

Selection of model selection criteria for multivariate ridge regression

Isamu NAGAI

(Received December 6, 2011)

(Revised January 11, 2012)

ABSTRACT. In the present study, we consider the selection of model selection criteria for multivariate ridge regression. There are several model selection criteria for selecting the ridge parameter in multivariate ridge regression, e.g., the C_p criterion and the modified C_p (MC_p) criterion. We propose the generalized C_p (GC_p) criterion, which includes C_p and MC_p criteria as special cases. The GC_p criterion is specified by a non-negative parameter λ , which is referred to as the penalty parameter. We attempt to select an optimal penalty parameter such that the predicted mean squared error (PMSE) of the predictor of the ridge regression after optimizing the ridge parameter is minimized. Through numerical experiments, we verify that the proposed optimization methods exhibit better performance than conventional optimization methods, i.e., optimizing only the ridge parameter by minimizing the C_p or MC_p criterion.

1. Introduction

In the present paper, we deal with a multivariate linear regression model with n observations of a p -dimensional vector of response variables and a k -dimensional vector of regressors (for more detailed information, see, for example, Srivastava, 2002, Chapter 9; Timm, 2002, Chapter 4). Let \mathbf{Y} , \mathbf{X} , and $\mathcal{E} = (\varepsilon_1, \dots, \varepsilon_n)'$ be the $n \times p$ matrix of response variables, the $n \times k$ matrix of non-stochastic centered explanatory variables (i.e., $\mathbf{X}'\mathbf{1}_n = \mathbf{0}_k$) of rank(\mathbf{X}) = k ($< n$), and the $n \times p$ matrix of error variables, respectively, where n is the sample size, $\mathbf{1}_n$ is an n -dimensional vector of ones, and $\mathbf{0}_k$ is a k -dimensional vector of zeros. Suppose that $n - k - p - 2 > 0$ and $\varepsilon_1, \dots, \varepsilon_n \stackrel{\text{i.i.d.}}{\sim} N_p(\mathbf{0}_p, \Sigma)$, where Σ is a $p \times p$ unknown covariance matrix. Then, the matrix form of the multivariate linear regression model is expressed as

$$\mathbf{Y} = \mathbf{1}_n \boldsymbol{\mu}' + \mathbf{X}\boldsymbol{\Xi} + \mathcal{E},$$

2010 *Mathematics Subject Classification.* Primary 62J07; Secondary 62H12.

Key words and phrases. Asymptotic expansion; Generalized C_p criterion; Model selection criterion; Multivariate linear regression model; Ridge regression; Selection of the model selection criterion.

where $\boldsymbol{\mu}$ is a p -dimensional unknown location vector, and $\boldsymbol{\Xi}$ is a $k \times p$ unknown regression coefficient matrix. This model can also be expressed as

$$\mathbf{Y} \sim N_{n \times p}(\mathbf{1}_n \boldsymbol{\mu}' + \mathbf{X} \boldsymbol{\Xi}, \boldsymbol{\Sigma} \otimes \mathbf{I}_n).$$

Note that \mathbf{X} is centered. The maximum likelihood or the least squares (LS) estimators of $\boldsymbol{\mu}$ and $\boldsymbol{\Xi}$ are given by $\hat{\boldsymbol{\mu}} = \mathbf{Y}' \mathbf{1}_n / n$ and $\hat{\boldsymbol{\Xi}} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}$, respectively. Since $\hat{\boldsymbol{\mu}}$ and $\hat{\boldsymbol{\Xi}}$ are simple, and the unbiased estimators of $\boldsymbol{\mu}$ and $\boldsymbol{\Xi}$, the LS estimators are widely used in actual data analysis (see e.g., Dien *et al.*, 2006; Sárbu *et al.*, 2008, Saxén and Sundell, 2006; Skagerberg, Macgregor, and Kiparissides, 1992; Yoshimoto, Yanagihara, and Ninomiya, 2005). However, when the multicollinearity occurs, the problems that an estimator of $\boldsymbol{\Xi}$ becomes unstable happens. In order to avoid this problem, a ridge regression was proposed by Hoerl and Kennard (1970) when $p = 1$. Several authors extended this univariate ridge regression to the multivariate case, e.g., Brown and Zidek (1980), Haitovsky (1987), and Yanagihara and Satoh (2010). The ridge regression estimator of $\boldsymbol{\Xi}$ is given as

$$\hat{\boldsymbol{\Xi}}_\theta = \mathbf{M}_\theta^{-1} \mathbf{X}' \mathbf{Y},$$

where $\mathbf{M}_\theta = \mathbf{X}' \mathbf{X} + \theta \mathbf{I}_k$, and θ is a nonnegative value, which is referred to as a ridge parameter. Since an estimate of $\hat{\boldsymbol{\Xi}}_\theta$ strongly depends on the value of the ridge parameter θ , the optimization of θ is an important problem in the ridge regression.

An optimal θ is commonly determined by minimizing the predicted mean squared error (PMSE) of the predictor of \mathbf{Y} which is defined by $\hat{\mathbf{Y}}_\theta = \mathbf{1}_n \hat{\boldsymbol{\mu}}' + \mathbf{X} \hat{\boldsymbol{\Xi}}_\theta$. However, we cannot directly use the PMSE to optimize θ , because unknown parameters are included in the PMSE. Hence, we adopt an optimization method using a model selection criterion, i.e., an estimator of the PMSE, instead of the unknown PMSE. As an estimator of the PMSE, Yanagihara and Satoh (2010) proposed a C_p criterion. This criterion includes C_p criteria for selecting variables in an univariate linear model, which was proposed by Mallows (1973; 1995), for selecting variables in a multivariate linear model, which was proposed by Sparks, Coutsourides, and Troskie (1983) as a special case. Yanagihara and Satoh (2010) also proposed the modified C_p (MC_p) criterion such that the bias of the C_p criterion for choosing the ridge parameter to the PMSE is completely corrected under a fixed θ . This criterion coincides with the bias-corrected C_p criterion proposed by Fujikoshi and Satoh (1997) when $\theta = 0$. The MC_p criterion has several desirable properties as the estimator of the PMSE as described by, e.g., Fujikoshi, Yanagihara, and Wakaki (2005), and Yanagihara and Satoh (2010).

Unfortunately, optimizing θ by minimizing MC_p , i.e., an unbiased estimator of PMSE, does not always minimize the PMSE of $\hat{\mathbf{Y}}_\theta$. This indicates

that there will be an optimal model selection criterion for selecting θ . Thus, we propose a generalized C_p (GC_p) criterion that includes the C_p and MC_p criteria as special cases (originally, the GC_p criterion was proposed by Atkinson (1980) for selecting variables in the univariate linear model). The GC_p criterion is specified by a non-negative parameter λ , which is referred to as the penalty parameter. From the viewpoint of making the PMSE of the predictor of Y after optimizing θ small, we select the optimal penalty parameter λ , which is basically the selection of the model selection criterion. In the present paper, we optimize λ by the following three methods:

- (Double optimization): We optimize θ and λ simultaneously by minimizing the GC_p and the penalty selection criteria, respectively.
- (Optimization of λ with an approximated value of an optimal θ): We optimize λ by minimizing the penalty selection criterion made from the approximated value of the optimal θ .
- (Asymptotic optimization of λ): We calculate an asymptotic optimal λ from an asymptotic expansion of the PMSE. We then estimate the asymptotic optimal λ .

From the optimization of the model selection criterion, we will perform a reasonable optimization of θ .

