# A Communication-Efficient And Privacy-Aware Distributed Algorithm For Sparse PCA

Lei Wang[*]        Xin Liu[†]        Yin Zhang[‡]

**Abstract**

Sparse principal component analysis (PCA) improves interpretability of the classic PCA by introducing sparsity into the dimension-reduction process. Optimization models for sparse PCA, however, are generally non-convex, non-smooth and more difficult to solve, especially on large-scale datasets requiring distributed computation over a wide network. In this paper, we develop a distributed and centralized algorithm called DSSAL1 for sparse PCA that aims to achieve low communication overheads by adapting a newly proposed subspace-splitting strategy to accelerate convergence. Theoretically, convergence to stationary points is established for DSSAL1. Extensive numerical results show that DSSAL1 requires far fewer rounds of communication than state-of-the-art peer methods. In addition, we make the case that since messages exchanged in DSSAL1 are well-masked, the possibility of private-data leakage in DSSAL1 is much lower than in some other distributed algorithms.

## 1 Introduction

The principal component analysis (PCA) is a fundamental and ubiquitous tool in statistics and data analytics. In particular, it frequently serves as a critical preprocessing step in numerous data science or machine learning tasks. In general, solutions produced by the classical PCA are dense combinations of all features. However, in practice, sparse combinations not only enhance the interpretability of the principal components but also reduce storage [47, 9], which motivates the idea of sparse PCA. More importantly, from a theoretical perspective as given in [34, 63], sparse PCA is able to remediate some inconsistency phenomenon present in the classical PCA under the high-dimensional setting. As a dimension reduction and feature extraction method, sparse PCA can be widely applied to application areas where PCA is normally used; such as medical imaging [47], biological knowledge mining [54], ecology study [24], cancer analysis [9], neuroscience research [7], to mention a few examples.

Let $A = [a_1, a_2, \ldots, a_m]$ be an $n \times m$ data matrix, where $n$ and $m$ denote the numbers of features and samples, respectively. Without loss of generality, we assume that each row of $A$ has a zero mean. Mathematically speaking, PCA finds an orthonormal basis $Z \in \mathcal{S}_{n,p} := \{Z \in \mathbb{R}^{n \times p} \mid Z^\top Z = I_p\}$ of a $p$-dimensional subspace such that the projection of samples $a_1, a_2, \ldots, a_m$ on this subspace has the most variance. The set $\mathcal{S}_{n,p}$ is usually referred to as the Stiefel manifold [48]. Then PCA can be formally expressed as the following optimization problem with the orthogonality constraints.

$$\min_{Z \in \mathbb{R}^{n \times p}} \quad f(Z) := -\frac{1}{2} \mathrm{tr} \left( Z^\top A A^\top Z \right)$$
$$\text{s.t.} \quad Z \in \mathcal{S}_{n,p}, \tag{1}$$

where $\mathrm{tr}(\cdot)$ represents the trace of a given square matrix, and the column of $Z$ are called loading vectors or simply loadings.

In the projected data $Z^\top A \in \mathbb{R}^{p \times m}$, the number of features is reduced from $n$ to $p$ and each feature (row of $Z^\top A$) is a linear combination of the original features (rows of $A$) with coefficients from $Z$. For a sufficiently sparse $Z$, each reduced feature depends only on a few original features instead of all of them, leading to better interpretability in many applications. For this purpose, we consider the following $\ell_1$-regularized optimization model for sparse PCA:

$$\begin{aligned} \min_{Z \in \mathbb{R}^{n \times p}} \quad & \bar{f}(Z) := f(Z) + r(Z) \\ \mathrm{s.\,t.} \quad & Z \in \mathcal{S}_{n,p}, \end{aligned} \tag{2}$$

where the (non-operator) $\ell_1$-norm regularizer $r(Z) := \mu \|Z\|_1 = \mu \sum_{i,j} |[Z]_{ij}|$ is imposed to promote sparsity in $Z$, and $\mu > 0$ is the parameter used to control the amount of sparseness. Here, $[Z]_{ij}$ denotes the $(i,j)$-th entry of the matrix $Z$. Our distributed approach can efficiently tackle (2) with $p > 1$ by pursuing sparsity and orthogonality at the same time.

The optimization problem (2) is a penalized version of the SCoTLASS model proposed in [35]. Evidently, there is a significant difficulty gap in going from the standard PCA to sparse PCA in terms of both problem complexity and solution methodology. The standard PCA is polynomial-time solvable, while the sparse PCA is NP-hard if the $\ell_0$-regularizer is used to enforce sparsity [39], though it is still unclear what is the computational time complexity of solving the $\ell_1$-regularized model (2). In this paper we will be content with a theoretical result of convergence to first-order stationary points of (2). There are other formulations for sparse PCA, such as regression model [62], semidefinite programming [14, 15], and matrix decompositions [46, 53]. A comparative study of these formulations is beyond the scope of this paper. We refer interested readers to [63] for a recent survey of theoretical and computational results for sparse PCA.

## 1.1 Distributed setting

We consider the following distributed setting. The data matrix $A$ is divided into $d$ blocks, each containing a number of samples; namely, $A = [A_1, A_2, \ldots, A_d]$ where $A_i \in \mathbb{R}^{n \times m_i}$ so that $m_1 + \cdots + m_d = m$. This is a natural setting since each sample is a column of $A$. These submatrices $A_i$, $i = 1, \ldots, d$, are stored locally in $d$ locations, possibly having been collected at different locations by different agents, and all the agents are connected through a communication network. According to the splitting scheme of $A$, the function $f(Z)$ can also be distributed into $d$ agents, namely,

$$f(Z) = \sum_{i=1}^{d} f_i(Z) \quad \text{with} \quad f_i(Z) = -\frac{1}{2} \mathrm{tr}\left(Z^\top A_i A_i^\top Z\right). \tag{3}$$

In terms of network topology, we only need to assume that the network allows global summation operations (say, through all-reduce type of communications [43]) which are required by our algorithm. In particular, our algorithm will operate well in the federated-learning setting [41] where all the agents are connected to a center server so that global summations can be readily achieved in the network. To this extent, we say our algorithm is a centralized algorithm.

For most distributed algorithms, since the amount of communications per iteration remains essentially the same, the total communication overhead is proportional to the required number of iterations regardless of the underlying network topology. Therefore, our goal is to devise an algorithm that converges fast in terms of iteration counts, while considerations on other network topologies and communication patterns are beyond the scope of the current paper.

In certain applications, such as those in healthcare and financial industry [38, 60], preserving privacy of local data is of primary importance. In this paper, we consider the following scenario: each agent $i$ wants to keep its local dataset (i.e., $A_i A_i^\top$) from being discovered by any other agents including the center. In this situation, it is not an option to implement a pre-agreed encryption or a coordinated masking operation. For convenience, we will call such a privacy situation of *intrinsic privacy*. For an algorithm to preserve intrinsic privacy, publicly exchanged information must be safe so that none of the

local-data matrices can be figured out by any means. We will show later that, in general, algorithms based on traditional methods with distributed matrix multiplications are not intrinsically private.

## 1.2 Related works

Optimization problems with orthogonality constraints have been actively investigated in recent decades, for which many algorithms and solvers have been developed, such as, gradient approaches [40, 42, 1, 3], conjugate gradient approaches [16, 45, 61], constraint preserving updating schemes [52, 32], Newton methods [30, 29], trust-region methods [2], multipliers correction frameworks [21, 49], and orthonormalization-free approaches [22, 55]. These aforementioned algorithms are designed for smooth objective functions, and are generally not suitable for problem (2).

There exist some algorithms specifically designed for solving non-smooth optimization problems with orthogonality constraints; most of which adopt certain non-smooth optimization techniques to the Stiefel manifold; for instances, Riemannian subgradient methods [17, 18], proximal point algorithms [6], non-smooth trust-region methods [25], gradient sampling methods [28], and proximal gradient methods [11, 31]. Out of these algorithms, proximal gradient methods are among the most efficient. Specifically, Chen et al. [11] propose a manifold proximal gradient method (ManPG) and its accelerated version ManPG-Ada for non-smooth optimization problems over the Stiefel manifold. Starting from a current iterate, say $Z^{(k)} \in \mathcal{S}_{n,p}$, ManPG and ManPG-Ada generate the next iterate by solving the following subproblem restricted to the tangent space $\mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p} = \{D \in \mathbb{R}^{n \times p} \mid D^\top Z^{(k)} + (Z^{(k)})^\top D = 0\}$:

$$\min_{D \in \mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p}} \quad -\left\langle \sum_{i=1}^{d} A_i A_i^\top Z^{(k)}, D \right\rangle + \frac{1}{2\eta} \|D\|_{\mathrm{F}}^2 + r(Z^{(k)} + D), \tag{4}$$

which consumes the most computational time in the algorithm. In ManPG [11], the semi-smooth Newton (SSN) method [57] is deployed to solve the above subproblem, and global convergence to stationary points is established with the convergence rate $O(1/\sqrt{k})$. Huang and Wei [31] later extend the framework of ManPG to general Riemannian manifolds beyond the Stiefel manifold. Another class of algorithms first introduces auxiliary variables to split the objective function and the orthogonality constraint and then utilizes alternating minimization techniques to solve the resulting model, including SOC [37], MADMM [36], and PAMAL [12]. It is worth mentioning that SOC and MADMM lack a convergence guarantee. As for PAMAL, although convergence is guaranteed, its numerical performance is heavily dependent on its parameters as discussed in [11].

In principle, all the aforementioned algorithms can be adapted to the distributed setting considered in this paper. We take the algorithm ManPG as an example. Under the distributed setting, each agent computes the local product $A_i A_i^\top Z^{(k)}$ individually, then one round of communications will gather the global sum $\sum_{i=1}^{d} A_i A_i^\top Z^{(k)}$ at a center. The center then solves subproblem (4) and scatters the updated $Z^{(k+1)}$ back to all agents. At each iteration, distributed computation is basically limited to the matrix-multiplication level.

We point out that, without prior data-masking, distributed algorithms at the matrix-multiplication level cannot preserve local-data privacy intrinsically. Specifically, local data $A_i A_i^\top$, privately owned by agent $i$, can be uncovered by anyone who has access to publicly exchanged products $A_i A_i^\top Z^{(k)}$, after collecting enough such products and then solving a system of linear equations for the "unknown" $A_i A_i^\top$. This idea works even for more complicated iterative procedures as long as a mapping from data $A_i A_i^\top$ to a publicly exchanged quantity is deterministic and known.

To illustrate this data leakage situation, we now present a simple experiment on algorithm ManPG-Ada [11]. The test environment well be described in Section 5. As mentioned earlier, at iteration $k$ the publicly shared quantity in ManPG-Ada by agent $i$ is $S_i^{(k)} := A_i A_i^\top Z^{(k)}$, where $Z^{(k)}$ is the $k$-th iterate of ManPG-Ada accessible to all agents. For instance, to recover the unknown local data $A_1 A_1^\top$ we construct the following linear system with unknown $Y \in \mathbb{R}^{n \times n}$:

$$Y[Z^{(1)}, \cdots, Z^{(k)}] = [S_1^{(1)}, \cdots, S_1^{(k)}]. \tag{5}$$

Let $Y^{(k)}$ be the minimum norm least-squares solution to the linear equation (5) at iteration $k$. We perform an experiment to illustrate that $Y^{(k)}$ will quickly converge to $A_1 A_1^\top$, when ManPG-Ada is

deployed to solve the sparse PCA problem on a randomly generated test matrix with $n = 100$ and $m = 1280$ (other parameters are $d = 10$, $p = 10$, and $\mu = 0.05$). Figure 1 depicts how the reconstruction error $\|Y^{(k)} - A_1 A_1^\top\|_F$ and the stationarity violation[1] vary, on a logarithmic scale, as the number of iterations increases. We observe that local data $A_1 A_1^\top$ is obtained with high accuracy much faster than solving the sparse PCA problem.



Figure 1: Local data uncovered by solving linear systems during ManPG-Ada iterations.

In order to handle distributed datasets, some distributed ADMM-type algorithms have been introduced to PCA and related problems. These methods achieve algorithm-level parallelization, which are generally more secure than those based on matrix-multiplication-level parallelization. For instance, [26] proposes a distributed algorithm for sparse PCA with convergence to stationary points, but only studies the special case of $p = 1$. Moreover, under the distributed setting, [51] develop a subspace splitting strategy to tackle the smooth optimization problem over the Grassmann manifold without the $\ell_1$-norm regularizer in (2).

Recently, decentralized optimization has attracted attentions partly due to its wide applications in wireless sensor networks. Only local communications are carried out in decentralized optimization, namely, agents exchange information only with their immediate neighbors. This is quite different from the distributed setting considered in the current paper. We refer interested readers to the references [23, 19, 20, 4, 59, 10, 50] for some decentralized PCA algorithms.

## 1.3 Main contributions

Recently, a subspace splitting strategy has been proposed [51] to accelerate convergence of distributed algorithms for optimization problems over Grassmann manifolds where objective functions are orthogonal-transformation invariant[2] and smooth. In regularized problems such as sparse PCA, orthogonal-transformation invariance and smoothness no longer hold. In this paper, we present a non-trivial extension of the subspace splitting strategy to more general optimization problems over Stiefel manifolds. In particular, by incorporating the subspace splitting strategy into the framework of ManPG [11], we have developed a distributed algorithm DSSAL1 to solve the $\ell_1$-regularized optimization problem (2) for sparse PCA.

The main step of DSSAL1 is to solve the subproblem:

$$\min_{D \in \mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p}} \left\langle \sum_{i=1}^{d} Q_i^{(k)} Z^{(k)}, D \right\rangle + \frac{1}{2\eta} \|D\|_F^2 + r(Z^{(k)} + D), \tag{6}$$

---

[1] Suppose $D^{(k)}$ is the solution to (4). According to Lemma 5.3 in [11], $X^{(k)}$ is a first-order stationary point if $D^{(k)} = 0$. Therefore, the stationarity violation is defined as $\|D^{(k)}\|_F$.

[2] A function $f(X)$ is called orthogonal-transformation invariant if $f(XO) = f(X)$ for any $X \in \mathcal{S}_{n,p}$ and $O \in \mathcal{S}_{p,p}$.

which is identical to subproblem (4) in ManPG except that each local data matrix $A_i A_i^\top$ is replaced, or masked, by a matrix $-Q_i^{(k)}$ whose expression will be derived later.

In a sense, the main contributions of this paper can be attributed to this remarkably simple replacement or masking, which brings two great benefits. Firstly, convergence will be significantly accelerated which will be shown empirically in Section 5. Secondly, publicly exchanged products $Q_i^{(k)} Z^{(k)}$ are no longer images of some fixed and known mapping of $A_i A_i^\top$, making it practically impossible to uncover $A_i A_i^\top$ from such products via equation-solving.

Several innovations have been made in the development of our algorithms. Firstly, in our algorithm, only the global variable can possibly converge to a solution, while the local variables will never do. The role of the local variables, which are generally dense, is to help identify the right subspace. Secondly, we devise an inexact-solution strategy to effectively handle the difficult subproblem for the global variable where orthogonality and sparsity are pursued simultaneously. Thirdly, we establish the global convergence to stationarity for our algorithm, overcoming a number of technical difficulties associated with the rather complex algorithm construction, as well as the non-convexity and non-smoothness in problem (2).

