

# Stein's unbiased risk estimate and adaptive singular value shrinkage for estimation problem of low-rank matrix mean with unknown covariance matrix

Yoshihiko KONNO

Osaka Metropolitan University

Conference  
Mathematical Methods of Modern Statistics 3

27 June- 1st July 2022

- 1 Mean matrix estimation when covariance is known
- 2 Sufficient conditions for Stein's identity
- 3 Case when the covariance matrix is unknown
- 4 Simulation Study

## Mean matrix estimation when covariance is known

- Assume that  $m$ ,  $p$  are positive integers s.t.  $m \geq p$ .
- Let

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_m^T \end{pmatrix}$$

be an  $m \times p$  data matrix whose row vectors are independently distributed as

$$\mathbf{x}_i : p \times 1 \sim N(\xi_i, \sigma^2 I_p), \quad (i = 1, 2, \dots, m)$$

Here  $\Xi^T := (\xi_1, \dots, \xi_m)$  and  $\sigma > 0$  are unknown.

- First consider the problem of estimating  $\Xi$  under a loss function and its risk

$$\mathbf{L}_1(\widehat{\Xi}, \Xi) = \text{tr} \{(\widehat{\Xi} - \Xi)(\widehat{\Xi} - \Xi)^\top\} =: \|\widehat{\Xi} - \Xi\|_F^2$$

and

$$\mathbf{R}_1(\widehat{\Xi}, \Xi) = \mathbb{E}[\mathbf{L}_1(\widehat{\Xi}, \Xi)].$$

- Here  $\widehat{\Xi}$  is an estimator based on  $\mathbf{X}$ .
- $\|\mathbf{A}\|_F$  is the Frobenius norm of a matrix  $\mathbf{A}$ .
- $\text{tr } \mathbf{A}$  and  $\mathbf{A}^\top$  stand for the trace and the transpose of a matrix  $\mathbf{A}$ , respectively.
- Low-rank mean matrix condition, i.e.,

$$\text{rank } \Xi = r < p; \quad r \text{ is unknown.}$$

## Eckart-Young approximation theorem

- Singular Value Decomposition: Decompose  $\mathbf{X}$  as

$$\begin{aligned}\mathbf{X} &= \mathbf{U}\mathbf{L}\mathbf{V}^T; \quad \mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_p), \quad \mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_p) \\ \mathbf{L} &= \mathbf{diag}(\ell_1, \ell_2, \dots, \ell_p) \quad \text{with } \ell_1 \geq \ell_2 \geq \dots \geq \ell_p \geq 0\end{aligned}$$

where  $\mathbf{u}_i \in \mathbb{R}^m$ ,  $\mathbf{v}_i \in \mathbb{R}^p$  ( $i = 1, \dots, p$ ) s.t.  
 $\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}_p$ .

- The total least squares (TLS) pseudo estimator is given by

$$\hat{\Xi}_{TLS} = \sum_{i=1}^r \ell_i \mathbf{u}_i \mathbf{v}_i^T. \quad \Leftrightarrow \quad \hat{\Xi}_{TLS} \in \underset{\Xi: \text{rank } \Xi \leq r}{\text{argmin}} \|\Xi - \mathbf{X}\|_F^2.$$

## A hard-shresholding rule

- Assume that  $\sigma^2$  is **known**.
- Solve

$$\mathbf{SVHT}_{\lambda}(\mathbf{X}) = \underset{\Xi}{\operatorname{argmin}} \left[ \|\Xi - \mathbf{X}\|_{\mathcal{F}}^2 + \lambda \operatorname{rank}(\Xi) \right]$$

where  $\lambda > \mathbf{0}$  is a tuning scalar parameter.

- Then the solution is given by

$$\mathbf{SVHT}_{\lambda}(\mathbf{X}) = \sum_{i=1}^p \ell_i \mathbb{1}\{\ell_i \geq \lambda\} \mathbf{u}_i \mathbf{v}_i^{\top}; \quad \mathbb{1}\{\ell_i \geq \lambda\} = \begin{cases} 1 & \ell_i \geq \lambda \\ 0 & \text{otherwise} \end{cases}$$

- The optimal shresholding is  $\frac{4}{\sqrt{3}} \sqrt{p} \sigma$  when  $p = m$ .

