
Improved Distributed Principal Component Analysis

Maria-Florina Balcan
School of Computer Science
Carnegie Mellon University
ninamf@cs.cmu.edu

Vandana Kanchanapally
School of Computer Science
Georgia Institute of Technology
vvandana@gatech.edu

Yingyu Liang
Department of Computer Science
Princeton University
yingyul@cs.princeton.edu

David Woodruff
Almaden Research Center
IBM Research
dpwoodru@us.ibm.com

Abstract

We study the distributed computing setting in which there are multiple servers, each holding a set of points, who wish to compute functions on the union of their point sets. A key task in this setting is Principal Component Analysis (PCA), in which the servers would like to compute a low dimensional subspace capturing as much of the variance of the union of their point sets as possible. Given a procedure for approximate PCA, one can use it to approximately solve problems such as k -means clustering and low rank approximation. The essential properties of an approximate distributed PCA algorithm are its communication cost and computational efficiency for a given desired accuracy in downstream applications. We give new algorithms and analyses for distributed PCA which lead to improved communication and computational costs for k -means clustering and related problems. Our empirical study on real world data shows a speedup of orders of magnitude, preserving communication with only a negligible degradation in solution quality. Some of these techniques we develop, such as a general transformation from a constant success probability subspace embedding to a high success probability subspace embedding with a dimension and sparsity independent of the success probability, may be of independent interest.

1 Introduction

Since data is often partitioned across multiple servers [20, 7, 18], there is an increased interest in computing on it in the distributed model. A basic tool for distributed data analysis is Principal Component Analysis (PCA). The goal of PCA is to find an r -dimensional (affine) subspace that captures as much of the variance of the data as possible. Hence, it can reveal low-dimensional structure in very high dimensional data. Moreover, it can serve as a preprocessing step to reduce the data dimension in various machine learning tasks, such as Non-Negative Matrix Factorization (NNMF) [15] and Latent Dirichlet Allocation (LDA) [4].

In the distributed model, approximate PCA was used by Feldman et al. [9] for solving a number of shape fitting problems such as k -means clustering, where the approximation is in the form of a *coreset*, and has the property that local coresets can be easily combined across servers into a global coreset, thereby providing an approximate PCA to the union of the data sets. Designing small coresets therefore leads to communication-efficient protocols. Coresets have the nice property that their size typically does not depend on the number n of points being approximated. A beautiful property of the coresets developed in [9] is that for approximate PCA their size also only depends linearly on the dimension d , whereas previous coresets depended quadratically on d [8]. This gives the best known communication protocols for approximate PCA and k -means clustering.

Despite this recent exciting progress, several important questions remain. First, can we improve the communication further as a function of the number of servers, the approximation error, and other parameters of the downstream applications (such as the number k of clusters in k -means clustering)? Second, while preserving optimal or nearly-optimal communication, can we improve the computational costs of the protocols? We note that in the protocols of Feldman et al. each server has to run a singular value decomposition (SVD) on her local data set, while additional work needs to be performed to combine the outputs of each server into a global approximate PCA. Third, are these algorithms practical and do they scale well with large-scale datasets? In this paper we give answers to the above questions. To state our results more precisely, we first define the model and the problems.

Communication Model. In the distributed setting, we consider a set of s nodes $\mathcal{V} = \{v_i, 1 \leq i \leq s\}$, each of which can communicate with a central coordinator v_0 . On each node v_i , there is a local data matrix $\mathbf{P}_i \in \mathbb{R}^{n_i \times d}$ having n_i data points in d dimension ($n_i > d$). The global data $\mathbf{P} \in \mathbb{R}^{n \times d}$ is then a concatenation of the local data matrix, i.e. $\mathbf{P}^\top = [\mathbf{P}_1^\top, \mathbf{P}_2^\top, \dots, \mathbf{P}_s^\top]$ and $n = \sum_{i=1}^s n_i$. Let p_i denote the i -th row of \mathbf{P} . Throughout the paper, we assume that the data points are centered to have zero mean, i.e., $\sum_{i=1}^n p_i = 0$. Uncentered data requires a rank-one modification to the algorithms, whose communication and computation costs are dominated by those in the other steps.

Approximate PCA and ℓ_2 -Error Fitting. For a matrix $\mathbf{A} = [a_{ij}]$, let $\|\mathbf{A}\|_F^2 = \sum_{i,j} a_{ij}^2$ be its Frobenius norm, and let $\sigma_i(\mathbf{A})$ be the i -th singular value of \mathbf{A} . Let $\mathbf{A}^{(t)}$ denote the matrix that contains the first t columns of \mathbf{A} . Let $L_{\mathbf{X}}$ denote the linear subspace spanned by the columns of \mathbf{X} . For a point p , let $\pi_L(p)$ be its projection onto subspace L and let $\pi_{L_{\mathbf{X}}}(p)$ be shorthand for $\pi_{L_{\mathbf{X}}}(p)$. For a point $p \in \mathbb{R}^d$ and a subspace $L \subseteq \mathbb{R}^d$, we denote the squared distance between p and L by

$$d^2(p, L) := \min_{q \in L} \|p - q\|_2^2 = \|p - \pi_L(p)\|_2^2.$$

Definition 1. The linear (or affine) r -Subspace k -Clustering on $\mathbf{P} \in \mathbb{R}^{n \times d}$ is

$$\min_{\mathcal{L}} d^2(\mathbf{P}, \mathcal{L}) := \sum_{i=1}^n \min_{L \in \mathcal{L}} d^2(p_i, L) \quad (1)$$

where \mathbf{P} is an $n \times d$ matrix whose rows are p_1, \dots, p_n , and $\mathcal{L} = \{L_j\}_{j=1}^k$ is a set of k centers, each of which is an r -dimensional linear (or affine) subspace.

PCA is a special case when $k = 1$ and the center is an r -dimensional subspace. This optimal r -dimensional subspace is spanned by the top r right singular vectors of \mathbf{P} , also known as the principal components, and can be found using the singular value decomposition (SVD). Another special case of the above is k -means clustering when the centers are points ($r = 0$). Constrained versions of this problem include NMF where the r -dimensional subspace should be spanned by positive vectors, and LDA which assumes a prior distribution defining a probability for each r -dimensional subspace. We will primarily be concerned with relative-error approximation algorithms, for which we would like to output a set \mathcal{L}' of k centers for which $d^2(\mathbf{P}, \mathcal{L}') \leq (1 + \epsilon) \min_{\mathcal{L}} d^2(\mathbf{P}, \mathcal{L})$.

For approximate distributed PCA, the following protocol is implicit in [9]: each server i computes its top $O(r/\epsilon)$ principal components \mathbf{Y}_i of \mathbf{P}_i and sends them to the coordinator. The coordinator stacks the $O(r/\epsilon) \times d$ matrices \mathbf{Y}_i on top of each other, forming an $O(sr/\epsilon) \times d$ matrix \mathbf{Y} , and computes the top r principal components of \mathbf{Y} , and returns these to the servers. This provides a relative-error approximation to the PCA problem. We refer to this algorithm as Algorithm `disPCA`.

Our Contributions. Our results are summarized as follows.

Improved Communication: We improve the communication cost for using distributed PCA for k -means clustering and similar ℓ_2 -fitting problems. The best previous approach is to use Corollary 4.5 in [9], which shows that given a data matrix \mathbf{P} , if we project the rows onto the space spanned by the top $O(k/\epsilon^2)$ principal components, and solve the k -means problem in this subspace, we obtain a $(1 + \epsilon)$ -approximation. In the distributed setting, this would require first running Algorithm `disPCA` with parameter $r = O(k/\epsilon^2)$, and thus communication at least $O(skd/\epsilon^3)$ to compute the $O(k/\epsilon^2)$ global principal components. Then one can solve a distributed k -means problem in this subspace, and an α -approximation in it translates to an overall $\alpha(1 + \epsilon)$ approximation.

Our Theorem 3 shows that it suffices to run Algorithm `disPCA` while only incurring $O(skd/\epsilon^2)$ communication to compute the $O(k/\epsilon^2)$ global principal components, preserving the k -means solution cost up to a $(1 + \epsilon)$ -factor. Our communication is thus a $1/\epsilon$ factor better, and illustrates that

for downstream applications it is sometimes important to “open up the box” rather than to directly use the guarantees of a generic PCA algorithm (which would give $O(skd/\epsilon^3)$ communication). One feature of this approach is that by using the distributed k -means algorithm in [3] on the projected data, the coordinator can sample points from the servers proportional to their local k -means cost solutions, which reduces the communication roughly by a factor of s , which would come from each server sending their local k -means coresets to the coordinator. Furthermore, before applying the above approach, one can first run any other dimension reduction to dimension d' so that the k -means cost is preserved up to certain accuracy. For example, if we want a $1+\epsilon$ approximation factor, we can set $d' = O(\log n/\epsilon^2)$ by a Johnson-Lindenstrauss transform; if we want a larger $2+\epsilon$ approximation factor, we can set $d' = O(k/\epsilon^2)$ using [5]. In this way the parameter d in the above communication cost bound can be replaced by d' . Note that unlike these dimension reductions, our algorithm for projecting onto principal components is deterministic and does not incur error probability.

Improved Computation: We turn to the computational cost of Algorithm `disPCA`, which to the best of our knowledge has not been addressed. A major bottleneck is that each player is computing a singular value decomposition (SVD) of its point set \mathbf{P}_i , which takes $\min(n_i d^2, n_i^2 d)$ time. We change Algorithm `disPCA` to instead have each server first sample an oblivious subspace embedding (OSE) [22, 6, 19, 17] matrix \mathbf{H}_i , and instead run the algorithm on the point set defined by the rows of $\mathbf{H}_i \mathbf{P}_i$. Using known OSEs, one can choose \mathbf{H}_i to have only a single non-zero entry per column and thus $\mathbf{H}_i \mathbf{P}_i$ can be computed in $\text{nnz}(\mathbf{P}_i)$ time. Moreover, the number of rows of \mathbf{H}_i is $O(d^2/\epsilon^2)$, which may be significantly less than the original n_i number of rows. This number of rows can be further reduced to $O(d \log^{O(1)} d/\epsilon^2)$ if one is willing to spend $O(\text{nnz}(\mathbf{P}_i) \log^{O(1)} d/\epsilon)$ time [19]. We note that the number of non-zero entries of $\mathbf{H}_i \mathbf{P}_i$ is no more than that of \mathbf{P}_i .

One technical issue is that each of s servers is locally performing a subspace embedding, which succeeds with only constant probability. If we want a single non-zero entry per column of \mathbf{H}_i , to achieve success probability $1 - O(1/s)$ so that we can union bound over all s servers succeeding, we naively would need to increase the number of rows of \mathbf{H}_i by a factor linear in s . We give a general technique, which takes a subspace embedding that succeeds with constant probability as a black box, and show how to perform a procedure which applies it $O(\log 1/\delta)$ times independently and from these applications finds one which is guaranteed to succeed with probability $1 - \delta$. Thus, in this setting the players can compute a subspace embedding of their data in $\text{nnz}(\mathbf{P}_i)$ time, for which the number of non-zero entries of $\mathbf{H}_i \mathbf{P}_i$ is no larger than that of \mathbf{P}_i , and without incurring this additional factor of s . This may be of independent interest.

It may still be expensive to perform the SVD of $\mathbf{H}_i \mathbf{P}_i$ and for the coordinator to perform an SVD on \mathbf{Y} in Algorithm `disPCA`. We therefore replace the SVD computation with a randomized approximate SVD computation with spectral norm error. Our contribution here is to analyze the error in distributed PCA and k -means after performing these speedups.

Empirical Results: Our speedups result in significant computational savings. The randomized techniques we use reduce the time by orders of magnitude on medium and large-scale data sets, while preserving the communication cost. Although the theory predicts a new small additive error because of our speedups, in our experiments the solution quality was only negligibly affected.

Related Work A number of algorithms for approximate distributed PCA have been proposed [21, 2, 14, 16, 9], but either without theoretical guarantees, or without considering communication. Most closely related to our work is [9, 12]. [9] observes that the top singular vectors of the local data is its summary and the union of these summaries is a summary of the global data, i.e., Algorithm `disPCA` discussed above. [12] studies algorithms in the arbitrary partition model in which each server holds a matrix \mathbf{P}_i and $\mathbf{P} = \sum_{i=1}^s \mathbf{P}_i$. More details and more related work can be found in the appendix.

