

# INFERENCE FOR HETEROSKEDASTIC PCA WITH MISSING DATA

BY YULING YAN<sup>1,a</sup>, YUXIN CHEN<sup>2,b</sup> AND JIANQING FAN<sup>3,c</sup>

<sup>1</sup>*Institute for Data, Systems, and Society, Massachusetts Institute of Technology, [yulingy@mit.edu](mailto:yulingy@mit.edu)*

<sup>2</sup>*Department of Statistics and Data Science, Wharton School, University of Pennsylvania, [yuxinc@wharton.upenn.edu](mailto:yuxinc@wharton.upenn.edu)*

<sup>3</sup>*Department of Operations Research and Financial Engineering, Princeton University, [cjfan@princeton.edu](mailto:cjfan@princeton.edu)*

This paper studies how to construct confidence regions for principal component analysis (PCA) in high dimension, a problem that has been vastly underexplored. While computing measures of uncertainty for nonlinear/nonconvex estimators is in general difficult in high dimension, the challenge is further compounded by the prevalent presence of missing data and heteroskedastic noise. We propose a novel approach to performing valid inference on the principal subspace, on the basis of an estimator called HeteroPCA (*Ann. Statist.* **50** (2022b) 53–80). We develop nonasymptotic distributional guarantees for HeteroPCA, and demonstrate how these can be invoked to compute both confidence regions for the principal subspace and entrywise confidence intervals for the spiked covariance matrix. Our inference procedures are fully data-driven and adaptive to heteroskedastic random noise, without requiring prior knowledge about the noise levels.

**1. Introduction.** The applications of modern data science frequently ask for succinct representations of high-dimensional data. At the core of this pursuit lies principal component analysis (PCA), which serves as an effective means of dimension reduction and has been deployed across a broad range of domains (Fan et al. (2021), Johnstone and Paul (2018), Jolliffe (1986), Vaswani, Chi and Bouwmans (2018)). In reality, data collection could often be far from ideal—for instance, the acquired data might be subject to random contamination and contain incomplete observations—which inevitably affects the fidelity of PCA and calls for additional care when interpreting the results. To enable informative assessment of the influence of imperfect data acquisition, it would be desirable to accompany the PCA estimators in use with valid measures of uncertainty or “confidence.”

**1.1. Problem formulation.** To allow for concrete and precise studies, the present paper concentrates on a tractable model that captures the effects of random heteroskedastic noise and missing data in PCA. In what follows, we start by formulating the problem, in the hope of facilitating more precise discussions.

*Model.* Imagine we are interested in  $n$  independent random vectors  $\mathbf{x}_j = [x_{1,j}, \dots, x_{d,j}]^\top \in \mathbb{R}^d$  drawn from the following distribution:<sup>1</sup>

$$(1.1) \quad \mathbf{x}_j \stackrel{\text{ind.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{S}^*), \quad 1 \leq j \leq n,$$

where the unknown covariance matrix  $\mathbf{S}^* \in \mathbb{R}^{d \times d}$  is assumed to be rank- $r$  ( $r < n$ ) with eigen-decomposition

$$(1.2) \quad \mathbf{S}^* = \mathbf{U}^* \mathbf{\Lambda}^* \mathbf{U}^{*\top}.$$

---

Received May 2022; revised January 2024.

*MSC2020 subject classifications.* Primary 62H25; secondary 62E17.

*Key words and phrases.* Principal component analysis, confidence regions, missing data, uncertainty quantification, heteroskedastic data, subspace estimation.

<sup>1</sup>All results in this paper continue to hold if the sample vectors are generated such that  $\mathbf{x}_j = \mathbf{U}^*(\mathbf{\Lambda}^*)^{1/2} \mathbf{f}_j$ , where  $\{\mathbf{f}_j\}_{j=1}^n$  are independent *sub-Gaussian* random vectors in  $\mathbb{R}^r$  satisfying  $\mathbb{E}[\mathbf{f}_j] = \mathbf{0}$ ,  $\mathbb{E}[\mathbf{f}_j \mathbf{f}_j^\top] = \mathbf{I}_r$  and  $\|\mathbf{f}_j\|_{\psi_2} = O(1)$ . Here,  $\|\cdot\|_{\psi_2}$  denotes the sub-Gaussian norm (Vershynin (2018)).

Here, the orthonormal columns of  $U^* \in \mathbb{R}^{d \times r}$  constitute the  $r$  leading eigenvectors of  $S^*$ , whereas  $\Lambda^* \in \mathbb{R}^{r \times r}$  is a diagonal matrix whose diagonal entries are composed of the nonzero eigenvalues of  $S^*$ . In other words, these vectors  $\{x_j\}_{1 \leq j \leq n}$  are randomly drawn from a low-dimensional subspace when  $r$  is small. What we have available are partial and randomly corrupted observations of the entries of the above vectors. Specifically, suppose that we only get to observe

$$(1.3) \quad y_{l,j} = x_{l,j} + \eta_{l,j} \quad \text{for all } (l, j) \in \Omega$$

over a subsampled index set  $\Omega \subseteq [d] \times [n]$  (with  $[n] := \{1, \dots, n\}$ ), where  $\eta_{l,j}$  represents the noise that contaminates the observation in this location. Throughout this paper, we focus on the following random sampling and random noise models.

- *Random sampling*: each index  $(l, j)$  is contained in  $\Omega$  independently with probability  $p$ ;
- *Heteroskedastic random noise with unknown variance*: the noise components  $\{\eta_{l,j}\}$  are independently generated sub-Gaussian random variables obeying

$$\mathbb{E}[\eta_{l,j}] = 0, \quad \mathbb{E}[\eta_{l,j}^2] = \omega_l^{*2}, \quad \text{and} \quad \|\eta_{l,j}\|_{\psi_2} = O(\omega_l^*),$$

where  $\{\omega_l^*\}_{1 \leq l \leq d}$  denote the standard deviations that are *a priori* unknown, and  $\|\cdot\|_{\psi_2}$  stands for the sub-Gaussian norm of a random variable (Vershynin (2018)). The noise levels  $\{\omega_l^*\}_{1 \leq l \leq d}$  are allowed to vary across locations, so as to model the so-called *heteroskedasticity* of noise.

This model can be viewed as a generalization of the spiked covariance model (Baik, Ben Arous and Pécché (2005), Bao et al. (2022a), Cai et al. (2021), Donoho, Gavish and Johnstone (2018), Johnstone (2001), Nadler (2008), Paul (2007)) to account for missing data and heteroskedastic noise. With the observed data  $\{y_{l,j} | (l, j) \in \Omega\}$  in hand, can we perform statistical inference on the orthonormal matrix  $U^*$ —which embodies the ground-truth  $r$ -dimensional principal subspace underlying the vectors  $\{x_j\}_{1 \leq j \leq n}$ —and make inference on the underlying covariance matrix  $S^*$ . Mathematically, the task can often be phrased as constructing valid confidence intervals/regions for both  $U^*$  and  $S^*$  based on the incomplete and corrupted observations  $\{y_{l,j} | (l, j) \in \Omega\}$ . Noteworthily, this model is frequently studied in econometrics and financial modeling under the name of factor models (Bai and Wang (2016), Fan, Li and Liao (2021), Fan et al. (2021), Fan and Yao (2017), Gagliardini, Ossola and Scaillet (2020)), and is closely related to the noisy matrix completion problem where we also quantify uncertainty of missing entries (Candès and Plan (2010), Candès and Recht (2009), Chi, Lu and Chen (2019), Keshavan, Montanari and Oh (2010a)).

*Inadequacy of prior works.* While methods for estimating principal subspace are certainly not in shortage (e.g., Balzano, Chi and Lu (2018), Cai et al. (2021), Cai and Zhang (2018), Li et al. (2021), Lounici (2014), Zhang, Cai and Wu (2022b), Zhu, Wang and Samworth (2022)), methods for constructing confidence regions for principal subspace remain vastly underexplored. The fact that the estimators in use for PCA are typically nonlinear and nonconvex presents a substantial challenge in the development of a distributional theory, let alone uncertainty quantification. As some representative recent attempts, Bao, Ding and Wang (2021), Xia (2021) established normal approximations of the distance between the true subspace and its estimate for the matrix denoising setting, while Koltchinskii, Löffler and Nickl (2020) further established asymptotic normality of some debiased estimator for linear functions of principal components. These distributional guarantees pave the way for the development of statistical inference procedures for PCA. However, it is noteworthy that these results required the noise components to either be i.i.d. Gaussian or at least exhibit matching moments (up to the 4th order), which fell short of accommodating heteroskedastic noise. The challenge is further compounded when statistical inference needs to be conducted in the face of missing data, a scenario that is beyond the reach of these prior works.

1.2. *Our contributions.* In light of the insufficiency of prior results, this paper takes a step toward developing data-driven inference and uncertainty quantification procedures for PCA, in the hope of accommodating both heteroskedastic noise and missing data. Our inference procedures are built on an estimator called HeteroPCA recently proposed by Zhang, Cai and Wu (2022b), which is an iterative algorithm in nature and will be detailed in Section 2. The main contributions of this paper are summarized as follows.

- *Distributional theory for PCA and covariance estimation.* We derive, in a nonasymptotic manner, rowwise distributional characterizations of the principal subspace estimate returned by HeteroPCA (see Theorem 1), as well as entrywise distributional guarantees of the estimate for the covariance matrix estimate of  $\{\mathbf{x}_l\}_{1 \leq l \leq n}$  (see Theorem 3). These distributional characterizations take the form of tractable Gaussian approximations centered at the ground truth.
- *Fine-grained confidence regions and intervals.* Our distributional theory in turns allows for construction of rowwise confidence region for the subspace  $\mathbf{U}^*$  (see Algorithm 3 and Theorem 2) as well as entrywise confidence intervals for the matrix  $\mathbf{S}^*$  (see Algorithm 4 and Theorem 4). The proposed inference procedures are fully data-driven and do not require prior knowledge of the noise levels.

Along the way, we have significantly strengthened the estimation guarantees for HeteroPCA in the presence of missing data. It is noteworthy that all of our theory allows the observed data to be highly incomplete and covers heteroskedastic noise, which is previously unavailable.

1.3. *Paper organization.* The remainder of the paper is organized as follows. In Section 2, we introduce the estimation algorithms available in prior literature. Section 3 develops a suite of distributional theory for HeteroPCA and demonstrates how to use it to construct fine-grained confidence regions and confidence intervals for the unknowns; the detailed proofs of our theorems are deferred to the Appendices. In Section 4, we carry out a series of numerical experiments to confirm the validity and applicability of our theoretical findings. Section 5 gives an overview of several related works. Section 6 takes a detour to analyze two intimately related problems, which will then be utilized to establish our main results. We conclude the paper with a discussion of future directions in Section 7. Most of the proof details are deferred to the Appendices.

1.4. *Notation.* Before proceeding, we introduce several notation that will be useful throughout. We let  $f(n) \lesssim g(n)$  or  $f(n) = O(g(n))$  represent the condition that  $|f(n)| \leq Cg(n)$  for some constant  $C > 0$  when  $n$  is sufficiently large; we use  $f(n) \gtrsim g(n)$  to denote  $f(n) \geq C|g(n)|$  for some constant  $C > 0$  when  $n$  is sufficiently large; and we let  $f(n) \asymp g(n)$  indicate that  $f(n) \lesssim g(n)$  and  $f(n) \gtrsim g(n)$  hold simultaneously. The notation  $f(n) \gg g(n)$  (resp.,  $f(n) \ll g(n)$ ) means that there exists some sufficiently large (resp., small) constant  $c_1 > 0$  (resp.,  $c_2 > 0$ ) such that  $f(n) \geq c_1g(n)$  (resp.,  $f(n) \leq c_2g(n)$ ). We also let  $f(n) = o(g(n))$  indicate that  $\lim_{n \rightarrow \infty} f(n)/g(n) = 0$ . For any real number  $a, b \in \mathbb{R}$ , we shall define  $a \wedge b := \min\{a, b\}$  and  $a \vee b := \max\{a, b\}$ .

For any matrix  $\mathbf{M} = [M_{i,j}]_{1 \leq i \leq n_1, 1 \leq j \leq n_2}$ , we let  $\mathbf{M}_{i,\cdot}$  and  $\mathbf{M}_{\cdot,j}$  stand for the  $i$ th row and the  $j$ th column of  $\mathbf{M}$ , respectively. We shall also let  $\|\mathbf{M}\|$ ,  $\|\mathbf{M}\|_F$ ,  $\|\mathbf{M}\|_{2,\infty}$  and  $\|\mathbf{M}\|_\infty$  denote the spectral norm, the Frobenius norm, the  $\ell_{2,\infty}$  norm (i.e.,  $\|\mathbf{M}\|_{2,\infty} := \max_i \|\mathbf{M}_{i,\cdot}\|_2$ ) and the entrywise  $\ell_\infty$  norm ( $\|\mathbf{M}\|_\infty := \max_{i,j} |M_{i,j}|$ ) of  $\mathbf{M}$ , respectively. For any index set  $\Omega$ , the notation  $\mathcal{P}_\Omega(\mathbf{M})$  represents the Euclidean projection of a matrix  $\mathbf{M}$  onto the subspace of matrices supported on  $\Omega$ , and define  $\mathcal{P}_{\Omega^c}(\mathbf{M}) := \mathbf{M} - \mathcal{P}_\Omega(\mathbf{M})$  as well. In addition, we denote by  $\mathcal{P}_{\text{diag}}(\mathbf{G})$  the Euclidean projection of a square matrix  $\mathbf{G}$  onto the subspace of matrices that vanish outside the diagonal, and define  $\mathcal{P}_{\text{off-diag}}(\mathbf{G}) := \mathbf{G} - \mathcal{P}_{\text{diag}}(\mathbf{G})$ . For a nonsingular

matrix  $\mathbf{H} \in \mathbb{R}^{k \times k}$  with SVD  $\mathbf{U}_H \mathbf{\Sigma}_H \mathbf{V}_H^\top$ , we denote by  $\text{sgn}(\mathbf{H})$  the following orthogonal matrix:

$$(1.4) \quad \text{sgn}(\mathbf{H}) := \mathbf{U}_H \mathbf{V}_H^\top.$$

Finally, we denote by  $\mathcal{C}^d$  the set of all convex sets in  $\mathbb{R}^d$ . For any Lebesgue measurable set  $\mathcal{A} \subseteq \mathbb{R}^d$ , we adopt the shorthand notation  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})\{\mathcal{A}\} := \mathbb{P}(\mathbf{z} \in \mathcal{A})$ , where  $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ . Throughout this paper, we let  $\Phi(\cdot)$  (resp.,  $\phi(\cdot)$ ) represent the cumulative distribution function (resp., probability distribution function) of the standard Gaussian distribution. We also denote by  $\chi_k^2$  the chi-square distribution with  $k$  degrees of freedom.

**2. Background: The estimation algorithm HeteroPCA.** In order to conduct statistical inference for PCA, the first step lies in selecting an algorithm to estimate the principal subspace and the covariance matrix of interest, which we discuss in this section. Before continuing, we introduce several useful matrix notation as follows:

$$(2.1a) \quad \mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n},$$

$$(2.1b) \quad \mathbf{Y} := \mathcal{P}_\Omega(\mathbf{X} + \mathbf{N}) \in \mathbb{R}^{d \times n},$$

where  $\mathcal{P}_\Omega$  has been defined in Section 1.4, and  $\mathbf{N} \in \mathbb{R}^{d \times n}$  represents the noise matrix such that the  $(l, j)$ -th entry of  $\mathbf{N}$  is given by  $\eta_{l,j}$ . In other words,  $\mathbf{Y}$  encapsulates all the observed data  $\{y_{l,j} | (l, j) \in \Omega\}$ , with any entry outside  $\Omega$  taken to be zero. If one has full access to the noiseless data matrix  $\mathbf{X}$ , then a natural strategy to estimate  $\mathbf{U}^*$  would be to return the top- $r$  eigenspace of the sample covariance matrix  $n^{-1} \mathbf{X} \mathbf{X}^\top$ , or equivalently, the top- $r$  left singular subspace of  $\mathbf{X}$ . In practice, however, one needs to extract information from the corrupted and incomplete data matrix  $\mathbf{Y}$ .