## A soft-thresholding rule

- Cèndes et al. define an adaptive soft-thresholding rule based on SURE:

$$\mathbf{SVT}_\lambda(\mathbf{X}) = \sum_{i=1}^p (\ell_i - \lambda) \mathbb{1}\{\ell_i \geq \lambda\} \mathbf{u}_i \mathbf{v}_i^\top =: \sum_{i=1}^p (\ell_i - \lambda)_+ \mathbf{u}_i \mathbf{v}_i^\top \quad (1)$$

which is obtained from

$$\min_{\mathbf{Y}} \left\{ \|\mathbf{X} - \mathbf{Y}\|_F^2 + (2\lambda) \sum_{j=1}^p \lambda_j \right\} \quad \mathbf{Y} = \mathbf{SVT}_\lambda(\mathbf{X}).$$

- The parameter  $\lambda$  in (1) is selected by minimizing SURE, Stein's unbiased risk estimate for (1).

- Gaussian integration-by-parts (=Stein's identity) and a bit of algebraic calculation lead to

$$R_B(\text{SVT}_\lambda, \Xi) = \mathbb{E}[\text{SURE}(\text{SVT}_\lambda)(\mathbf{X})],$$

$$\begin{aligned} \text{SURE}(\text{SVT}_\lambda)(\mathbf{X}) &= -m\rho\sigma^2 + \sum_{i=1}^p \min\{\ell_i^2, \lambda^2\} \\ &\quad + 2\sigma^2 \text{div}(\text{SVT}_\lambda(\mathbf{X})), \end{aligned}$$

$$\begin{aligned} \text{div}(\text{SVT}_\lambda(\mathbf{X})) &= (m-p) \sum_{i=1}^p \left(1 - \frac{\lambda}{\ell_i}\right)_+ + \sum_{i=1}^p \mathbb{1}\{\ell_i > \lambda\} \\ &\quad + 2 \sum_{i=1}^p \sum_{j \neq i}^p \frac{\ell_i(\ell_i - \lambda)_+}{\ell_i^2 - \ell_j^2} \end{aligned}$$

whenever  $\ell_1 > \ell_2 > \dots > \ell_p \geq 0$ .

- An adaptive estimator is given by

$$\mathbf{SVT}_{\widehat{\lambda}}(\mathbf{X}) = \sum_{i=1}^p (\ell_i - \widehat{\lambda})_+ \mathbf{u}_i \mathbf{v}_i^T, \quad (2)$$

$$\widehat{\lambda} = \underset{\lambda \geq 0}{\operatorname{argmin}} \left[ \sum_{i=1}^p \min\{\ell_i^2, \lambda^2\} + 2\sigma^2 \operatorname{div}(\mathbf{SVT}_{\lambda}(\mathbf{X})) \right].$$

- Numerical evaluation of the risk of (2) was carried out by Candés et. al. Obviously

$$\mathbf{R}(\mathbf{SVT}_{\widehat{\lambda}}(\mathbf{X}), \Xi) < \mathbf{R}(\mathbf{X}, \Xi) \quad \text{for } \forall \Xi.$$

- But it is not clear if  $\mathbf{R}(\mathbf{SVT}_{\widehat{\lambda}}(\mathbf{X}), \Xi)$  is close to  $\mathbf{R}(\widehat{\Xi}_{\text{TLS}}(\mathbf{X}), \Xi)$  for  $\forall \Xi$  s.t.  $\operatorname{rank} \Xi \leq r < p$ .

## Sufficient conditions for Stein's identity

- For  $\mathbf{s} \geq \mathbf{0}$  and  $\mathbf{A} \subset \mathbb{R}^p$ , let  $\mathcal{H}^{\mathbf{s}}$  be the Hausdorff measure, i.e.,

$$\mathcal{H}^{\mathbf{s}}(\mathbf{A}) = \lim_{\delta \rightarrow 0} \left[ \inf \left\{ \sum_{j=1}^{\infty} \alpha(\mathbf{s}) \left( \frac{\text{diam} \mathbf{C}_j}{2} \right)^{\mathbf{s}} ; \mathbf{A} \subset \cup_{j=1}^{\infty} \mathbf{C}_j, \text{diam} \mathbf{C}_j \leq \delta \right\} \right]$$

where

$$\alpha(\mathbf{s}) = \frac{\pi^{\mathbf{s}/2}}{\Gamma(1 + \mathbf{s}/2)}.$$

- Let  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_p)^{\top}$  be distributed as  $\mathbf{N}_p(\Xi, \sigma^2 \mathbf{I}_p)$  and  $\widehat{\Xi} = (\widehat{\xi}_1, \dots, \widehat{\xi}_p)^{\top}$  where  $\widehat{\xi}_i : \mathbb{R}^p \rightarrow \mathbb{R}$  ( $i = 1, 2, \dots, p$ ).

