

Reconstruction of a signal matrix for high-dimension, low-sample-size data

筑波大学・数理物質科学研究科 村山 航 (Wataru Murayama)
Graduate School of Pure and Applied Sciences
University of Tsukuba

筑波大学・数学域 矢田 和善 (Kazuyoshi Yata)
Institute of Mathematics
University of Tsukuba

筑波大学・数学域 青嶋 誠 (Makoto Aoshima)
Institute of Mathematics
University of Tsukuba

Abstract: We consider recovering a signal matrix in high-dimension, low-sample-size (HDLSS) situations. We first consider using the conventional PCA to estimate a signal matrix and show that the induced estimator holds consistency properties under several conditions. We show that the estimator is directly affected by noise structures. In order to overcome the difficulty, we apply the noise-reduction (NR) methodology to a recovery of the signal matrix. We show that the NR method gives a preferable estimator of the signal matrix which holds the consistency properties under mild conditions. The NR method improves the accuracy of the conventional PCA successfully. Finally, we give several simulation results to recover a signal matrix.

Key words and phrases: HDLSS, Noise-reduction methodology, Large p small n , PCA.

1 Introduction

High-dimension, low-sample-size (HDLSS) data situations occur in many areas of modern science such as genetic microarrays, medical imaging, text recognition, finance, chemometrics, and so on. The asymptotic studies of HDLSS data are becoming increasingly relevant. In recent years, substantial work has been done on HDLSS asymptotic theory. Hall et al. [7], Ahn et al. [1], and Yata and Aoshima [11] explored several types of geometric representations of HDLSS data. Jung and Marron [8] investigated the inconsistency of the eigenvalues and eigenvectors of the sample covariance matrix. Yata and Aoshima [11] gave consistent estimators for both the eigenvalues and eigenvectors together with the principal component (PC) scores by developing the *noise-reduction (NR) methodology*. The HDLSS asymptotic theory had been studied under the assumption that either the population distribution is Gaussian or the random variables in a

sphered data matrix have a ρ -mixing dependency. However, Yata and Aoshima [9] provided asymptotic theory without assuming either the Gaussian assumption or the ρ -mixing condition. Moreover, Yata and Aoshima [10] created a new principal component analysis (PCA) called the *cross-data-matrix methodology* that is applicable to constructing an unbiased estimator in non-parametric settings. Aoshima and Yata [3, 4] developed a variety of high-dimensional statistical inference based on the geometric representations by using the cross-data-matrix methodology. See Aoshima and Yata [5, 6] for a review covering this field of research.

In this paper, we address the problem of recovering an unknown $d \times n$ low-rank matrix, $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$. \mathbf{A} is called the signal matrix. Let $r = \text{rank}(\mathbf{A})$. We assume $r (\leq \min\{d, n\})$ is fixed. Suppose we have a $d \times n$ data matrix, $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$, where

$$\mathbf{X} = \sqrt{n}\mathbf{A} + \mathbf{W}. \quad (1)$$

Here, $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_n]$ is a $d \times n$ noise matrix, where \mathbf{w}_j , $j = 1, \dots, n$, are independent and identically distributed (i.i.d.) as a d -dimensional distribution with mean zero and covariance matrix $\Sigma_W (\geq \mathbf{O})$. Note that $\mathbf{x}_j - \sqrt{n}\mathbf{a}_j$, $j = 1, \dots, n$, are i.i.d. Let $\Sigma_A = \mathbf{A}\mathbf{A}^T$. Then, it holds that $E(\mathbf{X}\mathbf{X}^T)/n = \Sigma_A + \Sigma_W (= \Sigma, \text{ say})$. Andrey and Nobel [2] considered the model (1) in a high-dimensional setting, where the data dimension d and the sample size n increase at the same rate, i.e. $n/d \rightarrow c > 0$. They assumed that the elements of \mathbf{W} are i.i.d. standard normal random variables. Note that the conditions such as “ $n/d \rightarrow c > 0$ ” and the normality are quite strict in real high-dimensional analyses. In this paper, we consider the model (1) in HDLSS settings without assuming the severe conditions.

