# HETEROSKEDASTIC PCA: ALGORITHM, OPTIMALITY, AND APPLICATIONS

By Anru Zhang\*, T. Tony Cai<sup>†</sup>, and Yihong Wu<sup>‡</sup>

University of Wisconsin-Madison, University of Pennsylvania, and Yale University

Principal component analysis (PCA) and singular value decomposition (SVD) are widely used in statistics, econometrics, machine learning, and applied mathematics. It has been well studied in the case of homoskedastic noise, where the noise levels of the contamination are homogeneous.

In this paper, we consider PCA and SVD in the presence of heteroskedastic noise, which is a commonly used model for factor analysis and arises naturally in a range of applications. We introduce a general framework for heteroskedastic PCA and propose an algorithm called HeteroPCA, which involves iteratively imputing the diagonal entries to remove the bias due to heteroskedasticity. This procedure is computationally efficient and provably optimal under the generalized spiked covariance model. A key technical step is a deterministic robust perturbation analysis on singular subspaces, which can be of independent interest. The effectiveness of the proposed algorithm is demonstrated in a suite of applications, including heteroskedastic low-rank matrix denoising, Poisson PCA, and SVD based on heteroskedastic and incomplete data.

1. Introduction. Principal component analysis (PCA) and spectral methods are ubiquitous tools in many fields including statistics, econometrics, machine learning, and applied mathematics. They have been extensively studied and used in a wide range of applications. Recent examples include matrix denoising (Donoho and Gavish, 2014; Shabalin and Nobel, 2013), community detection (Donath and Hoffman, 2003; Newman, 2013), ranking from pairwise comparisons (Negahban et al., 2012; Chen and Suh,

 $<sup>^{*}\</sup>mathrm{The}$  research of Anru Zhang was supported in part by NSF Grant DMS-1811868 and NIH grant R01-GM131399-01.

 $<sup>^\</sup>dagger The$  research of Tony Cai was supported in part by NSF Grant DMS-1712735 and NIH grants R01-GM129781 and R01-GM123056.

<sup>&</sup>lt;sup>‡</sup>The research of Yihong Wu was supported in part by the NSF Grant CCF-1527105, an NSF CAREER award CCF-1651588, and an Alfred Sloan fellowship.

MSC 2010 subject classifications: Primary 62H12, 62H25; secondary 62C20

Keywords and phrases: factor analysis model, heteroskedasticity, low-rank matrix denoising, principal component analysis, singular value decomposition, spectral method

2015), matrix completion (Keshavan et al., 2010b; Sun and Luo, 2016), high-dimensional clustering (Jin et al., 2016), Markov process and reinforcement learning (Zhang and Wang, 2018), multidimensional scaling (Aflalo and Kimmel, 2013), topic modeling (Ke and Wang, 2017), phase retrieval (Candes et al., 2015; Cai et al., 2016), tensor PCA (Richard and Montanari, 2014; Zhang and Xia, 2018; Zhang and Han, 2018).

The central idea of PCA is to extract hidden low-rank structures from noisy observations. The following spiked covariance model has been well studied and used as a baseline for both methodological and theoretical developments for PCA (Johnstone, 2001; Baik et al., 2005; Baik and Silverstein, 2006; Paul, 2007; Donoho et al., 2018). Under this model, one observes  $Y_1, \ldots, Y_n \stackrel{iid}{\sim} N\left(\mu, \Sigma_0 + \sigma^2 I_p\right)$ , where  $\Sigma_0 = U\Lambda U^{\top}$  is a symmetric low-rank matrix and  $I_p$  is a p-dimensional identity matrix. The spiked covariance model can be equivalently written as

$$(1.1) \quad Y_k = X_k + \varepsilon_k, \quad X_k \stackrel{iid}{\sim} N(\mu, \Sigma_0), \quad \varepsilon_k \stackrel{iid}{\sim} N(0, \sigma^2 I_p), \quad k = 1, \dots, n.$$

The goal is either to recover  $\Sigma_0$  or factor loadings U. Let  $\hat{\Sigma}$  be the sample covariance matrix of  $Y_1, \ldots, Y_n$ . In literature, asymptotic properties of eigenvalues and eigenvectors of  $\hat{\Sigma}$  have been well established and estimators based on the eigen-decomposition of  $\hat{\Sigma}$  have been introduced and extensively studied. A key assumption here is that errors are homoskedastic, in the sense that each  $\varepsilon_k$  is assumed to be spherically symmetric Gaussian.

1.1. Heteroskedastic PCA. In many applications, the noise can be highly heteroskedastic, i.e., the magnitude of perturbation varies significantly from entry to entry in the data matrix. For example, heteroskedastic noise naturally appears in datasets with different types of variables. For various biological sequencing and photon imaging data, the observations are discrete counts that are commonly modeled by Poisson, multinomial, or negative binomial distributions (Salmon et al., 2014; Cao et al., 2017) and are naturally heteroskedastic. In network analysis and recommender systems, the observations are usually binary or ordinal, which are heteroskedastic. PCA is also used in the analysis of spectrophotometric data to determine the number of linearly independent species in rapid scanning wavelength kinetics experiments (Cochran and Horne, 1977). The spectrophotometric data often contain heteroskedastic noise since the measurements are based on averages over varying lengths of time intervals.

Motivated by these applications, it is natural to relax the homoskedasticity assumption in (1.1) and consider the following generalized spiked covari-

ance model (Bai and Yao, 2012; Yao et al., 2015):

$$(1.2) Y = X + \varepsilon, \mathbb{E}X = \mu, \operatorname{Cov}(X) = \Sigma_0,$$

$$\mathbb{E}\varepsilon = 0, \operatorname{Cov}(\varepsilon) = \operatorname{diag}(\sigma_1^2, \dots, \sigma_p^2),$$

$$\varepsilon = ((\varepsilon)_1, \dots, (\varepsilon)_p)^\top; X, (\varepsilon)_1, \dots, (\varepsilon)_p \text{ are independent.}$$

Here,  $\operatorname{Cov}(X)$  is rank-r and admits the following eigenvalue decomposition:  $\operatorname{Cov}(X) = \Sigma_0 = U\Lambda U^{\top}$  with  $U \in \mathbb{R}^{p \times r}$  and  $\Lambda \in \mathbb{R}^{r \times r}$ ;  $\sigma_1^2, \ldots, \sigma_p^2$  are unknown and need not be identical. This model is also widely used as the standard model in the *factor analysis* literature (see, e.g., Tipping and Bishop (1999); Ghosh and Dunson (2009) and the references therein). Given i.i.d. samples  $Y_1, \ldots, Y_n$  drawn from (1.2), the goal is to estimate the factor loadings U.

Due to the heteroskedasticity of the noise variances, it turns out that the classical PCA can lead to inconsistent estimates. To see this, note that performing PCA on  $\{Y_1, \ldots, Y_n\}$  amounts to applying the regular SVD on  $\mathbf{Y} = [Y_1, \ldots, Y_n]$ , i.e., estimating U by leading left singular vectors of the centralized sample matrix:

(1.3) 
$$\mathbf{Y} - \bar{Y} \mathbf{1}_n^{\mathsf{T}}, \quad \bar{Y} = \frac{1}{n} \sum_{k=1}^n Y_k.$$

Moreover, the left singular vectors of  $\mathbf{Y} - \bar{Y} \mathbf{1}_n^{\top}$  are identical to eigenvectors of the sample covariance matrix,

$$(1.4) \quad \hat{\Sigma} = \frac{1}{n-1} (\mathbf{Y} - \bar{Y} \mathbf{1}_n^{\top}) (\mathbf{Y} - \bar{Y} \mathbf{1}_n^{\top})^{\top} = \frac{1}{n-1} \sum_{k=1}^{n} (Y_k - \bar{Y}) (Y_k - \bar{Y})^{\top}.$$

Note that  $\mathbb{E}\hat{\Sigma} = \Sigma_0 + \operatorname{diag}(\sigma_1^2, \dots, \sigma_p^2)$ . When  $\sigma_1^2, \dots, \sigma_p^2$  are the same, the top eigenvectors of  $\mathbb{E}\hat{\Sigma}$  and  $\Sigma_0$  coincide; however, when  $\sigma_1^2, \dots, \sigma_p^2$  are not identical, the principal components of  $\mathbb{E}\hat{\Sigma}$  and those of  $\Sigma_0$  can differ significantly due to the bias of  $\hat{\Sigma}$  on diagonal entries. This shows the inadequacy of regular SVD in the case of heteroskedastic noise. In addition to PCA in the generalized spiked covariance model, this phenomenon similarly appears in other problems with heteroskedastic noise, such as heteroskedastic low-rank matrix denoising, Poisson PCA, SVD from incomplete and heteroskedastic values. See later Section 3 for details.

To better cope with the bias incurred on the diagonal, Florescu and Perkins (2016) introduced a method called the *diagonal-deletion SVD* in the context of bipartite stochastic block model. The idea is to set the diagonal of

Gram matrix to zero, then perform singular value decomposition. However, it is a priori unclear whether zeroing out the diagonals is always the best choice. In fact, we can construct explicit examples where diagonal-deletion SVD is inconsistent (c.f., forthcoming Proposition 1).

In the paper, we introduce a novel method, called *HeteroPCA*, for heteroskedastic principal component analysis. Instead of zeroing out diagonal entries of the sample covariance/Gram matrix, we propose to iteratively update diagonal entries based on off-diagonals, so that the bias incurred on the diagonal is significantly reduced and more accurate estimation can be achieved.

The performance of the proposed procedure is studied both theoretically and numerically. By proving matching minimax upper and lower bounds, we show that HeteroPCA achieves the optimal rate of convergence among a general class of settings under the generalized spiked covariance model. In particular, the procedure provably outperforms regular SVD and diagonal-deletion SVD.

The subspace perturbation bound plays a key role for theoretical analysis of various PCA methods. Classic tools, such as Davis-Kahan and Wedin's theorems (Davis and Kahan, 1970; Wedin, 1972), bound the subspace estimation error in terms of overall perturbation and spectral gap on the sample covariance matrix. Due to the aforementioned bias on the diagonal entries of the sample covariance matrix, these classic tools may not be suitable for analyzing heteroskedastic PCA here. To tackle this difficulty, we develop a new deterministic subspace perturbation bound (see Theorem 3 next), which provides the key technical tool for analyzing HeteroPCA procedure and may be of independent interest.

1.2. Applications and Related Literature. In addition to heteroskedastic PCA in the generalized spiked covariance model, the newly established HeteroPCA algorithm is applicable to a collection of high-dimensional statistical problems with heteroskedastic data. We also discuss in detail the applications of heteroskedastic low-rank matrix denoising, Poisson PCA, and SVD based on heteroskedastic and incomplete data in Section 3. Moreover, our result is also useful to a range of other applications, such as heteroskedastic canonical correlation analysis, heteroskedastic tensor SVD, exponential family PCA, community detection in bipartite stochastic network, and bipartite multidimensional scaling.

This paper is related to several recent works on PCA for heteroskedastic data. For example, Bai and Yao (2012); Yao et al. (2015) extended the

theory of regular spiked covariance model to generalized one and studied the eigenvalue limiting distribution of sample covariance matrix. Hong et al. (2016, 2018a,b) considered PCA with heteroskedastic noise in an alternative way. They introduced a model for heteroskedastic data, where the noise is non-uniform across different samples but uniform within each sample. They further studied the performance of regular SVD and established asymptotic distributions for both eigenvalue and eigenvector estimators. Our work is also closely related to a substantial body of literature on factor model analysis (Thomson, 1939; Lawley and Maxwell, 1962; Tipping and Bishop, 1999; Ghosh and Dunson, 2009; Bai and Li, 2012; Owen and Wang, 2016; Wang and Fan, 2017). Various approaches have been developed for estimating principal components in factor models, such as regression method (Thomson, 1939), weighted least squares (Bartlett, 1937) EM (Tipping and Bishop, 1999), and Bayesian MCMC (Ghosh and Dunson, 2009). The asymptotic theory was also extensively studied (Bai and Li, 2012; Wang and Fan, 2017). Departing from the previous results, this paper mainly concerns a non-asymptotic framework, providing computational algorithm with provable guarantees and allowing heteroskedastic noise within each sample; in fact, the noise variances can be different among all entries of data (see the heteroskedastic low-rank matrix denoising in Section 3.1). Since the regular SVD no longer achieves good performance when noise is heteroskedastic (c.f., the forthcoming Proposition 1), we instead propose and analyze the new procedure HeteroPCA. To the best of our knowledge, this is the first paper that provides a frequentist approach for subspace estimation with minimax optimal theoretical guarantees.

1.3. Organization of the Paper. The rest of the paper is organized as follows. After a brief introduction of notation and definitions (Section 2.1), we focus on the generalized spiked covariance model, present the HeteroPCA algorithm (Section 2.2), and develop matching minimax upper and lower bounds for the proposed procedure (Section 2.3). Then, we introduce a deterministic robust perturbation analysis that serves as a key technical step in our analysis (Section 2.4). We also illustrate main proof ideas of technical results in Section 2.5. Applications to other high-dimensional statistical problems, including heteroskedastic matrix denoising, Poisson PCA, and SVD based on heteroskedastic and incomplete data are discussed in Section 3. Numerical results are presented in Section 4 and other applications are briefly discussed in Section 5. The proofs of main results are given in Section 6. The additional proofs and technical lemmas are provided in the supplementary materials (Zhang et al., 2018).

#### 2. Optimal Heteroskedastic Principal Component Analysis.

2.1. Notation and Preliminaries. We use lowercase letters, e.g., x, y, z, to denote scalars or vectors; we use uppercase letters, e.g. U, M, N to denote matrices. For any sequences of positive numbers  $\{a_k\}$  and  $\{b_k\}$ , denote  $a \lesssim b$ and  $b \gtrsim a$  if there exists a uniform constant C > 0 such that  $a_k \leq Cb_k$  for all k. We also say  $a \approx b$  if  $a \lesssim b$  and  $a \gtrsim b$  both hold. For any matrix  $M \in \mathbb{R}^{p_1 \times p_2}$ , let  $\lambda_k(M)$  be the k-th largest singular value. Then, the SVD of M can be written as  $M = \sum_{k=1}^{p_1 \wedge p_2} \lambda_k(M) u_k v_k^{\top}$ . We also let  $SVD_r(M) =$  $[u_1 \cdots u_r]$  be the collection of leading r left singular vectors and QR(M) be the Q part of QR orthogonalization of M. The matrix spectral norm ||M|| = $\sup_{\|u\|_2=1} \|Mu\|_2 = \lambda_1(M)$  and Frobenius norm  $\|M\|_F = (\sum_{i,j} M_{ij}^2)^{1/2} =$  $(\sum_{k} \lambda_{k}^{2}(M))^{1/2}$  will be extensively used throughout the paper. Let  $I_{r}, 0_{m \times n}$ , and  $1_{m \times n}$  be the r-by-r identity,  $m \times n$  zero, and  $m \times n$  all-one matrices, respectively. Also let  $0_m$  and  $1_m$  denote the m-dimensional zero and allone column vectors. Denote  $\mathbb{O}_{p,r} = \{U \in \mathbb{R}^{p \times r} : U^{\top}U = I_r\}$  as the set of all p-by-r matrices with orthonormal columns. For any  $U \in \mathbb{O}_{p,r}$ , we note  $U_{\perp} \in \mathbb{O}_{p,p-r}$  as the orthogonal complement so that  $[U \ U_{\perp}] \in \mathbb{R}^{p \times p}$  is a complete orthogonal matrix.

Motivated by incoherence condition, a widely used assumption in the matrix completion literature (Candès and Recht, 2009), we define *incoherence* constant of  $U \in \mathbb{O}_{p,r}$  as

(2.1) 
$$I(U) = -\frac{p}{r} \max_{i \in [p]} \|e_i^\top U\|_2^2.$$

sin  $\Theta$  distance is used to quantify the distance between singular subspaces. Specifically for any  $U_1, U_2 \in \mathbb{O}_{p,r}$ , define  $\|\sin\Theta(U_1, U_2)\| \triangleq \|U_{1\perp}^\top U_2\| = \|U_{2\perp}^\top U_1\|$ . For any square matrix A, let  $\Delta(A)$  be A with all diagonal entries set to zero and D(A) be A with all off-diagonal entries set to zero. Then  $A = \Delta(A) + D(A)$ . We use  $C, C_1, \ldots, c, c_1, \cdots$  to respectively represent generic large and small constants, whose values may differ in different lines.

2.2. Methods for Heteroskedastic PCA. Now we are in position to investigate the heteroskedastic principal component analysis in detail. Suppose one observes i.i.d. copies  $Y_1, \ldots, Y_n$  of Y from the generalized spiked covariance model (1.2). Let  $\hat{\Sigma}$  be the sample covariance matrix defined as (1.3) (1.4). In order to estimate U, i.e., the leading principal components of  $\Sigma_0$ , the most natural estimator is  $\tilde{U} = \text{SVD}_r(\hat{\Sigma})$ , i.e., the subspace composed of first r left singular vectors of  $\hat{\Sigma}$ . This idea is widely referred to as singular

value thresholding (SVT) in literature (Donoho and Gavish, 2014; Chatterjee, 2015). The singular value shrinkage is another closely related method that has been proposed and studied previously (Nadakuditi, 2014; Gavish and Donoho, 2017; Donoho et al., 2018). By the well-known Eckart-Young-Mirsky Theorem (Golub et al., 1987), SVT, or the regular SVD estimator, is equivalent to the following optimization problem,

(2.2) 
$$\tilde{U} = \text{SVD}_r(\tilde{\Sigma}), \text{ where } \tilde{\Sigma} = \underset{\tilde{\Sigma}: \text{rank}(\tilde{\Sigma}) < r}{\arg \min} \left\| \tilde{\Sigma} - \hat{\Sigma} \right\|.$$

In particular, an important variant of Davis-Kahan's theorem (Davis and Kahan, 1970) given by Yu, Wang, and Samworth (Yu et al., 2014) yields

(2.3) 
$$\left\| \sin \Theta(\tilde{U}, U) \right\| \lesssim \frac{\|\hat{\Sigma} - (\Sigma_0 + \beta I_p)\|}{\lambda_r(\Lambda)} \wedge 1$$

for any scalar  $\beta \geq 0$ . Such a bound is sharp in the worst case. However, as discussed earlier,  $\mathbb{E}\hat{\Sigma} = \Sigma_0 + \operatorname{diag}(\sigma_1^2, \dots, \sigma_p^2)$  and the perturbation of  $\hat{\Sigma} - (\Sigma_0 + \beta I_p)$  need not be homogeneous for any scalar  $\beta \geq 0$  if  $\sigma_1^2, \dots, \sigma_p^2$  have different values. In other words, the diagonal entries of perturbation matrix  $(\hat{\Sigma} - \Sigma_0 + \beta I_p)$  may be significantly larger than the rest.

To achieve more robust estimates of U with provable guarantees, we propose the following simple and computationally feasible procedure. To be specific, if perturbation Z has higher amplitude on the diagonal, we ignore those diagonal entries of Z in (2.2) and consider

(2.4) 
$$\hat{U} = \text{SVD}_r(\hat{M}), \quad \text{where} \quad \hat{M} = \underset{\hat{M}: \text{rank}(\hat{M}) \leq r}{\arg \min} \left\| \Delta(\hat{M} - \hat{\Sigma}) \right\|.$$

Since (2.4) is non-convex, we instead consider the following procedure.

Step 1 Initialize by setting the diagonal of  $\hat{\Sigma}$  to zero:  $N^{(0)} = \Delta(\hat{\Sigma})$ .

Step 2 For t = 0, ..., perform SVD on  $N^{(t)}$  and let  $\tilde{N}^{(t)}$  be its best rank-rapproximation:

$$N^{(t)} = U^{(t)} \Sigma^{(t)} (V^{(t)})^{\top} = \sum_{i} \lambda_{i}^{(t)} u_{i}^{(t)} (v_{i}^{(t)})^{\top}, \quad \lambda_{1}^{(t)} \ge \dots \ge \lambda_{m}^{(t)} \ge 0,$$

$$\tilde{N}^{(t)} = \sum_{i=1}^{r} \lambda_i^{(t)} u_i^{(t)} (v_i^{(t)})^{\top}.$$

Step 3 Update  $N^{(t+1)} = D(\tilde{N}^{(t)}) + \Delta(N^{(t)})$ , i.e., replace the diagonal entries of  $N^{(t)}$  by those in  $\tilde{N}^{(t)}$ . In other words,

(2.5) 
$$N_{ij}^{(t+1)} = \begin{cases} N_{ij}^{(t)} = \hat{\Sigma}_{ij}, & i = j; \\ \tilde{N}_{ij}^{(t)}, & i \neq j. \end{cases}$$

Step 4 Repeat Steps 2-3 until convergence or the maximum number of iterations is reached.

The pseudo-code of proposed procedure is summarized as Algorithm 1.

## Algorithm 1 HeteroPCA

```
1: Input: matrix \hat{\Sigma}, rank r, number of iterations T.
```

- 2: Set  $N^{(0)} = \Delta(\hat{\Sigma})$ .
- 3: **for** t = 1, ..., T **do**
- Calculate SVD:  $N^{(t)} = \sum_{i} \lambda_{i}^{(t)} u_{i}^{(t)} (v_{i}^{(t)})^{\top}$ , where  $\lambda_{1}^{(t)} \geq \lambda_{2}^{(t)} \cdots \geq 0$ . Let  $\tilde{N}^{(t)} = \sum_{i=1}^{r} \lambda_{i}^{(t)} u_{i}^{(t)} (v_{i}^{(t)})^{\top}$ . Update diagonal entries:  $N^{(t+1)} = D(\tilde{N}^{(t)}) + \Delta(N^{(t)})$ .

- 8: Output:  $\hat{U} = U^{(T)} = [u_1^{(T)} \cdots u_r^{(T)}].$

### 2.3. Theoretical Analysis. Denote

$$\sigma_{\max}^2 \triangleq \max_i \sigma_i^2, \quad \sigma_{\sup}^2 \triangleq \sum_i \sigma_i^2.$$

Recall  $\Sigma_0 = U\Lambda U^{\top}$ ,  $\lambda_r(\Lambda)$  is the r-th singular value of  $\Lambda$ , which is also the r-th largest eigenvalue of  $\Sigma_0$ . We have the following theoretical guarantee for Algorithm 1.

Theorem 1 (Heteroskedastic PCA: upper bound). Consider the generalized spiked covariance model (1.2), where X and  $\varepsilon$  are sub-Gaussian in the sense that

$$\max_{q \ge 1, ||v||_2 = 1} q^{-1/2} \left( \mathbb{E} |v^{\top} \Lambda^{-1/2} U^{\top} X|^q \right)^{1/q} \le C,$$
$$\max_{q \ge 1, ||v||_2 = 1} q^{-1/2} \left( \mathbb{E} |\varepsilon_i / \sigma_i|^q \right)^{1/q} \le C.$$

Let  $Y_1, \ldots, Y_n$  be i.i.d. samples from (1.2). Assume that  $n \geq Cr$ ,  $\sigma_{sum}^2/\lambda_r(\Lambda) \geq$  $\exp(-Cn)$ , and  $\|\Lambda\|/\lambda_r(\Lambda) \leq C$  for constant C>0. Then, there exists constant  $c_I > 0$  such that if the incoherence constant  $I(U) = \max_i \frac{p}{r} ||e_i^{\top} U||_2^2$ satisfies  $I(U) \leq c_I p/r$ , then the output  $\hat{U}$  of Algorithm 1 applied to the sample covariance matrix  $\hat{\Sigma}$  satisfies:

$$(2.6) \mathbb{E}\left\|\sin\Theta(\hat{U},U)\right\| \lesssim \frac{1}{\sqrt{n}} \left(\frac{\sigma_{sum} + r^{1/2}\sigma_{max}}{\lambda_r^{1/2}(\Lambda)} + \frac{\sigma_{sum}\sigma_{max}}{\lambda_r(\Lambda)}\right) \wedge 1,$$

where  $\|\sin\Theta(\cdot,\cdot)\|$  is the  $\sin\Theta$  distance between two subspaces.

REMARK 1 (Interpretation of (2.6)). Let  $\tilde{p} = \sigma_{\text{sum}}^2/\sigma_{\text{max}}^2$ . The upper bound (2.6) can be rewritten as

$$(2.7) \mathbb{E}\left\|\sin\Theta(\hat{U},U)\right\| \lesssim \left(\sqrt{\frac{\tilde{p}\vee r}{n}}\frac{\sigma_{\max}}{\lambda_r^{1/2}(\Lambda)} + \sqrt{\frac{\tilde{p}}{n}}\frac{\sigma_{\max}^2}{\lambda_r(\Lambda)}\right) \wedge 1.$$

Consider the homoskedastic PCA setting where  $\sigma_1^2 = \cdots = \sigma_p^2 = \sigma_{\text{max}}^2$ . A special case of Theorem 1 yields:

(2.8) 
$$\mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\| \lesssim \sqrt{\frac{p}{n}} \left( \frac{\sigma_{\max}}{\lambda_r^{1/2}(\Lambda)} + \frac{\sigma_{\max}^2}{\lambda_r(\Lambda)} \right) \wedge 1.$$

Comparing (2.7) with (2.8), we see that a weighted average between  $\tilde{p} \vee r$  and  $\tilde{p}$  can be viewed as the "effective dimension" for heteroskedastic PCA.

Next, to establish the optimality of Theorem 1, we consider the following class of generalized spiked covariance matrices:

(2.9) 
$$\mathcal{F}_{p,n,r}(\sigma_{\text{sum}}, \sigma_{\text{max}}, \nu) = \left\{ \Sigma = U \Lambda U^{\top} + D : \right. \\ \left. D \text{ is non-negative diagonal, } \sum_{i} D_{ii} \leq \sigma_{\text{sum}}^{2}, \max_{i} D_{ii} \leq \sigma_{\text{max}}^{2}, \\ \left. U \in \mathbb{O}_{p,r}, I(U) \leq c_{I} p/r, \|\Lambda\|/\lambda_{r}(\Lambda) \leq C, \lambda_{r}(\Lambda) \geq \nu \right\}.$$

We establish the following minimax lower bound of heteroskedastic PCA for covariance matrices in  $\mathcal{F}_{p,n,r}(\sigma_{\text{sum}}, \sigma_{\text{max}}, \nu)$ .

THEOREM 2 (Heteroskedastic PCA: lower bound). Suppose  $\sqrt{p}\sigma_{max} \geq \sigma_{sum} \geq \sigma_{max} > 0$ . There exists constant C > 0, such that if  $p \geq Cr$ , the following lower bound holds, (2.10)

$$\inf_{\hat{U}} \sup_{\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{sum},\sigma_{max},\nu)} \mathbb{E} \left\| \sin \Theta(\hat{U},U) \right\| \gtrsim \frac{1}{\sqrt{n}} \left( \frac{\sigma_{sum} + r^{1/2}\sigma_{max}}{\nu^{1/2}} + \frac{\sigma_{sum}\sigma_{max}}{\nu} \right) \wedge 1.$$

REMARK 2. By combining Theorems 1 and 2, the proposed Algorithm 1 achieves the following optimal rate of convergence:

$$\inf_{\hat{U}} \sup_{\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{\text{sum}},\sigma_{\text{max}},\nu)} \mathbb{E} \left\| \sin \Theta(\hat{U},U) \right\| \approx \frac{1}{\sqrt{n}} \left( \frac{\sigma_{\text{sum}} + r^{1/2}\sigma_{\text{max}}}{\nu^{1/2}} + \frac{\sigma_{\text{sum}}\sigma_{\text{max}}}{\nu} \right) \wedge 1.$$

Remark 3 (Regular and diagonal-deletion SVDs in heteroskedastic PCA). Noting that  $\mathbb{E}\hat{\Sigma} = U\Lambda U^{\top} + \mathrm{diag}(\sigma_1^2,\ldots,\sigma_p^2)$  and  $\mathbb{E}\Delta(\hat{\Sigma}) = \Delta(U\Lambda U^{\top})$ , both  $\mathbb{E}\hat{\Sigma}$  and  $\mathbb{E}\Delta(\hat{\Sigma})$  may have different singular subspaces than U when heteroskedastic noise exists. Then, it may be less appropriate to estimate U based on  $\hat{\Sigma}$  itself or  $\hat{\Sigma}$  with diagonal entries simply replaced by zero. In fact, we can show that the regular SVD or diagonal-deletion SVD

$$\hat{U}^{\text{SVD}} = \text{SVD}_r(\hat{\Sigma}), \quad \hat{U}^{\text{DD}} = \text{SVD}_r(\Delta(\hat{\Sigma}))$$

may be inconsistent, even in the "fixed p, growing n" scenario, by the following lower bound argument.

Proposition 1. There exists a constant C > 0 such that if  $p \ge Cr$ , we have

$$(2.11) \sup_{\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{sum},\sigma_{max},\nu)} \mathbb{E} \| \sin \Theta(\hat{U}^{\text{SVD}}, U) \| \gtrsim \left( \frac{\sigma_{sum}/n^{1/2} + \sigma_{max}}{\nu^{1/2}} \right) \wedge 1,$$

$$(2.12) \sup_{\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{sum},\sigma_{max},\nu)} \mathbb{E} \| \sin \Theta(\hat{U}^{\text{DD}}, U) \| \gtrsim 1.$$

When the covariance matrix  $\Sigma_0$  is approximately rank-r and there exists a significant gap between  $\lambda_r(\Sigma_0)$  and  $\lambda_{r+1}(\Sigma_0)$ , the proposed HeteroPCA algorithm still achieves stable performance with provable guarantees.

PROPOSITION 2 (HeteroPCA for approximately low-rank covariance). Consider the generalized spiked covariance model (2.13). Suppose  $\Sigma_0 = \tilde{U}\Lambda\tilde{U}^{\top}$  is the eigenvalue decomposition, where  $\tilde{U} = [U\ U_{\perp}]$  and U is the collection of leading r singular vectors. Assume X and  $\varepsilon$  are sub-Gaussian in the sense that there exists a constant  $C_0 > 0$  such that

$$\max_{q \ge 1, \|v\|_2 = 1} q^{-1/2} \left( \mathbb{E} |v^{\top} \Lambda^{-1/2} \tilde{U}^{\top} X|^q \right)^{1/q} \le C_0,$$
$$\max_{q \ge 1, \|v\|_2 = 1} q^{-1/2} \left( \mathbb{E} |\varepsilon_i / \sigma_i|^q \right)^{1/q} \le C_0.$$

Also assume that  $n \geq Cr$ ,  $\sigma_{sum}^2/\lambda_r(\Lambda) \geq \exp(-Cn) + \exp(-Cp)$ , and  $\|\Lambda\|/\lambda_r(\Lambda) \leq C$  for some constant C>0. Then there exists some constant  $c_I>0$  such that if the incoherence constant  $I(U)=\max_i \frac{p}{r}\|e_i^\top U\|_2^2$  satisfies  $I(U) \leq c_I p/r$ , then the output  $\hat{U}$  of Algorithm 1 with input rank r

and matrix  $\hat{\Sigma}$  satisfies

$$\mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\|$$

$$\lesssim \left( \frac{\sigma_{sum} + \sqrt{r} \sigma_{max}}{n^{1/2} \lambda_r^{1/2}(\Lambda)} + \frac{\sigma_{sum} \sigma_{max}}{n^{1/2} \lambda_r(\Lambda)} + \frac{((np)^{1/2} + p) \lambda_{r+1}^{1/2}(\Lambda)}{n \lambda_r^{1/2}(\Lambda)} + \frac{\lambda_{r+1}(\Lambda)}{\lambda_r(\Lambda)} \right) \wedge 1.$$

2.4. A Deterministic Robust Perturbation Analysis. Next, we provide a low-rank subspace perturbation analysis in a deterministic framework. The theory developed here will play a key role for the heteroskedastic PCA problem in Section 1.1. Let N, M, Z be deterministic symmetric matrices that satisfy

$$(2.13) N = M + Z.$$

Here N is the observation,  $M \in \mathbb{R}^{p \times p}$  is the rank-r matrix of interest, and  $Z \in \mathbb{R}^{p \times p}$  is the perturbation. For the heteroskedastic PCA problem, N, M, Z correspond to the sample covariance matrix  $\hat{\Sigma}$ , the samples covariance matrix of the signal vectors  $\hat{\Sigma}_X$ , and their difference  $\hat{\Sigma} - \hat{\Sigma}_X$ , respectively. Let  $U \in \mathbb{O}_{p,r}$  consist of r singular vectors of M. As discussed earlier, applying the proposed HeteroPCA (Algorithm 1) to the matrix N provides an adaptive estimate of U that is unaffected by the significant corruption on the diagonal of perturbation matrix Z. Next, we analyze the theoretical property for the proposed Algorithm 1 under the general robust perturbation model (2.13).

THEOREM 3 (Robust  $\sin \Theta$  theorem). Suppose  $M \in \mathbb{R}^{p \times p}$  is a rank-r symmetric matrix and  $U \in \mathbb{O}_{p,r}$  consists of the eigenvectors of M. Then there exists a universal constant  $c_I > 0$  such that if

(2.14) 
$$\frac{I(U)\|M\|}{\lambda_r(M)} \le \frac{c_I p}{r},$$

(where  $I(U) = \max_i \frac{p}{r} ||e_i^\top U||_2^2$  is the incoherence constant defined in (2.1)), then the output  $\hat{U}$  of Algorithm 1 with input matrix  $\hat{\Sigma} = N$ , rank r, number of iterations  $T = \Omega\left(\log \frac{\lambda_r(M)}{||Z||} \vee 1\right)$  satisfies

(2.15) 
$$\left\| \sin \Theta(\hat{U}, U) \right\| \lesssim \frac{\|\Delta(Z)\|}{\lambda_r(M)} \wedge 1$$

where  $\Delta(Z)$  is the same as Z except that all diagonal entries are set to zero.