*A vanilla SVD-based approach.* Given that  $p^{-1} \mathbf{Y} = p^{-1} \mathcal{P}_\Omega(\mathbf{X} + \mathbf{N})$  is an unbiased estimate of  $\mathbf{X}$  (conditional on  $\mathbf{X}$ ), a natural idea that comes into mind is to resort to the top- $r$  left singular subspace of  $p^{-1} \mathbf{Y}$  when estimating  $\mathbf{U}^*$ . This simple procedure is summarized in Algorithm 1.

*An improved iterative estimator: HeteroPCA.* While Algorithm 1 returns reliable estimates of  $\mathbf{U}^*$  and  $\mathbf{S}^*$  in the regime of moderate-to-high signal-to-noise ratio (SNR), it might fail to be effective if either the missing rate  $1 - p$  or the noise levels are too large. To offer a high-level explanation, we find it helpful to compute the expectation of a properly rescaled sample covariance matrix:

$$(2.2) \quad \frac{1}{p^2} \mathbb{E}[\mathbf{Y} \mathbf{Y}^\top | \mathbf{X}] = \mathbf{X} \mathbf{X}^\top + \left(\frac{1}{p} - 1\right) \mathcal{P}_{\text{diag}}(\mathbf{X} \mathbf{X}^\top) + \frac{n}{p} \text{diag}\{[\omega_l^{*2}]_{1 \leq l \leq d}\},$$

where for any vector  $\mathbf{z} = [z_l]_{1 \leq l \leq d}$  we denote by  $\text{diag}(\mathbf{z}) \in \mathbb{R}^{d \times d}$  a diagonal matrix whose  $(l, l)$ -th entry equals  $z_l$ . Here, we rescale the sample covariance matrix by  $p^{-2}$  on the left-hand side, given that  $p^{-1} \mathbf{Y}$  is an unbiased estimate for  $\mathbf{X}$  and, therefore, we expect  $p^{-2} \mathbf{Y} \mathbf{Y}^\top$  to be close to  $\mathbf{X} \mathbf{X}^\top$ . If the sampling rate  $p$  is overly small and/or if the noise is of large

**Algorithm 1** A vanilla SVD-based approach

*Input:* data matrix  $\mathbf{Y}$  (cf. (2.1b)), sampling rate  $p$ , rank  $r$ .

*Compute* the truncated rank- $r$  SVD  $\mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$  of  $p^{-1} \mathbf{Y} / \sqrt{n}$ , where  $\mathbf{U} \in \mathbb{R}^{d \times r}$ ,  $\mathbf{\Sigma} \in \mathbb{R}^{r \times r}$  and  $\mathbf{V} \in \mathbb{R}^{n \times r}$ .

*Output:*  $\mathbf{U}$  as the subspace estimate,  $\mathbf{\Sigma}$  as an estimate of  $(\boldsymbol{\Lambda}^*)^{1/2}$ , and  $\mathbf{S} = \mathbf{U} \mathbf{\Sigma}^2 \mathbf{U}^\top$  as the covariance matrix estimate of  $\mathbf{x}$ .

---

**Algorithm 2** HeteroPCA (by Zhang, Cai and Wu (2022b))

---

*Input:* data matrix  $Y$  (cf. (2.1b)), sampling rate  $p$ , rank  $r$ , maximum number of iterations  $t_0$ .

*Initialization:* set  $G^0 = \frac{1}{np^2} \mathcal{P}_{\text{off-diag}}(YY^\top)$ .

*Updates:* for  $t = 0, 1, \dots, t_0$  do

$$(U^t, \Lambda^t) = \text{eigs}(G^t, r),$$

$$G^{t+1} = \mathcal{P}_{\text{off-diag}}(G^t) + \mathcal{P}_{\text{diag}}(U^t \Lambda^t U^{t\top}) = \frac{1}{np^2} \mathcal{P}_{\text{off-diag}}(YY^\top) + \mathcal{P}_{\text{diag}}(U^t \Lambda^t U^{t\top}).$$

Here, for any symmetric matrix  $G \in \mathbb{R}^{d \times d}$  and  $1 \leq r \leq d$ ,  $\text{eigs}(G, r)$  returns  $(U, \Lambda)$ , where  $U \Lambda U^\top$  is the top- $r$  eigendecomposition of  $G$ .

*Output:*  $U = U^{t_0}$  as the subspace estimate,  $\Sigma = (\Lambda^{t_0})^{1/2}$  as an estimate of  $(\Lambda^*)^{1/2}$  and  $S = U^{t_0} \Lambda^{t_0} U^{t_0\top}$  as the covariance matrix estimate.

---

size but heteroskedastic, then the second and the third terms on the right-hand side of (2.2) might result in significant bias on the diagonal of the matrix  $\mathbb{E}[YY^\top | X]$ , thus hampering the statistical accuracy of the eigenspace of  $p^{-2}YY^\top$  (or equivalently, the left singular space of  $p^{-1}Y$ ) when employed to estimate  $U^*$ . Viewed in this light, a more effective estimator would include procedures that properly handle the diagonal components of  $p^{-2}YY^\top$ .

To remedy this issue, several previous works (e.g., Cai et al. (2021), Florescu and Perkins (2016)) adopted a spectral method with diagonal deletion, which essentially discards any diagonal entry of  $p^{-2}YY^\top$  before computing its top- $r$  eigenspace. However, diagonal deletion comes at a price: while this operation mitigates the significant bias due to heteroskedasticity and missing data, it introduces another type of bias that might be nonnegligible if the goal is to enable efficient fine-grained inference. To address this bias issue, Zhang, Cai and Wu (2022b) proposed an iterative refinement scheme—termed HeteroPCA—that copes with the diagonal entries in a more refined manner. Informally, HeteroPCA starts by computing the rank- $r$  eigenspace of the diagonal-deleted version of  $p^{-2}YY^\top$ , and then alternates between imputing the diagonal entries of  $XX^\top$  and estimating the eigenspace of  $p^{-2}YY^\top$  with the aid of the imputed diagonal. A precise description of this procedure is summarized in Algorithm 2; here,  $\mathcal{P}_{\text{off-diag}}$  and  $\mathcal{P}_{\text{diag}}$  have been defined in Section 1.4.

**3. Distributional theory and inference procedures.** In this section, we augment the HeteroPCA estimator introduced in Section 2 by a suite of distributional theory, and demonstrate how to employ our distributional characterizations to perform inference on both the principal subspace represented by  $U^*$  and the covariance matrix  $S^*$ .

3.1. *Key quantities and assumptions.* Before continuing, we introduce several additional notation and assumptions that play a key role in our theoretical development. Recall that the eigendecomposition of the covariance matrix  $S^* \in \mathbb{R}^{d \times d}$  (see (1.2)) is assumed to be  $U^* \Lambda^* U^{*\top}$ . We assume the diagonal matrix  $\Lambda^*$  to be  $\Lambda^* = \text{diag}\{\lambda_1^*, \dots, \lambda_r^*\}$ , where the diagonal entries are given by the nonzero eigenvalues of  $S^*$  obeying

$$\lambda_1^* \geq \dots \geq \lambda_r^* > 0.$$

The condition number of  $S^*$  is denoted by

$$(3.1) \quad \kappa := \lambda_1^* / \lambda_r^*.$$

We also find it helpful to introduce the square root of  $\Lambda^*$  as follows:

$$(3.2) \quad \Sigma^* = \text{diag}\{\sigma_1^*, \dots, \sigma_r^*\} = (\Lambda^*)^{1/2}, \quad \text{where } \sigma_i^* = (\lambda_i^*)^{1/2}, \quad 1 \leq i \leq r.$$

Furthermore, we introduce an incoherence parameter commonly employed in prior literature (Candès (2014), Chi, Lu and Chen (2019)).

DEFINITION 1 (Incoherence). *The rank- $r$  matrix  $\mathbf{S}^* \in \mathbb{R}^{d \times d}$  defined in (1.2) is said to be  $\mu$ -incoherent if the following condition holds:*

$$(3.3) \quad \|\mathbf{U}^*\|_{2,\infty} \leq \sqrt{\frac{\mu r}{d}}.$$

Here, we recall that  $\|\mathbf{U}^*\|_{2,\infty}$  denotes the largest  $\ell_2$  norm of all rows of the matrix  $\mathbf{U}^*$ .

REMARK 1. When  $\mu$  is small (e.g.,  $\mu \asymp 1$ ), this condition essentially ensures that the energy of  $\mathbf{U}^*$  is nearly evenly dispersed across all of its rows. As a worthy note, the theory developed herein allows the incoherence parameter  $\mu$  to grow with the problem dimension.

In light of a global rotational ambiguity issue (i.e., for any  $r \times r$  rotation matrix  $\mathbf{R}$ , the matrices  $\mathbf{U} \in \mathbb{R}^{d \times r}$  and  $\mathbf{UR} \in \mathbb{R}^{d \times r}$  share the same column space), in general we can only hope to estimate  $\mathbf{U}^*$  up to global rotation (unless additional eigenvalue separation conditions are imposed). Consequently, our theoretical development focuses on characterizing the error distribution  $\mathbf{UR} - \mathbf{U}^*$  of an estimator  $\mathbf{U}$  when accounting for a proper rotation matrix  $\mathbf{R}$ . In particular, we shall pay particular attention to a specific way of rotation as follows:

$$\mathbf{U} \operatorname{sgn}(\mathbf{U}^\top \mathbf{U}^*) - \mathbf{U}^*,$$

where we recall that for any nonsingular matrix  $\mathbf{H} \in \mathbb{R}^{k \times k}$  with SVD  $\mathbf{U}_H \boldsymbol{\Sigma}_H \mathbf{V}_H^\top$ , the matrix  $\operatorname{sgn}(\mathbf{H})$  is defined to be the rotation matrix  $\mathbf{U}_H \mathbf{V}_H^\top$ . This particular choice aligns  $\mathbf{U}$  and  $\mathbf{U}^*$  in the following sense:

$$\operatorname{sgn}(\mathbf{U}^\top \mathbf{U}^*) = \arg \min_{\mathbf{R} \in \mathcal{O}^{r \times r}} \|\mathbf{UR} - \mathbf{U}^*\|_F,$$

where  $\mathcal{O}^{r \times r}$  indicates the set of all  $r \times r$  rotation matrices; see Ma et al. ((2020), Appendix D.2.1). This is often referred to as Wahba’s problem (Wahba (1965)) or the orthogonal Procrustes problem (Schönemann (1966)).

The last assumption is concerned with the noise levels, which are allowed to vary across different locations.

ASSUMPTION 1 (Noise levels). *The noise levels  $\{\omega_i^*\}_{1 \leq i \leq d}$  obey*

$$(3.4) \quad \frac{\omega_{\max}^2}{\omega_{\min}^2} \leq \kappa_\omega \quad \text{with} \quad \omega_{\max} := \max_{1 \leq i \leq d} \omega_i^* \quad \text{and} \quad \omega_{\min} := \min_{1 \leq i \leq d} \omega_i^*.$$

3.2. *Inferential procedure and theory for HeteroPCA.* We are now positioned to investigate how to assess the uncertainty of the estimator HeteroPCA. For simplicity of presentation, we shall abuse some notation (e.g.,  $\boldsymbol{\Sigma}_{U,l}^*$  and  $v_{i,j}^*$ ) whenever it is clear from the context.

3.2.1. *Distributional theory and inference for the principal subspace  $\mathbf{U}^*$ .* In this subsection, we shall begin by establishing a distributional theory for the subspace estimate  $\mathbf{U}$  returned by HeteroPCA (see Theorem 1), followed by a data-driven and provably valid method to construct fine-grained confidence regions for  $\mathbf{U}^*$  (see Algorithm 3 and Theorem 2). We shall also briefly discuss how our results improve upon prior estimation guarantees for HeteroPCA in the presence of missing data.

*Distributional guarantees.* As it turns out, the subspace estimate returned by Algorithm 2 is approximately unbiased and Gaussian under milder conditions, as posited in the following theorem. The general result beyond the case with  $\kappa, \mu, r, \kappa_\omega \asymp 1$  is postponed to Theorem 11 in Appendix B in the Supplementary Material (Yan, Chen and Fan (2024)).

**THEOREM 1.** *Assume that each column of the ground truth  $\mathbf{X}$  (cf. (2.1a)) is independently generated from  $\mathcal{N}(\mathbf{0}, \mathbf{S}^*)$ , and that the sampling set  $\Omega$  follows the random sampling model in Section 1.1. Suppose that  $p < 1 - \delta$  for some arbitrary constant  $0 < \delta < 1$  or  $p = 1$ , and  $\kappa, \mu, r, \kappa_\omega \asymp 1$ . Assume that Assumption 1 holds and  $d \gtrsim \log^5 n$ ,*

$$(3.5a) \quad \frac{\omega_{\max}^2}{p\sigma_r^{*2}} \sqrt{\frac{d}{n}} \lesssim \frac{1}{\log^{7/2}(n+d)}, \quad \frac{\omega_{\max}}{\sigma_r^*} \sqrt{\frac{d}{np}} \lesssim \frac{1}{\log^3(n+d)},$$

$$(3.5b) \quad ndp^2 \gtrsim \log^9(n+d), \quad np \gtrsim \log^7(n+d).$$

Suppose, in addition, that the number of iterations exceeds

$$(3.6) \quad t_0 \gtrsim \log \left[ \left( \frac{\log^2(n+d)}{\sqrt{ndp}} + \frac{\omega_{\max}^2}{p\sigma_r^{*2}} \sqrt{\frac{d}{n}} \log(n+d) + \frac{\log(n+d)}{\sqrt{np}} + \frac{\omega_{\max}}{\sigma_r^*} \sqrt{\frac{d \log(n+d)}{np}} \right)^{-1} \right].$$

Let  $\mathbf{R}$  be the  $r \times r$  rotation matrix  $\mathbf{R} = \text{sgn}(\mathbf{U}^\top \mathbf{U}^*)$ . Then the estimate  $\mathbf{U}$  returned by Algorithm 2 obeys: for all  $1 \leq l \leq d$ ,

$$\sup_{\mathcal{C} \in \mathcal{C}^r} |\mathbb{P}([\mathbf{UR} - \mathbf{U}^*]_{l,\cdot} \in \mathcal{C}) - \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{U,l}^*)\{\mathcal{C}\}| = o(1),$$

where  $\mathcal{C}^r$  represents the set of all convex sets in  $\mathbb{R}^r$ , and

$$(3.7) \quad \begin{aligned} \boldsymbol{\Sigma}_{U,l}^* := & \left( \frac{1-p}{np} \|\mathbf{U}_{l,\cdot}^* \boldsymbol{\Sigma}^*\|_2^2 + \frac{\omega_l^{*2}}{np} \right) (\boldsymbol{\Sigma}^*)^{-2} + \frac{2(1-p)}{np} \mathbf{U}_{l,\cdot}^{*\top} \mathbf{U}_{l,\cdot}^* \\ & + (\boldsymbol{\Sigma}^*)^{-2} \mathbf{U}^{*\top} \text{diag}\{[d_{l,i}^*]_{1 \leq i \leq d}\} \mathbf{U}^* (\boldsymbol{\Sigma}^*)^{-2} \end{aligned}$$

with

$$d_{l,i}^* := \frac{1}{np^2} [\omega_l^{*2} + (1-p) \|\mathbf{U}_{l,\cdot}^* \boldsymbol{\Sigma}^*\|_2^2] [\omega_i^{*2} + (1-p) \|\mathbf{U}_{i,\cdot}^* \boldsymbol{\Sigma}^*\|_2^2] + \frac{2(1-p)^2}{np^2} S_{l,i}^{*2}.$$

Theorem 1 asserts that each row of the estimate  $\mathbf{U}$  returned by HeteroPCA is nearly unbiased and admits a nearly tight Gaussian approximation, whose covariance matrix can be determined via the closed-form expression (3.7). Given that  $\mathbf{UR}$  and  $\mathbf{U}$  represent the same subspace, this theorem delivers a fine-grained rowwise distributional characterization for the estimator HeteroPCA.

