

# 高次元固有ベクトルの一貫性

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

## 1 はじめに

情報化の進展に伴い、高次元データの統計解析が、ますます重要になってきている。2000年以降、確率論と理論物理の方面から、ランダム行列の理論に基づく幾つかの重要な結果がもたらされた。Johnstone [7] や Paul [8] 等は、標本固有値の漸近分布を導出した。しかしながら、ここでは、データの次元数  $d$  と標本数  $n$  が  $n/d \rightarrow c > 0$  を満たす場合を考え、高次元において標本数は次元数と同程度を仮定した。例えば、次元数は優に 10,000 を超えるが標本数は高々 100 程度といった高次元小標本においては、標本数を次元数と同程度には仮定できない。それゆえ、 $n$  が  $d$  に依存しないような設定で、もしくは、 $n = n(d)$  であっても  $n/d \rightarrow 0$  となる設定で、高次元漸近理論を展開する必要がある。Yata and Aoshima [10] は、高次元小標本における PCA の性質を研究し、PCA が一貫性をもつための標本数  $n$  の  $d$  に関するオーダー条件を導き、高次元小標本において PCA が不適解を起こすことを示した。この問題を解決する策として、Yata and Aoshima [12] は、高次元小標本データ空間の幾何学的表現を研究し、それに基づいて“ノイズ掃き出し法”とよばれる方法論を考案した。一方で、Yata and Aoshima [13] は、高次元大標本も含む一般的な高次元データに対して、power spiked モデルと呼ばれる固有値モデルを考案し、高次元データに対する新しい PCA を構築した。最近、Aoshima and Yata [5] は、ノイズ掃き出し法による固有ベクトルの推定量を用いることで、新たな高次元二標本検定法を考案した。

本稿では、高次元固有ベクトルの一貫性について論じる。ノイズ掃き出し法による固有ベクトルの推定量を補正することで、緩い仮定のもとその一貫性を与える新たな

方法論を提案する。

## 2 高次元固有値の一致性

平均に  $d$  次のゼロベクトル, 共分散行列に  $d$  次の半正定値行列  $\Sigma$  をもつ母集団を考える. 母集団から  $n$  ( $\geq 2$ ) 個の  $d$  次データベクトル  $\mathbf{x}_1, \dots, \mathbf{x}_n$  を無作為に抽出して,  $d \times n$  データ行列  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$  を定義する. ただし,  $d > n$  と仮定する.  $\Sigma$  の固有値を  $\lambda_1 \geq \dots \geq \lambda_d (\geq 0)$  とし, 適当な直交行列  $\mathbf{H} = [\mathbf{h}_1, \dots, \mathbf{h}_d]$  で  $\Sigma$  を

$$\Sigma = \mathbf{H}\mathbf{\Lambda}\mathbf{H}^T, \quad \mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_d)$$

と分解する. そのとき  $\mathbf{Z} = \mathbf{\Lambda}^{-1/2}\mathbf{H}^T\mathbf{X}$  とおき,  $\mathbf{Z} = [z_1, \dots, z_d]^T$ ,  $z_s = (z_{s1}, \dots, z_{sn})^T$  と表記する. ただし,  $\mathbf{Z}$  の成分は, 4 次モーメントに一様有界性を仮定する. 標本共分散行列  $\mathbf{S} = \mathbf{X}\mathbf{X}^T/n$  の固有値を  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_d (\geq 0)$ ,  $\hat{\lambda}_j$  に対する固有ベクトルを  $\hat{\mathbf{h}}_j$  として, スペクトル分解を

$$\mathbf{S} = \sum_{s=1}^d \hat{\lambda}_s \hat{\mathbf{h}}_s \hat{\mathbf{h}}_s^T$$

とおく. 最近, Yata and Aoshima [13] は, power spiked モデルとよばれる固有値モデルを考案し, 高次元データに対する新しい PCA を研究した. いま,  $\Sigma_{(1)} = \sum_{s=1}^m \lambda_s \mathbf{h}_s \mathbf{h}_s^T$ ,  $\Sigma_{(2)} = \sum_{s=m+1}^d \lambda_s \mathbf{h}_s \mathbf{h}_s^T$  とおき,  $\Sigma = \Sigma_{(1)} + \Sigma_{(2)}$  という分解を考える. ただし,  $m < \infty$ . そのとき, 次の条件を満たすような  $\lambda_1 \geq \dots \geq \lambda_d$  を power spiked モデルと定義する.

