

# High-dimensional covariance matrix estimation under the SSE model

筑波大学大学院・数理物質科学研究科 小西 啓介 (Keisuke Konishi)

Graduate School of Pure and Applied Sciences

University of Tsukuba

筑波大学・数理物質系 矢田 和善 (Kazuyoshi Yata)

Institute of Mathematics

University of Tsukuba

筑波大学・数理物質系 青嶋 誠 (Makoto Aoshima)

Institute of Mathematics

University of Tsukuba

## Abstract

In this paper, we consider the estimation for the inverse matrix of a high-dimensional covariance matrix under the strongly spiked eigenvalue model. One of the well-known estimation methods is the principal orthogonal complement thresholding (POET) given by Fan et al. [5]. We show that the POET has consistency properties only under several severe conditions in high-dimensional settings. In order to overcome the difficulty, we consider applying the noise-reduction (NR) method given by Yata and Aoshima [8, 9] to the POET. We propose a new estimation of the inverse covariance matrix called the NR-POET. We compare the performance of the NR-POET with the POET by several simulations.

## 1 Introduction

One of the features of high-dimensional data is that the data dimension  $d$  is high, however, the sample size  $n$  is low. This is the so-called “HDLSS” or “large  $d$ , small  $n$ ” data. Such data situations appear in many fields of modern science such as genetic microarrays, medical imaging, text recognition, finance, chemometrics, and so on. For HDLSS data, the sample covariance matrix does not have its inverse matrix. The estimation for the inverse matrix of a covariance matrix is a crucial issue for high-dimensional data analyses, especially for pathway analysis and graphical modeling.

Bickel and Levina [4] gave a thresholding estimator for the inverse matrix of a covariance matrix when the covariance matrix is sparse and its eigenvalues are bounded. However, such sparsity conditions are severe for actual data and often out of touch with reality. In fact, Aoshima and Yata [1, 2, 3] and Yata and Aoshima [8, 9] showed that the first several eigenvalues diverge as  $d$  grows and the bounded-eigenvalues condition is quite strict for microarray data sets. Fan et al. [5] proposed a different thresholding estimator called the principal orthogonal complement thresholding (POET) under the assumption that the first

several eigenvalues diverge rapidly at the rate of  $d$ . Unfortunately, the assumption required in the POET cannot express the structure of actual data.

Suppose we have a  $d \times n$  data matrix  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ , where  $\mathbf{x}_i$ ,  $i = 1, \dots, n$  ( $< d$ ), are independent and identically distributed (i.i.d.) as a  $d$ -dimensional distribution with mean zero and covariance matrix  $\mathbf{\Sigma}$ . We denote the eigen-decomposition of  $\mathbf{\Sigma}$  by  $\mathbf{\Sigma} = \mathbf{H}\mathbf{\Lambda}\mathbf{H}^T$ , where  $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_d)$  with eigenvalues  $\lambda_1 \geq \dots \geq \lambda_d \geq 0$  and  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_d]$  is an orthogonal matrix with eigenvectors  $\mathbf{h}_1, \dots, \mathbf{h}_d$  corresponding to the  $\lambda_1, \dots, \lambda_d$ . Let  $\sigma_{ij}$  be the  $(i, j)$  element of  $\mathbf{\Sigma}$  for  $i, j = 1, \dots, d$ . We assume that  $\sigma_{jj} \in (0, \infty)$  as  $d \rightarrow \infty$  for all  $j$ . For a function,  $f(\cdot)$ , “ $f(d) \in (0, \infty)$  as  $d \rightarrow \infty$ ” implies that  $\liminf_{d \rightarrow \infty} f(d) > 0$  and  $\limsup_{d \rightarrow \infty} f(d) < \infty$ .