The eigen-decomposition of Σ_W is given by $\Sigma_W = \mathbf{U}_W \Lambda_W \mathbf{U}_W^T$, where Λ_W is a diagonal matrix of eigenvalues, $\lambda_{1(W)} \geq \dots \geq \lambda_{d(W)} (\geq 0)$, and \mathbf{U}_W is an orthogonal matrix of the corresponding eigenvectors. Let $\mathbf{W} = \mathbf{U}_W \Lambda_W^{1/2} \mathbf{Z}$. Then, \mathbf{Z} is a $d \times n$ sphered data matrix from a distribution with the identity covariance matrix, \mathbf{I}_n . Here, we write $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_d]^T$ and $\mathbf{z}_j = (z_{j1}, \dots, z_{jn})^T$, $j = 1, \dots, d$. Note that $E(z_{jk}z_{j'k}) = 0$ ($j \neq j'$) and $\text{Var}(\mathbf{z}_j) = \mathbf{I}_n$. We assume that the fourth moments of each variable in \mathbf{Z} are uniformly bounded. The singular value decomposition of \mathbf{A} is given by $\mathbf{A} = \sum_{j=1}^r \lambda_{j(A)}^{1/2} \mathbf{u}_{j(A)} \mathbf{v}_{j(A)}^T$, where $\lambda_{1(A)}^{1/2} \geq \dots \geq \lambda_{r(A)}^{1/2} (\geq 0)$ are singular values of \mathbf{A} and $\mathbf{u}_{j(A)}$ (or $\mathbf{v}_{j(A)}$) denotes a unit left- (or right-) singular vector corresponding to $\lambda_{j(A)}^{1/2}$ ($j = 1, \dots, r$). In this paper, we assume the following model.

$$\limsup_{d \rightarrow \infty} \frac{\lambda_{1(A)}}{\lambda_{r(A)}} < \infty \quad \text{and} \quad \lim_{d \rightarrow \infty} \frac{\text{tr}(\Sigma_W^2)}{\lambda_{r(A)}^2} = 0. \quad (2)$$

The model (2) is a special case of the power spiked model given by Yata and Aoshima [12]. Also, we assume that $\lambda_j(A)$ s are distinct in the sense that

$$\liminf_{d \rightarrow \infty} |\lambda_{j(A)}/\lambda_{j'(A)} - 1| > 0$$

for all $j \neq j' (\leq r)$.

In Section 2, we consider using the conventional PCA to estimate \mathbf{A} and show that the induced estimator holds consistency properties under several conditions. In Section 3, we apply the NR method instead and show that the NR method gives a preferable estimator which holds the consistency properties under mild conditions. The NR method improves the accuracy of the conventional PCA successfully. Finally, in Section 4, we give several simulation results to recover a signal matrix.

2 Reconstruction of \mathbf{A} by conventional PCA

In this section, we consider recovering the signal matrix \mathbf{A} by using the conventional PCA. The sample covariance matrix is given by $\mathbf{S} = n^{-1} \mathbf{X} \mathbf{X}^T$. We consider the dual sample covariance matrix defined by $\mathbf{S}_D = n^{-1} \mathbf{X}^T \mathbf{X}$. Note that \mathbf{S}_D and \mathbf{S} share non-zero eigenvalues and $\text{rank}(\mathbf{S}) = \text{rank}(\mathbf{S}_D) \leq \min\{d, n\}$. Let $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_{\min\{d, n\}} \geq 0$ be the eigenvalues of \mathbf{S}_D . The eigen-decompositions of \mathbf{S} and \mathbf{S}_D are given by $\mathbf{S} = \sum_{j=1}^{\min\{d, n\}} \hat{\lambda}_j \hat{\mathbf{u}}_j \hat{\mathbf{u}}_j^T$ and $\mathbf{S}_D = \sum_{j=1}^{\min\{d, n\}} \hat{\lambda}_j \hat{\mathbf{v}}_j \hat{\mathbf{v}}_j^T$ respectively, where $\hat{\mathbf{u}}_j$ (or $\hat{\mathbf{v}}_j$) denotes a unit left- (or right-) singular vector of \mathbf{X} corresponding to $\hat{\lambda}_j$. Note that $\hat{\mathbf{u}}_j$ can be calculated by $\hat{\mathbf{u}}_j = (n \hat{\lambda}_j)^{-1/2} \mathbf{X} \hat{\mathbf{v}}_j$.