$\lambda_m$  に対して,  $\lim_{d \rightarrow \infty} \text{tr}(\Sigma_{(2)}^{k_m}) / \lambda_m^{k_m} = 0$  なる (有界な) ある自然数  $k_m$  が存在する.

本稿では簡単のため,  $k_m = 2$  の場合を考える. すなわち,

$$\lim_{d \rightarrow \infty} \frac{\text{tr}(\Sigma_{(2)}^2)}{\lambda_j^2} = 0, \quad j = 1, \dots, m \quad (1)$$

なる spiked モデルを仮定する. いま,

$$\delta_j = \frac{\text{tr}(\Sigma_{(2)})}{n\lambda_j}, \quad j = 1, \dots, m$$

とおく. (1) のもと, 次の定理を得る.

定理 1 ([13]). 各  $j = 1, \dots, m$  について, 条件

$$(C-i) \quad \frac{\sum_{s,t=m+1}^d \lambda_s \lambda_t E\{(z_{sk}^2 - 1)(z_{tk}^2 - 1)\}}{n\lambda_j^2} = o(1)$$

のもと,  $d, n \rightarrow \infty$  のとき次が成り立つ.

$$\frac{\hat{\lambda}_j}{\lambda_j} = 1 + \delta_j + o_p(1).$$

注意. もし,  $z_{1k}, \dots, z_{dk}$  が互いに独立ならば, (1)のもと (C-i) を満たす.

定理 1 より,  $\delta_j \rightarrow \infty$ ,  $d, n \rightarrow \infty$  のとき,  $\hat{\lambda}_j$  は “ $\lambda_j/\hat{\lambda}_j = o_p(1)$ ” なる強不一致性をもつ. 一方で, Yata and Aoshima [12] は, 高次元小標本データ空間の幾何学的表現を研究し, それに基づいて “ノイズ掃き出し法” とよばれる方法論を考案し, 次のような固有値の推定量を提案した.

$$\tilde{\lambda}_j = \hat{\lambda}_j - \frac{\text{tr}(\mathbf{S}) - \sum_{i=1}^j \hat{\lambda}_i}{n-j} \quad (j = 1, \dots, n-1). \quad (2)$$

そのとき, (1)のもと, 次の定理を得る.

定理 2 ([13]). 各  $j = 1, \dots, m$  について, (C-i)のもと,  $d, n \rightarrow \infty$  のとき次が成り立つ.

$$\frac{\tilde{\lambda}_j}{\lambda_j} = 1 + o_p(1).$$

定理 2 より,  $\delta_j \rightarrow \infty$ ,  $d, n \rightarrow \infty$  の場合においても,  $\tilde{\lambda}_j$  は一致性をもつ.

一方で, (C-i) が仮定できない場合, Yata and Aoshima [11] は, 母集団分布の仮定を必要としないクロスデータ行列法という方法論を考案した. 詳細は Aoshima and Yata [1] や青嶋・矢田 [2, 3] を参照されたい.

### 3 高次元固有ベクトルの一致性

$\Sigma$  の固有ベクトルについて、ノイズ掃き出し法による推定を考える。推定量 (2) に基づいて、 $\Sigma$  の固有ベクトル  $\mathbf{h}_j$  を

$$\tilde{\mathbf{h}}_j = \sqrt{\frac{\hat{\lambda}_j}{\tilde{\lambda}_j}} \hat{\mathbf{h}}_j$$

で推定する。ただし、 $\mathbf{h}_j$  には符号の自由度があるため、各  $j$  で  $\tilde{\mathbf{h}}_j^T \mathbf{h}_j \geq 0$  を仮定する。ここで、 $\|\tilde{\mathbf{h}}_j\|^2 = \hat{\lambda}_j / \tilde{\lambda}_j > 1$  であることに注意する。ただし、 $\|\cdot\|$  はユークリッドノルムを表す。(1) のもと次の定理を得る。