Let us briefly mention the key error decomposition behind this theorem, which might help illuminate how Gaussian approximation emerges. Letting  $\mathbf{E} := n^{-1/2}(p^{-1}\mathbf{Y} - \mathbf{X})$  (which captures the randomness from both the noise and random subsampling), we can decompose

$$(3.8) \quad \begin{aligned} \mathbf{UR} - \mathbf{U}^* = & \underbrace{[\mathbf{E}\mathbf{X}^\top + \mathcal{P}_{\text{off-diag}}(\mathbf{E}\mathbf{E}^\top)] \mathbf{U}^* (\boldsymbol{\Sigma}^*)^{-2}}_{=: \mathbf{Z} \quad (\text{first- and second-order approximation})} + \underbrace{[\mathbf{UR} - \mathbf{U}^* - \mathbf{Z}]}_{=: \boldsymbol{\Psi} \quad (\text{residual term})}. \end{aligned}$$

Here,  $\mathbf{Z}$  contains not only a linear mapping of  $\mathbf{E}$  but also a certain quadratic mapping, the latter of which is crucial when coping with the regime  $n \gg d$ . As a consequence of the

central limit theorem (which will be solidified in the proof),  $\mathbf{Z}$  admits the following Gaussian approximation:

$$(3.9) \quad \mathbf{Z}_{l,\cdot} \stackrel{d}{\approx} \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{U,l}^*), \quad 1 \leq l \leq d.$$

At the same time, the  $\ell_2$  norm of the residual term  $\boldsymbol{\Psi}_{l,\cdot}$  is well controlled and provably negligible compared to the corresponding component in  $\mathbf{Z}_{l,\cdot}$ , thus ascertaining the tightness of the advertised Gaussian approximation.

**REMARK 2.** The decomposition (3.8) also sheds light on why our current theory assumes a finite  $\kappa_\omega$  (cf. (3.4)) when conducting statistical inference (which is unnecessary for the task of estimation). Consider, for example, a simple case when (i) there is no missing data ( $p = 1$ ) and (ii) for some  $1 \leq l \leq d$ , one has  $\omega_l^* = 0$  (and hence  $\kappa_\omega = \infty$ ) and  $\|\mathbf{U}_{l,\cdot}^*\|_2 > 0$ . In this case,  $\mathbf{Z}_{l,\cdot} = \mathbf{0}$  since  $\boldsymbol{\Sigma}_{U,l}^* = \mathbf{0}$ , although  $[\mathbf{U}\mathbf{R} - \mathbf{U}^*]_{l,\cdot}$  is in general nonzero. This implies that our Gaussian approximation—and the inference procedure developed based on this approximation—might fall short of efficacy when  $\kappa_\omega = \infty$ .

*Construction of confidence regions for the principal subspace.* With the above distributional theory in place, we are well equipped to construct fine-grained confidence regions for  $\mathbf{U}^*$ , provided that the covariance matrix  $\boldsymbol{\Sigma}_{U,l}^*$  can be estimated in a faithful manner. In Algorithm 3, we propose a procedure to estimate  $\boldsymbol{\Sigma}_{U,l}^*$ , which in turn allows us to build confidence regions. As before, our estimator for  $\boldsymbol{\Sigma}_{U,l}^*$  can be viewed as a sort of “plug-in” method in accordance with the expression (3.7).

The following theorem confirms the validity of the proposed inference procedure when  $\kappa, \mu, r, \kappa_\omega \asymp 1$ . The more general case will be studied in Theorem 12 in Appendix B in the Supplementary Material (Yan, Chen and Fan (2024)).

**Algorithm 3** Confidence regions for  $\mathbf{U}_{l,\cdot}^*$  ( $1 \leq l \leq d$ ) based on HeteroPCA

*Input:* output  $(\mathbf{U}, \boldsymbol{\Sigma}, \mathbf{S})$  of Algorithm 2, sampling rate  $p$ , coverage level  $1 - \alpha$ .

*Compute* estimates of the noise levels  $\{\omega_l^*\}_{1 \leq l \leq d}$  as follows:

$$\omega_l^2 := \frac{\sum_{j=1}^n y_{l,j}^2 \mathbb{1}_{(l,j) \in \Omega}}{\sum_{j=1}^n \mathbb{1}_{(l,j) \in \Omega}} - S_{l,l} \quad \text{for all } 1 \leq l \leq d.$$

*Compute* an estimate of  $\boldsymbol{\Sigma}_{U,l}^*$  (cf. (3.7)) as follows:

$$\begin{aligned} \boldsymbol{\Sigma}_{U,l} &:= \left( \frac{1-p}{np} \|\mathbf{U}_{l,\cdot} \boldsymbol{\Sigma}\|_2^2 + \frac{\omega_l^2}{np} \right) \boldsymbol{\Sigma}^{-2} + \frac{2(1-p)}{np} \mathbf{U}_{l,\cdot}^\top \mathbf{U}_{l,\cdot} \\ &\quad + (\boldsymbol{\Sigma})^{-2} \mathbf{U}^\top \text{diag}\{[d_{l,i}]_{1 \leq i \leq d}\} \mathbf{U} (\boldsymbol{\Sigma})^{-2}, \end{aligned}$$

where

$$d_{l,i} := \frac{1}{np^2} [\omega_l^2 + (1-p) \|\mathbf{U}_{l,\cdot} \boldsymbol{\Sigma}\|_2^2] [\omega_i^2 + (1-p) \|\mathbf{U}_{i,\cdot} \boldsymbol{\Sigma}\|_2^2] + \frac{2(1-p)^2}{np^2} S_{l,i}^2.$$

*Compute* the  $(1 - \alpha)$ -quantile  $\tau_{1-\alpha}$  of  $\chi_r^2$  and construct a Euclidean ball:

$$\mathcal{B}_{1-\alpha} := \{\mathbf{z} \in \mathbb{R}^r : \|\mathbf{z}\|_2^2 \leq \tau_{1-\alpha}\}.$$

*Output* the  $(1 - \alpha)$ -confidence region

$$\mathcal{CR}_{U,l}^{1-\alpha} := \mathbf{U}_{l,\cdot} + (\boldsymbol{\Sigma}_{U,l})^{1/2} \mathcal{B}_{1-\alpha} = \{\mathbf{U}_{l,\cdot} + (\boldsymbol{\Sigma}_{U,l})^{1/2} \mathbf{z} : \mathbf{z} \in \mathcal{B}_{1-\alpha}\}.$$



**THEOREM 2.** *Suppose that the conditions of Theorem 1 hold. Then there exists a  $r \times r$  rotation matrix  $\mathbf{R} = \text{sgn}(\mathbf{U}^\top \mathbf{U}^\star)$  such that the confidence regions  $\text{CR}_{U,l}^{1-\alpha}$  ( $1 \leq l \leq d$ ) computed in Algorithm 3 obey*

$$\sup_{1 \leq l \leq d} |\mathbb{P}(\mathbf{U}_{l,\cdot}^\star, \mathbf{R}^\top \in \text{CR}_{U,l}^{1-\alpha}) - (1 - \alpha)| = o(1).$$

In words, Theorem 2 uncovers that: a valid ground-truth subspace representation is contained—in a rowwise reliable manner—within the confidence regions  $\text{CR}_{U,l}^{1-\alpha}$  ( $1 \leq l \leq d$ ) we construct. In the special case with  $r = 1$ , this result leads to valid entrywise confidence intervals for the principal component.

*Interpretations and implications.* We now take a moment to interpret the conditions required in Theorem 1 and Theorem 2, and discuss some appealing attributes of our methods. As before, the discussion below focuses on the scenario where  $\mu, \kappa, r, \kappa_\omega \asymp 1$  for the sake of simplicity.

- *Missing data.* Both theorems accommodate the case when a large fraction of data are missing, namely they cover the range

$$p \geq \tilde{\Omega}\left(\frac{1}{n \wedge \sqrt{nd}}\right)$$

for both distributional characterizations and confidence region construction using HeteroPCA. In particular, if  $n \gg d$ , then the sampling rate  $p$  only needs to exceed

$$p \geq \tilde{\Omega}\left(\frac{1}{\sqrt{nd}}\right);$$

this range can include some sampling rate much smaller than  $1/d$  (with  $d$  the ambient dimension of each sample vector), and cannot be improved in general according to (Cai et al. ((2021), Theorem 3.4)).

- *Tolerable noise levels.* The noise condition required in both Theorem 1 and Theorem 2 is given by

$$\omega_{\max}^2 \leq \tilde{O}\left(\left(\frac{n}{d} \wedge \sqrt{\frac{n}{d}}\right) p \sigma_r^{\star 2}\right).$$

Note that when  $\kappa, \mu, r \asymp 1$ , the variance obeys

$$\max_{(l,j) \in \Omega} \text{var}(x_{l,j}) = \max_{l \in [d]} S_{l,l}^\star \asymp \max_{l \in [d]} \|\mathbf{U}_{l,\cdot}^\star\|_2^2 \sigma_1^{\star 2} \asymp \frac{1}{d} \sigma_1^{\star 2}.$$

This implies that: when  $p \geq \tilde{\Omega}(1/(n \wedge \sqrt{nd}))$ , our tolerable entrywise noise level  $\omega_{\max}^2$  is allowed to be significantly (i.e.,  $\tilde{\Omega}(np \wedge \sqrt{ndp^2})$  times) larger than the largest variance of  $x_{l,j}$  for all  $(l, j) \in \Omega$ , thereby accommodating a wide range of noise levels.

- *Adaptivity to heteroskedasticity and unknown noise levels.* Our proposed inferential procedure is fully data-driven: it is automatically adaptive to unknown heteroskedastic noise, without requiring prior knowledge of the noise levels.

*Comparison with prior estimation theory.* While the main purpose of the current paper is to enable efficient statistical inference for the principal subspace, our theory (see Lemmas 18 and 19 in the Supplementary Material (Yan, Chen and Fan (2024))) also enables improved estimation guarantees compared to prior works.

- Recall that the estimation algorithm HeteroPCA was originally proposed and studied by Zhang, Cai and Wu (2022b). Our results broaden the sample size range supported by their theory. More specifically, note that Zhang, Cai and Wu ((2022b), Theorem 6 and Remark 10) requires the sampling rate  $p$  to satisfy

$$ndp \gtrsim \max\{d^{1/3}n^{2/3}, d\} \text{polylog}(n, d)$$

in order to guarantee consistent estimation, while our theoretical guarantees only require

$$ndp \gtrsim \max\{\sqrt{nd}, d\} \text{polylog}(n, d).$$

When  $n \gg d$ , the sample size requirement in Zhang, Cai and Wu (2022b) is  $(n/d)^{1/6}$  times more stringent than the one imposed in our theory.

- Let us discuss the advantage of HeteroPCA compared to the diagonal-deleted spectral method studied in Cai et al. ((2021, Algorithms 1 and 3). Due to diagonal deletion, there is an additional bias term (see the last term  $\mu_{\text{ce}K_{\text{ce}R}}/d$  in equation (4.16) in Cai et al. (2021)), which turns out to negatively affect our capability of performing inference. In contrast, HeteroPCA eliminates this bias term by means of successive refining, thus facilitating the subsequent inference stage.

3.2.2. *Distributional theory and inference for the covariance matrix  $S^*$ .* As it turns out, the above distributional theory for  $U^*$  further hints at how to perform statistical inference for the covariance matrix  $S^*$ . In the sequel, we shall first develop an entrywise distributional theory for the estimate  $S$  returned by HeteroPCA (see Theorem 3), followed by a data-driven inference procedure to conduct entrywise confidence intervals for  $S^*$  (see Algorithm 4 and Theorem 4).

*Entrywise distributional guarantees.* We now focus attention on characterizing the distribution of the  $(i, j)$ -th entry of  $S$  returned by Algorithm 2, which in turn suggests how to construct entrywise confidence intervals for  $S^*$ . Before proceeding, let us define a set of variance parameters  $\{v_{i,j}^*\}_{1 \leq i, j \leq d}$  which, as we shall demonstrate momentarily, correspond to the (approximate) variance of the entries of  $S$ .

- For any  $1 \leq i, j \leq d$  obeying  $i \neq j$ , we define

$$\begin{aligned} v_{i,j}^* := & \frac{2-p}{np} S_{i,i}^* S_{j,j}^* + \frac{4-3p}{np} S_{i,j}^{*2} + \frac{1}{np} (\omega_i^{*2} S_{j,j}^* + \omega_j^{*2} S_{i,i}^*) \\ & + \frac{1}{np^2} \sum_{k=1}^d \{[\omega_i^{*2} + (1-p)S_{i,i}^*][\omega_k^{*2} + (1-p)S_{k,k}^*] \\ (3.10) \quad & + 2(1-p)^2 S_{i,k}^{*2}\} (U_{k,\cdot}^* \cdot U_{j,\cdot}^{*\top})^2 \\ & + \frac{1}{np^2} \sum_{k=1}^d \{[\omega_j^{*2} + (1-p)S_{j,j}^*][\omega_k^{*2} + (1-p)S_{k,k}^*] \\ & + 2(1-p)^2 S_{j,k}^{*2}\} (U_{k,\cdot}^* \cdot U_{i,\cdot}^{*\top})^2. \end{aligned}$$

- For any  $1 \leq i \leq d$ , we set

$$\begin{aligned} v_{i,i}^* := & \frac{12-9p}{np} S_{i,i}^{*2} + \frac{4}{np} \omega_i^{*2} S_{i,i}^* \\ (3.11) \quad & + \frac{4}{np^2} \sum_{k=1}^d \{[\omega_i^{*2} + (1-p)S_{i,i}^*][\omega_k^{*2} + (1-p)S_{k,k}^*] \\ & + 2(1-p)^2 S_{i,k}^{*2}\} (U_{k,\cdot}^* \cdot U_{i,\cdot}^{*\top})^2. \end{aligned}$$

We are now positioned to present our distributional theory for the scenario where  $\kappa, \mu, r, \kappa_\omega \asymp 1$ , with the more general version deferred to Theorem 13 in Appendix B in the Supplementary Material (Yan, Chen and Fan (2024)). Here and throughout,  $S_{i,j}$  (resp.,  $S_{i,j}^*$ ) represents the  $(i, j)$ -th entry of the matrix  $S$  (resp.,  $S^*$ ).

**THEOREM 3.** *Suppose that  $p < 1 - \delta$  for some arbitrary constant  $0 < \delta < 1$  or  $p = 1$ , and  $\kappa, \mu, r, \kappa_\omega \asymp 1$ . Consider any  $1 \leq i, j \leq d$ . Assume that  $U^*$  is  $\mu$ -incoherent and satisfies the following condition:*

$$(3.12) \quad \begin{aligned} & \|U_{i,\cdot}^*\|_2 + \|U_{j,\cdot}^*\|_2 \\ & \gtrsim \left[ \frac{\omega_{\max}}{\sigma_r^*} \sqrt{\frac{d \log^5(n+d)}{np}} + \frac{\omega_{\max}^2}{p\sigma_r^{*2}} \sqrt{\frac{d \log^5(n+d)}{n}} + \sqrt{\frac{\log^7(n+d)}{ndp^2}} \right] \sqrt{\frac{1}{d}}. \end{aligned}$$

In addition, suppose that Assumption 1 holds, and

$$\begin{aligned} d & \gtrsim \log^5 n, & np & \gtrsim \log^7(n+d), & ndp^2 & \gtrsim \log^7(n+d), \\ \frac{\omega_{\max}}{\sigma_r^*} \sqrt{\frac{d}{np}} & \lesssim \frac{1}{\log^3(n+d)}, & \frac{\omega_{\max}^2}{p\sigma_r^{*2}} \sqrt{\frac{d}{n}} & \lesssim \frac{1}{\log^{7/2}(n+d)}. \end{aligned}$$

Assume that the number of iterations satisfies (3.6). Then the matrix  $S$  computed by Algorithm 2 obeys

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P} \left( \frac{S_{i,j} - S_{i,j}^*}{\sqrt{v_{i,j}^*}} \leq t \right) - \Phi(t) \right| = o(1),$$

where  $\Phi(\cdot)$  denotes the CDF of the standard Gaussian distribution.

In words, the above theorem indicates that under the conditions in Theorem 1, if the sum of the  $\ell_2$  norm of the rows  $U_{i,\cdot}^*$  and  $U_{j,\cdot}^*$  are not exceedingly small, then the estimation error  $S_{i,j} - S_{i,j}^*$  of HeteroPCA is approximately a zero-mean Gaussian with variance  $v_{i,j}^*$ .