We reconstruct \mathbf{A} by $\hat{\lambda}_j$ s, $\hat{\mathbf{u}}_j$ s and $\hat{\mathbf{v}}_j$ s. We assume the following conditions as necessary:

$$(C-i) \quad \frac{\sum_{s,t=1}^d \lambda_s(W) \lambda_t(W) E\{(z_{sk}^2 - 1)(z_{tk}^2 - 1)\}}{n \lambda_{r(A)}^2} = o(1);$$

$$(C-ii) \quad \frac{\text{tr}(\boldsymbol{\Sigma}_W)}{n \lambda_{r(A)}} = o(1).$$

Let $\kappa_j = \text{tr}(\boldsymbol{\Sigma}_W)/(n \lambda_{j(A)})$ for $j = 1, \dots, r$. We have the following results.

Theorem 1. For $j = 1, \dots, r$, it holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\frac{\hat{\lambda}_j}{\lambda_{j(A)}} = 1 + \kappa_j + o_p(1), \quad \hat{\mathbf{u}}_j^T \mathbf{u}_{j(A)} = (1 + \kappa_j)^{-1/2} + o_p(1) \quad \text{and} \quad \hat{\mathbf{v}}_j^T \mathbf{v}_{j(A)} = 1 + o_p(1)$$

under (C-i).

Corollary 1. For $j = 1, \dots, r$, it holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\frac{\hat{\lambda}_j}{\lambda_{j(A)}} = 1 + o_p(1), \quad \hat{\mathbf{u}}_j^T \mathbf{u}_{j(A)} = 1 + o_p(1) \quad \text{and} \quad \hat{\mathbf{v}}_j^T \mathbf{v}_{j(A)} = 1 + o_p(1)$$

under (C-i) and (C-ii).

Based on the theoretical background, we consider recovering the signal matrix \mathbf{A} by $\hat{\mathbf{A}}_r = \sum_{j=1}^r \hat{\lambda}_j^{1/2} \hat{\mathbf{u}}_j \hat{\mathbf{v}}_j^T$. Then, we have the following results.

Theorem 2. It holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\|\hat{\mathbf{A}}_r - \mathbf{A}\|_F^2 = r \frac{\text{tr}(\boldsymbol{\Sigma}_W)}{n} + o_p(\lambda_{r(A)})$$

under (C-i), where $\|\cdot\|_F$ denotes the Frobenius norm.

Corollary 2. It holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\|\hat{\mathbf{A}}_r - \mathbf{A}\|_F^2 = o_p(\lambda_{r(A)})$$

under (C-i) and (C-ii).

It should be noted that (C-ii) is a necessary condition to claim the consistency property such as $\|\hat{\mathbf{A}}_r - \mathbf{A}\|_F^2 / \lambda_{r(A)} = o_p(1)$.

3 Reconstruction of A by NR method

We consider applying the NR method by Yata and Aoshima [11] to recover the signal matrix A . By using the NR method, we obtain an estimator of $\lambda_{i(A)}$ as

$$\hat{\lambda}_i = \hat{\lambda}_i - \frac{\text{tr}(\mathbf{S}_D) - \sum_{j=1}^i \hat{\lambda}_j}{n - i} \quad (i = 1, \dots, n - 1). \quad (3)$$

The following result claims that $\hat{\lambda}_i$ holds the consistency property without (C-ii). Remember that $\hat{\lambda}_i$ requires (C-ii) to hold the consistency property.

Theorem 3. For $j = 1, \dots, r$, it holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\frac{\hat{\lambda}_j}{\lambda_{j(A)}} = 1 + o_p(1)$$

under (C-i).

Now, we consider an adjustment of $\hat{\lambda}_i$ s to estimate the signal matrix A :

$$\hat{\lambda}_{i(r)} = \hat{\lambda}_i - \frac{\text{tr}(\mathbf{S}_D) - \sum_{j=1}^r \hat{\lambda}_j}{n - r} \quad (i = 1, \dots, r). \quad (4)$$

We consider recovering A by $\hat{A}_r = \sum_{j=1}^r \hat{\lambda}_{j(r)}^{1/2} \hat{u}_j \hat{v}_j^T$. Then, we have the following results.

Theorem 4. It holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\|\hat{A}_r - A\|_F^2 = 2 \sum_{i=1}^r \lambda_{i(A)} \left(1 - \frac{1}{(1 + \kappa_i)^{1/2}}\right) + o_p(\lambda_{r(A)})$$

under (C-i).

Corollary 3. It holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\|\hat{A}_r - A\|_F^2 = o_p(\lambda_{r(A)})$$

under (C-i) and (C-ii).

