{\Theta}) = -\infty. \tag{47}$$

Meanwhile, for all $t \geq 1$, the following inequalities hold:

$$\begin{aligned}
\Theta^t - \underline{\Theta} &\stackrel{\textcircled{1}}{\geq} \frac{1}{\sigma\beta^{t-1}} \langle \mathbf{z}^t - \mathbf{z}^{t-1}, \mathbf{z}^t \rangle \\
&\stackrel{\textcircled{2}}{=} \frac{1}{2\sigma} \left\{ \frac{1}{\beta^{t-1}} \|\mathbf{z}^t\|_2^2 - \frac{1}{\beta^{t-1}} \|\mathbf{z}^{t-1}\|_2^2 + \frac{1}{\beta^{t-1}} \|\mathbf{z}^t - \mathbf{z}^{t-1}\|_2^2 \right\} \\
&\stackrel{\textcircled{3}}{\geq} \frac{1}{2\sigma} \left\{ \frac{1}{\beta^t} \|\mathbf{z}^t\|_2^2 - \frac{1}{\beta^{t-1}} \|\mathbf{z}^{t-1}\|_2^2 + 0 \right\}, \tag{48}
\end{aligned}$$

where step ① uses Inequality (45) and $\mathbf{z}^{t+1} = \mathbf{z}^t + \sigma\beta^t(\mathbf{A}\mathbf{x}^{t+1} - \mathbf{b})$; step ② uses the Pythagoras relation in Lemma A.1; step ③ uses $\frac{1}{\beta^{t-1}} \geq \frac{1}{\beta^t}$.

Telescoping Inequality (48) over t from 1 to T , we have:

$$\sum_{t=1}^T (\Theta^t - \underline{\Theta}) \geq \frac{1}{2\sigma} \cdot \left\{ \frac{1}{\beta^T} \|\mathbf{z}^T\|_2^2 - \frac{1}{\beta^0} \|\mathbf{z}^0\|_2^2 \right\} \geq -\frac{1}{2\sigma\beta^0} \|\mathbf{z}^0\|_2^2. \tag{49}$$

The finiteness of the right-hand-side in (49) contradicts with (47).

Therefore, we conclude that $\Theta^t \geq \underline{\Theta}$ for all $t \geq 1$. \square

D.11 PROOF OF LEMMA 4.11

Proof. We define $C_\mu \triangleq \frac{3}{\beta^0} C_h^2 K_u$.

We define $\Gamma_\mu^t \triangleq C_h^2 \frac{K_u}{\beta^t} \cdot (\frac{\mu^{t-1}}{\mu^t} - 1)^2$, where $\beta^t = \beta^0(1 + \xi t^p)$, $\mu^t \propto \frac{1}{\beta^t}$.

Letting $T \in [1, \infty)$, we obtain:

$$\begin{aligned}
\sum_{t=1}^T (\frac{\mu^{t-1}}{\mu^t} - 1)^2 &\stackrel{\textcircled{1}}{=} \sum_{t=1}^T (\frac{\beta^t}{\beta^{t-1}} - 1)^2 = (\frac{\beta^1}{\beta^0} - 1)^2 + \sum_{t=2}^T (\frac{\beta^t}{\beta^{t-1}} - 1)^2 \\
&\stackrel{\textcircled{2}}{=} (1 + \xi 1^p - 1)^2 + \sum_{t=1}^{T-1} (\frac{\beta^{t+1}}{\beta^t} - 1)^2 \\
&\stackrel{\textcircled{3}}{\leq} 1 + \sum_{t=1}^{\infty} \frac{(\xi(t+1)^p - \xi t^p)^2}{(1+\xi t^p)^2} \\
&\stackrel{\textcircled{4}}{\leq} 1 + \sum_{t=1}^{\infty} (\frac{(t+1)^p - t^p}{t^p})^2 \\
&\stackrel{\textcircled{5}}{\leq} 1 + 2,
\end{aligned} \tag{50}$$

where step $\textcircled{1}$ uses $\mu^t \propto \frac{1}{\beta^t}$; step $\textcircled{2}$ uses $\beta^1 = \beta^0(1 + \xi 1^p)$; step $\textcircled{3}$ uses the definition of $\beta^t = \beta^0 + \beta^0 \xi t^p$; step $\textcircled{4}$ uses $\frac{1}{(1+\xi t^p)^2} \leq \frac{1}{(\xi t^p)^2}$; step $\textcircled{5}$ uses Lemma A.5.

We further obtain:

$$\sum_{t=1}^{\infty} \Gamma_{\mu}^t \stackrel{\textcircled{1}}{\leq} C_h^2 \frac{K_u}{\beta^0} \cdot \{\sum_{t=1}^{\infty} (\frac{\mu^{t-1}}{\mu^t} - 1)^2\} \stackrel{\textcircled{2}}{\leq} 3C_h^2 \frac{K_u}{\beta^0} \triangleq C_{\mu},$$

where step $\textcircled{1}$ uses $\beta^t \geq \beta^0$; step $\textcircled{2}$ uses Inequality (50).

□

D.12 PROOF OF THEOREM 4.12

For both conditions \mathbb{BI} and \mathbb{SU} , we have from Lemmas (4.6) and (4.9):

$$\mathcal{E}^{t+1} \leq \Theta^t - \Theta^{t+1} + \Gamma_{\mu}^t.$$

Telescoping this inequality over t from 1 to T , we have:

$$\sum_{t=1}^T \mathcal{E}^{t+1} \leq \Theta^1 - \Theta^{T+1} + \sum_{t=1}^T \Gamma_{\mu}^t \stackrel{\textcircled{1}}{\leq} \Theta^1 - \underline{\Theta} + C_{\mu} \triangleq K_e, \tag{51}$$

where step $\textcircled{1}$ uses Lemma 4.10 that $\Theta^t \geq \underline{\Theta}$ for all t , and Lemma 4.11.

D.13 PROOF OF LEMMA 4.13

Proof. Given $\sigma \in (0, 2)$, we define $\sigma_3 \triangleq \frac{\sigma}{1-|1-\sigma|} \in [1, \infty)$.

We define $\mathbf{w}_n^{t+1} \triangleq \nabla h_n(\mathbf{x}_n^{t+1}, \mu^t) + \nabla f_n(\mathbf{x}_n^t)$.

We define $\mathbf{u}_n^{t+1} \triangleq \mathbf{Q}^t(\mathbf{x}_n^{t+1} - \mathbf{x}_n^t)$, where $\mathbf{Q}^t \triangleq \theta_2 \mathbf{L}_n^t \mathbf{I} - \beta^t \mathbf{A}_n^T \mathbf{A}_n$.

We define $K_z \triangleq \frac{3}{\lambda} (\frac{1}{\beta^0} \bar{\lambda} \|\mathbf{z}^1\|_2^2 + 2\sigma_3 C_h^2 + 2\sigma_3 C_f^2 + \sigma_3 q^2 \bar{\lambda} \frac{K_e}{\varepsilon_2})$, and $\ddot{K}_z \triangleq K_e / \varepsilon_3$.

First, we have:

$$\begin{aligned}
\max_{i=1}^{\infty} \{\|\mathbf{w}_n^{i+1}\|_2^2\} &= \max_{i=1}^{\infty} \{\|\nabla h_n(\mathbf{x}_n^{i+1}, \mu^i) + \nabla f_n(\mathbf{x}_n^i)\|_2^2\} \\
&\stackrel{\textcircled{1}}{\leq} 2 \max_{i=1}^{\infty} \{\|\nabla h_n(\mathbf{x}_n^{i+1}, \mu^i)\|_2^2 + \|\nabla f_n(\mathbf{x}_n^i)\|_2^2\} \\
&\stackrel{\textcircled{2}}{\leq} 2C_h^2 + 2C_f^2
\end{aligned} \tag{52}$$

where step $\textcircled{1}$ uses $\|\mathbf{a} + \mathbf{b}\|_2^2 \leq 2\|\mathbf{a}\|_2^2 + 2\|\mathbf{b}\|_2^2$; step $\textcircled{2}$ uses Assumption 1.2.

Second, we have:

$$\begin{aligned}
\max_{i=1}^{\infty} \{\frac{1}{\beta^i} \|\mathbf{u}_n^{i+1}\|_2^2\} &= \max_{i=1}^{\infty} \{\frac{1}{\beta^i} \|\mathbf{Q}^i(\mathbf{x}_n^{i+1} - \mathbf{x}_n^i)\|_2^2\} \\
&\stackrel{\textcircled{1}}{\leq} \max_{i=1}^{\infty} \{\frac{1}{\beta^i} (q \bar{\lambda} \beta^i)^2 \|\mathbf{x}_n^{i+1} - \mathbf{x}_n^i\|_2^2\} \\
&\stackrel{\textcircled{2}}{\leq} q^2 \bar{\lambda} \sum_{i=1}^{\infty} \{\mathbf{L}_n^i \|\mathbf{x}_n^{i+1} - \mathbf{x}_n^i\|_2^2\} \\
&\stackrel{\textcircled{3}}{\leq} q^2 \bar{\lambda} \frac{K_e}{\varepsilon_2},
\end{aligned} \tag{53}$$

1566 where step ① uses $\|\mathbf{Q}^t\| \leq \beta^t \bar{\lambda} q$ for all $t \geq 0$, as shown in Lemma 4.3; step ② uses $\beta^i \bar{\lambda} \leq \mathsf{L}_n^i \triangleq$
 1567 $\beta^i \bar{\lambda} + L_n$; step ③ uses $K_e \geq \sum_{t=1}^{\infty} \mathcal{E}^{t+1} \geq \sum_{t=1}^{\infty} \mathcal{E}^{t+1} \geq \sum_{t=1}^{\infty} \varepsilon_2 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2$.

1568 (a) Using Part (b) of Lemma 4.2, we have:

$$1570 \quad \mathbf{A}_n^T \mathbf{z}^{t+1} = |1 - \sigma| \cdot \mathbf{A}_n^T \mathbf{z}^t + \sigma \{ \mathbf{w}_n^{t+1} + \mathbf{u}_n^{t+1} \}. \\ 1571$$

1572 Since $\|\cdot\|_2^2$ is convex, for all $t \geq 1$, we have:

$$1574 \quad \|\mathbf{A}_n^T \mathbf{z}^{t+1}\| - |1 - \sigma| \cdot \|\mathbf{A}_n^T \mathbf{z}^t\| \leq \sigma \{ \|\mathbf{w}_n^{t+1}\| + \|\mathbf{u}_n^{t+1}\| \}. \\ 1575$$