定理 3 ([13]). 各  $j = 1, \dots, m$  について、(C-i) と条件

$$(C\text{-ii}) \quad \liminf_{d \rightarrow \infty} \frac{\lambda_j}{\lambda_{j'}} > 1 \quad \text{for } j < j' \leq m$$

のもと、 $d, n \rightarrow \infty$  のとき次が成り立つ。

$$\mathbf{h}_j^T \hat{\mathbf{h}}_j = (1 + \delta_j)^{-1/2} + o_p(1) \quad \text{and} \quad \mathbf{h}_j^T \tilde{\mathbf{h}}_j = 1 + o_p(1).$$

それゆえ、 $\tilde{\mathbf{h}}_j$  は ( $\|\tilde{\mathbf{h}}_j\|^2 > 1$  であるが)  $\mathbf{h}_j$  の内積に関する一致性をもつ。ノイズ掃き出し法による推定量  $\tilde{\mathbf{h}}_j$  は、例えば、Aoshima and Yata [4, 5] では高次元二標本検定と高次元判別分析に応用され、Ishii et al. [6] では高次元共分散行列の同等性検定に応用されている。

ここで、定理 1 から 3 より、(C-i) と (C-ii) のもと、 $d, n \rightarrow \infty$  のとき次を得る。

$$\|\hat{\mathbf{h}}_j - \mathbf{h}_j\|^2 = 2\{1 - (1 + \delta_j)^{-1/2}\} + o_p(1) \quad \text{and} \quad \|\tilde{\mathbf{h}}_j - \mathbf{h}_j\|^2 = \delta_j + o_p(1).$$

すなわち、 $\liminf_{d, n \rightarrow \infty} \delta_j > 0$  のとき、 $\hat{\mathbf{h}}_j$  と  $\tilde{\mathbf{h}}_j$  はノルムに関する一致性をもたない。

### 4 高次元固有ベクトルのスパース推定

本節では、ノルムに関する一致性をもつように、 $\tilde{\mathbf{h}}_1$  を補正する。いま、 $\mathbf{h}_1 = (h_1, \dots, h_d)^T$  とおく。さらに、

$$D = \{s \mid h_s \neq 0, s = 1, \dots, d\}$$

とおく. そのとき, 次のモデルを仮定する.

$$\liminf_{d,n \rightarrow \infty} n^{1/2} |h_s| > 0 \text{ for all } s \in D. \quad (3)$$

ここで,  $(\sum_{s=2}^d \mathbf{h}_s \mathbf{h}_s^T) \mathbf{x}_j = (y_{j1}, \dots, y_{jd})^T$  とおく. そのとき, 各  $y_{js}$  の4次モーメントが有界であることを仮定する. いま,  $\tilde{\mathbf{h}}_1 = (\tilde{h}_1, \dots, \tilde{h}_d)^T$  とおき,  $\tilde{h}_1, \dots, \tilde{h}_d$  を絶対値の大きい順に並べ替えたものを  $\tilde{h}_{(1)}, \dots, \tilde{h}_{(d)}$  とおく. すなわち,  $|\tilde{h}_{(1)}| \geq \dots \geq |\tilde{h}_{(d)}|$  となる.  $\|\tilde{\mathbf{h}}_1\|^2 > 1$  であることに注意すれば,

$$\sum_{s=1}^{k-1} \tilde{h}_{(s)}^2 < 1 \quad \text{and} \quad \sum_{s=1}^k \tilde{h}_{(s)}^2 \geq 1$$

となる  $k$  が一意に定まる. そのとき,  $\tilde{\mathbf{h}}_1$  を次のようにスパース化する.

$$\hat{\mathbf{h}}_1 = (\hat{h}_1, \dots, \hat{h}_d)^T.$$

ただし,

$$\hat{h}_s = \begin{cases} \tilde{h}_s & (|\tilde{h}_s| \geq |\tilde{h}_{(k)}|) \\ 0 & (|\tilde{h}_s| < |\tilde{h}_{(k)}|) \end{cases}, \quad (s = 1, \dots, d)$$

とする. そのとき, (1) と (3) のもと次の定理を得る.