The remainder of the present paper is organized as follows: In Section 2, we propose the GC_p criterion, which includes criteria proposed by Yanagihara and Satoh (2010) as special cases. In Section 3, we propose three optimization methods for λ . In Section 4, we compare the optimization methods by conducting numerical studies. Finally, technical details are provided in Appendix.

2. Generalized C_p criterion

In this section, we propose the GC_p criterion for optimizing the ridge parameter, which includes C_p and MC_p criteria proposed by Yanagihara and Satoh (2010). Moreover, we present several mathematical properties of the optimal θ by minimizing the GC_p criterion.

The PMSE of \hat{Y}_θ is defined as

$$\text{PMSE}[\hat{Y}_\theta] = E_Y[E_U[\text{tr}\{(U - \hat{Y}_\theta)'(U - \hat{Y}_\theta)\Sigma^{-1}\}]],$$

where U is a random variable matrix that is independent of Y and has the same distribution as Y .

The C_p criterion proposed by Yanagihara and Satoh (2010) is a rough estimator of the PMSE of \hat{Y}_θ , which is defined by

$$C_p(\theta) = \text{tr}(W_\theta S^{-1}) + 2p \text{tr}(M_\theta^{-1} M_0),$$

where \mathbf{W}_θ is a residual sum of squares matrix defined by $\mathbf{W}_\theta = (\mathbf{Y} - \hat{\mathbf{Y}}_\theta)'(\mathbf{Y} - \hat{\mathbf{Y}}_\theta)$, and \mathbf{S} is an unbiased estimator of $\boldsymbol{\Sigma}$ defined by $\mathbf{S} = \mathbf{W}_\theta / (n - k - 1)$. From the definition of the C_p criterion, the first term of C_p measures the closeness of the ridge regression to the data, and the second term evaluates the penalty for the complexity of the ridge regression. However, the C_p criterion has the bias to the PMSE. The MC_p criterion proposed by Yanagihara and Satoh (2010) is an exact unbiased estimator of the PMSE. By neglecting the terms that are independent of θ , the MC_p criterion is defined as

$$MC_p(\theta) = c_M \text{tr}(\mathbf{W}_\theta \mathbf{S}^{-1}) + 2p \text{tr}(\mathbf{M}_\theta^{-1} \mathbf{M}_0),$$

where $c_M = 1 - (p + 1)/(n - k - 1)$. By comparing the two criteria, we can see that the difference between C_p and MC_p is a coefficient before $\text{tr}(\mathbf{W}_\theta \mathbf{S}^{-1})$.

Thus, we can generalize the model selection criterion for optimizing the ridge parameter as

$$GC_p(\theta, \lambda) = \lambda \text{tr}(\mathbf{W}_\theta \mathbf{S}^{-1}) + 2p \text{tr}(\mathbf{M}_\theta^{-1} \mathbf{M}_0), \quad (1)$$

where λ is a non-negative parameter. Note that $GC_p(\theta, 1) = C_p(\theta)$ and $GC_p(\theta, c_M) = MC_p(\theta)$. In this criterion, the penalty for the complexity of the model, which is in the second term of (1), becomes large when λ becomes small. This means that λ controls the penalty for the complexity of the model in the criterion (1). Hence, we can regard λ as a penalty parameter. In the present paper, we consider the optimization of λ to obtain the optimal θ , which further reduces the PMSE.

When λ is fixed, the optimized ridge parameter $\hat{\theta}(\lambda)$ is obtained by

$$\hat{\theta}(\lambda) = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, \lambda). \quad (2)$$

Since $\hat{\theta}(\lambda)$ is a minimizer of $GC_p(\theta, \lambda)$, the following equation holds:

$$\left. \frac{\partial GC_p(\theta, \lambda)}{\partial \theta} \right|_{\theta = \hat{\theta}(\lambda)} = 0. \quad (3)$$

Note that $\hat{\theta}(\lambda)$ changes with λ .

Here, we obtain the following mathematical properties of $\hat{\theta}(\lambda)$ (The proof is provided in Appendix A.1.):

THEOREM 1. *Let*

$$(\mathbf{z}_1, \dots, \mathbf{z}_k)' = \mathbf{Q}' \mathbf{X}' \mathbf{Y} \mathbf{S}^{-1/2}, \quad (4)$$

$$r_{\lambda, j} = \frac{\lambda \|\mathbf{z}_j\|^2 - p d_j}{p d_j^2}, \quad (j = 1, \dots, k), \quad (5)$$

where \mathbf{z}_i is a p -dimensional vector, \mathbf{Q} is a $k \times k$ orthogonal matrix which diagonalizes $\mathbf{X}'\mathbf{X}$, i.e., $\mathbf{Q}'\mathbf{X}'\mathbf{X}\mathbf{Q} = \mathbf{D} = \text{diag}(d_1, \dots, d_k)$ and d_i ($i = 1, \dots, k$) are eigenvalues of $\mathbf{X}'\mathbf{X}$, and $r_{\lambda,1}^+ \leq \dots \leq r_{\lambda,m}^+$ ($m \leq k$) are positive values of $r_{\lambda,1}, \dots, r_{\lambda,k}$. Then, $\hat{\theta}(\lambda)$ has the following properties:

- (1) $\hat{\theta}(\lambda)$ is a monotonic decreasing function with respect to λ .
- (2) $\hat{\theta}(\lambda)$ is not 0 when $\lambda \in [0, \infty)$.
- (3) $\hat{\theta}(\lambda) > (r_{\lambda,1}^+)^{-1}$ when $r_{\lambda,1}^+$ exists. $\hat{\theta}(\lambda) = \infty$ when $r_{\lambda,1}^+$ does not exist, i.e., $\max_{j=1, \dots, m} r_{\lambda,j} \leq 0$.
- (4) $\hat{\theta}(\infty) = 0$, $\hat{\theta}(0) = \infty$.
- (5) $\hat{\theta}(\lambda) = \infty$ for any $\lambda < \min_{j=1, \dots, k} pd_j^2 / \|\mathbf{z}_j\|^2$.

We suppose that $d_i = O(n)$. We must use an iterative computational algorithm to optimize θ because we cannot obtain $\hat{\theta}(\lambda)$ in closed form. In order to reduce computational tasks, we consider approximating $\hat{\theta}(\lambda)$ using an asymptotic expansion. By applying Taylor expansion to Equation (3), an asymptotic expansion of the GC_p criterion is derived. From this expansion, we obtain the asymptotic expansion of $\hat{\theta}(\lambda)$ as the follows:

THEOREM 2. *The $\hat{\theta}(\lambda)$ can be expanded as*

$$\hat{\theta}(\lambda) = \tilde{\theta}_{(L)}(\lambda) + O_p(n^{-L}),$$

where

$$\begin{aligned} \tilde{\theta}_{(L)}(\lambda) &= \frac{pb_1}{\lambda a_1} + \frac{1}{2\lambda a_1} \sum_{\ell=0}^{L-1} \frac{1}{n^\ell} (-1)^{\ell+1} (\ell+1) \tilde{\theta}_{(L-1)}^\ell(\lambda) \\ &\quad \times \{ \lambda(\ell+2)a_{\ell+1} \tilde{\theta}_{(L-1)}(\lambda) - 2pb_{\ell+1} \}. \end{aligned} \quad (6)$$

Here, $\tilde{\theta}_{(0)}(\lambda) = 0$, $\mathbf{V} = \mathbf{X}'\mathbf{Y}\mathbf{S}^{-1}\mathbf{Y}'\mathbf{X}$, $a_j = n^j \text{tr}(\mathbf{V}\mathbf{M}_0^{-(j+2)})$, $b_j = n^j \text{tr}(\mathbf{M}_0^{-j})$, and $\tilde{\theta}_{(L)}^\ell(\lambda)$ refers to $\{\tilde{\theta}_{(L)}(\lambda)\}^\ell$.