Let  $\mathbf{x}_j = \mathbf{H}\mathbf{\Lambda}^{1/2}\mathbf{z}_j$ , where  $\mathbf{z}_j = (z_{1j}, \dots, z_{dj})^T$  is considered as a sphered data vector having the zero mean vector and identity covariance matrix. We assume that

$$\text{(C-i)} \quad \limsup_{d \rightarrow \infty} E(z_{ri}^4) < \infty \text{ for all } r, \text{ and} \\ E(z_{ri}^2 z_{si}^2) = E(z_{ri}^2)E(z_{si}^2) = 1 \text{ and } E(z_{ri} z_{si} z_{ti} z_{ui}) = 0 \text{ for all } r \neq s, t, u.$$

When  $\mathbf{x}_j$ s are Gaussian, (C-i) naturally holds.

The sample covariance matrix is given by  $\mathbf{S} = n^{-1}\mathbf{X}\mathbf{X}^T$ . Let  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_d \geq 0$  be the eigenvalues of  $\mathbf{S}$ . Then, we denote the eigen-decomposition of  $\mathbf{S}$  by

$$\mathbf{S} = \sum_{i=1}^d \hat{\lambda}_i \hat{\mathbf{h}}_i \hat{\mathbf{h}}_i^T,$$

where  $\hat{\mathbf{h}}_i$  is a unit eigenvector corresponding to the  $\hat{\lambda}_i$ . The dual sample covariance matrix is given by  $\mathbf{S}_D = n^{-1}\mathbf{X}^T\mathbf{X}$ . We have the eigen-decomposition of  $\mathbf{S}_D$  by

$$\mathbf{S}_D = \sum_{i=1}^n \hat{\lambda}_i \hat{\mathbf{u}}_i \hat{\mathbf{u}}_i^T,$$

where  $\hat{\mathbf{u}}_i$  is a unit eigenvector corresponding to the  $\hat{\lambda}_i$ .

We assume the following spiked model for the eigenvalues of  $\mathbf{\Sigma}$ :

$$\text{(C-ii)} \quad \frac{\lambda_i}{d^{1/2}} \rightarrow \infty \text{ as } d \rightarrow \infty \text{ for } i = 1, \dots, m, \text{ and } \lambda_i \in (0, \infty) \text{ as } d \rightarrow \infty \text{ for all } i \geq m + 1.$$

Here,  $m$  is a positive and fixed integer. When  $m \geq 2$ ,  $\lambda_1, \dots, \lambda_m$  are distinct in the sense that

$$\liminf_{d \rightarrow \infty} (\lambda_i / \lambda_j - 1) > 0 \text{ for } 1 \leq i < j \leq m.$$

Note that (C-ii) is one of the strongly spiked eigenvalue models given by Aoshima and Yata [1]. See Remark 1.2. We divide  $\mathbf{\Sigma}$  into  $\mathbf{\Sigma}_1 = \sum_{j=1}^m \lambda_j \mathbf{h}_j \mathbf{h}_j^T$  and  $\mathbf{\Sigma}_2 = \sum_{j=m+1}^d \lambda_j \mathbf{h}_j \mathbf{h}_j^T$ , so that  $\mathbf{\Sigma} = \mathbf{\Sigma}_1 + \mathbf{\Sigma}_2$ . Here,  $\mathbf{\Sigma}_1$  is regarded as the signal and  $\mathbf{\Sigma}_2$  is regarded as the noise.

**Remark 1.1.** When we consider a spiked model such as

$$\lambda_i = a_i d^{\alpha_i} \quad (i = 1, \dots, m) \quad \text{and} \quad \lambda_i = c_i \quad (i = m + 1, \dots, d) \quad (1)$$

with positive and fixed constants,  $a_i$ s,  $c_i$ s and  $\alpha_i$ s. Note that (C-ii) is met when  $\alpha_m > 0.5$ . For instance, when we analyze a microarray data set, we find several gene networks and each network consists of genes that are highly correlated to each other. The high correlation is one of the reasons why strongly spiked eigenvalues appear in high-dimensional data analyses.

**Remark 1.2.** Aoshima and Yata [1] showed that the asymptotic normality of high-dimensional statistics cannot be established under the following model called the ‘‘strongly spiked eigenvalue (SSE) model’’:

$$\liminf_{d \rightarrow \infty} \left\{ \frac{\lambda_1}{\text{tr}(\Sigma^2)^{1/2}} \right\} > 0.$$

They gave a data transformation technique from the SSE model to the non-SSE model.