2 Tradeoff between Communication and Solution Quality

Algorithm `disPCA` for distributed PCA is suggested in [21, 9], which consists of a local stage and a global stage. In the local stage, each node performs SVD on its local data matrix, and communicates the first t_1 singular values $\Sigma_i^{(t_1)}$ and the first t_1 right singular vectors $\mathbf{V}_i^{(t_1)}$ to the central coordinator. Then in the global stage, the coordinator concatenates $\Sigma_i^{(t_1)} (\mathbf{V}_i^{(t_1)})^\top$ to form a matrix \mathbf{Y} , and performs SVD on it to get the first t_2 right singular vectors.

To get some intuition, consider the easy case when the data points actually lie in an r -dimensional subspace. We can run Algorithm `disPCA` with $t_1 = t_2 = r$. Since \mathbf{P}_i has rank r , its projection to

$$\mathbf{P} = \begin{bmatrix} \mathbf{P}_1 \\ \vdots \\ \mathbf{P}_s \end{bmatrix} \xrightarrow{\text{Local PCA}} \begin{bmatrix} \Sigma_1^{(t_1)} (\mathbf{V}_1^{(t_1)})^\top \\ \vdots \\ \Sigma_s^{(t_1)} (\mathbf{V}_s^{(t_1)})^\top \end{bmatrix} = \begin{bmatrix} \mathbf{Y}_1 \\ \vdots \\ \mathbf{Y}_s \end{bmatrix} = \mathbf{Y} \xrightarrow{\text{Global PCA}} \mathbf{V}^{(t_2)}$$

Figure 1: The key points of the algorithm disPCA.

the subspace spanned by its first $t_1 = r$ right singular vectors, $\hat{\mathbf{P}}_i = \mathbf{U}_i \Sigma_i^{(r)} (\mathbf{V}_i^{(r)})^\top$, is identical to \mathbf{P}_i . Then we only need to do PCA on $\hat{\mathbf{P}}$, the concatenation of $\hat{\mathbf{P}}_i$. Observing that $\hat{\mathbf{P}} = \tilde{\mathbf{U}} \mathbf{Y}$ where $\tilde{\mathbf{U}}$ is orthonormal, it suffices to compute SVD on \mathbf{Y} , and only $\Sigma_i^{(r)} \mathbf{V}_i^{(r)}$ needs to be communicated. In the general case when the data may have rank higher than r , it turns out that one needs to set t_1 sufficiently large, so that $\hat{\mathbf{P}}_i$ approximates \mathbf{P}_i well enough and does not introduce too much error into the final solution. In particular, the following *close projection* property about SVD is useful:

Lemma 1. *Suppose \mathbf{A} has SVD $\mathbf{A} = \mathbf{U} \Sigma \mathbf{V}$ and let $\hat{\mathbf{A}} = \mathbf{A} \mathbf{V}^{(t)} (\mathbf{V}^{(t)})^\top$ denote its SVD truncation. If $t = O(r/\epsilon)$, then for any $d \times r$ matrix \mathbf{X} with orthonormal columns,*

$$0 \leq \|\mathbf{A} \mathbf{X} - \hat{\mathbf{A}} \mathbf{X}\|_F^2 \leq \epsilon d^2(\mathbf{A}, L_{\mathbf{X}}), \quad \text{and} \quad 0 \leq \|\mathbf{A} \mathbf{X}\|_F^2 - \|\hat{\mathbf{A}} \mathbf{X}\|_F^2 \leq \epsilon d^2(\mathbf{A}, L_{\mathbf{X}}).$$

This means that the projections of $\hat{\mathbf{A}}$ and \mathbf{A} on any r -dimensional subspace are close, when the projected dimension t is sufficiently large compared to r . Now, note that the difference between $\|\mathbf{P} - \mathbf{P} \mathbf{X} \mathbf{X}^\top\|_F^2$ and $\|\hat{\mathbf{P}} - \hat{\mathbf{P}} \mathbf{X} \mathbf{X}^\top\|_F^2$ is only related to $\|\mathbf{P} \mathbf{X}\|_F^2 - \|\hat{\mathbf{P}} \mathbf{X}\|_F^2 = \sum_i [\|\mathbf{P}_i \mathbf{X}\|_F^2 - \|\hat{\mathbf{P}}_i \mathbf{X}\|_F^2]$. Each term in which is bounded by the lemma. So we can use $\hat{\mathbf{P}}$ as a proxy for \mathbf{P} in the PCA task. Again, computing PCA on $\hat{\mathbf{P}}$ is equivalent to computing SVD on \mathbf{Y} , as done in Algorithm disPCA. These lead to the following theorem, which is implicit in [9], stating that the algorithm can produce a $(1 + \epsilon)$ -approximation for the distributed PCA problem.

Theorem 2. *Suppose Algorithm disPCA takes parameters $t_1 \geq r + \lceil 4r/\epsilon \rceil - 1$ and $t_2 = r$. Then*

$$\|\mathbf{P} - \mathbf{P} \mathbf{V}^{(r)} (\mathbf{V}^{(r)})^\top\|_F^2 \leq (1 + \epsilon) \min_{\mathbf{X}} \|\mathbf{P} - \mathbf{P} \mathbf{X} \mathbf{X}^\top\|_F^2$$

where the minimization is over $d \times r$ orthonormal matrices \mathbf{X} . The communication is $O(\frac{sr d}{\epsilon})$ words.

2.1 Guarantees for Distributed ℓ_2 -Error Fitting

Algorithm disPCA can also be used as a pre-processing step for applications such as ℓ_2 -error fitting. In this section, we prove the correctness of Algorithm disPCA as pre-processing for these applications. In particular, we show that by setting t_1, t_2 sufficiently large, the objective value of any solution merely changes when the original data \mathbf{P} is replaced the projected data $\tilde{\mathbf{P}} = \mathbf{P} \mathbf{V}^{(t_2)} (\mathbf{V}^{(t_2)})^\top$. Therefore, the projected data serves as a proxy of the original data, i.e., any distributed algorithm can be applied on the projected data to get a solution on the original data. As the dimension is lower, the communication cost is reduced. Formally,

Theorem 3. *Let $t_1 = t_2 = O(rk/\epsilon^2)$ in Algorithm disPCA for $\epsilon \in (0, 1/3)$. Then there exists a constant $c_0 \geq 0$ such that for any set of k centers \mathcal{L} in r -Subspace k -Clustering,*

$$(1 - \epsilon) d^2(\mathbf{P}, \mathcal{L}) \leq d^2(\tilde{\mathbf{P}}, \mathcal{L}) + c_0 \leq (1 + \epsilon) d^2(\mathbf{P}, \mathcal{L}).$$

The theorem implies that any α -approximate solution \mathcal{L} on the projected data $\tilde{\mathbf{P}}$ is a $(1 + 3\epsilon)\alpha$ -approximation on the original data \mathbf{P} . To see this, let \mathcal{L}^* denote the optimal solution. Then

$$(1 - \epsilon) d^2(\mathbf{P}, \mathcal{L}) \leq d^2(\tilde{\mathbf{P}}, \mathcal{L}) + c_0 \leq \alpha d^2(\tilde{\mathbf{P}}, \mathcal{L}^*) + c_0 \leq \alpha(1 + \epsilon) d^2(\mathbf{P}, \mathcal{L}^*)$$

which leads to $d^2(\mathbf{P}, \mathcal{L}) \leq (1 + 3\epsilon)\alpha d^2(\mathbf{P}, \mathcal{L}^*)$. In other words, the distributed PCA step only introduces a small multiplicative approximation factor of $(1 + 3\epsilon)$.

The key to prove the theorem is the close projection property of the algorithm (Lemma 4): for any low dimensional subspace spanned by \mathbf{X} , the projections of \mathbf{P} and $\tilde{\mathbf{P}}$ on the subspace are close. In

Algorithm 1 Distributed k -means clustering

Input: $\{\mathbf{P}_i\}_{i=1}^s$, $k \in \mathbb{N}_+$ and $\epsilon \in (0, 1/2)$, a non-distributed α -approximation algorithm \mathcal{A}_α

- 1: Run Algorithm `disPCA` with $t_1 = t_2 = O(k/\epsilon^2)$ to get \mathbf{V} , and send \mathbf{V} to all nodes.
- 2: Run the distributed k -means clustering algorithm in [3] on $\{\mathbf{P}_i \mathbf{V} \mathbf{V}^\top\}_{i=1}^s$, using \mathcal{A}_α as a subroutine, to get k centers \mathcal{L} .

Output: \mathcal{L} .

particular, we choose \mathbf{X} to be the orthonormal basis of the subspace spanning the centers. Then the difference between the objective values of \mathbf{P} and $\tilde{\mathbf{P}}$ can be decomposed into two terms depending only on $\|\mathbf{P}\mathbf{X} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2$ and $\|\mathbf{P}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2$ respectively, which are small as shown by the lemma. The complete proof of Theorem 3 is provided in the appendix.

Lemma 4. *Let $t_1 = t_2 = O(k/\epsilon)$ in Algorithm `disPCA`. Then for any $d \times k$ matrix \mathbf{X} with orthonormal columns, $0 \leq \|\mathbf{P}\mathbf{X} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2 \leq \epsilon d^2(\mathbf{P}, L_{\mathbf{X}})$, and $0 \leq \|\mathbf{P}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \leq \epsilon d^2(\mathbf{P}, L_{\mathbf{X}})$.*

Proof Sketch: We first introduce some auxiliary variables for the analysis, which act as intermediate connections between \mathbf{P} and $\tilde{\mathbf{P}}$. Imagine we perform two kinds of projections: first project \mathbf{P}_i to $\hat{\mathbf{P}}_i = \mathbf{P}_i \mathbf{V}_i^{(t_1)} (\mathbf{V}_i^{(t_1)})^\top$, then project $\hat{\mathbf{P}}_i$ to $\bar{\mathbf{P}}_i = \hat{\mathbf{P}}_i \mathbf{V}^{(t_2)} (\mathbf{V}^{(t_2)})^\top$. Let $\hat{\mathbf{P}}$ denote the vertical concatenation of $\hat{\mathbf{P}}_i$ and let $\bar{\mathbf{P}}$ denote the vertical concatenation of $\bar{\mathbf{P}}_i$. These variables are designed so that the difference between \mathbf{P} and $\hat{\mathbf{P}}$ and that between $\hat{\mathbf{P}}$ and $\bar{\mathbf{P}}$ are easily bounded.

Our proof then proceeds by first bounding these differences, and then bounding that between $\bar{\mathbf{P}}$ and $\tilde{\mathbf{P}}$. In the following we sketch the proof for the second statement, while the other statement can be proved by a similar argument. See the appendix for details.

$$\|\mathbf{P}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 = \left[\|\mathbf{P}\mathbf{X}\|_F^2 - \|\hat{\mathbf{P}}\mathbf{X}\|_F^2 \right] + \left[\|\hat{\mathbf{P}}\mathbf{X}\|_F^2 - \|\bar{\mathbf{P}}\mathbf{X}\|_F^2 \right] + \left[\|\bar{\mathbf{P}}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \right].$$

The first term is just $\sum_{i=1}^s \left[\|\mathbf{P}_i \mathbf{X}\|_F^2 - \|\hat{\mathbf{P}}_i \mathbf{X}\|_F^2 \right]$, each of which can be bounded by Lemma 1, since $\hat{\mathbf{P}}_i$ is the SVD truncation of \mathbf{P}_i . The second term can be bounded similarly. The more difficult part is the third term. Note that $\bar{\mathbf{P}}_i = \hat{\mathbf{P}}_i \mathbf{Z}$, $\tilde{\mathbf{P}}_i = \mathbf{P}_i \mathbf{Z}$ where $\mathbf{Z} := \mathbf{V}^{(t_2)} (\mathbf{V}^{(t_2)})^\top \mathbf{X}$, leading to $\|\bar{\mathbf{P}}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 = \sum_{i=1}^s \left[\|\hat{\mathbf{P}}_i \mathbf{Z}\|_F^2 - \|\mathbf{P}_i \mathbf{Z}\|_F^2 \right]$. Although \mathbf{Z} is not orthonormal as required by Lemma 1, we prove a generalization (Lemma 7 in the appendix) which can be applied to show that the third term is indeed small. \square

Application to k -Means Clustering To see the implication, consider the k -means clustering problem. We can first perform any other possible dimension reduction to dimension d' so that the k -means cost is preserved up to accuracy ϵ , and then run Algorithm `disPCA` and finally run any distributed k -means clustering algorithm on the projected data to get a good approximate solution. For example, in the first step we can set $d' = O(\log n/\epsilon^2)$ using a Johnson-Lindenstrauss transform, or we can perform no reduction and simply use the original data.