Remark 4. We introduce the incoherence condition (2.14) here to avoid those M that are too "spiky". For example, consider  $M_1 = e_1 e_1^{\top}$  and  $M_2 = e_2 e_2^{\top}$ . Then  $\Delta(M_1) = \Delta(M_2)$  and there is no way to distinguish these two spiky matrices if one only has reliable off-diagonal observations. Similar conditions, such as the "delocalized condition," appear in recent work on PCA from noisy and linearly reduced data (Dobriban et al., 2016). The incoherence condition has been widely used in the matrix completion literature (e.g., (Recht, 2011, Assumption A0)), where  $I(U) \leq \mu_0$  is often assumed for some constant  $\mu_0$  independent of p. In comparison, in view of the trivial bound

 $I(U) = \frac{p}{r} \max_{1 \le i \le p} \|e_i^\top U\|_2^2 \le \frac{p}{r} \cdot 1,$ 

our assumption  $I(U) < c_I p/r$  is much looser than those that are prevalent in the matrix completion literature.

Compared to previous work on robust PCA (e.g., Candès et al. (2011)) that typically assumes corruptions have arbitrary amplitude but randomly selected sparse support, our robust perturbation analysis here is fully deterministic. In addition to Theorem 3, we also consider the subspace perturbation analysis where significant entries appear in a known subset  $\mathcal{G}$  rather than the diagonal of perturbation Z and provide a more general performance guarantee. See Section 6.2 for details.

REMARK 5. In addition to Theorem 3, we actually consider a more general situation where the corrupted entries lie in a more general set that need not be diagonal. We also prove the performance guarantee for a counterpart of Algorithm 1. See Section 6.2 for details.

The following lower bound shows that bounds for both the incoherence condition (2.14) and the estimation error (2.15) are rate-optimal.

PROPOSITION 3 (Robust  $\sin \Theta$  theorem: lower bound). Define the following collection of pairs of signal and perturbation matrices:

(2.16) 
$$\mathcal{D}_{p,r}(\nu,\delta,t) = \left\{ \begin{array}{c} M = U\Lambda U^{\top}, U \in \mathbb{O}_{p,r}, \\ (M,Z) : I(U) \|M\| / \lambda_r(M) \le t, \\ \|\Delta(Z)\| \le \delta, \lambda_r(M) \ge \nu \end{array} \right\}.$$

Suppose  $1 \le r \le p/2, t \ge 4$ , one observes  $N = M + Z \in \mathbb{R}^{p \times p}$ . Then

(2.17) 
$$\inf_{\hat{U}} \sup_{(M,Z) \in \mathcal{D}_{n,r}(\nu,\delta,t)} \left\| \sin \Theta(\hat{U},U) \right\| \ge c \left( \frac{\delta}{\nu} \wedge 1 \right).$$

If the incoherence constraint, i.e.,  $I(U)||M||/\lambda_r(M) \le t$ , is weak in the sense that  $t \ge p/r$ , then

(2.18) 
$$\inf_{\hat{U}} \sup_{(M,Z) \in \mathcal{D}_{p,r}(\nu,\delta,t)} \left\| \sin \Theta(\hat{U},U) \right\| \ge 1/2.$$

2.5. Proof Sketches of Main Technical Results. We briefly discuss the proofs of Theorems 1, 2, and 3 in this section. The detailed proofs are deferred to Section 6 and the supplementary materials.

The proof of Theorem 1 consists of three main steps. First, we define  $\hat{\Sigma}_X$  as the sample covariance matrix of signal vectors  $X_1, \ldots, X_n$  and  $E = [\varepsilon_1 \cdots \varepsilon_n]$  as the noise matrix. We aim to develop a concentration inequality for  $\Delta\left((n-1)(\hat{\Sigma}-\hat{\Sigma}_X)\right)$ , i.e., the off-diagonal part of perturbation. To this end, we decompose  $(n-1)(\hat{\Sigma}-\hat{\Sigma}_X)$  into  $(XE^\top+EX^\top), (EE^\top), n(\bar{X}\bar{E}^\top+\bar{E}\bar{X}^\top+\bar{E}\bar{E}^\top)$ , then bound them separately by heteroskedastic Wishart concentration inequality (Cai and Zhang, 2018a) and Lemma 3 in the supplementary materials. Second, we develop a lower bound for  $\lambda_r(\hat{\Sigma}_X)$ , i.e., the least non-trivial singular value of the signal covariance matrix. Finally, we apply the robust  $\sin\Theta$  theorem (Theorem 3), to complete the proof.

To show the lower bound in Theorem 2, it suffices to show the two terms in (2.10) separately; c.f., (A.1) and (A.2) in the detailed proof. To show each individual lower bound, we construct a series of "candidate matrices"  $\{U^{(k)}, \Sigma^{(k)}\}_{k=1}^{N}$  in  $\mathcal{F}_{p,n,r}(\sigma_{\text{sum}}, \sigma_{\text{max}}, \nu)$  so that  $\{U^{(k)}\}_{k=1}^{N}$  are well-separated while distinguishing them apart based on random sample  $Y_1, \ldots, Y_n \sim N(0, \Sigma^{(k)})$  is impossible. This implies the desired lower bound by applying Fano's method.

The proof of Theorem 3 is the main technical contribution of this paper. Specifically, we analyze how the estimation error  $K_t = ||N^{(t)} - M||$  at each iteration decays. First, we obtain an initialization error bound. Then for each t, we decompose  $K_t$  into four terms, bound them separately, and obtain an inequality that relates  $K_t$  to  $K_{t-1}$  (see (6.19)). By induction, this recursive inequality leads to the exponential decay of  $K_t$  and implies the desired upper bound. Note that Algorithm 1 can be viewed as successive compositions involving the projection operator  $P_U(\cdot)$  and the diagonal-deletion operator  $D(\cdot)$ . We thus introduce Lemma 1 to give sharp operator norm upper bounds for compositions of  $P_U(\cdot)$  and  $D(\cdot)$ . At the heart of the proof of Theorem 3, this lemma is useful for bounding the error at both the initialization and the subsequent iterations.

- 3. More Applications in High-dimensional Statistics. In this section, we consider a number of additional applications in high-dimensional statistics, including the heteroskedastic low-rank matrix denoising, Poisson PCA, and SVD based on heteroskedastic and incomplete data.
  - 3.1. Heteroskedastic Low-rank Matrix Denoising. Suppose one observes

$$(3.1) Y = X + E,$$

where X is the low-rank matrix of interest and the entries of noise E are independent, zero-mean, but need not have a common variance. The goal is to recover the singular subspace of X based on noisy observation Y. The problem arises naturally in a range of applications, such as magnetic resonance imaging (MRI) and relaxometry (Candes et al., 2013; Shabalin and Nobel, 2013). This model can also be viewed as a prototype of various problems in high-dimensional statistics and machine learning, including Poisson PCA (Salmon et al., 2014), bipartite stochastic block model (Florescu and Perkins, 2016), and exponential family PCA (Liu et al., 2016). Let the sample and population Gram matrices be  $N = YY^{\top}$  and  $M = XX^{\top}$ , respectively. Then,

$$(\mathbb{E}N)_{ij} = \begin{cases} M_{ij}, & i \neq j; \\ M_{ij} + \sum_{k=1}^{p_2} \operatorname{Var}(E_{ik}), & i = j. \end{cases}$$

Thus, entries of N are unbiased estimators only for the off-diagonal part of M. Under the heteroskedastic setting that  $\operatorname{Var}(E_{ij})$  are unequal, there can be significant differences between  $\mathbb{E}N$ ,  $\mathbb{E}\Delta(N)$ , and M on the diagonal entries, which may lead to significant perturbations on the diagonal of N-M. Here,  $\Delta(N)$  is the matrix N with diagonal entries set to zero. Since left singular vectors of Y and X are respectively identical to those of N and M, the regular SVD or diagonal-deletion SVD on Y can result in inconsistent estimates of the left singular subspace of X.

Compared to the regular or diagonal-deletion SVD on Y, the proposed HeteroPCA provides a better estimator. We have the following theoretical guarantees.

THEOREM 4 (Upper bound for heteroskedastic matrix denoising). Consider the model (3.1), where  $X \in \mathbb{R}^{p_1 \times p_2}$  is a fixed rank-r matrix and the noise matrix E consists of independent sub-Gaussian entries  $E_{ij}$  such that  $\mathbb{E}E_{ij} = 0$ ,  $\operatorname{Var}(E_{ij}) = \sigma_{ij}^2$ , and  $\max_{q \geq 1} q^{-1/2} (\mathbb{E}|E_{ij}/\sigma_{ij}|^q)^{1/q} \leq C$ . Suppose the left singular subspace of X is  $U \in \mathbb{O}_{p_1,r}$ . Assume that the condition

number of X is at most some absolute constant C, i.e.,  $||X|| \le C\lambda_r(X)$ . Denote

(3.2) 
$$\sigma_R^2 = \max_i \sum_{j=1}^{p_2} \sigma_{ij}^2, \quad \sigma_C^2 = \max_j \sum_{i=1}^{p_1} \sigma_{ij}^2, \quad \sigma_{max}^2 = \max_{ij} \sigma_{ij}^2$$

as the rowwise, columnwise, and entrywise noise variances. Then there exists a constant  $c_I > 0$  such that if U satisfies  $I(U) = \max_{1 \le i \le p_1} \frac{p_1}{r} \|e_i^\top U\|_2^2 \le c_I p_1/r$ , Algorithm 1 applied to  $YY^\top$  and rank r outputs an estimator  $\hat{U}$  that satisfies

(3.3)
$$\mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\|$$

$$\lesssim \left( \frac{\sigma_C + \sqrt{r}\sigma_{max}}{\lambda_r(X)} + \frac{\sigma_R \sigma_C + \sigma_R \sigma_{max} \sqrt{\log(p_1 \wedge p_2)} + \sigma_{max}^2 \log(p_1 \wedge p_2)}{\lambda_r^2(X)} \right) \wedge 1.$$

If  $\sigma_{max} \lesssim \sigma_C / \max\{\sqrt{r}, \sqrt{\log(p_1 \wedge p_2)}\}$  additionally holds, i.e., the variance array  $\{\sigma_{ij}^2\}$  is not too "spiky," we further have

(3.4) 
$$\mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\| \lesssim \left( \frac{\sigma_C}{\lambda_r(X)} + \frac{\sigma_R \sigma_C}{\lambda_r^2(X)} \right) \wedge 1.$$

Remark 6. Instead of HeteroPCA, one can directly apply the regular SVD or diagonal-deletion SVD:

$$\hat{U}^{\mathrm{SVD}} = \mathrm{SVD}(Y), \quad \text{and} \quad \hat{U}^{\mathrm{DD}} = \mathrm{SVD}\left(\Delta(Y^{\top}Y)\right).$$

Following the proof of Proposition 1, one can establish the lower bound to show that the proposed HeteroPCA outperforms the regular and diagonal-deletion SVDs. In particular, if  $\lambda_r(X) \gtrsim \sigma_R \vee \sigma_{\max} \sqrt{\log(p_1 \wedge p_2)}$  and  $\sigma_{\max}/\sigma_C \gtrsim \sqrt{r}$ , one can show that

$$\begin{split} & \mathbb{E} \left\| \sin \Theta \left( \hat{U}, U \right) \right\| \asymp \frac{\sigma_C}{\lambda_r(X)}, \\ & \mathbb{E} \left\| \sin \Theta \left( \hat{U}^{\text{SVD}}, U \right) \right\| \gtrsim \frac{\sigma_C + \sigma_R}{\lambda_r(X)}, \\ & \mathbb{E} \left\| \sin \Theta \left( \hat{U}^{\text{DD}}, U \right) \right\| \gtrsim 1, \end{split}$$

which illustrates the advantage of HeteroPCA.

REMARK 7. When  $\sigma_{ij} = \sigma_{\text{max}}$  for all  $1 \leq i \leq p_1, 1 \leq j \leq p_2$ , the upper bound of (3.3) reduces to

$$\mathbb{E}\left\|\sin\Theta(\hat{U}, U)\right\| \lesssim \left(\frac{\sqrt{p_1}\sigma_{\max}}{\lambda_r(X)} + \frac{\sqrt{p_1p_2}\sigma_{\max}}{\lambda_r^2(X)}\right),\,$$

which matches the optimal rate for homoskedastic matrix denoising in literature (Cai and Zhang, 2016, Theorems 3 and 4).

3.2. Poisson PCA. Poisson PCA is an important problem in statistics and engineering with a range of applications, including photon-limited imaging (Salmon et al., 2014) and biological sequencing data analysis (Cao et al., 2017). Suppose we observe  $Y \in \mathbb{R}^{p_1 \times p_2}$ , where  $Y_{ij} \stackrel{ind}{\sim} \text{Poisson}(X_{ij})$  and  $X \in \mathbb{R}^{p_1 \times p_2}$  is rank-r. Let  $X = U\Lambda V^{\top}$  be the singular value decomposition, where  $U \in \mathbb{O}_{p_1,r}, V \in \mathbb{O}_{p_2,r}$ . The goal is to estimate leading singular vectors of X, i.e., U or V, based on Y. HeteroPCA is an appropriate method for Poisson PCA since it can well handle the heteroskedasticity of Poisson distribution. Although the aforementioned heteroskedastic low-rank matrix denoising can be seen as a prototype problem of Poisson PCA, Theorem 4 is not directly applicable and more careful analysis is needed since the Poisson distribution has heavier tail than sub-Gaussian.

Theorem 5 (Poisson PCA). Suppose X is a nonnegative  $p_1$ -by- $p_2$  matrix,  $\operatorname{rank}(X) = r$ ,  $\lambda_1(X)/\lambda_r(X) \leq C$ ,  $X_{ij} \geq c$  for constant c > 0,  $U \in \mathbb{O}_{p_1,r}$  is the left singular subspace of X. Denote

(3.5) 
$$\sigma_R^2 = \max_i \sum_{j=1}^{p_2} X_{ij}, \quad \sigma_C^2 = \max_j \sum_{i=1}^{p_1} X_{ij}, \quad \sigma_*^2 = \max_{i,j} X_{ij}.$$

Suppose one observes  $Y \in \mathbb{R}^{p_1 \times p_2}, Y_{ij} \stackrel{ind}{\sim} \operatorname{Poisson}(X_{ij})$ . Then there exists constant  $c_I > 0$  such that if U satisfies  $I(U) = \max_i \frac{p_1}{r} \|e_i^\top U\|_2^2 \le c_I p_1/r$ , the proposed HeteroPCA procedure (Algorithm 1) on matrix  $YY^\top$  and rank r yield

$$(3.6) \lesssim \left(\frac{\sigma_C + r\sigma_{max}}{\lambda_r(X)} + \frac{\left\{\sigma_R + \sigma_C + \sigma_{max}\sqrt{\log(p_2)\log(p_1)}\right\}^2 - \sigma_R^2}{\lambda_r^2(X)}\right) \wedge 1.$$

In addition, if  $\sigma_{max} \leq \sigma_C / \max\{r, \sqrt{\log(p_1)\log(p_2)}\}$ , then

$$\mathbb{E}\left\|\sin\Theta(\hat{U}, U)\right\| \lesssim \left(\frac{\sigma_C}{\lambda_r(X)} + \frac{\sigma_R \sigma_C}{\lambda_r^2(X)}\right) \wedge 1.$$

REMARK 8. Similar results to Proposition 1 can be developed to show the advantage of HeteroPCA over the regular and diagonal-deletion SVD.

3.3. SVD based on Heteroskedastic and Incomplete Data. Missing data problems arise frequently in high-dimensional statistics. Let  $X \in \mathbb{R}^{p_1 \times p_2}$  be a rank-r unknown matrix. Suppose only a small fraction of entries of X, denoted by  $\Omega \subseteq [p_1] \times [p_2]$ , are observable with random noise,

$$Y_{ij} = X_{ij} + Z_{ij}, \quad (i,j) \in \Omega.$$

Here, each entry  $Y_{ij}$  is observed or missing with probability  $\theta$  or  $1 - \theta$  for some  $0 < \theta < 1$  and  $Z_{ij}$ 's are independent, zero-mean, and possibly heteroskedastic. Let  $R \in \mathbb{R}^{p_1 \times p_2}$  be the indicator of observable entries:

$$R_{ij} = \begin{cases} 1, & (i,j) \in \Omega; \\ 0, & (i,j) \notin \Omega, \end{cases}$$

and R and Y are independent. Assume  $X = U\Lambda V^{\top}$  is the singular value decomposition, where  $U \in \mathbb{O}_{p_1,r}$  and  $V \in \mathbb{O}_{p_2,r}$ . We specifically aim to estimate U based on  $\{Y_{ij}, (i,j) \in \Omega\}$ . Denote Y as the entry-wise product of Y and R, i.e.,  $\tilde{Y}_{ij} = Y_{ij}R_{ij}, \forall (i,j) \in [p_1] \times [p_2]$ . Since  $\mathbb{E}\tilde{Y}_{ij} = \theta X_{ij}$  and  $\mathrm{Var}(\tilde{Y}_{ij})$  are not necessarily identical for different (i,j) pairs, we can apply HeteroPCA on  $\tilde{Y}\tilde{Y}^{\top}$  to estimate U. The following theoretical guarantee holds.

THEOREM 6. Let X be a  $p_1$ -by- $p_2$  rank-r matrix, whose left singular subspace is denoted by  $U \in \mathbb{O}_{p_1,r}$ . Assume that  $\mathbb{E}Y = X$ . Let Y consist of sub-Gaussian entries in the sense that  $\max_{ij} \|Y_{ij}\|_{\psi_2} \leq C$ . Here,  $\|Y\|_{\psi_2} \triangleq \sup_{q \geq 1} q^{-1/2} (|Y|^q)^{1/q}$  is the Orlicz- $\psi_2$  norm of Y. Suppose  $0 < \theta \leq 1 - c$  for constant c > 0. There exists constant  $c_I > 0$  such that if  $U \in \mathbb{O}_{p_1 \times r}$  satisfies  $I(U) \|X\| / \lambda_r(X) \leq c_I p_1 / r$ , HeteroPCA applied to  $\tilde{Y}\tilde{Y}^{\top}$  outputs an estimator  $\hat{U}$  satisfying

$$(3.7) \qquad \left\| \sin \Theta(\hat{U}, U) \right\| \lesssim \frac{\max \left\{ \sqrt{p_2(\theta + \theta^3 p_1^2) \log(p_1)}, \theta p_1 \log^2(p_1) \right\}}{\theta^2 \lambda_r^2(X)} \wedge 1$$

with probability at least  $1 - p_1^{-C}$ .

REMARK 9 (Comparison with matrix completion). Our work is related to a substantial body of literature on low-rank matrix completion. For example, Candès and Recht (2009); Candès and Tao (2010); Recht (2011) analyzed the performance of nuclear norm minimization; Toh and Yun (2010)

developed the accelerated proximal gradient descent algorithm and software package  $NNLS^1$ ; Mazumder et al. (2010) introduced the spectral regularization algorithm for incomplete matrix learning and developed the software package  $SoftImpute^2$ ; Keshavan et al. (2010a,b); Keshavan (2012); Jain et al. (2013) analyzed the alternating gradient descent and spectral algorithm for matrix completion with/without noise; Vaswani and Narayanamurthy (2017, 2018) studied the performance of regular SVD estimator when there are missing values and the noise is non-isotropic and sparsely data-dependent; Robin et al. (2019) studied the low-rank model for count data with missing values; also see Cai and Wei (2018) for a recent survey of matrix completion. Different from previous matrix completion works, our goal here is to estimate the singular subspace  $U \in \mathbb{O}_{p_1,r}$  rather than the whole matrix  $X \in \mathbb{R}^{p_1 \times p_2}$ . We apply HeteroPCA to impute the diagonal entries of  $XX^{\top}$ , not the missing entries in X itself in most of the aforementioned matrix completion literature.

In addition, in the case that the average amplitude of all entries is a constant,  $||X||_F^2/(p_1p_2) \approx 1$ , and X has a good condition number,  $\lambda_1(X) \approx \lambda_r(X)$ , Theorem 6 implies that as long as the expected sample size satisfies (3.8)

$$\mathbb{E}|\Omega| \gg \max \left\{ p_1^{1/3} p_2^{2/3} r^{2/3} \log^{1/3}(p_1), p_1 r^2 \log(p_1), p_1 r \log(p_1) \log(p_1 p_2) \right\},$$

the HeteroPCA estimator is consistent. This requirement is weaker than ones in classic matrix completion literature (Candès and Tao, 2010; Keshavan et al., 2010a; Recht, 2011)

$$|\Omega| \gtrsim (p_1 + p_2)r \cdot \text{polylog}(p),$$

whose goal is to estimate the whole matrix X. When  $p_2 \gg p_1$ , (3.8) implies that HeteroPCA estimator is consistent, even if most columns of X are completely missing. To the best of our knowledge, we are among the first to give such a result.

REMARK 10 (Time complexity). If the target matrix X is  $p_1$ -by- $p_2$  and rank-r, the time complexity of HeteroPCA, regular SVD, diagonal-deletion SVD, OptShrink (Nadakuditi, 2014), and SoftImpute (Mazumder et al., 2010) are  $O(|\Omega|^2/p_2 + Tp_1^2r)$ ,  $O(T(|\Omega|r + r^3))$ ,  $O(|\Omega|^2/p_2 + Tp_1^2r)$ ,  $O(T(|\Omega|r + r^3))$ , and  $O(T(|\Omega| + p_1p_2r))$ , respectively. Here, T denotes the number of iterations in each method.

<sup>1</sup>http://www.math.nus.edu.sg/~mattohkc/NNLS.html

<sup>2</sup>https://cran.r-project.org/web/packages/softImpute/index.html

Remark 11. PCA based on heteroskedastic and incomplete data is another closely related problem. Although most existing literature on PCA with incomplete data focused on regular SVD methods under the homoskedastic noisy setting (see, e.g., Lounici et al. (2014); Cai and Zhang (2016)), we are able to achieve better performance by applying the proposed HeteroPCA algorithm if the noise is heteroskedastic. To be more specific, suppose one observes incomplete i.i.d. samples  $Y_1, \ldots, Y_n \in \mathbb{R}^p$  from the generalized spiked covariance model,

$$Y = X + \varepsilon \in \mathbb{R}^{p}, \quad \mathbb{E}X = \mu, \operatorname{Cov}(X) = U\Lambda U^{\top},$$

$$\mathbb{E}\varepsilon = 0, \operatorname{Cov}(\varepsilon) = D = \operatorname{diag}(\sigma_{1}^{2}, \dots, \sigma_{p}^{2}),$$

$$\varepsilon = ((\varepsilon)_{1}, \dots, (\varepsilon)_{p})^{\top}; \quad X, (\varepsilon)_{1}, \dots, (\varepsilon)_{p} \text{ are independent};$$

$$\forall k = 1, \dots, n, \quad Y_{k} = (Y_{1k}, \dots, Y_{pk})^{\top}, \quad Y_{1}, \dots, Y_{n} \text{ are i.i.d. copies of } Y;$$

$$\forall 1 \leq i \leq p, 1 \leq k \leq n, \quad R_{ik} = \begin{cases} 1, & Y_{ik} \text{ is observable}; \\ 0, & Y_{ik} \text{ is missing}, \end{cases}$$

and  $\{R_{ik}\}_{1\leq i\leq p, 1\leq k\leq n}$  are independent of  $Y_1, \ldots, Y_n$ . To estimate the leading factor loadings, i.e.,  $U\in \mathbb{O}_{p,r}$ , we can first evaluate the generalized sample covariance matrix,

$$\hat{\Sigma}^* = (\hat{\sigma}_{ij}^*)_{1 \le i, j \le p}, \quad \text{with} \quad \hat{\sigma}_{ij}^* = \frac{\sum_{k=1}^n (Y_{ik} - \bar{Y}_i^*)(Y_{ik} - \bar{Y}_j^*)R_{ik}R_{jk}}{\sum_{k=1}^n R_{ik}R_{jk}}$$
and 
$$\bar{Y}_i^* = \frac{\sum_{k=1}^n Y_{ik}R_{ik}}{\sum_{k=1}^n R_{ik}}.$$

Then estimate U by applying Algorithm 1 on  $\hat{\Sigma}^*$ . A similar consistent upper bound result to Theorem 6 can be developed for this procedure.

4. Numerical Results. In this section, we perform simulation studies to further illustrate the merit of proposed procedure in singular subspace estimation when heteroskedastic noise is in presence. All simulation results below are based on the average of 1000 repeated independent experiments. The average and the standard deviation of estimation errors are indicated by markers and error bars, respectively.

We first consider PCA under the generalized spiked covariance model (1.2). For various values of p, n, and r, we generate a p-by-r random matrix  $U_0$  with i.i.d. standard Gaussian entries,  $w_1, \ldots, w_p \stackrel{iid}{\sim} \text{Unif}[0,1]$ , and  $\sigma_1, \ldots, \sigma_p \stackrel{iid}{\sim} \text{Unif}[0,1]$ . The purpose of generating uniform random vectors  $w, \sigma$  is to introduce heteroskedasticity into observations. Then, we let

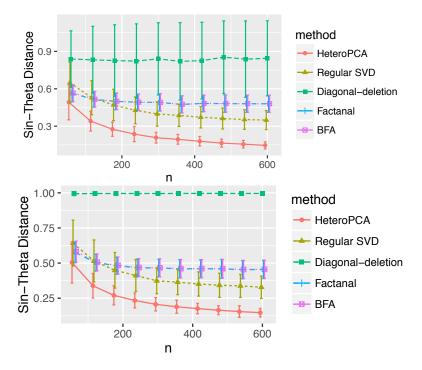


Fig 1. Average  $\sin \Theta$  loss of heteroskedastic PCA versus sample size n. Upper panel: r=3; lower panel: r=5.

 $U = \operatorname{QR}(U_0\operatorname{diag}(w)) \in \mathbb{O}_{p,r}$  and  $\Sigma_0 = UU^{\top} \in \mathbb{R}^{p \times p}$ . We aim to recover U based on i.i.d. observations  $\{Y_k = X_k + \varepsilon_k\}_{k=1}^n$ , where  $X_1, \ldots, X_n \stackrel{iid}{\sim} N(0, \Sigma_0), \varepsilon_1, \ldots, \varepsilon_n \stackrel{iid}{\sim} N(0, \operatorname{diag}(\sigma_1^2, \ldots, \sigma_n^2))$ . We implement the proposed HeteroPCA, diagonal-deletion, and regular SVD approaches and plot the average estimation errors and standard deviation in  $\sin \Theta$  distance. We also implement the classic factor analysis methods (Thomson, 1939; Lawley and Maxwell, 1962), factanal function in R stats package, and the Bayesian factor analysis method, MCMCfactanal function from R MCMCpack package (Martin et al., 2011). The simulation results are summarized to Figure 1. It can be seen that the proposed HeteroPCA estimator significantly outperforms other methods; the regular SVD yields larger estimation error; and the diagonal-deletion estimator performs unstably across different settings. This matches the theoretical findings in Section 2.

Next we study how the degree of heteroskedasticity affects the perfor-

mance. Let

$$v_1, \dots, v_p \stackrel{iid}{\sim} \text{Unif}[0, 1], \quad \sigma_k^2 = \frac{0.1 \cdot p \cdot v_k^{\alpha}}{\sum_{i=1}^p v_i^{\alpha}}, \quad k = 1, \dots, p.$$

In such case,  $\sigma_{\text{sum}}^2 = \sigma_1^2 + \cdots + \sigma_p^2$  always equals 0.1p and  $\alpha$  characterizes the degree of heteroskedasticity: the larger  $\alpha$  results more imbalanced distribution of  $(\sigma_1, \ldots, \sigma_p)$ ; if  $\alpha = 0$ ,  $\sigma_1 = \cdots = \sigma_p$  and the setting becomes homoskedastic. Now we generate  $U, \Sigma_0$  and  $\{Y_k, X_k, \varepsilon_k\}_{k=1}^n$  in the same way as the previous setting. Due to computational issues of factor analysis methods, we only compare HeteroPCA with regular SVD and diagonal-deletion estimator here. The average estimation errors for U are plotted in Figure 2. The results again suggest that the performance of diagonal-deletion estimator is unstable across different settings. When  $\alpha = 0$ , i.e., the noise is homoskedastic, the performance of HeteroPCA and regular SVD are comparable; but as  $\alpha$  increases, the estimation error of HeteroPCA grows significantly slower than that of the regular SVD, which is consistent with the theoretical results in Theorem 1.

Then, we consider the problem of denoising a low-rank matrix with heteroskedastic noise discussed in Section 3.1. Let  $U_0 \in \mathbb{R}^{p_1 \times r}$  and  $V_0 \in \mathbb{R}^{p_2 \times r}$  be i.i.d. Gaussian ensembles for  $(p_1, p_2) = (50, 200), (200, 1000)$  and r = 3. To introduce heteroskedasticity, we also randomly draw  $w, v_1 \in \mathbb{R}^{p_1}$ , and  $v_2 \in \mathbb{R}^{p_2}$  with i.i.d. Unif[0, 1] entries. Then we evaluate  $U = \operatorname{QR}(U_0 \cdot \operatorname{diag}(w)^4)$ ,  $V = \operatorname{QR}(V_0)$ , and construct the signal matrix  $X = (p_1p_2)^{1/4} \cdot UV^{\top}$ . The noise matrix is drawn as  $E_{ij} \stackrel{ind}{\sim} N(0, \sigma_0^2 \cdot \sigma_{ij}^2)$ , where  $\sigma_{ij} = (v_1)_i^4 \cdot (v_2)_j^4$ ,  $\sigma_0$  varies from 0 to 2,  $1 \leq i \leq p_1$ , and  $1 \leq j \leq p_2$ . Based on the  $p_1$ -by- $p_2$  observation Y = X + E, we implement HeteroPCA with input of  $YY^{\top}$ , regular-SVD, diagonal-deletion, and OptShrink (Nadakuditi, 2014) methods and plot the average  $\sin \Theta$  distance error in Figures 3 (a) - (d). For each of estimators  $\hat{U}$  and  $\hat{V}$ , we also estimate X by  $\hat{X} = \hat{U}\hat{U}^{\top}Y\hat{V}\hat{V}^{\top}$  and plot the Frobenius norm error in Figure 3 (e) and (f). As one can clearly see from Figure 3, the proposed HeteroPCA outperforms other methods in all estimations for U, V, and X, and the advantage of HeteroPCA is more significant when the noise level increases.

Next, we consider the Poisson PCA problem. We generate  $U_0 \in \mathbb{R}^{p_1 \times r}$  and  $V_0 \in \mathbb{R}^{p_2 \times r}$  with i.i.d. standard normal entries for  $(p_1, p_2, r) = (50, 500, 3)$  or (200, 1000, 3). Similarly to previous settings, we introduce heteroskedasticity by generating a vector  $w \in \mathbb{R}^{p_1}$  with i.i.d. Unif[0, 1] entries. Let  $U = |U_0 \cdot \operatorname{diag}(w)^4| \in \mathbb{R}^{p_1 \times r}, V = |V_0| \in \mathbb{R}^{p_2 \times r}, X = \lambda (p_1 p_2)^{1/4} UV^{\top} \in \mathbb{R}^{p_1 \times p_2}$ , and  $Y_{ij} \sim \operatorname{Poisson}(X_{ij})$  independently. Here,  $\lambda > 0$  measures the signal strength. The performance of HeteroPCA, regular SVD, diagonal-deletion,

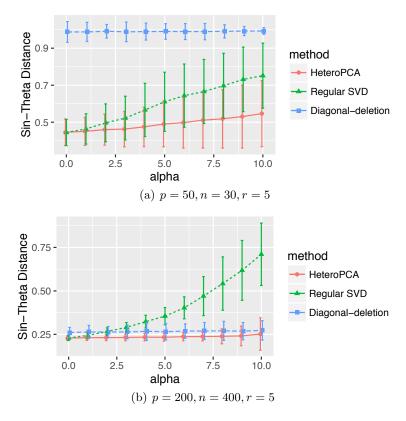


FIG 2. Average  $\sin \Theta$  loss for heteroskedastic PCA versus heteroskedastic level  $\alpha$ .

and OptShrink on estimation of left singular subspaces are provided in Figure 4. These plots again illustrate the merit of the proposed HeteroPCA method.