1577 Applying Lemma A.7 with $e^t \triangleq \|\mathbf{A}_n^T \mathbf{z}^t\|$ and $\Psi^t \triangleq \|\mathbf{w}_n^{t+1}\| + \|\mathbf{u}_n^{t+1}\|$, for all $t \geq 1$, we obtain:

$$\begin{aligned} 1579 \quad \|\mathbf{A}_n^T \mathbf{z}^t\|_2^2 &\leq (\|\mathbf{A}_n^T \mathbf{z}^1\| + \sigma_3 \max_{i=1}^{t-1} \{ \|\mathbf{w}_n^{i+1}\| + \|\mathbf{u}_n^{i+1}\| \ })^2 \\ 1580 &\stackrel{\textcircled{1}}{\leq} 3 \{ \bar{\lambda} \|\mathbf{z}^1\|_2^2 + \sigma_3 \max_{i=1}^{t-1} \|\mathbf{w}_n^{i+1}\|_2^2 + \sigma_3 \max_{i=1}^{t-1} \|\mathbf{u}_n^{i+1}\|_2^2 \ } \\ 1582 &\stackrel{\textcircled{2}}{\leq} 3 \beta^t \{ \frac{1}{\beta^0} \bar{\lambda} \|\mathbf{z}^1\|_2^2 + \sigma_3 \max_{i=1}^{\infty} \frac{1}{\beta^i} \|\mathbf{w}_n^{i+1}\|_2^2 + \sigma_3 \max_{i=1}^{\infty} \frac{1}{\beta^i} \|\mathbf{u}_n^{i+1}\|_2^2 \ } \\ 1584 &\stackrel{\textcircled{3}}{\leq} 3 \beta^t \{ \frac{1}{\beta^0} \bar{\lambda} \|\mathbf{z}^1\|_2^2 + 2\sigma_3 C_h^2 + 2\sigma_3 C_f^2 + \sigma_3 q^2 \bar{\lambda} \frac{K_e}{\varepsilon_2} \ } \\ 1586 &\stackrel{\textcircled{4}}{=} K_z \underline{\lambda} \beta^t, \end{aligned} \tag{54}$$

1588 where step ① use $(a + b + c)^2 \leq 3(a^2 + b^2 + c^2)$, Assumption 1.3 that $\|\mathbf{A}_n\|_2^2 \leq \bar{\lambda}$; step ② uses
 1589 $\beta^i \leq \beta^t$ for all $i \leq t$; ③ uses Inequalities (52) and (53); step ④ uses the definition of K_z . This
 1590 further leads to $\|\mathbf{z}^t\|_2^2 \leq \frac{1}{\underline{\lambda}} \|\mathbf{A}_n^T \mathbf{z}^t\|_2^2 = K_z \underline{\lambda} \beta^t$.

1591 (b) We have:

$$1593 \quad \ddot{K}_z \triangleq K_e / \varepsilon_3 \stackrel{\textcircled{1}}{\geq} \frac{1}{\varepsilon_3} \sum_{t=1}^{\infty} \mathcal{E}^{t+1} \stackrel{\textcircled{2}}{\geq} \frac{1}{\varepsilon_3} \sum_{t=1}^{\infty} \frac{\varepsilon_3}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2,$$

1596 where step ① uses Theorem 4.12; step ② uses the definition of $\mathcal{E}^{t+1} \triangleq [\varepsilon_1 \sum_{i=1}^{n-1} \mathsf{L}_i^t \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2] +$
 1597 $\varepsilon_2 \mathsf{L}_n^t \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \frac{\varepsilon_3}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2$ in Lemma 4.1.

1599 \square

1601 D.14 PROOF OF LEMMA 4.14

1603 *Proof.* We let $\sigma \in (0, 2)$.

1605 First, we derive the following inequalities:

$$\begin{aligned} 1606 \quad \langle \mathbf{A} \mathbf{x}^{t+1} - \mathbf{b}, \mathbf{z}^{t+1} \rangle &= \frac{1}{\sigma \beta^t} \langle \mathbf{z}^{t+1} - \mathbf{z}^t, \mathbf{z}^{t+1} \rangle \\ 1608 &\stackrel{\textcircled{1}}{=} \frac{1}{2\sigma} \{ \frac{1}{\beta^t} \|\mathbf{z}^{t+1}\|_2^2 - \frac{1}{\beta^t} \|\mathbf{z}^t\|_2^2 + \frac{1}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \ } \\ 1610 &\geq -\frac{1}{2\sigma \beta^t} \|\mathbf{z}^t\|_2^2, \end{aligned} \tag{55}$$

1612 where step ① uses the Pythagoras relation in Fact A.1.

1613 We consider Lemma 4.6 and Lemma 4.9. We let $i \geq 1$. Given $\mathcal{E}^{i+1} \geq 0$, it follows that:

$$1615 \quad 0 \leq \Theta^i - \Theta^{i+1} + \Gamma_{\mu}^i.$$

1617 Telescoping this inequality over i from 1 to t , we have:

$$1619 \quad 0 \leq \Theta^1 - \Theta^{t+1} + \sum_{i=1}^t \Gamma_{\mu}^i \stackrel{\textcircled{1}}{\leq} \Theta^1 - \Theta^{t+1} + C_{\mu},$$

1620 where step ① uses Lemma 4.11. For all $t \geq 1$, we derive the following results:

$$\begin{aligned}
1621 \quad \Theta^1 + C_\mu &\geq \Theta^{t+1} \\
1622 \quad &\stackrel{\textcircled{1}}{=} \Theta_L^{t+1} + \Theta_a^{t+1} + \Theta_u^{t+1} \\
1623 \quad &\stackrel{\textcircled{2}}{=} \mathcal{L}(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}; \beta^{t+1}, \mu^{t+1}) + \frac{1}{2}C_h\mu^{t+1} + \Theta_a^{t+1} + \Theta_u^{t+1} \\
1624 \quad &\stackrel{\textcircled{3}}{=} \sum_{i=1}^n f_i(\mathbf{x}_i^{t+1}) + \langle \mathbf{A}\mathbf{x}^{t+1} - \mathbf{b}, \mathbf{z}^{t+1} \rangle + \frac{\beta^{t+1}}{2} \|\mathbf{A}\mathbf{x}^{t+1} - \mathbf{b}\|_2^2 \\
1625 \quad &+ \{\sum_{i=1}^{n-1} h_i(\mathbf{x}_i^{t+1})\} + h_n(\mathbf{x}_n^{t+1}; \mu^{t+1}) + \frac{1}{2}C_h\mu^{t+1} + \Theta_a^{t+1} + \Theta_u^{t+1} \\
1626 \quad &\stackrel{\textcircled{4}}{\geq} \sum_{i=1}^n [f_i(\mathbf{x}_i^{t+1}) + h_i(\mathbf{x}_i^{t+1})] + \langle \mathbf{A}\mathbf{x}^{t+1} - \mathbf{b}, \mathbf{z}^{t+1} \rangle - \frac{1}{2}\mu^{t+1}C_h^2 \\
1627 \quad &\stackrel{\textcircled{5}}{\geq} \sum_{i=1}^n [f_i(\mathbf{x}_i^{t+1}) + h_i(\mathbf{x}_i^{t+1})] - \frac{1}{2\sigma\beta^t} \|\mathbf{z}^t\|_2^2 - \frac{1}{2}\mu^{t+1}C_h^2 \\
1628 \quad &\stackrel{\textcircled{6}}{\geq} \sum_{i=1}^n [f_i(\mathbf{x}_i^{t+1}) + h_i(\mathbf{x}_i^{t+1})] - \frac{1}{2\sigma\beta^t} \|\mathbf{z}^t\|_2^2 - \frac{1}{2}\mu^0C_h^2,
\end{aligned}$$

1629 where step ① uses the definition of Θ^{t+1} ; step ② uses the definition of Θ_L^{t+1} in Lemma 4.1;
1630 step ③ uses the definition of $\mathcal{L}(\mathbf{x}^{t+1}, \mathbf{z}^{t+1}; \beta^{t+1}, \mu^{t+1})$ in (4); step ④ uses $\frac{\beta^{t+1}}{2} \|\mathbf{A}\mathbf{x}^{t+1} - \mathbf{b}\|_2^2 \geq 0$,
1631 $\frac{1}{2}C_h\mu^{t+1} \geq 0$, $\Theta_a^{t+1} \geq 0$, $\Theta_u^{t+1} \geq 0$, and the fact that $h_n(\mathbf{x}_n^{t+1}; \mu^{t+1}) \geq h_n(\mathbf{x}_n^{t+1}) - \frac{1}{2}\mu^{t+1}C_h^2$;
1632 step ⑤ uses Inequality (55); step ⑥ uses $\mu^t \leq \mu^0$ for all t .

1633 We further obtain:

$$\begin{aligned}
1640 \quad \sum_{i=1}^n [f_i(\mathbf{x}_i^{t+1}) + h_i(\mathbf{x}_i^{t+1})] &\leq \Theta^1 + C_\mu + \frac{1}{2\sigma\beta^t} \|\mathbf{z}^t\|_2^2 + \frac{1}{2}\mu^0C_h^2 \\
1641 \quad &\stackrel{\textcircled{1}}{<} +\infty,
\end{aligned}$$

1642 where step ① uses the boundedness of $\frac{1}{\beta^t} \|\mathbf{z}^t\|_2^2$ for all $t \geq 0$, as shown in Lemma 4.13. According
1643 to Assumption 1.4, we have $\|\mathbf{x}_i^{t+1}\| < +\infty$ for all $i \in [n]$. □

1648 PROOF OF THEOREM 4.15

1649 *Proof.* We define $K_c \triangleq \frac{K_e}{K'_c}$, where $K'_c \triangleq \min\{\min(\varepsilon_1, \varepsilon_2)\underline{A}, \epsilon_3\}$, and $\underline{A} \triangleq \min_{i=1}^n \|\mathbf{A}_i\|_2^2$.

1650 We define $\mathcal{E}^{t+1} \triangleq [\varepsilon_1 \sum_{i=1}^{n-1} \mathsf{L}_i \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2 + \varepsilon_2 \mathsf{L}_n \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \frac{\varepsilon_3}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2]$.

1651 (a) We have:

$$\begin{aligned}
1652 \quad K_e &\stackrel{\textcircled{1}}{\geq} \sum_{t=1}^T \mathcal{E}^{t+1} \\
1653 \quad &\stackrel{\textcircled{2}}{=} \sum_{t=1}^T \{\varepsilon_1 \sum_{i=1}^{n-1} \mathsf{L}_i \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2 + \varepsilon_2 \mathsf{L}_n \|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 + \frac{\varepsilon_3}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2\} \\
1654 \quad &\stackrel{\textcircled{3}}{\geq} \frac{1}{\beta^T} \sum_{t=1}^T \{[\varepsilon_1 \sum_{i=1}^{n-1} \frac{\mathsf{L}_i}{\beta^t} \|\beta^t(\mathbf{x}_i^{t+1} - \mathbf{x}_i^t)\|_2^2] + \varepsilon_2 \frac{\mathsf{L}_n}{\beta^t} \|\beta^t(\mathbf{x}_n^{t+1} - \mathbf{x}_n^t)\|_2^2 + \varepsilon_3 \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2\} \\
1655 \quad &\stackrel{\textcircled{4}}{\geq} \frac{1}{\beta^T} \sum_{t=1}^T \{[\varepsilon_1 \sum_{i=1}^{n-1} \underline{A} \|\beta^t(\mathbf{x}_i^{t+1} - \mathbf{x}_i^t)\|_2^2] + \varepsilon_2 \underline{A} \|\beta^t(\mathbf{x}_n^{t+1} - \mathbf{x}_n^t)\|_2^2 + \varepsilon_3 \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2\} \\
1656 \quad &\stackrel{\textcircled{5}}{\geq} \frac{1}{\beta^T} \cdot K'_c \cdot \sum_{t=1}^T \{\sum_{i=1}^n \|\beta^t(\mathbf{x}_i^{t+1} - \mathbf{x}_i^t)\|_2^2 + \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2\} \\
1657 \quad &\stackrel{\textcircled{6}}{=} \frac{1}{\beta^T} \cdot K'_c \cdot \sum_{t=1}^T \{\|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2 + \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2\},
\end{aligned}$$

1658 where step ① uses Theorem (4.12); step ② uses the definition of \mathcal{E}^{t+1} ; step ③ uses $\beta^T \geq \beta^t$ for
1659 all $t \leq T$; step ④ uses $\frac{\mathsf{L}_i}{\beta^t} = \frac{\mathsf{L}_i + \beta^t \|\mathbf{A}_i\|_2^2}{\beta^t} \geq \|\mathbf{A}_i\|_2^2 \geq \underline{A}$; step ⑤ uses the definition of $K'_c \triangleq$
1660 $\min\{\min(\varepsilon_1, \varepsilon_2)\underline{A}, \epsilon_3\}$; step ⑥ uses $\sum_{i=1}^n \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2 = \|\mathbf{x}^{t+1} - \mathbf{x}^t\|_2^2$. Therefore, we obtain:

$$\sum_{t=1}^T \{\|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2 + \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2\} \leq \frac{K_e}{K'_c} \beta^T = K_c \beta^T.$$

1661 (c) By dividing both sides of the above inequality by T , we obtain:

$$\begin{aligned}
1662 \quad \frac{K_c \beta^T}{T} &\geq \frac{1}{T} \sum_{t=1}^T \{\|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2 + \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2\} \\
1663 &\geq \min_{t=1}^T \{\|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2 + \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2\}.
\end{aligned}$$

We conclude that there exists an index \bar{t} with $\bar{t} \leq T$ such that $\|\mathbf{z}^{\bar{t}+1} - \mathbf{z}^{\bar{t}}\|_2^2 + \|\beta^{\bar{t}}(\mathbf{x}^{\bar{t}+1} - \mathbf{x}^{\bar{t}})\|_2^2 \leq \frac{K_c \beta^T}{T}$.