The proof of this theorem is given in Appendix A.2. Note that $\tilde{\theta}_{(L)}(\lambda)$ can be used as an approximated value of $\hat{\theta}(\lambda)$. There is a one-to-one correspondence between $\tilde{\theta}_{(1)}(\lambda) = pb_1/(\lambda a_1)$ and λ , and $\tilde{\theta}_{(1)}(\lambda)$ satisfies the properties 1, 2, and 4 in THEOREM 1.

3. Optimization of penalty in the GC_p criterion

3.1. Double optimization of θ and λ . In the previous section, we considered the model selection criterion for selecting θ , which can be regarded as an estimator of the PMSE $[\hat{\mathbf{Y}}_\theta]$. By minimizing the estimator of the PMSE of $\hat{\mathbf{Y}}_\theta$, we expect to reduce the PMSE of $\hat{\mathbf{Y}}_\theta$. However, since the optimal ridge

parameter will be changed by the data, it is important to reduce not the PMSE of $\hat{\mathbf{Y}}_\theta$ but rather the PMSE of $\hat{\mathbf{Y}}_{\hat{\theta}}$, i.e., the predictor of \mathbf{Y} after optimizing θ . In this section, we consider optimizing λ using $\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}]$, where $\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)} = \mathbf{1}_n \hat{\boldsymbol{\mu}}' + \mathbf{X} \hat{\boldsymbol{\Sigma}}_{\hat{\theta}(\lambda)}$.

Without a loss of generality, we can assume that the covariance matrix of \mathbf{y}_i is \mathbf{I}_p in the $\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}]$. Therefore, from Efron (2004), we obtain the $\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}]$, which can be regarded as a function of λ , as follows:

$$\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}] = E_{\mathbf{Y}}[\text{tr}(\mathbf{W}_{\hat{\theta}(\lambda)} \boldsymbol{\Sigma}^{-1})] + 2E_{\mathbf{Y}} \left[\sum_{i=1}^n \sum_{j=1}^p \frac{\partial(\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)})_{ij}}{\partial(\mathbf{Y})_{ij}} \right],$$

where $(\mathbf{A})_{ij}$ is the (i, j) th element of \mathbf{A} . Since $\hat{\theta}(\lambda)$ depends on $(\mathbf{Y})_{ij}$, we can see that

$$\frac{\partial(\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)})_{ij}}{\partial(\mathbf{Y})_{ij}} = \frac{\partial(\hat{\mathbf{Y}}_\theta)_{ij}}{\partial(\mathbf{Y})_{ij}} \Big|_{\theta=\hat{\theta}(\lambda)} + \frac{\partial(\hat{\mathbf{Y}}_\theta)_{ij}}{\partial\theta} \Big|_{\theta=\hat{\theta}(\lambda)} \frac{\partial\hat{\theta}(\lambda)}{\partial(\mathbf{Y})_{ij}}. \quad (7)$$

The first term of the above equation is calculated as

$$\begin{aligned} \frac{\partial(\hat{\mathbf{Y}}_\theta)_{ij}}{\partial(\mathbf{Y})_{ij}} \Big|_{\theta=\hat{\theta}(\lambda)} &= \frac{\partial(\mathbf{1}_n \hat{\boldsymbol{\mu}}')_{ij}}{\partial(\mathbf{Y})_{ij}} + \frac{\partial(\mathbf{X} \mathbf{M}_\theta^{-1} \mathbf{X}' \mathbf{Y})_{ij}}{\partial(\mathbf{Y})_{ij}} \Big|_{\theta=\hat{\theta}(\lambda)} \\ &= \frac{\partial(\mathbf{1}_n \hat{\boldsymbol{\mu}}')_{ij}}{\partial(\mathbf{Y})_{ij}} + (\mathbf{X} \mathbf{M}_{\hat{\theta}(\lambda)}^{-1} \mathbf{X}')_{ii}. \end{aligned}$$

Note that $\sum_{i=1}^n \sum_{j=1}^p \partial(\mathbf{1}_n \hat{\boldsymbol{\mu}}')_{ij} / \partial(\mathbf{Y})_{ij} = p$ and $\sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\hat{\theta}(\lambda)}^{-1} \mathbf{X}')_{ii} = p \text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-1} \mathbf{M}_0)$.

Next, we consider obtaining the second term of (7). Note that

$$\frac{\partial(\hat{\mathbf{Y}}_\theta)_{ij}}{\partial\theta} = \frac{\partial(\mathbf{X} \mathbf{M}_\theta^{-1} \mathbf{X}' \mathbf{Y})_{ij}}{\partial\theta} = -(\mathbf{X} \mathbf{M}_\theta^{-2} \mathbf{X}' \mathbf{Y})_{ij}.$$

Hence, we derive

$$\begin{aligned} \text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}] &= E_{\mathbf{Y}}[\text{tr}(\mathbf{W}_{\hat{\theta}(\lambda)} \boldsymbol{\Sigma}^{-1})] + 2p\{E_{\mathbf{Y}}[\text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-1} \mathbf{M}_0)] + 1\} \\ &\quad - 2E_{\mathbf{Y}} \left[\sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\hat{\theta}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \frac{\partial\hat{\theta}(\lambda)}{\partial(\mathbf{Y})_{ij}} \right]. \quad (8) \end{aligned}$$

Based on this result, we need only to obtain $\partial\hat{\theta}(\lambda)/\partial(\mathbf{Y})_{ij}$ in order to calculate the $\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}]$. This derivative leads to the following theorem (The proof is given in Appendix A.3.):

THEOREM 3. The PMSE of $\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}$ is expressed as

$$\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}] = E_Y[\text{tr}(\mathbf{W}_{\hat{\theta}(\lambda)}\boldsymbol{\Sigma}^{-1}) + 2p\{\text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-1}\mathbf{M}_0) + 1\} + 4B(\hat{\theta}(\lambda))], \quad (9)$$

where

$$B(\theta) = \frac{\lambda\theta \text{tr}(\mathbf{V}\mathbf{M}_{\theta}^{-5}\mathbf{M}_0)}{\lambda \text{tr}(\mathbf{V}\mathbf{M}_{\theta}^{-3}) - 3\lambda\theta \text{tr}(\mathbf{V}\mathbf{M}_{\theta}^{-4}) + 2p \text{tr}(\mathbf{M}_{\theta}^{-3}\mathbf{M}_0)}.$$

By neglecting the terms that are independent of λ , we define the penalty selection criteria for optimizing λ as follows:

DEFINITION 1. The penalty selection criteria to optimize λ are defined as

$$C_p^\#(\lambda) = \text{tr}(\mathbf{W}_{\hat{\theta}(\lambda)}\mathbf{S}^{-1}) + 2p \text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-1}\mathbf{M}_0) + 4B(\hat{\theta}(\lambda)),$$

$$MC_p^\#(\lambda) = c_M \text{tr}(\mathbf{W}_{\hat{\theta}(\lambda)}\mathbf{S}^{-1}) + 2p \text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-1}\mathbf{M}_0) + 4B(\hat{\theta}(\lambda)),$$

where $\hat{\theta}(\lambda)$ is given by (2) and $c_M = 1 - (p+1)/(n-k-1)$.