Finally, in the following experiment we study SVD based on heteroskedastic and incomplete data in the setting of Section 3.3. Generate  $Y, X, Z \in \mathbb{R}^{p_1 \times p_2}$  in the same way as the previous heteroskedastic SVD setting with  $p_1 = 50, 100, r = 3, 5, \sigma_0 = .2$ , and  $p_2$  ranging from 800 to 3200. Each entry of Y is observed independently with probability  $\theta = 0.1$ . We aim to estimate U based on  $\{Y_{ij}: (i,j) \in \Omega\}$ . In addition to HeteroPCA, regular SVD, diagonal-deletion SVD, and OptShrink, we also apply the nuclear norm minimization (Mazumder et al., 2010, Soft-Impute package)

$$\hat{X}_* = \arg\min_{\hat{X} \in \mathbb{R}^{p_1 \times p_2}} \sum_{(i,j) \in \Omega} (\tilde{Y}_{ij} - \hat{X}_{ij})^2 + \nu ||\hat{X}||_*, \quad \hat{U} = \text{SVD}_r(\hat{X}).$$

To avoid the cumbersome issue of parameter  $\nu$  selection, we evaluate the

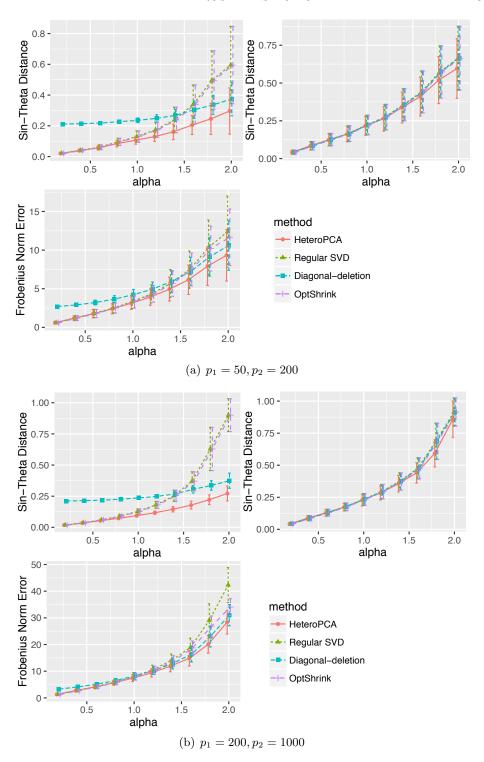


FIG 3. Estimation errors of  $\hat{U}$  (top left),  $\hat{V}$  (top right), and  $\hat{X}$  (bottom left) in heteroskedastic matrix denoising.

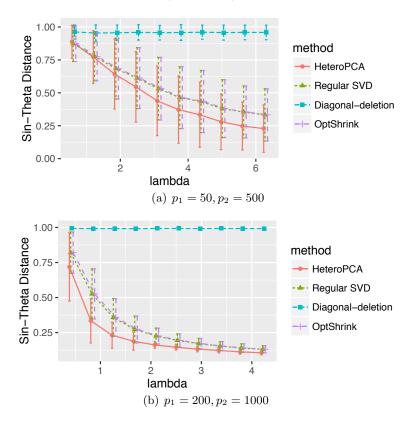


Fig 4. Estimation errors in Poisson PCA for a ranging value of signal strength  $\lambda$ 

above nuclear norm minimization estimator for a grid of values of  $\nu$ , then record the outcome with the minimum  $\sin\Theta$  distance error  $\|\sin\Theta(\hat{U},U)\|$ . From the results plotted in Figure 5, we can see that HeteroPCA significantly outperforms all other methods.

5. Discussions. We consider PCA in the presence of heteroskedastic noise in this paper. To alleviate the significant bias incurred on diagonal entries of the Gram matrix due to heteroskedastic noise, we introduced a new procedure named HeteroPCA that adaptively imputes diagonal entries to remove the bias. The proposed procedure achieves optimal rate of convergence in a range of settings. In addition, we discuss several applications of the proposed algorithm, including heteroskedastic low-rank matrix denoising, Poisson PCA, and SVD based on heteroskedastic and incomplete data.

The proposed HeteroPCA procedure can also be applied to many other

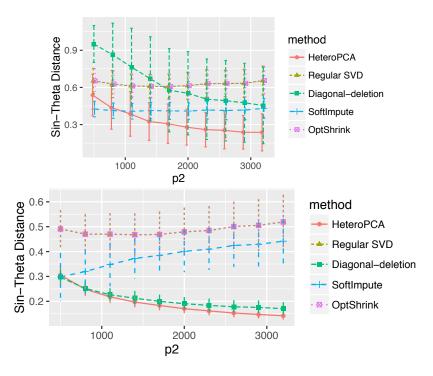


FIG 5. Average  $\sin \Theta$  distance error for SVD based on heteroskedastic and incomplete data. Here,  $p_1 = 50, r = 5, \theta = .2$  (Upper Panel) and  $p_1 = 100, r = 3, \theta = .2$  (Lower Panel);  $p_2$  varies from 800 to 3200.

problems where the noise is heteroskedastic. First, exponential family PCA is a commonly used technique for dimension reduction on non-real-valued datasets (Collins et al., 2002; Mohamed et al., 2009). As discussed in the introduction, the exponential family distributions, e.g., exponential, binomial, and negative binomial, may be highly heteroskedastic. As in the case of Poisson PCA considered in Section 3.2, the proposed HeteroPCA algorithm can be applied to exponential family PCA.

In addition, community detection in social network has attracted significant attention in the recent literature (Fortunato, 2010; Newman, 2013). Although most of existing results focused on unipartite graphs, bipartite graphs, i.e., all edges are between two groups of nodes, often appear in practice (Melamed, 2014; Florescu and Perkins, 2016; Alzahrani and Horadam, 2016; Zhou and Amini, 2018). The proposed HeteroPCA can also be applied to community detection for bipartite stochastic block model. Similarly to the analysis for heteroskedastic low-rank matrix denoising in Section 3.1, HeteroPCA can be shown to have advantages over other baseline methods.

The proposed framework is also applicable to solve the *heteroskedastic* tensor SVD problem, which aims to recover the low-rank structure from the tensorial observation corrupted by heteroskedastic noise. Suppose one observes  $\mathbf{Y} = \mathbf{X} + \mathbf{Z} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ , where  $\mathbf{X}$  is a Tucker low-rank signal tensor and  $\mathbf{Z}$  is the noise tensor with independent and zero-mean entries. If  $\mathbf{Z}$  is homoskedastic, the higher-order orthogonal iteration (HOOI) (De Lathauwer et al., 2000) was shown to achieve the optimal performance for recovering  $\mathbf{X}$  (Zhang and Xia, 2018). If  $\mathbf{Z}$  is heteroskedastic, we can apply HeteroPCA instead of the regular SVD to obtain a better initialization for HOOI. Similarly to the argument in this article, we are able to show that this modified HOOI yields more stable and accurate estimates than the regular HOOI.

Canonical correlation analysis (CCA) is one of the most important tools in multivariate analysis for exploring the relationship between two sets of vector samples (Hotelling, 1936). In the standard procedure of CCA, the core step is a regular SVD on the adjusted cross-covariance matrix between samples. When the observations contain heteroskedastic noise, one can replace the regular SVD procedure by HeteroPCA to achieve better performance.

Finally, it is interesting to further study PCA when the noise has heteroskedastic and dependent structure. Suppose one observes i.i.d. samples  $\{Y_k\}_{k=1}^n$  that admits a signal-noise decomposition:  $Y_k = X_k + \varepsilon_k$ . We still assume the covariance matrix of  $X_k$  is low-rank. To ensure that the noise and signal parts are distinguishable, some structural conditions on  $\text{Cov}(\varepsilon)$  are required. In addition to the focus of this paper that  $\text{Cov}(\varepsilon)$  is diagonal, it is interesting to explore other commonly-used high-dimensional covariance structures, such as sparsity, bandedness, etc. These structures result in bias in various index sets of the sample covariance matrix, other than the diagonal set discussed in this paper. HeteroPCA may be modified accordingly to alleviate the bias in the sample covariance matrix.

- **6. Proofs.** In this section, we prove the main results, namely, Theorems 1 and 3. For reasons of space, the remaining proofs are given in Section A of the supplementary materials (Zhang et al., 2018).
  - 6.1. Proofs for Heteroskedastic PCA.

PROOF OF THEOREM 1. Based on the generalized spiked covariance model, we introduce

$$E = [\varepsilon_1, \dots, \varepsilon_n] \in \mathbb{R}^{p \times n}, \quad \gamma_k = \Lambda^{1/2} U^\top (X_k - \mu) \in \mathbb{R}^r, \quad \Gamma = [\gamma_1, \dots, \gamma_n] \in \mathbb{R}^{r \times n}.$$

Then the observations can be written as

$$Y_k = X_k + \varepsilon_k = \mu + U\Lambda^{1/2}\gamma_k + \varepsilon_k$$
, or  $Y = \mu \mathbf{1}_n^{\top} + U\Lambda^{1/2}\Gamma + E$ ,

where  $\mu \in \mathbb{R}^p$  is a fixed vector,  $\mathbb{E}\gamma_k = 0$ ,  $\operatorname{Cov}(\gamma_k) = I$ , E has independent entries, and  $\Gamma$  has independent columns. We also denote  $\bar{X} \in \mathbb{R}^p$ ,  $\bar{E} \in \mathbb{R}^p$ ,  $\bar{\Gamma} \in \mathbb{R}^r$  as the averages of all columns of X, E, and  $\Gamma$ , respectively. Since  $\hat{\Sigma}$  is invariant after any translation on Y, we can assume  $\mu = 0$  without loss of generality. The rest of the proof is divided into three steps for the sake of presentation.

Step 1 We define  $\hat{\Sigma}_X = (XX^\top - n\bar{X}\bar{X}^\top)/(n-1)$  as the sample covariance of signal vectors. The aim of this step is to develop a concentration inequality for  $\hat{\Sigma} - \Sigma_X$ . To this end, we consider the following decomposition of  $n(\hat{\Sigma} - \hat{\Sigma}_X)$ ,

$$(6.1)$$

$$(n-1)(\hat{\Sigma} - \hat{\Sigma}_X) = (n-1)\hat{\Sigma} - (XX^{\top} - n\bar{X}\bar{X}^{\top})$$

$$= YY^{\top} - n\bar{Y}\bar{Y}^{\top} - (XX^{\top} - n\bar{X}\bar{X})$$

$$= (X+E)(X+E)^{\top} - (XX^{\top} - n\bar{X}\bar{X}^{\top}) - n\left(\bar{X}\bar{X}^{\top} + \bar{X}\bar{E}^{\top} + \bar{E}\bar{X}^{\top} + \bar{E}\bar{E}^{\top}\right)$$

$$= XE^{\top} + EX^{\top} + EE^{\top} - n\left(\bar{X}\bar{E}^{\top} + \bar{E}\bar{X}^{\top} + \bar{E}\bar{E}^{\top}\right).$$

We analyze each term of (6.1) separately as follows. Since E has independent entries and  $Var(E_{ij}) = \sigma_i^2$ , the rowwise structured heteroskedastic concentration inequality (c.f., Cai and Zhang (2018a)) implies

(6.2) 
$$\mathbb{E} \left\| E E^{\top} - \mathbb{E} E E^{\top} \right\| \lesssim \sqrt{n} \sigma_{\text{sum}} \sigma_{\text{max}} + \sigma_{\text{sum}}^2.$$

Since X is deterministic, E is random, and  $\mathbb{E}E = 0$ , we have  $\mathbb{E}EX^{\top} = 0$ . By Lemma 3 in the supplementary materials,

$$(6.3) \quad \mathbb{E}_{E}\left(\left\|EX^{\top} - \mathbb{E}EX^{\top}\right\| \middle| X\right) = \mathbb{E}_{E}\left(\left\|EX^{\top}\right\| \middle| X\right)$$

$$\lesssim \|X\| \left(\sigma_{C} + r^{1/4}\sigma_{\max}\sigma_{C} + \sqrt{r}\sigma_{\max}\right)^{\text{Cauchy-Schwarz}} \lesssim \|X\| \left(\sigma_{\text{sum}} + \sqrt{r}\sigma_{\max}\right).$$
Since  $\mathbb{E}\|\bar{E}\|_{2}^{2} = \sum_{i=1}^{p} \mathbb{E}(\bar{E})_{i}^{2} = \sum_{i=1}^{p} \sigma_{i}^{2}/n = \sigma_{\text{sum}}^{2}/n, \text{ we have}$ 

$$\mathbb{E}_{E}\left(n\left\|\bar{X}\bar{E}^{\top} + \bar{E}\bar{X}^{\top} + \bar{E}\bar{E}^{\top}\right\|\right)$$

$$\leq \mathbb{E}_{E}n\|\bar{X}\bar{E}^{\top}\| + \mathbb{E}_{E}n\|\bar{E}\bar{X}^{\top}\| + \mathbb{E}_{E}n\|\bar{E}\bar{E}^{\top}\|$$

$$\leq \mathbb{E}_{E}2n\|\bar{X}\|_{2}\|\bar{E}\|_{2} + \mathbb{E}n\|\bar{E}\|_{2}^{2}$$

$$\leq 2n\|\bar{X}\|_{2} \cdot (\mathbb{E}\|\bar{E}\|_{2}^{2})^{1/2} + \mathbb{E}n\|\bar{E}\|_{2}^{2} \leq 2n^{1/2}\sigma_{\text{sum}}\|\bar{X}\|_{2} + \sigma_{\text{sum}}^{2}.$$

Combining (6.2), (6.3), and (6.4), we have

$$\mathbb{E}_{E} \left\| (n-1)(\hat{\Sigma} - \hat{\Sigma}_{X}) - \mathbb{E}EE^{\top} \right\|$$

$$\lesssim \sqrt{n}\sigma_{\text{sum}}\sigma_{\text{max}} + \sigma_{\text{sum}}^{2} + \|X\|(\sigma_{\text{sum}} + \sqrt{r}\sigma_{\text{max}}) + n^{1/2}\|\bar{X}\|_{2}\sigma_{\text{sum}}.$$

Noting that  $\mathbb{E}EE^{\top} = n \operatorname{diag}(\sigma_1^2, \dots, \sigma_p^2)$  is diagonal and  $\Delta(\cdot)$  is the operator that sets all diagonal entries to zero, we further have

$$\mathbb{E}_{E} \left\| \Delta \left( (n-1)(\hat{\Sigma} - \hat{\Sigma}_{X}) \right) \right\| = \mathbb{E}_{E} \left\| \Delta \left( (n-1)(\hat{\Sigma} - \hat{\Sigma}_{X}) - \mathbb{E}EE^{\top} \right) \right\|$$

$$\stackrel{\text{Lemma 5}}{\leq} 2\mathbb{E}_{E} \left\| (n-1) \left( \hat{\Sigma} - \hat{\Sigma}_{X} \right) - \mathbb{E}EE^{\top} \right\|$$

$$\stackrel{\lesssim}{\sqrt{n}} \sigma_{\text{sum}} \sigma_{\text{max}} + \sigma_{\text{sum}}^{2} + \|X\| (\sigma_{\text{sum}} + \sqrt{r}\sigma_{\text{max}}) + n^{1/2} \|\bar{X}\|_{2} \sigma_{\text{sum}}.$$

Since  $\operatorname{rank}(\hat{\Sigma}_X) \leq r$ , the eigenvectors of  $\hat{\Sigma}_X$  are U, and U satisfies the incoherence condition:  $I(U) \leq c_I p/r$ , the robust  $\sin \Theta$  Theorem (Theorem 3) yields

(6.5)  $\mathbb{E}_{E} \left\| \sin \Theta(\hat{U}, U) \right\| \lesssim \frac{\mathbb{E}_{E} \left\| \Delta \left( (n-1)(\hat{\Sigma} - \hat{\Sigma}_{X}) \right) \right\|}{\lambda_{r} \left( (n-1)\hat{\Sigma}_{X} \right)} \wedge 1$   $\lesssim \frac{\sqrt{n}\sigma_{\text{sum}}\sigma_{\text{max}} + \sigma_{\text{sum}}^{2} + \|X\|(\sigma_{\text{sum}} + \sqrt{r}\sigma_{\text{max}}) + n^{1/2}\|\bar{X}\|_{2}\sigma_{\text{sum}}}{\lambda_{r} \left( (n-1)\hat{\Sigma}_{X} \right)} \wedge 1.$ 

Step 2 Next, we study the expectation of the target function with respect to X. We specifically need to study  $\lambda_r(n\hat{\Sigma}_X)$ , ||X||, and  $||\bar{X}||_2$ . Since  $\Gamma \in \mathbb{R}^{r \times n}$  has independent columns and each column is isotropic sub-Gaussian distributed, based on the random matrix theory (Vershynin, 2010, Corollary 5.35),

$$\mathbb{P}\left(\sqrt{n} + C\sqrt{r} + t \ge \|\Gamma\| \ge \lambda_r(\Gamma) \ge \sqrt{n} - C\sqrt{r} - t\right) \le \exp(-Ct^2/2).$$

In addition,  $\sqrt{n}\bar{\Gamma} \in \mathbb{R}^r$  is a sub-Gaussian vector satisfying

$$\max_{q \ge 1} \max_{\|v\|_2 \le 1} q^{-1/2} \left( \mathbb{E} |v^\top \cdot \sqrt{n} \bar{\Gamma}|^q \right)^{1/q} \le C$$

for any  $v \in \mathbb{R}^r$ . By the Bernstein-type concentration inequality (Vershynin, 2010, Proposition 5.16),

$$\mathbb{P}\left(\|\sqrt{n}\bar{\Gamma}\|_{2}^{2} \ge r + C\sqrt{rx} + Cx\right) \le C \exp(-cx).$$

If  $n \ge Cr$  for some large constant C > 0, by setting  $t = c\sqrt{n}$  and x = cn in the previous two inequalities, we have

(6.6) 
$$2\sqrt{n} \ge ||\Gamma|| \ge \lambda_r(\Gamma) \ge \sqrt{n}/2$$
, and  $||\sqrt{n}\bar{\Gamma}||_2 \le \sqrt{n}/3$ 

with probability at least  $1 - C \exp(-cn)$ . When (6.6) holds,

$$\lambda_{r}(n\hat{\Sigma}_{X}) = \lambda_{r} \left( n(XX^{\top} - n\bar{X}\bar{X}^{\top}) \right) = \lambda_{r} \left( nU\Lambda^{1/2}(\Gamma\Gamma^{\top} - n\bar{\Gamma}\bar{\Gamma}^{\top})\Lambda^{1/2}U^{\top} \right)$$

$$\geq \lambda_{r}(\Lambda) \cdot \lambda_{r} \left( \Gamma\Gamma^{\top} - n\bar{\Gamma}\bar{\Gamma}^{\top} \right) \geq \lambda_{r}(\Lambda) \left( \lambda_{r}^{2}(\Gamma) - \|\sqrt{n}\bar{\Gamma}\|_{2}^{2} \right)$$

$$\stackrel{(6.4)}{\geq} \lambda_{r}(\Lambda) \left( n/4 - n/9 \right) \gtrsim n\lambda_{r}(\Lambda);$$

(6.7) 
$$||X|| \le ||U\Lambda^{1/2}\Gamma|| \le ||\Lambda^{1/2}|| \cdot ||\Gamma|| \stackrel{(6.6)}{\le} 2\sqrt{n}||\Lambda^{1/2}|| \lesssim \sqrt{n}\lambda_r^{1/2}(\Lambda),$$

where the last inequality is due to the assumption that  $\|\Lambda\|/\lambda_r(\Lambda) \leq C$  for some constant C.

$$\|\bar{X}\|_2 = \|U\Lambda^{1/2}\bar{\Gamma}\|_2 \le \|\Lambda^{1/2}\| \cdot \|\bar{\Gamma}\|_2 \lesssim \lambda_r^{1/2}(\Lambda).$$

By combining the previous three inequalities and (6.5), we know if (6.6) holds,

$$(6.8) \quad \mathbb{E}_{E} \left\| \sin \Theta(\hat{U}, U) \right\| \\ \lesssim \frac{\sqrt{n} \sigma_{\text{sum}} \sigma_{\text{max}} + \sigma_{\text{sum}}^{2} + (n \lambda_{r}(\Lambda))^{1/2} (\sigma_{\text{sum}} + \sqrt{r} \sigma_{\text{max}}) + (n \lambda_{r}(\Lambda))^{1/2} (\Lambda) \sigma_{\text{sum}}}{n \lambda_{r}(\Lambda)} \wedge 1 \\ \lesssim \left( \frac{\sigma_{\text{sum}} + \sqrt{r} \sigma_{\text{max}}}{(n \lambda_{r}(\Lambda))^{1/2}} + \frac{\sqrt{n} \sigma_{\text{sum}} \sigma_{\text{max}} + \sigma_{\text{sum}}^{2}}{n \lambda_{r}(\Lambda)} \right) \wedge 1 \\ \lesssim \left( \frac{\sigma_{\text{sum}} + \sqrt{r} \sigma_{\text{max}}}{(n \lambda_{r}(\Lambda))^{1/2}} + \frac{\sigma_{\text{sum}} \sigma_{\text{max}}}{n^{1/2} \lambda_{r}(\Lambda)} \right) \wedge 1.$$

Here, the last " $\lesssim$ " is due to  $\sigma_{\text{sum}}^2/(n\lambda_r(\Lambda)) \wedge 1 \leq \sigma_{\text{sum}}/(n\lambda_r(\Lambda))^{1/2} \wedge 1$ . Step 3 Finally,

$$\mathbb{E}\|\sin\Theta(\hat{U},U)\| = \mathbb{E}\|\sin\Theta(\hat{U},U)\|1_{\{(6.6) \text{ holds}\}} + \mathbb{E}\|\sin\Theta(\hat{U},U)\|1_{\{(6.6) \text{ does not hold}\}}$$

$$\lesssim \left(\frac{\sigma_{\text{sum}} + \sqrt{r}\sigma_{\text{max}}}{(n\lambda_r(\Lambda))^{1/2}} + \frac{\sigma_{\text{sum}}\sigma_{\text{max}}}{n^{1/2}\lambda_r(\Lambda)}\right) \wedge 1 + \mathbb{P}\left((6.6) \text{ does not hold}\right)$$

$$\lesssim \left(\frac{\sigma_{\text{sum}} + \sqrt{r}\sigma_{\text{max}}}{(n\lambda_r(\Lambda))^{1/2}} + \frac{\sigma_{\text{sum}}\sigma_{\text{max}}}{n^{1/2}\lambda_r(\Lambda)}\right) \wedge 1 + C \exp(-cn)$$

$$\lesssim \left(\frac{\sigma_{\text{sum}} + \sqrt{r}\sigma_{\text{max}}}{(n\lambda_r(\Lambda))^{1/2}} + \frac{\sigma_{\text{sum}}\sigma_{\text{max}}}{n^{1/2}\lambda_r(\Lambda)}\right) \wedge 1.$$

The last inequality is due to the assumption that  $\lambda_r(\Lambda) \geq c \exp(-cn)$ . Therefore, we have finished the proof of this theorem.

6.2. Proof of Theorem 3. In this subsection we prove a more general version of Theorem 3, where the corrupted entries lie in a known set  $\mathcal{G} \subset [p] \times [p]$  which need not be the diagonal. Recall the model (2.13), where we observe a symmetric  $p \times p$  matrix N = M + Z, where M is a rank-r matrix of interest and Z is the perturbation. Our goal is to estimate  $U \in \mathbb{O}_{p,r}$ , consisting of the eigenvectors of M. Extending the ideas of Algorithm 1 for HeteroPCA, Algorithm 2 provides a robust estimate of U which iteratively impute the values in the corrupted entries in  $\mathcal{G}$ . In the special case where  $\mathcal{G}$  is the diagonal, i.e.,  $\mathcal{G} = \{(i,i) : 1 \leq i \leq p\}$ , Algorithm 2 reduces to Algorithm 1.

## Algorithm 2 Generalized HeteroPCA

```
1: Input: matrix \hat{\Sigma}, rank r, number of iterations T, corruption subset \mathcal{G} \subseteq [p] \times [p].

2: Set N^{(0)} = \Gamma(N).

3: for t = 1, \dots, T do

4: Calculate SVD: N^{(t)} = \sum_i \lambda_i^{(t)} u_i^{(t)} (v_i^{(t)})^{\top}, where \lambda_1^{(t)} \ge \lambda_2^{(t)} \dots \ge 0.

5: Let \tilde{N}^{(t)} = \sum_{i=1}^r \lambda_i^{(t)} u_i^{(t)} (v_i^{(t)})^{\top}.

6: Update corrupted entries: N^{(t+1)} = G(\tilde{N}^{(t)}) + \Gamma(\hat{\Sigma}).

7: end for

8: Output: \hat{U} = U^{(T)} = [u_1^{(T)} \dots u_r^{(T)}].
```

Next we give a performance guarantee for Algorithm 2. For any  $H \in \mathbb{R}^{p \times p}$ , let G(H) be the matrix H with all entries but those in  $\mathcal{G}$  set to zero and  $\Gamma(H) = H - G(H)$ . Define

(6.9) 
$$\eta = \max_{H \in \mathbb{R}^{m \times m, \operatorname{rank}(H) \le 2r}} \frac{\|G(H)\|}{\|H\|},$$

which essentially measures the maximum perturbations due to the entries in  $\mathcal{G}$  on the singular subspace. We also assume that the set of corrupted entries  $\mathcal{G}$  is b-sparse in the sense that

$$\max_{i} |\left\{j: (i,j) \in \mathcal{G}\right\}| \vee \max_{j} |\left\{i: (i,j) \in \mathcal{G}\right\}| \leq b,$$

i.e., the number of corrupted entries in each row and each column is at most b. To overcome the "spiky" issue discussed in Remark 4, we again assume the incoherence condition (6.10). We have the following theoretical results for Algorithm 2.

THEOREM 7 (General robust  $\sin \Theta$  theorem). Assume  $\mathcal{G} \in [p] \times [p]$  is b-sparse. Suppose one observes the symmetric matrix N = M + Z, where  $\operatorname{rank}(M) = r$  and Z is any symmetric perturbation, the eigenvectors of M are  $U \in \mathbb{O}_{p,r}$ . If  $\hat{U}$  is the output of Algorithm 1 with  $T = \Omega(\log \frac{\lambda_r(M)}{\eta \| \Gamma(Z) \|} \vee 1)$  iterations, there exists a universal constant c > 0 such that if the incoherence condition

(6.10) 
$$\frac{I(U)\|M\|}{\lambda_r(M)} \le \frac{cm}{\eta b r(b \wedge r)}$$

is satisfied and  $\eta \|\Gamma(Z)\| \leq c\lambda_r(M)$ , then the outcome of Algorithm 1 with corrupted index set  $\mathcal{G}$  satisfies

(6.11) 
$$\left\| \sin \Theta(\hat{U}, U) \right\| \lesssim \frac{\|\Gamma(Z)\|}{\lambda_r(M)} \wedge 1.$$

REMARK 12. Although the exact characterization of  $\eta$  may be difficult for general  $\mathcal{G}$ , Lemma 5 in the supplement shows that  $\eta \leq \sqrt{b \wedge (2r)}$  for all b-sparse  $\mathcal{G}$ .

Next, we prove Theorem 7, as the proof of Theorem 3 would follow by letting  $\mathcal{G}$  be the diagonal index set.

PROOFS OF THEOREM 7. To characterize how the proposed procedure refines the estimation by initialization and iterations, we define  $T_0 = \|\Gamma(N - M)\| = \|\Gamma(Z)\|$  and  $K_t = \|N^{(t)} - M\|$  for t = 0, 1, ... Since  $H = \Gamma(H) + G(H)$ , we have  $\|H\| \le \|G(H)\| + \|\Gamma(H)\|$  for all matrix  $H \in \mathbb{R}^{p \times p}$ .

Step 1. We first analyze the initial error satisfies

$$K_{0} = ||N^{(0)} - M|| = ||\Gamma(N) - M|| = ||\Gamma(N - M) - G(M)||$$

$$\leq ||\Gamma(N - M)|| + ||G(M)|| = ||\Gamma(Z)|| + ||G(M)|| = ||\Gamma(Z)|| + ||G(P_{U}MP_{U})||$$

$$\stackrel{\text{Lemma 1}}{\leq} ||\Gamma(Z)|| + \frac{I(U)rb}{p}||M|| = T_{0} + \frac{I(U)rb}{p}||M||.$$

Provided that  $\frac{I(U)rb}{p}||M|| \leq \lambda_r(M)/(16\eta)$  in the assumption, we have

(6.12) 
$$K_0 \le T_0 + \lambda_r(M)/(16\eta).$$

Step 2. Next, we analyze the evolution of iterations. By definitions,

(6.13) 
$$\tilde{N}^{(t-1)} = P_{U^{(t-1)}} N^{(t-1)}, \quad \Gamma(N^{(t)}) = \Gamma(N^{(t-1)}) = \dots = \Gamma(N),$$
$$G(N^{(t)}) = G(\tilde{N}^{(t-1)}).$$

Then for all  $t \geq 0$ ,

(6.14) 
$$\|\Gamma\left(N^{(t)} - M\right)\| = \|\Gamma(N - M)\| = \|\Gamma(Z)\| = T_0.$$

The analysis for  $||G(N^{(t)} - M)||$  is more complicated. Recall  $U^{(t-1)}$  is the leading r principal components of  $N^{(t-1)}$ . Then,

$$\begin{split} &\left\|G(N^{(t)}-M)\right\| \stackrel{(6.13)}{=} \left\|G(\tilde{N}^{(t-1)}-M)\right\| = \left\|G(P_{U^{(t-1)}}N^{(t-1)}-M)\right\| \\ &= \left\|G(P_{U^{(t-1)}}N^{(t-1)}-P_{U^{(t-1)}}M-P_{U_{\perp}^{(t-1)}}M)\right\| \\ &\leq \left\|G(P_{U^{(t-1)}}(N^{(t-1)}-M))\right\| + \left\|G(P_{U_{\perp}^{(t-1)}}M)\right\| \\ &\leq \left\|G(P_{U}(N^{(t-1)}-M))\right\| + \left\|G\left((P_{U^{(t-1)}}-P_{U})\cdot(N^{(t-1)}-M)\right)\right\| \\ &+ \left\|G\left(P_{U_{\perp}^{(t-1)}}M\right)\right\| \,. \end{split}$$

We bound these three terms separately:

- By Lemma 1,

(6.16) 
$$\left\| G\left(P_{U}(N^{(t-1)} - M)\right) \right\| \leq \sqrt{\frac{I(U)rb(b \wedge r)}{p}} \left\| N^{(t-1)} - M \right\|$$

$$= \sqrt{\frac{I(U)rb(b \wedge r)}{p}} K_{t-1}.$$

- Note that  $U^{(t-1)}(U^{(t-1)})^{\top}$  and  $UU^{\top}$  are both positive semi-definite and  $\|U^{(t-1)}(U^{(t-1)})^{\top}\|\vee\|UU^{\top}\| \leq 1$ , we have  $\|U^{(t-1)}(U^{(t-1)})^{\top}-UU^{\top}\| \leq 1$ . By Lemma 1 in Cai and Zhang (2018b),

$$\begin{aligned} \left\| U^{(t-1)}(U^{(t-1)})^{\top} - UU^{\top} \right\| &\leq 2\| \sin \Theta(U^{(t-1)}, U)\| \wedge 1 = 2\| (U_{\perp}^{(t-1)})^{\top} U\| \wedge 1 \\ &\leq \left( 2 \left\| (U_{\perp}^{(t-1)})^{\top} UU^{\top} M \right\| \cdot \lambda_{\min}^{-1} (U^{\top} M) \right) \wedge 1 \\ &\leq \left( 2 \left\| (U_{\perp}^{(t-1)})^{\top} M \right\| \cdot \lambda_{r}^{-1} (M) \right) \wedge 1 \\ &\leq \left( \frac{4\| N^{(t-1)} - M \|}{\lambda_{r}(M)} \right) \wedge 1 = \frac{4K_{t-1}}{\lambda_{r}(M)} \wedge 1, \end{aligned}$$

where the penultimate step follows from Lemma 7. Note that

$$rank((P_{U^{(t-1)}} - P_U)(N^{(t-1)} - M)) \le rank(P_{U^{(t-1)}} - P_U)$$
  
 
$$\le rank(P_{U^{(t-1)}}) + rank(P_U) \le 2r,$$

we have

(6.17) 
$$\left\| G \left( (P_{U^{(t-1)}} - P_U) \cdot (N^{(t-1)} - M) \right) \right\|$$

$$\leq \eta \cdot \|P_{U^{(t-1)}} - P_U\| \cdot \|N^{(t-1)} - M\| \leq \eta K_{t-1} \cdot \left( \frac{4K_{t-1}}{\lambda_r(M)} \wedge 1 \right).$$

- By Lemmas 1 and 7,

(6.18)

$$\|G(P_{U_{\perp}^{(t-1)}}M)\| = \|G(P_{U_{\perp}^{(t-1)}}MP_{U})\| \le \sqrt{\frac{I(U)rb(b \wedge r)}{p}} \|P_{U_{\perp}^{(t-1)}}M\|$$

$$\le 2\sqrt{\frac{I(U)rb(b \wedge r)}{p}} \|N^{(t-1)} - M\| = 2\sqrt{\frac{I(U)rb(b \wedge r)}{p}} K_{t-1}.$$

Combining (6.14)–(6.18), we have for all  $t \ge 1$ ,

(6.19) 
$$K_{t} \leq \|\Gamma(N^{(t)} - M)\| + \|G(N^{(t)} - M)\|$$

$$\leq T_{0} + 3\sqrt{\frac{I(U)rb(b \wedge r)}{p}}K_{t-1} + \frac{4\eta}{\lambda_{r}(M)}K_{t-1}^{2},$$

Step 3. Finally, we use induction to show that for all  $t \geq 0$ ,

(6.20) 
$$K_t \le 2T_0 + \frac{\lambda_r(M)}{\eta} 2^{-(t+4)}.$$

The base case of t = 0 is proved by (6.12). Next, suppose the statement (6.20) holds for t - 1. Then

$$K_{t} \overset{\text{(a)}}{\leq} T_{0} + 3\sqrt{\frac{I(U)rb(b \wedge r)}{p}} K_{t-1} + \frac{4\eta}{\lambda_{r}(M)} K_{t-1}^{2}$$

$$\overset{\text{(b)}}{\leq} T_{0} + \frac{K_{t-1}}{4} + K_{t-1} \left(\frac{8\eta T_{0}}{\lambda_{r}(M)} + \frac{1}{4}\right)$$

$$\overset{\text{(c)}}{\leq} T_{0} + \frac{K_{t-1}}{2} \overset{\text{(d)}}{\leq} T_{0} + \frac{2T_{0} + \lambda_{r}(M) \cdot (1/2)^{(t-1)+4}/\eta}{2}$$

$$= 2T_{0} + \lambda_{r}(M) \cdot (1/2)^{t+4}/\eta,$$

where (a) is (6.19); (b) is due to the assumption  $144I(U)rb(b \wedge r) \leq p$  and the induction hypothesis; (c) follows from the assumption  $T_0 \leq \lambda_r(M)/(64\eta)$ ; (d) is again by the induction hypothesis. Therefore, for all  $t \geq \Omega(\log \frac{\lambda_r(M)}{T_0\eta} \vee 1) = \Omega(\log \frac{\lambda_r(M)}{\eta \|\Gamma(Z)\|} \vee 1)$ , we have  $K_t \leq 3T_0$ . Finally, the desired (6.11) follows from Davis-Kahan's  $\sin \Theta$  theorem, completing the proof of Theorem 7.