□

D.16 PROOF OF THEOREM 4.18

To prove this theorem, we first provide the following lemma.

Lemma D.1. We define $\mathbf{q}^t \triangleq \{\mathbf{x}_1^t, \mathbf{x}_2^t, \dots, \mathbf{x}_{n-1}^t, \check{\mathbf{x}}_n^t\}$. We have:

$$(a) \|\mathbf{A}\mathbf{q}^{t+1} - \mathbf{b}\|_2^2 \leq B_1\|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + B_2(\beta^t)^{-2}.$$

$$(b) \text{dist}^2(\mathbf{0}, \partial h_n(\check{\mathbf{x}}_n^{t+1}) + \nabla_{\mathbf{x}_n} f_n(\check{\mathbf{x}}_n^{t+1}) + \mathbf{A}_n^\top \mathbf{z}^{t+1}) \leq B_3\|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + B_4\|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2 + B_5(\beta^t)^{-2}.$$

$$(c) \sum_{i=1}^{n-1} \text{dist}^2(\mathbf{0}, \partial h_i(\mathbf{x}_i^{t+1}) + \nabla_{\mathbf{x}_i} f_i(\mathbf{x}_i^{t+1}) + \mathbf{A}_i^\top \mathbf{z}^{t+1}) \leq B_6\|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + B_7\|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2.$$

Here, $B_1 = \frac{2}{\sigma^2(\beta^0)^2}$, $B_2 = 2\bar{A}(\frac{C_h}{\delta\lambda})^2$, $B_3 = 4(1 - \frac{1}{\sigma})^2\bar{A}$, $B_4 = 4q^2\bar{\lambda}^2 + \frac{4L_n^2}{(\beta^0)^2}$, $B_5 = \frac{4L_n^2 C_h^2}{(\delta\lambda)^2}$, $B_6 = 3(1 - \frac{1}{\sigma})^2\bar{A}(n-1)$, and $B_7 = \frac{3\bar{L}}{(\beta^0)^2} + 6\theta_1^2(\frac{\bar{L}}{\beta^0} + \bar{A})^2 + 6\bar{A}^2(n-1)$. Furthermore, $\bar{A} \triangleq \max_{i=1}^n \|\mathbf{A}_i\|_2^2$, $\bar{L} \triangleq \max_{i=1}^n L_i$.

Proof. We define $\mathbf{u}_i^{t+1} = \theta_1 \mathbf{L}_i^t (\mathbf{x}_i^{t+1} - \mathbf{x}_i^t) - \beta^t \mathbf{A}_i^\top [\sum_{j=i}^n \mathbf{A}_j (\mathbf{x}_j^{t+1} - \mathbf{x}_j^t)]$ with $i \in [n-1]$.

We define $\mathbf{u}_n^{t+1} \triangleq \mathbf{Q}^t(\mathbf{x}_n^{t+1} - \mathbf{x}_n^t)$ with $\mathbf{Q}^t \triangleq \theta_2 \mathbf{L}_n^t \mathbf{I} - \beta^t \mathbf{A}_n^\top \mathbf{A}_n$.

(a) We have:

$$\begin{aligned} & \|\mathbf{A}\mathbf{q}^{t+1} - \mathbf{b}\|_2^2 \\ &= \|\sum_{i=1}^n \mathbf{A}_i \mathbf{x}_i^{t+1} - \mathbf{A}_n \mathbf{x}_n^{t+1} + \mathbf{A}_n \check{\mathbf{x}}_n^{t+1} - \mathbf{b}\|_2^2 \\ &\stackrel{(1)}{\leq} 2\|\sum_{i=1}^n \mathbf{A}_i \mathbf{x}_i^{t+1} - \mathbf{b}\|_2^2 + 2\|\mathbf{A}_n(\mathbf{x}_n^{t+1} - \check{\mathbf{x}}_n^{t+1})\|_2^2 \\ &\stackrel{(2)}{\leq} 2\|\mathbf{A}\mathbf{x}^{t+1} - \mathbf{b}\|_2^2 + 2\bar{A}(\mu^t C_h)^2 \\ &\stackrel{(3)}{=} 2\|\frac{1}{\sigma\beta^t}(\mathbf{z}^{t+1} - \mathbf{z}^t)\|_2^2 + 2\bar{A}(\frac{C_h}{\delta\lambda\beta^t})^2, \\ &\stackrel{(4)}{\leq} \underbrace{\frac{2}{\sigma^2(\beta^0)^2}\|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2}_{\triangleq B_1} + \underbrace{2\bar{A}(\frac{C_h}{\delta\lambda})^2 \cdot (\beta^t)^{-2}}_{\triangleq B_2}, \end{aligned}$$

where step ① uses the inequality that $\|\mathbf{a} - \mathbf{b}\|_2^2 \leq 2\|\mathbf{a}\|_2^2 + 2\|\mathbf{b}\|_2^2$ for all \mathbf{a} and \mathbf{b} ; step ② uses $\|\mathbf{A}_n\|_2^2 \leq \bar{A}$ and Part (c) in Lemma 3.6; step ③ uses $\mathbf{z}^{t+1} = \mathbf{z}^t + \beta^t \sigma(\mathbf{A}\mathbf{x}^{t+1} - \mathbf{b})$; step ④ uses $\beta^0 \leq \beta^t$.

(b) We first have the following inequalities:

$$\begin{aligned} & 2\|\nabla f_n(\check{\mathbf{x}}_n^{t+1}) - \nabla f_n(\mathbf{x}_n^t)\|_2^2 \\ &= 2\|\nabla f_n(\check{\mathbf{x}}_n^{t+1}) - \nabla f_n(\mathbf{x}_n^{t+1}) + \nabla f_n(\mathbf{x}_n^{t+1}) - \nabla f_n(\mathbf{x}_n^t)\|_2^2 \\ &\stackrel{(1)}{\leq} 4\|\nabla f_n(\check{\mathbf{x}}_n^{t+1}) - \nabla f_n(\mathbf{x}_n^{t+1})\|_2^2 + 4\|\nabla f_n(\mathbf{x}_n^{t+1}) - \nabla f_n(\mathbf{x}_n^t)\|_2^2 \\ &\stackrel{(2)}{\leq} 4L_n^2\|\check{\mathbf{x}}_n^{t+1} - \mathbf{x}_n^{t+1}\|_2^2 + 4L_n^2\|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2 \\ &\stackrel{(3)}{\leq} 4L_n^2(\mu^t)^2 C_h^2 + 4L_n^2 \frac{1}{(\beta^t)^2} \|\beta^t(\mathbf{x}_n^{t+1} - \mathbf{x}_n^t)\|_2^2 \\ &\stackrel{(4)}{\leq} \underbrace{4L_n^2 \frac{1}{(\delta\lambda)^2} C_h^2 \cdot \frac{1}{(\beta^t)^2}}_{\triangleq B_5} + 4L_n^2\|\mathbf{x}_n^{t+1} - \mathbf{x}_n^t\|_2^2, \end{aligned} \tag{56}$$

where step ① uses the inequality that $\|\mathbf{a} - \mathbf{b}\|_2^2 \leq 2\|\mathbf{a}\|_2^2 + 2\|\mathbf{b}\|_2^2$ for all \mathbf{a} and \mathbf{b} ; step ② uses the fact that $f_n(\mathbf{x}_n)$ is L_n -smooth; step ③ uses Part (c) of Lemma 3.6 that: $\|\check{\mathbf{x}}_n^{t+1} - \mathbf{x}_n^{t+1}\| \leq \mu^t C_h$; step ④ uses $\mu^t \leq \frac{1}{\delta\lambda\beta^t}$.

1728 We further obtain:
1729
1730 $\text{dist}^2(\mathbf{0}, \partial h_n(\check{\mathbf{x}}_n^{t+1}) + \nabla f_i(\check{\mathbf{x}}_i^{t+1}) + \mathbf{A}_i^\top \mathbf{z}^{t+1})$
1731 $\stackrel{\textcircled{1}}{=} \|\theta_2 \mathsf{L}_n^t(\mathbf{c}^t - \mathbf{x}_n^{t+1}) + \nabla f_i(\check{\mathbf{x}}_i^{t+1}) + \mathbf{A}_i^\top \mathbf{z}^{t+1}\|_2^2$
1732 $\stackrel{\textcircled{2}}{=} \|\theta_2 \mathsf{L}_n^t(\mathbf{x}_n^t - \frac{1}{\theta_2 \mathsf{L}_n^t} \mathbf{g} - \mathbf{x}_n^{t+1}) + \nabla f_i(\check{\mathbf{x}}_i^{t+1}) + \mathbf{A}_i^\top \mathbf{z}^{t+1}\|_2^2$
1733 $\stackrel{\textcircled{3}}{=} \|(\theta_2 \mathsf{L}_n^t - \beta^t \mathbf{A}_n^\top \mathbf{A}_n)(\mathbf{x}_n^t - \mathbf{x}_n^{t+1}) + \nabla f_n(\check{\mathbf{x}}_n^{t+1}) - \nabla f_n(\mathbf{x}_n^t) + (1 - \frac{1}{\sigma}) \mathbf{A}_n^\top (\mathbf{z}^{t+1} - \mathbf{z}^t)\|_2^2$
1734 $\stackrel{\textcircled{4}}{\leq} 2\|\mathbf{Q}(\mathbf{x}_n^t - \mathbf{x}_n^{t+1}) + (1 - \frac{1}{\sigma}) \mathbf{A}_n^\top (\mathbf{z}^{t+1} - \mathbf{z}^t)\|_2^2 + 2\|\nabla f_n(\check{\mathbf{x}}_n^{t+1}) - \nabla f_n(\mathbf{x}_n^t)\|_2^2$
1735 $\stackrel{\textcircled{5}}{\leq} 4(1 - \frac{1}{\sigma})^2 \bar{\mathbf{A}} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + 4\|\mathbf{Q}(\mathbf{x}_n^t - \mathbf{x}_n^{t+1})\|_2^2 + 2\|\nabla f_n(\check{\mathbf{x}}_n^{t+1}) - \nabla f_n(\mathbf{x}_n^t)\|_2^2$
1736 $\stackrel{\textcircled{6}}{\leq} \underbrace{4(1 - \frac{1}{\sigma})^2 \bar{\mathbf{A}} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2}_{\triangleq B_3} + \underbrace{\{4q^2 \lambda^2 + 4L_n^2 \frac{1}{(\beta^0)^2}\} \cdot \|\beta^t(\mathbf{x}_n^t - \mathbf{x}_n^{t+1})\|_2^2}_{\triangleq B_4} + \frac{B_5}{(\beta^t)^2},$ (57)
1741
1742