This paper is organized as follows: In Section 2, we consider the POET to construct an estimator of  $\Sigma^{-1}$  and show that the POET has consistency properties under several severe conditions. In Section 3, we introduce the noise-reduction (NR) method that was given by Yata and Aoshima [8, 9]. The NR method is a new PCA having consistency properties for high-dimensional data. In Section 4, we consider applying the NR method to the POET for the inverse matrix estimation. We propose a new estimation of the inverse covariance matrix, called the NR-POET. Finally, in Section 5, we compare the performance of the NR-POET with POET by several simulations.

## 2 POET and its asymptotic properties

In this section, we introduce the principal orthogonal complement thresholding (POET) given by Fan et al. [5]. Let  $\sigma_{2ij}$  be the  $(i, j)$  element of  $\Sigma_2$ . Let  $\tau_{d,h} = \max_{1 \leq i \leq d} \sum_{j=1}^d |\sigma_{2ij}|^h$  for  $0 \leq h \leq 1$ . Here,  $\tau_{d,h}$  is the sparsity measure given by Bickel and Levina [4]. If  $\tau_{d,h}$  is much smaller than  $d$  for a constant  $h \in [0, 1)$ ,  $\Sigma_2$  is considered as sparse in the sense that many elements of  $\Sigma_2$  are very small. We assume  $\limsup_{d \rightarrow \infty} \tau_{d,h} < \infty$  for a constant  $h \in [0, 1)$ . Let  $\lambda_{\min}(\mathbf{M})$  and  $\lambda_{\max}(\mathbf{M})$  be the smallest and largest eigenvalues of any positive definite matrix,  $\mathbf{M}$ . Note that  $\lambda_{\max}(\Sigma_2) = O(\tau_{d,h})$  and  $\limsup_{d \rightarrow \infty} \lambda_{\max}(\Sigma_2) < \infty$ . We assume that  $\liminf_{d \rightarrow \infty} \lambda_{\min}(\Sigma_2) > 0$ .

Let

$$\mathbf{w}_l = (w_{1l}, \dots, w_{dl})^T = \left( \mathbf{I}_d - \sum_{j=1}^m \mathbf{h}_j \mathbf{h}_j^T \right) \mathbf{x}_l \quad \text{for } l = 1, \dots, n$$

and

$$\hat{\mathbf{w}}_l = (\hat{w}_{1l}, \dots, \hat{w}_{dl})^T = \left( \mathbf{I}_d - \sum_{j=1}^m \hat{\mathbf{h}}_j \hat{\mathbf{h}}_j^T \right) \mathbf{x}_l \quad \text{for } l = 1, \dots, n.$$

Let  $I(\cdot)$  be the indicator function. A thresholding operator is defined by

$$T(\mathbf{M}) = [m_{ij}\{I(i=j) + I(i \neq j)I(|m_{ij}| \geq t_{ij})\}]$$

for any symmetric matrix  $\mathbf{M} = (m_{ij})$  and  $t_{ij} > 0$ ,  $i \neq j$ . Fan et al. [5] considered estimating  $\Sigma_2$  by  $T(\widehat{\Sigma}_2)$  with  $\widehat{\Sigma}_2 = \sum_{l=1}^n \hat{\mathbf{w}}_l \hat{\mathbf{w}}_l^T / n$  and

$$t_{ij} = C \sqrt{\hat{\theta}_{ij}} \left( d^{-1/2} + \sqrt{n^{-1} \log d} \right) \text{ for all } i \neq j. \quad (2)$$

Here,  $C(> 0)$  is a sufficiently large constant and  $\hat{\theta}_{ij} = n^{-1} \sum_{s=1}^n (\hat{w}_{is} \hat{w}_{js} - n^{-1} \sum_{t=1}^n \hat{w}_{it} \hat{w}_{jt})^2$ . Then, they gave an estimator of  $\Sigma$  by

$$\widehat{\Sigma} = \sum_{j=1}^m \hat{\lambda}_j \hat{\mathbf{h}}_j \hat{\mathbf{h}}_j^T + T(\widehat{\Sigma}_2). \quad (3)$$