As a concrete example, we can use original data ($d' = d$), then run Algorithm `disPCA`, and finally run the distributed clustering algorithm in [3] which uses any non-distributed α -approximation algorithm as a subroutine and computes a $(1 + \epsilon)\alpha$ -approximate solution. The resulting algorithm is presented in Algorithm 1.

Theorem 5. *With probability at least $1 - \delta$, Algorithm 1 outputs a $(1 + \epsilon)^2 \alpha$ -approximate solution for distributed k -means clustering. The total communication cost of Algorithm 1 is $O(\frac{sk}{\epsilon^2})$ vectors in \mathbb{R}^d plus $O\left(\frac{1}{\epsilon^4} \left(\frac{k^2}{\epsilon^2} + \log \frac{1}{\delta}\right) + sk \log \frac{sk}{\delta}\right)$ vectors in $\mathbb{R}^{O(k/\epsilon^2)}$.*

3 Fast Distributed PCA

Subspace Embeddings One can significantly improve the time of the distributed PCA algorithms by using subspace embeddings, while keeping similar guarantees as in Lemma 4, which suffice for l_2 -error fitting. More precisely, a subspace embedding matrix $\mathbf{H} \in \mathbb{R}^{\ell \times n}$ for a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ has the property that for all vectors $y \in \mathbb{R}^d$, $\|\mathbf{H}\mathbf{A}y\|_2 = (1 \pm \epsilon)\|\mathbf{A}y\|_2$. Suppose independently,

each node v_i chooses a random subspace embedding matrix \mathbf{H}_i for its local data \mathbf{P}_i . Then, they run Algorithm `disPCA` on the embedded data $\{\mathbf{H}_i \mathbf{P}_i\}_{i=1}^s$ instead of on the original data $\{\mathbf{P}_i\}_{i=1}^s$.

The work of [22] pioneered subspace embeddings. The recent fast sparse subspace embeddings [6] and its optimizations [17, 19] are particularly suitable for large scale sparse data sets, since their running time is linear in the number of non-zero entries in the data matrix, and they also preserve the sparsity of the data. The algorithm takes as input an $n \times d$ matrix \mathbf{A} and a parameter ℓ , and outputs an $\ell \times d$ embedded matrix $\mathbf{A}' = \mathbf{H}\mathbf{A}$ (the embedded matrix \mathbf{H} does need to be built explicitly). The embedded matrix is constructed as follows: initialize $\mathbf{A}' = \mathbf{0}$; for each row in \mathbf{A} , multiply it by $+1$ or -1 with equal probability, then add it to a row in \mathbf{A}' chosen uniformly at random.

The success probability is constant, while we need to set it to be $1 - \delta$ where $\delta = \Theta(1/s)$. Known results which preserve the number of non-zero entries of \mathbf{H} to be 1 per column increase the dimension of \mathbf{H} by a factor of s . To avoid this, we propose an approach to boost the success probability by computing $O(\log \frac{1}{\delta})$ independent embeddings, each with only constant success probability, and then run a cross validation style procedure to find one which succeeds with probability $1 - \delta$. More precisely, we compute the SVD of all embedded matrices $\mathbf{H}_j \mathbf{A} = \mathbf{U}_j \Sigma_j \mathbf{V}_j^\top$, and find a $j \in [r]$ such that for at least half of the indices $j' \neq j$, all singular values of $\Sigma_j \mathbf{V}_j^\top \mathbf{V}_{j'} \Sigma_{j'}^\top$ are in $[1 \pm O(\epsilon)]$ (see Algorithm 4 in the appendix). The reason why such an embedding $\mathbf{H}_j \mathbf{A}$ succeeds with high probability is as follows. Any two successful embeddings $\mathbf{H}_j \mathbf{A}$ and $\mathbf{H}_{j'} \mathbf{A}$, by definition, satisfy that $\|\mathbf{H}_j \mathbf{A} x\|_2^2 = (1 \pm O(\epsilon)) \|\mathbf{H}_{j'} \mathbf{A} x\|_2^2$ for all x , which we show is equivalent to passing the test on the singular values. Since with probability at least $1 - \delta$, 9/10 fraction of the embeddings are successful, it follows that the one we choose is successful with probability $1 - \delta$.

Randomized SVD The exact SVD of an $n \times d$ matrix is impractical in the case when n or d is large. Here we show that the randomized SVD algorithm from [11] can be applied to speed up the computation without compromising the quality of the solution much. We need to use their specific form of randomized SVD since the error is with respect to the spectral norm, rather than the Frobenius norm, and so can be much smaller as needed by our applications.

The algorithm first probes the row space of the $\ell \times d$ input matrix \mathbf{A} with an $\ell \times 2t$ random matrix Ω and orthogonalizes the image of Ω to get a basis \mathbf{Q} (i.e., QR-factorize $\mathbf{A}^\top \Omega$); projects the data to this basis and computes the SVD factorization on the smaller matrix $\mathbf{A}\mathbf{Q}$. It also performs q power iterations to push the basis towards the top t singular vectors.

Fast Distributed PCA for ℓ_2 -Error Fitting We modify Algorithm `disPCA` by first having each node do a subspace embedding locally, then replace each SVD invocation with a randomized SVD invocation. We thus arrive at Algorithm 2. For ℓ_2 -error fitting problems, by combining approximation guarantees of the randomized techniques with that of distributed PCA, we are able to prove:

Theorem 6. *Suppose Algorithm 2 takes $\epsilon \in (0, 1/2]$, $t_1 = t_2 = O(\max\{\frac{k}{\epsilon^2}, \log \frac{s}{\delta}\})$, $\ell = O(\frac{d^2}{\epsilon^2})$, $q = O(\max\{\log \frac{d}{\epsilon}, \log \frac{sk}{\epsilon}\})$ as input, and sets the failure probability of each local subspace embedding to $\delta' = \delta/2s$. Let $\tilde{\mathbf{P}} = \mathbf{P}\mathbf{V}\mathbf{V}^\top$. Then with probability at least $1 - \delta$, there exists a constant $c_0 \geq 0$, such that for any set of k points \mathcal{L} ,*

$$(1 - \epsilon)d^2(\mathbf{P}, \mathcal{L}) - \epsilon \|\mathbf{P}\mathbf{X}\|_F^2 \leq d^2(\tilde{\mathbf{P}}, \mathcal{L}) + c_0 \leq (1 + \epsilon)d^2(\mathbf{P}, \mathcal{L}) + \epsilon \|\mathbf{P}\mathbf{X}\|_F^2$$

where \mathbf{X} is an orthonormal matrix whose columns span \mathcal{L} . The total communication is $O(sk d/\epsilon^2)$ and the total time is $O(\text{nnz}(\mathbf{P}) + s \left[\frac{d^3 k}{\epsilon^4} + \frac{k^2 d^2}{\epsilon^6} \right] \log \frac{d}{\epsilon} \log \frac{sk}{\delta \epsilon})$.

Proof Sketch: It suffices to show that $\tilde{\mathbf{P}}$ enjoys the close projection property as in Lemma 4, i.e., $\|\mathbf{P}\mathbf{X} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2 \approx 0$ and $\|\mathbf{P}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \approx 0$ for any orthonormal matrix whose columns span a low dimensional subspace. Note that Algorithm 2 is just running Algorithm `disPCA` (with randomized SVD) on $\mathbf{T}\mathbf{P}$ where $\mathbf{T} = \text{diag}(\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_s)$, so we first show that $\mathbf{T}\tilde{\mathbf{P}}$ enjoys this property. But now exact SVD is replaced with randomized SVD, for which we need to use the spectral error bound to argue that the error introduced is small. More precisely, for a matrix \mathbf{A} and its SVD truncation $\hat{\mathbf{A}}$ computed by randomized SVD, it is guaranteed that the spectral norm of $\mathbf{A} - \hat{\mathbf{A}}$ is small, then $\|(\mathbf{A} - \hat{\mathbf{A}})\mathbf{X}\|_F$ is small for any \mathbf{X} with small Frobenius norm, in particular, the orthonormal basis spanning a low dimensional subspace. This then suffices to guarantee $\mathbf{T}\tilde{\mathbf{P}}$ enjoys the close projection property. Given this, it suffices to show that $\tilde{\mathbf{P}}$ enjoys this property as $\mathbf{T}\tilde{\mathbf{P}}$, which follows from the definition of a subspace embedding. \square

Algorithm 2 Fast Distributed PCA for l_2 -Error Fitting

Input: $\{\mathbf{P}_i\}_{i=1}^s$; parameters t_1, t_2 for Algorithm `disPCA`; ℓ, q for randomized techniques.

1: **for** each node $v_i \in \mathcal{V}$ **do**

2: Compute subspace embedding $\mathbf{P}'_i = \mathbf{H}_i \mathbf{P}_i$.

3: **end for**

4: Run Algorithm `disPCA` on $\{\mathbf{P}'_i\}_{i=1}^s$ to get \mathbf{V} , where the SVD is randomized.

Output: \mathbf{V} .

4 Experiments

Our focus is to show the randomized techniques used in Algorithm 2 reduce the time taken significantly without compromising the quality of the solution. We perform experiments for three tasks: rank- r approximation, k -means clustering and principal component regression (PCR).

Datasets We choose the following real world datasets from UCI repository [1] for our experiments. For low rank approximation and k -means clustering, we choose two medium size datasets NewsGroups (18774×61188) and MNIST (70000×784), and two large-scale Bag-of-Words datasets: NYTimes news articles (BOWnytimes) (300000×102660) and PubMed abstracts (BOWpubmed) (8200000×141043). We use $r = 10$ for rank- r approximation and $k = 10$ for k -means clustering. For PCR, we use MNIST and further choose YearPredictionMSD (515345×90), CTslices (53500×386), and a large dataset MNIST8m (800000×784).

Experimental Methodology The algorithms are evaluated on a star network. The number of nodes is $s = 25$ for medium-size datasets, and $s = 100$ for the larger ones. We distribute the data over the nodes using a weighted partition, where each point is distributed to the nodes with probability proportional to the node’s weight chosen from the power law with parameter $\alpha = 2$.

For each projection dimension, we first construct the projected data using distributed PCA. For low rank approximation, we report the ratio between the cost of the obtained solution to that of the solution computed by SVD on the global data. For k -means, we run the algorithm in [3] (with Lloyd’s method as a subroutine) on the projected data to get a solution. Then we report the ratio between the cost of the above solution to that of a solution obtained by running Lloyd’s method directly on the global data. For PCR, we perform regression on the projected data to get a solution. Then we report the ratio between the error of the above solution to that of a solution obtained by PCR directly on the global data. We stop the algorithm if it takes more than 24 hours. For each projection dimension and each algorithm with randomness, the average ratio over 5 runs is reported.

Results Figure 2 shows the results for low rank approximation. We observe that the error of the fast distributed PCA is comparable to that of the exact solution computed directly on the global data. This is also observed for distributed PCA with one or none of subspace embedding and randomized SVD. Furthermore, the error of the fast PCA is comparable to that of normal PCA, which means that the speedup techniques merely affects the accuracy of the solution. The second row shows the computational time, which suggests a significant decrease in the time taken to run the fast distributed PCA. For example, on NewsGroups, the time of the fast distributed PCA improves over that of normal distributed PCA by a factor between 10 to 100. On the large dataset BOWpubmed, the normal PCA takes too long to finish and no results are presented, while the speedup versions produce good results in reasonable time. The use of the randomized techniques gives us a good performance improvement while keeping the solution quality almost the same.

Figure 3 and Figure 4 show the results for k -means clustering and PCR respectively. Similar to that for low rank approximation, we observe that the distributed solutions are almost as good as that computed directly on the global data, and the speedup merely affects the solution quality. We again observe a huge decrease in the running time by the speedup techniques.

Acknowledgments This work was supported in part by NSF grants CCF-0953192 and CCF-1101215, AFOSR grant FA9550-09-1-0538, ONR grant N00014-09-1-0751, a Google Research Award, and a Microsoft Research Faculty Fellowship. David Woodruff would like to acknowledge the XDATA program of the Defense Advanced Research Projects Agency (DARPA), administered through Air Force Research Laboratory contract FA8750-12-C0323, for supporting this work.