REMARK 13. In fact, Theorem 7 implies Theorem 3. To see this, note that if the corruption set  $\mathcal{G}$  is the diagonal, i.e.,  $\mathcal{G} = \{(i,i) : 1 \leq i \leq p\}$ , we have

$$b = \max_{i} \{j : (i, j) \in \mathcal{G}\} \vee \max_{j} \{i : (i, j) \in \mathcal{G}\} = 1,$$
$$\eta = \max_{M} \|D(M)\| / \|M\| = \max_{M} \max_{i} \frac{|M_{ii}|}{\|M\|} = 1.$$

The next Lemma 1 provides an important technical tool for the proof of robust  $\sin \Theta$  theorem (Theorem 3). It essentially shows that the operator norm of the composition of linear maps  $G(P_U \cdot)$  is much smaller than the product of individual operator norms  $||G(\cdot)||$  and  $||P_U||$ , provided that the basis U is incoherent; the same conclusion also applies to  $G(\cdot P_V)$ .

LEMMA 1. Assume  $\mathcal{G} \subseteq [m_1] \times [m_2]$  is b-sparse, i.e.,  $\max_j \{i : (i,j) \in \mathcal{G}\} \vee \max_i \{j : (i,j) \in \mathcal{G}\} \leq b$ . Suppose  $U \in \mathbb{O}_{m_1,r}$  and  $V \in \mathbb{O}_{m_2,r}$ . Recall that G(A) is the matrix A with all entries in  $\mathcal{G}^c$  set to zero,  $I(U) = \frac{m_1}{r} \max_i \|e_i^\top U\|_2^2$ ,  $I(V) = \frac{m_2}{r} \max_i \|e_i^\top V\|_2^2$ ,  $P_U = UU^\top$ , and  $P_V = VV^\top$ . Then for any matrix  $A \in \mathbb{R}^{p_1 \times p_2}$ , we have

$$||G(P_U A)|| \le \sqrt{\frac{I(U)rb(b \wedge r)}{m_1}} ||A||, \quad ||G(AP_V)|| \le \sqrt{\frac{I(V)rb(b \wedge r)}{m_2}} ||A||,$$

$$and \quad ||G(P_U A P_V)|| \le \frac{\sqrt{I(U)I(V)} \cdot rb}{\sqrt{m_1 m_2}} ||A||.$$

In particular, recall that D(A) is the matrix A with all off-diagonal entries set to zero. Suppose  $U \in \mathbb{O}_{m,r}$ . Then for any matrix  $A \in \mathbb{R}^{m \times m}$ ,

$$||D(P_U(D(A)))|| \le \frac{I(U)r}{m}||D(A)||, \quad ||D(P_UA)|| \le \sqrt{\frac{I(U)r}{m}}||A||.$$

Proof of Lemma 1.

$$||G(P_{U}A)|| = \max_{||v||_{2}=1} ||v^{\top}G(UU^{\top}A)||_{2} = \max_{||v||_{2}=1} \left(\sum_{j=1}^{m_{2}} \left(v^{\top}[G(UU^{\top}A)]_{\cdot j}\right)^{2}\right)^{1/2}$$

$$= \max_{||v||_{2}=1} \left(\sum_{j=1}^{m_{2}} \left(\sum_{i:(i,j)\in\mathcal{G}} v_{i}(UU^{\top}A)_{i,j}\right)^{2}\right)^{1/2}$$

$$\leq \max_{||v||_{2}=1} \left(\sum_{j=1}^{m_{2}} \left(\sum_{i:(i,j)\in\mathcal{G}} v_{i}^{2}\right) \left(\sum_{i:(i,j)\in\mathcal{G}} (UU^{\top}A)_{i,j}^{2}\right)\right)^{1/2},$$

where the inequality is due to Cauchy-Schwarz. In particular, for any  $1 \le j \le m_2$ ,

$$\sum_{i:(i,j)\in\mathcal{G}} (UU^{\top}A)_{ij}^{2} \leq \sum_{i:(i,j)\in\mathcal{G}} \left( U_{i\cdot} \cdot (U^{\top}A)_{\cdot j} \right)^{2} = \sum_{i:(i,j)\in\mathcal{G}} \|U_{i\cdot}\|_{2}^{2} \cdot \left\| (U^{\top}A)_{\cdot j} \right\|_{2}^{2} \\
\leq \sum_{i:(i,j)\in\mathcal{G}} \frac{I(U)r}{m_{1}} \|A\|^{2} \leq \frac{I(U)rb}{m_{1}} \|A\|^{2}.$$

Thus,

$$||G(P_{U}A)|| \leq \sqrt{\frac{I(U)rb}{m_{1}}} ||A|| \max_{||v||_{2}=1} \left( \sum_{j=1}^{m_{2}} \left( \sum_{i:(i,j)\in\mathcal{G}} v_{i}^{2} \right) \right)^{1/2}$$

$$= \sqrt{\frac{I(U)rb}{m_{1}}} ||A|| \cdot \max_{||v||_{2}=1} \left( \sum_{i=1}^{m_{1}} \sum_{j:(i,j)\in\mathcal{G}} v_{i}^{2} \right)^{1/2}$$

$$\leq \sqrt{\frac{I(U)rb}{m_{1}}} ||A|| \cdot \max_{||v||_{2}=1} \left( \sum_{i=1}^{m_{1}} v_{i}^{2}b \right)^{1/2} \leq \sqrt{\frac{I(U)rb^{2}}{m_{1}}} ||A||.$$

Additionally, since  $rank(A) \leq r$  and  $\mathcal{G}$  is b-sparse,

$$||G(P_{U}A)||^{2} \leq ||G(P_{U}A)||_{F}^{2}$$

$$= \sum_{i,j} \left(G(UU^{\top}A)\right)_{ij}^{2} = \sum_{j=1}^{m_{2}} \sum_{i:(i,j)\in\mathcal{G}} \left(U_{i\cdot}(U^{\top}A)_{\cdot j}\right)^{2}$$

$$\leq \sum_{j=1}^{m_{2}} \sum_{i:(i,j)\in\mathcal{G}} ||U_{i\cdot}||_{2}^{2} \cdot ||U^{\top}A_{\cdot j}||_{2}^{2} \leq \sum_{j=1}^{m_{2}} \frac{I(U)rb}{m_{1}} \cdot ||U^{\top}A_{\cdot j}||_{2}^{2}$$

$$= \frac{I(U)rb}{m_{1}} \cdot ||U^{\top}A||_{F}^{2} \leq \frac{I(U)rb}{m_{1}} \cdot r||U^{\top}A||^{2} \leq \frac{I(U)r^{2}b}{m_{1}}||A||^{2}.$$

Combining previous two inequalities, we have

$$||G(P_U A)|| \le \sqrt{\frac{I(U)rb(r \wedge b)}{m_1}} ||A||.$$

The proof for  $||G(AP_V)|| \leq \sqrt{I(V)rb(b \wedge r)/m_2}||A||$  similarly follows. Next, for any  $u \in \mathbb{R}^{m_1}$ ,  $v \in \mathbb{R}^{m_2}$  such that  $||u||_2 = ||v||_2 = 1$ , we have

$$u^{\top}G(P_{U}AP_{V})v = u^{\top}G(UU^{\top}AVV^{\top})v = \sum_{(i,j)\in\mathcal{G}} u_{i}v_{j} \left[ UU^{\top}AVV^{\top} \right]_{ij}$$

$$\leq \sum_{(i,j)\in\mathcal{G}} |u_{i}v_{j}| ||U_{i}.||_{2} \cdot ||U^{\top}AV|| \cdot ||V_{j}.||_{2}$$

$$\leq \sum_{(i,j)\in\mathcal{G}} |u_{i}v_{j}| \cdot \sqrt{\frac{rI(U)}{m_{1}}} \cdot ||A|| \cdot \sqrt{\frac{rI(V)}{m_{2}}}$$

$$\leq \frac{\sqrt{I(U)I(V)}r}{\sqrt{m_{1}m_{2}}} \cdot ||A|| \cdot \sum_{(i,j)\in\mathcal{G}} \frac{u_{i}^{2} + v_{j}^{2}}{2}$$

$$\leq \frac{\sqrt{I(U)I(V)}r}{\sqrt{m_{1}m_{2}}} \cdot ||A|| \cdot \left(\sum_{i} \sum_{j:(i,j)\in\mathcal{G}} \frac{u_{i}^{2}}{2} + \sum_{j} \sum_{i:(i,j)\in\mathcal{G}} \frac{v_{j}^{2}}{2}\right)$$

$$\leq \frac{br\sqrt{I(U)I(V)}}{\sqrt{m_{1}m_{2}}} ||A||,$$

which means

$$||G(P_UAP_V)|| \le \frac{br\sqrt{I(U)I(V)}}{\sqrt{m_1m_2}}||A||.$$

For the diagonal operator  $D(\cdot)$ , since D(A) is a diagonal matrix, we have  $D(A)e_i = D(A)_{ii}e_i$  and

$$||D(P_{U}(D(A)))|| = \max_{i} |\{P_{U}(D(A))\}_{ii}| = \max_{i} \left| e_{i}^{\top} P_{U} D(A) e_{i} \right|$$

$$= \max_{i} \left| e_{i}^{\top} P_{U} e_{i} \cdot A_{ii} \right| = \max_{i} ||U^{\top} e_{i}||_{2}^{2} \cdot |A_{ii}| \leq \frac{I(U)r}{m} ||D(A)||,$$

$$||D(P_{U}(A))|| = \max_{i} |(P_{U} A)_{ii}| = \max_{i} \left| e_{i}^{\top} U U^{\top} A e_{i} \right|$$

$$\leq ||e_{i}^{\top} U||_{2} \cdot ||A|| \leq \sqrt{\frac{I(U)r}{m}} ||A||.$$

#### References.

- Aflalo, Y. and Kimmel, R. (2013). Spectral multidimensional scaling. Proceedings of the National Academy of Sciences, page 201308708.
- Alzahrani, T. and Horadam, K. (2016). Community detection in bipartite networks: Algorithms and case studies. In *Complex Systems and Networks*, pages 25–50. Springer.
- Bai, J. and Li, K. (2012). Statistical analysis of factor models of high dimension. The Annals of Statistics, 40(1):436–465.
- Bai, Z. and Yao, J. (2012). On sample eigenvalues in a generalized spiked population model. *Journal of Multivariate Analysis*, 106:167–177.
- Baik, J., Arous, G. B., Péché, S., et al. (2005). Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. The Annals of Probability, 33(5):1643– 1697.
- Baik, J. and Silverstein, J. W. (2006). Eigenvalues of large sample covariance matrices of spiked population models. *Journal of multivariate analysis*, 97(6):1382–1408.
- Bartlett, M. S. (1937). The statistical conception of mental factors. British journal of Psychology, 28(1):97.
- Boucheron, S., Lugosi, G., and Massart, P. (2013). Concentration inequalities: A nonasymptotic theory of independence. Oxford university press.
- Cai, J.-F. and Wei, K. (2018). Exploiting the structure effectively and efficiently in low-rank matrix recovery. Processing, Analyzing and Learning of Images, Shapes, and Forms, 19:21.
- Cai, T. T., Li, X., and Ma, Z. (2016). Optimal rates of convergence for noisy sparse phase retrieval via thresholded wirtinger flow. *The Annals of Statistics*, 44(5):2221–2251.
- Cai, T. T. and Zhang, A. (2016). Minimax rate-optimal estimation of high-dimensional covariance matrices with incomplete data. *Journal of multivariate analysis*, 150:55–74.
- Cai, T. T. and Zhang, A. (2018a). Heteroskedastic wishart-type concentration inequalities. preprint.
- Cai, T. T. and Zhang, A. (2018b). Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics. The Annals of Statistics, 46(1):60-89.
- Candès, E. J., Li, X., Ma, Y., and Wright, J. (2011). Robust principal component analysis? Journal of the ACM (JACM), 58(3):11.
- Candes, E. J., Li, X., and Soltanolkotabi, M. (2015). Phase retrieval via wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007.
- Candès, E. J. and Recht, B. (2009). Exact matrix completion via convex optimization. Foundations of Computational mathematics, 9(6):717.
- Candes, E. J., Sing-Long, C. A., and Trzasko, J. D. (2013). Unbiased risk estimates for singular value thresholding and spectral estimators. *IEEE transactions on signal* processing, 61(19):4643–4657.
- Candès, E. J. and Tao, T. (2010). The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080.
- Cao, Y., Zhang, A., and Li, H. (2017). Multi-sample estimation of bacterial composition matrix in metagenomics data. arXiv preprint arXiv:1706.02380.
- Chatterjee, S. (2015). Matrix estimation by universal singular value thresholding. The Annals of Statistics, 43(1):177-214.
- Chen, Y. and Suh, C. (2015). Spectral mle: Top-k rank aggregation from pairwise comparisons. In *International Conference on Machine Learning*, pages 371–380.
- Cochran, R. N. and Horne, F. H. (1977). Statistically weighted principal component analysis of rapid scanning wavelength kinetics experiments. *Analytical Chemistry*, 49(6):846–

853.

- Collins, M., Dasgupta, S., and Schapire, R. E. (2002). A generalization of principal components analysis to the exponential family. In Advances in neural information processing systems, pages 617–624.
- Davis, C. and Kahan, W. M. (1970). The rotation of eigenvectors by a perturbation. iii. SIAM Journal on Numerical Analysis, 7(1):1–46.
- De Lathauwer, L., De Moor, B., and Vandewalle, J. (2000). On the best rank-1 and rank-(r 1, r 2,..., rn) approximation of higher-order tensors. SIAM journal on Matrix Analysis and Applications, 21(4):1324–1342.
- Dobriban, E., Leeb, W., and Singer, A. (2016). Pca from noisy, linearly reduced data: the diagonal case. arXiv preprint arXiv:1611.10333.
- Donath, W. E. and Hoffman, A. J. (2003). Lower bounds for the partitioning of graphs. In Selected Papers Of Alan J Hoffman: With Commentary, pages 437–442. World Scientific.
- Donoho, D. and Gavish, M. (2014). Minimax risk of matrix denoising by singular value thresholding. *The Annals of Statistics*, 42(6):2413–2440.
- Donoho, D. L., Gavish, M., and Johnstone, I. M. (2018). Optimal shrinkage of eigenvalues in the spiked covariance model. Annals of statistics, 46(4):1742.
- Florescu, L. and Perkins, W. (2016). Spectral thresholds in the bipartite stochastic block model. In *Conference on Learning Theory*, pages 943–959.
- Fortunato, S. (2010). Community detection in graphs. Physics reports, 486(3-5):75-174.
- Gavish, M. and Donoho, D. L. (2017). Optimal shrinkage of singular values. *IEEE Transactions on Information Theory*, 63(4):2137–2152.
- Ghosh, J. and Dunson, D. B. (2009). Default prior distributions and efficient posterior computation in bayesian factor analysis. *Journal of Computational and Graphical Statistics*, 18(2):306–320.
- Golub, G. H., Hoffman, A., and Stewart, G. W. (1987). A generalization of the eckartyoung-mirsky matrix approximation theorem. *Linear Algebra and its applications*, 88:317–327.
- Hao, B., Zhang, A., and Cheng, G. (2018). Sparse and low-rank tensor estimation via cubic sketchings. arXiv preprint arXiv:1801.09326.
- Hong, D., Balzano, L., and Fessler, J. A. (2016). Towards a theoretical analysis of pca for heteroscedastic data. In Communication, Control, and Computing (Allerton), 2016 54th Annual Allerton Conference on, pages 496–503. IEEE.
- Hong, D., Balzano, L., and Fessler, J. A. (2018a). Asymptotic performance of pca for high-dimensional heteroscedastic data. *Journal of Multivariate Analysis*.
- Hong, D., Fessler, J. A., and Balzano, L. (2018b). Optimally weighted pca for high-dimensional heteroscedastic data. arXiv preprint arXiv:1810.12862.
- Hotelling, H. (1936). Relations between two sets of variates. Biometrika, 28(3/4):321–377.
  Jain, P., Netrapalli, P., and Sanghavi, S. (2013). Low-rank matrix completion using alternating minimization. In Proceedings of the forty-fifth annual ACM symposium on Theory of computing, pages 665–674. ACM.
- Jin, J., Wang, W., et al. (2016). Influential features pca for high dimensional clustering. The Annals of Statistics, 44(6):2323–2359.
- Johnstone, I. M. (2001). On the distribution of the largest eigenvalue in principal components analysis. Annals of statistics, pages 295–327.
- Katznelson, Y. (2004). An introduction to harmonic analysis. Cambridge University Press. Ke, Z. T. and Wang, M. (2017). A new svd approach to optimal topic estimation. arXiv preprint arXiv:1704.07016.
- Keshavan, R. H. (2012). Efficient algorithms for collaborative filtering. PhD thesis, Stan-

- ford University.
- Keshavan, R. H., Montanari, A., and Oh, S. (2010a). Matrix completion from a few entries. IEEE Transactions on Information Theory, 56(6):2980–2998.
- Keshavan, R. H., Montanari, A., and Oh, S. (2010b). Matrix completion from noisy entries. Journal of Machine Learning Research, 11(Jul):2057–2078.
- Koltchinskii, V., Lounici, K., and Tsybakov, A. B. (2011). Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. The Annals of Statistics, 39(5):2302–2329.
- Laurent, B. and Massart, P. (2000). Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338.
- Lawley, D. N. and Maxwell, A. E. (1962). Factor analysis as a statistical method. *Journal of the Royal Statistical Society. Series D (The Statistician)*, 12(3):209–229.
- Liu, L. T., Dobriban, E., and Singer, A. (2016). e pca: High dimensional exponential family pca. arXiv preprint arXiv:1611.05550.
- Lounici, K. et al. (2014). High-dimensional covariance matrix estimation with missing observations. *Bernoulli*, 20(3):1029–1058.
- Martin, A. D., Quinn, K. M., Park, J. H., et al. (2011). Mcmcpack: Markov chain monte carlo in r. *Journal of Statistical Software*, 42(i09).
- Massart, P. (2007). Concentration inequalities and model selection.
- Mazumder, R., Hastie, T., and Tibshirani, R. (2010). Spectral regularization algorithms for learning large incomplete matrices. *Journal of machine learning research*, 11(Aug):2287–2322.
- Melamed, D. (2014). Community structures in bipartite networks: A dual-projection approach. *PloS one*, 9(5):e97823.
- Mohamed, S., Ghahramani, Z., and Heller, K. A. (2009). Bayesian exponential family pca. In *Advances in neural information processing systems*, pages 1089–1096.
- Nadakuditi, R. R. (2014). Optshrink: An algorithm for improved low-rank signal matrix denoising by optimal, data-driven singular value shrinkage. *IEEE Transactions on Information Theory*, 60(5):3002–3018.
- Negahban, S., Oh, S., and Shah, D. (2012). Rank centrality: Ranking from pair-wise comparisons. arXiv preprint arXiv:1209.1688.
- Newman, M. E. (2013). Spectral methods for community detection and graph partitioning. *Physical Review E*, 88(4):042822.
- Owen, A. B. and Wang, J. (2016). Bi-cross-validation for factor analysis. Statistical Science, 31(1):119–139.
- Paul, D. (2007). Asymptotics of sample eigenstructure for a large dimensional spiked covariance model. *Statistica Sinica*, pages 1617–1642.
- Recht, B. (2011). A simpler approach to matrix completion. *Journal of Machine Learning Research*, 12(Dec):3413–3430.
- Richard, E. and Montanari, A. (2014). A statistical model for tensor pca. In *Advances in Neural Information Processing Systems*, pages 2897–2905.
- Robin, G., Josse, J., Moulines, É., and Sardy, S. (2019). Low-rank model with covariates for count data with missing values. *Journal of Multivariate Analysis*, 173:416–434.
- Salmon, J., Harmany, Z., Deledalle, C.-A., and Willett, R. (2014). Poisson noise reduction with non-local pca. *Journal of mathematical imaging and vision*, 48(2):279–294.
- Shabalin, A. A. and Nobel, A. B. (2013). Reconstruction of a low-rank matrix in the presence of gaussian noise. *Journal of Multivariate Analysis*, 118:67–76.
- Sun, R. and Luo, Z.-Q. (2016). Guaranteed matrix completion via non-convex factorization. *IEEE Transactions on Information Theory*, 62(11):6535–6579.
- Thomson, G. (1939). The factorial analysis of human ability. British Journal of Educa-

- tional Psychology, 9(2):188-195.
- Tipping, M. E. and Bishop, C. M. (1999). Probabilistic principal component analysis. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 61(3):611–622.
- Toh, K.-C. and Yun, S. (2010). An accelerated proximal gradient algorithm for nuclear norm regularized linear least squares problems. *Pacific Journal of optimization*, 6(615-640):15.
- Vaswani, N. and Narayanamurthy, P. (2017). Finite sample guarantees for pca in non-isotropic and data-dependent noise. In 2017 55th Annual Allerton Conference on Communication, Control, and Computing (Allerton), pages 783–789. IEEE.
- Vaswani, N. and Narayanamurthy, P. (2018). Pca in sparse data-dependent noise. In 2018 IEEE International Symposium on Information Theory (ISIT), pages 641–645. IEEE.
- Vershynin, R. (2010). Introduction to the non-asymptotic analysis of random matrices. arXiv preprint arXiv:1011.3027.
- Vershynin, R. (2011). Spectral norm of products of random and deterministic matrices. *Probability theory and related fields*, 150(3-4):471–509.
- Wang, W. and Fan, J. (2017). Asymptotics of empirical eigenstructure for high dimensional spiked covariance. Annals of statistics, 45(3):1342.
- Wedin, P.-A. (1972). Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics*, 12(1):99–111.
- Yao, J., Zheng, S., and Bai, Z. (2015). Sample covariance matrices and high-dimensional data analysis. Cambridge University Press.
- Yu, B. (1997). Assouad, fano, and le cam. In Festschrift for Lucien Le Cam, pages 423–435. Springer.
- Yu, Y., Wang, T., and Samworth, R. J. (2014). A useful variant of the davis–kahan theorem for statisticians. *Biometrika*, 102(2):315–323.
- Zhang, A., Cai, T. T., and Wu, Y. (2018). Supplement to "Heteroskedastic PCA: Algorithm, optimality, and applications". Technical Report.
- Zhang, A. and Han, R. (2018). Optimal sparse singular value decomposition for high-dimensional high-order data. arXiv preprint arXiv:1809.01796.
- Zhang, A. and Wang, M. (2018). Spectral state compression of markov processes. arXiv preprint arXiv:1802.02920.
- Zhang, A. and Xia, D. (2018). Tensor svd: Statistical and computational limits. *IEEE Transactions on Information Theory*, to appear.
- Zhou, Z. and Amini, A. A. (2018). Optimal bipartite network clustering.  $arXiv\ preprint\ arXiv:1803.06031.$

# Supplement to "Heteroskedastic PCA: Algorithm,

# Optimality, and Applications"

Anru Zhang, T. Tony Cai, and Yihong Wu

### APPENDIX A: ADDITIONAL PROOFS

#### A.1. Additional Proofs for Heteroskedastic PCA.

PROOF OF THEOREM 2. We only need to show the following two inequalities to prove this theorem, (A.1)

$$\inf_{\hat{U}} \sup_{\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{\text{sum}},\sigma_{\text{max}},\nu)} \mathbb{E} \left\| \sin \Theta(\hat{U},U) \right\| \gtrsim \left( \frac{\sigma_{\text{sum}}}{(n\nu)^{1/2}} + \frac{\sigma_{\text{sum}}\sigma_{\text{max}}}{n^{1/2}\nu} \right) \wedge 1,$$

(A.2) 
$$\inf_{\hat{U}} \sup_{\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{\text{sum}},\sigma_{\text{max}},\nu)} \mathbb{E} \left\| \sin \Theta(\hat{U},U) \right\| \gtrsim \frac{\sqrt{r}\sigma_{\text{max}}}{(n\nu)^{1/2}} \wedge 1.$$

We first consider (A.1). Since all parameters can be rescaled, we assume  $\nu = 1$  without loss of generality. The proof is divided into three steps.

Step 1 In this step, we construct a series of "candidate covariance matrices" and prove that they belong to the subset of covariance matrices in the theorem statement. Let

(A.3) 
$$d = \lfloor \sigma_{\text{sum}}^2 / (8\sigma_{\text{max}}^2) \rfloor \vee 6, \quad L = 2\lceil 1/(dc_I) \rceil.$$

Now, we impose the assumption that

(A.4) 
$$p \ge 50 \vee \{2(r-1)(1+c_I)/c_I\} \vee \{8/c_I\}.$$

Since  $\sigma_{\text{sum}} \leq \sqrt{p}\sigma_{\text{max}}$ , we must have

(A.5) 
$$Ld \stackrel{\text{(A.3)}}{=} 2d \left\lceil \frac{1}{dc_I} \right\rceil < 2d \left( \frac{1}{dc_I} + 1 \right) = \frac{2}{c_I} + 2 \left( \left\lfloor \frac{\sigma_{\text{sum}}^2}{8\sigma_{\text{max}}^2} \right\rfloor \vee 6 \right)$$

$$\stackrel{\text{(A.4)}}{\leq} \frac{p}{4} + \frac{\sigma_{\text{sum}}^2}{4\sigma_{\text{max}}^2} \vee 12 \stackrel{\text{(A.4)}}{\leq} \frac{p}{4} + \max \left\{ \frac{p}{4}, \frac{p}{4} \right\} = \frac{p}{2}.$$

By Lemma 9, we can construct  $Q \in \mathbb{O}_{(p-Ld),(r-1)}$  with small incoherence constant:

(A.6) 
$$\max_{i} \|e_{i}^{\top} Q\|_{2}^{2} \leq \frac{1}{\lfloor \frac{p-Ld}{r-1} \rfloor} \leq \frac{1}{\frac{p-Ld}{r-1} - 1}$$

$$\leq \frac{1}{\frac{p/2}{r-1} - 1} \stackrel{\text{(A.4)}}{\leq} \frac{r - 1}{(r-1)(1 + c_{I})/c_{I} - (r-1)} \leq c_{I}.$$

By the Varshamov-Gilbert bound (Massart, 2007, Lemma 4.7), we can find series of vectors  $v^{(1)}, \ldots, v^{(N)} \subseteq \{-1,1\}^d$  with  $N \ge \exp(d/8)$ , such that

(A.7) 
$$||v^{(l)} - v^{(k)}||_2^2 \ge d$$
, for all  $1 \le k \ne l \le N$ 

Next, we construct a series of candidate covariance matrices for  $k = 1, \ldots, N$ ,

$$U^{(k)} = \begin{bmatrix} u^{(k)} & 0_{(Ld)\times(r-1)} \\ 0_{(p-Ld)\times 1} & Q \end{bmatrix} \in \mathbb{R}^{p\times r},$$

$$u^{(k)} = \begin{bmatrix} \frac{1}{\sqrt{Ld(1+\theta^2)}} \left( 1_d + \theta v^{(k)} \right) \\ \vdots \\ \frac{1}{\sqrt{Ld(1+\theta^2)}} \left( 1_d + \theta v^{(k)} \right) \\ \frac{1}{\sqrt{Ld(1+\theta^2)}} \left( 1_d - \theta v^{(k)} \right) \\ \vdots \\ \frac{1}{\sqrt{Ld(1+\theta^2)}} \left( 1_d - \theta v^{(k)} \right) \end{bmatrix} \in \mathbb{R}^{Ld};$$

$$D_{ij} = \begin{cases} \sigma_0^2, & 1 \le i = j \le Ld; \\ 0, & \text{otherwise,} \end{cases} \qquad \sigma_0^2 = \sigma_{\text{max}}^2 \wedge \{\sigma_{\text{sum}}^2/(Ld)\},$$
$$\Sigma^{(k)} = U^{(k)}(U^{(k)})^\top + D.$$

Here,  $0 \le \theta \le 1$  is a constant to be specified later; both  $\frac{1}{\sqrt{Ld(1+\theta^2)}}(1+\theta v^{(k)})$  and  $\frac{1}{\sqrt{Ld(1+\theta^2)}}(1-\theta v^{(k)})$  are repeated for (L/2) times in the first column of  $U^{(k)}$ . Then, all columns of  $U^{(k)}$  are orthonormal and

$$\max_{1 \le i \le p} \|e_i^{\top} U^{(k)}\|_2^2 \le \max \left\{ \frac{(1+\theta)^2}{Ld(1+\theta^2)}, \max_i \|e_i^{\top} Q\|_2^2 \right\} \\
\le \max \left\{ \frac{2}{Ld}, \max_i \|e_i^{\top} Q\|_2^2 \right\} \stackrel{\text{(A.3)(A.6)}}{\le} c_I.$$

Then  $U^{(k)}$  satisfies the incoherence constraint of the class  $\mathcal{F}_{p,r}(\sigma_{\text{sum}}, \sigma_{\text{max}}, \nu)$ ,

$$I\left(U^{(k)}\right) = \frac{p}{r} \max_{i} \|e_i^{\top} U^{(k)}\|_2^2 \le c_I p/r.$$

In addition,

$$\max_{1 \le i \le p} D_{ii} = \sigma_{\max}^2 \wedge \{\sigma_{\text{sum}}^2 / (Ld)\} \le \sigma_{\max}^2,$$

$$\sum_{i=1}^p D_{ii} = Ld\left(\sigma_{\max}^2 \wedge \{\sigma_{\text{sum}}^2 / (Ld)\}\right) \le \sigma_{\text{sum}}^2,$$

$$\lambda_r \left(U^{(k)} (U^{(k)})^\top\right) = 1 = \nu.$$

Therefore,  $\Sigma^{(1)}, \ldots, \Sigma^{(N)}$  truly belongs to the class in the theorem statement:

(A.8) 
$$\Sigma^{(1)}, \dots, \Sigma^{(N)} \subseteq \mathcal{F}_{p,n,r}(\sigma_{\text{sum}}, \sigma_{\text{max}}, \nu).$$

Step 2 Next for any  $k \neq l$ , we prove that  $U^{(k)}, U^{(l)}$  are well-separated and the KL-divergence of  $X^{(k)}$  and  $U^{(l)}$  are bounded if  $X^{(k)} \sim N(0, \Sigma^{(k)}), X^{(l)} \sim N(0, \Sigma^{(l)})$ . Since  $\sigma_{\text{sum}} \geq \sigma_{\text{max}}$ , we have

$$\begin{split} &(\mathrm{A}.9) \\ &\sigma_0^2 = \!\! \sigma_{\mathrm{max}}^2 \wedge \frac{\sigma_{\mathrm{sum}}^2}{Ld} \stackrel{(\mathrm{A}.3)}{\geq} \sigma_{\mathrm{max}}^2 \wedge \frac{\sigma_{\mathrm{sum}}^2}{2d\lceil 1/(dc_I) \rceil} \geq \sigma_{\mathrm{max}}^2 \wedge \frac{\sigma_{\mathrm{sum}}^2}{2d\left(\frac{1}{dc_I} + 1\right)} \\ &\geq \!\! \sigma_{\mathrm{max}}^2 \wedge \frac{\sigma_{\mathrm{sum}}^2}{\frac{2}{c_I} + 2\left(\lfloor \sigma_{\mathrm{sum}}^2/(8\sigma_{\mathrm{max}}^2) \rfloor \wedge 6\right)} \geq \sigma_{\mathrm{max}}^2 \wedge \frac{\sigma_{\mathrm{sum}}^2}{\frac{2}{c_I} + 12} \geq c\sigma_{\mathrm{max}}^2; \\ &d\sigma_0^2 \geq \!\! cd\sigma_{\mathrm{max}}^2 = c\left(\lfloor \sigma_{\mathrm{sum}}^2/(8\sigma_{\mathrm{max}}^2) \rfloor \vee 6\right) \sigma_{\mathrm{max}}^2 \geq c\left(\sigma_{\mathrm{sum}}^2/(16\sigma_{\mathrm{max}}^2)\right) \sigma_{\mathrm{max}}^2 \geq c'\sigma_{\mathrm{sum}}^2 \end{split}$$

for some constants c, c' > 0 that only rely on  $c_I$ .