## 1.4   Notations

The Euclidean inner product of two matrices $Y_1, Y_2$ of the same size is defined as $\langle Y_1, Y_2 \rangle = \mathrm{tr}(Y_1^\top Y_2)$. The Frobenius norm and spectral norm of a given matrix $C$ are denoted by $\|C\|_{\mathrm{F}}$ and $\|C\|_2$, respectively. Given a differentiable function $g(X) : \mathbb{R}^{n \times p} \to \mathbb{R}$, the gradient of $g$ with respect to $X$ is denoted by $\nabla g(X)$. And the subdifferential of a Lipschitz continuous function $h(X)$ is denoted by $\partial h(X)$. The tangent space to the Stiefel manifold $\mathcal{S}_{n,p}$ at $Z \in \mathcal{S}_{n,p}$, is represented by $\mathcal{T}_Z \mathcal{S}_{n,p} = \{Y \in \mathbb{R}^{n \times p} \mid Y^\top Z + Z^\top Y = 0\}$, and the orthogonal projection of $Y$ onto $\mathcal{T}_Z \mathcal{S}_{n,p}$ is denoted by $\mathrm{Proj}_{\mathcal{T}_Z \mathcal{S}_{n,p}} (Y) = \left(I_n - ZZ^\top\right) Y + Z \left(Z^\top Y - Y^\top Z\right)/2$. For $X \in \mathcal{S}_{n,p}$, we define $\mathbf{P}_X^\perp := I_n - XX^\top$ standing for the projection operator onto the null space of $X^\top$. Further notation will be introduced as it occurs.

## 1.5   Organization

The rest of this paper is organized as follows. In Section 2, we introduce a novel subspace-splitting model for sparse PCA, and investigate the structure of associated Lagrangian multipliers. Then a distributed algorithm is proposed to solve this model based on an ADMM-like framework in Section 3. Moreover, convergence properties of the proposed algorithm are studied in Section 4. Numerical experiments on a variety of test problems are presented in Section 5 to evaluate the performance of the proposed algorithm. We draw final conclusions and discuss some possible future developments in the last section.

# 2   Subspace-splitting model for sparse PCA

We consider the scenario that the $i$-th component function $f_i(X)$ of $f$ in (3) can be evaluated only by the $i$-th agent since local dadaset $A_i$ is accessible only to the $i$-th agent. In order to devise a distributed algorithm, the classic variable-splitting approach is to introduce a set of local variables $\{X_i\}$ to make the sum of local functions nominally separable. Then a (centralized) distributed algorithm would maintain a global variable $Z$ and impose variable-consensus constraints $X_i = Z$.

Despite the regularizer term $r$ in (2), all component functions $f_i(X)$ in $f$ are still invariant under orthogonal transformations. It should be natural for us to adapt the subspace-splitting idea introduced in [51], that is, to use the subspace-consensus constraints $X_i X_i^\top = ZZ^\top$ to accelerate convergence. In

this paper, we propose to solve the following optimization problem:

$$\min_{X_i, Z \in \mathbb{R}^{n \times p}} \quad \sum_{i=1}^{d} f_i(X_i) + r(Z) \tag{7a}$$

$$\text{s.t.} \qquad X_i^\top X_i = I_p, \qquad i = 1, \ldots, d, \tag{7b}$$

$$X_i X_i^\top = Z Z^\top, \quad i = 1, \ldots, d, \tag{7c}$$

$$Z^\top Z = I_p, \tag{7d}$$

which we will call the *subspace-splitting* model for problem (2), noting that both sides of (7c) are orthogonal projections onto subspaces. For brevity, we collect the global variable $Z$ and all local variables $\{X_i\}$ into a point $(Z, \{X_i\})$. A point $(Z, \{X_i\})$ is feasible if it satisfies the constraints (7b)-(7d).

## 2.1 Stationarity conditions

In this subsection, we aim to present the stationarity conditions of the sparse PCA problem (2). We first introduce the definition of Clarke subgradient [13] for non-smooth functions.

**Definition 2.1.** *Suppose $f : \mathbb{R}^{n \times p} \to \mathbb{R}$ is a Lipschitz continuous function. The generalized directional derivative of $f$ at the point $X \in \mathbb{R}^{n \times p}$ along the direction $H \in \mathbb{R}^{n \times p}$ is defined by:*

$$f^\circ(X; H) := \limsup_{Y \to X, \, t \to 0^+} \frac{f(Y + tH) - f(Y)}{t}.$$

*Based on generalized directional derivative of $f$, the (Clark) subgradient of $f$ is defined by:*

$$\partial f(X) := \{G \in \mathbb{R}^{n \times p} \mid \langle G, H \rangle \leq f^\circ(X; H)\}.$$

As discussed in [58, 11], the first-order stationarity condition of (2) can be stated as:

$$0 \in \mathrm{Proj}_{\mathcal{T}_Z \mathcal{S}_{n,p}} \left( -AA^\top Z + \partial r(Z) \right).$$

We provide an equivalent description of the above first-order stationarity condition, which will be used in the theoretical analysis.

**Lemma 2.2.** *A point $Z \in \mathcal{S}_{n,p}$ is a first-order stationary point of (2) if and only if there exists $R(Z) \in \partial r(Z)$ such that the following conditions hold:*

$$\begin{cases} \mathbf{P}_Z^\perp \left( -AA^\top Z + R(Z) \right) = 0, \\ Z^\top R(Z) - R(Z)^\top Z = 0. \end{cases} \tag{8}$$

The proof of Lemma 2.2 is put into Appendix A. Then we can characterize the first-order stationary points of (7) in the following manner.

**Definition 2.3.** *Suppose $X_i \in \mathbb{R}^{n \times p}(i = 1, \ldots, d)$ and $Z \in \mathbb{R}^{n \times p}$. A point $(Z, \{X_i\})$ is called a first-order stationary point of (7) if it is feasible and $Z$ satisfies the conditions in (8).*

## 2.2 Existence of low-rank multipliers

By associating dual variables $\Gamma_i$, $\Lambda_i$, and $\Theta$ to the constraints (7b), (7c), and (7d), respectively, we derive an equivalent description of first-order stationarity conditions of (7).

**Proposition 2.4.** *Suppose $X_i \in \mathbb{R}^{n \times p}(i = 1, \ldots, d)$ and $Z \in \mathbb{R}^{n \times p}$. A point $(Z, \{X_i\})$ is a first-order stationary point of (7) if and only if there exist symmetric matrices $\Lambda_i \in \mathbb{R}^{n \times n}$, $\Gamma_i \in \mathbb{R}^{p \times p}$, and*

$\Theta \in \mathbb{R}^{p \times p}$ such that $(Z, \{X_i\})$ satisfies the following condition:

$$
\begin{cases}
0 = A_i A_i^\top X_i + X_i \Gamma_i + \Lambda_i X_i, & i = 1, \dots, d, \\
0 \in \partial r(Z) + \displaystyle\sum_{i=1}^d \Lambda_i Z - Z\Theta, \\
0 = X_i^\top X_i - I_p, & i = 1, \dots, d, \\
0 = X_i X_i^\top - ZZ^\top, & i = 1, \dots, d, \\
0 = Z^\top Z - I_p.
\end{cases} \tag{9}
$$

The proof of Proposition 2.4 is relegated to Appendix B. Actually, the equations in (9) can be viewed as the KKT conditions of (7), while $\Gamma_i \in \mathbb{R}^{p \times p}$, $\Lambda_i \in \mathbb{R}^{n \times n}$, and $\Theta \in \mathbb{R}^{p \times p}$ are the Lagrangian multipliers corresponding to the constraints (7b), (7c), and (7d), respectively.

In [51], a low-rank and closed-form multiplier is devised with respect to the subspace constraint (7c) for the PCA problems, which can be expressed as:

$$
\Lambda_i = -X_i X_i^\top A_i A_i^\top \mathbf{P}_{X_i}^\perp - \mathbf{P}_{X_i}^\perp A_i A_i^\top X_i X_i^\top \ (i = 1, \dots, d). \tag{10}
$$

In Appendix (B), we further verify that the above formulation is also valid for the sparse PCA problems at any first-order stationary point $(Z, \{X_i\})$. In the next section, we will use (10) to update the multipliers in our framework of ADMM. This strategy simultaneously saves computational costs and storage requirements.

# 3    Algorithmic framework

Now we describe the proposed algorithm, based on an ADMM-like framework, to solve the subspace-splitting model (7). Note that there are three constraints (7b)-(7d). We only penalize the subspace constraints (7c) to the objective function, and obtain the corresponding augmented Lagrangian function:

$$
\mathcal{L}(Z, \{X_i\}, \{\Lambda_i\}) = \sum_{i=1}^d \mathcal{L}_i(Z, X_i, \Lambda_i) + r(Z), \tag{11}
$$

where

$$
\begin{aligned}
\mathcal{L}_i(Z, X_i, \Lambda_i) = &-\frac{1}{2}\mathrm{tr}\left(X_i^\top A_i A_i^\top X_i\right) - \frac{1}{2}\left\langle \Lambda_i, X_i X_i^\top - ZZ^\top \right\rangle \\
&+ \frac{\beta_i}{4}\left\| X_i X_i^\top - ZZ^\top \right\|_{\mathrm{F}}^2,
\end{aligned}
$$

and $\beta_i > 0 (i = 1, \dots, d)$ are penalty parameters. Schematically, we will follow the ADMM-like framework below to build an algorithm for solving subspace-splitting model (7), though we quickly add that the two optimization subproblems below in (12) and (13) are not "solved" in a normal sense since we may stop after a single iteration of an iterative scheme.

$$
\begin{cases}
Z^{(k+1)} \approx \underset{Z \in \mathcal{S}_{n,p}}{\arg\min} \ \mathcal{L}(Z, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}). & (12) \\
X_i^{(k+1)} \approx \underset{X_i \in \mathcal{S}_{n,p}}{\arg\min} \ \mathcal{L}_i(Z^{(k+1)}, X_i, \Lambda_i^{(k)}), \ \ i = 1, \dots, d. & (13) \\
W_i^{(k+1)} = -\mathbf{P}_{X_i^{(k+1)}}^\perp A_i A_i^\top X_i^{(k+1)}, \ \ i = 1, \dots, d. & (14) \\
\Lambda_i^{(k+1)} = X_i^{(k+1)}(W_i^{(k+1)})^\top + W_i^{(k+1)}(X_i^{(k+1)})^\top, \ \ i = 1, \dots, d. & (15)
\end{cases}
$$

In the above framework, the superscript $(k)$ counts the number of iterations, and the subscript $i$ indicates the number of agents.

A novel feature in our algorithm is the way to update the multipliers $\Lambda_i$ associated with the subspace constraints (7c). Generally speaking, in the augment Lagrangian based approach, the multipliers are updated by the dual ascent step

$$\Lambda_i^{(k+1)} = \Lambda_i^{(k)} - \tau\beta_i \left( X_i^{(k+1)}(X_i^{(k+1)})^\top - Z^{(k+1)}(Z^{(k+1)})^\top \right),$$

where $\tau > 0$ is the step size. This standard method would require to store and update an $n \times n$ matrix at each agent, which could be costly when $n$ is large. Instead, we use the low-rank updating formulas (14)-(15) based on the closed-form expression (10) derived in Section 2.2. In our iterative setting, these multiplier matrices are never stored but used in matrix multiplications, in which the required additional storage for agent $i$ is for the $n \times p$ matrix $W_i$.

## 3.1 Subproblem for global variable

We now describe how to approximate subproblems (12) and (13). By rearrangments, subproblem (12) reduces to:

$$\min_{Z \in \mathcal{S}_{n,p}} \quad q^{(k)}(Z) := \frac{1}{2}\mathrm{tr}\left( Z^\top Q^{(k)} Z \right) + r(Z) \tag{16}$$

where $Q^{(k)} = \sum_{i=1}^d Q_i^{(k)}$ and $Q_i^{(k)}$ is an $n \times n$ matrix defined by

$$Q_i^{(k)} = \Lambda_i^{(k)} - \beta_i X_i^{(k)}(X_i^{(k)})^\top. \tag{17}$$

We quickly point out here that it is not necessary to construct and store these $Q$-matrices explicitly since we will only use them to multiply $n \times p$ matrices in an iterative scheme to obtain approximate solutions to subproblem (16). More details will follow later.

Apparently, subproblem (16) pursues the orthogonality and sparsity simultaneously, and is not easier to solve than the the original problem (2). However, we will demonstrate later by both theoretical analysis and numerical experiments that inexactly solving (16) by conducting one proximal gradient step is adequate for the global convergence.

Starting from the current iterate $Z^{(k)}$, we first find a decent direction $D^{(k)}$ restricted to the tangent space $\mathcal{T}_{Z^{(k)}}\mathcal{S}_{n,p}$ by solving the following subproblem

$$\begin{aligned} \min_{D \in \mathbb{R}^{n \times p}} \quad & \left\langle Q^{(k)} Z^{(k)}, D \right\rangle + \frac{1}{2\eta}\|D\|_{\mathrm{F}}^2 + r(Z^{(k)} + D) \\ \text{s.t.} \quad & D^\top Z^{(k)} + (Z^{(k)})^\top D = 0, \end{aligned} \tag{18}$$

where $\eta > 0$ is the step size. Since $Z^{(k)} + D^{(k)}$ does not necessarily lie on the Stiefel manifold $\mathcal{S}_{n,p}$, we then perform a projection to bring it back to $\mathcal{S}_{n,p}$, which can be represented as:

$$Z^{(k+1)} = \mathrm{Proj}_{\mathcal{S}_{n,p}}\left( Z^{(k)} + D^{(k)} \right).$$

Here, the orthogonal projection of a matrix $C \in \mathbb{R}^{n \times p}$ onto $\mathcal{S}_{n,p}$ is denoted by $\mathrm{Proj}_{\mathcal{S}_{n,p}}(C) = U_C V_C^\top$, where $U_C \Sigma_C V_C^\top$ is the economic form of the singular value decomposition of $C$.

**Remark 1.** *We note that the subproblems solved by ManPG [11] are identical in form to our subproblem (18) but with $Q^{(k)}$ replaced by the data matrix $AA^\top$. That is, ManPG applies manifold proximal gradient steps to a fixed problem, while our algorithm computes steps of the same type using a sequence of matrices, each being updated to incorporate latest information. From this point of view, our algorithm can be interpreted as an acceleration scheme to reduce the number of outer-iterations, thereby reducing communication overheads.*

Since these $n \times n$ matrices $Q_i^{(k)}(i = 1, \ldots, d)$ are distributively maintained in $d$ agents, each agent is not able to independently solve subproblem (18). Fortunately, we only need to calculate

$$Q^{(k)} Z^{(k)} = \sum_{i=1}^d Q_i^{(k)} Z^{(k)}, \tag{19}$$

to solve this subproblem. In the distributed setting, the right hand side of (19) can be accomplished by calculating $Q_i^{(k)} Z^{(k)}$ in each agent and then invoking the all-reduce type of communication, where each agent just shares one $n \times p$ matrix. If the all-reduce type of communication is realized by the butterfly algorithm [43], the communication overhead per iteration is $O\left(np \log(d)\right)$. Furthermore, each local product $Q_i^{(k)} Z^{(k)}$ can be computed from

$$
\begin{aligned}
Q_i^{(k)} Z^{(k)} &= \Lambda_i^{(k)} Z^{(k)} - \beta_i X_i^{(k)} (X_i^{(k)})^\top Z^{(k)} \\
&= X_i^{(k)} (W_i^{(k)})^\top Z^{(k)} + W_i^{(k)} (X_i^{(k)})^\top Z^{(k)} \\
&\quad - \beta_i X_i^{(k)} (X_i^{(k)})^\top Z^{(k)},
\end{aligned}
\tag{20}
$$

with a computational cost in the order of $O(np^2)$. From the above formula, one observes that the $n \times n$ matrices $Q_i$ need not be stored explicitly.