1743 where step ① uses the optimality condition as shown in Part (b) of Lemma 3.6 that:

1744 $\rho(\mathbf{c}^t - \mathbf{x}_n^{t+1}) \in \partial h_n(\check{\mathbf{x}}_n^{t+1}), \text{ with } \rho = \theta_2 \mathsf{L}_n^t;$

1745 step ② uses $\mathbf{c}^t = \mathbf{x}_n^t - \mathbf{g}/\rho$ as shown in Algorithm 1; step ③ uses the fact that:

1746 $\mathbf{g} = \nabla f_n(\mathbf{x}_n^t) + \mathbf{A}_n^\top \mathbf{z}^t + \frac{1}{\sigma} \mathbf{A}_n^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) + \beta^t \mathbf{A}_n^\top \mathbf{A}_n(\mathbf{x}_n^t - \mathbf{x}_n^{t+1}),$

1747 step ④ uses the definition of \mathbf{Q} as in Lemma 4.2 and the inequality that $\|\mathbf{a} - \mathbf{b}\|_2^2 \leq 2\|\mathbf{a}\|_2^2 + 2\|\mathbf{b}\|_2^2$;

1748 step ⑤ uses the inequality that $\|\mathbf{a} - \mathbf{b}\|_2^2 \leq 2\|\mathbf{a}\|_2^2 + 2\|\mathbf{b}\|_2^2$ and $\|\mathbf{A}_n\| \leq \bar{\mathbf{A}}$; step ⑥ uses $\|\mathbf{Q}^t\| \leq \beta^t \lambda q$ as shown in Lemma 4.3, Inequality (56), and the fact that $\beta^0 \leq \beta^t$.

1749 (c) We first have the following inequalities:
1750

1751 $\sum_{i=1}^{n-1} \|\mathbf{u}_i^{t+1}\|_2^2$
1752 $\stackrel{\textcircled{1}}{=} \|\sum_{i=1}^{n-1} \{\theta_1 \mathsf{L}_i^t(\mathbf{x}_i^{t+1} - \mathbf{x}_i^t) - \beta^t \mathbf{A}_i^\top [\sum_{j=i}^n \mathbf{A}_j(\mathbf{x}_j^{t+1} - \mathbf{x}_j^t)]\}\|_2^2$
1753 $\stackrel{\textcircled{2}}{\leq} 2\|\sum_{i=1}^{n-1} \theta_1 \mathsf{L}_i^t(\mathbf{x}_i^{t+1} - \mathbf{x}_i^t)\|_2^2 + 2\|\sum_{i=1}^{n-1} \beta^t \mathbf{A}_i^\top [\sum_{j=i}^n \mathbf{A}_j(\mathbf{x}_j^{t+1} - \mathbf{x}_j^t)]\|_2^2$
1754 $\stackrel{\textcircled{3}}{\leq} 2(\theta_1 \mathsf{L}_i^t)^2 \sum_{i=1}^{n-1} \|\mathbf{x}_i^{t+1} - \mathbf{x}_i^t\|_2^2 + 2\bar{\mathbf{A}}^2(n-1) \sum_{j=1}^{n-1} \|\beta^t(\mathbf{x}_j^{t+1} - \mathbf{x}_j^t)\|_2^2$
1755 $\stackrel{\textcircled{4}}{\leq} 2(\theta_1 \mathsf{L}_i^t)^2 \|\mathbf{x}^{t+1} - \mathbf{x}^t\|_2^2 + 2\bar{\mathbf{A}}^2(n-1) \|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2$
1756 $\stackrel{\textcircled{5}}{\leq} 2\theta_1^2 (\frac{\bar{L}}{\beta^0} + \bar{\mathbf{A}})^2 \cdot \|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2 + 2\bar{\mathbf{A}}^2(n-1) \|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2$
1757 $= \{2\theta_1^2 (\frac{\bar{L}}{\beta^0} + \bar{\mathbf{A}})^2 + 2\bar{\mathbf{A}}^2(n-1)\} \cdot \|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2,$ (58)

1758 where step ① uses the definition of \mathbf{u}_i^{t+1} for all $i \in [n-1]$; step ② uses the inequality that $\|\mathbf{a} - \mathbf{b}\|_2^2 \leq$

1759 $2\|\mathbf{a}\|_2^2 + 2\|\mathbf{b}\|_2^2$; step ③ uses $\|\mathbf{A}_n\|_2^2 \leq \bar{\mathbf{A}}$; step ④ uses $\sum_{j=1}^{n-1} \|\mathbf{x}_j^{t+1} - \mathbf{x}_j^t\|_2^2 \leq \|\mathbf{x}^{t+1} - \mathbf{x}^t\|_2^2$; step

1760 ⑤ uses $\mathsf{L}_i^t = L_i + \beta^t \|\mathbf{A}_i\|_2^2 \leq \frac{\beta^t L_i}{\beta^0} + \beta^t \bar{\mathbf{A}} \leq \frac{\beta^t \bar{L}}{\beta^0} + \beta^t \bar{\mathbf{A}}$.

1761 We have:

1762 $\sum_{i=1}^{n-1} \text{dist}^2(\partial h_i(\mathbf{x}_i^{t+1}) + \nabla f_i(\mathbf{x}_i^{t+1}) + \mathbf{A}_i^\top \mathbf{z}^{t+1})$
1763 $\stackrel{\textcircled{1}}{=} \sum_{i=1}^{n-1} \|(1 - \frac{1}{\sigma}) \mathbf{A}_i^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) - \nabla f_i(\mathbf{x}_i^t) - \mathbf{u}_i^{t+1} + \nabla f_i(\mathbf{x}_i^{t+1})\|_2^2$
1764 $\stackrel{\textcircled{2}}{\leq} 3\sum_{i=1}^{n-1} \|(1 - \frac{1}{\sigma}) \mathbf{A}_i^\top (\mathbf{z}^{t+1} - \mathbf{z}^t)\|_2^2 + 3\sum_{i=1}^{n-1} \|\nabla f_i(\mathbf{x}_i^t) - \nabla f_i(\mathbf{x}_i^{t+1})\|_2^2 + 3\sum_{i=1}^{n-1} \|\mathbf{u}_i^{t+1}\|_2^2$
1765 $\stackrel{\textcircled{3}}{\leq} \underbrace{3(1 - \frac{1}{\sigma})^2 \bar{\mathbf{A}}(n-1)}_{\triangleq B_6} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + \underbrace{\frac{3\bar{L}}{(\beta^0)^2} \|\beta^t(\mathbf{x}^t - \mathbf{x}^{t+1})\|_2^2}_{\triangleq B_7} + 3\sum_{i=1}^{n-1} \|\mathbf{u}_i^{t+1}\|_2^2$
1766 $\stackrel{\textcircled{4}}{=} B_6 \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + \underbrace{6\theta_1^2 (\frac{\bar{L}}{\beta^0} + \bar{\mathbf{A}})^2 + 6\bar{\mathbf{A}}^2(n-1) \cdot \|\beta^t(\mathbf{x}^t - \mathbf{x}^{t+1})\|_2^2}_{\triangleq B_7},$

1782 where step ① uses Part (a) in Lemma 4.2 that:
 1783

$$1784 \quad i \in [n-1], \partial h_i(\mathbf{x}_i^{t+1}) \ni -\mathbf{u}_i^{t+1} - \mathbf{A}_i^\top \mathbf{z}^t - \frac{1}{\sigma} \mathbf{A}_i^\top (\mathbf{z}^{t+1} - \mathbf{z}^t) - \nabla f_i(\mathbf{x}_i^t);$$

1785 step ② uses the inequality that $\|\mathbf{a} + \mathbf{b} + \mathbf{c}\|_2^2 \leq 3\|\mathbf{a}\|_2^2 + 3\|\mathbf{b}\|_2^2 + 3\|\mathbf{c}\|_2^2$; step ③ uses $\|\mathbf{A}_i\|_2^2 \leq \bar{\mathbf{A}}$,
 1786 $f_i(\mathbf{x}_i)$ is L_i -smooth, $L_i \leq \bar{L}$, and $\beta^0 \leq \beta^t$; step ④ uses Inequality (58). \square
 1787

1788 Now, we proceed to prove the theorem.
 1789

1790 *Proof.* We define $\text{Crit}(\mathbf{x}, \mathbf{z}) \triangleq \|\mathbf{Ax} - \mathbf{b}\|_2^2 + \sum_{i=1}^n \text{dist}^2(\mathbf{0}, \nabla f_i(\mathbf{x}_i) + \partial h_i(\mathbf{x}_i) + \mathbf{A}_i^\top \mathbf{z})$.
 1791

1792 We define $\mathbf{q}^t \triangleq \{\mathbf{x}_1^t, \mathbf{x}_2^t, \dots, \mathbf{x}_{n-1}^t, \check{\mathbf{x}}_n^t\}$.
 1793

1794 Using lemma D.1, for all $t \geq 0$, we have:
 1795

$$1796 \quad \text{Crit}(\mathbf{q}^{t+1}, \mathbf{z}^{t+1}) \\ 1797 \leq \underbrace{(B_1 + B_3 + B_6)}_{\triangleq D_1} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + \underbrace{(B_4 + B_7)}_{\triangleq D_2} \|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2 + \underbrace{(B_2 + B_5)}_{\triangleq D_3} (\beta^t)^{-2} \quad (59)$$

1799 We further derive:
 1800

$$1801 \quad \frac{1}{T} \sum_{t=0}^T \text{Crit}(\mathbf{q}^{t+1}, \mathbf{z}^{t+1}) \\ 1802 \stackrel{(1)}{\leq} \frac{1}{T} \max(D_1, D_2) \sum_{t=0}^T \{\|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + \|\beta^t(\mathbf{x}^{t+1} - \mathbf{x}^t)\|_2^2\} + \frac{D_3}{T} (\beta^0)^{-2} + \frac{D_3}{T} \sum_{t=0}^T (\beta^t)^{-2} \\ 1803 \stackrel{(2)}{\leq} \frac{1}{T} \max(D_1, D_2) K_c \beta^T + \frac{D_3}{T} (\beta^0)^{-2} + \frac{D_3}{(\beta^0 \xi)^2} \cdot \frac{1}{T} \sum_{t=1}^T t^{-2p} \\ 1804 \stackrel{(3)}{\leq} \frac{1}{T} \max(D_1, D_2) K_c \beta^T + \frac{D_3}{T} (\beta^0)^{-2} + \frac{D_3}{(\beta^0 \xi)^2} \cdot \frac{1}{T} \cdot \frac{T^{1-2p}}{1-2p}, \text{ with } 2p \in (0, 1) \\ 1805 = \mathcal{O}(T^{p-1}) + \mathcal{O}(T^{-1}) + \mathcal{O}(T^{1-2p}), \text{ with } 2p \in (0, 1).$$