They assumed the following spiked model:

$$(C\text{-ii}') \quad \frac{\lambda_i}{d} \in (0, \infty) \text{ as } d \rightarrow \infty \text{ for } i = 1, \dots, m, \text{ and } \lambda_i \in (0, \infty) \text{ as } d \rightarrow \infty \text{ for all } i \geq m+1.$$

Note that (C-ii') is met when  $\alpha_m = 1$  in (1). Thus (C-ii') is much stricter than (C-ii). Let  $v = \min\{d, n\}$ . We denote the Frobenius and spectral norms by  $\|\cdot\|_F$  and  $\|\cdot\|$ , that is,  $\|\mathbf{M}\|_F = \{\text{tr}(\mathbf{M}^T \mathbf{M})\}^{1/2}$  and  $\|\mathbf{M}\| = \{\lambda_{\max}(\mathbf{M}^T \mathbf{M})\}^{1/2}$  for any  $d \times d$  matrix,  $\mathbf{M}$ . Then, the following result was given by Fan et al. [5].

**Theorem 2.1** ([5]). *Assume (C-i) and (C-ii'). Then, under some regularity conditions, for a sufficiently large constant  $C(> 0)$  it holds that as  $v \rightarrow \infty$*

$$\|\Sigma^{-1/2} \widehat{\Sigma} \Sigma^{-1/2} - \mathbf{I}_d\|_F / \sqrt{d} = o_P(1) \quad \text{and} \quad \|\widehat{\Sigma}^{-1} - \Sigma^{-1}\| = o_P(1).$$

It should be noted that (C-ii') is quite strict for high-dimensional data.

### 3 NR method and its asymptotic properties

In this section, we introduce the noise-reduction (NR) method given by Yata and Aoshima [8, 9]. Let  $\delta_j = \lambda_j^{-1} \sum_{s=m+1}^d \lambda_s / n$  for  $j = 1, \dots, m$ . Aoshima and Yata [1] and Yata and Aoshima [9] gave the following result when  $v \rightarrow \infty$ .

**Proposition 3.1** ([1, 9]). *Assume (C-i) and (C-ii). It holds for  $j = 1, \dots, m$ , that  $\hat{\lambda}_j / \lambda_j = 1 + \delta_j + O_P(n^{-1/2})$  and  $(\hat{\mathbf{h}}_j^T \mathbf{h}_j)^2 = (1 + \delta_j)^{-1} + O_P(n^{-1/2})$  as  $v \rightarrow \infty$ .*

If  $\delta_j \rightarrow \infty$  as  $v \rightarrow \infty$ ,  $\hat{\lambda}_j$  and  $\hat{\mathbf{h}}_j$  are strongly inconsistent in the sense that  $\lambda_j / \hat{\lambda}_j = o_P(1)$  and  $(\hat{\mathbf{h}}_j^T \mathbf{h}_j)^2 = o_P(1)$ . See Jung and Marron [6] for the concept of the strong inconsistency. In order to overcome the curse of dimensionality, Yata and Aoshima [8, 9] proposed an eigenvalue

estimation called the noise-reduction (NR) method, which was brought about by a geometric representation of  $\mathbf{S}_D$ . If one applies the NR method, the  $\lambda_j$  is estimated by

$$\tilde{\lambda}_j = \hat{\lambda}_j - \frac{\text{tr}(\mathbf{S}_D) - \sum_{l=1}^j \hat{\lambda}_l}{n-j} \quad (j = 1, \dots, n-1). \quad (4)$$

Note that  $\tilde{\lambda}_j \geq 0$  w.p.1 for  $j = 1, \dots, n-1$ , and the second term in (4) is an estimator of  $\lambda_j \delta_j$ . When applying the NR method to the PC direction vector, one obtains

$$\tilde{\mathbf{h}}_j = (n\tilde{\lambda}_j)^{-1/2} \mathbf{X} \hat{\mathbf{u}}_j$$

for  $j = 1, \dots, n-1$ . Then, we have the following result.