Here, $C_p^\#(\lambda)$ is obtained by substituting \mathbf{S}^{-1} into $\boldsymbol{\Sigma}^{-1}$ when we neglect the terms that are independent of λ in (9). However, there exists a bias because \mathbf{S}^{-1} is not an unbiased estimator of $\boldsymbol{\Sigma}^{-1}$ (see e.g., Siotani, Hayakawa, and Fujikoshi, 1985). Based on the results reported by Yanagihara and Satoh (2010), we may correct the bias of $E_Y[\text{tr}(\mathbf{W}_{\hat{\theta}(\lambda)}\mathbf{S}^{-1})]$ to $E_Y[\text{tr}(\mathbf{W}_{\hat{\theta}(\lambda)}\boldsymbol{\Sigma}^{-1})]$. Finally, we define $MC_p^\#(\lambda)$ by neglecting the terms that are independent of θ and λ . Using these criteria, λ and θ are optimized as follows:

$$\hat{\lambda}_C^\# = \arg \min_{\lambda \in [0, \infty]} C_p^\#(\lambda) \quad \text{and} \quad \hat{\theta}(\hat{\lambda}_C^\#) = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, \hat{\lambda}_C^\#),$$

$$\hat{\lambda}_M^\# = \arg \min_{\lambda \in [0, \infty]} MC_p^\#(\lambda) \quad \text{and} \quad \hat{\theta}(\hat{\lambda}_M^\#) = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, \hat{\lambda}_M^\#).$$

These optimization methods are similar to those reported by Ye (1998) and Shen and Ye (2002).

3.2. Optimization of λ with the approximated $\hat{\theta}(\lambda)$. In the previous subsection, we proposed the penalty selection criteria for selecting λ . These criteria are made from the optimal θ obtained by minimizing the GC_p criterion. This indicates that we need to repeat the optimization of θ until obtaining the optimal λ . Hence, many computational tasks are required for such an optimization. In this subsection, we try to reduce computational tasks by using the approximated $\hat{\theta}(\lambda)$, which is given by (6). Thus, we propose the

penalty selection criterion when the approximated $\hat{\theta}(\lambda)$ is used. As such, we calculate $\partial\tilde{\theta}_{(L)}(\lambda)/\partial(\mathbf{Y})_{ij}$. The following lemma is useful for obtaining such a derivative (The proof is provided in Appendix A.4):

LEMMA 1. *For any ℓ , the first derivative of a_ℓ with respect to $(\mathbf{Y})_{ij}$ is calculated as*

$$\frac{\partial a_\ell}{\partial(\mathbf{Y})_{ij}} = 2n^\ell \left(\mathbf{S}^{-1} \mathbf{Y}' \mathbf{X} \mathbf{M}_0^{-\ell-2} \mathbf{X}' \left(\mathbf{I}_n - \frac{1}{n-k-1} \mathbf{Y} \mathbf{S}^{-1} \mathbf{Y}' \mathbf{H} \right) \right)_{ji},$$

where $\mathbf{H} = \mathbf{I}_n - \mathbf{1}_n \mathbf{1}'_n / n - \mathbf{X} \mathbf{M}_0^{-1} \mathbf{X}'$.

By using this lemma and (6), we obtain the following theorem:

THEOREM 4. *The PMSE of $\hat{\mathbf{Y}}_{\tilde{\theta}_{(L)}(\lambda)}$ is expressed as*

$$\begin{aligned} \text{PMSE}[\hat{\mathbf{Y}}_{\tilde{\theta}_{(L)}(\lambda)}] &= E_{\mathbf{Y}}[\text{tr}(\mathbf{W}_{\tilde{\theta}_{(L)}(\lambda)} \boldsymbol{\Sigma}^{-1}) + 2p\{\text{tr}(\mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-1} \mathbf{M}_0) + 1\}] \\ &\quad + 2E_{\mathbf{Y}}[\mathbf{B}'(\tilde{\theta}_{(L)}(\lambda))], \end{aligned}$$

where

$$\begin{aligned} \mathbf{B}'(\tilde{\theta}_{(L)}(\lambda)) &= \frac{2n}{a_1} \tilde{\theta}_{(1)}(\lambda) \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{V}) \\ &\quad + \frac{1}{\lambda a_1^2} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{V}) \\ &\quad \times \sum_{\ell=0}^{L-1} \frac{1}{n^{\ell-1}} (-1)^{\ell+1} (\ell+1) \tilde{\theta}_{(L-1)}^{\ell}(\lambda) \{\lambda(\ell+2) a_{\ell+1} \tilde{\theta}_{(L-1)}(\lambda) - 2pb_{\ell+1}\} \\ &\quad - \frac{1}{2\lambda a_1} \sum_{\ell=0}^{L-1} \frac{1}{n^\ell} (-1)^{\ell+1} (\ell+1) \tilde{\theta}_{(L-1)}^{\ell-1}(\lambda) \\ &\quad \times \{\lambda(\ell+1)(\ell+2) a_{\ell+1} \tilde{\theta}_{(L-1)}(\lambda) - 2p\ell b_{\ell+1}\} \\ &\quad \times \sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \frac{\partial \tilde{\theta}_{(L-1)}(\lambda)}{\partial(\mathbf{Y})_{ij}} \\ &\quad - \frac{n}{a_1} \sum_{\ell=0}^{L-1} (-1)^{\ell+1} (\ell+1)(\ell+2) \tilde{\theta}_{(L-1)}^{\ell+1}(\lambda) \text{tr}(\mathbf{M}_0^{-(\ell+2)} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{V}). \end{aligned}$$

The proof of this theorem is presented in Appendix A.5. When $\tilde{\theta}_{(1)}(\lambda)$ is used, we obtain

$$\begin{aligned} \text{PMSE}[\hat{\mathbf{Y}}_{\tilde{\theta}_{(1)}(\lambda)}] &= E_{\mathbf{Y}}[\text{tr}(\mathbf{W}_{\tilde{\theta}_{(1)}(\lambda)} \boldsymbol{\Sigma}^{-1}) + 2p\{\text{tr}(\mathbf{M}_{\tilde{\theta}_{(1)}(\lambda)}^{-1} \mathbf{M}_0) + 1\}] \\ &\quad + E_{\mathbf{Y}} \left[\frac{4n\tilde{\theta}_{(1)}(\lambda)}{a_1} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(1)}(\lambda)}^{-2} \mathbf{V}) \right]. \end{aligned}$$

Thus, by neglecting the terms that are independent of λ , the penalty selection criteria with $\tilde{\theta}_{(1)}(\lambda)$ are defined as follows:

DEFINITION 2. The penalty selection criteria to optimize λ when $\tilde{\theta}_{(1)}(\lambda)$ is used are defined as follows:

$$\begin{aligned} C_p^{(1)}(\lambda) &= \text{tr}(\mathbf{W}_{\tilde{\theta}_{(1)}(\lambda)} \mathbf{S}^{-1}) + 2p \text{tr}(\mathbf{M}_{\tilde{\theta}_{(1)}(\lambda)}^{-1} \mathbf{M}_0) \\ &\quad + \frac{4n}{a_1} \tilde{\theta}_{(1)}(\lambda) \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(1)}(\lambda)}^{-2} \mathbf{V}), \\ MC_p^{(1)}(\lambda) &= c_M \text{tr}(\mathbf{W}_{\tilde{\theta}_{(1)}(\lambda)} \mathbf{S}^{-1}) + 2p \text{tr}(\mathbf{M}_{\tilde{\theta}_{(1)}(\lambda)}^{-1} \mathbf{M}_0) \\ &\quad + \frac{4n}{a_1} \tilde{\theta}_{(1)}(\lambda) \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(1)}(\lambda)}^{-2} \mathbf{V}), \end{aligned}$$

where $\tilde{\theta}_{(1)}(\lambda) = pb_1/(\lambda a_1)$ and $c_M = 1 - (p+1)/(n-k-1)$.