1806 Here, step ① uses Inequality 59; step ② uses Theorem 4.15, and $\frac{1}{(\beta^t)^2} = \frac{1}{(\beta^0 + \beta^0 \xi t^p)^2} \leq \frac{1}{(\beta^0 \xi t^p)^2}$;
 1807 step ③ uses the fact that $\sum_{t=1}^T t^{-p'} \leq \frac{T^{(1-p')}}{1-p'}$ if $p' \in (0, 1)$, as shown in Lemma A.6.
 1808

1809 In particular, with the choice $p = 1/3$, we have: $\frac{1}{T} \sum_{t=0}^T \text{Crit}(\mathbf{q}^{t+1}, \mathbf{z}^{t+1}) \leq \mathcal{O}(T^{-2/3})$.
 1810 \square

1811 D.17 PROOF OF LEMMA 4.20

1812 *Proof.* We let $\frac{\mathbf{z}^t}{\sqrt{\beta^t}} \triangleq \hat{\mathbf{z}}^t$ for all t .
 1813

1814 Initially, we derive:
 1815

$$1816 \quad \sum_{t=1}^{\infty} (1 - \sqrt{\frac{\beta^t}{\beta^{t+1}}})^2 \stackrel{(1)}{=} \sum_{t=1}^{\infty} (1 - \sqrt{\frac{1+\xi t^p}{1+\xi(t+1)^p}})^2 \\ 1817 \stackrel{(2)}{\leq} \sum_{t=1}^{\infty} (1 - \sqrt{\frac{t^p}{(t+1)^p}})^2 \\ 1818 = \sum_{t=1}^{\infty} \frac{\{(t+1)^{p/2} - t^{p/2}\}^2}{(t+1)^p} \\ 1819 \stackrel{(3)}{\leq} \sum_{t=1}^{\infty} \frac{\{\frac{p}{2} \cdot t^{(p/2-1)}\}^2}{t^p} \\ 1820 \stackrel{(4)}{\leq} \frac{1}{4} \sum_{t=1}^{\infty} \frac{t^{(p-2)}}{t^p} \\ 1821 \stackrel{(5)}{\leq} 1/2, \quad (60)$$

1822 where step ① uses $\beta^t = \beta^0(1 + \xi t^p)$ for all $t \geq 0$; step ② uses $\frac{1+\xi t^p}{1+\xi(t+1)^p} \leq \frac{\xi t^p}{\xi(t+1)^p}$; step ③ uses
 1823 Lemma A.4 that $(t+1)^{p/2} - t^{p/2} \leq \frac{p}{2} t^{(p/2-1)}$ for all $t \geq 1$ and $\frac{p}{2} \in (0, 1)$; step ④ uses $p \leq 1$ and
 1824 $\frac{1}{t+1} \leq \frac{1}{t}$; step ⑤ uses $\sum_{t=1}^{\infty} \frac{1}{t^2} = \frac{\pi^2}{6} < 2$.
 1825

1836 (a) We have: $\|\hat{\mathbf{z}}^t\|_2^2 = \|\frac{\mathbf{z}^t}{\sqrt{\beta^t}}\|_2^2 = \frac{1}{\beta^t}\|\mathbf{z}^t\|_2^2 \leq K_z < +\infty$, where the last step uses Lemma 4.13.
 1837

1838 (b) We have:

$$\begin{aligned}
 \sum_{t=1}^{\infty} \|\hat{\mathbf{z}}^{t+1} - \hat{\mathbf{z}}^t\|_2^2 &\stackrel{(1)}{=} \sum_{t=1}^{\infty} \left\| \frac{\mathbf{z}^{t+1}}{\sqrt{\beta^{t+1}}} - \frac{\mathbf{z}^t}{\sqrt{\beta^t}} \right\|_2^2 \\
 &= \sum_{t=1}^{\infty} \left\| \frac{\mathbf{z}^{t+1} - \mathbf{z}^t}{\sqrt{\beta^{t+1}}} - \mathbf{z}^t \left(\frac{1}{\sqrt{\beta^t}} - \frac{1}{\sqrt{\beta^{t+1}}} \right) \right\|_2^2 \\
 &\stackrel{(2)}{\leq} 2 \sum_{t=1}^{\infty} \left\| \frac{\mathbf{z}^{t+1} - \mathbf{z}^t}{\sqrt{\beta^{t+1}}} \right\|_2^2 + 2 \sum_{t=1}^{\infty} \left\| \mathbf{z}^t \left(\frac{1}{\sqrt{\beta^t}} - \frac{1}{\sqrt{\beta^{t+1}}} \right) \right\|_2^2 \\
 &\stackrel{(3)}{\leq} 2 \sum_{t=1}^{\infty} \frac{1}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 + 2 \sum_{t=1}^{\infty} \frac{1}{\beta^t} \left\| \left(1 - \sqrt{\frac{\beta^t}{\beta^{t+1}}} \right) \cdot \mathbf{z}^t \right\|_2^2 \\
 &\stackrel{(4)}{\leq} 2K_{zz} + \frac{2}{\beta^t} \|\mathbf{z}\|_2^2 \cdot \sum_{t=1}^{\infty} \left(1 - \sqrt{\frac{\beta^t}{\beta^{t+1}}} \right)^2 \\
 &\stackrel{(5)}{\leq} 2K_{zz} + 2K_z \cdot \frac{1}{2},
 \end{aligned}$$

1851 where step ① uses the definition $\frac{\mathbf{z}^t}{\sqrt{\beta^t}} \triangleq \hat{\mathbf{z}}^t$ for all t ; step ② uses $\|\mathbf{a} - \mathbf{b}\|_2^2 \leq 2\|\mathbf{a}\|_2^2 + 2\|\mathbf{b}\|_2^2$; step
 1852 ③ uses $\frac{1}{\beta^{t+1}} \leq \frac{1}{\beta^t}$; step ④ uses $\sum_{t=1}^{\infty} \frac{1}{\beta^t} \|\mathbf{z}^{t+1} - \mathbf{z}^t\|_2^2 \leq K_{zz}$ as shown in Lemma 4.13; step ⑤ uses
 1853 Inequality (60), and $\frac{1}{\beta^t} \|\mathbf{z}^t\|_2^2 \leq K_z$ as shown in Lemma 4.13.

1855 \square

1856 E ADDITIONAL EXPERIMENT DETAILS AND RESULTS

1860 We offer further experimental details in Sections E.1 and E.2, and include additional results in
 1861 Section E.3.

1863 E.1 DATASETS

1865 We incorporate four datasets in our experiments, including both randomly generated data and pub-
 1866 licly available real-world data. These datasets serve as our data matrices $\mathbf{D} \in \mathbb{R}^{m \times d}$. The
 1867 dataset names are as follows: ‘TDT2- m - d ’, ‘sector- m - d ’, ‘mnist- m - d ’, and ‘randn- m - d ’. Here,
 1868 randn(m, n) refers to a function that generates a standard Gaussian random matrix with dimensions
 1869 $m \times n$. The matrix $\mathbf{D} \in \mathbb{R}^{m \times d}$ is constructed by randomly selecting m examples and d dimen-
 1870 sions from the original real-world dataset (<http://www.cad.zju.edu.cn/home/dengcai/Data/TextData.html>,
 1871 <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>). We
 1872 normalize each column of \mathbf{D} to have a unit norm and center the data by subtracting the mean.

1873 E.2 PROJECTION ON ORTHOGONALITY CONSTRAINTS

1875 When $h(\mathbf{x}) = \iota_{\mathcal{M}}(\text{mat}(\mathbf{x}))$ with $\Omega \triangleq \{\mathbf{V} | \mathbf{V}^T \mathbf{V} = \mathbf{I}\}$, computing the proximal operator reduces
 1876 to the following optimization problem:
 1877

$$\bar{\mathbf{x}} \in \arg \min_{\mathbf{x}} \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2, \text{ s.t. } \text{mat}(\mathbf{x}) \in \mathcal{M} \triangleq \{\mathbf{V} | \mathbf{V}^T \mathbf{V} = \mathbf{I}\}.$$

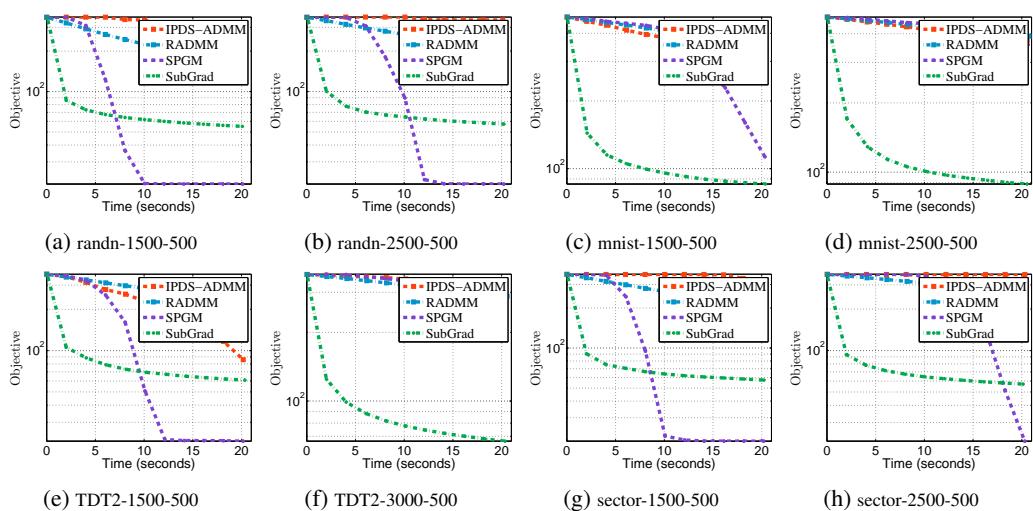
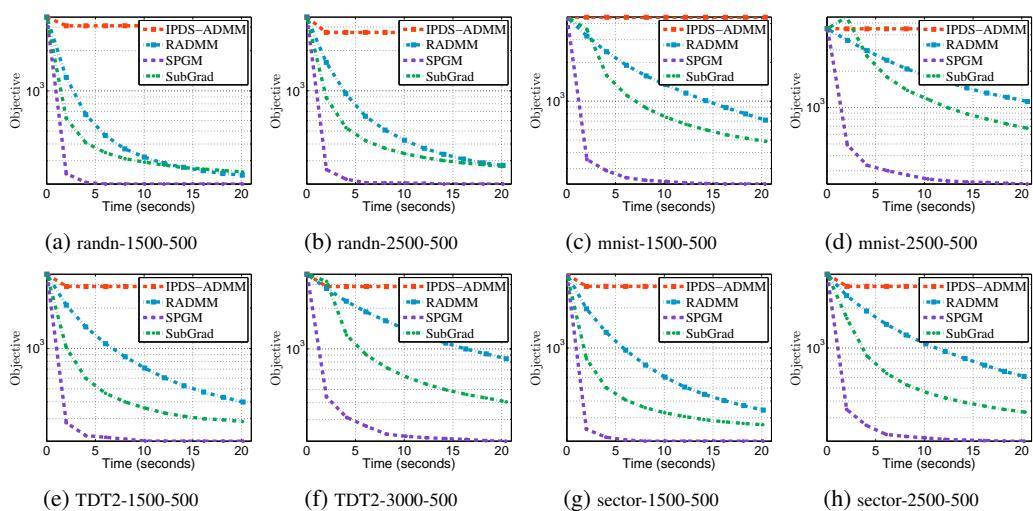
1879 This is the nearest orthogonality matrix problem, and the optimal solution can be computed as
 1880 $\bar{\mathbf{x}} = \text{vec}(\hat{\mathbf{U}}\hat{\mathbf{V}}^T)$, where $\text{mat}(\mathbf{x}') = \hat{\mathbf{U}}\text{Diag}(\mathbf{s})\hat{\mathbf{U}}^T$ is the singular value decompositon of the matrix
 1881 $\text{mat}(\mathbf{x}')$. Please refer to (Lai & Osher, 2014).