**Proposition 3.2** ([1, 9]). *Assume (C-i) and (C-ii). It holds for  $j = 1, \dots, k$ , that  $\tilde{\lambda}_j/\lambda_j = 1 + O_P(n^{-1/2})$  and  $(\tilde{\mathbf{h}}_j^T \mathbf{h}_j)^2 = 1 + O_P(n^{-1})$  as  $v \rightarrow \infty$ .*

Thus,  $\tilde{\lambda}_j$  and  $\tilde{\mathbf{h}}_j$  have the consistency properties even when  $\delta_j \rightarrow \infty$  and (C-ii) is met.

**Remark 3.1.** Wang and Fan [7] proposed the following estimator of  $\Sigma$  by using the NR method:

$$\hat{\Sigma} = \sum_{j=1}^m \tilde{\lambda}_j \hat{\mathbf{h}}_j \hat{\mathbf{h}}_j^T + T(\hat{\Sigma}_2). \quad (5)$$

They called this estimation method the Shrinkage-POET (S-POET).

For estimating  $\mathbf{x}_l^T \mathbf{h}_j$ , Aoshima and Yata [1] showed that  $\mathbf{x}_l^T \hat{\mathbf{h}}_j$  and even  $\mathbf{x}_l^T \tilde{\mathbf{h}}_j$  involve a huge bias. In order to overcome the inconvenience, Aoshima and Yata [1] gave the modified NR method. According to [1], we modify  $\mathbf{x}_l^T \tilde{\mathbf{h}}_j$  as  $\mathbf{x}_l^T \tilde{\mathbf{h}}_{jl}$  by

$$\tilde{\mathbf{h}}_{jl} = \left( \frac{n}{n-1} \right) \frac{\mathbf{X} \hat{\mathbf{u}}_{jl}}{(n\tilde{\lambda}_j)^{1/2}} = \frac{n^{1/2} \mathbf{X} \hat{\mathbf{u}}_{jl}}{(n-1)\tilde{\lambda}_j^{1/2}},$$

where

$$\hat{\mathbf{u}}_{jl} = (\hat{u}_{j1}, \dots, \hat{u}_{jl-1}, 0, \hat{u}_{jl+1}, \dots, \hat{u}_{jn})^T.$$

Note that  $\sum_{l=1}^n \hat{\mathbf{u}}_{jl}/n = \{(n-1)/n\} \hat{\mathbf{u}}_j$ . We give the following R-code to calculate  $\mathbf{x}_l^T \tilde{\mathbf{h}}_{jl}$ :

```

xTh <- function(X, m, MeanZero=F){
  d <- dim(X)[1]
  n <- dim(X)[2]
  if (MeanZero){
    r <- min(n-1, d)
    Sd <- t(X) %*% X / n
    eig <- eigen(Sd)
    dualval <- eig$values[1:m]
    dualvec <- eig$vectors
    ans <- matrix(0, n, m)
    c <- sqrt(n) / (n-1)
    for (i in 1:m){
      nrmval <- dualval[i] - (sum(diag(Sd)) - sum(dualval[1:i])) / (n-i)
      u_hat <- dualvec[, i]
      for (j in 1:n){
        u_hat[j] <- 0
        nrmvec_self <- c * X %*% u_hat / sqrt(nrmval)
        ans[j, i] <- as.numeric(t(nrmvec_self) %*% X[, j])
        u_hat <- dualvec[, i]
      }
    }
  } else {
    r <- min(n-2, d)
    X <- sweep(X, 1, apply(X, 1, mean), '-')
    Sd <- t(X) %*% X / (n-1)
    eig <- eigen(Sd)
    dualval <- eig$values[1:m]
    dualvec <- eig$vectors
    ans <- matrix(0, n, m)
    c <- sqrt(n-1) / (n-2)
    for (i in 1:m){
      nrmval <- dualval[i] - (sum(diag(Sd)) - sum(dualval[1:i])) / (n-i-1)
      u_hat <- dualvec[, i]
      for (j in 1:n){
        u_hat[j] <- 0
        nrmvec_self <- c * X %*% u_hat / sqrt(nrmval)
        ans[j, i] <- as.numeric(t(nrmvec_self) %*% X[, j])
        u_hat <- dualvec[, i]
      }
    }
  }
  return(ans)
}