Similar to $MC_p^\#(\lambda)$, $MC_p^{(1)}(\lambda)$ can be regarded as a simple bias-corrected $C_p^{(1)}(\lambda)$. At least, when $\tilde{\theta}_{(1)}(\lambda) = 0$, $MC_p^{(1)}(\lambda)$ completely corrects the bias of $C_p^{(1)}(\lambda)$. If we use a $\tilde{\theta}_{(L)}(\lambda)$ other than $\tilde{\theta}_{(1)}(\lambda)$, the penalty selection criteria becomes more complicated as the number of L increases. As an example, we describe the penalty selection criteria for $\tilde{\theta}_{(2)}(\lambda)$ in Appendix A.6. From the viewpoint of an application, $C_p^{(1)}(\lambda)$ and $MC_p^{(1)}(\lambda)$ are useful because these are the simplest criteria among all L . When we use $C_p^{(1)}(\lambda)$ and $MC_p^{(1)}(\lambda)$, the optimal θ and λ are given as follows:

$$\begin{aligned} \hat{\lambda}_C^{(1)} &= \arg \min_{\lambda \in [0, \infty]} C_p^{(1)}(\lambda) \quad \text{and} \quad \hat{\theta}(\hat{\lambda}_C^{(1)}) = \tilde{\theta}_{(1)}(\hat{\lambda}_C^{(1)}), \\ \hat{\lambda}_M^{(1)} &= \arg \min_{\lambda \in [0, \infty]} MC_p^{(1)}(\lambda) \quad \text{and} \quad \hat{\theta}(\hat{\lambda}_M^{(1)}) = \tilde{\theta}_{(1)}(\hat{\lambda}_M^{(1)}). \end{aligned}$$

3.3. Asymptotic optimization for λ . In the previous subsections, we proposed the penalty selection criteria with $\tilde{\theta}_{(1)}(\lambda)$. When such criteria are used to optimize λ , we must perform an iterative procedure. In this subsection, we consider the non-iterative optimization of λ . This requires the calculation of an asymptotic optimal λ , which minimizes an asymptotic expansion of the $\text{PMSE}[\hat{\mathbf{Y}}_{\tilde{\theta}(\lambda)}]$ among $\lambda \in [0, \infty]$. The following theorem gives such an asymptotic optimal value of λ (The proof is provided in Appendix A.7.):

THEOREM 5. *An asymptotic optimal λ^* minimizes the $\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}]$ asymptotically is given by*

$$\frac{1}{\lambda^*} = \frac{E_Y[a_1^{-1}]}{E_Y[a_1^*/a_1^2]} - \frac{2E_Y[a_2/a_1^2]}{pb_1E_Y[a_1^*/a_1^2]},$$

where $\mathbf{V}^* = \mathbf{X}'\mathbf{Y}\Sigma^{-1}\mathbf{Y}'\mathbf{X}$ and $a_j^* = n^j \text{tr}(\mathbf{V}^*\mathbf{M}_0^{-(j+2)})$.

By replacing a_1^* with a_1 , we estimate λ^* as follows:

$$\hat{\lambda}_0 = \left\{ 1 - \frac{2 \text{tr}(\mathbf{M}_0^{-4}\mathbf{V})}{p \text{tr}(\mathbf{M}_0^{-3}\mathbf{V}) \text{tr}(\mathbf{M}_0^{-1})} \right\}^{-1}. \quad (10)$$

Note that $E_Y[a_1] = c_M^{-1}E_Y[a_1^*]$ holds. Hence, we can estimate $E_Y[a_1^*]$ as $c_M a_1$. This implies the new estimator of λ^* is given by $\hat{\lambda}_M = c_M \hat{\lambda}_0$. When we use $\hat{\lambda}_0$ and $\hat{\lambda}_M$, optimal θ is given by

$$\hat{\theta}(\hat{\lambda}_0) = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, \hat{\lambda}_0) \quad \text{and} \quad \hat{\lambda}_0 \text{ is in (10),}$$

$$\hat{\theta}(\hat{\lambda}_M) = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, \hat{\lambda}_M) \quad \text{and} \quad \hat{\lambda}_M = c_M \hat{\lambda}_0.$$

4. Numerical study

In this section, we conduct numerical experiments to compare the PMSEs of predictors of \mathbf{Y} consisting of the ridge regression estimators with the optimized ridge and penalty parameters. Let \mathbf{R}_q and $\mathbf{A}_q(\rho)$ be $q \times q$ matrices defined by $\mathbf{R}_q = \text{diag}(1, \dots, q)$ and $(\mathbf{A}_q(\rho))_{ij} = \rho^{|i-j|}$. The explanatory matrix \mathbf{X} was generated from $\mathbf{X} = \mathbf{W}\Psi^{1/2}$, where $\Psi = \mathbf{R}_k^{1/2} \mathbf{A}_k(\rho_x) \mathbf{R}_k^{1/2}$, and \mathbf{W} is an $n \times k$ matrix, whose elements were generated independently from the uniform distribution on $(-1, 1)$. The $k \times p$ unknown regression coefficient matrix Ξ was defined by $\Xi = \delta \mathbf{F} \Xi_i$, where δ is a constant, \mathbf{F} is defined as $\mathbf{F} = \text{diag}(\mathbf{1}_\kappa, \mathbf{0}_{k-\kappa})$, which is a $k \times k$ matrix, and Ξ_i is defined by Ξ_0 when $k = 5$, Ξ_1 when $k = 10$, and Ξ_2 when $k = 15$, respectively. Here, Ξ_0 is defined by the first five rows of Ξ_1 , and Ξ_1 and Ξ_2 are given by

$$\Xi_1 = \begin{pmatrix} 0.8501 & 0.6571 & 0.2159 \\ -0.2753 & -0.2432 & -0.1187 \\ -0.3193 & -0.2926 & -0.1671 \\ 0.2754 & 0.2608 & 0.1766 \\ 0.2693 & 0.2164 & 0.2066 \\ -0.0676 & -0.0663 & -0.0561 \\ 0.2239 & 0.2197 & 0.1880 \\ -0.0352 & -0.0346 & -0.0305 \\ 0.3240 & 0.3199 & 0.2868 \\ -0.3747 & -0.3727 & -0.3554 \end{pmatrix},$$

$$\mathbf{\Xi}_2 = \begin{pmatrix} 1.3794 & 0.0645 & 0.0330 \\ -0.0766 & -0.0241 & -0.0143 \\ -0.2618 & -0.1396 & -0.0951 \\ -0.4619 & -0.2589 & -0.1798 \\ 0.2381 & 0.1488 & 0.1082 \\ 0.2140 & 0.1463 & 0.1112 \\ 0.3002 & 0.2364 & 0.1950 \\ 0.1155 & 0.0953 & 0.0812 \\ -0.2774 & -0.2395 & -0.2091 \\ 0.3392 & 0.3072 & 0.2807 \\ 0.0016 & 0.0107 & 0.0100 \\ 0.0438 & 0.0408 & 0.0381 \\ -0.3187 & -0.3039 & -0.2904 \\ 0.0529 & 0.0510 & 0.0493 \\ 0.2505 & 0.2451 & 0.2399 \end{pmatrix}.$$