From Theorems 2 and 4, we compare $2\lambda_{i(A)}\{1 - 1/(1 + \kappa_i)^{1/2}\}$ with $\lambda_{i(A)}\kappa_i$. It holds that $2\{1 - 1/(1 + \kappa_i)^{1/2}\} < \kappa_i$ ($i = 1, \dots, r$) for any $\kappa_i > 0$, so that $\|\hat{A}_r - A\|_F^2$ is smaller than $\|\hat{A}_r - A\|_F^2$ asymptotically. Thus, \hat{A}_r improves the error rate of \hat{A}_r .

4 Simulations

We used computer simulations to compare the performance of \hat{A}_r with \hat{A}_r . We set $r = 3$ and $\Sigma_A = \text{diag}(\lambda_{1(A)}, \lambda_{2(A)}, \lambda_{3(A)}, 0, \dots, 0)$ with $\lambda_{1(A)} = 2d^{3/5}$, $\lambda_{2(A)} = 1.5d^{3/5}$ and $\lambda_{3(A)} = d^{3/5}$. Note that the model (2) holds. We generated pseudo random vectors for w_j , $j = 1, \dots, n$, i.i.d. as a d -dimensional normal distribution with mean zero and covariance matrix $\Sigma_W (\geq O)$.

We considered two cases: (a) $\Sigma_W = I_d$ and (b) $\Sigma_W = (\sigma_{ij})$ with $\sigma_{ij} = (0.3)^{|i-j|^{1/3}}$. Let $F(\mathbf{M}) = \|\mathbf{M} - \mathbf{A}\|_F^2 / \lambda_r(\mathbf{A})$ for any $d \times n$ matrix, \mathbf{M} . The findings were obtained by averaging the outcomes from 2000 ($= K$, say) replications. Under a fixed scenario, suppose that the k -th replication ends with estimates, $F(\hat{\mathbf{A}}_r)_k$ and $F(\check{\mathbf{A}}_r)_k$, for $k = 1, \dots, K$. Let us simply write $F_{\hat{\mathbf{A}}} = K^{-1} \sum_{k=1}^K F(\hat{\mathbf{A}}_r)_k$ and $F_{\check{\mathbf{A}}} = K^{-1} \sum_{k=1}^K F(\check{\mathbf{A}}_r)_k$. We also considered the Monte Carlo variability. Let $\text{var}(F_{\hat{\mathbf{A}}}) = (K-1)^{-1} \sum_{k=1}^K (F(\hat{\mathbf{A}}_r)_k - F_{\hat{\mathbf{A}}})^2$ and $\text{var}(F_{\check{\mathbf{A}}}) = (K-1)^{-1} \sum_{k=1}^K (F(\check{\mathbf{A}}_r)_k - F_{\check{\mathbf{A}}})^2$. Fig. 1 shows the behaviors of $(F_{\hat{\mathbf{A}}}, F_{\check{\mathbf{A}}})$ and $(\text{var}(F_{\hat{\mathbf{A}}}), \text{var}(F_{\check{\mathbf{A}}}))$ for $(d, n) = (2^s, 3s)$, $s = 5, \dots, 11$, in the case of (a). Fig. 2 shows them in the case of (b). The dashed lines denote the simulation results. In the left panels of each figure, we gave the corresponding theoretical values, $\text{rtr}(\Sigma_W)/(n\lambda_r)$ and $2 \sum_{i=1}^r \lambda_{i(\mathbf{A})} \{1 - (1 + \kappa_i)^{-1/2}\} / \lambda_r$, that are denoted by the solid lines. See Theorems 2 and 4 for the details. The simulation results appeared close to the theoretical values and it seemed to be good approximations. As expected theoretically, we observed that the estimates by the NR method give more preferable performances both for (a) and (b) compared to the conventional PCA.