## Stein's lemma by Niels R. Hansen(SPL(2018))

- Let  $\mathbf{A} \subset \mathbb{R}^p$  be a closed subset s.t.  $\widehat{\Xi} : \mathbf{A}^c \rightarrow \mathbb{R}^p$  is continuously differentiable. If  $\mathcal{H}^{p-1}(\mathbf{A}) = \mathbf{0}$  and

$$\sum_{i=1}^p \mathbb{E} \left[ \left\| \frac{\partial \widehat{\xi}_i}{\partial \mathbf{x}_i} \right\| \right] < \infty, \quad (2)$$

then

$$\frac{1}{\sigma^2} \sum_{i=1}^p \text{Cov}(\widehat{\xi}_i, \mathbf{x}_i) = \frac{1}{\sigma^2} \mathbb{E}[(\mathbf{X} - \Xi)^\top \widehat{\Xi}] = \mathbb{E}[\text{div}(\widehat{\Xi})] = \mathbb{E} \left[ \sum_{i=1}^p \frac{\partial \widehat{\xi}_i}{\partial \mathbf{x}_i} \right].$$

- The condition (3) is replaced with

$$\text{div}(\widehat{\Xi}) \geq \mathbf{0}, \quad \text{Lebesgue almost everywhere.}$$

- Hansen (2018) showed that Stein's lemma is valid for

$$\hat{\Xi} = \sum_{i=1}^p h_i(\ell_i) \mathbf{u}_i \mathbf{v}_i^T,$$

if

$h_i : (0, \infty) \rightarrow \mathbb{R}$  is continuously differentiable

for  $i = 1, 2, \dots, p-1$  and

$h_p : [0, \infty) \rightarrow \mathbb{R}$  is continuously differentiable with  $h_p(0) = 0$

as well as (3) is satisfied for each  $h_i$  ( $i = 1, 2, \dots, p$ ).

## Case when the covarianc matrix is unknown

- Assume that  $\min\{m, n\} \geq p$  and that

$$\begin{matrix} m \\ n \end{matrix} \begin{matrix} p \\ \\ \end{matrix} \begin{pmatrix} X \\ Y \end{pmatrix} = \begin{pmatrix} \Xi \\ \mathbf{0} \end{pmatrix} + \mathbf{E}; \quad \mathbf{E} = \begin{pmatrix} \overleftrightarrow{\phantom{0}} \\ \overleftrightarrow{\phantom{0}} \\ \vdots \\ \overleftrightarrow{\phantom{0}} \\ \overleftrightarrow{\phantom{0}} \end{pmatrix} \begin{matrix} p \\ \\ \\ \\ \end{matrix}$$

- The  $m \times p$  mean matrix  $\Xi$  is of rank  $r < p$
- The error  $\mathbf{E}$  is an  $(m + n) \times p$  error matrix (unobservable) whose rows are identically distributed as  $\mathbf{N}_p(\mathbf{0}, \Sigma)$ .
- The covariance matrix  $\Sigma$  is a  $p \times p$  positive-definite and unknown.

- We consider the problem of estimating  $\Xi$  under low-rank mean matrix condition, i.e.,

$$\text{rank } \Xi = r < p; \quad r \text{ is unknown.}$$

- A loss fuction and its risk are given by

$$\mathbf{L}_2(\widehat{\Xi}, \Xi) = \text{tr} \{ (\widehat{\Xi} - \Xi) \Sigma^{-1} (\widehat{\Xi} - \Xi)^T \} =: \|\widehat{\Xi} - \Xi\|_{F, \Sigma}^2$$

and

$$\mathbf{R}_2(\widehat{\Xi}, \Xi) = \mathbb{E}[\mathbf{L}_2(\widehat{\Xi}, \Xi)]$$

where  $\widehat{\Xi}$  is an estimator based on  $(\mathbf{X}, \mathbf{S})$ .

- $\mathbf{S} = \mathbf{Y}^T \mathbf{Y} \sim \mathbf{W}_p(\Sigma, \mathbf{n})$ , which is the Wishart distribution with the degree of freedom  $\mathbf{n}$  and the scale matrix  $\Sigma$ .