Here, δ controls the scale of the regression coefficient matrix, and F controls the number of non-zero regression coefficients via κ (the dimension of the true model). The values of the elements of $\mathbf{\Xi}_1$ and $\mathbf{\Xi}_2$, which is an essential regression coefficient matrix, are the same as in Lawless (1981). Simulated data \mathbf{Y} were generated by $N_{n \times 3}(\mathbf{X}\mathbf{\Xi}, \mathbf{\Sigma} \otimes \mathbf{I}_n)$ repeatedly under several selections of n , k , κ , δ , ρ_y , and ρ_x , where $\mathbf{\Sigma} = \mathbf{R}_3^{1/2} \mathbf{A}_3(\rho_y) \mathbf{R}_3^{1/2}$, and the number of repetition was 1000. At each repetition, we evaluated $r(\mathbf{X}\mathbf{\Xi}, \hat{\mathbf{Y}}_{\hat{\theta}}) = \text{tr}\{(\mathbf{X}\mathbf{\Xi} - \hat{\mathbf{Y}}_{\hat{\theta}})'(\mathbf{X}\mathbf{\Xi} - \hat{\mathbf{Y}}_{\hat{\theta}})\mathbf{\Sigma}^{-1}\}$, where $\hat{\mathbf{Y}}_{\hat{\theta}} = \mathbf{1}_n \hat{\boldsymbol{\mu}}' + \mathbf{X}\hat{\boldsymbol{\Xi}}_{\hat{\theta}}$, which is the predicted value of \mathbf{Y} obtained from each method. The average of $np + r(\mathbf{X}\mathbf{\Xi}, \hat{\mathbf{Y}}_{\hat{\theta}})$ across 1000 repetitions was regarded as the PMSE of $\hat{\mathbf{Y}}_{\hat{\theta}}$. In the simulation, a standardized \mathbf{X} was used to estimate the regression coefficients.

Recall that $GC_p(\theta, \lambda)$ is defined in (1). Here, λ and θ are optimized by the following methods:

Method 1: $\hat{\theta} = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, \hat{\lambda})$ and $\hat{\lambda} = \hat{\lambda}_0$, where $\hat{\lambda}_0$ is defined in (10).

Method 2: $\hat{\theta} = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, \hat{\lambda})$ and $\hat{\lambda} = \hat{\lambda}_M = c_M \hat{\lambda}_0$, where $c_M = 1 - (p + 1)/(n - k - 1)$.

Method 3: $\hat{\theta} = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, \hat{\lambda})$ and $\hat{\lambda} = \arg \min_{\lambda \in [0, \infty]} C_p^\#(\lambda)$, where $C_p^\#(\lambda)$ is given by DEFINITION 1.

Method 4: $\hat{\theta} = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, \hat{\lambda})$ and $\hat{\lambda} = \arg \min_{\lambda \in [0, \infty]} MC_p^\#(\lambda)$, where $MC_p^\#(\lambda)$ is given by DEFINITION 1.

Method 5: $\hat{\theta} = \tilde{\theta}_{(1)}(\hat{\lambda})$ and $\hat{\lambda} = \arg \min_{\lambda \in [0, \infty]} C_p^{(1)}(\lambda)$, where $\tilde{\theta}_{(1)}(\lambda)$ and $C_p^{(1)}(\lambda)$ are defined in (6) and by DEFINITION 2.

Method 6: $\hat{\theta} = \tilde{\theta}_{(1)}(\hat{\lambda})$ and $\hat{\lambda} = \arg \min_{\lambda \in [0, \infty]} MC_p^{(1)}(\lambda)$, where $MC_p^{(1)}(\lambda)$ is given by DEFINITION 2.

For the purpose of comparison with the proposed methods, we prepare conventional optimization methods, which are obtained using the following methods:

Method 7: $\hat{\theta}_C = \arg \min_{\theta \in [0, \infty]} C_p(\theta) = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, 1)$.

Method 8: $\hat{\theta}_M = \arg \min_{\theta \in [0, \infty]} MC_p(\theta) = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, c_M)$, where $c_M = 1 - (p + 1)/(n - k - 1)$.

In Methods 3 through 8, the `fminsearch` function in Matlab is used to find the minimizer of the penalty selection criterion or model selection criterion. In the `fminsearch` function, the Nelder-Mead simplex method (see e.g., Lagarias *et al.*, 1998) is used to search the value that minimizes the function. When Methods 1 through 4 are used, an optimal θ is searched using the `fminsearch` function. We can see that the computational speeds of Methods 1 and 2 are the same as those of Methods 5 and 6. Furthermore, the computational speeds of Methods 5 and 6 are almost the same as those of Methods 7 and 8 because these six methods optimize one parameter. It is easy to predict that the computational speeds of Methods 3 and 4 are slower than the other methods because Methods 3 and 4 optimize two parameters simultaneously.

In this paper, we proposed Methods 1 through 6 as referred to above, and these methods can be regarded as the estimation methods for the optimal λ . To obtain the optimal λ , called λ^{**} , which minimizes the PMSE, we divided a range $[0, 2]$ into 100 parts and used each point. Then we compute $r(\mathbf{X}\boldsymbol{\Xi}, \hat{\mathbf{Y}}_{\hat{\theta}})$ for each point in each repetition. After 1000 repetitions, we compute the averages of these values for each point which are regarded as the main term of the PMSE of $\hat{\mathbf{Y}}_{\hat{\theta}}$. By comparing the average values, the λ^{**} is obtained. For comparing $\hat{\lambda}$ which is estimated λ by using each method in above Methods 1 through 6, we show the Figure 1 which is the boxplot of $\hat{\lambda}$ of 1000 repetitions in several situations. The horizontal line means the λ^{**} . Tables 1 through 4 show the averages of $(\hat{\lambda} - \lambda^{**})^2$ across 1000 repetitions, which is referred as the mean squared error (MSE) of $\hat{\lambda}$, for each method.

In THEOREM 2, we derived the expansion for $\hat{\theta}(\lambda)$ and we suggested to use the first term of the expansion which is referred as $\tilde{\theta}_{(1)}(\lambda)$. To compare the $\hat{\theta}(\lambda)$ and $\tilde{\theta}_{(1)}(\lambda)$, we show the scatter plots in several situations when we fix λ as 1 or 2 in Figure 2. In each scatter plot, the 45-degree line means the line of $\hat{\theta}(\lambda) = \tilde{\theta}_{(1)}(\lambda)$. When the scatter plot closes up this line, $\tilde{\theta}_{(1)}(\lambda)$ closes to $\hat{\theta}(\lambda)$.

Tables 7 through 12 show the simulation results obtained for $\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}}] / \{p(n + k + 1)\} \times 100$ for the cases in which $(k, n) = (5, 30), (5, 50), (10, 30), (10, 50), (15, 30)$, and $(15, 50)$, respectively, where $p(n + k + 1)$ is the PMSE

of the predictor of Y derived using the LS estimators. We note $p = 3$ in this numerical studies. In the tables, bold face indicates the minimized PMSE, and italic face indicates the second-smallest PMSE.