By the definition of (A.7), we have for any  $1 \le k \ne l \le N$ ,

$$\begin{aligned} & \left\| \sin \Theta \left( U^{(k)}, U^{(l)} \right) \right\| = \left( 1 - \lambda_r^2 \left( (U^{(k)})^\top U^{(l)} \right) \right)^{1/2} = \left( 1 - \left( u^{(k)\top} u^{(l)} \right)^2 \right)^{1/2} \\ &= \left( 1 - \frac{(L/2)^2}{L^2 d^2 (1 + \theta^2)^2} \left( (1_d + \theta v^{(k)})^\top (1_d + \theta v^{(l)}) + (1_d - \theta v^{(k)})^\top (1_d - \theta v^{(l)}) \right)^2 \right)^{1/2} \\ &= \left( 1 - \frac{1}{4d^2 (1 + \theta^2)^2} \left( 2d + 2\theta^2 v^{(k)\top} v^{(l)} \right)^2 \right)^{1/2} \\ &= \left( 1 - \left( \frac{1 + \theta^2 (v^{(k)})^\top v^{(l)} / d}{1 + \theta^2} \right)^2 \right)^{1/2} .\end{aligned}$$

By (A.7), for any  $k \neq l$ , we have  $d \leq ||v^{(k)} - v^{(l)}||_2^2 \leq 4d$  and

$$(v^{(k)})^{\top} v^{(l)} = \frac{1}{2} \left( \|v^{(k)}\|_{2}^{2} + \|v^{(l)}\|_{2}^{2} - \|v^{(k)} - v^{(l)}\|_{2}^{2} \right)$$
$$= \frac{1}{2} \left( 2d - \|v^{(k)} - v^{(l)}\|_{2}^{2} \right) \in [-d, d/2].$$

Consequently,

(A.10)

$$\left(1 - \left(\frac{1 + \theta^2/2}{1 + \theta^2}\right)^2\right)^{1/2} \le \left\|\sin\Theta(U^{(k)}, U^{(l)})\right\| \le \left(1 - \left(\frac{1 - \theta^2}{1 + \theta^2}\right)^2\right)^{1/2}.$$

Provided that  $0 < \theta \le 1$ ,

(A.11)

$$\left(1 - \left(\frac{1 - \theta^2}{1 + \theta^2}\right)^2\right)^{1/2} = \left(\frac{(1 + \theta^2)^2 - (1 - \theta^2)^2}{(1 + \theta^2)^2}\right)^{1/2} = \frac{2\theta}{1 + \theta^2} \le 2\theta,$$

(A.12) 
$$\left(1 - \left(\frac{1 + \theta^2/2}{1 + \theta^2}\right)^2\right)^{1/2} = \frac{\left(\theta^2 + (3/4)\theta^4\right)^{1/2}}{1 + \theta^2} \ge \frac{\theta}{2}.$$

Combining (A.10), (A.11), and (A.12), we have

(A.13) 
$$\frac{\theta}{2} \le \left\| \sin \Theta(U^{(k)}, U^{(l)}) \right\| \le 2\theta, \quad \forall 1 \le k \ne l \le N.$$

Suppose

$$X^{(k)} = \left[ X_1^{(k)} \dots X_n^{(k)} \right] \stackrel{iid}{\sim} N(0, \Sigma^{(k)}), \quad k = 1, \dots, N.$$

Next, we consider the Kullback-Leibler divergence between  $X^{(k)}$  and  $X^{(l)}$  for any  $1 \leq k \neq l \leq N$ . Note the following fact on the Kullback-Leibler divergence between multivariate Gaussians: suppose  $X = [X_1, \ldots, X_n] \stackrel{iid}{\sim} N(0, \Sigma)$  and  $X' = [X'_1, \ldots, X'_n] \stackrel{iid}{\sim} N(0, \Sigma')$  are p-dimensional vectors. If  $\Sigma$  and  $\Sigma'$  are non-degenerating, then

$$D_{KL}(X||X') = \frac{n}{2} \left( \operatorname{tr} \left( (\Sigma')^{-1} \Sigma \right) - p + \log \left( \frac{\det \Sigma'}{\det \Sigma} \right) \right).$$

Since  $\Sigma^{(k)}$  and  $\Sigma^{(l)}$  may be degenerating, one cannot directly apply the previous formula to calculate their KL divergence. Instead, denote the top

(Ld)-by-(Ld) sub-matrix of  $\Sigma^{(k)}$  as

$$\tilde{\Sigma}^{(k)} = u^{(k)}u^{(k)})^{\top} + \tilde{D} \in \mathbb{R}^{(Ld) \times (Ld)},$$
where 
$$u^{(k)} = \begin{bmatrix} \frac{1}{\sqrt{Ld(1+\theta^2)}} \left(1_d + \theta v^{(k)}\right) \\ \vdots \\ \frac{1}{\sqrt{Ld(1+\theta^2)}} \left(1_d + \theta v^{(k)}\right) \\ \frac{1}{\sqrt{Ld(1+\theta^2)}} \left(1_d - \theta v^{(k)}\right) \\ \vdots \\ \frac{1}{\sqrt{Ld(1+\theta^2)}} \left(1_d - \theta v^{(k)}\right) \end{bmatrix} \in \mathbb{R}^{Ld}, \quad \tilde{D} = \sigma_0^2 I.$$

By the structure of  $\Sigma^{(k)}$ , we know  $\det(\tilde{\Sigma}^{(k)}) = \det(\tilde{\Sigma}^{(l)})$  for all  $1 \leq k, l \leq N$ , and  $\Sigma^{(k)}_{[1:Ld,1:Ld]} = \tilde{\Sigma}^{(k)}$ ,  $\Sigma^{(k)}_{[(Ld+1):p,1:Ld]} = 0$ ,  $\Sigma^{(k)}_{[1:Ld,(Ld+1):p]} = 0$ ,  $\Sigma^{(k)}_{[(Ld+1):p,(Ld+1):p]} = QQ^{\top}$ . Here,  $\Sigma^{(k)}_{[1:Ld,1:Ld]}$  represents the submatrix formed by the first to Ld-th rows and first to Ld-th columns of  $\Sigma^{(k)}$ ;  $\Sigma^{(k)}_{[1:Ld,(Ld+1):p]}$  and  $\Sigma^{(k)}_{[(Ld+1):p,(Ld+1):p]}$  are defined in a similar fashion. Then, 1) for any  $1 \leq k \leq N$  and  $1 \leq i \leq n$ ,  $(X^{(k)}_i)_{[1:Ld]}$  and  $(X^{(k)}_i)_{[(Ld+1):p]}$ , i.e., the first Ld entries and the other entries of  $X_i$ , are two independent vectors; 2)  $(X^{(k)}_1)_{[(Ld+1):p]}, \ldots, (X^{(k)}_n)_{[(Ld+1):p]}$  are independent and identically distributed. Thus,

$$D_{KL}\left(X^{(l)}||X^{(k)}\right) = D_{KL}\left(X^{(l)}_{[1:Ld,:]}||X^{(k)}_{[1:Ld,:]}\right) = \frac{n}{2}\left(\operatorname{tr}\left((\tilde{\Sigma}^{(k)})^{-1}\tilde{\Sigma}^{(l)}\right) - Ld\right).$$

Here,  $X_{[1:LD,:]}^{(k)}$  and  $X_{[1:LD,:]}^{(l)}$  represent the first LD rows of  $X^{(k)}$  and  $X^{(l)}$ , respectively. Since  $u^{(k)}$  is a unit vector, one can verify that

$$(\tilde{\Sigma}^{(k)})^{-1} = \sigma_0^{-2} I_{Ld} + \left(\frac{1}{\sigma_0^2 + 1} - \sigma_0^{-2}\right) u^{(k)} (u^{(k)})^\top,$$

$$(\tilde{\Sigma}^{(k)})^{-1} \tilde{\Sigma}^{(l)} = I_{Ld} + \left(\frac{\sigma_0^2}{\sigma_0^2 + 1} - 1\right) u^{(k)} (u^{(k)})^\top + \sigma_0^{-2} u^{(l)} (u^{(l)})^\top + \left(\frac{1}{\sigma_0^2 + 1} - \sigma_0^{-2}\right) u^{(k)} (u^{(k)})^\top u^{(l)} (u^{(l)})^\top,$$

and

$$(A.14)$$

$$D_{KL}\left(X^{(l)}||X^{(k)}\right)$$

$$= \frac{n}{2}\left(Ld + \left(\frac{\sigma_0^2}{\sigma_0^2 + 1} - 1 + \sigma_0^{-2}\right) + \left(\frac{1}{\sigma_0^2 + 1} - \sigma_0^{-2}\right)\left((u^{(k)})^\top u^{(l)}\right)^2 - Ld\right)$$

$$= \frac{n}{2\sigma_0^2(\sigma_0^2 + 1)} \cdot \left(1 - \left((u^{(k)})^\top u^{(l)}\right)^2\right) = \frac{n}{2\sigma_0^2(\sigma_0^2 + 1)} \left\|\sin\Theta\left(U^{(k)}, U^{(l)}\right)\right\|^2$$

$$\stackrel{\text{(A.13)}}{\leq \frac{2n\theta^2}{\sigma_0^2(1 + \sigma_0^2)}}.$$

Step 3 We finalize the proof by the generalized Fano's lemma. Specifically by (Yu, 1997, Lemma 3), we have

$$\inf_{\hat{U}} \sup_{\Sigma \in \mathcal{F}_{p,r}(\sigma_{\text{sum}},\sigma_{\text{max}},\nu)} \mathbb{E} \| \sin \Theta(\hat{U},U) \| \stackrel{\text{(A.8)}}{\geq} \inf_{\hat{U}} \sup_{\Sigma \in \{\Sigma^{(l)}\}_{l=1}^{N}} \mathbb{E} \| \sin \Theta(\hat{U},U) \|$$

$$\stackrel{\text{(A.13)(A.14)}}{\geq} \frac{\theta}{4} \left( 1 - \frac{\frac{2n\theta^2}{\sigma_0^2(1+\sigma_0^2)} + \log(2)}{\log(N)} \right) \stackrel{N \geq 3}{\geq} \frac{\theta}{4} \left( 1 - \frac{\frac{2n\theta^2}{\sigma_0^2(1+\sigma_0^2)} + \log(2)}{(d/8) \vee \log(3)} \right)$$

$$\geq \frac{\theta}{4} \left( 1 - \frac{\frac{2n\theta^2}{\sigma_0^2(1+\sigma_0^2)}}{(d/8) \vee \log(3)} - \frac{\log(2)}{(d/8) \vee \log(3)} \right)$$

$$\geq \frac{\theta}{4} \left( 1 - \frac{\frac{2n\theta^2}{\sigma_0^2(1+\sigma_0^2)}}{d/8} - \frac{\log(2)}{\log(3)} \right).$$

Now we set  $\theta = \left(\frac{\sigma_0^2(1+\sigma_0^2)}{2n} \cdot \left(\frac{d}{32}\right)\right)^{1/2} \wedge 1$ . Then, for uniform constant c > 0, we have

$$\theta \ge c \left( \sqrt{\frac{d}{n}} (\sigma_0 + \sigma_0^2) \wedge 1 \right) \ge \frac{c \left( \sqrt{d\sigma_0^2} + \sqrt{d\sigma_0^2 \cdot \sigma_0^2} \right)}{\sqrt{n}} \wedge 1$$

$$\stackrel{\text{(A.9)}}{\ge} \frac{c \left( \sigma_{\text{sum}} + \sigma_{\text{max}} \sigma_{\text{sum}} \right)}{\sqrt{n}} \wedge 1.$$

Therefore,

$$\inf_{\hat{U}} \sup_{\Sigma \in \mathcal{F}_{p,r}(\sigma_{\text{sum}},\sigma_{\text{max}},\nu)} \mathbb{E} \| \sin \Theta(\hat{U}, U) \|$$

$$\geq c \left( \frac{\sigma_{\text{sum}} + \sigma_{\text{max}}\sigma_{\text{sum}}}{\sqrt{n}} \wedge 1 \right) \gtrsim \left( \frac{\sigma_{\text{sum}}}{(n\nu)^{1/2}} + \frac{\sigma_{\text{sum}}\sigma_{\text{max}}}{n^{1/2}\nu} \right) \wedge 1,$$

which has finished the proof for (A.1).

The proof of (A.2) is similar to (A.1): we still (a) first construct a series of candidate covariance matrices, (b) prove separateness of these covariance matrices and boundedness of KL divergence of random samples, and (c) apply generalized Fano's lemma to finalize the proof.

We still assume  $\nu=1$  without loss of generality. Since  $\sigma_{\max} \leq \sigma_{\text{sum}}$ , (A.2) is directly implied by (A.1) (which has been just proved) when r is a constant. Thus, we can assume  $r \geq 50$  in this part of proof without loss of generality. By the Varshamov-Gilbert bound (Massart, 2007, Lemma 4.7), we can find  $w^{(1)}, \ldots, w^{(N)} \subseteq \{\pm 1\}^r$ , such that

(A.15) 
$$\|w^{(l)} - w^{(k)}\|_{2}^{2} \ge r \text{ for all } 1 \le k \ne l \le N,$$

and  $N \ge \exp(r/8)$ . Consider the following set of covariance matrices for l = 1, ..., N,

$$A^{(l)} = \begin{bmatrix} (\theta w^{(l)})^\top \\ \frac{1}{\sqrt{L}} I_r \\ \vdots \\ \frac{1}{\sqrt{L}} I_r \\ 0_{(p-dr-1)\times r} \end{bmatrix}, \quad A^{(l)} = U^{(l)} R^{(l)} \text{ is the QR orthogonalization;}$$

$$\Sigma^{(l)} = A^{(l)} (A^{(l)})^{\top} + D \in \mathbb{R}^{p \times p}, \quad D_{ij} = \begin{cases} \sigma_{\max}^2, & i = j = 1; \\ 0, & \text{otherwise.} \end{cases}$$

Here,  $L = \lceil 1/c_I \rceil$ ;  $w^{(l)} \in \mathbb{R}^r$  has i.i.d. Rademacher entries;  $0 < \theta \le \sqrt{(c_I \wedge 1)/r}$  is some parameter to be determined later;  $\frac{1}{\sqrt{L}}I_r$  is repeated for L times; by design, the noise only appears in the first entry of the vector, so that the conditions

$$\max_{i} D_{ii} = D_{11} \le \sigma_{\max}^{2}$$
 and  $\sum_{i=1}^{p} D_{ii} = D_{11} \le \sigma_{\text{sum}}^{2}$ 

naturally hold, provided that  $\sigma_{\text{sum}} \geq \sigma_{\text{max}}$ .

By the relationship between singular values of the matrix and its subma-

trices (see (Cai and Zhang, 2018b, Lemma 2)), we have

$$\lambda_r \left( A^{(l)} \right) \ge \lambda_r \left( \begin{bmatrix} \frac{1}{\sqrt{L}} I_r \\ \vdots \\ \frac{1}{\sqrt{L}} I_r \end{bmatrix} \right) = 1,$$

$$\|A^{(l)}\| \le \left( \|\theta w^{(l)}\|_2^2 + \left\| \begin{bmatrix} \frac{1}{\sqrt{L}} I_r \\ \vdots \\ \frac{1}{\sqrt{L}} I_r \end{bmatrix} \right\|^2 \right)^{1/2},$$

which means

$$I(U^{(l)}) = \frac{p}{r} \max_{i} \|e_{i}^{\top} U^{(l)}\|_{2}^{2} \leq \frac{p}{r} \max_{i} \|e_{i}^{\top} A^{(l)} (R^{(l)})^{-1}\|_{2}^{2}$$

$$\leq \frac{p}{r} \max_{i} \|e_{i}^{\top} A^{(l)}\|_{2}^{2} \cdot \lambda_{r}^{-2} (R^{(l)})$$

$$\leq \frac{p}{r} \max \left\{ \theta^{2} r, \frac{1}{L} \right\} \cdot \lambda_{r}^{-2} (A^{(l)}) \leq c_{I} p / r.$$

Therefore,

(A.16) 
$$\Sigma^{(1)}, \dots, \Sigma^{(N)} \in \mathcal{F}_{p,n,r}(\sigma_{\text{sum}}, \sigma_{\text{max}}, \nu).$$

Again, suppose  $X^{(l)} = [X_1^{(l)}, \dots, X_n^{(l)}] \stackrel{iid}{\sim} N(0, \Sigma^{(l)})$  for  $l = 1, \dots, N$ . Next, we evaluate the  $\sin \Theta$  distances between each pair of  $(U^{(l)}, U^{(k)})$  and the KL divergence among  $X^{(l)}$ 's. Similarly to the proof for the first part of this theorem, we introduce a "condensed version" of  $A^{(l)}, \Sigma^{(l)}$ , and  $X^{(l)}$ .

$$\begin{split} \tilde{A}^{(l)} &= \begin{bmatrix} (\theta w^{(l)})^\top \\ I_r \end{bmatrix} \in \mathbb{R}^{(r+1)\times r}, \quad \tilde{A}^{(l)} = \tilde{U}^{(l)} R^{(l)} \text{ is the QR decomposition,} \\ \tilde{\Sigma}^{(l)} &= \tilde{A}^{(l)} (\tilde{A}^{(l)})^\top + \tilde{D} \in \mathbb{R}^{(r+1)\times (r+1)}, \quad \tilde{D}_{ij} = \left\{ \begin{array}{l} \sigma_{\max}^2, & i=j=1; \\ 0, & \text{otherwise,} \end{array} \right. \\ \tilde{X}^{(l)} &= [\tilde{X}_1^{(l)}, \dots, \tilde{X}_n^{(l)}], \quad X_i^{(l)} = T \tilde{X}_i^{(l)} \in \mathbb{R}^{r+1}, \\ \tilde{X}^{(l)} &= [\tilde{X}_1^{(l)}, \dots, \tilde{X}_n^{(l)}], \quad X_i^{(l)} = T \tilde{X}_i^{(l)} \in \mathbb{R}^{r+1}, \\ 0_{r\times 1} & \frac{1}{\sqrt{L}} I_r \\ \vdots & \vdots \\ 0_{r\times 1} & \frac{1}{\sqrt{L}} I_r \\ 0_{(p-Lr-1)\times 1} & 0_{(p-Lr-1)\times r} \\ \end{array} \right]. \end{split}$$

Then,  $\tilde{A}^{(l)},\,A^{(l)},\,\tilde{U}^{(l)},$  and  $U^{(l)}$  can be similarly related via T,

(A.17) 
$$T\tilde{A}^{(l)} = A^{(l)}, \quad T\tilde{U}^{(k)} = U^{(k)}.$$

One can also verify that  $\tilde{X}^{(l)} \stackrel{iid}{\sim} N(0, \tilde{\Sigma}^{(l)})$ . Noting that

$$v^{(l)} = \frac{1}{\sqrt{1 + r\theta^2}} \begin{pmatrix} 1 \\ -\theta w^{(l)} \end{pmatrix} \in \mathbb{R}^{r+1}$$

is the orthogonal complement to  $\tilde{A}^{(l)}$ , we have

$$\begin{split} & \left\| \sin \Theta(\tilde{U}^{(k)}, \tilde{U}^{(l)}) \right\| = \left\| (v^{(l)})^{\top} \tilde{U}^{(k)} \right\| = \left\| (v^{(l)})^{\top} \tilde{A}^{(k)} (R^{(l)})^{-1} \right\| \\ & \geq \left\| (v^{(l)})^{\top} \tilde{A}^{(k)} \right\| \cdot \lambda_r^{-1} (A^{(l)}) \geq \left\| \frac{\theta w^{(l)} - \theta w^{(k)}}{\sqrt{1 + r\theta^2}} \right\|_2 \frac{1}{\sqrt{1 + r\theta^2}} \\ & = \frac{\theta}{1 + r\theta^2} \left\| w^{(l)} - w^{(k)} \right\|_2. \end{split}$$

Since  $0 \le \theta \le \sqrt{(c_I \land 1)/r}$ , we additionally have

(A.18) 
$$\begin{aligned} & \left\| \sin(U^{(k)}, U^{(l)}) \right\| \stackrel{\text{(A.17)}}{=} \left\| \sin(\tilde{U}^{(k)}, \tilde{U}^{(l)}) \right\| \ge \frac{\theta}{1 + r\theta^2} \left\| w^{(l)} - w^{(k)} \right\|_2 \\ & \ge \frac{\theta}{2} \left\| w^{(l)} - w^{(k)} \right\|_2 \stackrel{\text{(A.15)}}{\geq} \frac{\sqrt{r}\theta}{2}, \quad \text{for all } 1 \le k \ne l \le N. \end{aligned}$$

Next, we consider the KL divergence among these samples. Given the linear relationship  $X^{(l)}=T\tilde{X}^{(l)}, X^{(k)}=T\tilde{X}^{(k)}$  with non-singular map T, we have

$$\begin{split} &D_{KL}\left(X^{(l)}||X^{(k)}\right) = D_{KL}\left(\tilde{X}^{(l)}||\tilde{X}^{(k)}\right) \\ &= \frac{n}{2}\left(\operatorname{tr}\left((\tilde{\Sigma}^{(k)})^{-1}\tilde{\Sigma}^{(l)}\right) - (r+1) + \log\left(\frac{\det(\tilde{\Sigma}^{(k)})}{\det(\tilde{\Sigma}^{(l)})}\right)\right). \end{split}$$

Noting that

$$\tilde{\Sigma}^{(k)} = \begin{bmatrix} \theta^2 r + \sigma_{\max}^2 & \theta(w^{(k)})^\top \\ \theta w^{(k)} & I_r \end{bmatrix},$$

 $\det(\tilde{\Sigma}^{(k)})=\det(\tilde{\Sigma}^{(l)})$  by symmetry. By the matrix inversion formula and calculation, one has

$$(\tilde{\Sigma}^{(k)})^{-1} = \begin{bmatrix} \sigma_{\max}^{-2} & -\sigma_{\max}^{-2} \nu(w^{(k)})^{\top} \\ -\sigma_{\max}^{-2} \nu w^{(k)} & I_r + \sigma_{\max}^{-2} \nu^2 w^{(k)} (w^{(k)})^{\top} \end{bmatrix},$$

and

$$D_{KL}\left(X^{(l)}||X^{(k)}\right) = \frac{n}{2}\left(\operatorname{tr}\left((\tilde{\Sigma}^{(k)})^{-1}\tilde{\Sigma}^{(l)}\right) - (r+1)\right)$$

$$= \frac{n}{2}\left((r+1) + 2\sigma_{\max}^{-2}\theta^{2}\left(r - (w^{(k)})^{\top}w^{(l)}\right) - (r+1)\right)$$

$$= \frac{n}{2}\left(2\sigma_{\max}^{-2}\theta^{2}(r - (w^{(k)})^{\top}w^{(l)})\right)$$

$$\leq n\sigma_{\max}^{-2}\theta^{2}r.$$

Finally, by generalized Fano's lemma (Yu, 1997, Lemma 3),

$$\inf_{\hat{U}} \sup_{\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{\text{sum}},\sigma_{\text{max}},\nu)} \left\| \sin \Theta(\hat{U}, U) \right\| \\
\stackrel{\text{(A.16)}}{\geq} \inf_{\hat{U}} \sup_{\Sigma \in \{\Sigma_k\}_{k=1}^N} \left\| \sin \Theta(\hat{U}, U) \right\| \\
\stackrel{\text{(A.18)(A.19)}}{\geq} \frac{\sqrt{r\theta}}{4} \left( 1 - \frac{n\sigma_{\text{max}}^{-2}\theta^2 r + \log(2)}{r/8} \right).$$

Set  $\theta = \sigma_{\max}/(32\sqrt{n}) \wedge \sqrt{(c_I \wedge 1)/r}$ . Given  $r \geq 50$ , we have

$$1 - \frac{n\sigma_{\max}^{-2}\theta^2r + \log(2)}{r/8} \ge 1 - \frac{r/32 + \log(2)}{r/8} \ge 1/3,$$

which means

$$\inf_{\hat{U}} \sup_{\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{\text{sum}},\sigma_{\text{max}},\nu)} \left\| \sin \Theta(\hat{U},U) \right\| \gtrsim c \left( \sqrt{\frac{r}{n}} \sigma_{\text{max}} \wedge 1 \right) = c \left( \frac{r^{1/2} \sigma_{\text{max}}}{(n\nu)^{1/2}} \wedge 1 \right).$$

for some constant c > 0 that only relies on  $c_I$ . Thus, we have finished the proof for (A.2).

PROOF OF PROPOSITION 1. For both lower bounds, it suffices to consider the rank-one case (r=1), where  $\hat{U}$  and U are unit vectors denoted by  $\hat{u}$  and u.

In order to prove the lower bound (2.11) on  $\hat{u}^{\text{SVD}}$ , it suffices to show

(A.20) 
$$\sup_{\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{\text{sum}},\sigma_{\text{max}},\nu)} \mathbb{E} \| \sin \Theta(\hat{u}^{\text{SVD}},u) \| \gtrsim \frac{\sigma_{\text{sum}}}{(n\nu)^{1/2}} \wedge 1,$$

and

(A.21) 
$$\sup_{\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{\text{sum}},\sigma_{\text{max}},\nu)} \mathbb{E} \| \sin \Theta(\hat{u}^{\text{SVD}},u) \| \gtrsim \frac{\sigma_{\text{max}}^2}{\nu} \wedge 1.$$

Note that (A.20) was implied by the minimax lower bound of Theorem 2. Thus we only need to show (A.21). By scaling  $\sigma_{\text{max}}$ , without loss of generality, we set  $\nu = 1$ .

Let  $L = \lceil \frac{1}{c_I} \rceil$ . Consider  $\Sigma = uu^\top + D$ , where

$$u = \begin{bmatrix} \frac{1}{\sqrt{L}} \mathbf{1}_L \\ \mathbf{0}_{p-L} \end{bmatrix}, \quad D_{ij} = \begin{cases} \sigma_{\max}^2 \wedge \frac{1}{L}, & i = j = 1; \\ 0, & \text{otherwise,} \end{cases}$$

Then  $I(u) \leq c_I p/r$ ,  $\lambda_r(uu^{\top}) = 1$ , and hence  $\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{\text{sum}}, \sigma_{\text{max}}, \nu)$  for any  $\sigma_{\text{sum}}$  that satisfies  $\sigma_{\text{max}}^2 \leq \sigma_{\text{sum}}^2 \leq p\sigma_{\text{max}}^2$ . Let  $X_1, \ldots, X_n \stackrel{iid}{\sim} N(0, \Sigma)$ . Then each  $X_k$  (ignoring the zero entries) can be expressed as

$$X_k = \frac{1}{\sqrt{L}} 1_L \gamma_k + e_1 \varepsilon_k,$$

where  $\gamma_k \sim N(0,1)$  and  $\varepsilon_k \sim N(0,\sigma_{\max}^2 \wedge \frac{1}{L})$  are independent. Thus,

$$n\hat{\Sigma} = XX^{\top} = \sum_{k=1}^{n} X_k X_k^{\top}$$

$$= \sum_{k=1}^{n} \frac{\gamma_k^2}{L} \mathbf{1}_L \mathbf{1}_L^{\top} + \left(\sum_{k=1}^{n} \varepsilon_k^2\right) e_1 e_1^{\top} + \left(\sum_{k=1}^{n} \frac{\gamma_k \varepsilon_k}{\sqrt{L}}\right) \left(e_1 \mathbf{1}_L^{\top} + \mathbf{1}_L e_1^{\top}\right).$$

Note that

$$\begin{split} & \sum_{k=1}^{n} \frac{\gamma_k^2}{L} \sim \frac{\chi_n^2}{L}, \quad \sum_{k=1}^{n} \varepsilon_k^2 \sim \left(\sigma_{\max}^2 \wedge \frac{1}{L}\right) \chi_n^2, \\ & \sum_{k=1}^{n} \frac{\gamma_k \varepsilon_k}{\sqrt{L}} \sim \sqrt{\frac{\sigma_{\max}^2}{L} \wedge \frac{1}{L^2}} \cdot N(0,1) \cdot \sqrt{\chi_n^2}. \end{split}$$

By the tail bound of the  $\chi^2$  distribution (c.f., (Laurent and Massart, 2000, Lemma 1)) and the symmetry of the normal distribution,

$$\sum_{k=1}^n \frac{\gamma_k^2}{L} \leq \frac{n+4\sqrt{n}}{L}, \quad \sum_{k=1}^n \varepsilon_k^2 \geq \left(\sigma_{\max}^2 \wedge \frac{1}{L}\right) (n-4\sqrt{n}), \quad \text{and} \quad \sum_{k=1}^n \frac{\gamma_k \varepsilon_k}{\sqrt{L}} \geq 0$$

hold with probability at least some constant c > 0. Denote

$$a \triangleq \frac{\sum_{k=1}^{n} \varepsilon_k^2 / \sqrt{L}}{\sum_{k=1}^{n} \gamma_k^2 / L}, \quad b \triangleq \frac{\sum_{k=1}^{n} \gamma_k \varepsilon_k}{\sum_{k=1}^{n} \gamma_k^2 / \sqrt{L}}.$$

For large n, we have

(A.22) 
$$a \ge \frac{1}{2} \left( \sqrt{L} \sigma_{\text{max}}^2 \wedge (1/\sqrt{L}) \right), \quad b \ge 0$$

with probability at least 0 < c < 1. Next, we introduce the following lemma to characterize the leading singular vector of  $\hat{\Sigma}$ .

LEMMA 2. Let  $L \geq 2$ . Consider  $A = 1_L 1_L^{\top} + a \sqrt{L} e_1 e_1^{\top} + b \left( e_1 1_L^{\top} + 1_L e_1^{\top} \right)$  and the leading singular vector of A is  $\hat{u}$ ,  $u = \frac{1}{\sqrt{L}} 1_L$ . If  $b \geq 0$ ,  $a \geq 0$ , we have

(A.23) 
$$\|\sin\Theta(\hat{u}, u)\| \ge c(a \land 1)/L.$$

PROOF OF LEMMA 2. The statement is clearly true when a=0. Then we can assume a>0 in the following analysis. Since A is symmetric and all entries of A are positive, the leading singular vector of A, i.e.,

$$\hat{u} = \underset{u \in \mathbb{R}^L, ||u||_2 = 1}{\arg\max} |u^{\top} A u|,$$

must have the same signs for all its entries. Since A has the block structure, such that  $A_{i.} = A_{i'.}$  and  $A_{.j} = A_{.j'}$  for any  $2 \le i, i', j, j' \le L$ ,  $\hat{u}$  must have the same structure in the sense that  $\hat{u}_i = \hat{u}_{i'}$  for any  $2 \le i, i' \le L$ . Based these, we can write  $\hat{u}$  as

$$\hat{u} = \left(z, \sqrt{(1-z^2)/(L-1)} \mathbf{1}_{L-1}\right)^{\top}$$

for some  $0 \le z \le 1$ . Then,

$$\hat{u}^{\top} A \hat{u} = \left(z + \sqrt{(1 - z^2)(L - 1)}\right)^2 + a\sqrt{L}z^2 + 2bz \left(z + \sqrt{(1 - z^2)(L - 1)}\right),$$

$$\frac{\partial}{\partial z} (\hat{u}^{\top} A \hat{u}) = 2\left(z + \sqrt{(1 - z^2)(L - 1)}\right) \left(1 - z\sqrt{\frac{L - 1}{1 - z^2}}\right) + 2\sqrt{L}az$$

$$+ 2b\left(2z + \sqrt{(1 - z^2)(L - 1)} - z^2\sqrt{\frac{L - 1}{1 - z^2}}\right).$$

We analyze 
$$2\left(z+\sqrt{(1-z^2)(L-1)}\right)\left(1-z\sqrt{\frac{L-1}{1-z^2}}\right)+2\sqrt{L}az$$
 and  $2b\left(2z+\sqrt{(1-z^2)(L-1)}-z^2\sqrt{\frac{L-1}{1-z^2}}\right)$  separately as below.

• If  $0 \le z \le 1/\sqrt{2}$ , we immediately have

$$2z + \sqrt{(1-z^2)(L-1)} - z^2\sqrt{\frac{L-1}{1-z^2}} = 2z + (1-2z^2)\sqrt{\frac{L-1}{1-z^2}} \ge 2z \ge 0;$$

if  $1/\sqrt{2} < z \le 1/\sqrt{L} + (a \wedge 1)/(10L)$ , we must have  $L = 2, z \le 0.8$ , and

$$2z + \sqrt{(1-z^2)(L-1)} - z^2 \sqrt{\frac{L-1}{1-z^2}} = 2z + (1-2z^2)\sqrt{\frac{L-1}{1-z^2}}$$
$$\geq \sqrt{2} + (1-2(.8)^2) \cdot \sqrt{\frac{1}{1-(.8)^2}} > 0.$$

Provided that  $b \geq 0$ , we always have

(A.24) 
$$2b\left(2z + \sqrt{(1-z^2)(L-1)} - z^2\sqrt{\frac{L-1}{1-z^2}}\right) \ge 0.$$

• Next, we consider  $Q \triangleq 2\left(z + \sqrt{(1-z^2)(L-1)}\right)\left(1 - z\sqrt{\frac{L-1}{1-z^2}}\right) + 2\sqrt{L}az$ . When  $0 \le z < 1/\sqrt{L}$ , by calculation we can verify that  $1 - z\sqrt{\frac{L-1}{1-z^2}} > 0$ , which means Q > 0.