```

## 4 NR-POET

In this section, we propose a new estimation of  $\Sigma^{-1}$  by applying the NR method to the POET. By using the modified NR method in Section 3, we estimate  $\mathbf{w}_l$  by

$$\tilde{\mathbf{w}}_l = (\tilde{w}_{1l}, \dots, \tilde{w}_{dl})^T = \left( \mathbf{I}_d - \sum_{j=1}^m \tilde{\mathbf{h}}_j \tilde{\mathbf{h}}_j^T \right) \mathbf{x}_l.$$

We consider estimating  $\Sigma_2$  by  $T(\tilde{\Sigma}_2)$ , where  $\tilde{\Sigma}_2 = \sum_{l=1}^n \tilde{\mathbf{w}}_l \tilde{\mathbf{w}}_l^T / n$  and

$$t_{ij} = C \sqrt{\tilde{\theta}_{ij}} \left( d^{-1/2} + \sqrt{n^{-1} \log d} \right) \text{ for all } i \neq j \quad (6)$$

with  $\tilde{\theta}_{ij} = n^{-1} \sum_{s=1}^n (\tilde{w}_{is} \tilde{w}_{js} - n^{-1} \sum_{t=1}^n \tilde{w}_{it} \tilde{w}_{jt})^2$ . We denote the eigen-decomposition of  $T(\tilde{\Sigma}_2)$  by

$$T(\tilde{\Sigma}_2) = \sum_{j=1}^d \acute{\lambda}_j \acute{\mathbf{h}}_j \acute{\mathbf{h}}_j^T,$$

where  $\acute{\lambda}_j$ s are eigenvalues of  $T(\tilde{\Sigma}_2)$  having  $\acute{\lambda}_1 \geq \dots \geq \acute{\lambda}_d \geq 0$  and  $\acute{\mathbf{h}}_j$  is a unit eigenvector corresponding to the  $\acute{\lambda}_j$ . Note that

$$\Sigma^{-1} = \sum_{j=1}^m \lambda_j^{-1} \mathbf{h}_j \mathbf{h}_j^T + \sum_{j=m+1}^d \lambda_j^{-1} \mathbf{h}_j \mathbf{h}_j^T.$$

Finally, by applying the NR method to the signal part and the noise part, we propose to estimate  $\Sigma^{-1}$  by

$$\tilde{\Sigma}^{-1} = \sum_{j=1}^m \tilde{\lambda}_j^{-1} \tilde{\mathbf{h}}_j \tilde{\mathbf{h}}_j^T + \left( \mathbf{I}_d - \sum_{j=1}^m \tilde{\mathbf{h}}_j \tilde{\mathbf{h}}_j^T \right) \left( \sum_{j=1}^{d-m} \acute{\lambda}_j^{-1} \acute{\mathbf{h}}_j \acute{\mathbf{h}}_j^T \right) \left( \mathbf{I}_d - \sum_{j=1}^m \tilde{\mathbf{h}}_j \tilde{\mathbf{h}}_j^T \right). \quad (7)$$

Note that the signal and noise parts are orthogonal. We call this new estimation method the ‘‘NR-POET’’. In the next section, we compare the performance of the NR-POET with the POET by several simulations.

## 5 Simulation studies

In this section, we compare the performance of the NR-POET with the POET and S-POET by using computer simulations. We considered the following covariance matrix:

$$\Sigma = \begin{pmatrix} \mathbf{\Gamma}_{d(1),\sigma(1)} & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{\Gamma}_{d(2),\sigma(2)} & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{\Omega}_{d(3)}(\rho) \end{pmatrix},$$