Figure 1, we can see the dispersion of $\hat{\lambda}$ and the differences between $\hat{\lambda}$ and λ^{**} . Methods 2, 4 and 6 give always smaller values than Methods 1, 3 and 5. The dispersions of $\hat{\lambda}$ obtained from Methods 2, 4 and 6 are smaller than those obtained from Methods 1, 3 and 5. This facts mean that the formal bias correction of each method makes $\hat{\lambda}$ and dispersion of $\hat{\lambda}$ smaller. We note that the dispersions of $\hat{\lambda}$ obtained from Methods 1 and 2 are smaller than other methods. When ρ_y and ρ_x are small, our optimization method gives nearly closer value to λ^{**} . From Tables 1 through 6, we can see the numerical evaluation for each method. When k and δ are zeros, Method 6 is the best and Method 5 is the second best. Methods 2 and 4 are the best and the second best when κ and δ is small. When κ is equal to $k = 10$ and δ is large, Method 6 is the best method. On the other hand, when $k = \kappa = 15$ and δ is large, Method 6 or 2 is the best in ρ_y is small or large. Consequently, Method 6 and 5 was, on average, the best and the second best method except $k = 15$. When $k = 15$, Method 5 was the best. Hence we recommend using Method 6 to optimize the penalty parameter λ .

Figure 2, we can see the difference between dispersions of $\tilde{\theta}_{(1)}(\lambda)$ and $\hat{\theta}(\lambda)$ in each situation. We note that λ becomes large, the difference between $\hat{\theta}(\lambda)$ and $\tilde{\theta}_{(1)}(\lambda)$ becomes small. Also when ρ_y or n becomes large, the difference between $\hat{\theta}(\lambda)$ and $\tilde{\theta}_{(1)}(\lambda)$ becomes small. On the other hand, the difference between $\hat{\theta}(\lambda)$ and $\tilde{\theta}_{(1)}(\lambda)$ becomes large when ρ_x is large. In almost case, $\tilde{\theta}_{(1)}(\lambda)$ is smaller than $\hat{\theta}(\lambda)$. This fact is corresponding the result $\hat{\theta}(\lambda) = \tilde{\theta}_{(1)}(\lambda) + O_p(n^{-1})$ in THEOREM 2. When ρ_y or δ becomes large, the dispersion of $\hat{\theta}(\lambda)$ becomes small. The each value of $\hat{\theta}(\lambda)$ and $\tilde{\theta}_{(1)}(\lambda)$ become small when ρ_y , ρ_x , δ or λ becomes large.

Based on the simulations, we can see that all of the methods improved the PMSEs of the LS estimators in almost all cases. All of the methods greatly improved the PMSE when n becomes small or k becomes large. Moreover, the improvement in the PMSE of the proposed method increases as ρ_y decreases. The improvement in the PMSE when $\kappa \neq 0$ and $\delta \neq 0$ of the proposed method increases as ρ_x increases. A comparison of several methods reveals that Methods 2 and 4 were better than Methods 1 and 3, respectively, in almost all cases when ρ_x is large. When k and ρ_x become large, Methods 5 and 6 provide a greater improvement in the PMSE than Methods 3 and 4. When k becomes small and n becomes large, Methods 4 and 6 improve the PMSE more than Methods 2 and 4 in most cases. Occasionally, Method 7 improves the PMSE more than Method 8, especially when κ and δ become large. Consequently, Method 6 was, on average, the best method. In

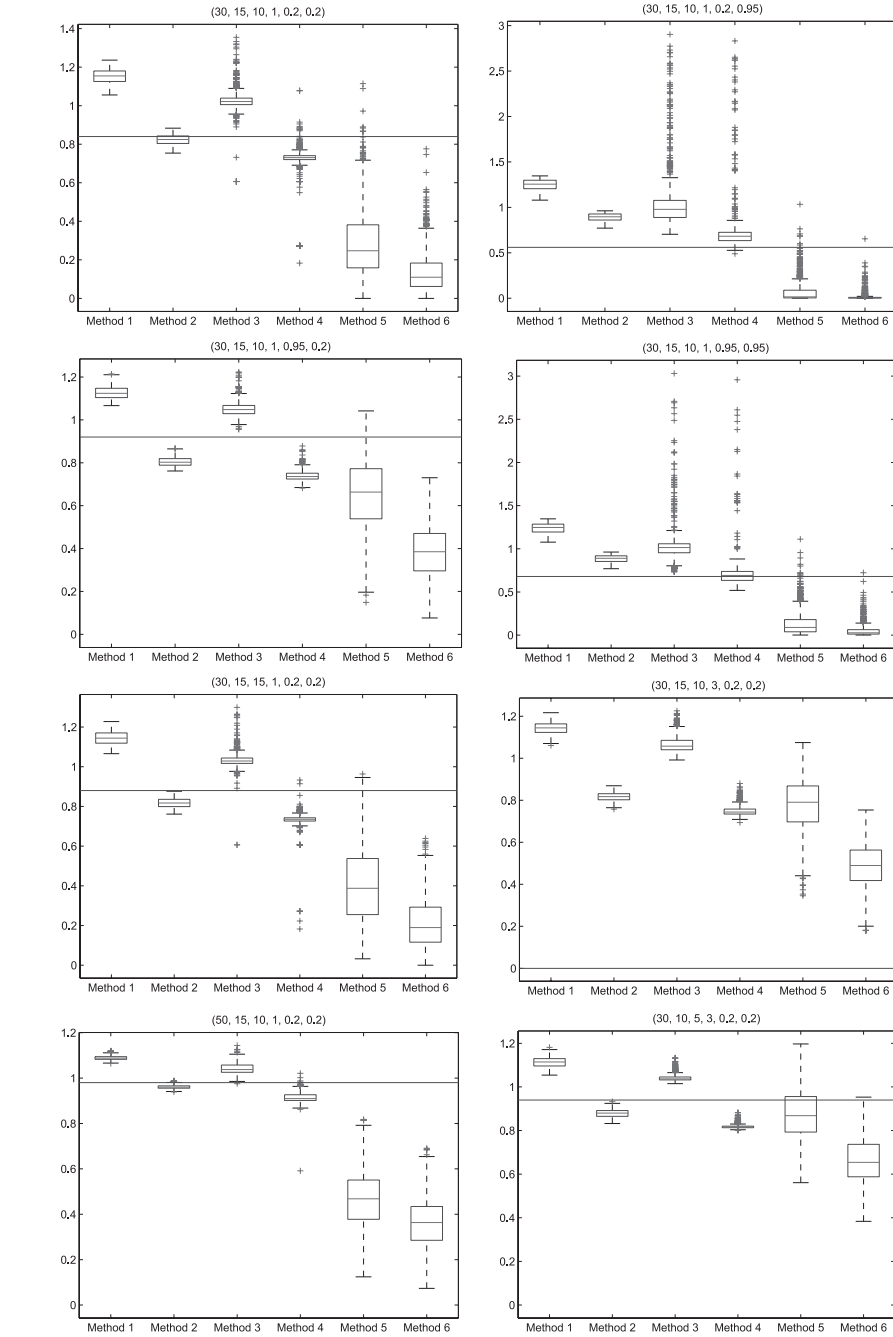


Fig. 1. The box plots for $\hat{\lambda}$ in the case of $(n, k, \kappa, \delta, \rho_y, \rho_x)$