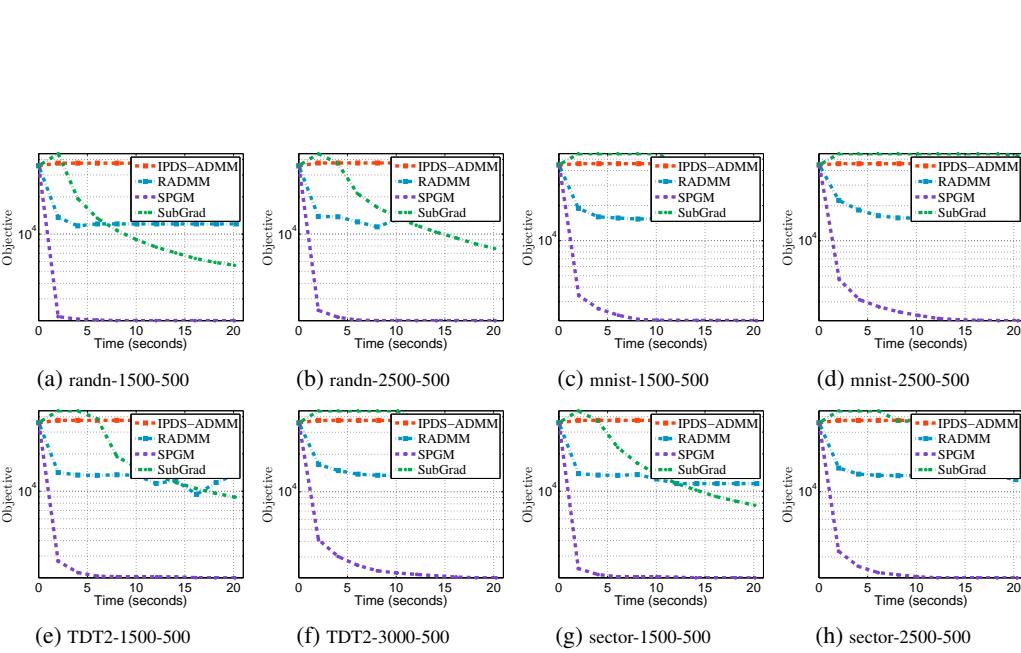
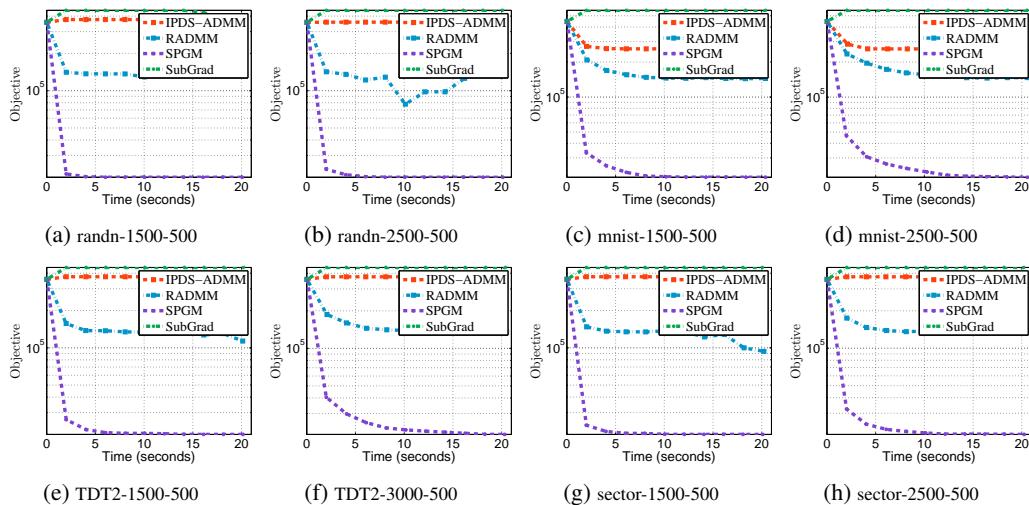
1883 E.3 ADDITIONAL EXPERIMENT RESULTS

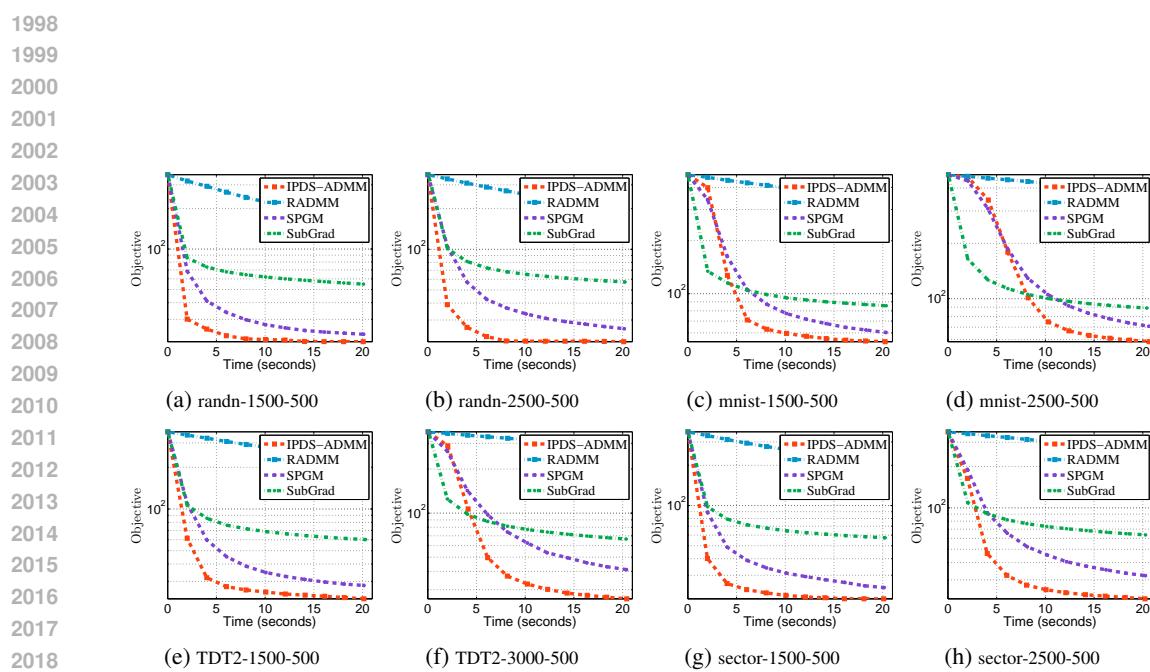
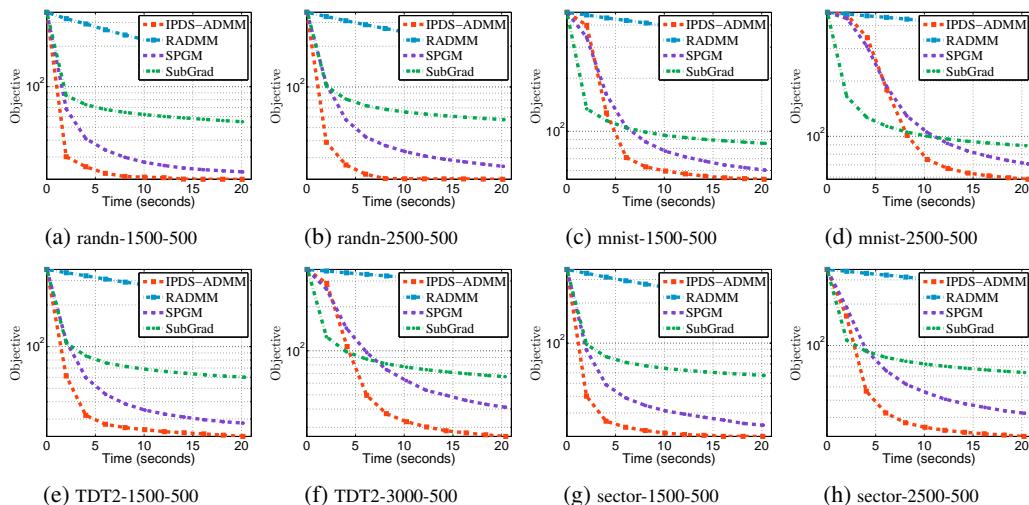
1885 We present the convergence curves of the compared methods for solving sparse PCA with varying
 1886 $\dot{\rho} = \{1, 10, 100, 1000\}$ and $\beta^0 = \{10\dot{\rho}, 50\dot{\rho}, 100\dot{\rho}, 500\dot{\rho}\}$, as shown in Figures 2 to 17. Please refer
 1887 to Table 2 for the mapping between $(\dot{\rho}, \beta^0)$ and the corresponding convergence curves. The results
 1888 demonstrate that the proposed IPDS-ADMM consistently outperforms other methods in terms of
 1889 speed for solving the sparse PCA problem, particularly for the ranges $\dot{\rho} = \{1, 10, 100, 1000\}$ and
 $\beta^0 = \{50\dot{\rho}, 100\dot{\rho}\}$.