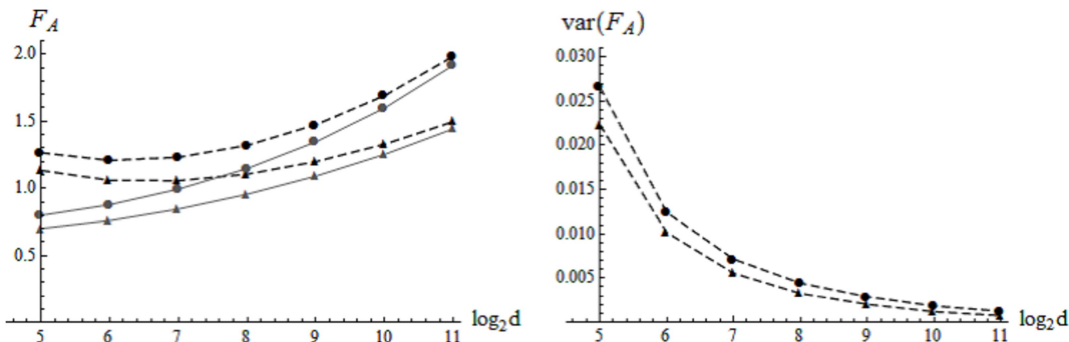


Figure 1: The behaviors of two estimates, $\hat{\mathbf{A}}_r$ denoted by \bullet and $\check{\mathbf{A}}_r$ denoted by \blacktriangle , when $\Sigma_W = I_d$. The values of $F_{\hat{\mathbf{A}}}$ and $F_{\check{\mathbf{A}}}$ are denoted by the dashed lines in the left panel. The values of their sample variances, $\text{var}(F_{\hat{\mathbf{A}}})$ and $\text{var}(F_{\check{\mathbf{A}}})$, are denoted by the dashed lines in the right panel. The theoretical values, $\text{rtr}(\Sigma_W)/(n\lambda_r)$ and $2 \sum_{i=1}^r \lambda_{i(\mathbf{A})} \{1 - (1 + \kappa_i)^{-1/2}\} / \lambda_r$, are denoted by the solid lines in the left panel.

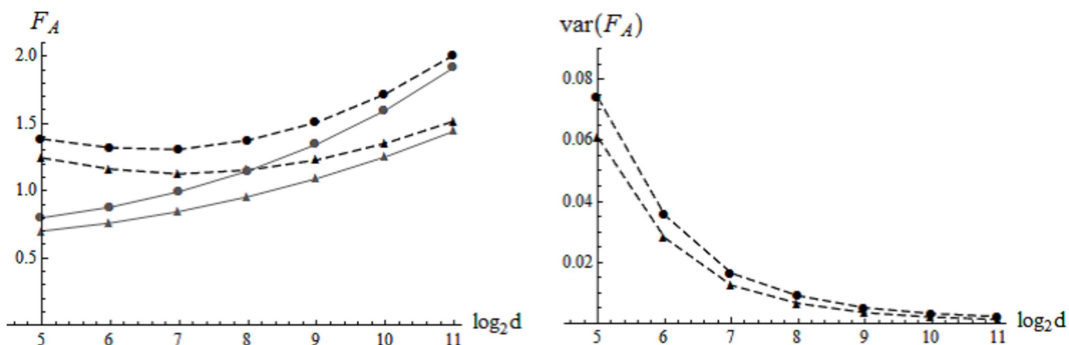


Figure 2: (Continued). When $\Sigma_W = (\sigma_{ij})$ with $\sigma_{ij} = (0.3)^{|i-j|^{1/3}}$.

A Appendix

Throughout, let $\mathbf{e}_n = (e_1, \dots, e_n)^T$ be an arbitrary unit random vector.

Lemma 1. *It holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$*

$$\mathbf{e}_n^T \frac{\mathbf{W}^T \mathbf{W}}{n\lambda_j(A)} \mathbf{e}_n = \frac{\text{tr}(\boldsymbol{\Sigma}_W)}{n\lambda_j(A)} + o_p(1)$$

for $j (\leq r)$ under (C-i).

Proof. We write that

$$\mathbf{e}_n^T \frac{\mathbf{W}^T \mathbf{W}}{n\lambda_j(A)} \mathbf{e}_n = \mathbf{e}_n^T (n^{-1} \sum_{s=1}^d \lambda_s(W) \mathbf{z}_s \mathbf{z}_s^T) \mathbf{e}_n = \mathbf{e}_n^T \{n^{-1} \sum_{s=1}^d \lambda_s(W) (\mathbf{z}_s \mathbf{z}_s^T - \mathbf{I}_n)\} \mathbf{e}_n + \frac{\text{tr}(\boldsymbol{\Sigma}_W)}{n}.$$