**REMARK 3.** When it comes to inference for  $S_{i,j}^*$ , our theorems impose the following condition (cf. (3.12)):

$$\|U_{i,\cdot}^*\|_2 + \|U_{j,\cdot}^*\|_2 \geq \tilde{\Omega} \left( \frac{1}{\sqrt{ndp^2}} + \frac{\omega_{\max}}{\sigma_r^*} \sqrt{\frac{d}{np}} + \frac{\omega_{\max}^2}{p\sigma_r^{*2}} \sqrt{\frac{d}{n}} \right) \cdot \sqrt{\frac{1}{d}} \|U^*\|_F.$$

Note that the typical  $\ell_2$  norm of a row of  $U^*$  is  $\|U^*\|_F/\sqrt{d}$  when the energy is uniformly spread out across all rows. This means that under our sampling rate condition, our results allow  $\|U_{i,\cdot}^*\|_2 + \|U_{j,\cdot}^*\|_2$  to be much smaller than its typical size. As it turns out, a lower bound on  $\|U_{i,\cdot}^*\|_2 + \|U_{j,\cdot}^*\|_2$  might be necessary for  $S_{i,j} - S_{i,j}^*$  to be approximately Gaussian. Consider, for example, the case when  $\|U_{i,\cdot}\|_2 = \|U_{j,\cdot}\|_2 = 0$ . It can be seen from our analysis that

$$S_{i,j} - S_{i,j}^* \approx Z_{i,\cdot} \Sigma^* Z_{j,\cdot}^\top + A_{i,j},$$

where  $Z_{i,\cdot}, Z_{j,\cdot}$  and  $A_{i,j}$  are all (approximately) Gaussian. This means that  $S_{i,j} - S_{i,j}^*$  might not follow the (approximate) Gaussian distribution claimed in Theorem 3 if  $\|U_{i,\cdot}^*\|_2 + \|U_{j,\cdot}^*\|_2$  is too small.

**Algorithm 4** Confidence intervals for  $S_{i,j}^*$  ( $1 \leq i, j \leq d$ ) based on HeteroPCA

*Input:* output  $(\mathbf{U}, \Sigma, \mathbf{S})$  of Algorithm 2, sampling rate  $p$ , coverage level  $1 - \alpha$ .

*Compute* estimates of the noise level  $\omega_l^*$  as follows:

$$\omega_l^2 := \frac{\sum_{j=1}^n y_{l,j}^2 \mathbb{1}_{(l,j) \in \Omega}}{\sum_{j=1}^n \mathbb{1}_{(l,j) \in \Omega}} - S_{l,l}.$$

*Compute* an estimate of  $v_{i,j}^*$  (cf. (3.10) or (3.11)) as follows: if  $i \neq j$ , then

$$\begin{aligned} v_{i,j} := & \frac{2-p}{np} S_{i,i} S_{j,j} + \frac{4-3p}{np} S_{i,j}^2 + \frac{1}{np} (\omega_i^2 S_{j,j}^* + \omega_j^2 S_{i,i}^*) \\ & + \frac{1}{np^2} \sum_{k=1}^d \{ [\omega_i^2 + (1-p) S_{i,i}] [\omega_k^2 + (1-p) S_{k,k}] + 2(1-p)^2 S_{i,k}^2 \} (\mathbf{U}_k, \mathbf{U}_{j,\cdot}^\top)^2 \\ & + \frac{1}{np^2} \sum_{k=1}^d \{ [\omega_j^2 + (1-p) S_{j,j}] [\omega_k^2 + (1-p) S_{k,k}] + 2(1-p)^2 S_{j,k}^2 \} (\mathbf{U}_k, \mathbf{U}_{i,\cdot}^\top)^2. \end{aligned}$$

If  $i = j$ , then

$$\begin{aligned} v_{i,i} := & \frac{12-9p}{np} S_{i,i}^2 + \frac{4}{np} \omega_i^2 S_{i,i} \\ & + \frac{4}{np^2} \sum_{k=1}^d \{ [\omega_i^2 + (1-p) S_{i,i}] [\omega_k^2 + (1-p) S_{k,k}] + 2(1-p)^2 S_{i,k}^2 \} (\mathbf{U}_k, \mathbf{U}_{i,\cdot}^\top)^2. \end{aligned}$$

*Output* the  $(1 - \alpha)$ -confidence interval

$$\text{CI}_{i,j}^{1-\alpha} := [S_{i,j} \pm \Phi^{-1}(1 - \alpha/2) \sqrt{v_{i,j}}].$$

*Construction of entrywise confidence intervals.* The distributional characterization in Theorem 3 enables valid construction of entrywise confidence intervals for  $\mathbf{S}^*$ , as long as we can obtain reliable estimate of the variance  $v_{i,j}^*$ . In what follows, we come up with an algorithm—as summarized in Algorithm 4—that attempts to estimate  $v_{i,j}^*$  and build confidence intervals in a data-driven manner, as confirmed by the following theorem for the scenario with  $\kappa, \mu, r, \kappa_\omega \asymp 1$ . The more general result is postponed to Theorem 14 in Appendix B in the Supplementary Material (Yan, Chen and Fan (2024)).

**THEOREM 4.** *Suppose that the conditions in Theorem 3 hold. Assume that  $ndp^2 \gtrsim \log^8(n+d)$ . Then the confidence interval computed in Algorithm 4 obeys*

$$\mathbb{P}(S_{i,j}^* \in \text{CI}_{i,j}^{1-\alpha}) = 1 - \alpha + o(1).$$

Compared with Cai et al. ((2021), Corollary 2), we can see that when consistent estimation is possible—namely, under the sampling rate condition  $p \geq \tilde{\Omega}((n \wedge \sqrt{nd})^{-1})$  and the noise conditions  $\omega_{\max} \leq \tilde{\Omega}((\sqrt{n/d} \wedge \sqrt[4]{n/d}) \sqrt{p} \sigma_r^*)$ —it is plausible to construct fine-grained confidence interval for  $S_{i,j}^*$ , provided that the size of  $\|\mathbf{U}_{i,\cdot}^*\|_2 + \|\mathbf{U}_{j,\cdot}^*\|_2$  is not exceedingly small.

**REMARK 4.** Before concluding this subsection, we note that the sampling rate  $p$  might be unknown *a priori* in practice. If this is the case, then one plausible strategy is to replace  $p$

in Algorithms 2, 3 and 4 with the following empirical estimate:

$$\widehat{p} = \frac{\sum_{l=1}^d \sum_{j=1}^n \mathbb{1}\{(l, j) \in \Omega\}}{nd}.$$

In view of the standard concentration results  $\widehat{p} = (1 + o(1))p$ , it is straightforward to verify that all of these inference procedure and the accompanying theory remain valid. We omit the details for the sake of brevity.

3.3. *A glimpse of key technical ingredients.* Let us take a moment to highlight several technical ingredients of the current theory, which might be applicable to other high-dimensional statistical problems beyond the analysis of HeteroPCA.

*Second-order perturbation theory for principal subspace.* At the core of our analysis lies a “second-order” perturbation theory tailored to general subspace estimation problems, to be presented in Section 6. More concretely, we establish a second-order expansion of the subspace perturbation error (see Theorem 5) that makes explicit the following two parts: (i) nearly tight first- and second-order terms, which can be expressed succinctly as linear and quadratic mappings of the perturbation matrix; (ii) the remaining higher-order terms that are provably negligible. Given that HeteroPCA is an iterative algorithm, developing such a refined perturbation theory for HeteroPCA becomes substantially more challenging than the vanilla SVD-based approach. Our refined perturbation theory allows us to tighten prior estimation theory (e.g., the Davis–Kahan  $\sin \Theta$  Theorem (Davis and Kahan (1970)) or recent  $\ell_{2,\infty}$ -type perturbation bounds (Abbe et al. (2020), Cai et al. (2021), Chen et al. (2021))), the latter of which focused mainly on providing orderwise estimation error bounds.

*Fine-grained distributional characterizations for the principal subspace  $U^*$ .* As alluded to previously, we establish the distributional characterization for the principal subspace (i.e., Theorem 1) based on a key error decomposition

$$(3.13) \quad UR - U^* = \underbrace{[EX^\top + P_{\text{off-diag}}(EE^\top)]U^*(\Sigma^*)^{-2}}_{=: Z \quad (\text{first- and second-order approximation})} + \underbrace{[UR - U^* - Z]}_{=: \Psi \quad (\text{residual term})},$$

where  $E := n^{-1/2}(p^{-1}Y - X)$ . For each  $l \in [d]$ , the multivariate Berry–Esseen theorem reveals the approximate Gaussianity of  $Z_{l,\cdot}$ , while at the same time, our second-order perturbation theory (cf. Theorem 5) ensures that  $\Psi_{l,\cdot}$  is stochastically dominated by  $Z_{l,\cdot}$ . Additionally, rather than providing general  $\ell_{2,\infty}$  bounds (as in the prior work Cai et al. (2021)), our proof relies crucially on more delicate row-dependent error control (so that the size of  $\Psi_{l,\cdot}$  is carefully bounded in accordance with the  $l$ th row of  $U^*$  and  $N$ ).

*Entrywise distributional characterizations when estimating the covariance matrix  $S^*$ .* Moving one step further, we derive the following key error decomposition w.r.t. the covariance matrix  $S^*$ :

$$(3.14) \quad S - S^* = \underbrace{U^*\Sigma^{*2}Z^\top + Z\Sigma^{*2}U^*}_{=: W} + \underbrace{n^{-1}XX^\top - S^*}_{=: A} + \underbrace{[S - S^* - W - A]}_{=: \Phi \quad (\text{residual term})},$$

where  $Z$  is defined in (3.13) and approximately Gaussian. Here,  $W$  serves as the main component as induced by the subspace estimation error,  $A$  indicates the discrepancy between the empirical covariance (using clean and fully observed data) and the true covariance, whereas  $\Phi$  is some higher-order term that is provably negligible in a strong entrywise sense. This in turn allows us to pin down tight entrywise distributional characterizations for  $S - S^*$ .

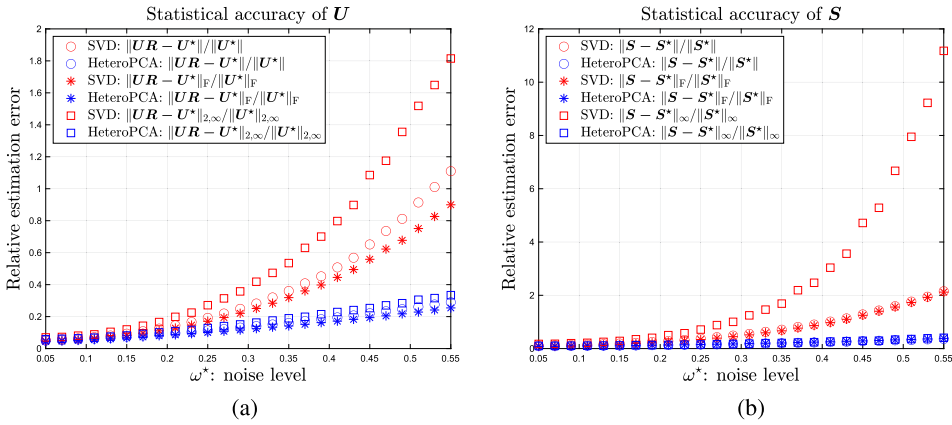


FIG. 1. The relative estimation error of  $U$  and  $S$  returned by both SVD-based approach (cf. Algorithm 1) and HeteroPCA (cf. Algorithm 2) over different noise level  $\omega^*$ . (a) Relative estimation errors of  $UR - U^*$  measured by  $\|\cdot\|$ ,  $\|\cdot\|_F$  and  $\|\cdot\|_{2,\infty}$  vs. the noise level  $\omega^*$ ; (b) Relative estimation errors of  $S - S^*$  measured by  $\|\cdot\|$ ,  $\|\cdot\|_F$  and  $\|\cdot\|_{\infty}$  vs. the noise level  $\omega^*$ . The results are reported over 200 independent trials for  $r = 3$  and  $p = 0.6$ .

#### 4. Numerical experiments.

*Setup.* This section conducts a series of numerical experiments to validate our distributional and inference theory developed in Section 3. Throughout this section, unless otherwise noted, we fix the dimension to be  $d = 100$  and the number of sample vectors to be  $n = 2000$ , and we generate the covariance matrix as  $S^* = U^*U^{*\top}$  with  $U^* \in \mathbb{R}^{n \times r}$  being a random orthonormal matrix following the Haar distribution over the Grassmann manifold  $G_{d,r}$  (Vershynin (2018), Section 5.2.6). In each Monte Carlo trial, the observed data are produced according to the model described in Section 1.1. For the purpose of introducing heteroskedasticity, we will introduce a parameter  $\omega^*$  that controls the noise level: in each independent trial, each noise level  $\omega_l^*$  ( $1 \leq l \leq d$ ) is independently drawn from  $\text{Uniform}[0.1\omega^*, 2\omega^*]$ ; the random noise component  $\eta_{l,j}$  is then drawn from  $\mathcal{N}(0, \omega_l^{*2})$  independently for every  $l \in [d]$  and  $j \in [n]$ .

*Superiority of HeteroPCA to the SVD-based approach in estimation.* To begin with, we first compare the empirical estimation accuracy of the SVD approach (cf. Algorithm 1) and HeteroPCA (cf. Algorithm 2). Figure 1 displays the relative estimation errors—including the ones tailored to the principal subspace:  $\|UR - U^*\|/\|U^*\|$ ,  $\|UR - U^*\|_F/\|U^*\|_F$ ,  $\|UR - U^*\|_{2,\infty}/\|U^*\|_{2,\infty}$ , and the ones tailored to the covariance matrix:  $\|S - S^*\|/\|S^*\|$ ,  $\|S - S^*\|_F/\|S^*\|_F$ ,  $\|S - S^*\|_{\infty}/\|S^*\|_{\infty}$ —of both algorithms as the noise level  $\omega^*$  varies, with  $r = 3$  and  $p = 0.6$ . Similarly, Figure 2 shows the relative numerical estimation errors of both algorithms versus the sampling rate  $p$ , with  $r = 3$  and  $\omega^* = 0.05$ . As we shall see from both figures, HeteroPCA uniformly outperforms the SVD-based approach in all experiments, and is able to achieve appealing performance for a much wider range of noise levels and sampling rates.

*Superiority of HeteroPCA to diagonal-deleted PCA in estimation.* Let us also compare the empirical estimation accuracy of the diagonal-deleted spectral method (Cai et al. (2021)) and HeteroPCA (cf. Algorithm 2). Recall from Section 3.2 that the main difference between the estimation error bounds of these two algorithms lies in an additional bias term due to the diagonal deletion operation (see the last term  $\mu_{\text{ce}}\kappa_{\text{ce}}r/d$  in equation (4.16) in Cai et al. (2021)). Figure 3 displays the relative estimation errors for estimating the principal subspace  $\|UR - U^*\|/\|U^*\|$  and for estimating the covariance matrix  $\|S - S^*\|/\|S^*\|$  as the dimension  $d$  varies with  $r = 3$ ,  $\omega^* = 0.05$  and  $p = 0.6$ . As can be seen from the plots, HeteroPCA uniformly outperforms the diagonal-deleted spectral method, especially when  $d$  is not too

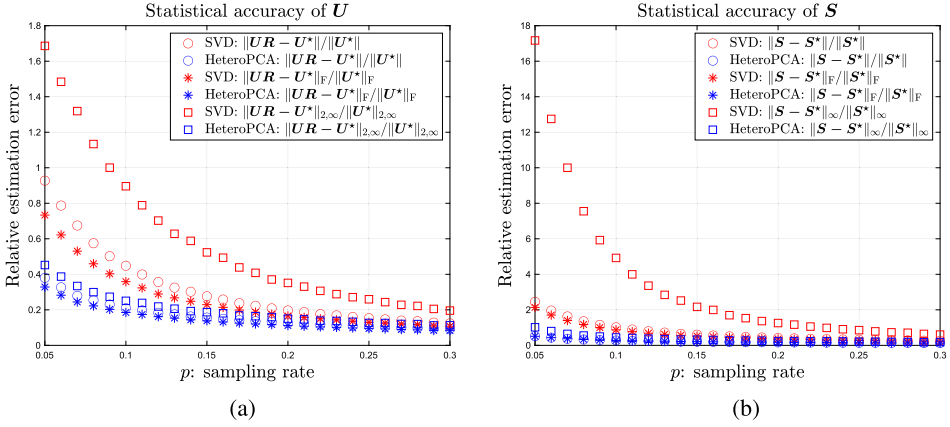


FIG. 2. The relative estimation error of  $U$  and  $S$  returned by both SVD-based approach (cf. Algorithm 1) and HeteroPCA (cf. Algorithm 2) across different missing probability  $p$ . (a) Relative estimation errors of  $UR - U^*$  measured by  $\|\cdot\|$ ,  $\|\cdot\|_F$  and  $\|\cdot\|_{2,\infty}$  vs. the missing rate  $p$ ; (b) Relative estimation errors of  $S - S^*$  measured by  $\|\cdot\|$ ,  $\|\cdot\|_F$  and  $\|\cdot\|_{\infty}$  vs. the missing rate  $p$ . The results are reported over 200 independent trials for  $r = 3$  and  $\omega^* = 0.05$ .

large. This numerical evidence corroborates the efficacy of the diagonal refinement scheme adopted in HeteroPCA.