When  $1/\sqrt{L} \le z \le 1/\sqrt{L} + (a \land 1)/(10L) \le 1/\sqrt{L} + 1/(10L)$ , we have

(A.25) 
$$0 \le z + \sqrt{(1-z^2)(L-1)} \le \frac{1}{\sqrt{L}} + \frac{1}{10L} + \sqrt{(1-1/L)(L-1)} = \sqrt{L} + \frac{1}{10L}.$$

For any  $1/\sqrt{L} \le z \le 1/\sqrt{L} + \frac{(a \land 1)}{10L}$ , we have

$$\left(1 - z\sqrt{\frac{L-1}{1-z^2}}\right)' = -\sqrt{\frac{L-1}{1-z^2}} - \frac{z^2}{1-z^2}\sqrt{\frac{L-1}{1-z^2}} = -\sqrt{\frac{(L-1)}{(1-z^2)^3}}$$
$$\geq -\sqrt{\frac{(L-1)}{(1-z^2)^3}} \geq -4.63\sqrt{L-1}.$$

Thus, by Taylor's theorem, there exists  $\xi \in \left[\frac{1}{\sqrt{L}}, \frac{1}{\sqrt{L}} + \frac{a \wedge 1}{10L}\right]$  such that

$$\left(1 - z\sqrt{\frac{L-1}{1-z^2}}\right) = \left(1 - z\sqrt{\frac{L-1}{1-z^2}}\right) \Big|_{z=\frac{1}{\sqrt{L}}}$$

$$+ \left(z - \frac{1}{\sqrt{L}}\right) \left(1 - z\sqrt{\frac{L-1}{1-z^2}}\right)' \Big|_{z=\xi},$$

$$\ge 0 - 4.63\sqrt{L-1} \left(z - \frac{1}{\sqrt{L}}\right) \ge -4.63 \cdot \frac{(a \wedge 1)\sqrt{L-1}}{10L}.$$

Therefore, for any a>0 and  $1/\sqrt{L}\leq z\leq 1/\sqrt{L}+(a\wedge 1)/(10L),$  we have

$$Q \stackrel{\text{(A.25)(A.26)}}{\geq} - \frac{2 \cdot 4.65(a \wedge 1)\sqrt{L-1}}{10L} \cdot \left(\sqrt{L} + \frac{1}{10L}\right) + 2\sqrt{L}az$$

$$\geq - \frac{9.3(a \wedge 1)\sqrt{L-1}}{10L} \left(\sqrt{L} + \frac{1}{\sqrt{L}}\right) + \frac{2\sqrt{L}a}{\sqrt{L}} > 0.$$

In summary of the previous two bullet points, for all  $0 \le z \le 1/\sqrt{L} + (a \land 1)/(10L)$ , we have

$$\frac{\partial}{\partial z}(\hat{u}^{\top}A\hat{u}) > 0.$$

Now, if  $\hat{u}$  is truly the singular vector of A, z must be a stationary point of  $\hat{u}^{\top}A\hat{u}$ , i.e.,  $\frac{\partial}{\partial z}(\hat{u}^{\top}A\hat{u})=0$ , which additionally means  $z\geq 1/\sqrt{L}+(a\wedge 1)/(10L)$ . Finally, if  $\hat{u}$  is the leading singular vector of A, we have  $z\geq 1/\sqrt{L}+(a\wedge 1)/(10L)$  and

$$\begin{aligned} &\|\sin\Theta(\hat{u},u)\| \overset{\text{(Cai and Zhang, 2018b, Lemma 1)}}{\geq} \frac{1}{\sqrt{2}} \min \|\hat{u} \pm u\|_2 = \frac{1}{\sqrt{2}} \|\hat{u} - u\|_2 \\ &\geq \frac{1}{\sqrt{2}} |\hat{u}_1 - u_1| = \frac{1}{\sqrt{2}} |z - 1/\sqrt{L}| \geq (a \wedge 1)/(10\sqrt{2}L). \end{aligned}$$

By Lemma 2 and L is of the constant order, the leading eigenvalue of  $n\hat{\Sigma}$ , i.e.,  $\hat{u} \in \mathbb{O}_{p,r}$ , satisfies

$$\mathbb{E}\|\sin\Theta(\hat{u}, u)\| \ge \mathbb{E}\left[\|\sin\Theta(\hat{u}, u)\| 1_{\{(\mathbf{A}.23) \text{ holds}\}}\right]$$

$$\ge \frac{c}{L} \left(\frac{\sqrt{L}\sigma_{\max}^2 \wedge (1/\sqrt{L})}{2} \wedge 1\right) \cdot \mathbb{P}((\mathbf{A}.23) \text{ holds})$$

$$\ge c \left(\sigma_{\max}^2 \wedge 1\right) \gtrsim \frac{\sigma_{\max}^2}{\nu} \wedge 1,$$

which proves (A.21) and hence also (2.11).

Then we consider the lower bound (2.12) for  $\hat{u}^{\mathrm{DD}}$ . Without loss of generality, we assume  $\nu = 1$ , again. Let  $L = \lceil 1/c_I \rceil$  and set

$$u = \begin{bmatrix} \frac{2}{\sqrt{5L}} \mathbf{1}_L \\ \frac{1}{\sqrt{5L}} \mathbf{1}_L \\ \mathbf{0}_{(p-2L)} \end{bmatrix}, \quad D = 0, \quad \Sigma = uu^\top + D.$$

Then  $I(u) = (p/r) \max_i ||u_i||_2^2 \le c_I \cdot p/r, \ \sum_i D_{ii} \le \sigma_{\text{sum}}^2, \max_i D_{ii} \le \sigma_{\text{max}}^2$ for any  $\sigma_{\text{sum}}^2, \sigma_{\text{sum}}^2 \geq 0$ . Thus,  $\Sigma \in \mathcal{F}_{p,n,r}(\sigma_{\text{sum}}, \sigma_{\text{max}}, \nu)$ . If  $X_1, \ldots, X_n \stackrel{iid}{\sim}$ 

$$X_k = u\gamma_k, \quad \gamma_k \sim N(0, 1).$$

$$n\hat{\Sigma} = \sum_{k=1}^{n} X_k X_k^{\top} = \sum_{k=1}^{n} \gamma_k^2 \begin{bmatrix} \frac{4}{5L} \mathbf{1}_{L \times L} & \frac{2}{5L} \mathbf{1}_{L \times L} & \mathbf{0}_{L \times (p-2L)} \\ \frac{2}{5L} \mathbf{1}_{L \times L} & \frac{1}{5L} \mathbf{1}_{L \times L} & \mathbf{0}_{L \times (p-2L)} \\ \mathbf{0}_{(p-2L) \times L} & \mathbf{0}_{(p-2L) \times L} & \mathbf{0}_{(p_2L) \times (p-2L)} \end{bmatrix}$$

Then we can write

$$\Delta(n\hat{\Sigma})/\sum_{k=1}^{n}\gamma_{k}^{2} = \begin{bmatrix} \frac{4}{5L}(1_{L\times L} - I_{L}) & \frac{2}{5L}1_{L\times L} & 0_{L\times(p-2L)} \\ \frac{2}{5L}1_{L\times L} & \frac{1}{5L}(1_{L\times L} - I_{L}) & 0_{L\times(p-2L)} \\ 0_{(p-2L)\times L} & 0_{(p-2L)\times L} & 0_{(p-2L)\times(p-2L)} \end{bmatrix} \triangleq R.$$

Since R is symmetric and  $u^{\top}R$  has a different direction than u, we know u is not an eigenvalue of R. Consequently, u is not a singular vector of R. In other words, if  $\hat{u}^{\mathrm{DD}}$  is the leading singular vector of the diagonal-deletion matrix of  $\hat{\Sigma}$ , we have

$$\|\sin\Theta(\hat{u}^{DD}, u)\| = \sqrt{1 - \hat{u}^{\top}u} \ge c > 0.$$

Here, c only depends on the constant  $c_I$ . We have thus finished the proof of this proposition.

PROOF OF PROPOSITION 2. Since  $\hat{\Sigma}$  is invariant after translation on Y, we can assume that the mean vector  $\mu = 0$  without loss of generality. Let

$$\Sigma_0 = \tilde{U}\Lambda \tilde{U}^{\top} = \begin{bmatrix} U & U_{\perp} \end{bmatrix} \begin{bmatrix} \Lambda_1 & \\ & \Lambda_2 \end{bmatrix} \begin{bmatrix} U^{\top} \\ U_{\perp}^{\top} \end{bmatrix}$$

be the full eigenvalue decomposition of  $\Sigma_0$ . Here,  $\tilde{U} = [U \ U_{\perp}]$  is the p-byp orthogonal matrix comprised of all eigenvectors of  $\Sigma_0$ ,  $U = [U \ U_{\perp}], \Lambda_1$  and  $\Lambda_2$  are r-by-r and (p-r)-by-(p-r) non-negative diagonal matrices containing the first r and the other (p-r) eigenvalues of  $\Sigma_0$ , respectively. We can also decompose  $Y_k$  based on its principal components as

$$Y_k = X_k + \varepsilon_k = U\Lambda_1^{1/2}\gamma_{1k} + U_\perp \Lambda_2^{1/2}\gamma_{2k} + \varepsilon_k,$$

where the random scores satisfy  $\mathbb{E}(\gamma_{1k}^{\top}, \gamma_{2k}^{\top}) = 0$ ,  $\text{Cov}((\gamma_{1k}^{\top}, \gamma_{2k}^{\top})) = I$ . We can further write this decomposition in a matrix form,

$$\Gamma_1 = [\gamma_{11} \cdots \gamma_{1n}], \quad \Gamma_2 = [\gamma_{21} \cdots \gamma_{2n}], \quad \Gamma = \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} = \begin{bmatrix} \gamma_{11} \cdots \gamma_{1n} \\ \gamma_{21} \cdots \gamma_{2n} \end{bmatrix},$$

$$Y = X^{(1)} + X^{(2)} + E, \quad X^{(1)} = U\Lambda_1^{1/2}\Gamma_1, \quad X^{(2)} = U_{\perp}\Lambda_2^{1/2}\Gamma_2.$$

We divide the rest of the proof in three steps.

Step 1 Define  $\hat{\Sigma}_X = (XX^\top - n\bar{X}\bar{X}^\top)/(n-1)$  and  $\hat{\Sigma}_{X^{(1)}} = (X^{(1)}X^{(1)\top} - n\bar{X}^{(1)}\bar{X}^{(1)\top})/(n-1)$ . By the same argument as the proof of Theorem 1, we can prove the following average perturbation inequality for  $\hat{\Sigma} - \hat{\Sigma}_X$ ,

$$\mathbb{E}_E \left\| \Delta \left( (n-1)(\hat{\Sigma} - \hat{\Sigma}_X) \right) \right\| \lesssim \sqrt{n} \sigma_{\text{sum}} \sigma_{\text{max}} + \sigma_{\text{sum}}^2 + \|X\| (\sigma_{\text{sum}} + \sqrt{r} \sigma_{\text{max}}) + n^{1/2} \|\bar{X}\|_2 \sigma_{\text{sum}}.$$

Here,  $\mathbb{E}_E$  means the expectation with respect to the noise part E. In addition, we can decompose  $(n-1)(\hat{\Sigma}-\hat{\Sigma}_{X^{(1)}})$  in the similar way as (6.1):

$$(n-1)(\hat{\Sigma}_X - \hat{\Sigma}_{X^{(1)}}) = X^{(1)}X^{(2)\top} + X^{(2)}X^{(1)\top} + X^{(2)}X^{(2)\top} - n\left(\bar{X}^{(1)}\bar{X}^{(2)\top} + \bar{X}^{(2)}\bar{X}^{(1)\top} + \bar{X}^{(2)}\bar{X}^{(2)\top}\right).$$

Therefore,

$$\mathbb{E}_{E} \left\| \Delta (n\hat{\Sigma} - n\hat{\Sigma}_{X}) \right\| \lesssim \sqrt{n} \sigma_{\text{sum}} \sigma_{\text{max}} + \sigma_{\text{sum}}^{2} + \|X\| (\sigma_{\text{sum}} + \sqrt{r} \sigma_{\text{max}}) + n^{1/2} \|\bar{X}\|_{2} \sigma_{\text{sum}} + 2\|X^{(1)}X^{(2)\top}\| + \|X^{(2)}\|^{2} + 2n\|\bar{X}^{(1)}\|_{2} \|\bar{X}^{(2)}\|_{2} + n\|\bar{X}^{(2)}\|_{2}^{2}.$$

Noting that  $\hat{\Sigma}^{(1)}$  is rank-r and has singular subspace U, by the robust  $\sin \Theta$  theorem (Theorem 3),

$$\mathbb{E}_{E} \left\| \sin \Theta(\hat{U}, U) \right\| \lesssim \left( \frac{\sqrt{n} \sigma_{\text{sum}} \sigma_{\text{max}} + \sigma_{\text{sum}}^{2} + \|X\| (\sigma_{\text{sum}} + \sqrt{r} \sigma_{\text{max}}) + n^{1/2} \|\bar{X}\|_{2} \sigma_{\text{sum}}}{(n-1) \lambda_{r}(\hat{\Sigma}^{(1)})} + \frac{2\|X^{(1)} X^{(2)\top}\| + \|X^{(2)}\|^{2} + 2n \|\bar{X}^{(1)}\|_{2} \|\bar{X}^{(2)}\|_{2} + n \|\bar{X}^{(2)}\|_{2}^{2}}{(n-1) \lambda_{r}(\hat{\Sigma}^{(1)})} \right) \wedge 1.$$

We analyze each term above as follows. Specifically, we introduce the following desirable probability event A, which happens if the inequalities (A.29) - (A.31) all hold:

(A.29) 
$$\sqrt{n} + C\sqrt{p} \ge \lambda_1 \left( \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} \right), \lambda_1 \left( \begin{bmatrix} \Gamma_1 \\ -\Gamma_2 \end{bmatrix} \right), \lambda_p \left( \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} \right), \lambda_p \left( \begin{bmatrix} \Gamma_1 \\ -\Gamma_2 \end{bmatrix} \right) \ge (\sqrt{n} - C\sqrt{p}) \lor 0,$$

(A.30) 
$$\lambda_r \left( \hat{\Sigma}_{X^{(1)}} \right) \ge \frac{5}{36} \lambda_r(\Lambda), \quad \|X^{(1)}\| \le 2\sqrt{n} \|\Lambda^{1/2}\|, \quad \|\sqrt{n}\bar{\Gamma}^{(1)}\|_2 \le \sqrt{n}/3,$$

(A.31) 
$$\|\Gamma_2\| \le C(\sqrt{n} + \sqrt{p}), \quad \|\bar{\Gamma}^{(2)}\|_2 \le C\sqrt{p/n}.$$

In the next two steps, we analyze each term in (A.28) given  $\mathcal{A}$  hold, then evaluate the probability that  $\mathcal{A}$  holds.

Step 2 Now we assume  $\mathcal{A}$  happens and (A.29)–(A.31) all hold. We plug in  $X^{(1)} = U\Lambda_1^{1/2}\Gamma_1$  and  $X^{(2)} = U_{\perp}\Lambda_2^{1/2}\Gamma_2$  and obtain

$$\begin{split} \|X^{(1)}X^{(2)\top}\| &= \|X^{(2)}X^{(1)\top}\| \leq \|\Lambda_1^{1/2}\|\|\Lambda_2^{1/2}\|\|\Gamma_1\Gamma_2^\top\| \\ &= \frac{1}{2}\lambda_1^{1/2}(\Lambda)\lambda_{r+1}^{1/2}(\Lambda) \left\| \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} \begin{bmatrix} \Gamma_1^\top & \Gamma_2^\top \end{bmatrix} - \begin{bmatrix} \Gamma_1 \\ -\Gamma_2 \end{bmatrix} \begin{bmatrix} \Gamma_1^\top & -\Gamma_2^\top \end{bmatrix} \right\|, \\ \left\| \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} \begin{bmatrix} \Gamma_1^\top & \Gamma_2^\top \end{bmatrix} - \begin{bmatrix} \Gamma_1 \\ -\Gamma_2 \end{bmatrix} \begin{bmatrix} \Gamma_1^\top & -\Gamma_2^\top \end{bmatrix} \right\| \\ &= \max_{w \in \mathbb{R}^p: \|w\|_2 \leq 1} \left| \left( \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} \right) - \left( \begin{bmatrix} \Gamma_1 \\ -\Gamma_2 \end{bmatrix} \right) \right\} - \min \left\{ \lambda_p^2 \left( \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} \right), \lambda_p^2 \left( \begin{bmatrix} \Gamma_1 \\ -\Gamma_2 \end{bmatrix} \right) \right\} \\ &\leq \max \left\{ \lambda_1^2 \left( \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} \right), \lambda_1^2 \left( \begin{bmatrix} \Gamma_1 \\ -\Gamma_2 \end{bmatrix} \right) \right\} - \min \left\{ \lambda_p^2 \left( \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} \right), \lambda_p^2 \left( \begin{bmatrix} \Gamma_1 \\ -\Gamma_2 \end{bmatrix} \right) \right\} \\ &\leq \left( \sqrt{n} + C\sqrt{p} \right)^2 - \left\{ \left( \sqrt{n} - C\sqrt{p} \right) \vee 0 \right\}^2 \leq C(\sqrt{np} + p). \end{split}$$

Thus,

$$||X^{(1)}X^{(2)\top}|| \le C(\sqrt{np} + p)\lambda_1^{1/2}(\Lambda)\lambda_{r+1}^{1/2}(\Lambda) \lesssim (\sqrt{np} + p)\lambda_r^{1/2}(\Lambda)\lambda_{r+1}^{1/2}(\Lambda).$$

Similarly,

$$||X^{(2)}X^{(2)\top}|| = ||\Lambda_2^{1/2}\Gamma_2\Gamma_2^{\top}\Lambda_2^{1/2}|| = ||\Lambda_2||||\Gamma_2||^2 \stackrel{(A.31)}{\leq} C(n+p)\lambda_{r+1}(\Lambda);$$

$$2n\|\bar{X}^{(1)}\|_{2}\|\bar{X}^{(2)}\|_{2} + n\|\bar{X}^{(2)}\|_{2}^{2}$$

$$=2n\|U\Lambda_{1}\bar{\Gamma}_{1}\|_{2}\|U_{\perp}\Lambda_{2}\bar{\Gamma}_{2}\|_{2} + n\|U\Lambda_{1}\bar{\Gamma}_{2}\|_{2}^{2}$$

$$\leq 2n\lambda_{1}^{1/2}(\Lambda)\lambda_{r+1}^{1/2}(\Lambda)\|\bar{\Gamma}_{1}\|_{2}\|\bar{\Gamma}_{2}\|_{2} + n\lambda_{r+1}(\Lambda)\|\bar{\Gamma}_{2}\|_{2}^{2}$$

$$(A.30)(A.31)$$

$$\lesssim n\lambda_{r}^{1/2}(\Lambda)\lambda_{r+1}^{1/2}(\Lambda)\sqrt{p/n} + n\lambda_{r+1}(\Lambda)p/n,$$

$$\|\bar{X}^{(1)}\|_{2} = \|U\Lambda_{1}\bar{\Gamma}_{1}\|_{2} \overset{(A.30)}{\leq} \lambda_{1}^{1/2}(\Lambda)\|\bar{\Gamma}^{(1)}\|_{2} \lesssim \lambda_{r}^{1/2}(\Lambda),$$

$$\|\hat{X}^{(2)}\|_{2} \leq \|\Lambda_{2}^{1/2}\|\|\bar{\Gamma}_{2}\|_{2} \lesssim \lambda_{r+1}(\Lambda)\sqrt{p/n},$$

$$\|\bar{X}\|_{2} \leq \|\bar{X}^{(1)}\|_{2} + \|\bar{X}^{(2)}\|_{2} \leq \lambda_{r}^{1/2}(\Lambda) + \lambda_{r+1}^{1/2}(\Lambda)\sqrt{p/n}.$$

Summarizing all previous bounds, when A holds, we have

$$\begin{split} &\mathbb{E}_{E} \left\| \sin \Theta(\hat{U}, U) \right\| \lesssim \left( \frac{\sqrt{n} \sigma_{\text{sum}} \sigma_{\text{max}} + \sigma_{\text{sum}}^{2} + \|X\| (\sigma_{\text{sum}} + \sqrt{r} \sigma_{\text{max}}) + n^{1/2} \|\bar{X}\|_{2} \sigma_{\text{sum}}}{n \lambda_{r} (\hat{\Sigma}_{X^{(1)}})} \right. \\ &+ \frac{2 \|X^{(1)} X^{(2)\top}\| + \|X^{(2)}\|^{2} + 2n \|\bar{X}^{(1)}\|_{2} \|\bar{X}^{(2)}\|_{2} + n \|\bar{X}^{(2)}\|_{2}^{2}}{n \lambda_{r} (\hat{\Sigma}_{X^{(1)}})} \wedge 1 \\ \lesssim & \left( \frac{\sqrt{n} \sigma_{\text{sum}} \sigma_{\text{max}} + \sigma_{\text{sum}}^{2} + (\sqrt{n} \lambda_{r}^{1/2} + \sqrt{p} \lambda_{r+1}^{1/2} (\Lambda)) (\sigma_{\text{sum}} + \sqrt{r} \sigma_{\text{max}})}{n \lambda_{r} (\Lambda)} \right. \\ &+ \frac{n^{1/2} \sigma_{\text{sum}} (\lambda_{r}^{1/2} (\Lambda) + \lambda_{r+1}^{1/2} (\Lambda) \sqrt{p/n})}{n \lambda_{r} (\Lambda)} \\ &+ \frac{(\sqrt{np} + p) \lambda_{r}^{1/2} (\Lambda) \lambda_{r+1}^{1/2} (\Lambda) + (n + p) \lambda_{r+1} (\Lambda) + n \lambda_{r}^{1/2} (\Lambda) \lambda_{r+1}^{1/2} (\Lambda) \sqrt{p/n} + n \lambda_{r+1} (\Lambda) p/n}{n \lambda_{r} (\Lambda)} \right) \wedge 1 \\ \lesssim & \left( \frac{\sigma_{\text{sum}} + \sqrt{r} \sigma_{\text{max}}}{n^{1/2} \lambda_{r}^{1/2} (\Lambda)} + \frac{\sigma_{\text{sum}} \sigma_{\text{max}}}{n^{1/2} \lambda_{r} (\Lambda)} + \frac{(\sqrt{np} + p) \lambda_{r+1}^{1/2} (\Lambda)}{n \lambda_{r}^{1/2} (\Lambda)} + \frac{\lambda_{r+1} (\Lambda)}{\lambda_{r} (\Lambda)} \right) \wedge 1. \end{split}$$

Here, the penultimate inequality is due to the following facts:

$$- \frac{\sigma_{\text{sum}}^2}{n\lambda_r(\Lambda)} \wedge 1 \leq \frac{\sigma_{\text{sum}}}{n^{1/2}\lambda_r^{1/2}(\Lambda)} \wedge 1;$$

– Since  $ab \wedge 1 \leq (a+b) \wedge 1$  for any  $a, b \geq 0$ ,

$$\frac{\sqrt{p}\lambda_{r+1}^{1/2}(\Lambda)(\sigma_{\text{sum}} + \sqrt{r}\sigma_{\text{max}})}{n\lambda_r(\Lambda)} \wedge 1 \leq \left(\frac{p^{1/2}\lambda_{r+1}^{1/2}(\Lambda)}{n^{1/2}\lambda_r^{1/2}(\Lambda)} + \frac{\sigma_{\text{sum}} + \sqrt{r}\sigma_{\text{max}}}{n^{1/2}\lambda_r^{1/2}(\Lambda)}\right) \wedge 1;$$

- By the same reason above,

$$\frac{n^{1/2}\sigma_{\text{sum}}\lambda_{r+1}^{1/2}(\Lambda)\sqrt{p/n}}{n\lambda_r(\Lambda)}\wedge 1 \leq \left(\frac{\sigma_{\text{sum}}}{n^{1/2}\lambda_r^{1/2}(\Lambda)} + \left(\frac{p\lambda_{r+1}(\Lambda)}{n\lambda_r(\Lambda)}\right)^{1/2}\right)\wedge 1,$$

$$- \left( \frac{p\lambda_{r+1}(\Lambda)}{n\lambda_r(\Lambda)} \right) \wedge 1 \le \left( \frac{p\lambda_{r+1}(\Lambda)}{n\lambda_r(\Lambda)} \right)^{1/2} \wedge 1,$$

- Step 3 In this step, we evaluate the probability that the event  $\mathcal{A}$  holds by giving probability upper bounds for (A.29) (A.31) as follows.
  - Noting that  $\begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} \in \mathbb{R}^{p \times n}$  and  $\begin{bmatrix} \Gamma_1 \\ -\Gamma_2 \end{bmatrix} \in \mathbb{R}^{p \times n}$  are random matrices with i.i.d. columns, by (Vershynin, 2010, Corollary 5.35),

$$\mathbb{P}\left(\sqrt{n} + C\sqrt{p} + t \ge \text{all singular values of } \begin{bmatrix} \Gamma_1 \\ \Gamma_2 \end{bmatrix} \begin{bmatrix} \Gamma_1 \\ -\Gamma_2 \end{bmatrix} \ge \sqrt{n} - C\sqrt{p} - t \right) \le \exp(-Ct^2/2).$$

By setting  $t = C\sqrt{p}$  for large constant C > 0, we know (A.29) holds with probability at least  $1 - C \exp(-Cp)$ .

- Since  $\Gamma_1 \in \mathbb{R}^{r \times n}$  has isotropic sub-Gaussian columns, based on the argument of (6.6) in the proof of Theorem 1, (A.30) holds with probability at least  $1 C \exp(-Cn)$ .
- Noting that  $\Gamma_2$  is a (p-r)-by-n random matrix with i.i.d. isotropic sub-Gaussian rows, by (Vershynin, 2011, Corollary 5.35),

$$\|\Gamma_2\| \le C(\sqrt{n} + \sqrt{p})$$

with probability at least  $1 - \exp(-C(n+p))$ ; by Bernstein-type concentration inequality (Vershynin, 2011, Proposition 5.16),

$$\mathbb{P}\left(\|\sqrt{n}\overline{\Gamma}^{(2)}\|_{2}^{2} \ge p + C\sqrt{px} + Cx\right) \le C\exp(-cx).$$

By setting x = Cp, we conclude that (A.31) holds with probability at least  $1 - C \exp(-Cp)$ .

To sum up, the event  $\mathcal{A}$  happens, i.e., (A.29) - (A.31) all hold, with probability at least  $1 - C \exp(-Cn) - C \exp(-Cp)$ .

Step 4 We finalize the proof in this step.

$$\mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\| = \mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\| 1_{\mathcal{A}} + \mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\| 1_{\mathcal{A}^{c}}$$

$$\lesssim \frac{\sigma_{\text{sum}} + \sqrt{r}\sigma_{\text{max}}}{n^{1/2}\lambda_{r}^{1/2}(\Lambda)} + \frac{\sigma_{\text{sum}}\sigma_{\text{max}}}{n^{1/2}\lambda_{r}(\Lambda)} + \frac{((np)^{1/2} + p)\lambda_{r+1}^{1/2}(\Lambda)}{n\lambda_{r}^{1/2}(\Lambda)} + \frac{\lambda_{r+1}(\Lambda)}{\lambda_{r}(\Lambda)}$$

$$+ C \exp(-Cn) + C \exp(-Cp)$$

$$\lesssim \frac{\sigma_{\text{sum}} + \sqrt{r}\sigma_{\text{max}}}{n^{1/2}\lambda_{r}^{1/2}(\Lambda)} + \frac{\sigma_{\text{sum}}\sigma_{\text{max}}}{n^{1/2}\lambda_{r}(\Lambda)} + \frac{((np)^{1/2} + p)\lambda_{r+1}^{1/2}(\Lambda)}{n\lambda_{r}^{1/2}(\Lambda)} + \frac{\lambda_{r+1}(\Lambda)}{\lambda_{r}(\Lambda)},$$

where the last inequality is due to  $\sigma_{\text{sum}}^2/\lambda_r(\Lambda) \ge \exp(-Cp) + \exp(-Cn)$  in the assumption. Finally, the trivial upper bound 1 always holds for  $\mathbb{E}\left\|\sin\Theta(\hat{U},U)\right\|$ . We thus have finished this proof.

### A.2. Additional Proofs for Robust $\sin \Theta$ Theorem.

PROOF OF PROPOSITION 3. We first develop the lower bound with the incoherence constraint. We first assume  $\delta/\nu \leq 1/\sqrt{2}$ . Let  $d = 2\lfloor p/(2r) \rfloor$ ,  $\alpha, \beta \in \mathbb{R}^d$  be unit vectors such that

$$\alpha = \frac{1}{\sqrt{d}} (1, \dots, 1), \quad \beta = \frac{1}{\sqrt{d(1+\theta^2)}} (1+\theta, \dots, 1+\theta, 1-\theta, \dots, 1-\theta).$$

Clearly,  $f(\theta) \triangleq \|\alpha\alpha^{\top} - \beta\beta^{\top}\|$  is a continuous function of  $\theta$ . One can verify that f(0) = 0;  $f(1) = 1/\sqrt{2}$ , then there exists  $0 \leq \theta \leq 1$  to ensure that

Based on (A.32), we additionally construct

(A.33) 
$$U^{(1)} = \begin{bmatrix} \alpha_1 I_r \\ \vdots \\ \alpha_d I_r \\ 0_{(p-rd),r} \end{bmatrix}, \quad U^{(2)} = \begin{bmatrix} \beta_1 I_r \\ \vdots \\ \beta_d I_r \\ 0_{(p-rd),r} \end{bmatrix}.$$

Here,  $\frac{1}{\sqrt{d}}I_r$  is repeated for d times in  $U^{(1)}$ ; both  $\frac{1+\theta}{\sqrt{d(1+\theta^2)}}I_r$  and  $\frac{1-\theta}{\sqrt{d(1+\theta^2)}}I_r$  are repeated for d/2 times in  $U^{(2)}$ . Let  $M^{(1)} = \nu U^{(1)}(U^{(1)})^{\top}$ ,  $M^{(2)} = \nu U^{(2)}(U^{(2)})^{\top}$ ,

$$Z^{(1)} = \frac{1}{2}(M^{(2)} - M^{(1)}), Z^{(2)} = \frac{1}{2}(M^{(1)} - M^{(2)}).$$
 By such the construction,  $\lambda_r(M^{(1)}) = \lambda_r(M^{(2)}) = \nu, \ \|M^{(1)}\|/\lambda_r(M^{(1)}) = \|M^{(2)}\|/\lambda_r(M^{(2)}) = 1,$ 

$$I(U^{(1)}) = \frac{p}{r} \max_i \|e_i^\top U^{(1)}\|_2^2 \le \frac{p}{rd} = \frac{p}{r \cdot 2|p/(2r)|} < \frac{p}{r \cdot 2(p/(2r) - 1)} \le 2,$$

$$\begin{split} I(U^{(2)}) = & \frac{p}{r} \max_{i} \|e_{i}^{\top} U^{(2)}\|_{2}^{2} \leq \frac{p(1+\theta)^{2}}{r \cdot (d(1+\theta^{2}))} \leq \frac{p}{r} \cdot \frac{2}{d} = \frac{p}{r} \cdot \frac{1}{\lfloor p/(2r) \rfloor} \\ \leq & \begin{cases} 4 \cdot 1 \leq 4, & \text{if } 2r \leq p \leq 4r; \\ \frac{p}{r} \cdot \frac{1}{p/(2r) - 1} = \frac{2p}{p - 2r} \leq 4, & \text{if } 4r + 1 \leq p. \end{cases} \end{split}$$

$$\|\Delta(Z^{(1)})\| = \|\Delta(Z^{(2)})\| \stackrel{\text{Lemma 5}}{\leq} 2\|Z^{(2)}\| \leq 2 \left\| \frac{1}{2} \left( M^{(2)} - M^{(1)} \right) \right\|$$
$$= \nu \left\| \alpha \alpha^{\top} - \beta \beta^{\top} \right\| = \delta,$$

which means  $(M^{(1)}, Z^{(1)}), (M^{(2)}, Z^{(2)}) \in \mathcal{D}_{p,r}(\nu, \delta, t)$  for  $t \geq 4$ . On the other hand, by (Cai and Zhang, 2018b, Lemma 1),

$$\left\| \sin \Theta(U^{(1)}, U^{(2)}) \right\| \ge \frac{1}{2} \| U^{(1)}(U^{(1)})^{\top} - U^{(2)}(U^{(2)})^{\top} \|$$

$$\stackrel{\text{(A.33)}}{=} \frac{1}{2} \left\| \begin{bmatrix} (\alpha_1^2 - \beta_1^2) I_r & \cdots & (\alpha_1 \alpha_d - \beta_1 \beta_d) I_r \\ \vdots & & \vdots \\ (\alpha_d \alpha_1 - \beta_d \beta_1) I_r & \cdots & (\alpha_d^2 - \beta_d^2) I_r \end{bmatrix} \right\|$$

$$= \frac{1}{2} \| \alpha \alpha^{\top} - \beta \beta^{\top} \| = \delta/(2\nu).$$

Given  $M^{(1)} + Z^{(1)} = M^{(2)} + Z^{(2)}$ , we have

$$\inf_{\hat{U}} \sup_{(M,Z) \in \mathcal{D}_{p,r}(\nu,\delta,t)} \left\| \sin \Theta(\hat{U},U) \right\| 
\geq \inf_{\hat{U}} \sup_{(M,Z) \in \left\{ (M^{(1)},Z^{(1)}),(M^{(2)},Z^{(2)}) \right\}} \left\| \sin \Theta(\hat{U},U) \right\| 
\geq \inf_{\hat{U}} \frac{1}{2} \left( \left\| \sin \Theta(\hat{U},U^{(1)}) \right\| + \left\| \sin \Theta(\hat{U},U^{(2)}) \right\| \right) \geq \frac{1}{2} \left\| \sin \Theta(U^{(1)},U^{(2)}) \right\| = \frac{\delta}{4\nu}.$$

Next, if  $\delta/\nu \geq \sqrt{2}/2$ , let  $\delta_0 = \nu \cdot \sqrt{2}/2$ . By the previous argument, one can show

$$\inf_{\hat{U}} \sup_{(M,Z) \in \mathcal{D}_{p,r}(\nu,\delta,t)} \left\| \sin \Theta(\hat{U},U) \right\| \ge \frac{\delta_0}{4\nu} = \frac{\sqrt{2}}{8} \ge \frac{\sqrt{2}}{8} \left( \frac{\delta}{\nu} \wedge 1 \right).$$

In summary, we must have

$$\inf_{\hat{U}} \sup_{(M,Z) \in \mathcal{D}_{p,r}(\nu,\delta,t)} \left\| \sin \Theta(\hat{U}, U) \right\| \ge \frac{\sqrt{2}}{8} \left( \frac{\delta}{\nu} \wedge 1 \right)$$

in the first scenario that  $t \geq 4$ .