Now we consider how to efficiently solve subproblem (18). By associating a multiplier $\Upsilon \in \mathbb{R}^{p \times p}$ to the linear equality constraint, the Lagrangian function of (18) can be written as:

$$
\begin{aligned}
\mathfrak{L}(D, \Upsilon) =& \left\langle Q^{(k)} Z^{(k)}, D \right\rangle + \frac{1}{2\eta} \|D\|_{\mathrm{F}}^2 + r(Z^{(k)} + D) \\
& - \frac{1}{2} \left\langle \Upsilon, D^\top Z^{(k)} + (Z^{(k)})^\top D \right\rangle.
\end{aligned}
$$

We choose to apply the Uzawa method [5] to solve subproblem (18). At the $j$-th inner iteration, we first minimize the above Lagrangian function with respect to $D$ for a fixed $\Upsilon = \Upsilon(j)$:

$$
\begin{aligned}
D(j+1) =& \arg\min_{D \in \mathbb{R}^{n \times p}} \mathfrak{L}(D, \Upsilon(j)) \\
=& \operatorname{Prox}_{\eta r} \left( Z^{(k)} - \eta \left( Q^{(k)} Z^{(k)} - Z^{(k)} \Upsilon(j) \right) \right) - Z^{(k)},
\end{aligned}
\tag{21}
$$

where $D(j)$ and $\Upsilon(j)$ denote the $j$-th inner iterate of $D$ and $\Upsilon$, respectively. Here, we use $\operatorname{Prox}_g (X)$ to denote the proximal mapping of a given function $g : \mathbb{R}^{n \times p} \to \mathbb{R}$ at the point $X \in \mathbb{R}^{n \times p}$, which is defined by:

$$
\operatorname{Prox}_g (X) = \arg\min_{Y \in \mathbb{R}^{n \times p}} \ g(Y) + \frac{1}{2} \|Y - X\|_{\mathrm{F}}^2.
$$

For the $\ell_1$-norm regularizer term $r(X) = \mu \|X\|_1$, the proximal mapping in (21) admits a closed-form solution:

$$
[\operatorname{Prox}_{\eta r} (X)]_{ij} = \begin{cases} [X]_{ij} - \eta\mu, & \text{if } [X]_{ij} > \eta\mu, \\ 0, & \text{if } -\eta\mu \le [X]_{ij} \le \eta\mu, \\ [X]_{ij} + \eta\mu, & \text{if } [X]_{ij} < -\eta\mu, \end{cases}
$$

where the subscript $[\cdot]_{ij}$ represents the $(i,j)$-th entry of a matrix. Then the multiplier is updated by a dual ascent step:

$$
\Upsilon(j+1) = \Upsilon(j) - \tau \left( D(j+1)^\top Z^{(k)} + (Z^{(k)})^\top D(j+1) \right),
\tag{22}
$$

where $\tau > 0$ is the step size. These two steps are repeated until convergence. The complete framework is summarized in Algorithm 1.

The Uzawa method can be viewed as a special case of the primal-dual hybrid gradient algorithm (PDHG) developed in [27] with the convergence rate $O(1/k)$ in the ergodic sense under mild conditions. Moreover, as an inner solver it can bring a higher overall efficiency than the SSN method used by ManPG [11] (see [56] for a recent study).

## 3.2 Subproblems for local variables

In this subsection, we focus on for the $i$-th local variable $X_i$, which can be rearranged as the following equivalent problem:

$$
\min_{X_i \in \mathcal{S}_{n,p}} \ h_i^{(k)}(X_i) := -\frac{1}{2} \operatorname{tr}\left( X_i^\top H_i^{(k)} X_i \right).
\tag{23}
$$

---

**Algorithm 1:** Uzawa method for subproblem (18).

---
**1 Input:** $Z^{(k)}$, $Q^{(k)}Z^{(k)}$, and $\eta$ in subproblem (18).
**2** Set $j := 0$, and choose the step size $\tau > 0$ as well as the initial variable $\Upsilon(0)$.
**3 while** *not converged* **do**
**4** $\quad$ Compute $D(j+1)$ by (21).
**5** $\quad$ Update $\Upsilon(j+1)$ by (22).
**6** $\quad$ Set $j := j + 1$.
**7 Output:** $D(j)$.

---

Here, $H_i^{(k)}$ is an $n \times n$ real symmetric matrix:

$$H_i^{(k)} = A_i A_i^\top + \Lambda_i^{(k)} + \beta_i Z^{(k+1)}(Z^{(k+1)})^\top, \tag{24}$$

which is only related to the local data $A_i$. This is a standard eigenvalue problem where one needs to compute a $p$-dimensional dominant eigenspace of $H_i^{(k)}$.

As a subproblem, it is not necessary to solve (23) to high precision. In practice, we just need to find a point $X_i^{(k+1)} \in \mathcal{S}_{n,p}$ satisfying the following two conditions, which suffices to be a good inexact solution empirically, and to guarantee the global convergence of the whole algorithm. The first condition demands a sufficient decrease in function value:

$$h_i^{(k)}\left(X_i^{(k)}\right) - h_i^{(k)}\left(X_i^{(k+1)}\right) \geq \frac{c_i}{c_i' \|A_i\|_2^2 + \beta_i} \left\| \mathbf{P}_{X_i^{(k)}}^\perp H_i^{(k)} X_i^{(k)} \right\|_{\mathrm{F}}^2, \tag{25}$$

where $c_i > 0$ and $c_i' > 0$ are two constants independent of $\beta_i$. The second condition is a sufficient decrease in KKT violation:

$$\left\| \mathbf{P}_{X_i^{(k+1)}}^\perp H_i^{(k)} X_i^{(k+1)} \right\|_{\mathrm{F}} \leq \delta_i \left\| \mathbf{P}_{X_i^{(k)}}^\perp H_i^{(k)} X_i^{(k)} \right\|_{\mathrm{F}}, \tag{26}$$

where $\delta_i \in [0,1)$ is a constant independent of $\beta_i$. It turns out that these two rather weak termination conditions for subproblem (23) are sufficient for us to derive global convergence of our ADMM-like algorithm framework (12)-(15).

In practice, we can combine a warm-start strategy with a single iteration of SSI [44] to generate the next iterate $X_i^{(k+1)} = \mathrm{Proj}_{\mathcal{S}_{n,p}}\left(H_i^{(k)} X_i^{(k)}\right)$, i.e.,

$$\begin{aligned}
X_i^{(k+1)} &= \mathrm{Proj}_{\mathcal{S}_{n,p}}\left(A_i A_i^\top X_i^{(k)} + \Lambda_i^{(k)} X_i^{(k)} + \beta_i Z^{(k+1)}(Z^{(k+1)})^\top X_i^{(k)}\right) \\
&= \mathrm{Proj}_{\mathcal{S}_{n,p}}\left(A_i A_i^\top X_i^{(k)} + W_i^{(k)} + \beta_i Z^{(k+1)}(Z^{(k+1)})^\top X_i^{(k)}\right),
\end{aligned} \tag{27}$$

which can be computed in the order of $O(np^2)$ floating-point operations, given that the term $A_i A_i^\top X_i^{(k)}$ is inherited from the last iteration as a result of updating $W_i^{(k)}$, see (14).

## 3.3 Algorithm description

We formally present the detailed algorithmic framework as Algorithm 2 below, named *distributed subspace splitting algorithm with $\ell_1$ regularization* and abbreviated to DSSAL1. In the distributed environment, all agents are initiated from the same point $Z^{(0)} \in \mathcal{S}_{n,p}$. And the initial guess of multipliers are computed by (15). After initialization, all agents first solve the common subproblem for $Z$ collaboratively by certain communication strategy. Then each agent solves its subproblem for $X_i$ and updates its multiplier $\Lambda_i$. These two steps only involve the local data privately stored at each agent, and hence can be carried out in $d$ agents concurrently. This procedure is repeated until convergence.

---

**Algorithm 2:** Distributed Subspace Splitting Algorithm with $\ell_1$ regularization (DSSAL1).

---

**1 Input:** functions $f_i(i = 1, \ldots, d)$ and $r$.

**2** Set $k := 0$, choose penalty parameters $\{\beta_i\}$, and initialize $Z^{(0)}$.

**3** Set $X_i^{(0)} = Z^{(0)}$ for $i = 1, \ldots, d$.

**4** Compute the initial multipliers $\{\Lambda_i^{(0)}\}$ by (15).

**5 while** *not converged* **do**

**6**   Solve (18) to obtain $D^{(k)}$ by Algorithm 1.

**7**   Set $Z^{(k+1)} = \mathrm{Proj}_{\mathcal{S}_{n,p}}\left(Z^{(k)} + D^{(k)}\right)$.

**8**   **for** $i = 1, \ldots, d$ **do**

**9**     Find $X_i^{(k+1)} \in \mathcal{S}_{n,p}$ satisfying (25) and (26).

**10**     Update the multipliers $\Lambda_i^{(k+1)}$ by (15).

**11**   Set $k := k + 1$.

**12 Output:** $Z^{(k)}$.

---

## 3.4 Data privacy

We claim that DSSAL1 can naturally protect the intrinsic privacy of local data. To form the global sum in (19), the shared information in DSSAL1 at iteration $k$ from the $i$-th agent is $S_i^{(k)} := Q_i^{(k)} Z^{(k)}$ where $Z^{(k)}$ is known to all agents. However, the $n \times p$ system of equations, $Q_i^{(k)} Z^{(k)} = S_i^{(k)}$, is insufficient for obtaining the $n \times n$ mask matrix $Q_i^{(k)}$ which changes from iteration to iteration. Secondly, even if a few mask matrices $Q_i^{(k)}$ were unveiled, it would still be impossible to derive the local data matrix $A_i A_i^\top$ from these $Q_i^{(k)}$ without knowing corresponding $X_i^{(k)}$ (and $\beta_i$) which are always kept privately by the $i$-th agent. Finally, consider the ideal "converged" case where $X_i X_i^\top = Z Z^\top$ held at iteration $k$, and $\beta_i$ were known. In this case, $Q_i^{(k)}$ would be a known linear function of $A_i A_i^\top$ parameterized by $Z^{(k)}$. Still, the $n \times p$ system $Q_i^{(k)} Z^{(k)} = S_i^{(k)}$ would not be sufficient to uncover the $n \times n$ local data matrix $A_i A_i^\top$ (strictly speaking, one only needs to recover $n(n+1)/2$ entries since $A_i A_i^\top$ is symmetric). Based on this discussion, we call DSSAL1 a privacy-aware method.

## 3.5 Computational cost

We conclude this section by discussing the computational cost of our algorithm per iteration. We first compute the matrix multiplication $Q^{(k)} Z^{(k)}$ by (19) and (20), whose computational cost for each agent is $O(np^2)$ as mentioned earlier. Then, at the center, the Uzawa method is applied to solving subproblem (18) which has a per-iteration complexity $O(np^2)$. In practice, it usually takes very few iterations to generate $Z^{(k+1)}$. Next, each agent uses a single iteration of SSI to generate $X_i^{(k+1)}$ by (27) with the computational cost $O(np^2)$ as discussed before. Finally, agent $i$ updates $W_i^{(k+1)}$ by (14) with the computational cost $4npm_i + O(np^2)$ (which represents the multiplier matrix $\Lambda_i^{(k+1)}$ implicitly). Overall, for each agent, the computational cost of our algorithm is $4npm_i + O(np^2)$ per iteration. At the center, the computational cost for approximately solving (18) is, empirically speaking, $O(np^2)$.

# 4 Convergence analysis

In this section, we analyze the global convergence of Algorithm 2 and prove that a sequence $\{Z^{(k)}\}$ generated by Algorithm 2 has at least one accumulation point, and any accumulation point is a first-order stationary point. A global, sub-linear convergence rate is also established.

We start with a property that a feasible point is first-order stationary if no progress can be made by solving (18).

**Lemma 4.1.** *Let $(Z^{(k)}, \{X_i^{(k)}\})$ be feasible. Then $Z^{(k)}$ is a first-order stationary point of the sparse PCA problem (2) if $D^{(k)} = 0$ is the minimizer of subproblem (18).*

The proof of Lemma 4.1 is deferred to Appendix C. It motivates the following definition of an $\epsilon$-stationary point for the subspace-splitting model (7).

**Definition 4.2.** *Suppose $Z^{(k)}$ is the $k$-th iterate of Algorithm 2. Then $Z^{(k)}$ is called an $\epsilon$-stationary point if the following condition holds:*

$$\frac{1}{d}\sum_{i=1}^{d}\left\|Z^{(k)}(Z^{(k)})^{\top} - X_i^{(k)}(X_i^{(k)})^{\top}\right\|_{\mathrm{F}}^{2} + \left\|D^{(k)}\right\|_{\mathrm{F}}^{2} \le \epsilon^{2},$$

*where $\epsilon > 0$ is a small constant.*

In order to prove the convergence of our algorithm, we need to impose some mild conditions on algorithm parameters, which are summarized below.

**Condition 1.** *The algorithm parameters $\eta > 0$, $c_i > 0$, $c_i' > 0$, $\delta_i \in [0,1)$, as well as two auxiliary parameters $\rho \ge 1$ and $\underline{\sigma} \in (0,1)$, satisfy the following conditions:*

$$0 < \eta < \frac{1}{2\bar{M}},\ 0 \le \delta_i < \frac{\underline{\sigma}}{2\sqrt{\rho d}},\ 0 < \underline{\sigma} < \min\left\{1, \frac{1}{\sqrt{c_i}}\right\},\ i = 1,\dots,d,$$

*where $\bar{M} = \mu\sqrt{np}/2 + 2\|A\|_{\mathrm{F}}^{2} + \sqrt{p}\sum_{i=1}^{d}\beta_i > 0$ is a constant.*

**Condition 2.** *Each penalty parameter $\beta_i(i = 1,\dots,d)$ in (11) has a lower bound*

$$\max\left\{\xi_i\|A_i\|_2^2,\ \frac{8(\mu np + \sqrt{p}\|A\|_{\mathrm{F}}^2)}{(1 - \underline{\sigma}^2)},\ \frac{6\sqrt{p}\|A_i\|_{\mathrm{F}}^2}{c_i\underline{\sigma}^2(1 - \underline{\sigma}^2)}\right\},$$

*where $\xi_i = \max\{c_i',\ 4\sqrt{2}/\underline{\sigma},\ 4(2\sqrt{\rho d} + \sqrt{2})/(\underline{\sigma} - 2\sqrt{\rho d}\delta_i),\ 4(\sqrt{2\rho d} + 1)/(c_i\underline{\sigma}^2\rho d)\} > 0$ is a constant. In addition, $\beta_i \le \rho\beta_j$ holds for any $i, j \in \{1,\dots,d\}$.*

We note that the above technical conditions are not necessary in a practical implementation, They are introduced purely for the purpose of theoretical analysis to facilitate obtaining a global convergence rate and corresponding worst-case complexity for Algorithm 2.

**Theorem 4.3.** *Suppose $\{Z^{(k)}\}$ is an iterate sequence generated by Algorithm 2, starting from an arbitrary orthogonal matrix $Z^{(0)} \in \mathcal{S}_{n,p}$, with parameters satisfying Conditions 1 and 2. Then $\{Z^{(k)}\}$ has at least one accumulation point and any accumulation point must be a first-order stationary point of the sparse PCA problem (2). Moreover, for any integer $K > 1$, it holds that*

$$\min_{k=1,\dots,K}\left\{\left\|D^{(k)}\right\|_{\mathrm{F}}^{2} + \frac{1}{d}\sum_{i=1}^{d}\left\|Z^{(k)}(Z^{(k)})^{\top} - X_i^{(k)}(X_i^{(k)})^{\top}\right\|_{\mathrm{F}}^{2}\right\} \le \frac{C}{K},$$

*where $C > 0$ is a constant.*

The proof of Theorem 4.3, which is rather complicated and lengthy, will be given in Appendix D. The global sub-linear convergence rate in Theorem 4.3 guarantees that DSSAL1 is able to return an $\epsilon$-stationary point in at most $O(1/\epsilon^2)$ iterations. Since DSSAL1 performs one round of communication per iteration, the number of communication rounds required to obtain an $\epsilon$-stationary point is also $O(1/\epsilon^2)$ at the most.