*Confidence regions for the principal subspace  $U^*$ .* Next, we carry out a series of experiments to corroborate the practical validity of the confidence regions constructed using the SVD-based approach (Yan, Chen and Fan ((2021), Algorithm 3)) and HeteroPCA (cf. Algorithm 3). To this end, we define  $\widehat{\text{Cov}}_U(i)$  to be the empirical probability that the constructed confidence interval  $\text{CR}_{U,i}^{0.95}$  covers  $U^*_i \cdot \text{sgn}(U^{*\top}U)$  over 200 Monte Carlo trials, where  $U$  is the estimate returned by either algorithm. We also let  $\text{Mean}(\widehat{\text{Cov}}_U)$  (resp.,  $\text{std}(\widehat{\text{Cov}}_U)$ ) be the empirical mean (resp., standard deviation) of  $\widehat{\text{Cov}}_U(i)$  over  $i \in [d]$ . Table 1 gathers  $\text{Mean}(\widehat{\text{Cov}})$  and  $\text{std}(\widehat{\text{Cov}})$  for  $r = 3$  and different choices of  $(p, \omega^*)$  for both algorithms. Encouragingly, the empirical coverage rates are all close to 95% for both methods when  $p$  is not too small and  $\omega^*$  is not too large. When  $p$  becomes smaller or  $\omega^*$  grows larger, HeteroPCA is still capable of performing valid statistical inference, while the SVD-based approach fails.

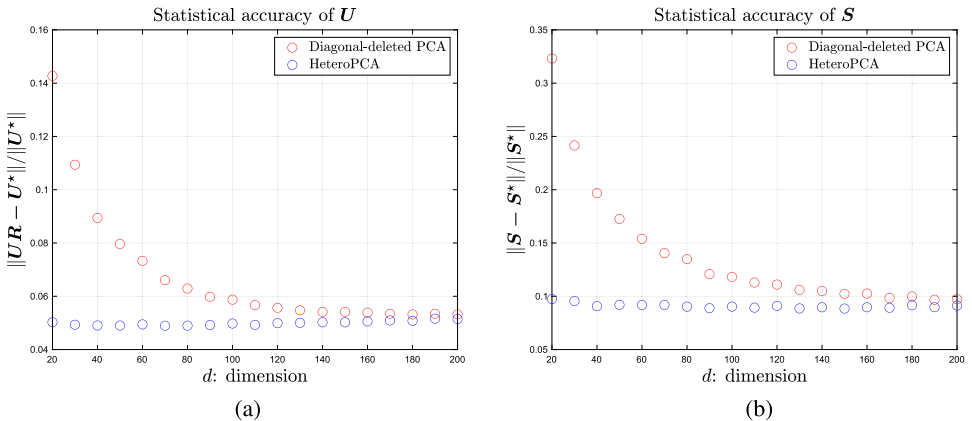


FIG. 3. The relative estimation error of  $U$  and  $S$  returned by both diagonal-deleted spectral method (Cai et al. (2021)) and HeteroPCA (cf. Algorithm 2). (a) Relative estimation error  $\|UR - U^*\|/\|U^*\|$  vs. dimension  $d$ ; (b) Relative estimation error  $\|S - S^*\|/\|S^*\|$  vs. the dimension  $d$ . The results are reported over 200 independent trials for  $r = 3$ ,  $\omega^* = 0.05$  and  $p = 0.6$ .

TABLE 1  
*Empirical coverage rates of  $U^* \text{sgn}(U^{*\top} U)$  for different  $(p, \omega^*)$ 's over 200 Monte Carlo trials*

$p$	$\omega^*$	The SVD-based approach		HeteroPCA	
		Mean( $\widehat{\text{Cov}}$ )	Std( $\widehat{\text{Cov}}$ )	Mean( $\widehat{\text{Cov}}$ )	Std( $\widehat{\text{Cov}}$ )
0.6	0.05	0.9270	0.0292	0.9523	0.0157
0.6	0.1	0.8989	0.0521	0.9484	0.0154
0.4	0.05	0.8849	0.0501	0.9448	0.0184
0.4	0.1	0.8458	0.0853	0.9405	0.0182
0.2	0.05	0.7370	0.1196	0.9287	0.0204
0.2	0.1	0.6856	0.1569	0.9219	0.0204

This provides another empirical evidence on the advantage and broader applicability of HeteroPCA compared to the SVD-based approach. In addition, for the rank-1 case ( $r = 1$ ), we define  $T_i := [U - \text{sign}(U^\top U^*)U^*]_i / \sqrt{\Sigma_{U,i}}$ . Figure 4 displays the Q-Q (quantile-quantile) plot of  $T_1 := [U - \text{sign}(U^\top U^*)U^*]_1 / \sqrt{\Sigma_{U,1}}$  versus the standard Gaussian random variable over 2000 Monte Carlo simulations for both algorithms (when  $p = 0.6$  and  $\omega^* = 0.05$ ); the near-Gaussian empirical distribution of  $T_1$  also corroborates our distributional guarantees.

*Entrywise confidence intervals for  $S^*$ .* Finally, we provide numerical evidence that confirms the validity of the confidence interval constructed on the basis of the SVD-based approach (Yan, Chen and Fan ((2021), Algorithm 4)) and HeteroPCA (cf. Algorithm 4). Define  $\widehat{\text{Cov}}_S(i, j)$  to be the empirical probability that the 95% confidence interval  $[S_{i,j} \pm 1.96\sqrt{v_{i,j}}]$  covers  $S_{i,j}^*$  over 200 Monte Carlo trials, where  $S_{i,j}$  is the  $(i, j)$ -th entry of the estimate  $S$  returned by either algorithm. Let  $\text{Mean}(\widehat{\text{Cov}}_S)$  (resp.,  $\text{std}(\widehat{\text{Cov}}_S)$ ) be the empirical mean (resp., standard deviation) of  $\widehat{\text{Cov}}_S(i, j)$  over all  $i, j \in [d]$ . Table 2 collects  $\text{Mean}(\widehat{\text{Cov}})$  and  $\text{std}(\widehat{\text{Cov}})$  for  $r = 3$  and accounts for different choices of  $(p, \omega^*)$  for both algorithms. Similar to previous experiments, HeteroPCA uniformly outperforms the SVD-based approach, which again suggests that HeteroPCA is the method of choice. In addition, we define  $Z_{i,j} := (S_{i,j} - S_{i,j}^*) / \sqrt{v_{i,j}}$ . For both algorithms, Figure 5 and Figure 6 depict the Q-Q (quantile-quantile) plot of  $Z_{1,1}$  and  $Z_{1,2}$  versus standard Gaussian distributions over 2000 Monte Carlo trials for the case with  $r = 3$ ,  $p = 0.6$  and  $\omega^* = 0.05$ , which again confirm the practical validity of our distributional theory.

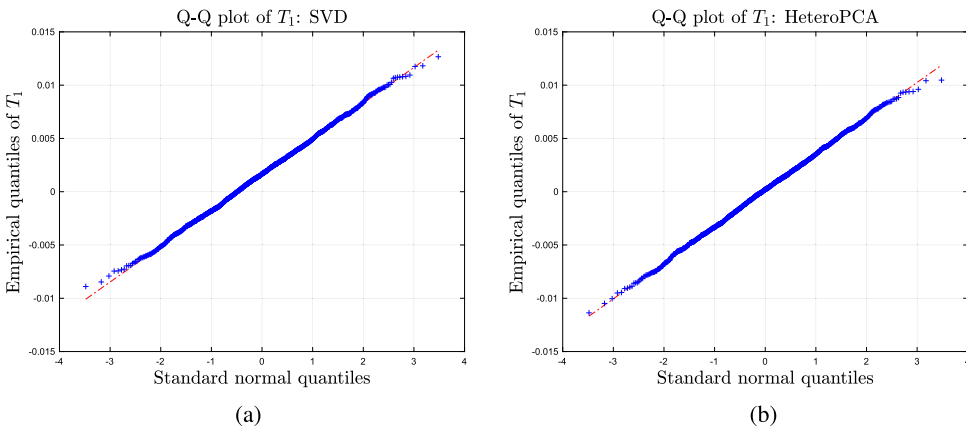


FIG. 4. (a) Q-Q (quantile-quantile) plot of  $T_1$  vs. the standard normal distribution for the SVD-based approach; (b) Q-Q (quantile-quantile) plot of  $T_1$  vs. the standard normal distribution for HeteroPCA. The results are reported over 2000 independent trials for  $r = 1$ ,  $p = 0.6$  and  $\omega^* = 0.05$ .



TABLE 2  
*Empirical coverage rates of  $S_{i,j}^*$  for different  $(\omega^*, p)$ 's over 200 Monte Carlo trials*

$p$	$\omega^*$	The SVD-based approach		HeteroPCA	
		Mean( $\widehat{\text{Cov}}$ )	Std( $\widehat{\text{Cov}}$ )	Mean( $\widehat{\text{Cov}}$ )	Std( $\widehat{\text{Cov}}$ )
0.6	0.05	0.9380	0.0244	0.9475	0.0153
0.6	0.1	0.9243	0.0425	0.9484	0.0151
0.4	0.05	0.9200	0.0509	0.9485	0.0156
0.4	0.1	0.9027	0.0713	0.9490	0.0153
0.2	0.05	0.8657	0.1031	0.9494	0.0164
0.2	0.1	0.8488	0.1186	0.9491	0.0162

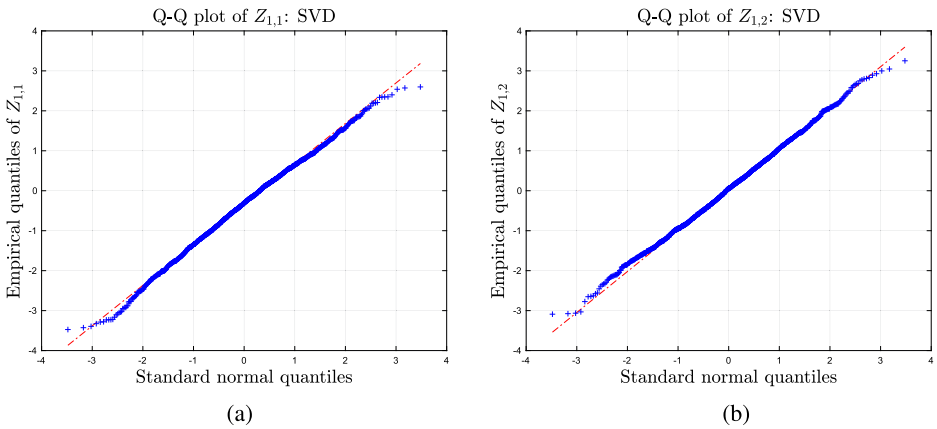


FIG. 5. (a)  $Q-Q$  (quantile–quantile) plot of  $Z_{1,1}$  vs. the standard normal distribution for the SVD-based approach; (b)  $Q-Q$  (quantile–quantile) plot of  $Z_{1,2}$  vs. a standard Gaussian distribution for the SVD-based approach. The results are reported over 2000 independent trials for  $r = 3, p = 0.6, \omega^* = 0.05$ .

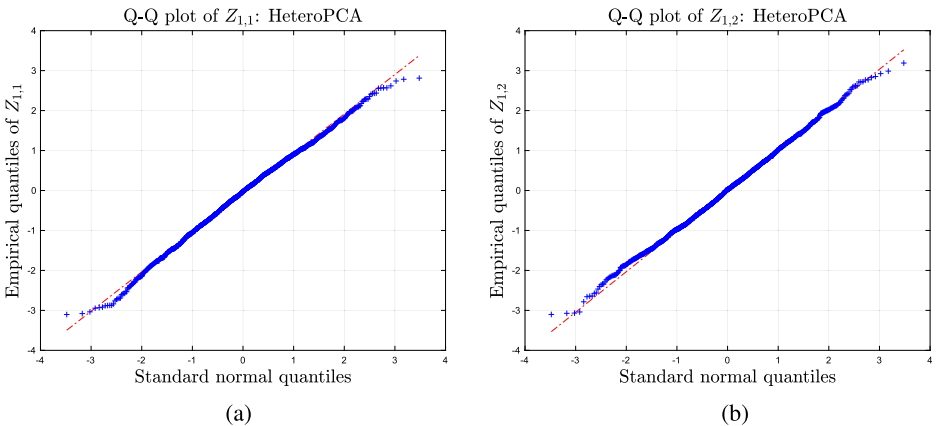


FIG. 6. (a)  $Q-Q$  (quantile–quantile) plot of  $Z_{1,1}$  vs. the standard normal distribution for HeteroPCA; (b)  $Q-Q$  (quantile–quantile) plot of  $Z_{1,2}$  vs. a standard Gaussian distribution for HeteroPCA. The results are reported over 2000 independent trials for  $r = 3, p = 0.6, \omega^* = 0.05$ .

**5. Other related works.** Low-rank matrix denoising serves as a common model to study the effectiveness of spectral methods (Chen et al. (2021)), and has been the main subject of many prior works including Abbe et al. (2020), Agterberg, Lubberts and Priebe (2022), Bao, Ding and Wang (2021), Cai and Zhang (2018), Cape, Tang and Priebe (2019), Chen, Cheng and Fan (2021), Ding (2020), Lei (2019), Montanari, Ruan and Yan (2018), Xia (2021), among others. Several recent works began to pursue a distributional theory for the eigenvector or singular vectors of the observed data matrix (Bao, Ding and Wang (2021), Cheng, Wei and Chen (2021), Fan et al. (2022), Xia (2021)). To name a few examples, Bao et al. (2022b) studied the limiting distribution of the inner product between an empirical singular vector and the corresponding ground truth, assuming that the associated spectral gap is sufficient large and that the noise components are homoskedastic; Xia (2021) established nonasymptotic Gaussian approximation for certain projection distance in the presence of i.i.d. Gaussian noise. Furthermore, the presence of missing data forms another source of technical challenges, leading to a problem often dubbed as noisy low-rank matrix completion (Candès and Plan (2010), Chen et al. (2020), Negahban and Wainwright (2012)).

Spectral methods have been successfully applied to tackle noisy matrix completion (Chen, Liu and Li (2020), Chen and Wainwright (2015), Cho, Kim and Rohe (2017), Keshavan, Montanari and Oh (2010b), Ma et al. (2020), Sun and Luo (2016), Zheng and Lafferty (2016)), which commonly serve as an effective initialization scheme for nonconvex optimization methods (Chi, Lu and Chen (2019)). While statistical inference for noisy matrix completion has been investigated recently (Chen et al. (2019b), Chernozhukov et al. (2023), Xia and Yuan (2021)), these prior works focused on performing inference based on optimization-based estimators. How to construct fine-grained confidence intervals based on spectral methods remains previously out of reach for noisy matrix completion. It is also noteworthy that the inferential procedures proposed in Chen et al. (2019b), Xia and Yuan (2021) (for noisy matrix completion) were developed for the regime where reliable estimation of the full low-rank matrix is feasible. This, however, falls short of covering the most challenging regime considered herein (where one might only be able to estimate the column subspace but not the row subspace). This crucial difference in the regimes of interest leads to substantial challenges unaddressed by these prior works.