Then we consider the second part that  $t \geq p/r$ . Let

$$U^{(1)} = \begin{bmatrix} I_r \\ 0_{(p-r)\times r} \end{bmatrix}, \quad U^{(2)} = \begin{bmatrix} 0_{r\times r} \\ I_r \\ 0_{(p-2r)\times r} \end{bmatrix}$$

be two orthogonal matrices,  $M^{(1)} = \nu U^{(1)}(U^{(1)})^{\top}, M^{(2)} = \nu U^{(2)}(U^{(2)})^{\top}, Z^{(1)} = -M^{(1)}, Z^{(2)} = -M^{(2)}$ . Then clearly,  $M^{(1)} + Z^{(1)} = M^{(2)} + Z^{(2)}, \lambda_r(M^{(1)}) = \lambda_r(M^{(2)}) \ge \nu, \|\Delta(Z^{(1)})\| = \|\Delta(Z^{(2)})\| = 0,$ 

$$\|\sin\Theta(U^{(1)}, U^{(2)})\| = \left(1 - \lambda_r((U^{(1)})^\top U^{(2)})\right)^{1/2} = (1 - 0)^{1/2} = 1.$$

Moreover, for any  $t \geq p/r$ ,

$$I(U^{(1)}) = \frac{p}{r} \|e_i^\top U^{(1)}\|_2^2 = \frac{p}{r} \leq t, \quad I(U^{(2)}) = \frac{p}{r} \|e_i^\top U^{(2)}\|_2^2 = \frac{p}{r} \leq t.$$

We thus have

$$(M^{(1)}, Z^{(1)}), (M^{(2)}, Z^{(2)}) \in \mathcal{D}_{p,r}(\nu, \delta, t)$$

if  $t \ge p/r$ . Given  $M^{(1)} + Z^{(1)} = M^{(2)} + Z^{(2)}$ , we have

$$\inf_{\hat{U}} \sup_{(M,Z) \in \mathcal{D}_{p,r}(\nu,\delta,t)} \left\| \sin \Theta(\hat{U}, U) \right\|$$

$$\geq \inf_{\hat{U}} \sup_{(M,Z) \in \left\{ (M^{(1)}, Z^{(1)}), (M^{(2)}, Z^{(2)}) \right\}} \left\| \sin \Theta(\hat{U}, U) \right\|$$

$$\geq \inf_{\hat{U}} \frac{1}{2} \left( \left\| \sin \Theta(\hat{U}, U^{(1)}) \right\| + \left\| \sin \Theta(\hat{U}, U^{(2)}) \right\| \right) \geq \frac{1}{2} \left\| \sin \Theta(U^{(1)}, U^{(2)}) \right\| = \frac{1}{2},$$

which has finished the proof of this theorem.

## A.3. Proofs in Heteroskedastic Low-rank Matrix Denoising.

PROOF OF THEOREM 4. First, we assume  $\delta \in \mathbb{R}^{p_1}$ ,  $\delta_i = \sum_{j=1}^{p_2} \sigma_{ij}^2$  as the row-wise summation of variances. Note that

$$YY^{\top} = XX^{\top} + XE^{\top} + EX^{\top} + EE^{\top}.$$

Then,  $\mathbb{E}YY^{\top} = XX^{\top} + \operatorname{diag}(\delta)$ . By the Wishart-type heteroskedastic concentration inequality (c.f., Cai and Zhang (2018a)),

$$\mathbb{E}\left\|EE^{\top} - \operatorname{diag}(\delta)\right\| \lesssim \sigma_C^2 + \sigma_C \sigma_R + \sigma_R \sigma_{\max} \sqrt{\log(p_1 \wedge p_2)} + \sigma_{\max}^2 \log(p_1 \wedge p_2),$$

By Lemma 3 and  $||X|| \leq C\lambda_r(X)$ ,

(A.35) 
$$\mathbb{E} \|XE^{\top}\| \lesssim \|X\| \left(\sigma_C + \sqrt{r}\sigma_{\max}\right) \lesssim \lambda_r(X) \left(\sigma_C + \sqrt{r}\sigma_{\max}\right).$$

By Lemma 5,

$$\left\| \Delta (YY^{\top} - XX^{\top}) \right\| = \left\| \Delta (YY^{\top} - XX^{\top} - \operatorname{diag}(\delta)) \right\|$$

$$(A.36) \leq 2 \left\| YY^{\top} - XX^{\top} - \operatorname{diag}(\delta) \right\| \leq \left\| XE^{\top} + EX^{\top} + EE^{\top} - \operatorname{diag}(\delta) \right\|$$

$$\leq 2 \left\| EE^{\top} - \operatorname{diag}(\delta) \right\| + 4 \|EX^{\top}\|.$$

Combining (A.34), (A.35), and (A.36), we have

(A.37)
$$\mathbb{E} \left\| \Delta (YY^{\top} - XX^{\top}) \right\|$$

$$\lesssim \sigma_C^2 + \sigma_C \sigma_R + \sigma_R \sigma_{\max} \sqrt{\log(p_1 \wedge p_2)} + \sigma_{\max}^2 \log(p_1 \wedge p_2) + \lambda_r(X) \left( \sigma_C + \sqrt{r} \sigma_{\max} \right).$$

Note that the eigen-subspace of  $XX^{\top}$  is the same as U, i.e., the left singular subspace of X. Since  $I(U) \leq c_I p/r$ , the robust  $\sin \Theta$  theorem (Theorem 3) implies

$$\mathbb{E} \left\| \sin \Theta \left( \hat{U}, U \right) \right\| \leq \frac{C \mathbb{E} \| \Delta (YY^{\top} - XX^{\top}) \|}{\lambda_r^2(X)} \wedge 1$$

$$\lesssim \left( \frac{\sigma_C^2 + \sigma_C \sigma_R + \sigma_R \sigma_{\max} \sqrt{\log(p_1 \wedge p_2)} + \sigma_{\max}^2 \log(p_1 \wedge p_2) + \lambda_r(X) \left( \sigma_C + \sqrt{r} \sigma_{\max} \right)}{\lambda_r^2(X)} \right) \wedge 1$$

$$\lesssim \left( \frac{\sigma_C + \sqrt{r} \sigma_{\max}}{\lambda_r(X)} + \frac{\sigma_C \sigma_R + \sigma_R \sigma_{\max} \sqrt{\log(p_1 \wedge p_2)} + \sigma_{\max}^2 \log(p_1 \wedge p_2)}{\lambda_r^2(X)} \right) \wedge 1.$$

The last inequality is due to the fact that  $\sigma_C/\lambda_r(X) \wedge 1 \leq \sigma_C^2/\lambda_r^2(X) \wedge 1$ . In particular when  $\sigma_{\max} \lesssim \sigma_C/\max\{\sqrt{r}, \sqrt{\log(p_1 \wedge p_1)}\}$ , we have

$$\frac{\sqrt{r}\sigma_{\max}}{\lambda_r(X)} \lesssim \frac{\sigma_C}{\lambda_r(X)}, \quad \frac{\sigma_R\sigma_{\max}\sqrt{\log(p_1 \wedge p_2)}}{\lambda_r^2(X)} \lesssim \frac{\sigma_C\sigma_R}{\lambda_r^2(X)},$$

$$\frac{\sigma_{\max}^2 \log(p_1 \wedge p_2)}{\lambda_r^2(X)} \wedge 1 \lesssim \frac{\sigma_C^2}{\lambda_r^2(X)} \wedge 1 \leq \frac{\sigma_C}{\lambda_r(X)} \wedge 1.$$

We thus have

$$\mathbb{E}\left\|\sin\Theta(\hat{U}, U)\right\| \lesssim \left(\frac{\sigma_C}{\lambda_r(X)} + \frac{\sigma_C \sigma_R}{\lambda_r^2(X)}\right) \wedge 1.$$

### A.4. Proofs in Poisson PCA.

PROOF OF THEOREM 5. Denote  $E = Y - X \in \mathbb{R}^{p_1 \times p_2}$ . Recall the following tail probability bound of Poisson distribution (see, e.g., (Boucheron et al., 2013, Pages 22-23)),

$$\mathbb{P}\left(|Y_{ij} - X_{ij}| \ge t\right) \le 2\exp\left(-(t + X_{ij})\log\left(1 + \frac{t}{X_{ij}}\right) + t\right), \quad \forall t \ge 0.$$

Next, we aim to show

(A.39) 
$$\mathbb{P}(|Y_{ij} - X_{ij}| \ge t) \le 2\exp\left(1 - ct/\sqrt{X_{ij}}\right), \quad \forall t > 0.$$

for some uniform constant c > 0.

• If  $t < \sqrt{X_{ij}}$ ,

$$\mathbb{P}(|Y_{ij} - X_{ij}| \ge t) \le 1 \le 2\exp(1 - t/\sqrt{X_{ij}}), \quad \forall t > 0.$$

• If  $\sqrt{X_{ij}} \le t \le X_{ij}/2$ , by Taylor expansion for  $\log(1 + x/X_{ij})$ ,

$$(X_{ij} + t) \log \left( 1 + \frac{t}{X_{ij}} \right) - t \ge (X_{ij} + t) \left( \frac{t}{X_{ij}} - \frac{t^2}{2X_{ij}^2} \right) - t$$
$$= \frac{t^2}{X_{ij}} - \frac{t^2}{X_{ij}} \cdot \frac{X_{ij} + t}{2X_{ij}} \ge \frac{t^2}{X_{ij}} - \frac{t^2}{X_{ij}} \cdot \frac{3}{4} = \frac{t^2}{4X_{ij}} \ge \frac{t}{4\sqrt{X_{ij}}} - \frac{1}{16}.$$

Thus,

$$\mathbb{P}\left(|Y_{ij} - X_{ij}| \ge t\right) \le 2 \exp\left(\frac{1}{16} - \frac{t}{4\sqrt{X_{ij}}}\right).$$

• If  $X_{ij}/2 \le t \le 2X_{ij}$ , we shall note that

$$\frac{\partial}{\partial X_{ij}} \left( (t + X_{ij}) \log \left( 1 + \frac{t}{X_{ij}} \right) \right) = \log \left( 1 + \frac{t}{X_{ij}} \right) - \frac{t}{X_{ij}} \le 0,$$

then  $(t + X_{ij}) \log \left(1 + \frac{t}{X_{ij}}\right)$  is a decreasing function of  $X_{ij}$ . Thus,

$$\mathbb{P}\left(|Y_{ij} - X_{ij}| \ge t\right) \le 2\exp\left(-(t + X_{ij})\log\left(1 + \frac{t}{X_{ij}}\right) + t\right)$$

$$\le 2\exp\left(-(t + 2t)\log\left(1 + \frac{t}{2t}\right) + t\right)$$

$$\le 2\exp\left(-(3\log(1.5) - 1)t\right)$$

$$\le 2\exp\left(-\sqrt{c}(3\log(1.5) - 1)t/\sqrt{X_{ij}}\right)$$

$$(X_{ij} + t) \log \left(1 + \frac{t}{X_{ij}}\right) - t \ge (t+t) \log(1+1/2) - t.$$

• If  $t \geq 2X_{ij}$ ,

$$\mathbb{P}(|Y_{ij} - X_{ij}| \ge t) \le 2 \exp(-t \log(1+2) + t)$$
  
 
$$\le 2 \exp(-t/\sqrt{X_{ij}} \cdot (\sqrt{c}(\log(3) - 1))).$$

In summary, (A.39) always hold, which means  $E_{ij}/\sqrt{X_{ij}}$  is a sub-exponential distributed random variable. By the sub-exponential Wishart-type concentration inequality (c.f., Cai and Zhang (2018a)),

$$\mathbb{E}\left\|EE^{\top} - \mathbb{E}EE^{\top}\right\| \lesssim \sigma_{C}\sigma_{R} + \sigma_{C}^{2} + \sigma_{R}\sigma_{\max}\sqrt{\log(p_{1})\log(p_{2})} + \sigma_{\max}\log(p_{1})\log(p_{2}).$$

Suppose the right singular subspace of X is  $V \in \mathbb{O}_{p_2,r}$ . By Lemma 4 and  $||X|| \leq C\lambda_r(X)$ ,

$$\mathbb{E} \left\| X E^{\top} \right\| \lesssim \|X\| \mathbb{E} \|EV\| \lesssim \|X\| (\sigma_C + r\sigma_{\max}) \lesssim \lambda_r(X) (\sigma_C + \sqrt{r}\sigma_{\max}).$$

Now, the rest of the proof follows from the Inequality A.36 and the arguments below in proof of Theorem 4. We can finally prove that

$$\mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\|$$

$$\lesssim \left( \frac{\sigma_C + r\sigma_{\max}}{\lambda_r(X)} + \frac{\sigma_C \sigma_R + \sigma_R \sigma_{\max} \sqrt{\log(p_1) \log(p_2)} + \sigma_{\max}^2 \log(p_1) \log(p_2)}{\lambda_r^2(X)} \right) \wedge 1.$$

# A.5. Proofs in SVD Based on Heteroskedastic and Incomplete Data.

PROOF OF THEOREM 6.

Step 1 We first derive bounds for some key quantities, including  $\sigma_B^2$  and  $|||\mathbf{B}_k|||_{\psi_1}$  to be defined later, for the application of matrix concentration in the next step. Based on the property of sub-Gaussian random variable and  $||Y_{ij}||_{\psi_2} \leq C$ ,  $Y_{ij}$  has bounded moments

$$\mathbb{E}|Y_{ij}|^{\alpha} \le C, \quad \alpha = 1, 2, 3, 4.$$

Since

(A.40) 
$$\begin{pmatrix} \mathbb{E}\tilde{Y}\tilde{Y}^{\top} \end{pmatrix}_{ij} = \sum_{k=1}^{p_2} \mathbb{E}\tilde{Y}_{ik}\tilde{Y}_{jk} = \begin{cases} \sum_{k=1}^{n} \theta \mathbb{E}Y_{ik}^2, & i=j; \\ \sum_{k=1}^{n} \theta^2 \mathbb{E}Y_{ik}Y_{jk}, & i\neq j. \end{cases}$$
$$= \begin{cases} \theta(XX^{\top})_{ii} + \theta \sum_{k=1}^{p_2} \operatorname{Var}(Z_{ik}), & i=j; \\ \theta^2(XX^{\top})_{ij}, & i\neq j, \end{cases}$$

we know  $\Delta(\mathbb{E}\tilde{Y}\tilde{Y}^{\top}) = \Delta(\theta^2XX^{\top})$ , i.e.,  $\mathbb{E}\tilde{Y}\tilde{Y}^{\top}$  and  $\theta^2XX^{\top}$  share the off-diagonal part. Recall  $D(\cdot)$  and  $\Delta(\cdot)$  represent the diagonal and off-diagonal part of the matrix, respectively.

Next, we establish a concentration inequality for  $\|\tilde{Y}\tilde{Y}^{\top} - \mathbb{E}\tilde{Y}\tilde{Y}^{\top}\|$ . Note the following decomposition, (A.41)

$$\tilde{Y}\tilde{Y}^{\top} - \mathbb{E}\tilde{Y}\tilde{Y}^{\top} = \sum_{k=1}^{p_2} \left( \tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top} - \mathbb{E}\tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top} \right) \triangleq \sum_{k=1}^{p_2} B_k, \quad B_k = \tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top} - \mathbb{E}\tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top}.$$

Based on the assumption,

$$|\mathbb{E}\tilde{Y}_{ij}|^{\alpha} = \theta \mathbb{E}|Y_{ij}|^{\alpha} \le C\theta, \quad \alpha = 1, 2, 3, 4.$$

Then.

(A.42) 
$$0 \leq \mathbb{E}B_k B_k^{\top} = \mathbb{E}\left(\tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top} - \mathbb{E}\tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top}\right)^2$$
$$= \mathbb{E}\tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top} \tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top} - (\mathbb{E}\tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top})^2 \leq \mathbb{E}\tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top} \tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top},$$

$$\left| \left( \mathbb{E} \tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top} \tilde{Y}_{\cdot k} \tilde{Y}_{\cdot k}^{\top} \right)_{ij} \right| = \left| \mathbb{E} \tilde{Y}_{ik} \left( \sum_{s=1}^{p_1} \tilde{Y}_{sk}^2 \right) \tilde{Y}_{jk} \right|.$$

If  $i \neq j$ ,

$$\left| \mathbb{E}\tilde{Y}_{ik} \left( \sum_{s=1}^{p_1} \tilde{Y}_{sk}^2 \right) \tilde{Y}_{jk} \right| = \left| \mathbb{E}\tilde{Y}_{ik}^3 \tilde{Y}_{jk} + \mathbb{E}\tilde{Y}_{ik} \tilde{Y}_{jk}^3 + \sum_{s \neq i,j} \mathbb{E}\tilde{Y}_{ik} \tilde{Y}_{sk}^2 \tilde{Y}_{jk} \right|$$

$$(A.44) \leq \mathbb{E}|\tilde{Y}_{ik}|^3 \cdot \mathbb{E}|\tilde{Y}_{jk}| + \mathbb{E}|\tilde{Y}_{ik}| \cdot \mathbb{E}|\tilde{Y}_{jk}|^3 + \sum_{s \neq i,j} \mathbb{E}|\tilde{Y}_{ik}| \cdot \mathbb{E}|\tilde{Y}_{sk}|^2 \cdot \mathbb{E}|\tilde{Y}_{jk}|$$

$$\leq C(\theta^3(p_1 - 2) + 2\theta^2);$$

if i = j,

$$\left| \mathbb{E} \tilde{Y}_{ik} \left( \sum_{s=1}^{p_1} \tilde{Y}_{sk}^2 \right) \tilde{Y}_{jk} \right| = \left| \mathbb{E} \tilde{Y}_{ik}^4 + \sum_{s \neq i} \mathbb{E} \tilde{Y}_{ik}^2 \tilde{Y}_{sk}^2 \right| \le C(\theta^2 (p_1 - 1) + \theta).$$

Then,

$$\sigma_{B}^{2} \triangleq \left\| \sum_{k=1}^{p_{2}} \mathbb{E}B_{k}^{2} \right\| \leq \sum_{k=1}^{p_{2}} \left\| \mathbb{E}B_{k}^{2} \right\| \stackrel{(A.42)}{\leq} \sum_{k=1}^{p_{2}} \left\| \mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top} \right\|$$

$$\leq \sum_{k=1}^{p_{2}} \left( \left\| D(\mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}) \right\| + \left\| \Delta(\mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}) \right\| \right)$$

$$\leq \sum_{k=1}^{p_{2}} \left( \max_{i} \left( \mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top} \right)_{ii} + \left\| \Delta(\mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}) \right\|_{F} \right)$$

$$\stackrel{(A.43)}{\leq} Cp_{2} \left( \theta^{2}p_{1} + \theta + \left\{ \sum_{1 \leq i \neq j \leq p_{1}} \left( \mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top} \right)_{ij}^{2} \right\}^{1/2} \right\}$$

$$\stackrel{(A.44)}{\leq} Cp_{2} \left( \theta^{2}p_{1} + \theta + p_{1} (\theta^{3}p_{1} + \theta^{2}) \right)$$

$$= Cp_{2} (\theta + \theta^{2}p_{1} + \theta^{3}p_{1}^{2}) \leq Cp_{2} \left( \theta + \theta^{3}p_{1}^{2} \right).$$

On the other hand,

$$\sigma_B^2 \ge \max_{1 \le i \le p_1} \left( \sum_{k=1}^{p_2} \mathbb{E}B_k^2 \right)_{ii},$$

where 
$$(\mathbb{E}B_k^2)_{ii} = (\mathbb{E}\tilde{Y}_{.k}\tilde{Y}_{.k}^{\top}\tilde{Y}_{.k}\tilde{Y}_{.k}^{\top})_{ii} - ((\mathbb{E}\tilde{Y}_{.k}\tilde{Y}_{.k}^{\top})^2)_{ii}$$
  

$$= \mathbb{E}\tilde{Y}_{ik}^2 \left(\sum_{s=1}^{p_1} \tilde{Y}_{sk}^2\right) - \sum_{s=1}^{p_1} (\mathbb{E}\tilde{Y}_{ik}\tilde{Y}_{sk})^2$$

$$= \mathbb{E}\tilde{Y}_{ik}^4 + \sum_{s\neq i} \mathbb{E}\tilde{Y}_{ik}^2 \cdot \mathbb{E}\tilde{Y}_{sk}^2 - (\mathbb{E}\tilde{Y}_{ik}^2)^2 - \sum_{s\neq i} (\mathbb{E}\tilde{Y}_{ik})^2 (\mathbb{E}\tilde{Y}_{sk})^2$$

$$\geq \mathbb{E}\tilde{Y}_{ik}^4 - (\mathbb{E}\tilde{Y}_{ik}^2)^2 = \theta \mathbb{E}Y_{ik}^4 - \theta^2 \mathbb{E}Y_{ik}^2 \geq (\theta - \theta^2) \mathbb{E}Y_{ik}^4.$$

Provided that  $\theta \leq 1 - c$  for constant c > 0, we have

(A.46)

$$\sigma_B^2 \ge \max_i \left( \sum_{k=1}^{p_2} \mathbb{E} B_k^2 \right)_{ii} \ge (\theta - \theta^2) \max_i \sum_{k=1}^{p_2} \mathbb{E} Y_{ik}^4 \ge \frac{c\theta}{p_2} \max_i \left( \mathbb{E} \sum_{k=1}^{p_2} Y_{ik}^2 \right)^2$$

$$\ge \frac{c\theta}{p_1^2 p_2} \left( \mathbb{E} \sum_{i=1}^{p_1} \sum_{k=1}^{p_2} Y_{ik}^2 \right)^2 \ge \frac{c\theta}{p_1^2 p_2} (\mathbb{E} \|X\|_F^2)^2 \ge \frac{c\theta r^2}{p_1^2 p_2} \lambda_r^4(X).$$

Next, we give an upper bound for  $|||B_k|||_{\psi_1}$ . Note that

$$||B_{k}|| = ||\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top} - \mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}|| \le ||\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}|| + ||\mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}|| \le ||\tilde{Y}_{\cdot k}||_{2}^{2} + ||\mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}||$$

In particular, we set  $t = C_1 \theta p_1$  for sufficiently large constant  $C_1 > 0$ . Then,

$$\mathbb{E} \exp\left(\|B_{k}\|/t\right) \leq \mathbb{E} \exp\left\{\left(\|\tilde{Y}_{\cdot k}\|_{2}^{2} + \|\mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}\|\right)/t\right\}$$

$$=\mathbb{E} \exp\left(\|\tilde{Y}_{\cdot k}\|_{2}^{2}/t\right) \cdot \exp\left(\|\mathbb{E}\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}\|/t\right)$$

$$\leq \mathbb{E} \exp\left(\|\tilde{Y}_{\cdot k}\|_{2}^{2}/t\right) \cdot \mathbb{E} \exp\left(\|\tilde{Y}_{\cdot k}\tilde{Y}_{\cdot k}^{\top}\|/t\right) \quad \text{(by Jensen's inequality)}$$

$$= \left(\mathbb{E} \exp\left(\|\tilde{Y}_{\cdot k}\|_{2}^{2}/t\right)\right)^{2} = \left(\mathbb{E} \prod_{i=1}^{p_{1}} \exp\left(\tilde{Y}_{ik}^{2}/t\right)\right)^{2}$$

$$= \prod_{i=1}^{p_{1}} \left(\mathbb{E} \exp\left(\tilde{Y}_{ik}^{2}/t\right)\right)^{2} \leq \prod_{i=1}^{p_{1}} \left(\mathbb{E} \exp(0/t)1_{\{R_{ik}=0\}} + \mathbb{E} \exp(Y_{ik}^{2}/t)1_{\{R_{ik}=1\}}\right)^{2}$$

$$\stackrel{\text{Lemma } 8}{\leq} \prod_{i=1}^{p_{1}} \left((1-\theta) + \theta(1+C/t)\right)^{2} = (1+C\theta/t)^{2p_{1}} \leq 1 + C\theta p_{1}/t \leq 1 + C/C_{1} \leq 2,$$

which means

(A.47) 
$$U_B^{(1)} \triangleq \|\|B_k\|\|_{\psi_1} = \inf\{b > 0 : \mathbb{E}\exp(\|B_k\|/b) \le 2\} \le C_1 \theta p_1.$$

Step 2 Next, we derive an upper bound for  $\|\Delta(\tilde{Y}\tilde{Y}^{\top} - \theta^2XX^{\top})\|$  based on the results of the previous step. By the Bernstein-type matrix concentration inequality (c.f., Proposition 2 in Koltchinskii et al. (2011)), (A.45), (A.46), and (A.47), we have

$$\left\| \sum_{k=1}^{p_2} B_k \right\| \le C \max \left\{ \sigma_B \sqrt{\log(p_1)}, U_B^{(1)} \log(p_1) \log \left( \frac{U_B^{(1)}}{\sigma_B / \sqrt{p_2}} \right) \right\}$$

$$\le C \max \left\{ \sqrt{p_2(\theta + \theta^3 p_1^2) \log(p_1)}, \theta p_1 \log(p_1) \log \left( \frac{C\theta^{1/2} p_1^2 p_2}{r \lambda_r^2(X)} \right) \right\}$$

with probability at least  $1 - p_1^{-C}$ . By (A.40) and (A.41), we further have  $P(\mathcal{A}) \geq 1 - p_1^{-C}$ , where  $\mathcal{A}$  is the event such that

$$\mathcal{A} = \left\{ \left\| \Delta \left( \tilde{Y} \tilde{Y}^{\top} - \theta^2 X X^{\top} \right) \right\|$$

$$\leq C \max \left\{ \sqrt{p_2(\theta + \theta^3 p_1^2) \log(p_1)}, \theta p_1 \log(p_1) \log \left( \frac{C \theta^{1/2} p_1^2 p_2}{r \lambda_r^2(X)} \right) \right\} \right\}$$

Step 3 Finally, we finalize the proof by using the robust  $\sin \Theta$  theorem. When the event  $\mathcal{A}$  holds, by Theorem 3, we have the following theoretical guarantee for the HeteroPCA estimator applying to  $\tilde{Y}\tilde{Y}^{\top}$ ,

(A.48)

$$\left\| \sin \Theta(\hat{U}, U) \right\| \leq \frac{C \|\Delta(\tilde{Y}\tilde{Y}^{\top} - \theta^2 X X^{\top})\|}{\lambda_r(\theta^2 X X^{\top})} \wedge 1$$

$$\leq \frac{C \max \left\{ \sqrt{p_2(\theta + \theta^3 p_1^2) \log(p_1)}, \theta p_1 \log(p_1) \log\left(\frac{C \theta^{1/2} p_1^2 p_2}{r \lambda_r^2(X)}\right) \right\}}{\theta^2 \lambda_r^2(X)} \wedge 1.$$

We discuss the bound above in two cases: first, if  $\lambda_r^2(X) \ge \sqrt{p_2 p_1^2/\theta}$ ,

$$\log \left( \frac{C\theta^{1/2}p_1^2p_2}{r\lambda_r^2(X)} \right) \le C \log (p_1p_2);$$

second, if  $\lambda_r^2(X) \leq \sqrt{p_2 p_1^2/\theta}$ , we have

$$\left\|\sin\Theta(\hat{U}, U)\right\| \le 1 \le \frac{C\sqrt{p_2(\theta + \theta^3 p_1^2)\log(p_1)}}{\theta^2 \lambda_r^2(X)} \wedge 1.$$

Thus, if A holds, we always have

$$\left\|\sin\Theta(\hat{U}, U)\right\| \leq \frac{C \max\left\{\sqrt{p_2(\theta + \theta^3 p_1^2)\log(p_1)}, \theta p_1\log(p_1)\log(p_1p_2)\right\}}{\theta^2 \lambda_r^2(X)} \wedge 1.$$

PROOF OF THE CONSISTENCY RESULT IN REMARK 9. If  $||X|| \le C\lambda_r(X)$  and  $||X||_F^2 \ge cp_1p_2$ , we have

$$\lambda_r^2(X) \geq \frac{1}{C} \|X\|^2 \geq \frac{1}{Cr} \sum_{i=1}^r \lambda_i^2(X) \geq \frac{1}{Cr} \|X\|_F^2 \geq \frac{p_1 p_2}{Cr}.$$

If

$$\theta \gg \max \left\{ \frac{r^{2/3} \log^{1/3}(p_1)}{p_1^{2/3} p_2^{1/3}}, \frac{r^2 \log(p_1)}{p_2}, \frac{r \log(p_1) \log(p_1 p_2)}{p_2} \right\},$$

or equivalently

$$\mathbb{E}|\Omega| \gg \max\left\{p_1^{1/3}p_2^{2/3}r^{2/3}\log^{1/3}(p_1), p_1r^2\log(p_1), p_1r\log(p_1)\log(p_1p_2)\right\},\,$$

we have

$$\begin{split} & \mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\| = \mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\| 1_{\mathcal{A}} + \mathbb{E} \left\| \sin \Theta(\hat{U}, U) \right\| 1_{\mathcal{A}^c} \\ & \leq \frac{C \max \left\{ \sqrt{p_2(\theta + \theta^3 p_1^2) \log(p_1)}, \theta p_1 \log(p_1) \log(p_1 p_2) \right\}}{\theta^2 \lambda_r^2(X)} \wedge 1 + \mathbb{P}(\mathcal{A}^c) \\ & \leq \frac{C \max \left\{ \sqrt{p_2(\theta + \theta^3 p_1^2) \log(p_1)}, \theta p_1 \log(p_1) \log(p_1 p_2) \right\}}{C \theta^2 p_1 p_2 / r} \wedge 1 + p_1^{-C} = o(1). \end{split}$$

as 
$$p_1, p_2 \to \infty$$
.