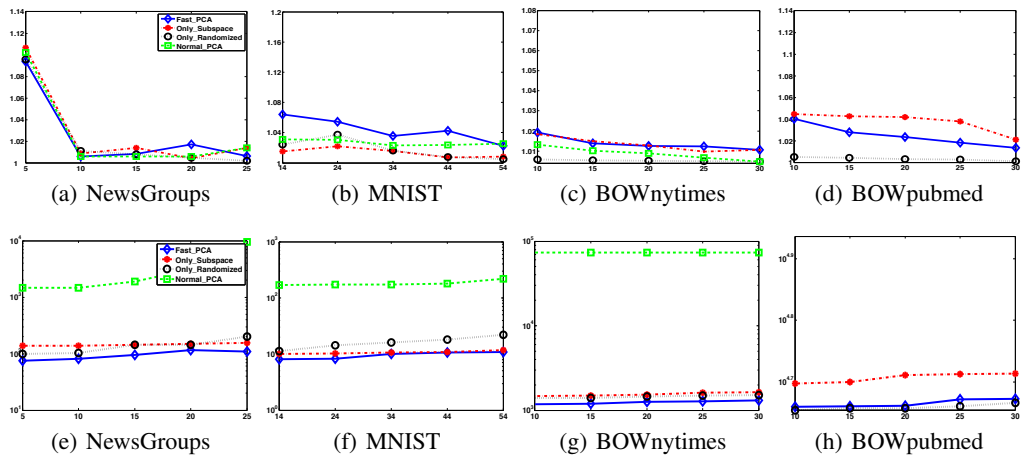


Figure 2: Low rank approximation. First row: error (normalized by baseline) v.s. projection dimension. Second row: time v.s. projection dimension.

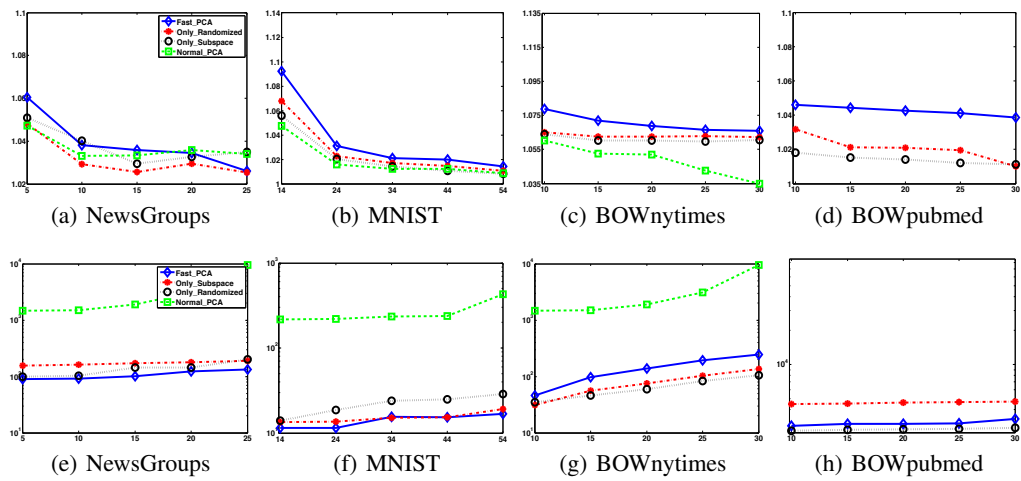


Figure 3: k -means clustering. First row: cost (normalized by baseline) v.s. projection dimension. Second row: time v.s. projection dimension.

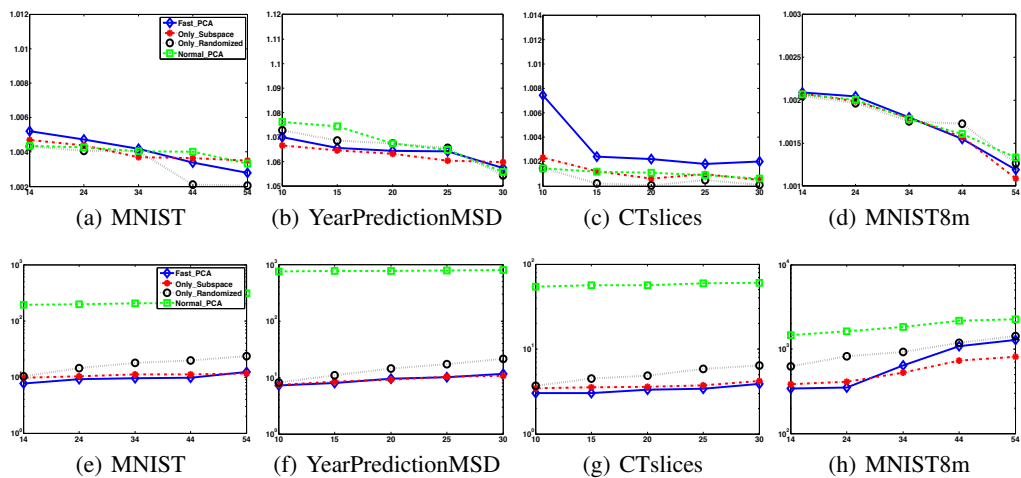


Figure 4: PCR. First row: error (normalized by baseline) v.s. projection dimension. Second row: time v.s. projection dimension.

References

- [1] K. Bache and M. Lichman. UCI machine learning repository, 2013.
- [2] Z.-J. Bai, R. H. Chan, and F. T. Luk. Principal component analysis for distributed data sets with updating. In *Proceedings of the International Conference on Advanced Parallel Processing Technologies*, 2005.
- [3] M.-F. Balcan, S. Ehrlich, and Y. Liang. Distributed k-means and k-median clustering on general communication topologies. In *Advances in Neural Information Processing Systems*, 2013.
- [4] D. M. Blei, A. Y. Ng, and M. I. Jordan. Latent dirichlet allocation. *the Journal of machine Learning research*, 2003.
- [5] C. Boutsidis, A. Zouzias, M. W. Mahoney, and P. Drineas. Stochastic dimensionality reduction for k-means clustering. *CoRR*, abs/1110.2897, 2011.
- [6] K. L. Clarkson and D. P. Woodruff. Low rank approximation and regression in input sparsity time. In *Proceedings of the 45th Annual ACM Symposium on Theory of Computing*, 2013.
- [7] J. C. Corbett, J. Dean, M. Epstein, A. Fikes, C. Frost, J. Furman, S. Ghemawat, A. Gubarev, C. Heiser, P. Hochschild, et al. Spanner: Googles globally-distributed database. In *Proceedings of the USENIX Symposium on Operating Systems Design and Implementation*, 2012.
- [8] D. Feldman and M. Langberg. A unified framework for approximating and clustering data. In *Proceedings of the Annual ACM Symposium on Theory of Computing*, 2011.
- [9] D. Feldman, M. Schmidt, and C. Sohler. Turning big data into tiny data: Constant-size coresets for k-means, pca and projective clustering. In *Proceedings of the Annual ACM-SIAM Symposium on Discrete Algorithms*, 2013.
- [10] M. Ghashami and J. M. Phillips. Relative errors for deterministic low-rank matrix approximations. *CoRR*, abs/1307.7454, 2013.
- [11] N. Halko, P.-G. Martinsson, and J. A. Tropp. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. *SIAM review*, 2011.
- [12] R. Kannan, S. Vempala, and D. Woodruff. Nimble algorithms for cloud computing. *arXiv preprint arXiv:1304.3162*, 2013.
- [13] N. Karampatziakis and P. Mineiro. Combining structured and unstructured randomness in large scale pca. *CoRR*, abs/1310.6304, 2013.
- [14] Y.-A. Le Borgne, S. Raybaud, and G. Bontempi. Distributed principal component analysis for wireless sensor networks. *Sensors*, 2008.
- [15] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. *Advances in Neural Information Processing Systems*, 2001.
- [16] S. V. Macua, P. Belanovic, and S. Zazo. Consensus-based distributed principal component analysis in wireless sensor networks. In *Proceedings of the IEEE International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*, 2010.
- [17] X. Meng and M. W. Mahoney. Low-distortion subspace embeddings in input-sparsity time and applications to robust linear regression. In *Proceedings of the Annual ACM symposium on Symposium on theory of computing*, 2013.
- [18] S. Mitra, M. Agrawal, A. Yadav, N. Carlsson, D. Eager, and A. Mahanti. Characterizing web-based video sharing workloads. *ACM Transactions on the Web*, 2011.
- [19] J. Nelson and H. L. Nguyễn. Osnap: Faster numerical linear algebra algorithms via sparser subspace embeddings. *arXiv preprint arXiv:1211.1002*, 2012.
- [20] C. Olston, J. Jiang, and J. Widom. Adaptive filters for continuous queries over distributed data streams. In *Proceedings of the ACM SIGMOD International Conference on Management of Data*, 2003.
- [21] Y. Qu, G. Ostrouchov, N. Samatova, and A. Geist. Principal component analysis for dimension reduction in massive distributed data sets. In *Proceedings of IEEE International Conference on Data Mining*, 2002.
- [22] T. Sarlós. Improved approximation algorithms for large matrices via random projections. In *FOCS*, pages 143–152, 2006.

A Related Work

A number of algorithms for approximate distributed PCA have been proposed [21, 2, 14, 16, 9], but either without theoretical guarantees, or without considering communication. [21] proposed an algorithm but provided no analysis on the tradeoff between communication and approximation. Most closely related to our work is [9], which observes that the top singular vectors of the local point set can be viewed as its summary and the union of the local summaries can be viewed as a summary of the global data, i.e., Algorithm `disPCA` discussed above.

In [12] the authors study algorithms in the arbitrary partition model in which each server holds a matrix \mathbf{P}_i and $\mathbf{P} = \sum_{i=1}^s \mathbf{P}_i$. Thus, each row of \mathbf{P} is additively shared across the s servers, whereas in our model each row of \mathbf{P} belongs to a single server, though duplicate rows are allowed. Our model is motivated by applications in which points are indecomposable entities. As our model is a special case of the arbitrary partition model, we can achieve more efficient algorithms. For instance, our distributed PCA algorithms provide much stronger guarantees, see, e.g., Lemma 4, which are needed for the downstream k -means application. Moreover, our k -means algorithms are more general, in the sense that they do not make a well-separability assumption, and more efficient in that the communication of [12] is $O(sd^2) + s(k/\epsilon)^{O(1)}$ words as opposed to our $O(sdk/\epsilon^2) + sk + (k/\epsilon)^{O(1)}$.

After the announce of this work, [Michael Cohen, Sam Elder, Cameron Musco, Christopher Musco, Madalina Persu] improve the guarantee for the k -means application in two ways. First, they tighten the result in [9], showing that projecting to just the $O(k/\epsilon)$ rather than $O(k/\epsilon^2)$ top singular vectors is sufficient to approximate k -means with $(1 + \epsilon)$ error. Second, they show that performing a Johnson-Lindenstrauss transformation down to $O(k/\epsilon^2)$ dimension gives $(1 + \epsilon)$ approximation without requiring a $\log(n)$ dependence. This can be used as a preprocessing step before our algorithm, replacing d with $O(k/\epsilon^2)$ in our communication bounds. They further show how to reduce the dimension to $O(k/\epsilon)$ using only $O(sk/\epsilon)$ vectors, but by a technique different from distributed PCA.

Other related work includes the recent [10] (see also the references therein), who give a deterministic streaming algorithm for low rank approximation in which each point of \mathbf{P} is seen one at a time and uses $O(dk/\epsilon)$ words of communication. Their algorithm naturally gives an $O(sdk/\epsilon)$ communication algorithm for low rank approximation in the distributed model. However, their algorithm for PCA doesn't satisfy the stronger guarantees of Lemma 4, and therefore it is unclear how to use it for k -means clustering. It also involves an SVD computation for each point, making the overall computation per server $O(n_i dr^2/\epsilon^2)$, which is slower than what we achieve, and it is not clear how their algorithm can exploit sparsity.

Speeding up large scale PCA using different versions of subspace embeddings was also considered in [13], though not in a distributed setting and not for ℓ_2 -error shape fitting problems. Also, their error guarantees are in terms of the r -th singular value gap, and are incomparable to ours.

B Guarantees for Distributed PCA

B.1 Proof of Lemma 1

We first prove a generalization of Lemma 1.

Lemma 7. *Let $\mathbf{A} \in \mathbb{R}^{n \times d}$ be an $n \times d$ matrix with singular value decomposition $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$. Let $\epsilon \in (0, 1]$ and $r, t \in \mathbb{N}_+$ with $d - 1 \geq t \geq r + \lceil r/\epsilon \rceil - 1$, and let $\widehat{\mathbf{A}} = \mathbf{A}\mathbf{V}^{(t)}(\mathbf{V}^{(t)})^\top$. Then for any matrix \mathbf{X} with d rows and $\|\mathbf{X}\|_F^2 \leq r$, we have*

$$\|(\mathbf{A} - \widehat{\mathbf{A}})\mathbf{X}\|_F^2 = \|\mathbf{A}\mathbf{X}\|_F^2 - \|\widehat{\mathbf{A}}\mathbf{X}\|_F^2 \leq \epsilon \sum_{i=r+1}^d \sigma_i^2(\mathbf{A}).$$

Proof. The proof follows the idea in the proof of Lemma 6.1 in [9].