Additionally, the recent work (Xia (2019)) tackled the confidence regions for spectral estimators tailored to the low-rank matrix regression problem, without accommodating the noisy matrix completion context. Most importantly, while the SVD-based vanilla spectral method often works well for the balanced case (such that the column dimension and the row dimension are on the same order), suboptimality has been well recognized when estimating the column subspace of interest in the highly unbalanced case (so that the column dimension far exceeds the row dimension); this issue is also present when it comes to existing optimization-based methods like nuclear norm minimization. As a result, all prior schemes mentioned in this paragraph failed to tackle the highly unbalanced case in an statistically efficient manner.

Turning to PCA or subspace estimation, there has been an enormous literature dedicated to this topic; see Balzano, Chi and Lu (2018), Johnstone and Paul (2018) for an overview of prior development. Noteworthy, the need to handle the diagonals of the sample covariance matrix in the presence of heteroskedastic noise and/or missing data has been pointed out in many prior works, for example, Cai et al. (2021), Florescu and Perkins (2016), Loh and Wainwright (2012), Lounici (2014), Montanari and Sun (2018). The iterative refinement scheme proposed by Zhang, Cai and Wu (2022b) turns out to be among the most effective and adaptive schemes in handling the diagonals. Aimed at designing fine-grained estimators for the principal components, Koltchinskii, Löffler and Nickl (2020), Li et al. (2021) proposed statistically efficient debiased estimators for linear functionals of principal components, and moreover, the estimator proposed in Koltchinskii, Löffler and Nickl (2020) has also been

shown to exhibit asymptotic normality in the presence of i.i.d. Gaussian noise. Bloemendal et al. (2016) also pinned down the asymptotic distributions of certain principal components under a spiked covariance model. However, these papers fell short of presenting valid and data-driven uncertainty quantification methods for the proposed estimators, and their results operates under the assumptions of homoskedastic noise without any missing data, a scenario that is remarkably more restricted than ours. Under the spiked covariance model, Bao et al. (2022b) studied the limiting distribution of the angle between the eigenvectors of the sample covariance matrix and any fixed vector, under the “balanced” scenario where the aspect ratio  $n/d$  is a constant. In addition, recent years have witnessed much activity in high-dimensional PCA in the face of missing data (Cai et al. (2021), Pavez and Ortega (2021), Zhang, Cai and Wu (2022b), Zhu, Wang and Samworth (2022)); these works, however, focused primarily on developing estimation guarantees, which did not provide either distributional guarantees for the estimators or concrete procedures that allow for confidence region construction. Additionally, the HeteroPCA algorithm has been further extended by two follow-up works Zhou and Chen (2023a), Zhou and Chen (2023b) to accommodate the scenario with large condition numbers as well as tensor clustering in the presence of heteroskedastic noise.

From a technical viewpoint, it is worth mentioning that the  $\ell_\infty$  and  $\ell_{2,\infty}$  perturbation theory has been an active research direction in recent years (Agterberg, Lubberts and Priebe (2022), Cape, Tang and Priebe (2019), Chen, Cheng and Fan (2021), Eldridge, Belkin and Wang (2018), Fan, Wang and Zhong (2017), Xie (2021)). Among multiple existing technical frameworks, the leave-one-out analysis idea—which has been applied to a variety of statistical estimation problems (Cai et al. (2022), Cai, Poor and Chen (2023), Chen, Gao and Zhang (2022), Chen et al. (2019), Chen et al. (2021), El Karoui (2018), El Karoui et al. (2013), Ling (2022), Zhong and Boumal (2018))—provides a powerful and flexible framework that enables  $\ell_\infty$  and  $\ell_{2,\infty}$  statistical guarantees for spectral methods (Abbe et al. (2020), Cai et al. (2021), Chen et al. (2019a)); see (Chen et al. ((2021), Chapter 4)) for an accessible introduction of this powerful framework. Our analysis for the HeteroPCA approach is influenced by the one in Cai et al. (2021). Note, however, that Cai et al. (2021) didn’t come with any distributional guarantees for spectral methods, which we seek to accomplish in this paper.

It is important to note that although the current version of this paper focuses primarily on the HeteroPCA method, a preliminary version available on arXiv (Yan, Chen and Fan (2021)) includes a discussion on distributional theory and inferential procedures for PCA using the SVD-based approach (cf. Algorithm 1). This content was subsequently omitted during the revision phase based on editorial suggestions. Interested readers are referred to Yan, Chen and Fan (2021) for a set of inferential results developed for the SVD-based approach, in parallel to Theorems 1 to 4 in this paper.

Finally, we note in passing that constructing confidence intervals for sparse regression (based on, say, the Lasso estimator or other sparsity-promoting estimator) has attracted a flurry of research activity in the past few years (Cai and Guo (2017), Celentano, Montanari and Wei (2023), Javanmard and Montanari (2014), Ning and Liu (2017), Ren et al. (2015), van de Geer et al. (2014), Zhang and Zhang (2014)). The methods derived therein, however, are not directly applicable to perform statistical inference for PCA and/or other low-rank models.

**6. A detour: Subspace estimation.** We now take a detour to look at an intimately related problem, which we shall refer to as *subspace estimation* and will play a crucial role in understanding the HeteroPCA approach. We will set out to develop a fine-grained statistical theory for HeteroPCA when applied to this subspace estimation setting. The resulting theory will be invoked in Appendix D in the Supplementary Material (Yan, Chen and Fan (2024)) to analyze the PCA context.

6.1. *Model and algorithm.*

*Model and assumptions.* Suppose that we are interested in a rank- $r$  matrix  $\mathbf{M}^\natural \in \mathbb{R}^{n_1 \times n_2}$ , whose SVD is given by

$$(6.1) \quad \mathbf{M}^\natural = \sum_{i=1}^r \sigma_i^\natural \mathbf{u}_i^\natural \mathbf{v}_i^{\natural\top} = \mathbf{U}^\natural \mathbf{\Sigma}^\natural \mathbf{V}^{\natural\top} \in \mathbb{R}^{n_1 \times n_2}.$$

Here,  $\mathbf{U}^\natural = [\mathbf{u}_1^\natural, \dots, \mathbf{u}_r^\natural]$  (resp.,  $\mathbf{V}^\natural = [\mathbf{v}_1^\natural, \dots, \mathbf{v}_r^\natural]$ ) consists of orthonormal columns that correspond to the left (resp., right) singular vectors of  $\mathbf{M}^\natural$ , and  $\mathbf{\Sigma}^\natural = \text{diag}\{\sigma_1^\natural, \dots, \sigma_r^\natural\}$  is a diagonal matrix consisting of the singular values of  $\mathbf{M}^\natural$ . Without loss of generality, we assume that

$$n = \max\{n_1, n_2\}.$$

It is assumed that the singular values are sorted (in magnitude) in descending order, namely

$$(6.2) \quad \sigma_1^\natural \geq \dots \geq \sigma_r^\natural \geq 0,$$

with the condition number denoted by

$$(6.3) \quad \kappa^\natural := \sigma_1^\natural / \sigma_r^\natural.$$

What we have observed is a noisy copy of  $\mathbf{M}^\natural$ , namely

$$(6.4) \quad \mathbf{M} = \mathbf{M}^\natural + \mathbf{E},$$

where  $\mathbf{E} = [E_{i,j}]_{1 \leq i, j \leq n}$  stands for a noise matrix. We focus on *estimating the column subspace* represented by  $\mathbf{U}^\natural$  and the singular values encapsulated in  $\mathbf{\Sigma}^\natural$ , but not the row space  $\mathbf{V}^\natural$ . An important special scenario one should bear in mind is the highly unbalanced case where the column dimension  $n_2$  far exceeds the row dimension  $n_1$ ; in this case, it is common to encounter situations where reliable estimation of  $\mathbf{M}^\natural$  and  $\mathbf{V}^\natural$  is infeasible but that of  $\mathbf{U}^\natural$  shows promise. For this reason, we refer to this setting as *subspace estimation* in order to differentiate it from matrix denoising, emphasizing that we are only interested in column subspace estimation.

With the new aim in mind, we shall modify our incoherence and noise assumptions accordingly. Here, we abuse the notation with the understanding that the following set of assumptions will be used only when analyzing the approach based on HeteroPCA. We shall also denote  $n := \max\{n_1, n_2\}$ .

**ASSUMPTION 2 (Incoherence).** *The rank- $r$  matrix  $\mathbf{M}^\natural \in \mathbb{R}^{n_1 \times n_2}$  defined in (6.1) is said to be  $\mu^\natural$ -incoherent if the following holds:*

$$\|\mathbf{U}^\natural\|_{2,\infty} \leq \sqrt{\frac{\mu^\natural r}{n_1}}, \quad \|\mathbf{V}^\natural\|_{2,\infty} \leq \sqrt{\frac{\mu^\natural r}{n_2}}, \quad \text{and} \quad \|\mathbf{M}^\natural\|_\infty \leq \sqrt{\frac{\mu^\natural}{n_1 n_2}} \|\mathbf{M}^\natural\|_F.$$

**ASSUMPTION 3 (Heteroskedastic random noise).** *Assume that the  $E_{i,j}$ 's are independently generated, and suppose that there exist nonnegative quantities  $\{\sigma_i\}_{i=1}^{n_1}$ ,  $\{B_i\}_{i=1}^{n_1}$ ,  $\sigma$  and  $B$  obeying*

$$\forall (i, j) \in [n_1] \times [n_2]: \quad \mathbb{E}[E_{i,j}] = 0, \quad \text{var}(E_{i,j}^2) = \sigma_{i,j}^2 \leq \sigma_i^2 \leq \sigma^2, \quad |E_{i,j}| \leq B_i \leq B,$$

where for all  $i \in [n_1]$ ,

$$(6.5) \quad B_i \lesssim \frac{\sigma_i \min\{\sqrt{n_2}, \sqrt[4]{n_1 n_2}\}}{\sqrt{\log n}}, \quad \text{and} \quad B \lesssim \frac{\sigma \min\{\sqrt{n_2}, \sqrt[4]{n_1 n_2}\}}{\sqrt{\log n}}.$$

---

**Algorithm 5** HeteroPCA for general subspace estimation (HeteroPCA)

---

*Initialization:* set  $\mathbf{G}^0 = \mathcal{P}_{\text{off-diag}}(\mathbf{M}\mathbf{M}^\top)$ .

*Updates:* for  $t = 0, 1, \dots, t_0$  do

$$(6.8a) \quad (\mathbf{U}^t, \mathbf{\Lambda}^t) = \text{eigs}(\mathbf{G}^t, r);$$

$$(6.8b) \quad \mathbf{G}^{t+1} = \mathcal{P}_{\text{off-diag}}(\mathbf{M}\mathbf{M}^\top) + \mathcal{P}_{\text{diag}}(\mathbf{U}^t \mathbf{\Lambda}^t \mathbf{U}^{t\top}).$$

Here,  $\text{eigs}(\mathbf{G}, r)$  returns  $(\mathbf{U}, \mathbf{\Lambda})$  where  $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^\top$  is the top- $r$  eigendecomposition of  $\mathbf{G}$ .

*Output:*  $\mathbf{U} = \mathbf{U}^{t_0}$ ,  $\mathbf{\Lambda} = \mathbf{\Lambda}^{t_0}$ ,  $\mathbf{\Sigma} = (\mathbf{\Lambda}^{t_0})^{1/2}$ ,  $\mathbf{S} = \mathbf{U}^{t_0} \mathbf{\Lambda}^{t_0} \mathbf{U}^{t_0\top}$ .

---

*Algorithm: HeteroPCA for subspace estimation.* The paradigm HeteroPCA can naturally be applied to tackle the above subspace estimation task. Let us introduce the ground-truth gram matrix as follows:

$$(6.6) \quad \mathbf{G}^\natural := \mathbf{M}^\natural \mathbf{M}^{\natural\top}.$$

Given that  $\mathbf{M} = \mathbf{M}^\natural + \mathbf{E}$  is an unbiased estimate of  $\mathbf{M}^\natural$ , one might naturally attempt to estimate the column space of  $\mathbf{M}$  by looking at the eigenspace of the sample Gram matrix  $\mathbf{M}\mathbf{M}^\top$ . It can be easily seen that

$$(6.7) \quad E[\mathbf{M}\mathbf{M}^\top] = \mathbf{M}^\natural \mathbf{M}^{\natural\top} + \text{diag} \left\{ \left[ \sum_{j=1}^{n_2} \sigma_{i,j}^2 \right]_{1 \leq i \leq n_1} \right\},$$

where the diagonal term on the right-hand side of (6.7) might incur significant bias in the most challenging regime. The HeteroPCA algorithm seeks to handle the diagonal part in an iterative manner, alternating between imputing the values of the diagonal entries and eigendecomposition of  $\mathbf{M}\mathbf{M}^\top$  with the diagonal replaced by the imputed values. The procedure is summarized in Algorithm 5.

6.2. *Fine-grained statistical guarantees for HeteroPCA.* We now move on to present our theoretical guarantees for Algorithm 5. In order to account for the potential global rotational ambiguity, we introduce the following rotation matrix as before:

$$(6.9) \quad \mathbf{R}\mathbf{U} := \arg \min_{\mathbf{O} \in \mathcal{O}^{r \times r}} \|\mathbf{U}\mathbf{O} - \mathbf{U}^\natural\|_{\mathbb{F}}^2,$$

where we recall that  $\mathcal{O}^{r \times r}$  represents the set of  $r \times r$  orthonormal matrices. It is also helpful to define the following quantities: for all  $m \in [n_1]$ ,

$$(6.10a) \quad \zeta_{\text{op}} := \sigma^2 \sqrt{n_1 n_2} \log n + \sigma \sigma_1^\natural \sqrt{n_1 \log n},$$

$$(6.10b) \quad \zeta_{\text{op},m} := \sigma \sigma_m \sqrt{n_1 n_2} \log n + \sigma_m \sigma_1^\natural \sqrt{n_1 \log n}.$$

Our result is as follows with the proof postponed to Appendix C in the Supplementary Material (Yan, Chen and Fan (2024)).

**THEOREM 5.** *Suppose that Assumptions 2–3 hold. Assume that*

$$(6.11) \quad n_1 \gtrsim \kappa^{\natural 4} \mu^{\natural} r + \mu^{\natural 2} r \log^2 n, \quad n_2 \gtrsim r \log^4 n, \quad \text{and} \quad \zeta_{\text{op}} \ll \frac{\sigma_r^{\natural 2}}{\kappa^{\natural 2}},$$

*and that the algorithm is run for  $t_0 \geq \log(\frac{\sigma_1^{\natural 2}}{\zeta_{\text{op}}})$  iterations. With probability exceeding  $1 - O(n^{-10})$ , there exist two matrices  $\mathbf{Z}$  and  $\mathbf{\Psi}$  such that the estimates returned by HeteroPCA*

obey

$$(6.12) \quad \mathbf{U}\mathbf{R}_U - \mathbf{U}^{\natural} = \mathbf{Z} + \mathbf{\Psi},$$

where

$$(6.13a) \quad \mathbf{Z} := \mathbf{E}\mathbf{V}^{\natural}(\mathbf{\Sigma}^{\natural})^{-1} + \mathcal{P}_{\text{off-diag}}(\mathbf{E}\mathbf{E}^{\top})\mathbf{U}^{\natural}(\mathbf{\Sigma}^{\natural})^{-2},$$

$$(6.13b) \quad \|\mathbf{\Psi}\|_{2,\infty} \lesssim \kappa^{\natural 2} \frac{\mu^{\natural} r}{n_1} \frac{\zeta_{\text{op}}}{\sigma_r^{\natural 2}} + \kappa^{\natural 2} \frac{\zeta_{\text{op}}^2}{\sigma_r^{\natural 4}} \sqrt{\frac{\mu^{\natural} r}{n_1}}.$$

In fact, for each  $m \in [n_1]$ , we further have

$$(6.13c) \quad \|\mathbf{Z}_{m,\cdot}\|_2 \lesssim \frac{\zeta_{\text{op},m}}{\sigma_r^{\natural 2}} \sqrt{\frac{\mu^{\natural} r}{n_1}} + \|\mathbf{U}_{m,\cdot}^{\natural}\|_2 \left( \kappa^{\natural 2} \sqrt{\frac{\mu^{\natural} r}{n_1}} \frac{\zeta_{\text{op}}}{\sigma_r^{\natural 2}} + \kappa^{\natural 2} \frac{\zeta_{\text{op}}^2}{\sigma_r^{\natural 4}} \right),$$

$$(6.13d) \quad \|\mathbf{\Psi}_{m,\cdot}\|_2 \lesssim \kappa^{\natural 2} \frac{\zeta_{\text{op}} \zeta_{\text{op},m}}{\sigma_r^{\natural 4}} \sqrt{\frac{\mu^{\natural} r}{n_1}} + \|\mathbf{U}_{m,\cdot}^{\natural}\|_2 \left( \kappa^{\natural 2} \sqrt{\frac{\mu^{\natural} r}{n_1}} \frac{\zeta_{\text{op}}}{\sigma_r^{\natural 2}} + \kappa^{\natural 2} \frac{\zeta_{\text{op}}^2}{\sigma_r^{\natural 4}} \right).$$

REMARK 5. The interested reader might wonder why Theorem 5 is not valid for small  $n_1$  (e.g., (6.11) does not hold when  $n_1 = 2$ ), and we provide some intuition here. Recall from (6.7) that the diagonal of the sample Gram matrix can be significantly biased, and the HeteroPCA algorithm uses the off-diagonal information to iteratively estimate and refine the diagonal. When  $n_1$  is small, the (untrustworthy) diagonal entries account for a nonnegligible fraction of all entries of the entire sample Gram matrix, and as a result, we cannot hope to debias the diagonal reliably by HeteroPCA using only off-diagonal observations.