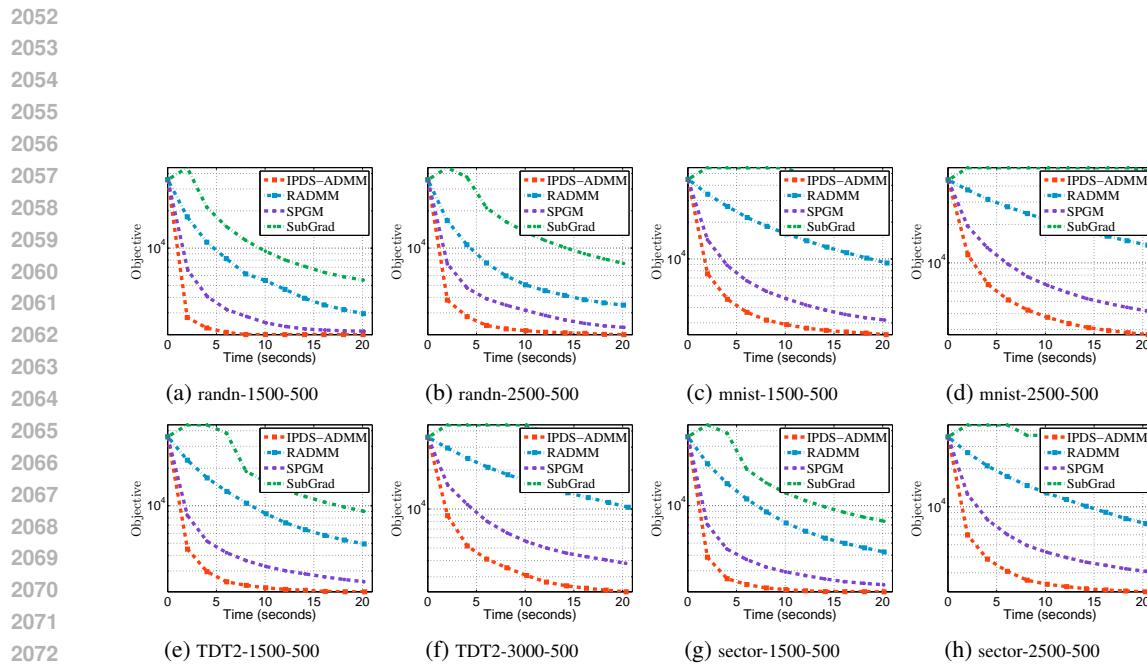
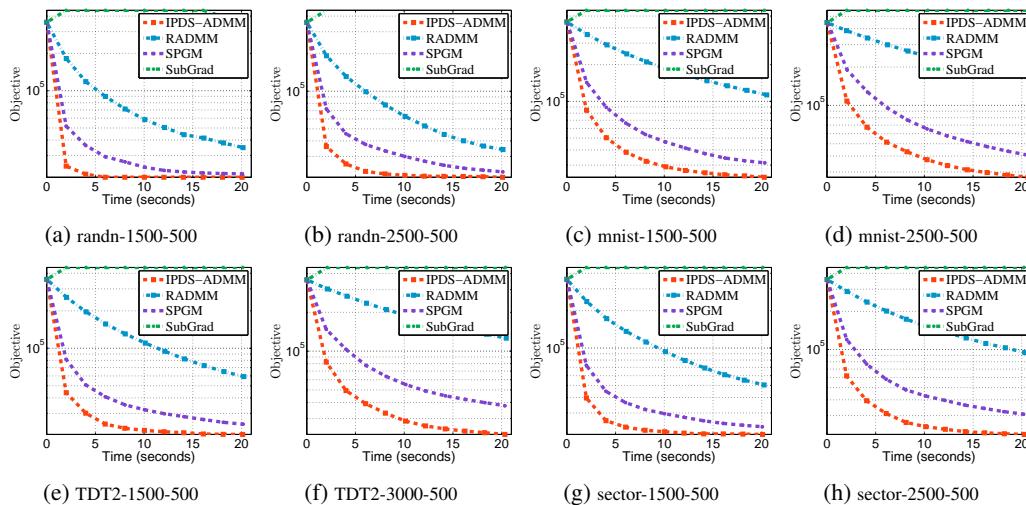
Table 2: The mapping between $(\dot{\rho}, \beta^0)$ and the corresponding convergence curves for sparse PCA.

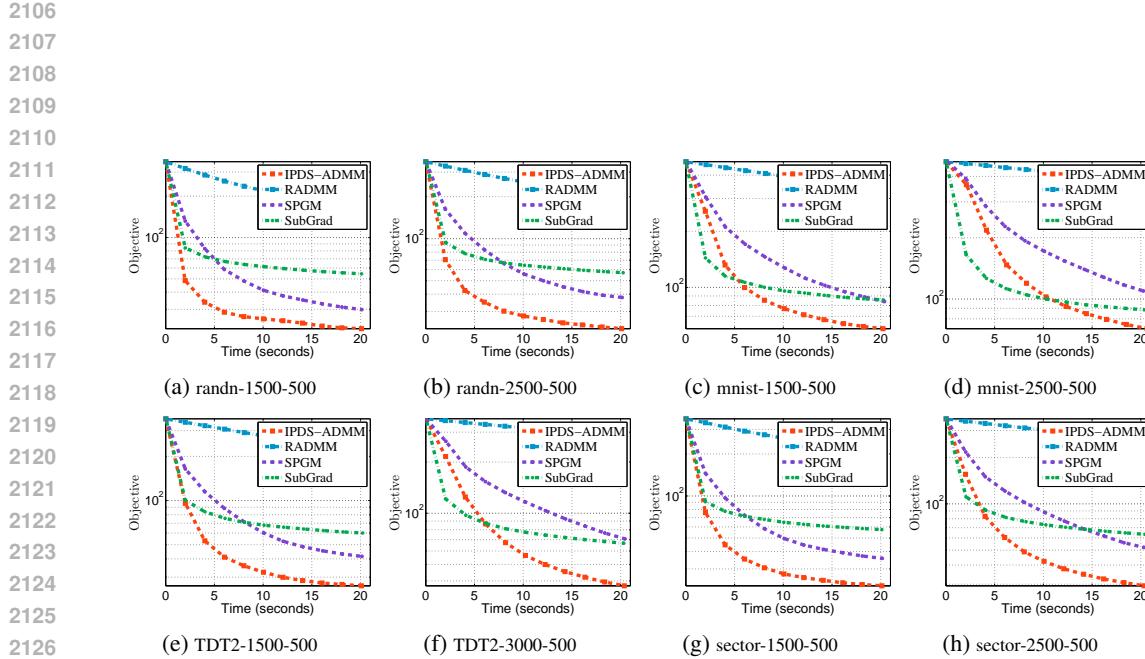
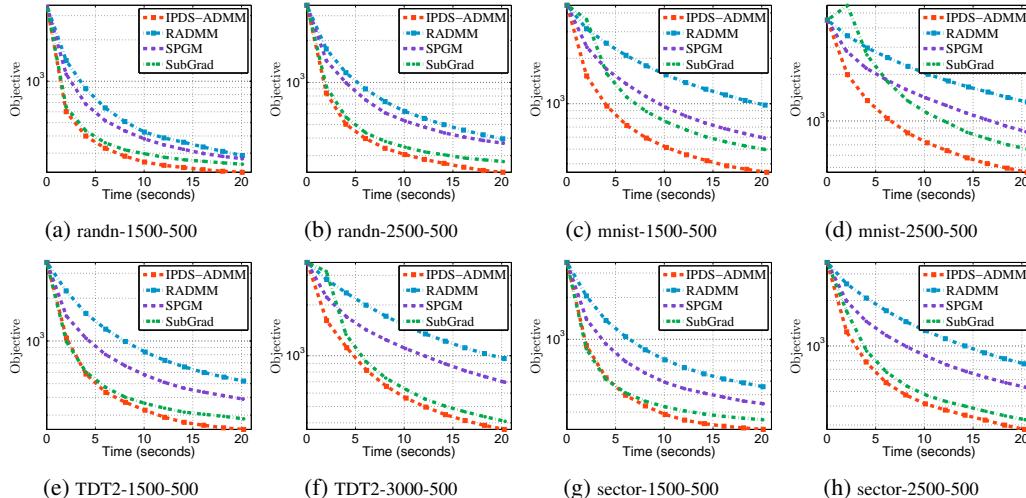
	$10\dot{\rho}$	$50\dot{\rho}$	$100\dot{\rho}$	$500\dot{\rho}$
$\dot{\rho} = 1$	Figure 2	Figure 6	Figure 10	Figure 14
$\dot{\rho} = 10$	Figure 3	Figure 7	Figure 11	Figure 15
$\dot{\rho} = 100$	Figure 4	Figure 8	Figure 12	Figure 16
$\dot{\rho} = 1000$	Figure 5	Figure 9	Figure 13	Figure 17

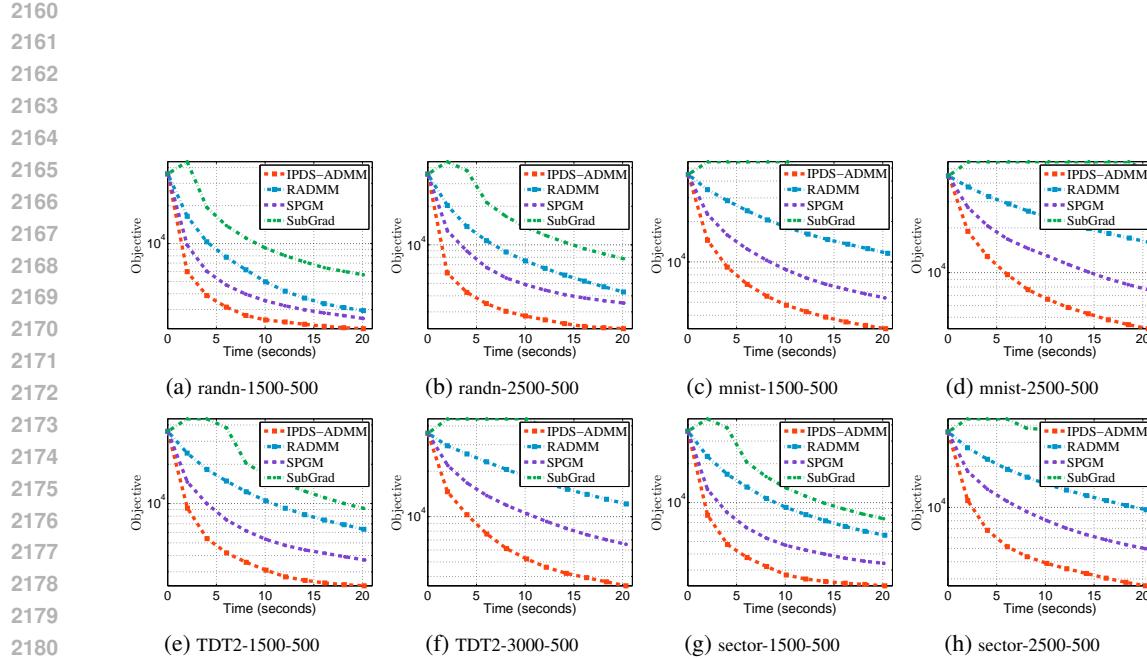
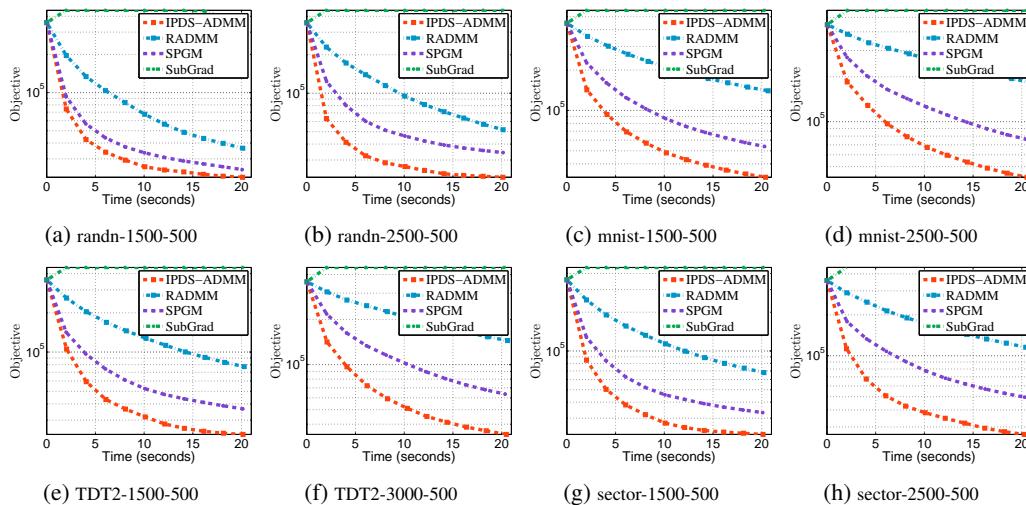
Figure 2: Convergence curves of methods for sparse PCA with $\dot{\rho} = 1$ and $\beta^0 = 10\dot{\rho}$.Figure 3: Convergence curves of methods for sparse PCA with $\dot{\rho} = 10$ and $\beta^0 = 10\dot{\rho}$.