定理 4.  $j = 1$  に対して (C-i) と (C-ii) を仮定する. 条件

$$(C\text{-iii}) \quad d/\lambda_1^2 \rightarrow 0, \quad d \rightarrow \infty$$

のもと,  $d, n \rightarrow \infty$  のとき次が成り立つ.

$$\|\hat{\mathbf{h}}_1 - \mathbf{h}_1\|^2 = o_p(1). \quad (4)$$

注意.  $\mathbf{h}_j$  ( $j \geq 2$ ) についてもノルムに関する一致性を証明できるが, 本稿では割愛する. Shen et al. [9] はあるチューニングパラメータを用いた  $\mathbf{h}_1$  のスパースな推定量を与え, その高次元一致性を議論した. しかしながら, その推定量がチューニングパラメータに大きく依存することに注意する. 一方で,  $\hat{\mathbf{h}}_1$  はそのようなパラメータに依存せず, 自動的にスパースな推定量を与えることができる.

## 5 シミュレーション

本節では、高次元小標本のもとで、 $\hat{h}_1$  と  $\hat{h}_1$  の精度を数値的に検証する。母集団分布には、 $d$  次元正規分布  $N_d(\mathbf{0}, \Sigma)$  を考え、次の 2 つの設定を考える。

(a)  $d = 2^s$ ,  $s = 6, \dots, 11$ ,  $n = \lceil d^{1/3} \rceil$  とおく。ただし、 $\lceil x \rceil$  は  $x$  以上の最小の整数を表す。 $\Sigma = \text{diag}(d^{2/3}, 1, \dots, 1)$  とおく。すなわち、 $\mathbf{h}_1 = (1, 0, \dots, 0)$  である。

(b)  $d = 500$ ,  $n = 2^s$ ,  $s = 3, \dots, 8$  とおく。

$$\Sigma = \begin{pmatrix} \Gamma_{\lceil d^{2/3} \rceil} & \mathbf{O} \\ \mathbf{O} & I_{d - \lceil d^{2/3} \rceil} \end{pmatrix} (= \Sigma_b)$$

とおく。ただし、 $\Gamma_t = I_t + \mathbf{1}_t \mathbf{1}_t^T$  であり、 $I_t$  は  $t$  次の単位行列である。このとき、 $\lambda_1 \approx d^{2/3}$  であり、最初の  $\lceil d^{2/3} \rceil$  個の成分が非ゼロである  $\mathbf{h}_1 = (\lceil d^{2/3} \rceil^{-1/2}, \dots, \lceil d^{2/3} \rceil^{-1/2}, 0, \dots, 0)^T$  となる。

設定 (a) と (b) において、 $A: \|\hat{h}_1 - \mathbf{h}_1\|^2$  と  $B: \|\hat{h}_1 - \mathbf{h}_1\|^2$  をそれぞれ 1000 回発生させ、その平均を図 1 と図 2 にプロットした。

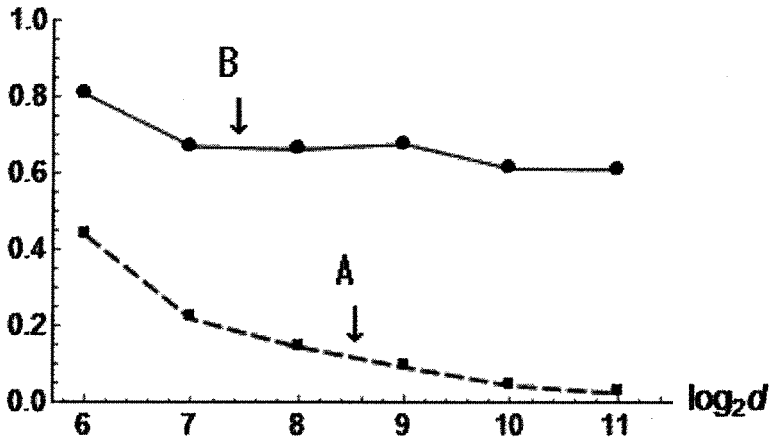


図 1 (a)  $d = 2^s$ ,  $s = 6, \dots, 11$ ,  $n = \lceil d^{1/3} \rceil$ ,  $\Sigma = \text{diag}(d^{2/3}, 1, \dots, 1)$ .