# 5 Numerical results

In this section, we evaluate the empirical performance of DSSAL1 through a set of comprehensive numerical experiments. All the experiments throughout this section are performed on a high-performance computing cluster [3], called LSSC-IV which is maintained at the State Key Laboratory of Scientific and Engineering Computing (LSEC), Chinese Academy of Sciences. The LSSC-IV cluster has 408 nodes, each consisting of two Inter(R) Xeon(R) Gold 6140 processors (at 2.30GHz ×18) with 192GB memory, running under the operating system Red Hat Enterprise Linux Server 7.3.

We compare the performance of DSSAL1 with two state-of-the-art algorithms: (1) an ADMM-type algorithm called SOC [37] and (2) a manifold proximal gradient method called ManPG-Ada [11]. Since open-source, parallel codes for the above two algorithms are not available, to conduct experiments under the distributed environment of the LSSC-IV cluster, we implemented the two existing algorithms and our own algorithm DSSAL1 in C++ with MPI for inter-process communications[4]. For the two existing algorithms, we set all parameters to their default values as described in [37, 11]. The linear algebra library Eigen[5] (version 3.3.8) is adopted for matrix computation tasks.

## 5.1 DSSAL1 Implementation details

In Algorithm 2, we set the penalty parameters to $\beta_i = 0.1(\|\nabla f_i(X_i^{(0)})\|_{\mathrm{F}} + \mu)$, and in subproblem (18) we set the hyperparameter to $\eta = 1/(\sum_{i=1}^d \beta_i)$ . In Algorithm 1, we set the step size to $\tau = 1/(2\eta)$ and terminate the algorithm whenever

$$\left\|D(j)^\top Z^{(k)} + (Z^{(k)})^\top D(j)\right\|_{\mathrm{F}} \leq \left\|D^{(k-1)}\right\|_{\mathrm{F}}$$

is satisfied or the number of iterations reaches 10.

We use the well-known SSI method [44] to obtain a very rough solution to subproblem (23) for $X_i$. More specifically, we initialize $X_i$ to the previous iterate $X_i^{(k)}$ and perform a single SSI iteration, as given by (27), to generate the next iterate $X_i^{(k+1)}$.

The stopping criteria used in Algorithm 2 are

$$\frac{1}{d}\sum_{i=1}^d \left\|Z^{(k)}(Z^{(k)})^\top - X_i^{(k)}(X_i^{(k)})^\top\right\|_{\mathrm{F}} \leq \epsilon_c \quad \text{and} \quad \left\|D^{(k)}\right\|_{\mathrm{F}} \leq \epsilon_g, \tag{28}$$

where $\epsilon_c$ and $\epsilon_g$ are two small positive constants. Unless otherwise specified, $\epsilon_c$ and $\epsilon_g$ are set to $10^{-6}$ and $10^{-8}np$, respectively. Algorithm 2 is also terminated once the iteration count reaches `MaxIter` = 50000.

## 5.2 Synthetic data generation

In our experiments, a synthetic data matrix $A \in \mathbb{R}^{n \times m}$ is constructed into the form of (economy-size) SVD:

$$A = U\Sigma V^\top, \tag{29}$$

where $U \in \mathbb{R}^{n \times n}$ and $V \in \mathbb{R}^{m \times n}$ satisfy $U^\top U = V^\top V = I_n$ and $\Sigma \in \mathbb{R}^{n \times n}$ is nonnegative and diagonal. Specifically, $U$ and $V$ are results of orthonormalization of random matrices whose entries are drawn independently and identically from $[-1, 1]$ under the uniform distribution, and

$$\Sigma_{ii} = \xi^{1-i}, \quad i = 1, \ldots, n,$$

where the parameter $\xi \geq 1$ determines the decay rate of singular values. Finally, we apply the standard PCA pre-processing operations to a data matrix $A = U\Sigma V^\top$ by subtracting the sample-mean from each sample and then normalizing the rows of the resulting matrix to make them of unit $\ell_2$-norm.

---

[3] More information at http://lsec.cc.ac.cn/chinese/lsec/LSSC-IVintroduction.pdf
[4] Our code is downloadable from http://lsec.cc.ac.cn/~liuxin/code.html
[5] Available from https://eigen.tuxfamily.org/index.php?title=Main_Page

For our synthetic data matrices, such pre-processing will only slightly perturb the decay rate of the singular-values before the pre-processing which is uniformly equal to $1/\xi$ by construction. Unless specified otherwise, the default value for the decay-rate parameter is $\xi = 1.1$.

In the numerical experiments, all the algorithms are started from the same initial points. Since the optimization problem is non-convex, different solvers may still occasionally return different solutions when starting from a common initial point at random. As suggested in [11], to increase the chance that all solvers find the same solution, we first run the Riemannian subgradient method [8, 18] for 500 iterations and then use the resulting output as the common starting point.

## 5.3 Comprehensive comparison on synthetic data

In order to do a thorough evaluation on the empirical performance of DSSAL1, we design four groups of test data matrices, generated as in Subsection 5.2. In each group, there is only one parameter varying while all the others are fixed. Specifically, the varying and fixed parameters for the four groups are as follows:

1. varying sample dimension $n = 1000 + 500j$ for $j = 0, 1, 2, 3, 4$, while $m = 128000$, $p = 10$, $\mu = 0.5$, $d = 128$;

2. varying number of computed loading vectors $p = 10 + 5j$ for $j = 0, 1, 2, 3, 4$, while $n = 1000$, $m = 128000$, $\mu = 0.3$, $d = 128$;

3. varying regularization parameter $\mu = 0.2 + 0.2j$ for $j = 0, 1, 2, 3, 4$, while $n = 1000$, $m = 128000$, $p = 10$, $d = 128$;

4. varying number of cores $d = 16 \times 2^j$ for $j = 0, 1, 2, 3, 4$, while $n = 1000$, $m = 256000$, $p = 10$, $\mu = 0.5$.

Additional experimental results for varying $\xi$ will be presented in the next subsection. Numerical results obtained for the above four test groups are presented in Tables 1 to 4, respectively, where we record wall-clock times in seconds and total rounds of communication. The average function values and sparsity levels for the four groups of tests are provided in Table 5. When computing the sparsity of a solution matrix (i.e., the percentage of zero elements), we set a matrix element to zero when its absolute value is less than $10^{-5}$.

From Table 5, we see that all three algorithms have attained comparable solution qualities with similar function values and sparsity levels in the four groups of testing problems.

Table 1: Comparison of DSSAL1, ManPG-Ada, and SOC for different $n$.

| | Wall-clock time in seconds | | | Rounds of communication | | |
|---|---|---|---|---|---|---|
| $n$ | DSSAL1 | ManPG-Ada | SOC | DSSAL1 | ManPG-Ada | SOC |
| 1000 | **10.97** | 26.76 | 72.07 | **655** | 1794 | 7569 |
| 1500 | **6.61** | 16.12 | 57.75 | **223** | 642 | 2201 |
| 2000 | **49.22** | 172.32 | 725.63 | **1238** | 5054 | 20444 |
| 2500 | **181.97** | 412.51 | 2296.26 | **3753** | 10971 | 45707 |
| 3000 | **34.58** | 153.40 | 860.76 | **680** | 2789 | 11480 |

It should be evident from Tables 1 to 4 that, in all four test groups and in terms of both wall-clock time and round of communication, DSSAL1 clearly outperforms ManPG-Ada which in turn outperforms SOC by large margins. Since the amount of communication per round for the three algorithms are essentially the same, their total communication overhead is proportional to the rounds of communication required. Because DSSAL1 takes far fewer rounds of communication than the other two algorithms (often by considerable margins), we conclude that DSSAL1 is a more communication-efficient algorithm than the other two. For example, in Table 3 for the case of $\mu = 0.8$, the number of communication rounds taken by DSSAL1 is less than a quarter of that by ManPG-Ada and one tenth of that by COS.

Table 2: Comparison of DSSAL1, ManPG-Ada, and SOC for different $p$.

| | Wall-clock time in seconds | | | Rounds of communication | | |
|---|---|---|---|---|---|---|
| $p$ | DSSAL1 | ManPG-Ada | SOC | DSSAL1 | ManPG-Ada | SOC |
| 10 | **9.74** | 26.42 | 75.25 | **629** | 1622 | 6652 |
| 15 | **29.99** | 56.31 | 153.98 | **1728** | 3586 | 15865 |
| 20 | **110.26** | 239.51 | 466.22 | **6144** | 14107 | 44086 |
| 25 | **68.79** | 148.34 | 334.58 | **3030** | 6153 | 27060 |
| 30 | **110.16** | 173.11 | 204.87 | **5133** | 6966 | 14621 |

Table 3: Comparison of DSSAL1, ManPG-Ada, and SOC for different $\mu$.

| | Wall-clock time in seconds | | | Rounds of communication | | |
|---|---|---|---|---|---|---|
| $\mu$ | DSSAL1 | ManPG-Ada | SOC | DSSAL1 | ManPG-Ada | SOC |
| 0.2 | **25.05** | 59.46 | 430.37 | **1537** | 4140 | 50000 |
| 0.4 | **22.16** | 46.76 | 115.36 | **1393** | 3061 | 13680 |
| 0.6 | **11.61** | 29.08 | 81.82 | **838** | 1899 | 8278 |
| 0.8 | **10.09** | 40.07 | 66.46 | **733** | 3369 | 8121 |
| 1.0 | **9.12** | 19.80 | 56.33 | **655** | 1348 | 6396 |

## 5.4 Empirical convergence rate

In this subsection, we examine empirical convergence rates of iterates produced by DSSAL1 and ManPG-Ada for comparison, while SOC is excluded from this experiment given its obvious non-competitiveness in previous experiments.

In the following experiments, we fix $n = 1000$, $m = 128000$, $p = 5$, $\mu = 0.2$, and $d = 128$. Three synthetic matrices $A \in \mathbb{R}^{n \times m}$ is randomly generated by (29) with $\xi$ taking three different values 1.15, 1.1, and 1.05, respectively, on which DSSAL1 and ManPG-Ada return $Z_{\mathrm{D}}^*$ and $Z_{\mathrm{M}}^*$, respectively, with smaller-than-usual termination tolerances $\epsilon_c = 10^{-8}$ and $\epsilon_g = 10^{-10} np$ in (28). We use the average of the two, $Z^* = (Z_{\mathrm{D}}^* + Z_{\mathrm{M}}^*)/2$, as a "ground truth" solution. Then we rerun the two algorithms on the same $A$ with the termination condition $\|Z^{(k)} - Z^*\|_{\mathrm{F}} \leq 3 \times 10^{-4}$ and record the quantity $\|Z^{(k)} - Z^*\|_{\mathrm{F}}$ at each iteration.

In Figure 2, we plot the iterate-error sequences $\{\|Z^{(k)} - Z^*\|_{\mathrm{F}}\}$ for both DSSAL1 and ManPG-Ada and observe that both algorithms appear to converge asymptotically at linear rates. Overall, however, the convergence of DSSAL1 is several times faster than that of ManPG-Ada. We also provide the ratio between the iteration number of ManPG-Ada and that of DSSAL1 for different values of $\xi$ in Table 6. In general, the closer to one $\xi$ is, the slower the singular values of $A$ decay, and the more difficult the problem tends to be. Table 6 demonstrates that in our test the advantage of DSSAL1 becomes more and more pronounced as the test instances become more and more difficult to solve.

# 6 Conclusion

In this paper, we propose a distributed algorithm, called DSSAL1, for solving the $\ell_1$-regularized optimization model (2) for sparse PCA computation. DSSAL1 has the following features.

1. The algorithm successfully extends the subspace-splitting strategy from orthogonal-transformation invariant objective function $f = \sum f_i$ to the sum $f + r$ where $r$ can be non-invariant and non-smooth.

2. The algorithm has built-in mechanism for local-data masking and hence naturally protects local-data privacy without requiring a special privacy preservation process.

Table 4: Comparison of DSSAL1, ManPG-Ada, and SOC for different $d$.

| | Wall-clock time in seconds | | | Rounds of communication | | |
|---|---|---|---|---|---|---|
| $d$ | DSSAL1 | ManPG-Ada | SOC | DSSAL1 | ManPG-Ada | SOC |
| 16 | **161.72** | 168.36 | 383.14 | **897** | 1169 | 5175 |
| 32 | **106.95** | 135.81 | 312.20 | **808** | 1169 | 5175 |
| 64 | **54.33** | 68.89 | 167.27 | **753** | 1169 | 5175 |
| 128 | **24.28** | 35.54 | 96.04 | **683** | 1169 | 5175 |
| 256 | **9.17** | 14.74 | 47.61 | **660** | 1169 | 5175 |

Table 5: Average function values and sparsity levels of DSSAL1, ManPG-Ada, and SOC for different tests.

| | Function value | | | Sparsity level | | |
|---|---|---|---|---|---|---|
| Test | DSSAL1 | ManPG-Ada | SOC | DSSAL1 | ManPG-Ada | SOC |
| Varying $n$ | -672.26 | -672.26 | -672.26 | 16.54% | 16.51% | 16.50% |
| Varying $p$ | -360.52 | -360.51 | -360.46 | 40.42% | 40.46% | 40.32% |
| Varying $\mu$ | -282.65 | -282.65 | -282.64 | 26.28% | 26.29% | 26.28% |
| Varying $d$ | -302.02 | -301.97 | -301.97 | 22.13% | 22.21% | 22.21% |

3. The algorithm is storage-efficient in that beside local data, it only requires storing $n \times p$ matrices (usually $p \ll n$) by each agent, thanks to a low-rank multiplier formula.

4. The algorithm has a global convergence guarantee to stationary points and a worst-case complexity under mild conditions, in spite of the nonlinear equality constraints for subspace consensus.

Comprehensive numerical simulations are conducted under a distributed environment to evaluate the performance of our algorithm in comparison to two state-of-the-art approaches. Remarkably, the communication rounds required by DSSAL1 are often over one order of magnitude smaller than existing methods. These results indicate that DSSAL1 has a great potential in solving large-scale application problems in distributed environments where data privacy is a primary concern.

Finally, we mention two related topics worthy of future studies. One is the possibility of developing asynchronous approaches for sparse PCA to address load balance issues in distributed environments. Another is to extend the subspace splitting strategy to decentralized networks so that a wider range of applications can benefit from the effectiveness of this approach.

# Appendices

## A  Proof of Lemma 2.2

*Proof of Lemma 2.2.* According to the definition of $\mathrm{Proj}_{\mathcal{T}_Z \mathcal{S}_{n,p}} (\cdot)$, it follows that

$$\left\| \mathrm{Proj}_{\mathcal{T}_Z \mathcal{S}_{n,p}} \left( -AA^\top Z + R(Z) \right) \right\|_{\mathrm{F}}^2$$
$$= \frac{1}{4} \left\| Z^\top \left( -AA^\top Z + R(Z) \right) - \left( -AA^\top Z + R(Z) \right)^\top Z \right\|_{\mathrm{F}}^2$$
$$+ \left\| \mathbf{P}_Z^\perp \left( -AA^\top Z + R(Z) \right) \right\|_{\mathrm{F}}^2$$
$$= \left\| \mathbf{P}_Z^\perp \left( -AA^\top Z + R(Z) \right) \right\|_{\mathrm{F}}^2 + \frac{1}{4} \left\| Z^\top R(Z) - R(Z)^\top Z \right\|_{\mathrm{F}}^2 ,$$

(a) $\xi = 1.15$     (b) $\xi = 1.1$     (c) $\xi = 1.05$

Figure 2: Comparison between DSSAL1 and ManPG-Ada of empirical convergence rates.