where  $d_{(1)} + d_{(2)} + d_{(3)} = d$ ,  $\mathbf{\Gamma}_{t,\sigma} = \sigma(\mathbf{I}_t + \mathbf{1}_t \mathbf{1}_t^T)$  with  $\mathbf{1}_t = (1, \dots, 1)^T$  and  $\mathbf{\Omega}_t(\rho) = \mathbf{B}_t(\rho^{|i-j|^{1/3}}) \mathbf{B}_t$  with  $\mathbf{B}_t = \text{diag}[\{0.5 + 1/(t+1)\}^{1/2}, \dots, \{0.5 + t/(t+1)\}^{1/2}]$ . Note that  $\lambda_{\max}(\mathbf{\Gamma}_{t,\sigma}) = \sigma(t+1)$  and its other eigenvalues are  $\sigma$ . Also, note that  $[\lambda_{\max}\{\mathbf{\Omega}_t(\rho)\}]^2 / \text{tr}\{\{\mathbf{\Omega}_t(\rho)\}^2\} = o(1)$  as  $t \rightarrow \infty$  for  $|\rho| < 1$ . We set  $(d_{(1)}, d_{(2)}) = (\lceil d/3 \rceil, \lceil d/3 \rceil)$ , where  $\lceil x \rceil$  denotes the smallest integer  $\geq x$ . We set  $(\sigma_{(1)}, \sigma_{(2)}) = (1, d^{-1/3})$ . Note that  $(\lambda_1, \lambda_2) \approx (d/3, d^{2/3}/3)$ , so that (C-ii) is satisfied while (C-ii') is not. We set  $d = 200(200)1400$ ,  $n = \lceil d^{3/5} \rceil$  and  $\rho = 0.3$ .

We fixed  $C = 5$  in (2) and (6). We considered two cases:

(S-i)  $\mathbf{x}_i$ s are generated from  $N_d(\mathbf{0}, \mathbf{\Sigma})$ ;

(S-ii)  $\mathbf{x}_i$ s are generated from  $z_{ri} = (y_{ri} - 2)/2$  ( $r = 1, \dots, d$ ) in which  $y_{ri}$ s are i.i.d as the chi-squared distribution with 2 degrees of freedom.

For each case, we estimated  $\mathbf{\Sigma}^{-1}$  by using the POET in (3), S-POET in (5) and NR-POET in (7). We calculated the average loss and its standard deviation by 2000 ( $= R$ , say) times replications for each estimator. Under a fixed scenario, let  $\widehat{\mathbf{\Sigma}}_r^{-1}$  be the  $r$ -th estimation of  $\mathbf{\Sigma}^{-1}$ . Let  $\|\mathbf{M}\|_{\Sigma} = \|\mathbf{\Sigma}^{1/2} \mathbf{M} \mathbf{\Sigma}^{1/2}\|_F / \sqrt{d}$ . In the above simulations,  $E\|\widehat{\mathbf{\Sigma}}_r^{-1} - \mathbf{\Sigma}^{-1}\|_{\Sigma}$  was estimated by  $R^{-1} \sum_{r=1}^R \|\widehat{\mathbf{\Sigma}}_r^{-1} - \mathbf{\Sigma}^{-1}\|_{\Sigma}$ . We denote the standard deviation of  $\|\widehat{\mathbf{\Sigma}}_r^{-1} - \mathbf{\Sigma}^{-1}\|_{\Sigma}$  by  $SD\|\widehat{\mathbf{\Sigma}}_r^{-1} - \mathbf{\Sigma}^{-1}\|_{\Sigma}$ . In the above simulations,  $SD\|\widehat{\mathbf{\Sigma}}_r^{-1} - \mathbf{\Sigma}^{-1}\|_{\Sigma}$  was estimated by the standard deviation of  $\|\widehat{\mathbf{\Sigma}}_r^{-1} - \mathbf{\Sigma}^{-1}\|_{\Sigma}$ ,  $r = 1, \dots, R$ . We displayed the results for (S-i) in Figure 1 and for (S-ii) in Figure 2. We observed that the NR-POET gave better performances than the POET and S-POET both for (S-i) and (S-ii).