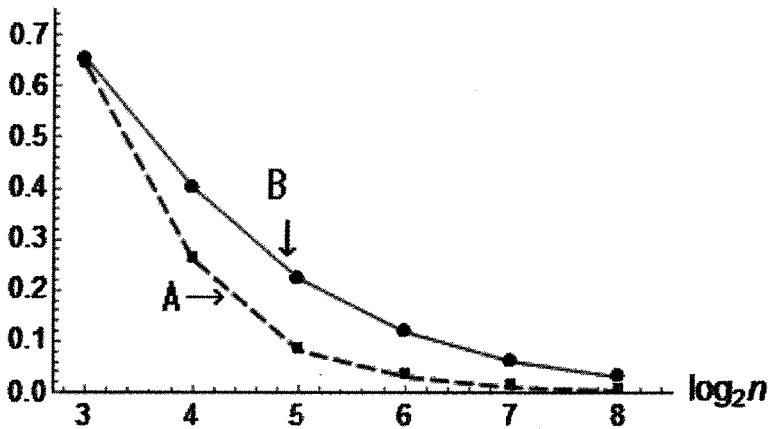


図2 (b)  $d = 500, n = 2^s, s = 3, \dots, 8, \Sigma = \Sigma_b$ .

これらの図からも分かるように、スパースな推定量  $\hat{h}_1$  が従来の推定量  $\hat{h}_1$  に比べ、高次元小標本のもと非常に良い推定量となっている。ここでは割愛するが、設定を変えて実験をしたときにも、 $\hat{h}_1$  が  $\hat{h}_1$  に比べ、高次元小標本のもと優れた推定量となっていることが確認されている。

## 6 定理4の証明

まず、 $j = 1$  に対して条件 (C-i), (C-ii) と (C-iii) を仮定する。いま、 $S_D = n^{-1} \mathbf{X}^T \mathbf{X}$  とおく。  $S_D$  は  $S$  と正の固有値を共有する双対標本共分散行列という。  $S_D$  の固有値を  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_n$  とし、  $\hat{\lambda}_j$  に対する固有ベクトルを  $\hat{u}_j$  として、スペクトル分解を

$$S_D = \sum_{s=1}^n \hat{\lambda}_s \hat{u}_s \hat{u}_s^T$$

とおく。そのとき、

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

と書ける。ここで、Yata and Aoshima [13] の補題9より、 $d, n \rightarrow \infty$  のとき次が成り立つ。

$$\hat{u}_1^T \frac{z_1/n^{1/2}}{\|z_1/n^{1/2}\|} = 1 + o_p(n^{-1/2}).$$

すなわち,

$$\hat{\mathbf{u}}_1 = \{1 + o_p(1)\} \mathbf{z}_1 / n^{1/2} + \boldsymbol{\omega}$$

と書ける. ただし,  $\mathbf{z}_1^T \boldsymbol{\omega} = 0$ ,  $\|\boldsymbol{\omega}\|^2 = o_p(n^{-1/2})$  である. いま,  $\mathbf{y}_{(s)} = (y_{1s}, \dots, y_{ns})^T$ ,  $s = 1, \dots, d$  とおく. そのとき, 定理 2 より次が成り立つ.

$$\begin{aligned} \tilde{\mathbf{h}}_1 &= \{1 + o_p(1)\} \mathbf{h}_1 + \frac{\{1 + o_p(1)\}}{n\lambda_1^{1/2}} (\mathbf{y}_{(1)}^T \mathbf{z}_1, \dots, \mathbf{y}_{(d)}^T \mathbf{z}_1)^T \\ &\quad + \frac{\{1 + o_p(1)\}}{n^{1/2}\lambda_1^{1/2}} (\mathbf{y}_{(1)}^T \boldsymbol{\omega}, \dots, \mathbf{y}_{(d)}^T \boldsymbol{\omega})^T. \end{aligned} \quad (5)$$

ここで, マルコフの不等式を用いると, 任意の  $\tau > 0$  について次が成り立つ.

$$\begin{aligned} \sum_{s=1}^d P\left(|\mathbf{y}_{(s)}^T \mathbf{z}_1| / (n\lambda_1^{1/2}) > \tau n^{-1/2}\right) &= \sum_{s=1}^d P\left(|\mathbf{y}_{(s)}^T \mathbf{z}_1|^4 / (n\lambda_1)^2 > \tau^4\right) \\ &= O(d/\lambda_1^2) \rightarrow 0. \end{aligned} \quad (6)$$