For convenience, let $\overline{\Sigma}^{(t)}$ denote the diagonal matrix that contains the first t diagonal entries in Σ and is 0 otherwise. Then $\widehat{\mathbf{A}} = \mathbf{U}\overline{\Sigma}^{(t)}\mathbf{V}^\top$. We first have

$$\begin{aligned}\|\mathbf{A}\mathbf{X}\|_F^2 - \|\widehat{\mathbf{A}}\mathbf{X}\|_F^2 &= \|\mathbf{U}\Sigma\mathbf{V}^\top\mathbf{X}\|_F^2 - \|\mathbf{U}\overline{\Sigma}^{(t)}\mathbf{V}^\top\mathbf{X}\|_F^2 \\ &= \|\Sigma\mathbf{V}^\top\mathbf{X}\|_F^2 - \|\overline{\Sigma}^{(t)}\mathbf{V}^\top\mathbf{X}\|_F^2 \\ &= \|(\Sigma - \overline{\Sigma}^{(t)})\mathbf{V}^\top\mathbf{X}\|_F^2 \\ &= \|\mathbf{U}(\Sigma - \overline{\Sigma}^{(t)})\mathbf{V}^\top\mathbf{X}\|_F^2 \\ &= \|\mathbf{A}\mathbf{X} - \widehat{\mathbf{A}}\mathbf{X}\|_F^2.\end{aligned}$$

where the second and fourth equalities follow since \mathbf{U} has orthonormal columns, and the third equality follows since for $\mathbf{M} = \mathbf{V}^\top\mathbf{X}$ we have

$$\begin{aligned}\|\Sigma\mathbf{M}\|_F^2 - \|\overline{\Sigma}^{(t)}\mathbf{M}\|_F^2 &= \sum_{i=1}^d \sum_{j=1}^d \sigma_i^2(\mathbf{A})m_{ij}^2 - \sum_{i=1}^t \sum_{j=1}^d \sigma_i^2(\mathbf{A})m_{ij}^2 \\ &= \sum_{i=t+1}^d \sum_{j=1}^d \sigma_i^2(\mathbf{A})m_{ij}^2 = \|(\Sigma - \overline{\Sigma}^{(t)})\mathbf{M}\|_F^2.\end{aligned}$$

Next, we bound $\|\mathbf{A}\mathbf{X} - \widehat{\mathbf{A}}\mathbf{X}\|_F^2$. We have

$$\|\mathbf{A}\mathbf{X} - \widehat{\mathbf{A}}\mathbf{X}\|_F^2 = \|(\Sigma - \overline{\Sigma}^{(t)})\mathbf{V}^\top\mathbf{X}\|_F^2 \leq \|(\Sigma - \overline{\Sigma}^{(t)})\|_S^2 \|\mathbf{X}\|_F^2 = r\sigma_{t+1}^2(\mathbf{A})$$

where the inequality follows because the spectral norm is consistent with the Euclidean norm. This implies the lemma since

$$r\sigma_{t+1}^2(\mathbf{A}) \leq \epsilon(t-r+1)\sigma_{t+1}^2(\mathbf{A}) \leq \epsilon \sum_{i=r+1}^{t+1} \sigma_i^2(\mathbf{A}) \leq \epsilon \sum_{i=r+1}^d \sigma_i^2(\mathbf{A}). \quad (2)$$

where the first inequality follows for our choice of t . \square

Then Lemma 1 immediately follows from Lemma 7 since any $d \times r$ orthonormal matrix \mathbf{A} has $\|\mathbf{A}\|_F^2 \leq r$, and $\sum_{i=r+1}^d \sigma_i^2(\mathbf{A}) \leq d^2(\mathbf{A}, L_{\mathbf{X}})$ by the property of the singular value decomposition.

B.2 Proof of Theorem 2

Theorem 2. *Suppose Algorithm disPCA takes parameters $t_1 \geq r + \lceil 4r/\epsilon \rceil - 1$ and $t_2 = r$, and outputs $\mathbf{V}^{(r)}$. Then*

$$\|\mathbf{P} - \mathbf{P}\mathbf{V}^{(r)}(\mathbf{V}^{(r)})^\top\|_F^2 \leq (1 + \epsilon) \min_{\mathbf{X}} d^2(\mathbf{P}, L_{\mathbf{X}})$$

where the minimization is over $d \times r$ orthonormal matrices \mathbf{X} . The communication is $O(\frac{sr d}{\epsilon})$ words.

Proof. Let $\widehat{\mathbf{P}}_i := \mathbf{P}_i \mathbf{V}_i^{(t)} (\mathbf{V}_i^{(t)})^\top$, and let $\widehat{\mathbf{P}}$ be the concatenation of $\widehat{\mathbf{P}}_i$.

First, we show that $\widehat{\mathbf{P}}$ serves as a proxy of \mathbf{P} for optimizing $d^2(\mathbf{P}, L_{\mathbf{X}})$. By Pythagorean Theorem, for any orthonormal matrix \mathbf{X} of size $d \times r$,

$$\begin{aligned}d^2(\widehat{\mathbf{P}}, L_{\mathbf{X}}) - d^2(\mathbf{P}, L_{\mathbf{X}}) &= (\|\widehat{\mathbf{P}}\|_F^2 - \|\widehat{\mathbf{P}}\mathbf{X}\|_F^2) - (\|\mathbf{P}\|_F^2 - \|\mathbf{P}\mathbf{X}\|_F^2) \\ &= \Delta(\mathbf{X}) - c_0\end{aligned} \quad (3)$$

where $\Delta(\mathbf{X}) := \|\mathbf{P}\mathbf{X}\|_F^2 - \|\widehat{\mathbf{P}}\mathbf{X}\|_F^2$ and $c_0 := \|\mathbf{P}\|_F^2 - \|\widehat{\mathbf{P}}\|_F^2$. Since $\Delta(\mathbf{X})$ is small by Lemma 1 and c_0 is a constant, $\widehat{\mathbf{P}}$ approximates \mathbf{P} for optimizing $d^2(\mathbf{P}, L_{\mathbf{X}})$.

Next, we note that the optimal principal components for $\widehat{\mathbf{P}}$ are $\mathbf{V}^{(r)}$. This is because $\widehat{\mathbf{P}} = \widetilde{\mathbf{U}}\mathbf{Y}$ where $\widetilde{\mathbf{U}}$ is a block-diagonal matrix with blocks $\mathbf{U}_1, \dots, \mathbf{U}_s$, and thus the right singular vectors of \mathbf{Y} are also the right singular vectors of $\widehat{\mathbf{P}}$.

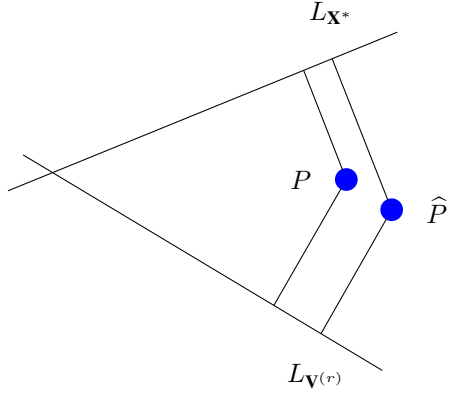


Figure 5: Illustration of low rank approximation.

Now, we are ready to bound $\|\mathbf{P} - \mathbf{P}\mathbf{V}^{(r)}(\mathbf{V}^{(r)})^\top\|_F^2 = d^2(\mathbf{P}, L_{\mathbf{V}^{(r)}})$. Suppose the r optimal loadings for \mathbf{P} are \mathbf{X}^* . See Figure 5 for an illustration. Then

$$\begin{aligned} \|\mathbf{P} - \mathbf{P}\mathbf{V}^{(r)}(\mathbf{V}^{(r)})^\top\|_F^2 &= d^2(\hat{\mathbf{P}}, L_{\mathbf{V}^{(r)}}) + c_0 - \Delta(\mathbf{V}^{(r)}) \\ &\leq d^2(\hat{\mathbf{P}}, L_{\mathbf{X}^*}) + c_0 - \Delta(\mathbf{V}^{(r)}) \\ &= d^2(\mathbf{P}, L_{\mathbf{X}^*}) + \Delta(\mathbf{X}^*) - \Delta(\mathbf{V}^{(r)}) \end{aligned} \quad (4)$$

where the first and third line follow from (3) and the second follows from the fact that $\mathbf{V}^{(r)}$ are the optimal principal loadings for $\hat{\mathbf{P}}$. By Lemma 1, $\Delta(\mathbf{V}^{(r)}) \geq 0$ and $\Delta(\mathbf{X}^*) \leq \epsilon d^2(\mathbf{P}, L_{\mathbf{X}^*})$. Combining these with (4) leads to the theorem. \square

Note A refinement of the proof of Lemma 1 leads to the following data dependent bound.

Lemma 8. *The statement in Lemma 7 holds if $t > \tau(\mathbf{A}, r, \epsilon)$ where*

$$\tau(\mathbf{A}, r, \epsilon) := \operatorname{argmin}_t \left\{ \sigma_t^2(\mathbf{A}) \leq \frac{\epsilon}{r} \sum_{i>r} \sigma_i^2(\mathbf{A}) \right\}.$$

Furthermore, $\tau(\mathbf{A}, r, \epsilon) = O(\frac{r}{\epsilon})$.

Proof. Note that the bound on t is only used in proving (2), for which $t > \tau(\mathbf{A}, r, \epsilon)$ suffices. $\tau(\mathbf{A}, r, \epsilon) = O(\frac{r}{\epsilon})$ follows by definition. \square

Theorem 9. *Suppose Algorithm disPCA takes parameters $t_1 \geq \max_i \tau(\mathbf{P}_i, r, \epsilon)$ and $t_2 = r$, and outputs $\mathbf{V}^{(r)}$. Then*

$$\|\mathbf{P} - \mathbf{P}\mathbf{V}^{(r)}(\mathbf{V}^{(r)})^\top\|_F^2 \leq (1 + \epsilon) \min_{\mathbf{X}} d^2(\mathbf{P}, L_{\mathbf{X}})$$

where the minimization is over orthonormal matrices $\mathbf{X} \in \mathbb{R}^{d \times r}$. The total communication cost is $O(sd \max_i \tau(\mathbf{P}_i, r, \epsilon))$ words.

$\tau(\mathbf{P}_i, r, \epsilon)$ is typically much less than $O(r/\epsilon)$ in practice. This provides an explanation for the fact that t_1 much smaller than $O(r/\epsilon)$ can still lead to good solution for many practical instances. Similar data dependent bounds can be derived for the other theorems in our paper.