Table 6: The iteration number ratio between ManPG-Ada and DSSAL1 for different values of $\xi$.

|  | $\xi = 1.15$ | $\xi = 1.1$ | $\xi = 1.05$ |
|---|---|---|---|
| $\dfrac{\text{It}_{\text{ManPG-Ada}}}{\text{It}_{\text{DSSAL1}}}$ | $\dfrac{1202}{296} \approx 4.06$ | $\dfrac{1417}{335} \approx 4.23$ | $\dfrac{2014}{338} \approx 5.96$ |

where $R(Z) \in \partial r(Z)$. The proof is completed. $\qquad\square$

# B    Proof of Proposition 2.4

*Proof of Proposition 2.4.* To begin with, we assume that $(Z, \{X_i\})$ is a first-order stationary point. Then there exists $R(Z) \in \partial r(Z)$ such that

$$\mathbf{P}_Z^{\perp} \left( -AA^{\top}Z + R(Z) \right) = 0,$$

and $Z^{\top}R(Z)$ is symmetric. Let $\Theta = Z^{\top}R(Z) \in Z^{\top}\partial r(Z)$, $\Gamma_i = -X_i^{\top}A_i A_i^{\top}X_i$, and

$$\Lambda_i = -\mathbf{P}_{X_i}^{\perp} A_i A_i^{\top} X_i X_i^{\top} - X_i X_i^{\top} A_i A_i^{\top} \mathbf{P}_{X_i}^{\perp}$$

with $i = 1, \ldots, d$. Then the matrices $\Theta$, $\Gamma_i$ and $\Lambda_i$ are symmetric and $\operatorname{rank}(\Lambda_i) \leq 2p$. Moreover, we can deduce that

$$A_i A_i^{\top} X_i + X_i \Gamma_i + \Lambda_i X_i = A_i A_i^{\top} X_i - X_i X_i^{\top} A_i A_i^{\top} X_i - \mathbf{P}_{X_i}^{\perp} A_i A_i^{\top} X_i = 0,$$

and

$$R(Z) + \sum_{i=1}^{d} \Lambda_i Z - Z\Theta = R(Z) - \sum_{i=1}^{d} \mathbf{P}_Z^{\perp} A_i A_i^{\top} Z - ZZ^{\top}R(Z)$$
$$= \mathbf{P}_Z^{\perp} \left( -AA^{\top}Z + R(Z) \right) = 0.$$

Hence, $(Z, \{X_i\})$ satisfies the conditions in (9) under these specific choices of $\Theta$, $\Gamma_i$ and $\Lambda_i$.

Conversely, we now assume that there exist $R(Z) \in \partial r(Z)$ and symmetric matrices $\Theta$, $\Gamma_i$ and $\Lambda_i$ such that $(Z, \{X_i\})$ satisfies the conditions in (9). It follows from the first and second equality in (9) that

$$\sum_{i=1}^{d} \mathbf{P}_{X_i}^{\perp} A_i A_i^{\top} X_i X_i^{\top} = -\sum_{i=1}^{d} \mathbf{P}_{X_i}^{\perp} \left( X_i \Gamma_i + \Lambda_i X_i \right) X_i^{\top} = -\sum_{i=1}^{d} \mathbf{P}_{X_i}^{\perp} \Lambda_i X_i X_i^{\top}$$
$$= -\mathbf{P}_Z^{\perp} \left( \sum_{i=1}^{d} \Lambda_i Z \right) Z^{\top} = \mathbf{P}_Z^{\perp} \left( R(Z) - Z\Theta \right) Z^{\top} = \mathbf{P}_Z^{\perp} R(Z) Z^{\top}.$$

17

At the same time, since $X_i X_i^\top = ZZ^\top$, we have

$$\sum_{i=1}^{d} \mathbf{P}_{X_i}^\perp A_i A_i^\top X_i X_i^\top = \sum_{i=1}^{d} \mathbf{P}_{Z}^\perp A_i A_i^\top ZZ^\top = \mathbf{P}_{Z}^\perp A A^\top ZZ^\top.$$

Combining the above two equalities and orthogonality of $Z$, we arrive at

$$\mathbf{P}_{Z}^\perp \left( -AA^\top Z + R(Z) \right) = 0.$$

Left-multiplying both sides of the second equality in (9) by $Z^\top$, we obtain that

$$Z^\top R(Z) = \Theta - \sum_{i=1}^{d} Z^\top \Lambda_i Z,$$

which together with the symmetry of $\Lambda_i$ and $\Theta$ implies that $Z^\top R(Z)$ is also symmetric. This completes the proof. □

## C  Proof of Lemma 4.1

*Proof of Lemma 4.1.* Since $(Z^{(k)}, \{X_i^{(k)}\})$ is feasible, we know $X_i^{(k)}(X_i^{(k)})^\top = Z^{(k)}(Z^{(k)})^\top$ for $i = 1, \dots, d$. Thus, it can be readily verified that

$$
\begin{aligned}
Q^{(k)} Z^{(k)} &= \sum_{i=1}^{d} \left( \Lambda_i^{(k)} - \beta_i X_i^{(k)}(X_i^{(k)})^\top \right) Z^{(k)} \\
&= \sum_{i=1}^{d} \left( -\mathbf{P}_{Z^{(k)}}^\perp A_i A_i^\top Z^{(k)}(Z^{(k)})^\top - \beta_i Z^{(k)}(Z^{(k)})^\top \right) Z^{(k)} \\
&= -\mathbf{P}_{Z^{(k)}}^\perp A_i A_i^\top Z^{(k)} - \left( \sum_{i=1}^{d} \beta_i \right) Z^{(k)},
\end{aligned}
$$

which implies that

$$\mathrm{Proj}_{\mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p}} \left( Q^{(k)} Z^{(k)} \right) = \mathrm{Proj}_{\mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p}} \left( -A_i A_i^\top Z^{(k)} \right).$$

According to Theorem 4.1 in [58], the first-order optimality condition of (18) can be stated as:

$$0 \in \mathrm{Proj}_{\mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p}} \left( Q^{(k)} Z^{(k)} + \frac{1}{\eta} D^{(k)} + \partial r(Z^{(k)} + D^{(k)}) \right).$$

Since $D^{(k)} = 0$ is the global minimizer of (18), we have

$$
\begin{aligned}
0 &\in \mathrm{Proj}_{\mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p}} \left( Q^{(k)} Z^{(k)} + \partial r(Z^{(k)}) \right) \\
&= \mathrm{Proj}_{\mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p}} \left( -A_i A_i^\top Z^{(k)} + \partial r(Z^{(k)}) \right).
\end{aligned}
$$

We obtain the assertion of this lemma. □

## D  Convergence of Algorithm 2

Now we prove Theorem 4.3 to establish the global convergence of Algorithm 2. In addition to the notations introduced in Section 1, we further adopt the followings throughout the theoretical analysis. The notations $\mathrm{rank}(C)$ and $\sigma_{\min}(C)$ represent the rank and the smallest singular value of $C$, respectively.

For $X, Y \in \mathcal{S}_{n,p}$, we define $\mathbf{D_p}(X,Y) := XX^\top - YY^\top$ and $\mathbf{d_p}(X,Y) := \|\mathbf{D_p}(X,Y)\|_F$, standing for, respectively, the projection distance matrix and its measurement.

To begin with, we provide a sketch of our proof. Suppose $\{Z^{(k)}\}$ is the iteration sequence generated by Algorithm 2, with $X_i^{(k)}$ and $\Lambda_i^{(k)}$ being the local variable and multiplier of the $i$-th agent at the $k$-th iteration, respectively. The proof includes the following main steps.

1. The sequence $\{Z^{(k)}\}$ is bounded and the sequence $\{\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\})\}$ is bounded from below.

2. The sequence $\{Z^{(k)}\}$ satisfies $\mathbf{d_p^2}\left(Z^{(k+1)}, X_i^{(k)}\right) \leq 2(1 - \underline{\sigma}^2)$, and $\underline{\sigma}$ is a unified lower bound of the smallest singular values of the matrices $(X_i^k)^\top Z^{k+1} (i = 1, \ldots, d)$.

3. The sequence $\{\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\})\}$ is monotonically non-increasing, and hence is convergent.

4. The sequence $\{Z^{(k)}\}$ has at least one accumulation point, and any accumulation point is a first-order stationary point of the sparse PCA problem (2).

Next we verify all the items in the above sketch by proving the following lemmas and corollaries.

**Lemma D.1.** *Suppose $\{Z^{(k)}\}$ is the iterate sequence generated by Algorithm 2. Let*

$$g^{(k)}(D) = \left\langle Q^{(k)} Z^{(k)}, D \right\rangle + \frac{1}{2\eta} \|D\|_F^2 + r(Z^{(k)} + D).$$

*Then the following relationship holds for any $k \in \mathbb{N}$,*

$$g^{(k)}(0) - g^{(k)}(D^{(k)}) \geq \frac{1}{2\eta} \left\| D^{(k)} \right\|_F^2.$$

*Proof.* Since $g^{(k)}$ is strongly convex with modulus $\dfrac{1}{\eta}$, we have

$$g^{(k)}(\hat{D}) \geq g^{(k)}(D) + \left\langle \partial g^{(k)}(D), \hat{D} - D \right\rangle + \frac{1}{2\eta} \left\| \hat{D} - D \right\|_F^2, \tag{30}$$

for any $D, \hat{D} \in \mathbb{R}^{n \times p}$. In particular, if $\hat{D}, D \in \mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p}$, it holds that

$$\left\langle \partial g^{(k)}(D), \hat{D} - D \right\rangle = \left\langle \mathrm{Proj}_{\mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p}} \left( \partial g^{(k)}(D) \right), \hat{D} - D \right\rangle.$$

It follows from the first-order optimality condition of (18) that $0 \in \mathrm{Proj}_{\mathcal{T}_{Z^{(k)}} \mathcal{S}_{n,p}} \left( \partial g^{(k)}(D^{(k)}) \right)$. Finally, taking $\hat{D} = 0$ and $D = D^{(k)}$ in (30) yields the assertion of this lemma. $\square$

**Lemma D.2.** *Suppose $Z \in \mathcal{S}_{n,p}$ and $D \in \mathcal{T}_Z \mathcal{S}_{n,p}$. Then it holds that*

$$\left\| \mathrm{Proj}_{\mathcal{S}_{n,p}} (Z + D) - Z \right\|_F \leq \|D\|_F,$$

*and*

$$\left\| \mathrm{Proj}_{\mathcal{S}_{n,p}} (Z + D) - Z - D \right\|_F \leq \frac{1}{2} \|D\|_F^2.$$

*Proof.* The proof can be found in, for example, [33]. For the sake of completeness, we provide a proof here. It follows from the orthogonality of $Z$ and the skew-symmetry of $Z^\top D$ that $Z + D$ has full column rank. This yields that $\mathrm{Proj}_{\mathcal{S}_{n,p}}(Z + D) = (Z + D)F^{-1}$, where $F = (I_p + D^\top D)^{1/2}$. Since $\mathrm{Proj}_{\mathcal{S}_{n,p}}(Z + D) - Z = (Z(I_p - F) + D)F^{-1}$, we have

$$\left\| \mathrm{Proj}_{\mathcal{S}_{n,p}} (Z + D) - Z \right\|_F^2 = 2\mathrm{tr}\left( I_p - F^{-1} \right) - 2\mathrm{tr}\left( F^{-1} Z^\top D \right) = 2\mathrm{tr}\left( I_p - F^{-1} \right)$$

$$= 2\sum_{j=1}^d \left( 1 - \left(1 + \tilde{\sigma}_i^2\right)^{-1/2} \right) \leq \sum_{j=1}^d \tilde{\sigma}_i^2 = \|D\|_F^2,$$

19

where $\tilde{\sigma}_1 \geq \cdots \geq \tilde{\sigma}_d \geq 0$ are the singular values of $D$. Similarly, it follows from the relationship $\text{Proj}_{\mathcal{S}_{n,p}}(Z + D) - Z - D = (Z + D)(F^{-1} - I_p)$ that

$$\left\| \text{Proj}_{\mathcal{S}_{n,p}}(Z + D) - Z - D \right\|_{\text{F}} = \text{tr}\left( (I_p - F)^2 \right) = \sum_{j=1}^{d} \left( 1 - (1 + \tilde{\sigma}_i^2)^{1/2} \right)^2$$

$$\leq \frac{1}{4} \sum_{j=1}^{d} \tilde{\sigma}_i^4 = \frac{1}{4} \|D\|_{\text{F}}^4,$$

which completes the proof. $\qquad\square$

**Corollary D.3.** *Suppose* $\{Z^{(k)}\}$ *is the iterate sequence generated by Algorithm 2 with the parameters satisfying Condition 1. Then for any* $k \in \mathbb{N}$, *it holds that*

$$\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) - \mathcal{L}(Z^{(k+1)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) \geq \bar{M} \left\| D^{(k)} \right\|_{\text{F}}^2, \tag{31}$$

*where* $\bar{M} > 0$ *is a constant defined in Section 4.*

*Proof.* Firstly, it can be readily verified that

$$\left\| Q^{(k)} \right\|_{\text{F}} \leq \sum_{i=1}^{d} \left\| Q_i^{(k)} \right\|_{\text{F}} \leq \sum_{i=1}^{d} \left( 2 \|A_i\|_{\text{F}}^2 + \sqrt{p}\beta_i \right).$$

Let $\bar{q}^{(k)}(Z) = \text{tr}(Z^\top Q^{(k)} Z)/2$ be the smooth part of the objective function $q^{(k)}(Z)$ in (16). Since $\nabla \bar{q}^{(k)}$ is Lipschitz continuous with the corresponding Lipschitz constant $\left\| Q^{(k)} \right\|_{\text{F}}$, we have

$$\bar{q}^{(k)}(Z^{(k+1)}) - \bar{q}^{(k)}(Z^{(k)}) \leq \left\langle Q^{(k)} Z^{(k)}, Z^{(k+1)} - Z^{(k)} \right\rangle$$
$$+ \frac{1}{2} \left\| Q^{(k)} \right\|_{\text{F}} \left\| Z^{(k+1)} - Z^{(k)} \right\|_{\text{F}}^2.$$

It follows from Lemma D.2 that

$$\left\langle Q^{(k)} Z^{(k)}, Z^{(k+1)} - Z^{(k)} - D^{(k)} \right\rangle \leq \left\| Q^{(k)} Z^{(k)} \right\|_{\text{F}} \left\| Z^{(k+1)} - Z^{(k)} - D^{(k)} \right\|_{\text{F}}$$
$$\leq \sum_{i=1}^{d} \left( \|A_i\|_{\text{F}}^2 + \frac{\sqrt{p}}{2}\beta_i \right) \left\| D^{(k)} \right\|_{\text{F}}^2,$$

and

$$\frac{1}{2} \left\| Q^{(k)} \right\|_{\text{F}} \left\| Z^{(k+1)} - Z^{(k)} \right\|_{\text{F}}^2 \leq \sum_{i=1}^{d} \left( \|A_i\|_{\text{F}}^2 + \frac{\sqrt{p}}{2}\beta_i \right) \left\| D^k \right\|_{\text{F}}^2.$$

Combing the above three inequalities, we can obtain that

$$\bar{q}^{(k)}(Z^{(k+1)}) - \bar{q}^{(k)}(Z^{(k)}) \leq \left\langle Q^{(k)} Z^{(k)}, D^{(k)} \right\rangle + \sum_{i=1}^{d} \left( 2 \|A_i\|_{\text{F}}^2 + \beta_i \sqrt{p} \right) \left\| D^{(k)} \right\|_{\text{F}}^2.$$