## Acknowledgements

Research of the second author was partially supported by Grant-in-Aid for Scientific Research (C), Japan Society for the Promotion of Science (JSPS), under Contract Number 18K03409. Research of the third author was partially supported by Grants-in-Aid for Scientific Research (A) and Challenging Research (Exploratory), JSPS, under Contract Numbers 15H01678 and 19K22837. This work was supported by the Research Institute for Mathematical Sciences, an International Joint Usage/Research Center located in Kyoto University.



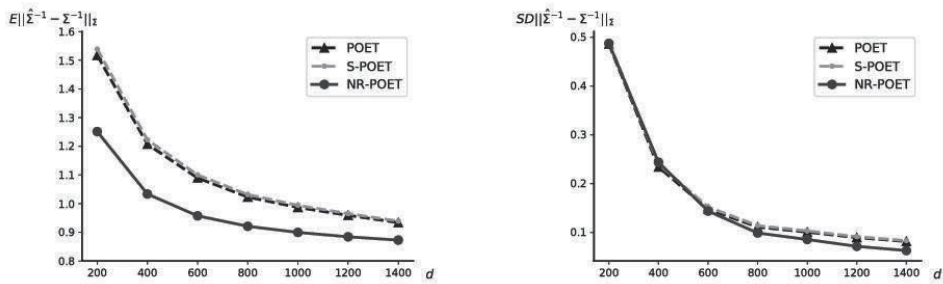


Figure 1: For Case (S-i), the left panel displays the average loss and the right panel displays its standard deviation.

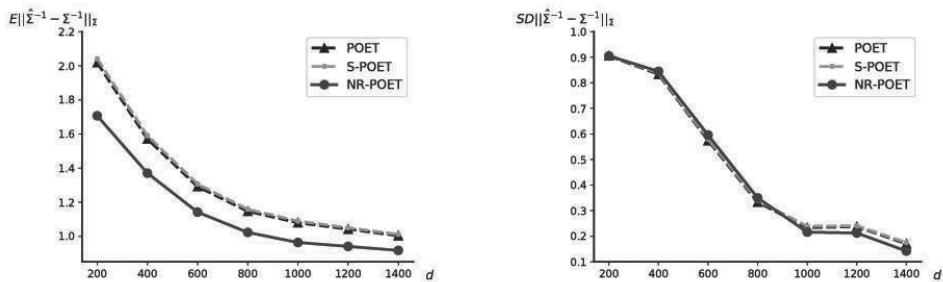


Figure 2: For Case (S-ii), the left panel displays the average loss and the right panel displays its standard deviation.

## References

- [1] Aoshima, M. and Yata, K. (2018). Two-sample tests for high-dimension, strongly spiked eigenvalue models. *Statistica Sinica*, **28**, 43-62.
- [2] Aoshima, M. and Yata, K. (2019a). Distance-based classifier by data transformation for high-dimension, strongly spiked eigenvalue models. *Annals of the Institute of Statistical Mathematics*, **71**, 473-503.
- [3] Aoshima, M. and Yata, K. (2019b). High-dimensional quadratic classifiers in non-sparse settings. *Methodology and Computing in Applied Probability*, **21**, 663-682.
- [4] Bickel, P. and Levina, E. (2008). Covariance regularization by thresholding. *Annals of Statistics*, **36**, 2577-2604.

- [5] Fan, J., Liao, Y. and Mincheva, M. (2013). Large covariance estimation by thresholding principal orthogonal complements. *Journal of the Royal Statistical Society, Series B.*, **75**, 603-80.
- [6] Jung, S. and Marron, J.S. (2009). PCA consistency in high dimension, low sample size context. *Annals of Statistics*, **37**, 4104-4130.
- [7] Wang, W. and Fan, J. (2017). Asymptotics of empirical eigenstructure for high dimensional spiked covariance. *Annals of Statistics*, **45**, 1342-1374.
- [8] Yata, K. and Aoshima, M. (2012). Effective PCA for high-dimension, low-sample-size data with noise reduction via geometric representations. *Journal of Multivariate Analysis*, **105**, 193-215.
- [9] Yata, K. and Aoshima, M. (2013). PCA consistency for the power spiked model in high-dimensional settings. *Journal of Multivariate Analysis*, **122**, 334-354.