C Guarantees for Distributed ℓ_2 -Error Fitting

C.1 Proof of Lemma 4

We first introduce some intermediate variables for our analysis. Imagine we perform two projections: first project \mathbf{P}_i to $\hat{\mathbf{P}}_i = \mathbf{P}_i \mathbf{V}_i^{(t)} (\mathbf{V}_i^{(t)})^\top$, then project $\hat{\mathbf{P}}_i$ to $\bar{\mathbf{P}}_i = \hat{\mathbf{P}}_i \mathbf{V}^{(t)} (\mathbf{V}^{(t)})^\top$ where

$t = t_1 = t_2$. Let $\widehat{\mathbf{P}}$ denote the vertical concatenation of $\widehat{\mathbf{P}}_i$ and let $\overline{\mathbf{P}}$ denote the vertical concatenation of $\overline{\mathbf{P}}_i$, i.e.

$$\widehat{\mathbf{P}} = \begin{bmatrix} \widehat{\mathbf{P}}_1 \\ \vdots \\ \widehat{\mathbf{P}}_s \end{bmatrix} \quad \text{and} \quad \overline{\mathbf{P}} = \begin{bmatrix} \overline{\mathbf{P}}_1 \\ \vdots \\ \overline{\mathbf{P}}_s \end{bmatrix}$$

Lemma 4. *Let $t_1 = t_2 \geq k + \lceil 8k/\epsilon \rceil - 1$ in Algorithm disPCA for $k \in \mathbb{N}_+$ and $\epsilon \in (0, 1)$. Then for any $d \times k$ matrix \mathbf{X} with orthonormal columns,*

$$\begin{aligned} 0 &\leq \|\mathbf{P}\mathbf{X} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2 &&\leq \epsilon d^2(\mathbf{P}, L_{\mathbf{X}}), \\ 0 &\leq \|\mathbf{P}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 &&\leq \epsilon d^2(\mathbf{P}, L_{\mathbf{X}}). \end{aligned}$$

Proof. For the first statement, we have

$$\|\mathbf{P}\mathbf{X} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2 \leq 2\|\mathbf{P}\mathbf{X} - \widehat{\mathbf{P}}\mathbf{X}\|_F^2 \quad (5)$$

$$+ 2\|\widehat{\mathbf{P}}\mathbf{X} - \overline{\mathbf{P}}\mathbf{X}\|_F^2 \quad (6)$$

$$+ 2\|\overline{\mathbf{P}}\mathbf{X} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2. \quad (7)$$

For (5), we have by Lemma 7

$$\|\mathbf{P}\mathbf{X} - \widehat{\mathbf{P}}\mathbf{X}\|_F^2 = \sum_{i=1}^s \|\mathbf{P}_i\mathbf{X} - \widehat{\mathbf{P}}_i\mathbf{X}\|_F^2 \leq \sum_{i=1}^s \frac{\epsilon}{4} d^2(\mathbf{P}_i, L_{\mathbf{X}}) = \frac{\epsilon}{8} d^2(\mathbf{P}, L_{\mathbf{X}}). \quad (8)$$

Similarly, for (6) we have by Lemma 7

$$\|\widehat{\mathbf{P}}\mathbf{X} - \overline{\mathbf{P}}\mathbf{X}\|_F^2 \leq \frac{\epsilon}{8} d^2(\widehat{\mathbf{P}}, L_{\mathbf{X}}). \quad (9)$$

To bound (7), let $\mathbf{Y} = \mathbf{V}^{(t)}(\mathbf{V}^{(t)})^\top \mathbf{X}$. Then by definition, $\overline{\mathbf{P}}_i\mathbf{X} = \widehat{\mathbf{P}}_i\mathbf{Y}$ and $\tilde{\mathbf{P}}_i\mathbf{X} = \mathbf{P}_i\mathbf{Y}$. By Lemma 7, we have

$$\|\overline{\mathbf{P}}\mathbf{X} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2 = \sum_{i=1}^s \|\widehat{\mathbf{P}}_i\mathbf{Y} - \mathbf{P}_i\mathbf{Y}\|_F^2 \quad (10)$$

$$\leq \sum_{i=1}^s \frac{\epsilon}{8} \sum_{i=r+1}^s \sigma_i^2(\mathbf{P}_i) \leq \frac{\epsilon}{8} \sum_{i=1}^s d^2(\mathbf{P}_i, L_{\mathbf{X}}) = \frac{\epsilon}{8} d^2(\mathbf{P}, L_{\mathbf{X}}). \quad (11)$$

Combining (8)(9) and (11) leads to

$$\|\mathbf{P}\mathbf{X} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2 \leq \frac{\epsilon}{2} d^2(\mathbf{P}, L_{\mathbf{X}}) + \frac{\epsilon}{4} d^2(\widehat{\mathbf{P}}, L_{\mathbf{X}}). \quad (12)$$

We now only need to bound $d^2(\widehat{\mathbf{P}}, L_{\mathbf{X}})$ is similar to $d^2(\mathbf{P}, L_{\mathbf{X}})$, which is done in Lemma 10. The first statement then follows.

For the second statement, we have a similar argument.

$$\|\mathbf{P}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 = \|\mathbf{P}\mathbf{X}\|_F^2 - \|\widehat{\mathbf{P}}\mathbf{X}\|_F^2 \quad (13)$$

$$+ \|\widehat{\mathbf{P}}\mathbf{X}\|_F^2 - \|\overline{\mathbf{P}}\mathbf{X}\|_F^2 \quad (14)$$

$$+ \|\overline{\mathbf{P}}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2. \quad (15)$$

For (13), we have by Lemma 7

$$\|\mathbf{P}\mathbf{X}\|_F^2 - \|\widehat{\mathbf{P}}\mathbf{X}\|_F^2 = \sum_{i=1}^s \left[\|\mathbf{P}_i\mathbf{X}\|_F^2 - \|\widehat{\mathbf{P}}_i\mathbf{X}\|_F^2 \right] \leq \sum_{i=1}^s \frac{\epsilon}{4} d^2(\mathbf{P}_i, L_{\mathbf{X}}) = \frac{\epsilon}{4} d^2(\mathbf{P}, L_{\mathbf{X}}). \quad (16)$$

Similarly, for (14) we have by Lemma 7

$$\|\widehat{\mathbf{P}}\mathbf{X}\|_F^2 - \|\overline{\mathbf{P}}\mathbf{X}\|_F^2 \leq \frac{\epsilon}{4} d^2(\widehat{\mathbf{P}}, L_{\mathbf{X}}). \quad (17)$$

By Lemma 7, we have

$$\begin{aligned}\|\bar{\mathbf{P}}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 &= \sum_{i=1}^s \left[\|\hat{\mathbf{P}}_i \mathbf{Y}\|_F^2 - \|\mathbf{P}_i \mathbf{Y}\|_F^2 \right] \\ &\leq \sum_{i=1}^s \frac{\epsilon}{4} \sum_{i=r+1}^s \sigma_i^2(\mathbf{P}_i) \leq \frac{\epsilon}{4} \sum_{i=1}^s d^2(\mathbf{P}_i, L_{\mathbf{X}}) = \frac{\epsilon}{4} d^2(\mathbf{P}, L_{\mathbf{X}}).\end{aligned}\quad (18)$$

Combining (16)(17) and (18) leads to

$$\|\mathbf{P}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \leq \frac{\epsilon}{2} d^2(\mathbf{P}, L_{\mathbf{X}}) + \frac{\epsilon}{4} d^2(\hat{\mathbf{P}}, L_{\mathbf{X}}).\quad (19)$$

The second statement then follows from (19) and Lemma 10. \square

The following is a technical lemma that will be used in the proof of Lemma 4.

Lemma 10.

$$d^2(\hat{\mathbf{P}}, L_{\mathbf{X}}) \leq (1 + \epsilon) d^2(\mathbf{P}, L_{\mathbf{X}}).$$

Proof. We have

$$\begin{aligned}d^2(\hat{\mathbf{P}}, L_{\mathbf{X}}) - d^2(\mathbf{P}, L_{\mathbf{X}}) &= \|\hat{\mathbf{P}} - \hat{\mathbf{P}}\mathbf{X}\mathbf{X}^\top\|_F^2 - \|\mathbf{P} - \mathbf{P}\mathbf{X}\mathbf{X}^\top\|_F^2 \\ &= \|\hat{\mathbf{P}}\|_F^2 - \|\hat{\mathbf{P}}\mathbf{X}\mathbf{X}^\top\|_F^2 - (\|\mathbf{P}\|_F^2 - \|\mathbf{P}\mathbf{X}\mathbf{X}^\top\|_F^2) \\ &= \sum_{i=1}^s \left[\|\hat{\mathbf{P}}_i\|_F^2 - \|\mathbf{P}_i\|_F^2 \right] + \sum_{i=1}^s \left[\|\mathbf{P}_i\mathbf{X}\mathbf{X}^\top\|_F^2 - \|\hat{\mathbf{P}}_i\mathbf{X}\mathbf{X}^\top\|_F^2 \right].\end{aligned}$$

By the Pythagorean Theorem, $\|\hat{\mathbf{P}}_i\|_F^2 \leq \|\mathbf{P}_i\|_F^2$. Also, since \mathbf{X} is orthonormal, $\|\mathbf{P}_i\mathbf{X}\mathbf{X}^\top\|_F^2 = \|\mathbf{P}_i\mathbf{X}\|_F^2$ and $\|\hat{\mathbf{P}}_i\mathbf{X}\mathbf{X}^\top\|_F^2 = \|\hat{\mathbf{P}}_i\mathbf{X}\|_F^2$. Then

$$d^2(\hat{\mathbf{P}}, L_{\mathbf{X}}) - d^2(\mathbf{P}, L_{\mathbf{X}}) \leq \sum_{i=1}^s \left[\|\mathbf{P}_i\mathbf{X}\|_F^2 - \|\hat{\mathbf{P}}_i\mathbf{X}\|_F^2 \right] \leq \sum_{i=1}^s \epsilon d^2(\mathbf{P}_i, L_{\mathbf{X}}) = \epsilon d^2(\mathbf{P}, L_{\mathbf{X}})\quad (20)$$

where the second inequality follows from Lemma 1. \square

C.2 Proof of Theorem 3

The following weak triangle inequality is useful for our analysis.

Fact 1. For any $a, b \in \mathbb{R}$ and $\epsilon \in (0, 1)$, $|a^2 - b^2| \leq \frac{3(a-b)^2}{\epsilon} + 2\epsilon a^2$.

Proof. Either $|a| \leq \frac{|a-b|}{\epsilon}$ or $|a-b| \leq \epsilon|a|$, so we have $|a||a-b| \leq \frac{(a-b)^2}{\epsilon} + \epsilon a^2$. This leads to

$$|a^2 - b^2| = |a-b||a+b| \leq |a-b|(|2a| + |b-a|) = 2|a||a-b| + (a-b)^2 \leq \frac{2(a-b)^2}{\epsilon} + 2\epsilon a^2 + (a-b)^2$$

which completes the proof. \square

We first prove the theorem for the special case of k -means clustering, and the same argument leads to the guarantee for general l_2 -error fitting problems.

Theorem 11. Let $t_1 = t_2 \geq k + \lceil 4k/\epsilon^2 \rceil - 1$ in Algorithm disPCA. Then there exists a constant $c_0 \geq 0$, such that for any set of k points \mathcal{L} ,

$$(1 - \epsilon) d^2(\mathbf{P}, \mathcal{L}) \leq d^2(\tilde{\mathbf{P}}, \mathcal{L}) + c_0 \leq (1 + \epsilon) d^2(\mathbf{P}, \mathcal{L}).$$

Algorithm 3 Fast Sparse Subspace Embedding [6]

Input: parameters $n, \ell \in \mathbb{N}_+$.

- 1: Let $h : [n] \mapsto [\ell]$ be a random map, so that for each $i \in [n]$, $h(i) = j$ for $j \in [\ell]$ with probability $1/\ell$.
- 2: Let Φ be an $\ell \times n$ binary matrix with $\Phi_{h(i),i} = 1$, and all remaining entries 0.
- 3: Let Σ be an $n \times n$ diagonal matrix, with each diagonal entry independently chosen as $+1$ or -1 with equal probability.

Output: $\mathbf{H} = \Phi \Sigma$.

Proof. The proof follows that in [9], with slight modification for the distributed setting.

Let $\mathbf{X} \in \mathbb{R}^{d \times k}$ has orthonormal columns that span \mathcal{L} . Let \tilde{p}_i be the point in $\tilde{\mathbf{P}}$ corresponding to p_i in \mathbf{P} . Let $c_0 = \|\mathbf{P}\|_F^2 - \|\tilde{\mathbf{P}}\|_F^2$. Then by Pythagorean theorem we have

$$|d^2(\mathbf{P}, \mathcal{L}) - d^2(\tilde{\mathbf{P}}, \mathcal{L}) - c_0| \leq \left| d^2(\mathbf{P}, L(\mathbf{X})) - d^2(\tilde{\mathbf{P}}, L_{\mathbf{X}}) - c_0 \right| + \left| \sum_{i=1}^{|\mathbf{P}|} [d(\pi_{\mathbf{X}}(p_i), \mathcal{L})^2 - d(\pi_{\mathbf{X}}(\tilde{p}_i), \mathcal{L})^2] \right|.$$

For the first part, we have by Pythagorean theorem

$$d^2(\mathbf{P}, L(\mathbf{X})) - d^2(\tilde{\mathbf{P}}, L_{\mathbf{X}}) - c_0 = (\|\mathbf{P}\|_F^2 - \|\mathbf{P}\mathbf{X}\|_F^2) - (\|\tilde{\mathbf{P}}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2) - c_0 = \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 - \|\mathbf{P}\mathbf{X}\|_F^2. \quad (21)$$

For the second part, by Fact 1 we have

$$\begin{aligned} \sum_{i=1}^{|\mathbf{P}|} |d(\pi_{\mathbf{X}}(p_i), \mathcal{L})^2 - d(\pi_{\mathbf{X}}(\tilde{p}_i), \mathcal{L})^2| &\leq \sum_{i=1}^{|\mathbf{P}|} \left[\frac{12d(\pi_{\mathbf{X}}(p_i), \pi_{\mathbf{X}}(\tilde{p}_i))^2}{\epsilon} + \frac{\epsilon}{2} d(\pi_{\mathbf{X}}(p_i), \mathcal{L})^2 \right] \\ &= \frac{12}{\epsilon} \|(\mathbf{P} - \tilde{\mathbf{P}})\mathbf{X}\|_F^2 + \frac{\epsilon}{2} \sum_{i=1}^{|\mathbf{P}|} d(\pi_{\mathbf{X}}(p_i), \mathcal{L})^2 \\ &\leq \frac{12}{\epsilon} \|(\mathbf{P} - \tilde{\mathbf{P}})\mathbf{X}\|_F^2 + \frac{\epsilon}{2} \sum_{i=1}^{|\mathbf{P}|} d(p_i, \mathcal{L})^2. \end{aligned} \quad (22)$$

Combining (21)(22) with Lemma 4 leads to the theorem, since $d^2(\mathbf{P}, L_{\mathbf{X}}) \leq d^2(\mathbf{P}, \mathcal{L})$. \square

The general statement for ℓ_2 -error geometric fitting problems follows from the same argument.