It follows from Lemma D.1 that

$$\left\langle Q^{(k)} Z^{(k)}, D^{(k)} \right\rangle + r(Z^{(k)} + D^{(k)}) - r(Z^{(k)})$$
$$= g^{(k)}(D^{(k)}) - g^{(k)}(0) - \frac{1}{2\eta} \left\| D^{(k)} \right\|_{\text{F}}^2 \leq -\frac{1}{\eta} \left\| D^{(k)} \right\|_{\text{F}}^2,$$

which infers that

$$\bar{q}^{(k)}(Z^{(k+1)}) - \bar{q}^{(k)}(Z^{(k)}) + r(Z^{(k)} + D^{(k)}) - r(Z^{(k)})$$
$$\leq \sum_{i=1}^{d} \left( 2 \|A_i\|_{\text{F}}^2 + \beta_i \sqrt{p} \right) \left\| D^{(k)} \right\|_{\text{F}}^2 - \frac{1}{\eta} \left\| D^{(k)} \right\|_{\text{F}}^2.$$

This together with the Lipschitz continuity of $r(Z)$ yields that

$$q^{(k)}(Z^{(k+1)}) - q^{(k)}(Z^{(k)})$$
$$= \bar{q}^{(k)}(Z^{(k+1)}) - \bar{q}^{(k)}(Z^{(k)}) + r(Z^{(k+1)}) - r(Z^{(k)})$$
$$= \bar{q}^{(k)}(Z^{(k+1)}) - \bar{q}^{(k)}(Z^{(k)}) + r(Z^{(k)} + D^{(k)}) - r(Z^{(k)})$$
$$\quad + r(Z^{(k+1)}) - r(Z^{(k)} + D^{(k)})$$
$$\leq \bar{q}^{(k)}(Z^{(k+1)}) - \bar{q}^{(k)}(Z^{(k)}) + r(Z^{(k)} + D^{(k)}) - r(Z^{(k)})$$
$$\quad + \mu\sqrt{np} \left\| Z^{(k+1)} - Z^{(k)} - D^{(k)} \right\|_{\mathrm{F}}$$
$$\leq \sum_{i=1}^{d} \left( 2 \left\| A_i \right\|_{\mathrm{F}}^2 + \beta_i \sqrt{p} \right) \left\| D^{(k)} \right\|_{\mathrm{F}}^2 - \frac{1}{\eta} \left\| D^{(k)} \right\|_{\mathrm{F}}^2 + \frac{\mu}{2}\sqrt{np} \left\| D^{(k)} \right\|_{\mathrm{F}}^2$$
$$= \left( \bar{M} - \frac{1}{\eta} \right) \left\| D^{(k)} \right\|_{\mathrm{F}}^2.$$

Here, $\bar{M} > 0$ is a constant defined in Section 4. According to Condition 1, we know that $\bar{M} - 1/\eta \leq -\bar{M}$. Hence, we finally arrive at

$$\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) - \mathcal{L}(Z^{(k+1)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) = q^{(k)}(Z^{(k)}) - q^{(k)}(Z^{(k+1)})$$
$$\geq \bar{M} \left\| D^{(k)} \right\|_{\mathrm{F}}^2.$$

This completes the proof. $\qquad\square$

**Lemma D.4.** *Suppose $\{Z^{(k)}\}$ is the iterate sequence generated by Algorithm 2 with the parameters satisfying Condition 1. Then for any $k \in \mathbb{N}$, it can be verified that*

$$\mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k)} \right) \leq \rho \sum_{j=1}^{d} \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k)}, X_j^{(k)} \right) + \frac{8}{\beta_i} \left( \sqrt{p} \left\| A \right\|_{\mathrm{F}}^2 + \mu np \right), \tag{32}$$

*where $\rho \geq 1$ is a constant defined in Section 4.*

*Proof.* The inequality (31) directly results in the following relationship.

$$q^{(k)}(Z^{(k)}) - q^{(k)}(Z^{(k+1)}) \geq 0.$$

According to the definition of $q^{(k)}$, it follows that

$$0 \leq \frac{1}{2}\mathrm{tr}\left( (Z^{(k)})^\top Q^{(k)} Z^{(k)} \right) - \frac{1}{2}\mathrm{tr}\left( (Z^{(k+1)})^\top Q^{(k)} Z^{(k+1)} \right) + r(Z^{(k)}) - r(Z^{(k+1)})$$
$$\leq \frac{1}{2}\sum_{j=1}^{d}\mathrm{tr}\left( \left( \beta_j X_j^{(k)}(X_j^{(k)})^\top - \Lambda_j^{(k)} \right) \mathbf{D}_{\mathbf{p}}\left( Z^{(k+1)}, Z^{(k)} \right) \right) + 2\mu np.$$

By straightforward calculations, we can deduce that

$$\sum_{j=1}^{d}\mathrm{tr}\left( \Lambda_j^{(k)} \mathbf{D}_{\mathbf{p}}\left( Z^{(k)}, Z^{(k+1)} \right) \right) \leq \sum_{j=1}^{d} \left\| \Lambda_j^{(k)} \right\|_{\mathrm{F}} \mathbf{d}_{\mathbf{p}}\left( Z^{(k+1)}, Z^{(k)} \right)$$
$$\leq 4\sqrt{p} \sum_{j=1}^{d} \left\| A_j \right\|_{\mathrm{F}}^2 = 4\sqrt{p} \left\| A \right\|_{\mathrm{F}}^2,$$

and

$$\sum_{j=1}^{d} \beta_j \mathrm{tr}\left( X_j^{(k)}(X_j^{(k)})^\top \mathbf{D}_{\mathbf{p}}\left( Z^{(k+1)}, Z^{(k)} \right) \right)$$
$$= \frac{1}{2}\sum_{j=1}^{d}\beta_j \mathbf{d}_{\mathbf{p}}^2\left( Z^{(k)}, X_j^{(k)} \right) - \frac{1}{2}\sum_{j=1}^{d}\beta_j \mathbf{d}_{\mathbf{p}}^2\left( Z^{(k+1)}, X_j^{(k)} \right).$$

The above three inequalities yield that

$$\sum_{j=1}^{d} \beta_j \mathbf{d}_{\mathbf{p}}^2\left(Z^{(k+1)}, X_j^{(k)}\right) \leq \sum_{j=1}^{d} \beta_j \mathbf{d}_{\mathbf{p}}^2\left(Z^{(k)}, X_j^{(k)}\right) + 8\sqrt{p}\|A\|_{\mathrm{F}}^2 + 8\mu n p,$$

which further implies that

$$\mathbf{d}_{\mathbf{p}}^2\left(Z^{(k+1)}, X_i^{(k)}\right) \leq \frac{1}{\beta_i} \sum_{j=1}^{d} \beta_j \mathbf{d}_{\mathbf{p}}^2\left(Z^{(k+1)}, X_j^{(k)}\right)$$

$$\leq \rho \sum_{j=1}^{d} \mathbf{d}_{\mathbf{p}}^2\left(Z^{(k)}, X_j^{(k)}\right) + \frac{8}{\beta_i}\left(\sqrt{p}\|A\|_{\mathrm{F}}^2 + \mu n p\right).$$

This completes the proof. $\qquad\square$

**Lemma D.5.** *Suppose $Z^{(k+1)}$ is the $(k+1)$-th iterate generated by Algorithm 2 and satisfies the following condition:*

$$\mathbf{d}_{\mathbf{p}}^2\left(Z^{(k+1)}, X_i^{(k)}\right) \leq 2\left(1 - \underline{\sigma}^2\right),$$

*where $\underline{\sigma} \in (0, 1)$ is a constant defined in Condition 1. Let the algorithm parameters satisfy Conditions 1 and 2. Then for any $i = 1, \ldots, d$, it holds that*

$$h_i^{(k)}(X_i^{(k)}) - h_i^{(k)}(X_i^{(k+1)}) \geq \frac{1}{4}\underline{\sigma}^2 c_i \beta_i \mathbf{d}_{\mathbf{p}}^2\left(Z^{(k+1)}, X_i^{(k)}\right), \tag{33}$$

*and*

$$\mathbf{d}_{\mathbf{p}}^2\left(Z^{(k+1)}, X_i^{(k+1)}\right) \leq \left(1 - c_i\underline{\sigma}^2\right) \mathbf{d}_{\mathbf{p}}^2\left(Z^{(k+1)}, X_i^{(k)}\right) + \frac{12}{\beta_i}\sqrt{p}\|A_i\|_{\mathrm{F}}^2. \tag{34}$$

*Proof.* It follows from Condition 2 that $\beta_i > c_i'\|A_i\|_2^2$, which together with (25) yields that

$$h_i^{(k)}(X_i^{(k)}) - h_i^{(k)}(X_i^{(k+1)}) \geq \frac{c_i}{2\beta_i}\left\|\mathbf{P}_{X_i^{(k)}}^{\perp} H_i^{(k)} X_i^{(k)}\right\|_{\mathrm{F}}^2, \tag{35}$$

And it can be checked that

$$\begin{aligned}
\mathbf{P}_{X_i^{(k)}}^{\perp} H_i^{(k)} X_i^{(k)} &= \mathbf{P}_{X_i^{(k)}}^{\perp}\left(A_i A_i^{\top} X_i^{(k)} + \Lambda_i^{(k)} X_i^{(k)} + \beta_i Z^{(k+1)}(Z^{(k+1)})^{\top} X_i^{(k)}\right) \\
&= \mathbf{P}_{X_i^{(k)}}^{\perp}\left(A_i A_i^{\top} X_i^{(k)} - \mathbf{P}_{X_i^{(k)}}^{\perp} A_i A_i^{\top} X_i^{(k)}\right) - \beta_i \mathbf{P}_{X_i^{(k)}}^{\perp} Z^{(k+1)}(Z^{(k+1)})^{\top} X_i^{(k)} \\
&= -\beta_i \mathbf{P}_{X_i^{(k)}}^{\perp} Z^{(k+1)}(Z^{(k+1)})^{\top} X_i^{(k)}.
\end{aligned} \tag{36}$$

Suppose $\hat{\sigma}_1, \ldots, \hat{\sigma}_p$ are the singular values of $(X_i^{(k)})^{\top} Z^{(k+1)}$. It is clear that $0 \leq \hat{\sigma}_i \leq 1$ for any $i = 1, \ldots, p$ due to the orthogonality of $X_i^{(k)}$ and $Z^{(k+1)}$. On the one hand, we have

$$\mathbf{d}_{\mathbf{p}}^2\left(Z^{(k+1)}, X_i^{(k)}\right) = \left\|X_i^{(k)}(X_i^{(k)})^{\top} - Z^{(k+1)}(Z^{(k+1)})^{\top}\right\|_{\mathrm{F}}^2 = 2\sum_{j=1}^{p}\left(1 - \hat{\sigma}_j^2\right).$$

On the other hand, it follows from $\mathbf{d}_{\mathbf{p}}^2\left(Z^{(k+1)}, X_i^{(k)}\right) \leq 2\left(1 - \underline{\sigma}^2\right)$ that

$$\sigma_{\min}\left((X_i^{(k)})^{\top} Z^{(k+1)}\right) \geq \underline{\sigma}.$$

Let $Y_i^{(k)} = (X_i^{(k)})^{\top} Z^{(k+1)}(Z^{(k+1)})^{\top} X_i^{(k)}$. By straightforward calculations, we can derive that

$$\begin{aligned}
\left\|\mathbf{P}_{X_i^{(k)}}^{\perp} Z^{(k+1)}(Z^{(k+1)})^{\top} X_i^{(k)}\right\|_{\mathrm{F}}^2 &= \operatorname{tr}\left(Y_i^{(k)}\right) - \operatorname{tr}\left((Y_i^{(k)})^2\right) = \sum_{j=1}^{p} \hat{\sigma}_j^2\left(1 - \hat{\sigma}_j^2\right) \\
&\geq \sum_{j=1}^{p} \sigma_{\min}^2\left((X_i^{(k)})^{\top} Z^{(k+1)}\right)\left(1 - \hat{\sigma}_j^2\right) \geq \frac{1}{2}\underline{\sigma}^2 \mathbf{d}_{\mathbf{p}}^2\left(Z^{(k+1)}, X_i^{(k)}\right).
\end{aligned} \tag{37}$$

Combining (35), (36) and (37), we acquire the assertion (33). Then it follows from the definition of $h_i^{(k)}$ that

$$c_i \underline{\sigma}^2 \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k)} \right) \leq 2\mathrm{tr}\left( Z^{(k+1)} (Z^{(k+1)})^\top \mathbf{D}_{\mathbf{p}} \left( X_i^{(k+1)}, X_i^{(k)} \right) \right)$$
$$+ \frac{2}{\beta_i} \mathrm{tr}\left( \left( A_i A_i^\top + \Lambda_i^{(k)} \right) \mathbf{D}_{\mathbf{p}} \left( X_i^{(k+1)}, X_i^{(k)} \right) \right).$$

By straightforward calculations, we can obtain that

$$\mathrm{tr}\left( \left( A_i A_i^\top + \Lambda_i^{(k)} \right) \mathbf{D}_{\mathbf{p}} \left( X_i^{(k+1)}, X_i^{(k)} \right) \right) \leq \left\| A_i A_i^\top + \Lambda_i^{(k)} \right\|_{\mathrm{F}} \mathbf{d}_{\mathbf{p}} \left( X_i^{(k+1)}, X_i^{(k)} \right)$$
$$\leq 6\sqrt{p} \left\| A_i \right\|_{\mathrm{F}}^2 ,$$

and

$$\mathrm{tr}\left( Z^{(k+1)} (Z^{(k+1)})^\top \mathbf{D}_{\mathbf{p}} \left( X_i^{(k+1)}, X_i^{(k)} \right) \right)$$
$$= \frac{1}{2} \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k)} \right) - \frac{1}{2} \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k+1)} \right).$$

The above three relationships yield (34). We complete the proof. □

**Lemma D.6.** *Suppose $\{Z^{(k)}\}$ is the iterate sequence generated by Algorithm 2 initiated from $Z^{(0)} \in \mathcal{S}_{n,p}$ with the parameters satisfying Conditions 1 and 2. Then for any $i = 1, \ldots, d$ and $k \in \mathbb{N}$, it holds that*

$$\mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k)} \right) \leq 2 \left( 1 - \underline{\sigma}^2 \right). \tag{38}$$

*Proof.* We use mathematical induction to prove this lemma. To begin with, it follows from the inequality (32) that

$$\mathbf{d}_{\mathbf{p}}^2 \left( Z^{(1)}, X_i^{(0)} \right) \leq \rho \sum_{j=1}^{d} \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(0)}, X_j^{(0)} \right) + \frac{8}{\beta_i} \left( \sqrt{p} \left\| A \right\|_{\mathrm{F}}^2 + \mu n p \right)$$
$$= \frac{8}{\beta_i} \left( \sqrt{p} \left\| A \right\|_{\mathrm{F}}^2 + \mu n p \right) \leq 2 \left( 1 - \underline{\sigma}^2 \right),$$

under the relationship $\beta_i > 4(\sqrt{p} \left\| A \right\|_{\mathrm{F}}^2 + \mu n p)/(1 - \underline{\sigma}^2)$ in Condition 2. Thus, the argument (38) directly holds for $(Z^{(1)}, \{X_i^{(0)}\})$. Now, we assume the argument holds at $(Z^{(k+1)}, \{X_i^{(k)}\})$, and investigate the situation at $(Z^{(k+2)}, \{X_i^{(k+1)}\})$.