- To derive a class of estimators, first assume that  $\Sigma$  is known.
- Then we have

$$\mathbf{X}\Sigma^{-1/2} \sim \mathbf{N}_{m \times p}(\tilde{\Xi}, \mathbf{I}_m \otimes \mathbf{I}_p), \quad \tilde{\Xi} = \Xi \Sigma^{-1/2}$$

which leads to an estimator of  $\tilde{\Xi}$  given by

$$\hat{\tilde{\Xi}}_{\text{TLS}} = \underset{\text{rank } \Xi \leq r}{\operatorname{argmin}} \|\mathbf{X}\Sigma^{-1/2} - \Xi\|_F^2 \implies \hat{\Xi} = \hat{\tilde{\Xi}}_{\text{TLS}} \Sigma^{1/2}.$$

- Hence we consider a class of estimators of the form

$$\hat{\Xi}_H = \left( \sum_{i=1}^p h_i(\ell_i) \mathbf{u}_i \mathbf{v}_i^T \right) \mathbf{S}^{1/2}; \quad \mathbf{X}\mathbf{S}^{-1/2} = \mathbf{U}\mathbf{L}\mathbf{V}^T$$

where  $\mathbf{L} = \operatorname{diag}(\ell_1, \dots, \ell_p)$ ,  $\mathbf{H} = \operatorname{diag}(h_1, \dots, h_p)$ ,

$\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_p)$  and  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_p)$  s.t.

$\mathbf{U}^T \mathbf{U} = \mathbf{V}^T \mathbf{V} = \mathbf{I}_p$ .

■ If

$$h_i(\ell_i) = \ell_i - \frac{\mathbf{c}}{\ell_i} \quad (i = 1, 2, \dots, \mathbf{p});$$

$\mathbf{c}$  is a known positive constant,

then it results in the Efron-Morris estimator which is given by

$$\begin{aligned}\widehat{\Xi}_H &= \mathbf{X}\mathbf{S}^{-1/2} \left[ \mathbf{I}_p - \mathbf{c} \{ (\mathbf{X}\mathbf{S}^{-1/2})^\top (\mathbf{X}\mathbf{S}^{-1/2}) \}^{-1} \right] \mathbf{S}^{1/2} \\ &= \mathbf{X} - \mathbf{c}\mathbf{X}\{\mathbf{X}^\top\mathbf{X}\}^{-1}\mathbf{S}.\end{aligned}$$

- On the other hand, Tsukuma and Kubokawa (2015) considered estimators of the form

$$\widehat{\Xi}_T = \mathbf{X} - \mathbf{U}\mathbf{T}\mathbf{U}^\top\mathbf{X},$$

where  $\mathbf{T} = \text{diag}(\mathbf{t}_1(\ell_1^2), \dots, \mathbf{t}_p(\ell_p^2))$ .

- Recall that

$$\mathbf{X}\mathbf{S}^{-1/2} = \mathbf{U}\mathbf{L}\mathbf{V}^T \iff \mathbf{L}^{-1}\mathbf{U}^T\mathbf{X} = \mathbf{V}^T\mathbf{S}^{1/2}.$$

From a simple calculation we get

$$\widehat{\Xi}_H = \mathbf{U}\mathbf{H}\mathbf{V}^T\mathbf{S}^{1/2} = \mathbf{U}\mathbf{H}\mathbf{L}^{-1}\mathbf{U}^T\mathbf{X} = \mathbf{U}\mathbf{L}^{-1}\mathbf{H}\mathbf{X}.$$

- If we set  $\mathbf{I}_p - \mathbf{T} = \mathbf{L}^{-1}\mathbf{H}(t_i(\mathbf{x}) = h_i(\sqrt{\mathbf{x}}))$ , then we have

$$\widehat{\Xi}_H = \widehat{\Xi}_T.$$

- From this we can see that

$$\mathbf{L}_M(\widehat{\Xi}_H, \Xi) = \mathbf{L}_M(\widehat{\Xi}_T, \Xi).$$

- Furthermore, using the result due to Tsukuma and Kubokawa (2015), we have

$$\mathbf{R}_M(\widehat{\Xi}_T, \Xi) = \mathbb{E}[\mathbf{SURE}(\mathbf{T})];$$

$$\mathbf{SURE}(\mathbf{T}) = \sum_{i=1}^p \left[ m + a\ell_i^2 t_i^2 - 2bt_i - 4\ell_i^2 t_i t'_i - 4\ell_i^2 t'_i \right. \\ \left. - 2 \sum_{j \neq i}^p \frac{\ell_i^4 t_i^2 - \ell_j^4 t_j^2}{\ell_i^2 - \ell_j^2} - 4 \sum_{j \neq i}^p \frac{\ell_i^2 t_i - \ell_j^2 t_j}{\ell_i^2 - \ell_j^2} \right];$$

$$t_i = 1 - \frac{h_i(\ell_i)}{\ell_i}; \quad t'_i = -\frac{1}{2\ell_i^2} \left( h'_i(\ell_i) + \frac{h(\ell_i)}{\ell_i} \right),$$

$a, b$  : known positive constants.