Theorem 3. *Let $t_1 = t_2 = O(rk/\epsilon^2)$ in Algorithm disPCA for $\epsilon \in (0, 1/3)$. Then there exists a constant $c_0 \geq 0$ such that for any set of k centers \mathcal{L} in r -Subspace k -Clustering,*

$$(1 - \epsilon)d^2(\mathbf{P}, \mathcal{L}) \leq d^2(\tilde{\mathbf{P}}, \mathcal{L}) + c_0 \leq (1 + \epsilon)d^2(\mathbf{P}, \mathcal{L}).$$

D Fast Distributed PCA

D.1 Proofs for Subspace Embedding

The construction of the embedding matrix \mathbf{H} is presented in Algorithm 3. Note that the embedding matrix \mathbf{H} does not need to be built explicitly; we can compute the embedding $\mathbf{H}\mathbf{A}$ for an given matrix \mathbf{A} in a direct and faster way. Algorithm 3 has the following guarantee.

Theorem 12. [6, 17, 19] *Suppose $n > d$ and $\ell = O(\frac{d^2}{\epsilon^2})$. With probability at least 99/100, $\|\mathbf{H}\mathbf{A}y\|_2 = (1 \pm \epsilon)\|\mathbf{A}y\|_2$ for all vectors $y \in \mathbb{R}^d$. Moreover, $\mathbf{H}\mathbf{A}$ can be computed in time $O(\text{nnz}(\mathbf{A}))$ where $\text{nnz}(\mathbf{A})$ is the number of non-zero entries in \mathbf{A} .*

Lemma 13. *Let $\epsilon \in (0, 1/2]$ and $k, t \in \mathbb{N}_+$ with $d - 1 \geq t \geq k + \lceil 4k/\epsilon \rceil - 1$. Suppose Algorithm disPCA takes input $\{\mathbf{H}_i \mathbf{P}_i\}_{i=1}^s$ and outputs $\mathbf{V}^{(t)}$. Let $\tilde{\mathbf{P}} = \mathbf{P}\mathbf{V}^{(t)}(\mathbf{V}^{(t)})^\top$. Then for any $d \times k$ matrix \mathbf{X} with orthonormal columns,*

$$\begin{aligned} \|\mathbf{P}\mathbf{X} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2 &\leq \epsilon d^2(\mathbf{P}, L_{\mathbf{X}}), \\ \left| \|\mathbf{P}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \right| &\leq 3\epsilon \|\mathbf{P}\mathbf{X}\|_F^2 + \epsilon d^2(\mathbf{P}, L_{\mathbf{X}}). \end{aligned}$$

Algorithm 4 Boosting success probability of embedding

Input: $\mathbf{A} \in \mathbb{R}^{n \times d}$, parameters ϵ, δ .

- 1: Construct $r = O(\log \frac{1}{\delta})$ independent subspace embeddings $\mathbf{H}_j \mathbf{A}$, each having accuracy $\epsilon/9$ and success probability $99/100$.
- 2: Compute SVD $\mathbf{H}_j \mathbf{A} = \mathbf{U}_j \mathbf{\Sigma}_j \mathbf{V}_j^\top$ for $j \in [r]$.
- 3: **for** $j \in [r]$ **do**
- 4: Check if for at least half $j' \neq j$,

$$\sigma_i(\mathbf{\Sigma}_{j'} \mathbf{V}_{j'}^\top \mathbf{V}_j \mathbf{\Sigma}_j^{-1}) \in [1 \pm \epsilon/3], \forall i.$$

- 5: If so, output $\mathbf{H}_j \mathbf{A}$.
 - 6: **end for**
-

Proof. First note that the input to Algorithm disPCA is \mathbf{TP} where \mathbf{T} is a block-diagonal matrix with blocks $\mathbf{H}_1, \dots, \mathbf{H}_s$. Then the projection of the input to $\mathbf{V}^{(t)}$ is $\mathbf{TPV}^{(t)}(\mathbf{V}^{(t)})^\top = \mathbf{TP}\tilde{\mathbf{P}}$. By Lemma 4, for any $d \times k$ matrix \mathbf{X} with orthonormal columns, we have

$$0 \leq \|\mathbf{TPX} - \mathbf{TP}\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \leq \frac{\epsilon}{4} d^2(\mathbf{TP}, L_{\mathbf{X}}), \quad (23)$$

$$0 \leq \|\mathbf{TPX}\|_F^2 - \|\mathbf{TP}\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \leq \frac{\epsilon}{4} d^2(\mathbf{TP}, L_{\mathbf{X}}). \quad (24)$$

By properties of \mathbf{T} , we have

$$\|\mathbf{TPX} - \mathbf{TP}\tilde{\mathbf{P}}\mathbf{X}\|_F^2 = \|\mathbf{T}(\mathbf{PX} - \tilde{\mathbf{P}}\mathbf{X})\|_F^2 \geq (1 - \epsilon)\|\mathbf{PX} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2$$

and

$$d^2(\mathbf{TP}, L_{\mathbf{X}}) = \|\mathbf{TP} - \mathbf{TPXX}^\top\|_F^2 \leq (1 + \epsilon)\|\mathbf{P} - \mathbf{PXX}^\top\|_F^2 = (1 + \epsilon)d^2(\mathbf{P}, L_{\mathbf{X}}).$$

Combined with (23), these lead to the first claim.

Similarly, we also have $\|\mathbf{TPX}\|_F^2 = (1 \pm \epsilon)\|\mathbf{PX}\|_F^2$ and $\|\mathbf{TP}\tilde{\mathbf{P}}\mathbf{X}\|_F^2 = (1 \pm \epsilon)\|\tilde{\mathbf{P}}\mathbf{X}\|_F^2$. Plugging these into (24), we obtain

$$-3\epsilon\|\mathbf{PX}\|_F^2 \leq \|\mathbf{PX}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \leq 3\epsilon\|\mathbf{PX}\|_F^2 + \epsilon d^2(\mathbf{P}, L_{\mathbf{X}})$$

which establishes the lemma. \square

Theorem 14. *Algorithm 4 outputs a subspace embedding with probability at least $1 - \delta$. In expectation Step 3 is run only a constant number of times with expected time $O(d^3 r^2 / \epsilon^2)$.*

Proof. For each j , $\mathbf{H}_j \mathbf{A}$ succeeds with probability $99/100$, meaning that for all x we have $\|\mathbf{H}_j \mathbf{A} x\|_2 = (1 \pm \epsilon/9)\|\mathbf{A} x\|_2$. Suppose for some $j \neq j'$, $\mathbf{H}_j \mathbf{A}$ and $\mathbf{H}_{j'} \mathbf{A}$ are both successful. By definition we have

$$\|\mathbf{H}_j \mathbf{A} x\|_2 = (1 \pm \epsilon/3)\|\mathbf{H}_{j'} \mathbf{A} x\|_2$$

for all x . Taking the SVD of the embeddings, this is equivalent to

$$\|\mathbf{\Sigma}_j \mathbf{V}_j^\top x\|_2 = (1 \pm \epsilon/3)\|\mathbf{\Sigma}_{j'} \mathbf{V}_{j'}^\top x\|_2$$

for all x . Making the change of variable $y := \mathbf{\Sigma}_j \mathbf{V}_j^\top x$, this is equivalent to

$$\|y\|_2 = (1 \pm \epsilon/3)\|\mathbf{\Sigma}_{j'} \mathbf{V}_{j'}^\top \mathbf{V}_j \mathbf{\Sigma}_j^{-1} y\|_2$$

for all y , which is true if and only if all singular values of $\mathbf{\Sigma}_{j'} \mathbf{V}_{j'}^\top \mathbf{V}_j \mathbf{\Sigma}_j^{-1}$ are in $[1 - \epsilon/3, 1 + \epsilon/3]$.

Conversely, if all singular values of $\mathbf{\Sigma}_{j'} \mathbf{V}_{j'}^\top \mathbf{V}_j \mathbf{\Sigma}_j^{-1}$ are in $[1 - \epsilon/3, 1 + \epsilon/3]$, one can trace the steps backward to conclude that $\|\mathbf{H}_j \mathbf{A} x\|_2 = (1 \pm \epsilon/3)\|\mathbf{H}_{j'} \mathbf{A} x\|_2$ for all x .

Since with probability at least $1 - \delta$, a $9/10$ fraction of the embeddings succeed with accuracy $\epsilon/9$, there exists a j that can pass the test. It follows that any index j which passes the test in the algorithm with a majority of the $j' \neq j$ is a successful subspace embedding with accuracy ϵ .

Algorithm 5 Randomized SVD [11]

Input: matrix $\mathbf{A} \in \mathbb{R}^{\ell \times d}$; parameters $t, q \in \mathbb{N}_+$.

1: \triangleright Stage A

2: Generate an $\ell \times 2t$ Gaussian test matrix Ω .

3: Set $\mathbf{Y} = (\mathbf{A}^\top \mathbf{A})^q \mathbf{A}^\top \Omega$, and compute QR-factorization: $\mathbf{Y} = \mathbf{Q}\mathbf{R}$.

4: \triangleright Stage B

5: Set $\mathbf{B} = \mathbf{A}\mathbf{Q}$, and compute SVD: $\mathbf{B} = \mathbf{U}\Sigma\tilde{\mathbf{V}}^\top$.

6: Set $\mathbf{V} = \mathbf{Q}\tilde{\mathbf{V}}$.

Output: Σ, \mathbf{V} .

Moreover, if we choose a random j to compare to the remaining j' , the expected number of choices of j until the test passes is only constant. Then finding the index j only takes an expected $O(r)$ SVDs.

The time to do the SVD naively is $O(d^4/\epsilon^2)$. We can improve this by letting \mathbf{T} be a fast Johnson-Lindenstrauss transform matrix of dimension $O(dr/\epsilon^2) \times O(d^2/\epsilon^2)$, then we can replace $\mathbf{H}_j \mathbf{A}$ with $\mathbf{T}\mathbf{H}_j \mathbf{A}$ for all $j \in [d]$. Then the verification procedure would only take $O(d^3 r^2/\epsilon^2)$ time. \square

D.2 Proofs for Randomized SVD

The details of randomized SVD are presented in Algorithm 5, rephrased in our notations. We have the following analog of Lemma 1.

Lemma 15. *Let $\mathbf{A} \in \mathbb{R}^{\ell \times d}$ be an $\ell \times d$ matrix ($\ell > d$). Let $\epsilon \in (0, 1]$, $k, t \in \mathbb{N}_+$ with $d - 1 \geq t \geq k + \lceil 6k/\epsilon^2 \rceil - 1$. Let $\hat{\mathbf{A}} = \mathbf{A}\mathbf{V}\mathbf{V}^\top$ where \mathbf{V} is computed by Algorithm 5 with $q = O(\log \max\{\ell, d\})$. Then with probability at least $1 - 3e^{-t}$, for any matrix \mathbf{X} with d rows and $\|\mathbf{X}\|_F^2 \leq k$, we have*

$$\begin{aligned} \|(\mathbf{A} - \hat{\mathbf{A}})\mathbf{X}\|_F^2 &\leq \frac{\epsilon^2}{3} \sum_{i=k+1}^d \sigma_i^2(\mathbf{A}), \\ \|\mathbf{A}\mathbf{X}\|_F^2 - \|\hat{\mathbf{A}}\mathbf{X}\|_F^2 &\leq \epsilon \sum_{i=k+1}^d \sigma_i^2(\mathbf{A}) + 2\epsilon \|\mathbf{A}\mathbf{X}\|_F^2. \end{aligned}$$

The algorithm runs in time $O(qt\ell d + t^2(\ell + d))$.