According to Condition 2, we have $12\sqrt{p} \left\| A_i \right\|_{\mathrm{F}}^2 / \beta_i < 2 \left( 1 - \underline{\sigma}^2 \right) c_i \underline{\sigma}^2$. Since we assume that $\mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k)} \right) \leq 2 \left( 1 - \underline{\sigma}^2 \right)$, it follows from the relationship (34) that

$$\mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k+1)} \right) \leq \left( 1 - c_i \underline{\sigma}^2 \right) \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k)} \right) + \frac{12}{\beta_i} \sqrt{p} \left\| A_i \right\|_{\mathrm{F}}^2$$
$$\leq 2 \left( 1 - \underline{\sigma}^2 \right) \left( 1 - c_i \underline{\sigma}^2 \right) + 2 \left( 1 - \underline{\sigma}^2 \right) c_i \underline{\sigma}^2 = 2 \left( 1 - \underline{\sigma}^2 \right),$$

which infers that $\sigma_{\min} \left( (X_i^{(k+1)})^\top Z^{(k+1)} \right) \geq \underline{\sigma}$. Similar to the proof of Lemma D.5, we can acquire that

$$\left\| \mathbf{P}_{X_i^{(k+1)}}^\perp Z^{(k+1)} (Z^{(k+1)})^\top X_i^{(k+1)} \right\|_{\mathrm{F}}^2 \geq \frac{1}{2} \underline{\sigma}^2 \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k+1)} \right). \tag{39}$$

Combining the condition (26) and the equality (36), we have

$$\left\| \mathbf{P}_{X_i^{(k+1)}}^\perp H_i^{(k)} X_i^{(k+1)} \right\|_{\mathrm{F}} \leq \delta_i \left\| \mathbf{P}_{X_i^{(k)}}^\perp H_i^{(k)} X_i^{(k)} \right\|_{\mathrm{F}}$$
$$= \delta_i \beta_i \left\| \mathbf{P}_{X_i^{(k)}}^\perp Z^{(k+1)} (Z^{(k+1)})^\top X_i^{(k)} \right\|_{\mathrm{F}} \leq \delta_i \beta_i \mathbf{d}_{\mathbf{p}} \left( Z^{(k+1)}, X_i^{(k)} \right). \tag{40}$$

On the other hand, it follows from the triangular inequality that

$$\left\| \mathbf{P}^{\perp}_{X_i^{(k+1)}} H_i^{(k)} X_i^{(k+1)} \right\|_{\mathrm{F}}$$
$$\geq \left\| \mathbf{P}^{\perp}_{X_i^{(k+1)}} \left( A_i A_i^{\top} + \Lambda_i^{(k+1)} + \beta_i Z^{(k+1)} (Z^{(k+1)})^{\top} \right) X_i^{(k+1)} \right\|_{\mathrm{F}}$$
$$- \left\| \mathbf{P}^{\perp}_{X_i^{(k+1)}} \left( \Lambda_i^{(k+1)} - \Lambda_i^{(k)} \right) X_i^{(k+1)} \right\|_{\mathrm{F}}$$

Combing the inequality (39), it can be verified that

$$\left\| \mathbf{P}^{\perp}_{X_i^{(k+1)}} \left( A_i A_i^{\top} + \Lambda_i^{(k+1)} + \beta_i Z^{(k+1)} (Z^{(k+1)})^{\top} \right) X_i^{(k+1)} \right\|_{\mathrm{F}}$$
$$= \beta_i \left\| \mathbf{P}^{\perp}_{X_i^{(k+1)}} Z^{(k+1)} (Z^{(k+1)})^{\top} X_i^{(k+1)} \right\|_{\mathrm{F}} \geq \frac{\sqrt{2}}{2} \underline{\sigma} \beta_i \mathbf{d_p} \left( Z^{(k+1)}, X_i^{(k+1)} \right).$$

Moreover, according to Lemma B.4 in [51], we have

$$\left\| \mathbf{P}^{\perp}_{X_i^{(k+1)}} \left( \Lambda_i^{(k+1)} - \Lambda_i^{(k)} \right) X_i^{(k+1)} \right\|_{\mathrm{F}} \leq \left\| \Lambda_i^{(k+1)} - \Lambda_i^{(k)} \right\|_{\mathrm{F}}$$
$$\leq 4 \left\| A_i \right\|_2^2 \mathbf{d_p} \left( X_i^{(k+1)}, X_i^{(k)} \right).$$

Combing the above three inequalities, we further obtain that

$$\left\| \mathbf{P}^{\perp}_{X_i^{(k+1)}} H_i^{(k)} X_i^{(k+1)} \right\|_{\mathrm{F}} \geq \frac{\sqrt{2}}{2} \underline{\sigma} \beta_i \mathbf{d_p} \left( Z^{(k+1)}, X_i^{(k+1)} \right)$$
$$- 4 \left\| A_i \right\|_2^2 \mathbf{d_p} \left( X_i^{(k+1)}, X_i^{(k)} \right).$$

Together with (40), this yields that

$$\frac{\sqrt{2}}{2} \underline{\sigma} \beta_i \mathbf{d_p} \left( Z^{(k+1)}, X_i^{(k+1)} \right)$$
$$\leq \delta_i \beta_i \mathbf{d_p} \left( Z^{(k+1)}, X_i^{(k)} \right) + 4 \left\| A_i \right\|_2^2 \mathbf{d_p} \left( X_i^{(k+1)}, X_i^{(k)} \right)$$
$$\leq \left( \delta_i \beta_i + 4 \left\| A_i \right\|_2^2 \right) \mathbf{d_p} \left( Z^{(k+1)}, X_i^{(k)} \right) + 4 \left\| A_i \right\|_2^2 \mathbf{d_p} \left( Z^{(k+1)}, X_i^{(k+1)} \right).$$

According to Conditions 1 and 2, we have $\sqrt{2}\underline{\sigma}\beta_i - 8\left\| A_i \right\|_2^2 > 0$ and $\underline{\sigma} - 2\sqrt{\rho d}\delta_i > 0$. Thus, it can be verified that

$$\mathbf{d_p} \left( Z^{(k+1)}, X_i^{(k+1)} \right) \leq \frac{2(\delta_i \beta_i + 4 \left\| A_i \right\|_2^2)}{\sqrt{2}\underline{\sigma}\beta_i - 8 \left\| A_i \right\|_2^2} \mathbf{d_p} \left( Z^{(k+1)}, X_i^{(k)} \right)$$
$$\leq \sqrt{\frac{1}{2\rho d}} \mathbf{d_p} \left( Z^{(k+1)}, X_i^{(k)} \right), \tag{41}$$

where the last inequality follows from the relationship $\beta > \dfrac{4 \left( 2\sqrt{\rho d} + \sqrt{2} \right) \left\| A_i \right\|_2^2}{\underline{\sigma} - 2\sqrt{\rho d}\delta_i}$ in Condition 2. This together with (32) and (38) yields that

$$\mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+2)}, X_i^{(k+1)} \right) \leq \rho \sum_{j=1}^{d} \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_j^{(k+1)} \right) + \frac{8}{\beta_i} \left( \sqrt{p} \left\| A \right\|_{\mathrm{F}}^2 + \mu n p \right)$$
$$\leq \frac{1}{2d} \sum_{j=1}^{d} \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_j^{(k)} \right) + \left( 1 - \underline{\sigma}^2 \right) \leq \left( 1 - \underline{\sigma}^2 \right) + \left( 1 - \underline{\sigma}^2 \right) = 2 \left( 1 - \underline{\sigma}^2 \right),$$

since we assume that $\beta_i > 8(\sqrt{p}\left\| A \right\|_{\mathrm{F}}^2 + \mu n p)/(1 - \underline{\sigma}^2)$ in Condition 2. The proof is completed. $\quad\square$

**Corollary D.7.** *Suppose $\{Z^{(k)}\}$ is the iterate sequence generated by Algorithm 2 initiated from $Z^{(0)} \in \mathcal{S}_{n,p}$, and the problem parameters satisfy Conditions 1 and 2. Then for any $k \in \mathbb{N}$, we can obtain that*

$$
\mathcal{L}(Z^{(k+1)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) - \mathcal{L}(Z^{(k+1)}, \{X_i^{(k+1)}\}, \{\Lambda_i^{(k)}\})
$$

$$
\geq \frac{1}{4}\underline{\sigma}^2 \sum_{i=1}^{d} c_i \beta_i \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k)} \right).
$$

*Proof.* This corollary directly follows from Lemma D.5 and Lemma D.6. $\qquad\square$

**Corollary D.8.** *Suppose $\{Z^{(k)}\}$ is the iterate sequence generated by Algorithm 2 initiated from $Z^{(0)} \in \mathcal{S}_{n,p}$, and problem parameters satisfy Conditions 1 and 2. Then for any $k \in \mathbb{N}$, we can acquire that*

$$
\mathcal{L}(Z^{(k+1)}, \{X_i^{(k+1)}\}, \{\Lambda_i^{(k)}\}) - \mathcal{L}(Z^{(k+1)}, \{X_i^{(k+1)}\}, \{\Lambda_i^{(k+1)}\})
$$

$$
\geq -\frac{\sqrt{2\rho d} + 1}{\rho d} \sum_{i=1}^{d} \|A_i\|_2^2 \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k)} \right).
$$

*Proof.* According to the Cauchy–Schwarz inequality, we can show that

$$
\left| \left\langle \Lambda_i^{(k+1)} - \Lambda_i^{(k)}, \mathbf{D}_{\mathbf{p}} \left( X_i^{(k+1)}, Z^{(k+1)} \right) \right\rangle \right| \leq \left\| \Lambda_i^{(k+1)} - \Lambda_i^{(k)} \right\|_{\mathrm{F}} \mathbf{d}_{\mathbf{p}} \left( Z^{(k+1)}, X_i^{(k+1)} \right)
$$

$$
\leq \sqrt{\frac{8}{\rho d}} \|A_i\|_2^2 \mathbf{d}_{\mathbf{p}} \left( X_i^{(k+1)}, X_i^{(k)} \right) \mathbf{d}_{\mathbf{p}} \left( Z^{(k+1)}, X_i^{(k)} \right),
$$

where the last inequality follows from Lemma B.4 in [51] and (41). In addition, we have

$$
\mathbf{d}_{\mathbf{p}} \left( X_i^{(k+1)}, X_i^{(k)} \right) \leq \mathbf{d}_{\mathbf{p}} \left( Z^{(k+1)}, X_i^{(k+1)} \right) + \mathbf{d}_{\mathbf{p}} \left( Z^{(k+1)}, X_i^{(k)} \right)
$$

$$
\leq \frac{\sqrt{2\rho d} + 1}{\sqrt{2\rho d}} \mathbf{d}_{\mathbf{p}} \left( Z^{(k+1)}, X_i^{(k)} \right),
$$

which implies that

$$
\left\langle \Lambda_i^{(k+1)} - \Lambda_i^{(k)}, \mathbf{D}_{\mathbf{p}} \left( X_i^{(k+1)}, Z^{(k+1)} \right) \right\rangle
$$

$$
\geq -\frac{2\left(\sqrt{2\rho d} + 1\right)}{\rho d} \|A_i\|_2^2 \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k)} \right).
$$

Combing the fact that

$$
\mathcal{L}(Z^{(k+1)}, \{X_i^{(k+1)}\}, \{\Lambda_i^{(k)}\}) - \mathcal{L}(Z^{(k+1)}, \{X_i^{(k+1)}\}, \{\Lambda_i^{(k+1)}\})
$$

$$
= \frac{1}{2} \sum_{i=1}^{d} \left\langle \Lambda_i^{(k+1)} - \Lambda_i^{(k)}, \mathbf{D}_{\mathbf{p}} \left( X_i^{(k+1)}, Z^{(k+1)} \right) \right\rangle,
$$

we complete the proof. $\qquad\square$

Now based on these lemmas and corollaries, we can demonstrate the monotonic non-increasing of $\{\mathcal{L}(\{X_i^k\}, Z^k, \{\Lambda_i^k\})\}$, which results in the global convergence of our algorithm.

**Proposition D.9.** *Suppose $\{Z^{(k)}\}$ is the iteration sequence generated by Algorithm 2 initiated from $Z^{(0)} \in \mathcal{S}_{n,p}$, and problem parameters satisfy Conditions 1 and 2. Then the sequence of augmented Lagrangian functions $\{\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\})\}$ is monotonically non-increasing, and for any $k \in \mathbb{N}$, it satisfies the following sufficient descent property:*

$$
\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) - \mathcal{L}(Z^{(k+1)}, \{X_i^{(k+1)}\}, \{\Lambda_i^{(k+1)}\})
$$

$$
\geq \sum_{i=1}^{d} J_i \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k+1)} \right) + \bar{M} \left\| D^{(k)} \right\|_{\mathrm{F}}^2, \tag{42}
$$

*where $J_i = \frac{1}{2}\rho d \underline{\sigma}^2 c_i \beta_i - 2(\sqrt{2\rho d} + 1) \|A_i\|_2^2 > 0$ is a constant.*

*Proof.* Combining Corollary D.3, Corollary D.7, and Corollary D.8, we obtain that

$$\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) - \mathcal{L}(Z^{(k+1)}, \{X_i^{(k+1)}\}, \{\Lambda_i^{(k+1)}\})$$

$$\geq \sum_{i=1}^{d} \left( \frac{1}{4}\underline{\sigma}^2 c_i \beta_i - \frac{\sqrt{2\rho d}+1}{\rho d} \|A_i\|_2^2 \right) \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k+1)}, X_i^{(k)} \right) + \bar{M} \left\| D^{(k)} \right\|_{\mathrm{F}}^2.$$

Recalling the relationship $\beta_i > 4(\sqrt{2\rho d} + 1) \|A_i\|_2^2/(\rho d \underline{\sigma}^2 c_i)$ in Condition 2, we can conclude that $\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) \geq \mathcal{L}(Z^{(k+1)}, \{X_i^{(k+1)}\}, \{\Lambda_i^{(k+1)}\})$. Hence, the sequence $\{\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\})\}$ is monotonically non-increasing. Finally, the above relationship together with (41) yields the assertion (42). The proof is finished. $\qquad\square$

Based on the above properties, we are ready to prove Theorem 4.3, which establishes the global convergence rate of our proposed algorithm.

*Proof of Theorem 4.3.* The whole sequence $\{Z^{(k)}, \{X_i^{(k)}\}\}$ is naturally bounded, since each of $X_i^{(k)}$ or $Z^{(k)}$ is orthogonal. Then it follows from the Bolzano-Weierstrass theorem that this sequence exists an accumulation point $\{Z^*, \{X_i^*\}\}$, where $Z^* \in \mathcal{S}_{n,p}$ and $X_i^* \in \mathcal{S}_{n,p}$. Moreover, the boundedness of $\{\Lambda_i^{(k)}\}$ results from the multipliers updating formula (15). Hence, the lower boundedness of $\{\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\})\}$ is owing to the continuity of the augmented Lagrangian function. Namely, there exists a constant $\underline{L}$ such that $\mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) \geq \underline{L}$ for all $k \in \mathbb{N}$.