From Lemma 5 given in Yata and Aoshima [12], it holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\mathbf{e}_n^T \frac{n^{-1} \sum_{s=1}^d \lambda_s(W) (\mathbf{z}_s \mathbf{z}_s^T - \mathbf{I}_n)}{\lambda_j} \mathbf{e}_n = o_p(1)$$

for $j (\leq r)$ under (C-i). Thus it concludes the result. \square

Lemma 2. *It holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$*

$$\frac{\mathbf{u}_{i(A)}^T \mathbf{W} \mathbf{e}_n}{n^{1/2}} = o_p(\lambda_{r(A)}^{1/2}), \quad i = 1, \dots, r.$$

Proof. We write that $\mathbf{u}_{i(A)}^T \mathbf{W} \mathbf{e}_n = \sum_{k=1}^n e_k \mathbf{w}_k^T \mathbf{u}_{i(A)}$. Note that $\lambda_1(W) = o(\lambda_{r(A)})$ as $d \rightarrow \infty$ from (2). By using Markov's inequality, for any $\tau > 0$ and $i = 1, \dots, r$, we have that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\begin{aligned} P\left(\sum_{k=1}^n (\mathbf{w}_k^T \mathbf{u}_{i(A)})^2 / n \geq \tau \lambda_{r(A)}\right) &\leq \frac{E\{\sum_{k=1}^n (\mathbf{w}_k^T \mathbf{u}_{i(A)})^2\}}{\tau n \lambda_{r(A)}} = \frac{\mathbf{u}_{i(A)}^T \boldsymbol{\Sigma}_W \mathbf{u}_{i(A)}}{\tau \lambda_{r(A)}} \\ &\leq \frac{\lambda_1(W)}{\tau \lambda_{r(A)}} = o(1) \end{aligned}$$

from the fact that $\mathbf{u}_{i(A)}^T \boldsymbol{\Sigma}_W \mathbf{u}_{i(A)} \leq \lambda_1(W)$. Then, by noting that

$$\begin{aligned} \left| \sum_{k=1}^n e_k (\mathbf{w}_k^T \mathbf{u}_{i(A)}) / n^{1/2} \right| &\leq \left\{ \sum_{k=1}^n e_k^2 \right\}^{1/2} \left\{ \sum_{k=1}^n (\mathbf{w}_k^T \mathbf{u}_{i(A)})^2 / n \right\}^{1/2} \\ &= \left\{ \sum_{k=1}^n (\mathbf{w}_k^T \mathbf{u}_{i(A)})^2 / n \right\}^{1/2}, \end{aligned}$$

we can conclude the result. \square

Proof of Theorem 1. We write that for $j = 1, \dots, r$

$$\begin{aligned} \frac{\hat{\lambda}_j}{\lambda_{j(A)}} &= \hat{\mathbf{v}}_j^T \frac{\mathbf{S}_D}{\lambda_{j(A)}} \hat{\mathbf{v}}_j = \hat{\mathbf{v}}_j^T \frac{(\mathbf{A} + \mathbf{W}/n^{1/2})^T (\mathbf{A} + \mathbf{W}/n^{1/2})}{\lambda_{j(A)}} \hat{\mathbf{v}}_j \\ &= \hat{\mathbf{v}}_j^T \frac{\mathbf{A}^T \mathbf{A}}{\lambda_{j(A)}} \hat{\mathbf{v}}_j + 2\hat{\mathbf{v}}_j^T \frac{\mathbf{A}^T \mathbf{W}}{n^{1/2} \lambda_{j(A)}} \hat{\mathbf{v}}_j + \hat{\mathbf{v}}_j^T \frac{\mathbf{W}^T \mathbf{W}}{n \lambda_{j(A)}} \hat{\mathbf{v}}_j. \end{aligned} \quad (5)$$

We note that $|\hat{\mathbf{v}}_j^T \mathbf{A}^T \mathbf{W} \hat{\mathbf{v}}_j| \leq \sum_{i=1}^r \lambda_{i(A)}^{1/2} |\mathbf{u}_{i(A)}^T \mathbf{W} \hat{\mathbf{v}}_j|$ for $j = 1, \dots, r$. From Lemma 2 and (2), it holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\hat{\mathbf{v}}_j^T \mathbf{A}^T \mathbf{W} \hat{\mathbf{v}}_j / n^{1/2} = o_p(\lambda_{r(A)}^{1/2}) \quad (6)$$

for $j = 1, \dots, r$. Then, by combining (5) with Lemma 1 and (6), it holds that for $j = 1, \dots, r$