- Then we have an adaptive soft-thresholding rule

$$\hat{\Xi}_{\hat{\lambda}} = \mathbf{S} \mathbf{V} \mathbf{T}_{\hat{\lambda}} (\mathbf{X} \mathbf{S}^{-1/2}) \mathbf{S}^{1/2} = \left( \sum_{i=1}^p (\ell_i - \hat{\lambda})_+ \mathbf{u}_i \mathbf{v}_i^T \right) \mathbf{S}^{1/2}$$

where

$$\hat{\lambda} = \underset{\lambda \geq 0}{\operatorname{argmin}} \operatorname{SURE}(\mathbf{S} \mathbf{V} \mathbf{T}_{\lambda}) (\mathbf{X} \mathbf{S}^{-1/2});$$

$$\operatorname{SURE}(\mathbf{S} \mathbf{V} \mathbf{T}_{\lambda}) (\mathbf{X} \mathbf{S}^{-1/2}) = \sum_{i=1}^p \left[ m + a \ell_i^2 t_i^2 - 2 b t_i - 4 \ell_i^2 t_i t'_i - 4 \ell_i^2 t'_i - 2 \sum_{j \neq i}^p \frac{\ell_i^4 t_i^2 - \ell_j^4 t_j^2}{\ell_i^2 - \ell_j^2} - 4 \sum_{j \neq i}^p \frac{\ell_i^2 t_i - \ell_j^2 t_j}{\ell_i^2 - \ell_j^2} \right];$$

$$t_i = 1 - \frac{(\ell_i - \lambda)_+}{\ell_i} \quad (i = 1, \dots, p);$$

$$t'_i = -\frac{1}{\ell_i^2} \left( \mathbb{1}\{\ell_i > \lambda\} + \frac{(\ell_i - \lambda)_+}{\ell_i} \right).$$

## Special case

- $\Sigma = \sigma^2 I_p$  where  $\sigma$  is postive but unknown.
- Let  $\mathbf{s}^2 = \text{tr}(\mathbf{Y}^\top \mathbf{Y})/p$ .
- Then an adaptive soft-thresholding rule for this case is given by  $\hat{\Xi}_{\hat{\lambda}} = \sum_{i=1}^p (\ell_i - \hat{\lambda} \mathbf{s}^2)_+ \mathbf{u}_i \mathbf{v}_i^\top$ ;  $\mathbf{X} = \mathbf{U} \mathbf{L} \mathbf{V}^\top$ , with  $\hat{\lambda} = \underset{\lambda \geq 0}{\text{argmin}} \text{SURE}(\text{SVT}_\lambda)(\mathbf{X})$  and

$$\begin{aligned} \text{SURE}(\text{SVT}_\lambda)(\mathbf{X}) &= \sum_{i=1}^p \left[ m \mathbf{s}^2 + a \ell_i^2 t_i^2 - 4 \ell_i t_i' - 2 \sum_{j \neq i}^p \frac{\ell_i^4 t_i^2 - \ell_j^4 t_j^2}{\ell_i^2 - \ell_j^2} \right. \\ &\quad \left. + \mathbf{s}^2 \left( a \ell_i^2 t_i^2 - 4 \ell_i^2 t_i t_i' - 4 \sum_{j \neq i}^p \frac{\ell_i^2 t_i - \ell_j t_j}{\ell_i^2 - \ell_j^2} \right) \right] \\ t_i &= 1 - \frac{(\ell_i - \lambda \mathbf{s}^2)_+}{\ell_i}; \quad t_i' = -\frac{1}{\ell_i^2} \left( \mathbb{1}\{\ell_i > \lambda \mathbf{s}^2\} + \frac{(\ell_i - \lambda \mathbf{s}^2)_+}{\ell_i} \right). \end{aligned}$$

- **Remark:** It is routine to convert this result to case for complex normal distribution.

## Simulation Study

## References

- 1 Candés, E.J., Sing-Long, C.A., and Trzasko, J.D. (2013): IEEE on Signal Processing **61** 4643–4657.
- 2 Chételat, D. and Wells, M.T. (2012): AOS **40** 3137–3160.
- 3 Efron, B. (2004): JASA **99** 619–642.
- 4 Hansen, N.R. (2018): SPL **135** 76–88.
- 5 Josse, J. and Sardy, S. (2016): Stat. Comput **26** 715–724.
- 6 Mukherjee, A., Chen, K., Wang, N., and Zhu, L. (2015): Biometrika **102** 457–477.
- 7 Tsukuma, H. and Kubokawa, T. (2015): JMVA **139** 312–328.