The expressions (6.12) and (6.13) make apparent a key decomposition of the estimation error. As we shall see, the term  $\mathbf{Z}$  is often the dominant term, which captures both the first-order and second-order approximation (w.r.t. the noise matrix  $\mathbf{E}$ ) of the estimation error. Unless the noise level  $\sigma$  is very small, we cannot simply ignore the second-order term  $\mathcal{P}_{\text{off-diag}}(\mathbf{E}\mathbf{E}^{\top})\mathbf{U}^{\natural}(\mathbf{\Sigma}^{\natural})^{-2}$ , as it is not necessarily dominated in size by the linear mapping term  $\mathbf{E}\mathbf{V}^{\natural}(\mathbf{\Sigma}^{\natural})^{-1}$ . The simple and closed-form expression of  $\mathbf{Z}$ —in conjunction with the fact that  $\mathbf{\Psi}$  is well controlled—plays a crucial role when developing a nonasymptotic distributional theory.

While Theorem 5 is established mainly to help derive distributional characterizations for PCA, we remark that our analysis also delivers  $\ell_{2,\infty}$  statistical guarantees in terms of estimating  $\mathbf{U}^{\natural}$  (see Lemma 6 in the Supplementary Material (Yan, Chen and Fan (2024))). More specifically, our analysis asserts that

$$(6.14) \quad \|\mathbf{U}\mathbf{R}_U - \mathbf{U}^{\natural}\|_{2,\infty} \lesssim \frac{\zeta_{\text{op}}}{\sigma_r^{\natural 2}} \sqrt{\frac{\mu^{\natural} r}{n_1}}$$

with high probability, under the conditions of Theorem 5. It is perhaps helpful to compare (6.14) with prior  $\ell_{2,\infty}$  theory concerning estimation of  $\mathbf{U}^{\natural}$ .

- We first compare Theorem 5 with the recent work (Agterberg, Lubberts and Priebe ((2022), Theorem 2), which focused on the regime  $n_2 \gtrsim n_1$  and showed that

$$\inf_{\mathbf{O} \in \mathcal{O}^{r \times r}} \|\mathbf{U}\mathbf{O} - \mathbf{U}^{\natural}\|_{2,\infty} \lesssim \left( \frac{\sigma^2}{\sigma_r^{\natural 2}} \sqrt{rn_1 n_2} \log n + \kappa^{\natural} \frac{\sigma}{\sigma_r^{\natural}} \sqrt{rn_1 \log n} \right) \sqrt{\frac{\mu^{\natural} r}{n_1}} \asymp \frac{\zeta_{\text{op}}}{\sigma_r^{\natural 2}} \sqrt{\frac{\mu^{\natural} r^2}{n_1}}$$

under the noise condition  $\sigma \sqrt{n_2} \ll \sigma_r^{\natural} / (\kappa^{\natural} \sqrt{r \log n})$  (in addition to a few other conditions omitted here). Note that when  $\kappa^{\natural}, \mu^{\natural}, r \asymp 1$ , their  $\ell_{2,\infty}$  error bound resembles (6.14), but

the condition  $\sigma\sqrt{n_2} \ll \sigma_r^\natural/\sqrt{\log n}$  required therein is much stronger than the noise condition  $\zeta_{\text{op}} \ll \sigma_r^{\natural 2}$ —which is equivalent to  $\sigma\sqrt[4]{n_1 n_2} \ll \sigma_r^\natural/\sqrt{\log n}$  when  $n_2 \gtrsim n_1$ —imposed by our theory (see (6.11)). It is also worth emphasizing that the theory of [Agterberg, Lubberts and Priebe \(2022\)](#) is capable of accommodating dependent data (i.e., they only require the rows of  $\mathbf{E}$  to be independent and allow dependence within rows), which is beyond the scope of the present paper.

- Compared with the  $\ell_{2,\infty}$  estimation error guarantees for the diagonal-deleted spectral method in [Cai et al. \(\(2021\), Theorem 1\)](#), our bound (6.14) is able to get rid of the bias term incurred by diagonal deletion (see [Cai et al. \(\(2021\), equation \(17\)\)](#)), thus improving upon this prior result.

It should be noted that fine-grained perturbation results akin to [Theorem 5](#) were also developed for the SVD algorithm in an earlier version of this paper on arXiv, as detailed in [Yan, Chen and Fan \(\(2021\), Section 6.1\)](#). Subsequently, [Yan and Wainwright \(2024\)](#) presented more refined results for cases where the entries of the noise matrix  $\mathbf{E}$  follow a sub-Gaussian distribution, with further information available in [Appendix F](#) therein.

Before concluding this section, it is natural to ask whether [Theorem 5](#) can be used to conduct subspace inference when every entry of  $\mathbf{E}$  is allowed to have completely difference variance. To begin with, for a broad class of  $\mathbf{E}$  with independent and heteroskedastic components, we can readily apply [Theorem 5](#) to obtain a distributional theory for HeteroPCA when estimating  $\mathbf{U}^\natural$ . Caution needs to be exercised, however, when it comes to confidence interval construction. On closer inspection, evaluating  $\mathbf{Z}$  (i.e., the first- and second-order approximation of the subspace estimation error) in [Theorem 5](#) requires knowledge about the right singular subspace  $\mathbf{V}^\natural$  of  $\mathbf{M}^\natural$ , which might sometimes be difficult or even infeasible to estimate in the unbalanced regime where  $n_2 \gg n_1$ . As a result, our theory is not guaranteed to deliver useful inferential methods for such cases, unless additional information about  $\mathbf{V}^\natural$  is available.

**7. Discussion.** In this paper, we have developed a suite of statistical inference procedures to construct confidence regions for PCA in the presence of missing data and heterogeneous corruption, which should be easy-to-use in practice due to their data-driven nature. Compared to other prior algorithms like the SVD-based approach and the diagonal-deleted spectral method, the solution developed based on HeteroPCA enjoys a broadened applicability range without compromising statistical efficiency. The fine-grained distributional characterizations we have developed are nonasymptotic, which naturally lend themselves to high-dimensional settings.

Moving forward, there are a variety of directions that are worthy of further investigation.

- *Improved dependency on  $\kappa$ ,  $\mu$ ,  $r$  and  $\kappa_\omega$ .* In our general theorems (see [Theorems 11–14](#) in the [Supplementary Material \(Yan, Chen and Fan \(2024\)\)](#)), we allow  $\kappa$ ,  $\mu$ ,  $r$  and  $\kappa_\omega$  to grow. However, our theoretical results scale suboptimally with these problem parameters. It remains unclear how to sharpen the dependency on these parameters, which might require developing more refined analysis techniques.
- *Approximate low-rank structure.* Our results assume exact low-rank structure of the spiked component  $\mathbf{S}^\star$  of the covariance matrix. In reality, there is no shortage of applications where  $\mathbf{S}^\star$  is at best approximately low rank. How to develop trustworthy inference procedures in the presence of approximate low-rank structure? Unfortunately, our current leave-one-out analysis framework relies heavily on the exact rank- $r$  structure (unless  $\sigma_{r+1}^\star$  is extremely small); new analysis ideas are needed in order to tackle approximate low-rank structure.

- *General missing pattern.* Uncertainty quantification in the face of heterogeneous missing patterns is another important topic of practical value. Consider, for example, the case where the entries in the same row of  $X$  are sampled with the same rate (i.e., the  $(l, j)$ -th entry of  $X$  is observed with probability  $p_l$ ). Then by constructing the following data matrix via inverse probability weighting:

$$[\text{diag}(p_1, p_2, \dots, p_d)]^{-1} Y,$$

we obtain an unbiased estimate of  $X$ , and the theory developed can be readily extended to perform valid inference. Note that we can also replace  $\{p_l\}$  via their empirical estimates in the inference procedures. Nevertheless, in the more general case where the sampling rates are allowed to vary across all locations, it is unclear how to construct an unbiased estimate of  $X$  without knowing the per-entry sampling rates in advance; hence, our theory fails to accommodate this general scenario. Extending our current results to such general sampling patterns might call for new analysis tools.

- *Inference for individual principal components.* Moving beyond inference and uncertainty quantification for the principal subspace and the spiked covariance matrix, it is interesting to investigate how to conduct valid inference on individual principal components, particularly when the associated eigengap is vanishingly small (Li et al. (2021)).
- *Extension to unknown mean, dependent or adversarial noise.* If the observed data are inherently biased with *a priori* unknown means, how to properly compensate for the bias? What if the noise components are interdependent, and what if the observed data samples are further corrupted by a nonnegligible fraction of adversarial outliers?
- *Minimax-optimal estimation and inference.* As recognized in the matrix completion literature (Chen, Liu and Li (2020), Keshavan, Montanari and Oh (2010b), Ma et al. (2020)), spectral methods alone are in general unable to yield minimax-optimal statistical accuracy in the presence of missing data, given that spectral methods inherently treat the missingness effect as some sort of “noise.” The same message—namely, suboptimality of HeteroPCA in the face of missing data—carries over to the PCA setting considered herein. We conjecture that a subsequent refinement procedure (e.g., gradient descent tailored to compute the maximum likelihood estimate) is needed in order to reach minimax optimality, and we leave this for future investigation.
- *Applications in financial econometrics.* In addition to applications to the uncertainty quantification in the matrix completion problems in recommender system, the inferential procedure and analysis tools we have developed in this paper have applications in finance and econometrics. For example, our analysis and results for principal subspace are useful in testing factor structures in famous Fama–French factor models, and can also be used in sector/industry clustering using stock returns (Porter et al. (1998)); our results on uncertainty quantification for the spiked covariance matrix could also shed light on how to better quantify the risk in portfolio optimization that takes into account on the uncertainty in the risk estimation.

**Acknowledgments.** Yuxin Chen is the corresponding author.

**Funding.** Y. Chen is supported in part by the Alfred P. Sloan Research Fellowship, the Google Research Scholar Award, the AFOSR YIP award FA9550-19-1-0030, by the ONR grant N00014-19-1-2120, by the ARO grant W911NF-20-1-0097 and by the NSF grants CCF-1907661, IIS-2218713, IIS-2218773, DMS-2014279, CCF-2221009.

J. Fan is supported in part by the ONR grants N00014-19-1-2120, N00014-22-1-2340, by the NSF grants DMS-1662139, DMS-1712591, DMS-2052926, DMS-2053832, DMS-2210833 and by the NIH grant 2R01-GM072611-15.

Y. Yan is supported in part by the Charlotte Elizabeth Procter Honorific Fellowship from Princeton University and the Norbert Wiener Postdoctoral Fellowship from MIT.



## SUPPLEMENTARY MATERIAL

**Supplement to “Inference for heteroskedastic PCA with missing data”** (DOI: 10.1214/24-AOS2366SUPP; .pdf). Supplementary information.

## REFERENCES

- ABBE, E., FAN, J., WANG, K. and ZHONG, Y. (2020). Entrywise eigenvector analysis of random matrices with low expected rank. *Ann. Statist.* **48** 1452–1474. MR4124330 <https://doi.org/10.1214/19-AOS1854>
- AGTERBERG, J., LUBBERTS, Z. and PRIEBE, C. E. (2022). Entrywise estimation of singular vectors of low-rank matrices with heteroskedasticity and dependence. *IEEE Trans. Inf. Theory* **68** 4618–4650. MR4449064
- BAI, J. and WANG, P. (2016). Econometric analysis of large factor models. *Ann. Rev. Econ.* **8** 53–80.
- BAIK, J., BEN AROUS, G. and PÉCHÉ, S. (2005). Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *Ann. Probab.* **33** 1643–1697. MR2165575 <https://doi.org/10.1214/009117905000000233>
- BALZANO, L., CHI, Y. and LU, Y. M. (2018). Streaming PCA and subspace tracking: The missing data case. *Proc. IEEE* **106** 1293–1310. <https://doi.org/10.1109/JPROC.2018.2847041>
- BAO, Z., DING, X., WANG, J. and WANG, K. (2022a). Statistical inference for principal components of spiked covariance matrices. *Ann. Statist.* **50** 1144–1169. MR4404931 <https://doi.org/10.1214/21-aos2143>
- BAO, Z., DING, X., WANG, J. and WANG, K. (2022b). Statistical inference for principal components of spiked covariance matrices. *Ann. Statist.* **50** 1144–1169. MR4404931 <https://doi.org/10.1214/21-aos2143>
- BAO, Z., DING, X. and WANG, K. (2021). Singular vector and singular subspace distribution for the matrix denoising model. *Ann. Statist.* **49** 370–392. MR4206682 <https://doi.org/10.1214/20-AOS1960>
- BLOEMENDAL, A., KNOWLES, A., YAU, H.-T. and YIN, J. (2016). On the principal components of sample covariance matrices. *Probab. Theory Related Fields* **164** 459–552. MR3449395 <https://doi.org/10.1007/s00440-015-0616-x>
- CAI, C., LI, G., CHI, Y., POOR, H. V. and CHEN, Y. (2021). Subspace estimation from unbalanced and incomplete data matrices:  $\ell_{2,\infty}$  statistical guarantees. *Ann. Statist.* **49** 944–967. MR4255114 <https://doi.org/10.1214/20-aos1986>
- CAI, C., LI, G., POOR, H. V. and CHEN, Y. (2022). Nonconvex low-rank tensor completion from noisy data. *Oper. Res.* **70** 1219–1237. MR4409613
- CAI, C., POOR, H. V. and CHEN, Y. (2023). Uncertainty quantification for nonconvex tensor completion: Confidence intervals, heteroscedasticity and optimality. *IEEE Trans. Inf. Theory* **69** 407–452. MR4544966
- CAI, T. T. and GUO, Z. (2017). Confidence intervals for high-dimensional linear regression: Minimax rates and adaptivity. *Ann. Statist.* **45** 615–646. MR3650395 <https://doi.org/10.1214/16-AOS1461>
- CAI, T. T. and ZHANG, A. (2018). Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics. *Ann. Statist.* **46** 60–89. MR3766946 <https://doi.org/10.1214/17-AOS1541>
- CANDÈS, E. and PLAN, Y. (2010). Matrix completion with noise. *Proc. IEEE* **98** 925–936.
- CANDÈS, E. J. (2014). Mathematics of sparsity (and a few other things). In *Proceedings of the International Congress of Mathematicians—Seoul 2014, Vol. 1* 235–258. Kyung Moon Sa, Seoul. MR3728471
- CANDÈS, E. J. and RECHT, B. (2009). Exact matrix completion via convex optimization. *Found. Comput. Math.* **9** 717–772. MR2565240 <https://doi.org/10.1007/s10208-009-9045-5>
- CAPE, J., TANG, M. and PRIEBE, C. E. (2019). The two-to-infinity norm and singular subspace geometry with applications to high-dimensional statistics. *Ann. Statist.* **47** 2405–2439. MR3988761 <https://doi.org/10.1214/18-AOS1752>
- CELENTANO, M., MONTANARI, A. and WEI, Y. (2023). The Lasso with general Gaussian designs with applications to hypothesis testing. *Ann. Statist.* **51** 2194–2220. MR4678801 <https://doi.org/10.1214/23-aos2327>
- CHEN, J., LIU, D. and LI, X. (2020). Nonconvex rectangular matrix completion via gradient descent without  $\ell_{2,\infty}$  regularization. *IEEE Trans. Inf. Theory* **66** 5806–5841. MR4158648 <https://doi.org/10.1109/TIT.2020.2992234>
- CHEN, P., GAO, C. and ZHANG, A. Y. (2022). Partial recovery for top- $k$  ranking: Optimality of MLE and suboptimality of the spectral method. *Ann. Statist.* **50** 1618–1652. MR4441134 <https://doi.org/10.1214/21-aos2166>
- CHEN, Y., CHENG, C. and FAN, J. (2021). Asymmetry helps: Eigenvalue and eigenvector analyses of asymmetrically perturbed low-rank matrices. *Ann. Statist.* **49** 435–458. MR4206685 <https://doi.org/10.1214/20-AOS1963>
- CHEN, Y., CHI, Y., FAN, J. and MA, C. (2019). Gradient descent with random initialization: Fast global convergence for nonconvex phase retrieval. *Math. Program.* **176** 5–37. MR3960803 <https://doi.org/10.1007/s10107-019-01363-6>
- CHEN, Y., CHI, Y., FAN, J., MA, C. et al. (2021). Spectral methods for data science: A statistical perspective. *Found. Trends Mach. Learn.* **14** 566–806.