#### APPENDIX B: TECHNICAL LEMMAS

LEMMA 3. Assume that  $E \in \mathbb{R}^{p_1 \times p_2}$  has independent sub-Gaussian entries,  $\operatorname{Var}(E_{ij}) = \sigma_{ij}^2$ ,  $\sigma_C^2 = \max_j \sum_i \sigma_{ij}^2$ ,  $\sigma_R^2 = \max_i \sum_j \sigma_{ij}^2$ ,  $\sigma_{max}^2 = \max_{i,j} \sigma_{ij}^2$ . Assume that

$$||E_{ij}/\sigma_{ij}||_{\psi_2} = \max_{q>1} q^{-1/2} \{ \mathbb{E}(E_{ij}/\sigma_{ij})^q \}^{1/q} \le \kappa.$$

Let  $V \in \mathbb{O}_{p_2,r}$  be a fixed orthogonal matrix. Then

(B.1) 
$$\mathbb{P}(\|EV\| \ge 2(\sigma_C + x)) \le 2 \exp\left(5r - \min\left\{\frac{x^4}{\kappa^4 \sigma_{max}^2 \sigma_C^2}, \frac{x^2}{\kappa^2 \sigma_{max}^2}\right\}\right)$$
,

(B.2) 
$$\mathbb{E}||EV|| \lesssim \sigma_C + \kappa r^{1/4} (\sigma_{max} \sigma_C)^{1/2} + \kappa r^{1/2} \sigma_{max}.$$

PROOF OF LEMMA 3. We first construct  $W \subseteq \mathcal{B}_r = \{w \in \mathbb{R}^r : ||w||_2 \le 1\}$  as the  $\ell_2$  distance  $\varepsilon$ -net in r-dimensional space, such that  $|\mathcal{W}| \le (1 + 2/\varepsilon)^r$  (Vershynin, 2011, Lemma 2.5). Since  $E \in \mathbb{R}^{p_1 \times p_2}$  has independent entries, for each fixed  $w \in \mathcal{W}$ ,  $EVw \in \mathbb{R}^{p_1}$  has independent entries and

$$\operatorname{Var}((EVw)_i) = \sum_{j=1}^{p_2} \operatorname{Var}(E_{ij}) \cdot (Vw)_j^2 \le \sum_{j=1}^{p_2} \sigma_{\max}^2 (Vw)_j^2 \le \sigma_{\max}^2 ||Vw||_2^2 \le \sigma_{\max}^2,$$

$$\sum_{i=1}^{p_1} \operatorname{Var}((EVw)_i) = \sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \operatorname{Var}(E_{ij}) \cdot (Vw)_j^2 \le \sum_{j=1}^{p_2} \sigma_C^2(Vw)_j^2 \le \sigma_C^2.$$

Thus, we can rewrite the centralized  $||EVw||_2^2$  as

$$||EVw||_2^2 - \sum_{i=1}^{p_1} \text{Var}((EVw)_i) = \sum_{i=1}^{p_1} ((EVw)_i^2/\text{Var}((EVw)_i) - 1) \cdot \text{Var}((EVw)_i)$$

Here,

$$\sum_{i=1}^{p_1} \operatorname{Var}((EVw)_i) \le \sigma_C^2, \quad \max_i \operatorname{Var}((EVw)_i) \le \sigma_{\max}^2,$$

$$\sum_{i=1}^{p_1} \operatorname{Var}^2((EVw)_i) \le \sigma_{\max}^2 \sum_i \operatorname{Var}((EVw)_i) \le \sigma_{\max}^2 \sigma_C^2.$$

By Bernstein-type concentration inequality (Vershynin, 2010, Proposition 5.16),

$$\mathbb{P}\left(\|EVw\|_{2}^{2} \ge \sigma_{C}^{2} + t\right) \le 2\exp\left(-\min\left\{\frac{t^{2}}{\kappa^{4}\sigma_{\max}^{2}\sigma_{C}^{2}}, \frac{t}{\kappa^{2}\sigma_{\max}^{2}}\right\}\right).$$

Applying the union bound for all  $w \in \mathcal{W}$ , we obtain

$$\mathbb{P}\left(\max_{w \in \mathcal{W}} \|EVw\|_{2}^{2} \ge \sigma_{C}^{2} + t\right) \le 2\left(1 + 2/\varepsilon\right)^{r} \exp\left(-\min\left\{\frac{t^{2}}{\kappa^{2}\sigma_{C}^{2}\sigma_{\max}^{2}}, \frac{t}{\kappa^{2}\sigma_{\max}^{2}}\right\}\right).$$

Next, suppose  $u^* = \arg\max_{\substack{u \in \mathbb{R}^r \\ \|u\|_2 \leq 1}} \|EVu\|_2$ . By definition of  $\varepsilon$ -net, there exists  $w \in \mathcal{W}$ , such that  $\|u^* - w\|_2 \leq \varepsilon$  and

$$||EV|| = ||EVu^*||_2 \le ||EVw||_2 + ||EV(u^* - w)||_2$$
  
$$\le \varepsilon ||EV|| + \max_{w \in \mathcal{W}} ||EVw||_2.$$

Namely,  $||EV|| \le \max_{w \in \mathcal{W}} ||EVw||_2/(1-\varepsilon)$ . Setting  $\varepsilon = 1/2$ , we have

(B.3) 
$$\mathbb{P}(\|EV\| \ge 2(\sigma_C + x)) \le 2 \exp\left(5r - \min\left\{\frac{x^4}{\kappa^4 \sigma_{\max}^2 \sigma_C^2}, \frac{x^2}{\kappa^2 \sigma_{\max}^2}\right\}\right)$$

which has proved (B.1).

Next, we consider the expectation upper bound. For any  $x \geq 0$ ,  $\mathbb{P}(\|EV\| \geq x) \leq 1$ ; for any  $x \geq 2\sigma_C + 10\kappa\sqrt{r}\sigma_{\max} + 10\kappa r^{1/4}(\sigma_{\max}\sigma_C)^{1/2}$ ,

$$\mathbb{P}\left(\|EV\| \ge x\right) \le 2\exp\left(5\log(r) - \min\left\{\frac{(x/2 - \sigma_C)^4}{\kappa^4 \sigma_{\max}^2 \sigma_C^2}, \frac{(x/2 - \sigma_C)^2}{\kappa^2 \sigma_{\max}^2}\right\}\right)$$

$$\le 2\exp\left(5\log(r) - \frac{(x/2 - \sigma_C)^4}{\kappa^4 \sigma_{\max}^2 \sigma_C^2}\right) + 2\exp\left(5\log(r) - \frac{(x/2 - \sigma_C)^2}{\kappa^2 \sigma_{\max}^2}\right)$$

$$\le 2\exp\left(-\frac{(x/2 - \sigma_C)^4}{2\kappa^4 \sigma_{\max}^2 \sigma_C^2}\right) + 2\exp\left(-\frac{(x/2 - \sigma_C)^2}{2\kappa^2 \sigma_{\max}^2}\right).$$

Thus,

$$\mathbb{E}||EV|| = \int_{0}^{\infty} \mathbb{P}\left(||EV|| \ge x\right) dx$$

$$= \int_{0}^{2\sigma_{C} + 10\kappa\sqrt{r}\sigma_{\max} + 10\kappa r^{1/4}(\sigma_{\max}\sigma_{C})^{1/2}} \mathbb{P}\left(||EV|| \ge x\right) dx$$

$$+ \int_{2\sigma_{C} + 10\kappa\sqrt{r}\sigma_{\max} + 10\kappa r^{1/4}(\sigma_{\max}\sigma_{C})^{1/2}} \mathbb{P}\left(||EV|| \ge x\right) dx$$

$$\le 2\sigma_{C} + 10\kappa\sqrt{r}\sigma_{\max} + 10\kappa r^{1/4}(\sigma_{\max}\sigma_{C})^{1/2}$$

$$+ \int_{0}^{\infty} \left\{ 2\exp\left(-\frac{(x/2)^{4}}{\kappa^{4}\sigma_{\max}^{2}\sigma_{C}^{2}}\right) + 2\exp\left(-\frac{(x/2)^{2}}{\kappa^{2}\sigma_{\max}^{2}}\right) \right\} dx$$

$$\le 2\sigma_{C} + 10\kappa\sqrt{r}\sigma_{\max} + 10\kappa r^{1/4}(\sigma_{\max}\sigma_{C})^{1/2}$$

$$+ 4\kappa(\sigma_{\max}\sigma_{C})^{1/2} \int_{0}^{\infty} e^{-x^{4}} dx + 4\kappa\sigma_{\max}\int_{0}^{\infty} -x^{2} dx$$

$$\le C\left(\sigma_{C} + \kappa r^{1/4}(\sigma_{\max}\sigma_{C})^{1/2} + \kappa\sigma_{\max}\sqrt{r}\right).$$

We thus have finished the proof of (B.2).

LEMMA 4 (Spectral Norm of Projected Random Matrix with independent Sub-exponential Entries). Suppose  $E \in \mathbb{R}^{p_1 \times p_2}$  has independent sub-exponential entries,  $\operatorname{Var}(E_{ij}) = \sigma_{ij}^2$ ,  $\sigma_C^2 = \max_j \sum_i \sigma_{ij}^2$ ,  $\sigma_{max}^2 = \max_{i,j} \sigma_{ij}^2$ . Assume that

$$||E_{ij}/\sigma_{ij}||_{\psi_1} = \max_{q \ge 1} q^{-1} \{ \mathbb{E}(E_{ij}/\sigma_{ij})^q \}^{1/q} \le C.$$

Suppose  $V \in \mathbb{O}_{p_2,r}$  is a fixed orthogonal matrix. Then,

$$\mathbb{E}||EV||^2 \lesssim \sigma_C^2 + r^2 \sigma_{max}^2.$$

PROOF OF LEMMA 4. We divide the proof into four steps.

Step 1 First, we introduce an  $\varepsilon$ -net to reduce the matrix concentration problem to a simpler vector one. Let  $\mathcal{W} \subseteq \mathcal{B}_r = \{w \in \mathbb{R}^r : \|w\|_2 \le 1\}$  be the  $\ell_2$  distance  $\varepsilon$ -net in r-dimensional space, such that  $\varepsilon = 1/2$  and  $|\mathcal{W}| \le (1+2/(1/2))^r = 5^r$  (Vershynin, 2011, Lemma 2.5). Since E is a random matrix with independent entries, for any fixed  $w \in \mathcal{W}$ , the vector EVw has independent entries,

(B.4) 
$$\mathbb{E}||EVw||_{2}^{2} = \sum_{i=1}^{p_{1}} \mathbb{E}(EVw)_{i}^{2} = \sum_{i=1}^{p_{1}} \operatorname{Var}((EVw)_{i})$$
$$= \sum_{i=1}^{p_{1}} \sum_{j=1}^{p_{2}} \operatorname{Var}(E_{ij}) \cdot (Vw)_{j}^{2} \leq \sum_{j=1}^{p_{2}} \sigma_{C}^{2}(Vw)_{j}^{2} = \sigma_{C}^{2}.$$

Step 2 Then we establish the concentration for each entry of EVw, say  $(EVw)_i$ . Denote  $f_{ij} = \sigma_{ij}(Vw)_j$ . We have (B.5)

$$\mathbb{E}(EVw)_i^2 = \operatorname{Var}\left(\sum_{j=1}^{p_2} E_{ij}(Vw)_j\right) = \sum_{j=1}^{p_2} \sigma_{ij}^2(Vw)_j^2 = \sum_{j=1}^{p_2} f_{ij}^2 = ||f_{i\cdot}||_2^2,$$

(B.6) 
$$\max_{i} \|f_{i\cdot}\|_{2}^{2} = \max_{i} \sum_{j=1}^{p_{2}} \sigma_{ij}^{2} (Vw)_{j}^{2} \le \left(\max_{i,j} \sigma_{ij}^{2}\right) \cdot \sum_{j=1}^{p_{2}} (Vw)_{j}^{2} \le \sigma_{\max}^{2},$$

(B.7) 
$$\sum_{i=1}^{p_1} \|f_{i\cdot}\|_2^2 = \sum_{i=1}^{p_1} \sum_{j=1}^{p_2} \sigma_{ij}^2 (Vw)_j^2 \le \sum_{i=1}^{p_2} \sigma_C^2 (Vw)_j^2 \le \sigma_C^2.$$

Note that

$$\sum_{j=1}^{p_2} \frac{E_{ij}}{\sigma_{ij}} \cdot \sigma_{ij}(Vw)_j = \sum_{j=1}^{p_2} E_{ij}(Vw)_j = (EVw)_i.$$

By Bernstein-type concentration inequality (c.f., (Vershynin, 2010, Proposition 5.16)),

$$\mathbb{P}(|(EVw)_{i}| \geq t) = \mathbb{P}\left(\left|\sum_{j=1}^{p_{2}} \frac{E_{ij}}{\sigma_{ij}} \cdot \sigma_{ij}(Vw)_{j}\right| \geq t\right) \\
\leq 2 \exp\left(-c \min\left\{\frac{t^{2}}{\|f_{i\cdot}\|_{2}^{2}}, \frac{t}{\|f_{i\cdot}\|_{\infty}}\right\}\right) \\
\leq 2 \exp\left(-c \min\left\{\frac{t^{2}}{\|f_{i\cdot}\|_{2}^{2}}, \frac{t}{\|f_{i\cdot}\|_{2}}\right\}\right) \\
\leq 2 \exp\left(-c \min\left\{\frac{t}{\|f_{i\cdot}\|_{2}} - \frac{1}{4}, \frac{t}{\|f_{i\cdot}\|_{2}}\right\}\right) \\
\leq 2 \exp\left(\frac{c}{4} - ct/\|f_{i\cdot}\|_{2}\right).$$

Next, we consider  $T_i \triangleq (EVw)_i^2 - ||f_{i\cdot}||_2^2$ ,  $i = 1, ..., p_1$  and aim to establish the tail property of  $T_i$ . Suppose  $C_1$  and  $\tilde{C}$  and two constants to be determined later. Then,

$$\mathbb{E} \exp\left(\frac{|T_{i}|^{1/2}}{C_{1}||f_{i\cdot}||_{2}}\right) = \mathbb{E} \exp\left(\frac{|(EVw)_{i}^{2} - ||f_{i\cdot}||_{2}^{2}|^{1/2}}{C_{1}||f_{i\cdot}||_{2}}\right)$$

$$\leq \mathbb{E} \exp\left(\frac{|(EVw)_{i}| + ||f_{i\cdot}||_{2}}{C_{1}||f_{i\cdot}||_{2}}\right)$$

$$= \int_{0}^{\infty} \frac{\partial}{\partial t} \left(\mathbb{E} \exp\left(\frac{t + ||f_{i\cdot}||_{2}}{C_{1}||f_{i\cdot}||_{2}}\right)\right) \mathbb{P}\left(|(EVw)_{i}| \geq t\right) dt$$

$$\stackrel{\text{(B.8)}}{\leq} \int_{0}^{\tilde{C}||f_{i\cdot}||_{2}} \frac{1}{C_{1}||f_{i\cdot}||_{2}} \exp\left(\frac{t + ||f_{i\cdot}||_{2}}{C_{1}||f_{i\cdot}||_{2}}\right) dt$$

$$+ \int_{\tilde{C}||f_{i\cdot}||_{2}}^{\infty} \frac{1}{C_{1}||f_{i\cdot}||_{2}} \exp\left(\frac{t + ||f_{i\cdot}||_{2}}{C_{1}||f_{i\cdot}||_{2}}\right) 2 \exp\left(\frac{c}{4} - \frac{ct}{||f_{i\cdot}||_{2}}\right) dt$$

$$\leq \frac{\tilde{C}}{C_{1}} \exp\left(\frac{\tilde{C} + 1}{C_{1}}\right) + \int_{\tilde{C}||f_{i\cdot}||_{2}}^{\infty} \frac{2 \exp(\frac{c}{4} + \frac{1}{C_{1}})}{C_{1}||f_{i\cdot}||_{2}} \exp\left(-\left(c - \frac{1}{C_{1}}\right) \frac{t}{||f_{i\cdot}||_{2}}\right) dt$$

$$= \frac{\tilde{C}}{C_{1}} \exp\left(\frac{\tilde{C} + 1}{C_{1}}\right) + \frac{2 \exp\left(\frac{c}{4} + \frac{1}{C_{1}} - \left(c - \frac{1}{C_{1}}\right)\tilde{C}\right)}{cC_{1} - 1}.$$

Let  $\tilde{C} = \sqrt{C_1}$ . We can see for large constant  $C_1$ ,  $\mathbb{E} \exp\left(\frac{|T_i|^{1/2}}{C_1||f_{i\cdot}||_2}\right) \leq 2$ . Then,

$$||T_i||_{\psi_{1/2}} \triangleq \inf \left\{ \alpha > 0 : \mathbb{E} \exp\left(|T_i/\alpha|^{1/2}\right) \le 2 \right\} \le C_1^2 ||f_i||_2^2.$$

Step 3 In this step we establish the concentration inequality for the  $\ell_2$  norm of the vector EVw. Noting that  $\mathbb{E}T_i = 0$ , by the tail inequality for sum of heavy tail random variables (c.f., Lemma 6 in Hao et al. (2018)), we have for any  $q \geq 2$ ,

$$\left(\mathbb{E} \left| \|EVw\|_{2}^{2} - \mathbb{E} \|EVw\|_{2}^{2} \right|^{q}\right)^{1/q} = \left(\mathbb{E} \left| \sum_{i=1}^{p_{1}} T_{i} \right|^{q}\right)^{1/q} \\
\leq C \left(\sqrt{q} \left( \sum_{i=1}^{p_{1}} \|f_{i \cdot}\|_{2}^{4} \right)^{1/2} + q^{2} \left( \sum_{i=1}^{p_{1}} \|f_{i \cdot}\|_{2}^{2q} \right)^{1/q}\right) \\
\leq C \left(\sqrt{q} \left( \sum_{i=1}^{p_{1}} \|f_{i \cdot}\|_{2}^{2} \cdot \max_{i} \|f_{i \cdot}\|_{2}^{2} \right)^{1/2} \\
+ q^{2} \left( \sum_{i=1}^{p_{1}} \|f_{i \cdot}\|_{2}^{2} \cdot \max_{i} \|f_{i \cdot}\|_{2}^{2q-2} \right)^{1/q} \right) \\
3.7)$$

$$\stackrel{\text{(B.6)(B.7)}}{\leq} C\sqrt{q}\sigma_C\sigma_{\max} + Cq^2\sigma_C^{2/q}\sigma_{\max}^{(2q-2)/q} = C\sqrt{q}\sigma_C\sigma_{\max} + Cq^2(\sigma_C/\sigma_{\max})^{2/q}\sigma_{\max}.$$

Set

(B.9) 
$$q = 2(r+1) + \log(\sigma_C/\sigma_{\text{max}}),$$

we have

$$\left(\mathbb{E}\left|\|EVw\|_{2}^{2} - \mathbb{E}\|EVw\|_{2}^{2}\right|^{q}\right)^{1/q} \leq C\sqrt{q}\sigma_{C}\sigma_{\max} + Cq^{2}\sigma_{\max} \triangleq G.$$

By Markov inequality,

$$\mathbb{P}\left(\left|\|EVw\|_{2}^{2} - \mathbb{E}\|EVw\|_{2}^{2}\right| \ge t\right) = \mathbb{P}\left(\left|\|EVw\|_{2}^{2} - \mathbb{E}\|EVw\|_{2}^{2}\right|^{q} \ge t^{q}\right) \\
\le \frac{\mathbb{E}\left|\|EVw\|_{2}^{2} - \mathbb{E}\|EVw\|_{2}^{2}\right|^{q}}{t^{q}} = \frac{G^{q}}{t^{q}}.$$

Step 4 Finally, we apply the  $\varepsilon$ -net technique to derive the upper bound for  $||EV||_2^2 = \max_{||w||_2 \le 1} ||EVw||_2^2$  from the concentration inequality of  $||EVw||_2^2$  with fixed w. Applying the union bound, we have

(B.10) 
$$\mathbb{P}\left(\max_{w \in \mathcal{W}} \left| \|EVw\|_2^2 - \mathbb{E}\|EVw\|_2^2 \right| \ge t \right)$$

$$\leq |\mathcal{W}| \cdot \mathbb{P}\left( \left| \|EVw\|_2^2 - \mathbb{E}\|EVw\|_2^2 \right| \ge t \right) \le \frac{G^q 5^r}{t^q}.$$

Suppose  $u^* = \arg\max_{\substack{u \in \mathbb{R}^r \\ \|u\|_2 \le 1}} \|EVu\|_2$ . By definition of  $\varepsilon$ -net, there exists  $w \in \mathcal{W}$ , such that  $\|u^* - w\|_2 \le 1/2$ . Then,

$$||EV|| = ||EVu^*||_2 \le ||EVw||_2 + ||EV(u^* - w)||_2 \le \max_{w \in \mathcal{W}} ||EVw||_2 + ||EV||/2,$$

which means  $||EV|| \le \max_{w \in \mathcal{W}} ||EVw||_2/(1-1/2) = 2 \max_{w \in \mathcal{W}} ||EVw||_2$ . Therefore,

$$\begin{split} \mathbb{E} \|EV\|_{2}^{2} &\leq 4\mathbb{E} \max_{w \in \mathcal{W}} \|EVw\|_{2}^{2} \leq 4\max_{w \in \mathcal{W}} \left(\mathbb{E} \|EVw\|_{2}^{2} + \left| \|EVw\|_{2}^{2} - \mathbb{E} \|EVw\|_{2}^{2} \right| \right) \\ &\leq 4\sigma_{C}^{2} + 4\int_{0}^{\infty} \mathbb{P} \left( \left| \|EVw\|_{2}^{2} - \mathbb{E} \|EVw\|_{2}^{2} \right| \geq t \right) dt \\ &\leq 4\sigma_{C}^{2} + 4\int_{0}^{5G} 1 \cdot dt + 4\int_{5G}^{\infty} \mathbb{P} \left( \left| \|EVw\|_{2}^{2} - \mathbb{E} \|EVw\|_{2}^{2} \right| \geq t \right) dt \\ &\leq 4\sigma_{C}^{2} + 20G + 4\int_{5G}^{\infty} \frac{G^{q}5^{r}}{t^{q}} dt = 4\sigma_{C}^{2} + 20G + \frac{G(q-1)}{5^{q-1-r}} \\ &\leq 4\sigma_{C}^{2} + 20G + \frac{G(q-1)}{5^{(q-1)/2}} \leq 4\sigma_{C}^{2} + CG \\ &= 4\sigma_{C}^{2} + C \left( r + \log(\sigma_{C}/\sigma_{\max}) \right) \sigma_{C}\sigma_{\max} + C \left( r + \log(\sigma_{C}/\sigma_{\max}) \right)^{2} \sigma_{\max}^{2}. \end{split}$$

Finally, by arithmetic-geometric inequality,

$$(r + \log(\sigma_C/\sigma_{\max}))\sigma_C\sigma_{\max} \le \frac{1}{2}\sigma_C^2 + \frac{1}{2}(r + \log(\sigma_C/\sigma_{\max}))^2\sigma_{\max}^2$$
  
 
$$\lesssim \sigma_C^2 + r^2 + \log^2(\sigma_C/\sigma_{\max})\sigma_{\max}^2,$$

$$\log^2(\sigma_C/\sigma_{\max})\sigma_{\max} \lesssim (\sigma_C/\sigma_{\max})\sigma_{\max} = \sigma_C^2$$

we have

$$\mathbb{E}||EV||_2^2 \lesssim \sigma_C^2 + r^2 \sigma_{\max}^2.$$

The following lemma provides a sharp bound for the operator norm of matrix sparsification.

LEMMA 5. If  $M \in \mathbb{R}^{m_1 \times m_2}$ , rank(M) = r,  $\mathcal{G} \subseteq [m_1] \times [m_2]$ , max<sub>i</sub>  $|\{j : (i,j) \in \mathcal{G}\}| \leq b$ , max<sub>j</sub>  $|\{i : (i,j) \in \mathcal{G}\}| \leq b$ , then we have

$$\|G(M)\| \leq \sqrt{b \wedge r} \|M\|, \quad \|\Gamma(M)\| \leq (\sqrt{b \wedge r} + 1) \|M\|.$$

In particular, if  $M \in \mathbb{R}^{p \times p}$  is any square matrix and  $\Delta(M)$  is the matrix M with diagonal entries set to 0, then

$$||\Delta(M)|| \le 2||M||.$$

Here, the factor "2" in the statement above is sharp in the sense that such a statement does not hold if one replaces it by any smaller value.

PROOF OF LEMMA 5. If  $M \in \mathbb{R}^{m_1 \times m_2}$ , M can be seen as a linear operator from  $\mathbb{R}^{m_2}$  to  $\mathbb{R}^{m_1}$ . Note that

$$||G(M)||_{\infty} \triangleq \max_{x \in \mathbb{R}^{m_2}} \frac{||G(M)x||_{\infty}}{||x||_{\infty}} = \max_{i} \sum_{j=1}^{m_2} |G(M_{ij})| = \max_{i} \sum_{j:(i,j) \in \mathcal{G}} |M_{ij}|$$

$$\leq \max_{i} \sqrt{b} \left( \sum_{j:(i,j) \in \mathcal{G}} |M_{ij}|^2 \right)^{1/2} \leq \sqrt{b} \max_{i} ||M_{i\cdot}||_{2} \leq \sqrt{b} ||M||;$$

$$||G(M)||_{1} \triangleq \max_{x \in \mathbb{R}^{m_2}} \frac{||G(M)x||_{1}}{||x||_{1}} = \max_{j} \sum_{i=1}^{m_1} |G(M_{ij})| = \max_{j} \sum_{i:(i,j) \in \mathcal{G}} |M_{ij}|$$

$$\leq \max_{j} \sqrt{b} \left( \sum_{i:(i,j) \in \mathcal{G}} |M_{ij}|^2 \right)^{1/2} \leq \sqrt{b} \max_{j} ||M_{\cdot j}||_{2} \leq \sqrt{b} ||M||.$$

By Riesz-Thorin interpolation theorem (Katznelson, 2004, Chapter 4, Section  $1.2)^3$ ,

$$||G(M)|| \le (||G(M)||_{\infty} \cdot ||G(M)||_1)^{1/2} \le \sqrt{b}||M||.$$

Since rank(M) = r, we also have

$$||G(M)|| \le ||G(M)||_F \le ||M||_F \le \sqrt{r}||M||.$$

The previous two inequalities yield

$$||G(M)|| \le \sqrt{b \wedge r} ||M||.$$

Finally,

$$\|\Gamma(M)\| = \|M - G(M)\| \le (\sqrt{b \wedge r} + 1)\|M\|.$$

<sup>&</sup>lt;sup>3</sup>Also see https://en.wikipedia.org/wiki/Riesz-Thorin\_theorem.

In particular, note that  $\Delta(M) = M - D(M)$ ,  $||D(M)|| = \max_i |M_{ii}| \le ||M||$ , we have

$$||\Delta M|| = ||M - D(M)|| \le ||M|| + ||D(M)|| \le 2||M||.$$

Finally we provide an example to illustrate that the factor "2" above is sharp. Suppose  $p \geq 2$ ,  $1_p$  is the p-dimensional all-one vector. Set  $M = 1_p 1_p^\top - \frac{p}{2} I_p$ . Then,  $\Delta(M) = 1_p 1_p^\top - I_p$ . Since the eigenvalues of  $1_p 1_p^\top$  are  $\{p,0,\ldots,0\}$ , the eigenvalues of  $(\Delta(M) = 1_p 1_p^\top - I_p)$  and  $(M = 1_p 1_p^\top - \frac{p}{2} \cdot I_p)$  are  $\{p-1,-1,\ldots,-1\}$  and  $\{p/2,-p/2,\ldots,-p/2\}$ , respectively. At this point,

$$\frac{\|\Delta M\|}{\|M\|} = \frac{p-1}{p/2} = 2 - \frac{2}{p}.$$

As  $p \to \infty$ , we can see the statement  $||\Delta M|| \le (2 - \varepsilon)||M||$  does not hold in general for any  $\varepsilon > 0$ .

LEMMA 6. Suppose  $E_{p_1} \subseteq \mathbb{S}^{p_1-1}$ ,  $E_{p_2} \subseteq \mathbb{S}^{p_2-1}$  are  $\varepsilon$ -net in  $p_1$ - and  $p_2$ dimensional spheres,  $\varepsilon < 1/2$ , then for any symmetric matrix  $A \in \mathbb{R}^{p_1 \times p_1}$ and general  $B \in \mathbb{R}^{p_1 \times p_2}$ ,

$$||A|| \le \frac{\max_{v \in E_{p_1}} |v^{\top} A v|}{1 - 2\varepsilon}, \quad ||B|| \le \frac{\max_{u \in E_{p_1}, v \in E_n} u^{\top} B v}{1 - 2\varepsilon}.$$

PROOF OF LEMMA 6. Suppose  $\tilde{v} \in \mathbb{S}^{p_1-1}$  is the eigenvector of A corresponding to the eigenvalue with largest absolute value, then  $\tilde{v}$  satisfies  $\tilde{v}^{\top} A \tilde{v} = ||A||$ . Since  $E_{p_1}$  is an  $\varepsilon$ -net of  $\mathbb{S}^{p_1-1}$ , there exists  $u \in E_{p_1}$  such that  $||u - \tilde{v}|| \leq \varepsilon$ . Thus,

$$\begin{split} \|A\| &= \left| \tilde{v}^{\top} A \tilde{v} \right| \leq \left| \tilde{v}^{\top} A (\tilde{v} - v) \right| + \left| (\tilde{v} - v)^{\top} A v \right| + \left| v^{\top} A v \right| \\ &\leq \|\tilde{v}\|_{2} \cdot \|\tilde{v} - v\|_{2} \cdot \|A\| + \|v\|_{2} \cdot \|\tilde{v} - v\|_{2} \cdot \|A\| + \max_{v \in E_{p_{1}}} \left| v^{\top} A v \right| \\ &\leq 2\varepsilon \|A\| + \max_{v \in E_{p_{1}}} \left| v^{\top} A v \right|, \end{split}$$

which implies  $||A|| \leq \frac{1}{1-2\varepsilon} \max_{v \in E_{p_1}} |v^\top Av|$ . Similarly, suppose  $\bar{u}$  and  $\bar{v} \in \mathbb{S}^{p_1-1}$  are the left and right singular vectors of B corresponding to its largest singular value. Then B satisfies  $\bar{u}^\top B \bar{v}_2 = ||B||$ , and there exists  $u \in E_{p_1}$  and  $v \in E_{p_2}$  such that  $||\bar{u} - u||_2 \leq \varepsilon$ ,  $||\bar{v} - u||_2 \leq \varepsilon$ . Therefore,

$$\begin{split} \|B\| &= \tilde{u}^{\top} B \tilde{v} \leq u^{\top} B v + (\tilde{u} - u)^{\top} B v + \tilde{u}^{\top} B (\tilde{v} - v) \\ &\leq \max_{u \in E_{p_1}, v \in E_{p_2}} u^{\top} B v + \|\tilde{u} - u\|_2 \|B\| \cdot \|v\| + \|\tilde{u}\|_2 \cdot \|B\| \cdot \|\tilde{v} - v\|_2 \\ &\leq 2\varepsilon \|B\| + \max_{u \in E_{p_1}, v \in E_{p_2}} u^{\top} B v, \end{split}$$

which implies  $||B|| \leq \frac{1}{1-2\varepsilon} \max_{u \in E_{p_1}, v \in E_{p_2}} u^{\top} B v$ .

The following technical tool characterizes the spectral and Frobenius norm of projections after SVD. The proof is provided in (Zhang and Xia, 2018, Lemma 6).

LEMMA 7. Suppose  $M, E \in \mathbb{R}^{p_1 \times p_2}$ , rank(M) = r. If  $\hat{U} = \text{SVD}_r(M+E)$  and  $\hat{U}_{\perp}$  is the orthogonal complement of  $\hat{U}$ , then

$$\left\|P_{\hat{U}_\perp}M\right\| \leq 2\|E\|, \quad \left\|P_{\hat{U}_\perp}M\right\|_F \leq 2\min\{\sqrt{r}\|E\|,\|E\|_F\}.$$

The following lemma gives a upper bound for  $\mathbb{E} \exp(X^2/t)$  for sub-Gaussian distributed random variable X.

Lemma 8. Suppose X is a sub-Gaussian distributed random variable such that

$$||X||_{\psi_2} \triangleq \max_{q \ge 1} q^{-1/2} (\mathbb{E}|X|^q)^{1/q} \le B.$$

Then whenever  $t \geq 4eB^2$ , we have

$$\mathbb{E}\exp(X^2/t) \le 1 + \sqrt{\frac{8}{\pi}}eB^2/t.$$

Proof of Lemma 8. If  $t \ge 4eB^2$ ,

$$\mathbb{E} \exp(X^{2}/t)$$
=1 +  $\sum_{k=1}^{\infty} \mathbb{E} \frac{X^{2k}}{t^{k}k!} \le 1 + \sum_{k=1}^{\infty} \frac{(2k)^{k} \cdot B^{2k}}{t^{k} \cdot \sqrt{2\pi}k^{k+.5} \cdot e^{-k}}$  (Stirling's Formula)  

$$\le 1 + \sum_{k=1}^{\infty} \left(\frac{2eB^{2}}{t}\right)^{k} \frac{1}{\sqrt{2\pi k}} \le 1 + \frac{1}{\sqrt{2\pi}} \sum_{k=1}^{\infty} \left(\frac{2eB^{2}}{t}\right)^{k}$$

$$\le 1 + \frac{2eB^{2}/t}{\sqrt{2\pi}(1 - 2eB^{2}/t)} \le 1 + \sqrt{\frac{8}{\pi}}eB^{2}/t.$$

The following lemma gives a simple construction of orthogonal matrix of arbitrary dimension that satisfies incoherence constraint.

LEMMA 9. Suppose  $p \geq r \geq 1$ . There exists a p-by-r matrix Q with orthonormal columns, i.e.,  $Q \in \mathbb{O}_{p,r}$ , such that

$$\max_{1 \le i \le p} \|e_i^\top Q\|_2^2 \le \frac{1}{\lfloor p/r \rfloor}.$$

PROOF. Let  $\alpha = \lfloor p/r \rfloor$ ,  $\beta = p - \alpha r$ . Construct

$$Q = \begin{bmatrix} I_r \\ \vdots \\ I_r \\ I_{\beta} \ 0_{\beta \times (p-\beta)} \end{bmatrix} R,$$

where the  $I_r$  block is repeated for  $\alpha$  times in Q; R is the r-by-r diagonal matrix with first  $\beta$  diagonal entries equal  $1/\sqrt{\alpha+1}$  and the other diagonal entries equal  $1/\sqrt{\alpha}$ . It is easy to check that all columns of Q are orthonormal, i.e.,  $Q \in \mathbb{O}_{p,r}$ . Moreover,

$$\max_{1 \le i \le p} \|e_i^{\top} Q\|_2^2 \le \min_{1 \le i \le r} R_{ii}^2 = \frac{1}{\alpha} = \frac{1}{\lfloor p/r \rfloor}.$$

DEPARTMENT OF STATISTICS UNIVERSITY OF WISCONSIN-MADISON MADISON, WI, 53706.

E-mail: anruzhang@stat.wisc.edu URL: www.stat.wisc.edu/~anruzhang/ DEPARTMENT OF STATISTICS
THE WHARTON SCHOOL
UNIVERSITY OF PENNSYLVANIA
PHILADELPHIA, PA, 19104.
E-MAIL: tcai@wharton.upenn.edu

URL: http://www-stat.wharton.upenn.edu/~tcai/

DEPARTMENT OF STATISTICS AND DATA SCIENCE YALE UNIVERSITY
NEW HEAVEN, CT, 06511.
E-MAIL: yihong.wu@yale.edu

URL: http://www.stat.yale.edu/~yw562