It follows from the sufficient descent property (42) that

$$\sum_{k=1}^{K} \left\| D^{(k)} \right\|_{\mathrm{F}}^2$$

$$\leq \bar{M}^{-1} \sum_{k=1}^{K} \left( \mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) - \mathcal{L}(Z^{(k+1)}, \{X_i^{(k+1)}\}, \{\Lambda_i^{(k+1)}\}) \right) \tag{43}$$

$$= \bar{M}^{-1} \left( \mathcal{L}(Z^{(1)}, \{X_i^{(1)}\}, \{\Lambda_i^{(1)}\}) - \mathcal{L}(Z^{(K+1)}, \{X_i^{(K+1)}\}, \{\Lambda_i^{(K+1)}\}) \right)$$

$$\leq \bar{M}^{-1} \left( \mathcal{L}(Z^{(1)}, \{X_i^{(1)}\}, \{\Lambda_i^{(1)}\}) - \underline{L} \right),$$

and

$$\sum_{k=1}^{K} \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k)}, X_i^{(k)} \right)$$

$$\leq J_i^{-1} \sum_{k=1}^{K} \left( \mathcal{L}(Z^{(k-1)}, \{X_i^{(k-1)}\}, \{\Lambda_i^{(k-1)}\}) - \mathcal{L}(Z^{(k)}, \{X_i^{(k)}\}, \{\Lambda_i^{(k)}\}) \right) \tag{44}$$

$$= J_i^{-1} \left( \mathcal{L}(Z^{(0)}, \{X_i^{(0)}\}, \{\Lambda_i^{(0)}\}) - \mathcal{L}(Z^{(K)}, \{X_i^{(K)}\}, \{\Lambda_i^{(K)}\}) \right)$$

$$\leq J_i^{-1} \left( \mathcal{L}(Z^{(0)}, \{X_i^{(0)}\}, \{\Lambda_i^{(0)}\}) - \underline{L} \right).$$

Upon taking the limit as $K \to \infty$, we obtain that

$$\sum_{k=1}^{\infty} \left\| D^{(k)} \right\|_{\mathrm{F}}^2 < \infty \quad \text{and} \quad \sum_{k=1}^{\infty} \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k)}, X_i^{(k)} \right) < \infty,$$

which further implies that

$$\lim_{k \to \infty} \left\| D^{(k)} \right\|_{\mathrm{F}} = 0 \quad \text{and} \quad \lim_{k \to \infty} \mathbf{d}_{\mathbf{p}} \left( Z^{(k)}, X_i^{(k)} \right) = 0,$$

respectively. Combing this with Lemma 4.1, we know that any accumulation point $Z^*$ of sequence $\{Z^{(k)}\}$ is a first-order stationary point of the problem (2).

Eventually, we prove the sublinear convergence rate. Indeed, it follows from the inequalities (43) and (44) that

$$
\min_{k=1,\dots,K} \left\{ \left\| D^{(k)} \right\|_{\mathrm{F}}^2 + \frac{1}{d} \sum_{i=1}^{d} \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k)}, X_i^{(k)} \right) \right\}
$$

$$
\leq \frac{1}{K} \sum_{k=1}^{K} \left\{ \left\| D^{(k)} \right\|_{\mathrm{F}}^2 + \frac{1}{d} \sum_{i=1}^{d} \mathbf{d}_{\mathbf{p}}^2 \left( Z^{(k)}, X_i^{(k)} \right) \right\} \leq \frac{C}{K},
$$

where

$$
C = \bar{M}^{-1} \left( \mathcal{L}(Z^{(1)}, \{X_i^{(1)}\}, \{\Lambda_i^{(1)}\}) - \underline{L} \right)
$$

$$
+ \left( \sum_{i=1}^{d} J_i^{-1} \right) d^{-1} \left( \mathcal{L}(Z^{(0)}, \{X_i^{(0)}\}, \{\Lambda_i^{(0)}\}) - \underline{L} \right)
$$

is a positive constant. This completes the proof. $\qquad\square$

# References

[1] T. E. ABRUDAN, J. ERIKSSON, AND V. KOIVUNEN, *Steepest descent algorithms for optimization under unitary matrix constraint*, IEEE Transactions on Signal Processing, 56 (2008), pp. 1134–1147.

[2] P.-A. ABSIL, C. G. BAKER, AND K. A. GALLIVAN, *Trust-region methods on Riemannian manifolds*, Foundations of Computational Mathematics, 7 (2006), pp. 303–330.

[3] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization algorithms on matrix manifolds*, Princeton University Press, 2008.

[4] F. L. ANDRADE, M. A. FIGUEIREDO, AND J. XAVIER, *Distributed Picard iteration: Application to distributed EM and distributed PCA*, arXiv:2106.10665, (2021).

[5] K. J. ARROW, H. AZAWA, L. HURWICZ, AND H. UZAWA, *Studies in linear and non-linear programming*, vol. 2, Stanford University Press, 1958.

[6] M. BACÁK, R. BERGMANN, G. STEIDL, AND A. WEINMANN, *A second order nonsmooth variational model for restoring manifold-valued images*, SIAM Journal on Scientific Computing, 38 (2016), pp. A567–A597.

[7] T. BADEN, P. BERENS, K. FRANKE, M. R. ROSÓN, M. BETHGE, AND T. EULER, *The functional diversity of retinal ganglion cells in the mouse*, Nature, 529 (2016), pp. 345–350.

[8] G. C. BENTO, O. P. FERREIRA, AND J. G. MELO, *Iteration-complexity of gradient, subgradient and proximal point methods on Riemannian manifolds*, Journal of Optimization Theory and Applications, 173 (2017), pp. 548–562.

[9] G. CHEN, P. F. SULLIVAN, AND M. R. KOSOROK, *Biclustering with heterogeneous variance*, Proceedings of the National Academy of Sciences, 110 (2013), pp. 12253–12258.

[10] S. CHEN, A. GARCIA, M. HONG, AND S. SHAHRAMPOUR, *Decentralized Riemannian gradient descent on the Stiefel manifold*, in Proceedings of the 38th International Conference on Machine Learning, vol. 139, 2021, pp. 1594–1605.

[11] S. CHEN, S. MA, A. MAN-CHO SO, AND T. ZHANG, *Proximal gradient method for nonsmooth optimization over the Stiefel manifold*, SIAM Journal on Optimization, 30 (2020), pp. 210–239.

[12] W. CHEN, H. JI, AND Y. YOU, *An augmented lagrangian method for $\ell_1$-regularized optimization problems with orthogonality constraints*, SIAM Journal on Scientific Computing, 38 (2016), pp. B570–B592.

[13] F. H. Clarke, *Optimization and nonsmooth analysis*, SIAM, 1990.

[14] A. d'Aspremont, F. Bach, and L. El Ghaoui, *Optimal solutions for sparse principal component analysis*, Journal of Machine Learning Research, 9 (2008), pp. 1269–1294.

[15] A. d'Aspremont, L. El Ghaoui, M. I. Jordan, and G. R. Lanckriet, *A direct formulation for sparse PCA using semidefinite programming*, SIAM Review, 49 (2007), pp. 434–448.

[16] A. Edelman, T. A. Arias, and S. T. Smith, *The geometry of algorithms with orthogonality constraints*, SIAM Journal on Matrix Analysis and Applications, 20 (1998), pp. 303–353.

[17] O. Ferreira and P. Oliveira, *Subgradient algorithm on Riemannian manifolds*, Journal of Optimization Theory and Applications, 97 (1998), pp. 93–104.

[18] O. P. Ferreira, M. S. Louzeiro, and L. F. Prudente, *Iteration-complexity of the subgradient method on Riemannian manifolds with lower bounded curvature*, Optimization, 68 (2019), pp. 713–729.

[19] A. Gang and W. U. Bajwa, *A linearly convergent algorithm for distributed principal component analysis*, arXiv:2101.01300, (2021).

[20] ——, *FAST-PCA: A fast and exact algorithm for distributed principal component analysis*, arXiv:2108.12373, (2021).

[21] B. Gao, X. Liu, X. Chen, and Y.-X. Yuan, *A new first-order algorithmic framework for optimization problems with orthogonality constraints*, SIAM Journal on Optimization, 28 (2018), pp. 302–332.

[22] B. Gao, X. Liu, and Y.-X. Yuan, *Parallelizable algorithms for optimization problems with orthogonality constraints*, SIAM Journal on Scientific Computing, 41 (2019), pp. A1949–A1983.

[23] I. Gemp, B. McWilliams, C. Vernade, and T. Graepel, *Eigengame: PCA as a nash equilibrium*, arXiv:2010.00554, (2020).

[24] K. Gravuer, J. J. Sullivan, P. A. Williams, and R. P. Duncan, *Strong human association with plant invasion success for Trifolium introductions to New Zealand*, Proceedings of the National Academy of Sciences, 105 (2008), pp. 6344–6349.

[25] P. Grohs and S. Hosseini, *Nonsmooth trust region algorithms for locally Lipschitz functions on Riemannian manifolds*, IMA Journal of Numerical Analysis, 36 (2016), pp. 1167–1192.

[26] D. Hajinezhad and M. Hong, *Nonconvex alternating direction method of multipliers for distributed sparse principal component analysis*, in 2015 IEEE Global Conference on Signal and Information Processing, 2015, pp. 255–259.

[27] B. He, Y. You, and X. Yuan, *On the convergence of primal-dual hybrid gradient algorithm*, SIAM Journal on Imaging Sciences, 7 (2014), pp. 2526–2537.

[28] S. Hosseini and A. Uschmajew, *A Riemannian gradient sampling algorithm for nonsmooth optimization on manifolds*, SIAM Journal on Optimization, 27 (2017), pp. 173–189.

[29] J. Hu, B. Jiang, L. Lin, Z. Wen, and Y.-X. Yuan, *Structured quasi-Newton methods for optimization with orthogonality constraints*, SIAM Journal on Scientific Computing, 41 (2019), pp. A2239–A2269.

[30] J. Hu, A. Milzarek, Z. Wen, and Y. Yuan, *Adaptive quadratically regularized Newton method for Riemannian optimization*, SIAM Journal on Matrix Analysis and Applications, 39 (2018), pp. 1181–1207.

[31] W. Huang and K. Wei, *Riemannian proximal gradient methods*, Mathematical Programming, (2021), pp. 1–43.

[32] B. JIANG AND Y.-H. DAI, *A framework of constraint preserving update schemes for optimization on Stiefel manifold*, Mathematical Programming, 153 (2015), pp. 535–575.

[33] B. JIANG, S. MA, A. M.-C. SO, AND S. ZHANG, *Vector transport-free SVRG with general retraction for Riemannian optimization: Complexity analysis and practical implementation*, arXiv:1705.09059, (2017).

[34] I. M. JOHNSTONE AND A. Y. LU, *On consistency and sparsity for principal components analysis in high dimensions*, Journal of the American Statistical Association, 104 (2009), pp. 682–693.

[35] I. T. JOLLIFFE, N. T. TRENDAFILOV, AND M. UDDIN, *A modified principal component technique based on the LASSO*, Journal of Computational and Graphical Statistics, 12 (2003), pp. 531–547.

[36] A. KOVNATSKY, K. GLASHOFF, AND M. M. BRONSTEIN, *MADMM: A generic algorithm for non-smooth optimization on manifolds*, in European Conference on Computer Vision, Springer, 2016, pp. 680–696.

[37] R. LAI AND S. OSHER, *A splitting method for orthogonality constrained problems*, Journal of Scientific Computing, 58 (2014), pp. 431–449.

[38] Y. LOU, L. YU, S. WANG, AND P. YI, *Privacy preservation in distributed subgradient optimization algorithms*, IEEE Transactions on Cybernetics, 48 (2017), pp. 2154–2165.

[39] M. MAGDON-ISMAIL, *NP-hardness and inapproximability of sparse PCA*, Information Processing Letters, 126 (2017), pp. 35–38.

[40] J. H. MANTON, *Optimization algorithms exploiting unitary constraints*, IEEE Transactions on Signal Processing, 50 (2002), pp. 635–650.

[41] B. MCMAHAN, E. MOORE, D. RAMAGE, S. HAMPSON, AND B. A. Y. ARCAS, *Communication-efficient learning of deep networks from decentralized data*, in Proceedings of the 20th International Conference on Artificial Intelligence and Statistics, vol. 54, PMLR, 2017, pp. 1273–1282.

[42] Y. NISHIMORI AND S. AKAHO, *Learning algorithms utilizing quasi-geodesic flows on the Stiefel manifold*, Neurocomputing, 67 (2005), pp. 106–135.

[43] P. S. PACHECO, *An introduction to parallel programming*, Elsevier, 2011.

[44] H. RUTISHAUSER, *Simultaneous iteration method for symmetric matrices*, Numerische Mathematik, 16 (1970), pp. 205–223.

[45] H. SATO, *A Dai–Yuan-type Riemannian conjugate gradient method with the weak Wolfe conditions*, Computational Optimization and Applications, 64 (2016), pp. 101–118.

[46] H. SHEN AND J. Z. HUANG, *Sparse principal component analysis via regularized low rank matrix approximation*, Journal of Multivariate Analysis, 99 (2008), pp. 1015–1034.

[47] K. SJOSTRAND, E. ROSTRUP, C. RYBERG, R. LARSEN, C. STUDHOLME, H. BAEZNER, J. FERRO, F. FAZEKAS, L. PANTONI, D. INZITARI, ET AL., *Sparse decomposition and modeling of anatomical shape variation*, IEEE Transactions on Medical Imaging, 26 (2007), pp. 1625–1635.

[48] E. STIEFEL, *Richtungsfelder und fernparallelismus in n-dimensionalen mannigfaltigkeiten*, Commentarii Mathematici Helvetici, 8 (1935), pp. 305–353.

[49] L. WANG, B. GAO, AND X. LIU, *Multipliers correction methods for optimization problems over the Stiefel manifold*, CSIAM Transactions on Applied Mathematics, 2 (2021), pp. 508–531.

[50] L. WANG AND X. LIU, *Decentralized optimization over the Stiefel manifold by an approximate augmented Lagrangian function*, IEEE Transactions on Signal Processing, 70 (2022), pp. 3029–3041.

[51] L. WANG, X. LIU, AND Y. ZHANG, *A distributed and secure algorithm for computing dominant SVD based on projection splitting*, arXiv:2012.03461, (2020).

[52] Z. WEN AND W. YIN, *A feasible method for optimization with orthogonality constraints*, Mathematical Programming, 142 (2013), pp. 397–434.

[53] D. M. WITTEN, R. TIBSHIRANI, AND T. HASTIE, *A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis*, Biostatistics, 10 (2009), pp. 515–534.

[54] T. WU, E. HU, S. XU, M. CHEN, P. GUO, Z. DAI, T. FENG, L. ZHOU, W. TANG, L. ZHAN, ET AL., *clusterProfiler 4.0: A universal enrichment tool for interpreting omics data*, The Innovation, 2 (2021), p. 100141.

[55] N. XIAO, X. LIU, AND Y.-X. YUAN, *A class of smooth exact penalty function methods for optimization problems with orthogonality constraints*, Optimization Methods and Software, (2020), pp. 1–37.

[56] N. XIAO, X. LIU, AND Y.-x. YUAN, *A penalty-free infeasible approach for a class of nonsmooth optimization problems over the Stiefel manifold*, arXiv:2103.03514, (2021).

[57] X. XIAO, Y. LI, Z. WEN, AND L. ZHANG, *A regularized semi-smooth Newton method with projection steps for composite convex programs*, Journal of Scientific Computing, 76 (2018), pp. 364–389.

[58] W. H. YANG, L.-H. ZHANG, AND R. SONG, *Optimality conditions for the nonlinear programming problems on Riemannian manifolds*, Pacific Journal of Optimization, 10 (2014), pp. 415–434.

[59] H. YE AND T. ZHANG, *DeEPCA: Decentralized exact PCA with linear convergence rate*, Journal of Machine Learning Research, 22 (2021), pp. 1–27.

[60] C. ZHANG, M. AHMAD, AND Y. WANG, *ADMM based privacy-preserving decentralized optimization*, IEEE Transactions on Information Forensics and Security, 14 (2018), pp. 565–580.

[61] X. ZHU, *A Riemannian conjugate gradient method for optimization on the Stiefel manifold*, Computational Optimization and Applications, 67 (2017), pp. 73–110.

[62] H. ZOU, T. HASTIE, AND R. TIBSHIRANI, *Sparse principal component analysis*, Journal of Computational and Graphical Statistics, 15 (2006), pp. 265–286.

[63] H. ZOU AND L. XUE, *A selective overview of sparse principal component analysis*, Proceedings of the IEEE, 106 (2018), pp. 1311–1320.