- CHEN, Y., CHI, Y., FAN, J., MA, C. and YAN, Y. (2020). Noisy matrix completion: Understanding statistical guarantees for convex relaxation via nonconvex optimization. *SIAM J. Optim.* **30** 3098–3121. MR4167625 <https://doi.org/10.1137/19M1290000>
- CHEN, Y., FAN, J., MA, C. and WANG, K. (2019a). Spectral method and regularized MLE are both optimal for top- $K$  ranking. *Ann. Statist.* **47** 2204–2235. MR3953449 <https://doi.org/10.1214/18-AOS1745>
- CHEN, Y., FAN, J., MA, C. and YAN, Y. (2019b). Inference and uncertainty quantification for noisy matrix completion. *Proc. Natl. Acad. Sci. USA* **116** 22931–22937. MR4036123 <https://doi.org/10.1073/pnas.1910053116>
- CHEN, Y., FAN, J., MA, C. and YAN, Y. (2021). Bridging convex and nonconvex optimization in robust PCA: Noise, outliers and missing data. *Ann. Statist.* **49** 2948–2971. MR4338899 <https://doi.org/10.1214/21-aos2066>
- CHEN, Y. and WAINWRIGHT, M. J. (2015). Fast low-rank estimation by projected gradient descent: General statistical and algorithmic guarantees. Available at [arXiv:1509.03025](https://arxiv.org/abs/1509.03025).
- CHENG, C., WEI, Y. and CHEN, Y. (2021). Tackling small eigen-gaps: Fine-grained eigenvector estimation and inference under heteroscedastic noise. *IEEE Trans. Inf. Theory* **67** 7380–7419. MR4345128 <https://doi.org/10.1109/TIT.2021.3111828>
- CHERNOZHUKOV, V., HANSEN, C., LIAO, Y. and ZHU, Y. (2023). Inference for low-rank models. *Ann. Statist.* **51** 1309–1330. MR4630950 <https://doi.org/10.1214/23-aos2293>
- CHI, Y., LU, Y. M. and CHEN, Y. (2019). Nonconvex optimization meets low-rank matrix factorization: An overview. *IEEE Trans. Signal Process.* **67** 5239–5269. MR4016283 <https://doi.org/10.1109/TSP.2019.2937282>
- CHO, J., KIM, D. and ROHE, K. (2017). Asymptotic theory for estimating the singular vectors and values of a partially-observed low rank matrix with noise. *Statist. Sinica* **27** 1921–1948. MR3726772
- DAVIS, C. and KAHAN, W. M. (1970). The rotation of eigenvectors by a perturbation. III. *SIAM J. Numer. Anal.* **7** 1–46. MR0264450 <https://doi.org/10.1137/0707001>
- DING, X. (2020). High dimensional deformed rectangular matrices with applications in matrix denoising. *Bernoulli* **26** 387–417. MR4036038 <https://doi.org/10.3150/19-BEJ1129>
- DONOHO, D., GAVISH, M. and JOHNSTONE, I. (2018). Optimal shrinkage of eigenvalues in the spiked covariance model. *Ann. Statist.* **46** 1742–1778. MR3819116 <https://doi.org/10.1214/17-AOS1601>
- EL KAROUI, N. (2018). On the impact of predictor geometry on the performance on high-dimensional ridge-regularized generalized robust regression estimators. *Probab. Theory Related Fields* **170** 95–175. MR3748322 <https://doi.org/10.1007/s00440-016-0754-9>
- EL KAROUI, N., BEAN, D., BICKEL, P. J., LIM, C. and YU, B. (2013). On robust regression with high-dimensional predictors. *Proc. Natl. Acad. Sci. USA* **110** 14557–14562.
- ELDRIDGE, J., BELKIN, M. and WANG, Y. (2018). Unperturbed: Spectral analysis beyond Davis–Kahan. In *Proceedings of Algorithmic Learning Theory. Proc. Mach. Learn. Res. (PMLR)* **83** 38. PMLR. MR3857310
- FAN, J., FAN, Y., HAN, X. and LV, J. (2022). Asymptotic theory of eigenvectors for random matrices with diverging spikes. *J. Amer. Statist. Assoc.* **117** 996–1009. MR4436328 <https://doi.org/10.1080/01621459.2020.1840990>
- FAN, J., LI, K. and LIAO, Y. (2021). Recent developments on factor models and its applications in econometric learning. *Annu. Rev. Financ. Econ.* **13** 401–430.
- FAN, J., WANG, K., ZHONG, Y. and ZHU, Z. (2021). Robust high-dimensional factor models with applications to statistical machine learning. *Statist. Sci.* **36** 303–327. MR4255196 <https://doi.org/10.1214/20-sts785>
- FAN, J., WANG, W. and ZHONG, Y. (2017). An  $\ell_\infty$  eigenvector perturbation bound and its application to robust covariance estimation. *J. Mach. Learn. Res.* **18** Paper No. 207, 42. MR3827095
- FAN, J. and YAO, Q. (2017). *The Elements of Financial Econometrics*. Cambridge Univ. Press, Cambridge.
- FLORESCU, L. and PERKINS, W. (2016). Spectral thresholds in the bipartite stochastic block model. In *Conference on Learning Theory* 943–959.
- GAGLIARDINI, P., OSSOLA, E. and SCAILLET, O. (2020). Estimation of large dimensional conditional factor models in finance. In *Handbook of Econometrics, Vol. 7A. Handbooks in Econom.* 219–282. Elsevier, Amsterdam. MR4254181 <https://doi.org/10.1016/bs.hoe.2020.10.001>
- JAVANMARD, A. and MONTANARI, A. (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *J. Mach. Learn. Res.* **15** 2869–2909. MR3277152
- JOHNSTONE, I. M. (2001). On the distribution of the largest eigenvalue in principal components analysis. *Ann. Statist.* **29** 295–327. MR1863961 <https://doi.org/10.1214/aos/1009210544>
- JOHNSTONE, I. M. and PAUL, D. (2018). PCA in high dimensions: An orientation. *Proc. IEEE* **106** 1277–1292. <https://doi.org/10.1109/JPROC.2018.2846730>
- JOLLIFFE, I. T. (1986). *Principal Component Analysis. Springer Series in Statistics*. Springer, New York. MR0841268 <https://doi.org/10.1007/978-1-4757-1904-8>
- KESHAVAN, R. H., MONTANARI, A. and OH, S. (2010a). Matrix completion from noisy entries. *J. Mach. Learn. Res.* **11** 2057–2078. MR2678022

- KESHAVAN, R. H., MONTANARI, A. and OH, S. (2010b). Matrix completion from a few entries. *IEEE Trans. Inf. Theory* **56** 2980–2998. MR2683452 <https://doi.org/10.1109/TIT.2010.2046205>
- KOLTCHINSKII, V., LÖFFLER, M. and NICKL, R. (2020). Efficient estimation of linear functionals of principal components. *Ann. Statist.* **48** 464–490. MR4065170 <https://doi.org/10.1214/19-AOS1816>
- LEI, L. (2019). Unified  $\ell_{2,\infty}$  eigenspace perturbation theory for symmetric random matrices. arXiv preprint. Available at arXiv:1909.04798.
- LI, G., CAI, C., GU, Y., POOR, H. V. and CHEN, Y. (2021). Minimax Estimation of Linear Functions of Eigenvectors in the Face of Small Eigen-Gaps. arXiv preprint. Available at arXiv:2104.03298.
- LING, S. (2022). Near-optimal performance bounds for orthogonal and permutation group synchronization via spectral methods. *Appl. Comput. Harmon. Anal.* **60** 20–52. MR4387245 <https://doi.org/10.1016/j.acha.2022.02.003>
- LOH, P.-L. and WAINWRIGHT, M. J. (2012). High-dimensional regression with noisy and missing data: Provable guarantees with nonconvexity. *Ann. Statist.* **40** 1637–1664. MR3015038 <https://doi.org/10.1214/12-AOS1018>
- LOUNICI, K. (2014). High-dimensional covariance matrix estimation with missing observations. *Bernoulli* **20** 1029–1058. MR3217437 <https://doi.org/10.3150/12-BEJ487>
- MA, C., WANG, K., CHI, Y. and CHEN, Y. (2020). Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution. *Found. Comput. Math.* **20** 451–632. MR4099988 <https://doi.org/10.1007/s10208-019-09429-9>
- MONTANARI, A., RUAN, F. and YAN, J. (2018). Adapting to unknown noise distribution in matrix denoising. arXiv preprint. Available at arXiv:1810.02954.
- MONTANARI, A. and SUN, N. (2018). Spectral algorithms for tensor completion. *Comm. Pure Appl. Math.* **71** 2381–2425. MR3862094 <https://doi.org/10.1002/cpa.21748>
- NADLER, B. (2008). Finite sample approximation results for principal component analysis: A matrix perturbation approach. *Ann. Statist.* **36** 2791–2817. MR2485013 <https://doi.org/10.1214/08-AOS618>
- NEGAHBAN, S. and WAINWRIGHT, M. J. (2012). Restricted strong convexity and weighted matrix completion: Optimal bounds with noise. *J. Mach. Learn. Res.* **13** 1665–1697. MR2930649
- NING, Y. and LIU, H. (2017). A general theory of hypothesis tests and confidence regions for sparse high dimensional models. *Ann. Statist.* **45** 158–195. MR3611489 <https://doi.org/10.1214/16-AOS1448>
- PAUL, D. (2007). Asymptotics of sample eigenstructure for a large dimensional spiked covariance model. *Statist. Sinica* **17** 1617–1642. MR2399865
- PAVEZ, E. and ORTEGA, A. (2021). Covariance matrix estimation with non uniform and data dependent missing observations. *IEEE Trans. Inf. Theory* **67** 1201–1215. MR4232009 <https://doi.org/10.1109/tit.2020.3039118>
- PORTER, M. E. et al. (1998). *Clusters and the New Economics of Competition* 76. Harvard Business Review, Boston, MA.
- REN, Z., SUN, T., ZHANG, C.-H. and ZHOU, H. H. (2015). Asymptotic normality and optimalities in estimation of large Gaussian graphical models. *Ann. Statist.* **43** 991–1026. MR3346695 <https://doi.org/10.1214/14-AOS1286>
- SCHÖNEMANN, P. H. (1966). A generalized solution of the orthogonal Procrustes problem. *Psychometrika* **31** 1–10. MR0215870 <https://doi.org/10.1007/BF02289451>
- SUN, R. and LUO, Z.-Q. (2016). Guaranteed matrix completion via non-convex factorization. *IEEE Trans. Inf. Theory* **62** 6535–6579. MR3565131 <https://doi.org/10.1109/TIT.2016.2598574>
- VAN DE GEER, S., BÜHLMANN, P., RITOV, Y. and DEZEURE, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *Ann. Statist.* **42** 1166–1202. MR3224285 <https://doi.org/10.1214/14-AOS1221>
- VASWANI, N., CHI, Y. and BOUWMANS, T. (2018). Rethinking PCA for modern data sets: Theory, algorithms, and applications. *Proc. IEEE* **106** 1274–1276.
- VERSHYNIN, R. (2018). *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics **47**. Cambridge Univ. Press, Cambridge. With a foreword by Sara van de Geer. MR3837109 <https://doi.org/10.1017/9781108231596>
- WAHBA, G. (1965). A least squares estimate of satellite attitude. *SIAM Rev.* **7** 409–409.
- XIA, D. (2019). Confidence region of singular subspaces for low-rank matrix regression. *IEEE Trans. Inf. Theory* **65** 7437–7459. MR4030894 <https://doi.org/10.1109/TIT.2019.2924900>
- XIA, D. (2021). Normal approximation and confidence region of singular subspaces. *Electron. J. Stat.* **15** 3798–3851. MR4298986 <https://doi.org/10.1214/21-ejs1876>
- XIA, D. and YUAN, M. (2021). Statistical inferences of linear forms for noisy matrix completion. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **83** 58–77. MR4220984 <https://doi.org/10.1111/rssb.12400>
- XIE, F. (2021). Entrywise limit theorems of eigenvectors and their one-step refinement for sparse random graphs. arXiv preprint. Available at arXiv:2106.09840.
- YAN, Y., CHEN, Y. and FAN, J. (2021). Inference for heteroskedastic PCA with missing data (full version). arXiv preprint. Available at arXiv:2107.12365.

- YAN, Y., CHEN, Y. and FAN, J. (2024). Supplement to “Inference for heteroskedastic PCA with missing data.” <https://doi.org/10.1214/24-AOS2366SUPP>
- YAN, Y. and WAINWRIGHT, M. J. (2024). Entrywise inference for causal panel data: A simple and instance-optimal approach. arXiv preprint. Available at [arXiv:2401.13665](https://arxiv.org/abs/2401.13665).
- ZHANG, A. R., CAI, T. T. and WU, Y. (2022b). Heteroskedastic PCA: Algorithm, optimality, and applications. *Ann. Statist.* **50** 53–80. [MR4382008 https://doi.org/10.1214/21-aos2074](https://doi.org/10.1214/21-aos2074)
- ZHANG, C.-H. and ZHANG, S. S. (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **76** 217–242. [MR3153940 https://doi.org/10.1111/rssb.12026](https://doi.org/10.1111/rssb.12026)
- ZHENG, Q. and LAFFERTY, J. (2016). Convergence analysis for rectangular matrix completion using Burer–Monteiro factorization and gradient descent. Available at [arXiv:1605.07051](https://arxiv.org/abs/1605.07051).
- ZHONG, Y. and BOUMAL, N. (2018). Near-optimal bound for phase synchronization. *SIAM J. Optim.* [MR3566919 https://doi.org/10.1137/16M105808X](https://doi.org/10.1137/16M105808X)
- ZHOU, Y. and CHEN, Y. (2023a). Deflated HeteroPCA: Overcoming the curse of ill-conditioning in heteroskedastic PCA. arXiv preprint. Available at [arXiv:2303.06198](https://arxiv.org/abs/2303.06198).
- ZHOU, Y. and CHEN, Y. (2023b). Heteroskedastic tensor clustering. arXiv preprint. Available at [arXiv:2311.02306/3](https://arxiv.org/abs/2311.02306).
- ZHU, Z., WANG, T. and SAMWORTH, R. J. (2022). High-dimensional principal component analysis with heterogeneous missingness. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **84** 2000–2031. [MR4515564](https://doi.org/10.1111/rssb.12026)