$$\frac{\hat{\lambda}_j}{\lambda_{j(A)}} = \hat{\mathbf{v}}_j^T \frac{\sum_{s=1}^r \lambda_{s(A)} \mathbf{v}_{s(A)} \mathbf{v}_{s(A)}^T}{\lambda_{j(A)}} \hat{\mathbf{v}}_j + \frac{\text{tr}(\boldsymbol{\Sigma}_W)}{n \lambda_{j(A)}} + o_p(1) = 1 + \frac{\text{tr}(\boldsymbol{\Sigma}_W)}{n \lambda_{j(A)}} + o_p(1)$$

under (C-i). Thus, we have that $\hat{\mathbf{v}}_j^T \mathbf{v}_{j(A)} = 1 + o_p(1)$ for $j = 1, \dots, r$. For $\hat{\mathbf{u}}_j$ s, from Lemma 2, under (C-i), it holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$\mathbf{u}_{j(A)}^T \hat{\mathbf{u}}_j = \hat{\lambda}_j^{-1/2} \lambda_{j(A)}^{1/2} \mathbf{v}_{j(A)}^T \hat{\mathbf{v}}_j + (n \hat{\lambda}_j)^{-1/2} \mathbf{u}_{j(A)}^T \mathbf{W} \hat{\mathbf{v}}_j = \{1 + \text{tr}(\boldsymbol{\Sigma}_W) / (n \lambda_{j(A)})\}^{-1/2} + o_p(1)$$

for $j = 1, \dots, r$. Thus it concludes the results. \square

Proof of Corollary 1. Note that $\text{tr}(\boldsymbol{\Sigma}_W)/n = o(\lambda_{r(A)})$ under (C-ii). From Theorem 1, we can conclude the result. \square

Proofs of Theorem 2 and Corollary 2. From Theorem 1, under (C-i), we have that

$$\|\hat{\mathbf{A}}_r - \mathbf{A}\|_F^2 = \left\| \sum_{s=1}^r (\hat{\lambda}_s^{1/2} \hat{\mathbf{u}}_s \hat{\mathbf{v}}_s^T - \lambda_{s(A)}^{1/2} \mathbf{u}_{s(A)} \mathbf{v}_{s(A)}^T) \right\|_F^2 = r \frac{\text{tr}(\boldsymbol{\Sigma}_W)}{n} + o_p(\lambda_{r(A)}).$$

It concludes the result. \square

Lemma 3. *It holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$*

$$\frac{\text{tr}(\mathbf{S}_D) - \sum_{i=1}^j \hat{\lambda}_i}{n - j} = \frac{\text{tr}(\boldsymbol{\Sigma}_W)}{n} + o_p(\lambda_{r(A)}) \quad \text{for } j = 1, \dots, r$$

under (C-i).

Proof. We write that $\text{tr}(\mathbf{W}^T \mathbf{W}/n) - \text{tr}(\boldsymbol{\Sigma}_W) = \sum_{s=1}^d \lambda_{s(W)} \sum_{k=1}^n (z_{sk}^2 - 1)/n$. Then, it holds that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$E[\{\text{tr}(\mathbf{W}^T \mathbf{W}/n) - \text{tr}(\boldsymbol{\Sigma}_W)\}^2] = \sum_{r,s=1}^d \lambda_{r(W)} \lambda_{s(W)} E\{(z_{rk}^2 - 1)(z_{sk}^2 - 1)\}/n = o(\lambda_{r(A)}^2)$$

under (C-i). Hence, by using Chebyshev's inequality, for any $\tau > 0$, we have that as $d \rightarrow \infty$ and $n \rightarrow \infty$

$$P\left(|\text{tr}(\mathbf{W}^T \mathbf{W}/n) - \text{tr}(\boldsymbol{\Sigma}_W)| \geq \tau \lambda_{r(A)}\right) \leq \frac{E[\{\text{tr}(\mathbf{W}^T \mathbf{W}/n) - \text{tr}(\boldsymbol{\Sigma}_W)\}^2]}{\tau^2 \lambda_{r(A)}^2} = o(1)$$

under (C-i). Thus it follows that $\text{tr}(\mathbf{W}^T \mathbf{W}/n) = \text{tr}(\boldsymbol{\Sigma}_W) + o_p(\lambda_{r(A)})$. On the other hand, from (2), we have that