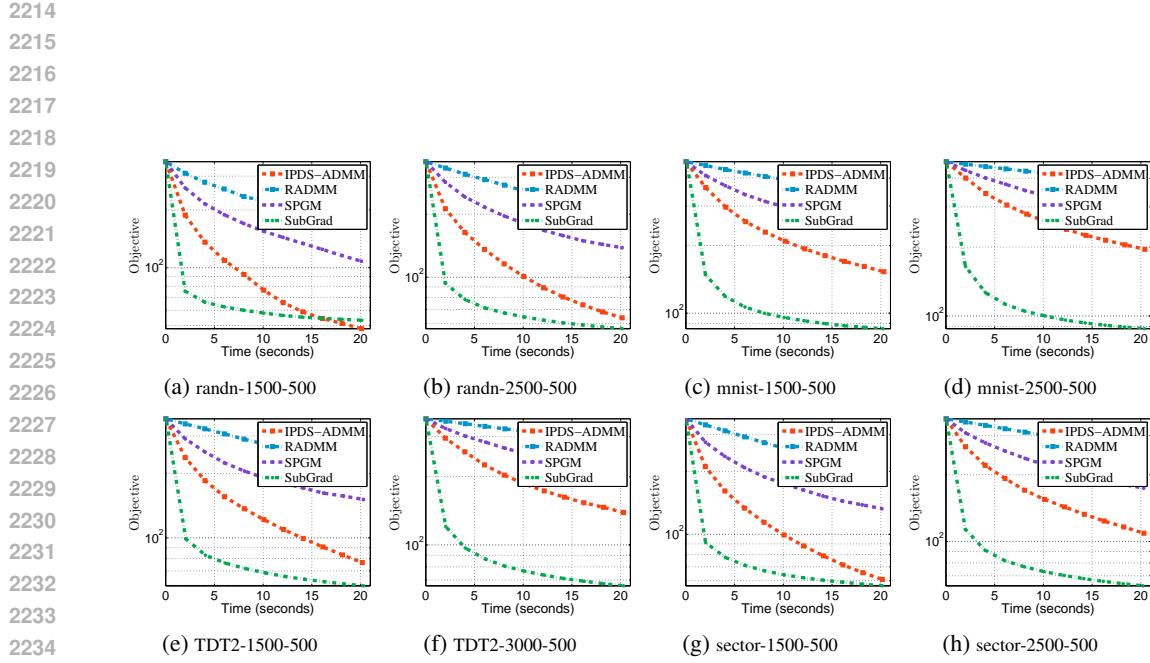
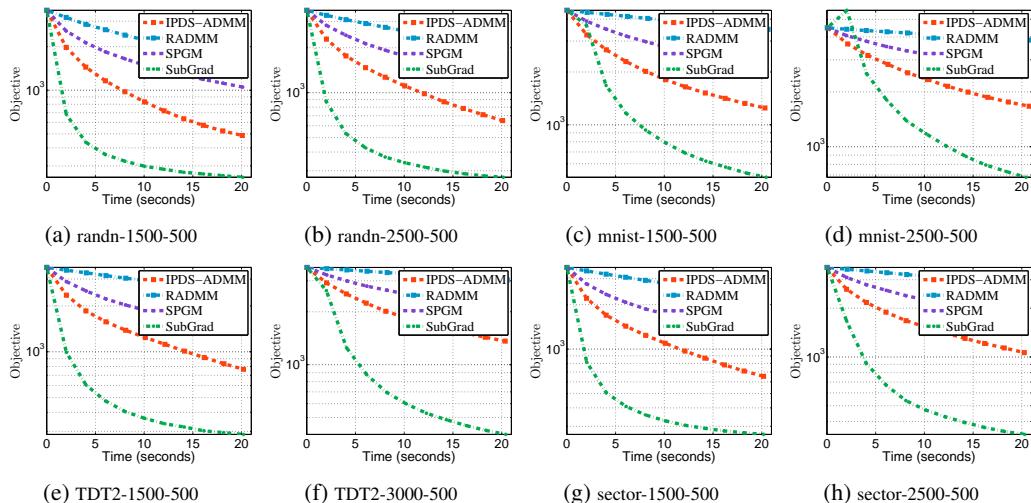
Figure 4: Convergence curves of methods for sparse PCA with $\dot{\rho} = 100$ and $\beta^0 = 10\dot{\rho}$.Figure 5: Convergence curves of methods for sparse PCA with $\dot{\rho} = 1000$ and $\beta^0 = 10\dot{\rho}$.

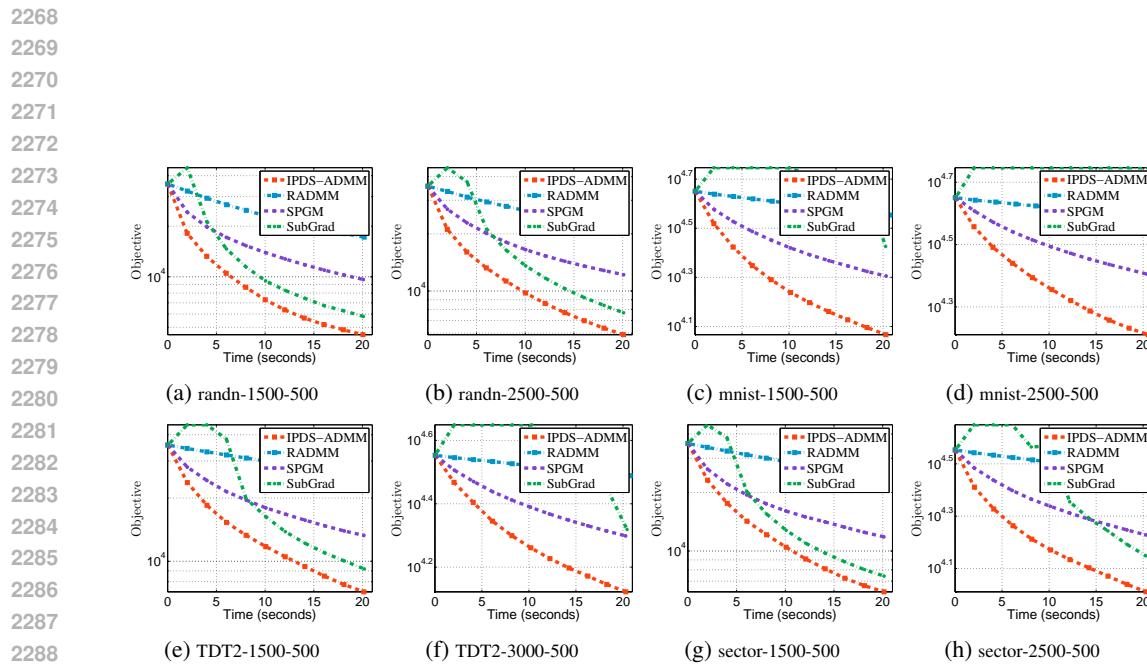
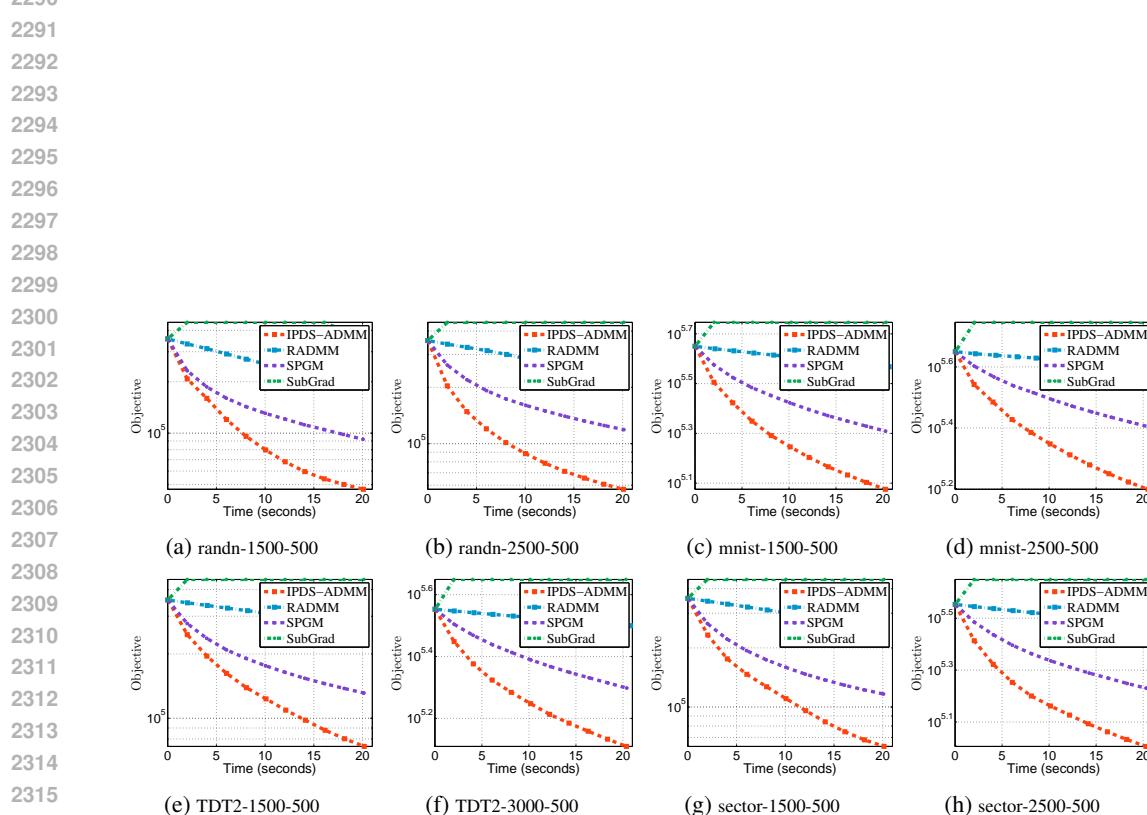
Figure 6: Convergence curves of methods for sparse PCA with $\hat{\rho} = 1$ and $\beta^0 = 50\hat{\rho}$.Figure 7: Convergence curves of methods for sparse PCA with $\hat{\rho} = 10$ and $\beta^0 = 50\hat{\rho}$.

Figure 8: Convergence curves of methods for sparse PCA with $\dot{\rho} = 100$ and $\beta^0 = 50\dot{\rho}$.Figure 9: Convergence curves of methods for sparse PCA with $\dot{\rho} = 1000$ and $\beta^0 = 50\dot{\rho}$.

Figure 10: Convergence curves of methods for sparse PCA with $\dot{\rho} = 1$ and $\beta^0 = 100\dot{\rho}$.Figure 11: Convergence curves of methods for sparse PCA with $\dot{\rho} = 10$ and $\beta^0 = 100\dot{\rho}$.

Figure 12: Convergence curves of methods for sparse PCA with $\dot{\rho} = 100$ and $\beta^0 = 100\dot{\rho}$.Figure 13: Convergence curves of methods for sparse PCA with $\dot{\rho} = 1000$ and $\beta^0 = 100\dot{\rho}$.

Figure 14: Convergence curves of methods for sparse PCA with $\hat{\rho} = 1$ and $\beta^0 = 500\hat{\rho}$.Figure 15: Convergence curves of methods for sparse PCA with $\hat{\rho} = 10$ and $\beta^0 = 500\hat{\rho}$.

Figure 16: Convergence curves of methods for sparse PCA with $\dot{\rho} = 100$ and $\beta^0 = 500\dot{\rho}$.Figure 17: Convergence curves of methods for sparse PCA with $\dot{\rho} = 1000$ and $\beta^0 = 500\dot{\rho}$.

2317
2318
2319
2320
2321