さらに, 各  $s$  で  $\sigma_{(s)} = E(y_{js}^2)$  とおき, 任意の  $\tau > 0$  について次が成り立つ.

$$\begin{aligned} &\sum_{s=1}^d P\left(|(\|\mathbf{y}_{(s)}\|^2/n - \sigma_{(s)})/\lambda_1| > \tau n^{-1/2}\right) \\ &= \sum_{s=1}^d P\left(|(\|\mathbf{y}_{(s)}\|^2/n - \sigma_{(s)})|^2/\lambda_1^2 > \tau^2 n^{-1}\right) = O(d/\lambda_1^2) \rightarrow 0. \end{aligned}$$

それゆえ,  $d > n$  に注意し, すべての  $s$  について次が成り立つ.

$$\frac{\|\mathbf{y}_{(s)}\|^2}{n\lambda_1} = \frac{\sigma_{(s)}}{\lambda_1} + o_p(n^{-1/2}) = o_p(n^{-1/2}). \quad (7)$$

それゆえ, (5) から (7) より, 次が成り立つ.

$$\tilde{\mathbf{h}}_1 = \{1 + o_p(1)\} \mathbf{h}_1 + \left(o_p(n^{-1/2}), \dots, o_p(n^{-1/2})\right)^T.$$

よって, 仮定 (3) より題意を得る.  $\square$

謝辞 本研究は, 科学研究費補助金 基盤研究 (A) 15H01678 研究代表者: 青嶋 誠「大規模複雑データの理論と方法論の総合的研究」, および, 若手研究 (B) 26800078 研究代表者: 矢田 和善「高次元漸近理論の統一的研究」から研究助成を受けています.



## 参考文献

- [1] Aoshima, M. and Yata, K. (2011). Two-stage procedures for high-dimensional data, *Sequential Analysis (Editor's special invited paper)*, **30**, 356-399.
- [2] 青嶋 誠, 矢田和善 (2013a). 論説：高次元小標本における統計的推測, *数学*, **65**, 225-247.
- [3] 青嶋 誠, 矢田和善 (2013b). 日本統計学会研究業績賞受賞者特別寄稿論文：高次元データの統計的方法論, *日本統計学会誌*, **43**, 123-150.
- [4] Aoshima, M. and Yata, K. (2016). A distance-based classifier for high-dimensional data under the strongly spiked eigenvalue model, submitted.
- [5] Aoshima, M. and Yata, K. (2017). Two-sample tests for high-dimension, strongly spiked eigenvalue models, *Statistica Sinica*, in press (arXiv:1602.02491).
- [6] Ishii, A., Yata, K. and Aoshima, M. (2016). Asymptotic properties of the first principal component and equality tests of covariance matrices in high-dimension, low-sample-size context, *Journal of Statistical Planning and Inference*, **170**, 186-199.
- [7] Johnstone, I.M. (2001). On the distribution of the largest eigenvalue in principal components analysis, *The Annals of Statistics*, **29**, 295-327.
- [8] Paul, D. (2007). Asymptotics of sample eigenstructure for a large dimensional spiked covariance model, *Statistica Sinica*, **17**, 1617-1642.
- [9] Shen, D., Shen, H. and Marron, J.S. (2013). Consistency of sparse PCA in high dimension, low sample size contexts, *Journal of Multivariate Analysis*, **115** 317-333.
- [10] Yata, K. and Aoshima, M. (2009). PCA consistency for non-Gaussian data in high dimension, low sample size context, *Communications in Statistics. Theory and Methods, Special Issue: Honoring Zacks, S. (ed. Mukhopadhyay, N.)*, **38**, 2634-2652.
- [11] Yata, K. and Aoshima, M. (2010). Effective PCA for high-dimension, low-sample-size data with singular value decomposition of cross data matrix,

*Journal of Multivariate Analysis*, **101**, 2060-2077.

- [12] 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.
- [13] Yata, K. and Aoshima, M. (2013). PCA consistency for the power spiked model in high-dimensional settings, *Journal of Multivariate Analysis*, **122**, 334-354.