$$\begin{aligned} E\{\text{tr}(\mathbf{A}^T \mathbf{W})^2\} &= E\left\{\left(\sum_{i=1}^n \mathbf{a}_i^T \mathbf{w}_i\right)^2\right\} = \sum_{i=1}^n \mathbf{a}_i^T \boldsymbol{\Sigma}_W \mathbf{a}_i = \text{tr}(\boldsymbol{\Sigma}_W \boldsymbol{\Sigma}_A) \\ &\leq \sqrt{\text{tr}(\boldsymbol{\Sigma}_W^2) \text{tr}(\boldsymbol{\Sigma}_A^2)} = O(\text{tr}(\boldsymbol{\Sigma}_W^2)^{1/2} r \lambda_{1(A)}) = o(\lambda_{r(A)}^2), \end{aligned}$$

so that $\text{tr}(\mathbf{A}^T \mathbf{W}) = o_p(\lambda_{r(A)})$. Then, we have that

$$\begin{aligned} \text{tr}(\mathbf{S}_D) &= \text{tr}(\mathbf{A}^T \mathbf{A}) + 2\text{tr}(\mathbf{A}^T \mathbf{W})/n^{1/2} + \text{tr}(\mathbf{W}^T \mathbf{W})/n \\ &= \sum_{i=1}^r \lambda_{i(A)} + \text{tr}(\boldsymbol{\Sigma}_W) + o_p(\lambda_{r(A)}). \end{aligned} \quad (7)$$

under (C-i). From Theorem 1 and (7), it holds that for $j (\leq r)$

$$\frac{\text{tr}(\mathbf{S}_D) - \sum_{i=1}^j \hat{\lambda}_i}{n - j} = \frac{\text{tr}(\boldsymbol{\Sigma}_W)}{n} + o_p(\lambda_{r(A)})$$

under (C-i). It concludes the result. \square

Proof of Theorem 3. By combining Theorem 1 with Lemma 3, we can conclude the result. \square

Proofs of Theorem 4 and Corollary 3. By combining Theorems 1 and 3, we can conclude the results. \square

Acknowledgment

Research of the first author was partially supported by Grant-in-Aid for Young Scientists (B), Japan Society for the Promotion of Science (JSPS), under Contract Number 26800078. Research of the second author was partially supported by Grants-in-Aid for Scientific Research (B) and Challenging Exploratory Research, JSPS, under Contract Numbers 22300094 and 26540010.

References

- [1] J. Ahn, J.S. Marron, K.M. Muller, Y.-Y. Chi, The high-dimension, low-sample-size geometric representation holds under mild conditions, *Biometrika* 94 (2007) 760-766.
- [2] A. Andrey, A. Nobel, Reconstruction of a low-rank matrix in the presence of Gaussian noise, *J. Multivariate Anal.* 118 (2013) 67-76.
- [3] M. Aoshima, K. Yata, Two-stage procedures for high-dimensional data, *Sequential Anal. (Editor's special invited paper)* 30 (2011) 356-399.
- [4] M. Aoshima, K. Yata, Authors' response, *Sequential Anal.* 30 (2011) 432-440.
- [5] M. Aoshima, K. Yata, Invited review article: Statistical inference in high-dimension, low-sample-size settings, *Sugaku* 65 (2013) 225-247.
- [6] M. Aoshima, K. Yata, The JSS research prize lecture: Effective methodologies for high-dimensional data, *J. Japan Statist. Soc. Ser. J* 43 (2013) 123-150.
- [7] P. Hall, J.S. Marron, A. Neeman, Geometric representation of high dimension, low sample size data, *J. R. Statist. Soc. Ser. B* 67 (2005) 427-444.
- [8] S. Jung, J.S. Marron, PCA consistency in high dimension, low sample size context, *Ann. Statist.* 37 (2009) 4104-4130.
- [9] K. Yata, M. Aoshima, PCA consistency for non-Gaussian data in high dimension, low sample size context, *Commun. Statist. Theory Methods, Special Issue Honoring S. Zacks (Ed. N. Mukhopadhyay)* 38 (2009) 2634-2652.
- [10] K. Yata, M. Aoshima, Effective PCA for high-dimension, low-sample-size data with singular value decomposition of cross data matrix, *J. Multivariate Anal.* 101 (2010) 2060-2077.
- [11] K. Yata, M. Aoshima, Effective PCA for high-dimension, low-sample-size data with noise reduction via geometric representations, *J. Multivariate Anal.* 105 (2012) 193-215.
- [12] K. Yata, M. Aoshima, PCA consistency for the power spiked model in high-dimensional settings, *J. Multivariate Anal.* 122 (2013) 334-354.