Proof. As stated in Section 10.4 in [11], with probability at least $1 - 3e^{-t}$, we have

$$\|\mathbf{A} - \hat{\mathbf{A}}\|_S \leq 2\sigma_{t+1}(\mathbf{A}). \quad (25)$$

Then we have

$$\|(\mathbf{A} - \hat{\mathbf{A}})\mathbf{X}\|_F^2 \leq \|\mathbf{X}\|_F^2 \|\mathbf{A} - \hat{\mathbf{A}}\|_S^2 \leq 2k\sigma_{t+1}^2(\mathbf{A})$$

where the first inequality follows because the spectral norm is consistent with the Euclidean norm, and the second inequality follows from (25). For our choice of t , we have

$$k\sigma_{t+1}^2(\mathbf{A}) \leq \frac{\epsilon^2}{6}(t - k + 1)\sigma_{t+1}^2(\mathbf{A}) \leq \frac{\epsilon^2}{6} \sum_{i=k+1}^{t+1} \sigma_i^2(\mathbf{A}) \leq \frac{\epsilon^2}{6} \sum_{i=k+1}^d \sigma_i^2(\mathbf{A}) \leq \frac{\epsilon^2}{6} d^2(\mathbf{A}, L_{\mathbf{X}}),$$

which leads to the first claim in the lemma.

To prove the second claim, first note that

$$\|\|\mathbf{A}\mathbf{X}\|_F - \|\hat{\mathbf{A}}\mathbf{X}\|_F\|^2 \leq \|(\mathbf{A} - \hat{\mathbf{A}})\mathbf{X}\|_F^2 \leq \frac{\epsilon^2}{3} d^2(\mathbf{A}, L_{\mathbf{X}}).$$

Then by Fact 1, we have

$$\|\|\mathbf{A}\mathbf{X}\|_F^2 - \|\hat{\mathbf{A}}\mathbf{X}\|_F^2\| \leq \frac{3}{\epsilon} \|\|\mathbf{A}\mathbf{X}\|_F - \|\hat{\mathbf{A}}\mathbf{X}\|_F\|^2 + 2\epsilon \|\mathbf{A}\mathbf{X}\|_F^2 \leq \epsilon d^2(\mathbf{A}, L_{\mathbf{X}}) + 2\epsilon \|\mathbf{A}\mathbf{X}\|_F^2$$

which completes the proof. \square

D.3 Proof of Theorem 6

Let \mathbf{T} to be a diagonal block matrix with $\mathbf{H}_1, \mathbf{H}_2, \dots, \mathbf{H}_s$ on the diagonal. Then Algorithm 2 is just to run Algorithm `disPCA` on \mathbf{TP} to get the principal components \mathbf{V} . Recall that the goal is to show $\tilde{\mathbf{P}} = \mathbf{P}\mathbf{V}\mathbf{V}^\top$ is a good proxy for the original data \mathbf{P} with respect to ℓ_2 error fitting problems. It suffices to show that $\tilde{\mathbf{P}}$ satisfies enjoys properties similar to those stated in Lemma 4.

To prove this, we begin with a lemma saying that $\mathbf{T}\tilde{\mathbf{P}}$ enjoys such properties, i.e. such properties are approximately preserved when replacing exact SVD with randomized SVD in Algorithm `disPCA` (Lemma 16). Then we can show that $\tilde{\mathbf{P}}$ enjoys similar properties as $\mathbf{T}\tilde{\mathbf{P}}$, i.e. these properties are approximately preserved under subspace embedding (Lemma 18).

Lemma 16. *For any $d \times k$ matrix \mathbf{X} with orthonormal columns,*

$$\begin{aligned} \|\mathbf{TPX} - \mathbf{T}\tilde{\mathbf{P}}\mathbf{X}\|_F^2 &\leq O(\epsilon^2)d^2(\mathbf{TP}, L_{\mathbf{X}}) + O(\epsilon^3)\|\mathbf{TPX}\|_F^2, \\ \left| \|\mathbf{TPX}\|_F^2 - \|\mathbf{T}\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \right| &\leq O(\epsilon)d^2(\mathbf{TP}, L_{\mathbf{X}}) + O(\epsilon)\|\mathbf{TPX}\|_F^2. \end{aligned}$$

Proof. The proof follows that of Lemma 4 to \mathbf{TP} . But now exact SVD is replaced with randomized SVD, so we need to argue that randomized SVD produces similar result as exact SVD in the sense of Lemma 7. This is already proved in Lemma 15. Also note that we need a technical lemma bounding the small error terms incurred on the intermediate result $\mathbf{T}\tilde{\mathbf{P}}$. This is done by Lemma 17. \square

Lemma 17.

$$\begin{aligned} \|\mathbf{T}\hat{\mathbf{P}}\mathbf{X}\|_F^2 &\leq \epsilon d^2(\mathbf{TP}, L_{\mathbf{X}}) + (1 + 2\epsilon)\|\mathbf{TPX}\|_F^2, \\ d^2(\mathbf{T}\hat{\mathbf{P}}, L_{\mathbf{X}}) &\leq (1 + \epsilon)d^2(\mathbf{TP}, L_{\mathbf{X}}) + \epsilon\|\mathbf{TPX}\|_F^2. \end{aligned}$$

Proof. For the first statement, by Lemma 15, we have

$$\begin{aligned} \left| \|\mathbf{T}\hat{\mathbf{P}}\mathbf{X}\|_F^2 - \|\mathbf{TPX}\|_F^2 \right| &\leq \sum_{i=1}^s \left| \|\mathbf{TP}_i\mathbf{X}\|_F^2 - \|\mathbf{T}\hat{\mathbf{P}}_i\mathbf{X}\|_F^2 \right| \\ &\leq \epsilon \sum_{i=1}^s d^2(\mathbf{TP}_i, L_{\mathbf{X}}) + 2\epsilon \sum_{i=1}^s \|\mathbf{TP}_i\mathbf{X}\|_F^2 \\ &\leq \epsilon d^2(\mathbf{TP}, L_{\mathbf{X}}) + 2\epsilon\|\mathbf{TPX}\|_F^2. \end{aligned} \quad (26)$$

For the second statement, by Pythagorean Theorem,

$$\begin{aligned} d^2(\mathbf{T}\hat{\mathbf{P}}, L_{\mathbf{X}}) - d^2(\mathbf{TP}, L_{\mathbf{X}}) &= \left[\|\mathbf{T}\hat{\mathbf{P}}\|_F^2 - \|\mathbf{T}\hat{\mathbf{P}}\mathbf{X}\|_F^2 \right] - \left[\|\mathbf{TP}\|_F^2 - \|\mathbf{TPX}\|_F^2 \right] \\ &= \left[\|\mathbf{T}\hat{\mathbf{P}}\|_F^2 - \|\mathbf{TP}\|_F^2 \right] + \left[\|\mathbf{TPX}\|_F^2 - \|\mathbf{T}\hat{\mathbf{P}}\mathbf{X}\|_F^2 \right] \\ &\leq \|\mathbf{TPX}\|_F^2 - \|\mathbf{T}\hat{\mathbf{P}}\mathbf{X}\|_F^2. \end{aligned}$$

The second statement then follows from the last inequality and (26). \square

Lemma 18. *For any $d \times k$ matrix \mathbf{X} with orthonormal columns,*

$$\begin{aligned} \|\mathbf{PX} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2 &\leq O(\epsilon^2)d^2(\mathbf{P}, L_{\mathbf{X}}) + O(\epsilon^3)\|\mathbf{PX}\|_F^2, \\ \left| \|\mathbf{PX}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \right| &\leq O(\epsilon)d^2(\mathbf{P}, L_{\mathbf{X}}) + O(\epsilon)\|\mathbf{PX}\|_F^2. \end{aligned}$$

Proof. By the property of subspace embedding, we have $\|\mathbf{TPX} - \mathbf{T}\tilde{\mathbf{P}}\mathbf{X}\|_F^2 = (1 \pm \epsilon)\|\mathbf{PX} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2$, $\|\mathbf{TPX}\|_F^2 = (1 \pm \epsilon)\|\mathbf{PX}\|_F^2$ and $d^2(\mathbf{TP}, L_{\mathbf{X}}) = \|\mathbf{TP} - \mathbf{TP}\mathbf{X}\mathbf{X}^\top\|_F^2 = (1 \pm \epsilon)\|\mathbf{P} - \mathbf{P}\mathbf{X}\mathbf{X}^\top\|_F^2 = (1 \pm \epsilon)d^2(\mathbf{P}, L_{\mathbf{X}})$. Then

$$\begin{aligned} (1 + \epsilon)\|\mathbf{PX} - \tilde{\mathbf{P}}\mathbf{X}\|_F^2 &\leq \|\mathbf{TPX} - \mathbf{T}\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \\ &\leq O(\epsilon^2)d^2(\mathbf{TP}, L_{\mathbf{X}}) + O(\epsilon^3)\|\mathbf{TPX}\|_F^2 \\ &\leq O(\epsilon^2)d^2(\mathbf{P}, L_{\mathbf{X}}) + O(\epsilon^3)\|\mathbf{PX}\|_F^2 \end{aligned}$$

where the second inequality is from Lemma 16. This then leads to the first statement.

For the second statement, we have

$$\begin{aligned}
(1 + \epsilon)\|\mathbf{P}\mathbf{X}\|_F^2 - (1 - \epsilon)\|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 &\leq \|\mathbf{TP}\mathbf{X}\|_F^2 - \|\mathbf{T}\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \\
&\leq O(\epsilon)d^2(\mathbf{TP}, L_{\mathbf{X}}) + O(\epsilon)\|\mathbf{TP}\mathbf{X}\|_F^2 \\
&\leq O(\epsilon)d^2(\mathbf{P}, L_{\mathbf{X}}) + O(\epsilon)\|\mathbf{P}\mathbf{X}\|_F^2
\end{aligned}$$

which leads to

$$\|\mathbf{P}\mathbf{X}\|_F^2 - \|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 \leq O(\epsilon)d^2(\mathbf{P}, L_{\mathbf{X}}) + O(\epsilon)\|\mathbf{P}\mathbf{X}\|_F^2.$$

A similar argument bounds $\|\tilde{\mathbf{P}}\mathbf{X}\|_F^2 - \|\mathbf{P}\mathbf{X}\|_F^2$, which completes the proof. \square

We represent Theorem 6 in a general form for ℓ_2 -error geometric fitting problems.

Theorem 6. *Suppose Algorithm 2 takes $\epsilon \in (0, 1/2]$, $t_1 = t_2 = O(\max\{\frac{k}{\epsilon^2}, \log \frac{s}{\delta}\})$, $\ell = O(\frac{d^2}{\epsilon^2})$, $q = O(\max\{\log \frac{d}{\epsilon}, \log \frac{sk}{\epsilon}\})$ as input, and sets the failure probability of each local subspace embedding to $\delta' = \delta/2s$. Let $\tilde{\mathbf{P}} = \mathbf{P}\mathbf{V}\mathbf{V}^\top$. Then with probability at least $1 - \delta$, there exists a constant $c_0 \geq 0$, such that for any set of k points \mathcal{L} ,*

$$(1 - \epsilon)d^2(\mathbf{P}, \mathcal{L}) - \epsilon\|\mathbf{P}\mathbf{X}\|_F^2 \leq d^2(\tilde{\mathbf{P}}, \mathcal{L}) + c_0 \leq (1 + \epsilon)d^2(\mathbf{P}, \mathcal{L}) + \epsilon\|\mathbf{P}\mathbf{X}\|_F^2$$

where \mathbf{X} is an orthonormal matrix whose columns span \mathcal{L} . The total communication is $O(sk d/\epsilon^2)$ and the total time is $O\left(\text{nnz}(\mathbf{P}) + s \left[\frac{d^3 k}{\epsilon^4} + \frac{k^2 d^2}{\epsilon^6}\right] \log \frac{d}{\epsilon} \log \frac{sk}{\delta \epsilon}\right)$.

Proof. The proof of correctness follows the proof of Theorem 3, replacing the use of Lemma 4 with Lemma 18.

On each node v_i , the subspace embedding takes time $O(\text{nnz}(\mathbf{P}_i))$, and the randomized SVD takes time $O(qt_1 \ell d + t_1^2(\ell + d))$; on the central coordinator, the randomized SVD takes time $O(qt_1(st_1)d + t_1^2(st_1 + d))$ since \mathbf{Y} has $O(st_1)$ non-zero rows. The total running time then follows from the choice of the parameters. The total communication cost follows from the fact that the algorithm only sends $\Sigma_i^{(t_1)}, \mathbf{V}_i^{(t_1)}$ from each node to the central coordinator. \square