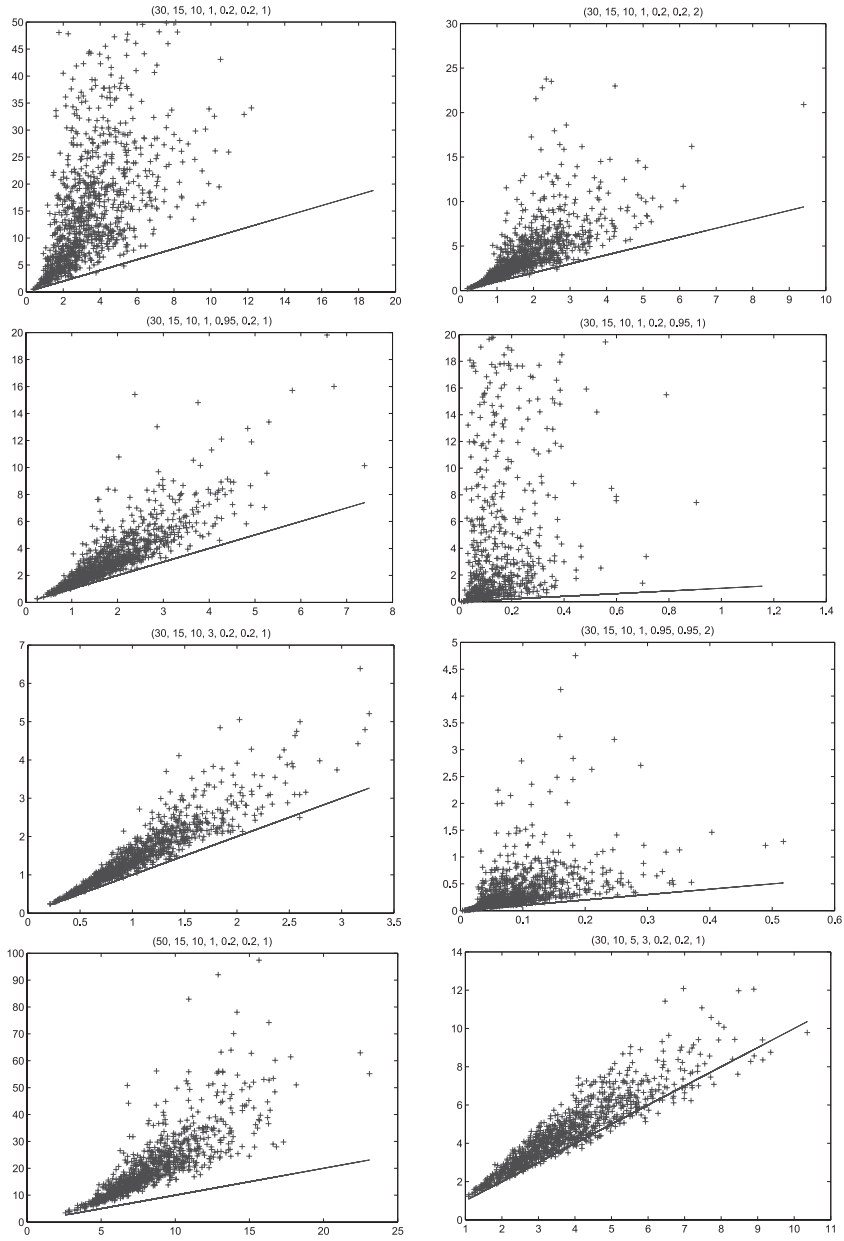


Fig. 2. The scatter plots for $\tilde{\theta}_{(1)}(\lambda)$ and $\hat{\theta}(\lambda)$ in the case of $(n, k, \kappa, \delta, \rho_y, \rho_x, \lambda)$

Table 1. MSE of $\hat{\lambda}$ in which $(k, n) = (5, 30)$

κ	δ	ρ_y	ρ_x	Method					
				1	2	3	4	5	6
0	0	0.2	0.2	1.470	1.014	1.475	1.145	0.039	0.014
			0.95	1.644	1.135	2.599	2.151	0.036	0.013
	0.95	0.2	1.327	0.896	1.239	0.907	0.033	0.013	
		0.95	1.593	1.092	2.482	2.050	0.035	0.012	
5	1	0.2	0.2	0.030	0.002	0.009	0.015	0.181	0.331
			0.95	0.305	0.114	0.424	0.266	0.378	0.445
		0.95	0.2	0.006	0.016	0.001	0.032	0.055	0.167
			0.95	0.340	0.137	0.327	0.169	0.270	0.338
	3	0.2	0.2	0.008	0.012	0.004	0.017	0.001	0.030
			0.95	1.577	1.095	1.261	0.815	0.303	0.154
		0.95	0.2	1.322	0.918	1.293	0.894	1.253	0.861
			0.95	1.516	1.053	1.359	0.927	0.677	0.397
Average				0.928	0.624	1.039	0.782	0.272	0.231

Table 2. MSE of $\hat{\lambda}$ in which $(k, n) = (5, 50)$

κ	δ	ρ_y	ρ_x	Method					
				1	2	3	4	5	6
0	0	0.2	0.2	1.394	1.148	1.504	1.277	0.027	0.015
			0.95	1.654	1.363	2.645	2.397	0.022	0.012
	0.95	0.2	1.164	0.941	1.218	0.998	0.023	0.017	
		0.95	1.656	1.364	2.680	2.436	0.020	0.010	
5	1	0.2	0.2	0.010	0.001	0.003	0.003	0.076	0.136
			0.95	0.288	0.176	0.413	0.310	0.406	0.442
		0.95	0.2	0.003	0.003	0.001	0.008	0.015	0.048
			0.95	0.232	0.133	0.187	0.121	0.322	0.377
	3	0.2	0.2	0.000	0.013	0.001	0.016	0.002	0.021
			0.95	1.550	1.281	1.281	1.040	0.460	0.337
		0.95	0.2	1.289	1.065	1.276	1.054	1.253	1.033
			0.95	1.483	1.226	1.335	1.100	0.928	0.731
Average				0.894	0.726	1.045	0.897	0.296	0.265

Table 3. MSE of $\hat{\lambda}$ in which $(k, n) = (10, 30)$

κ	δ	ρ_y	ρ_x	Method					
				1	2	3	4	5	6
0	0	0.2	0.2	1.293	0.798	0.977	0.773	0.016	0.003
			0.95	1.375	0.849	1.601	1.261	0.016	0.004
		0.95	0.2	1.292	0.798	0.930	0.701	0.017	0.004
			0.95	1.332	0.815	1.539	1.111	0.017	0.004
5	1	0.2	0.2	0.070	0.001	0.028	0.011	0.311	0.461
			0.95	0.328	0.104	0.375	0.202	0.311	0.355
		0.95	0.2	0.056	0.000	0.020	0.007	0.145	0.284
			0.95	0.303	0.090	0.283	0.127	0.308	0.361
	3	0.2	0.2	0.030	0.004	0.011	0.015	0.017	0.088
			0.95	0.181	0.031	0.110	0.029	0.329	0.438
		0.95	0.2	0.018	0.011	0.007	0.022	0.002	0.054
			0.95	1.369	0.853	1.130	0.673	0.241	0.106
10	1	0.2	0.2	0.035	0.003	0.010	0.019	0.267	0.450
			0.95	0.283	0.079	0.277	0.141	0.339	0.392
		0.95	0.2	0.043	0.002	0.015	0.014	0.121	0.277
			0.95	0.303	0.090	0.215	0.071	0.282	0.348
	3	0.2	0.2	1.299	0.810	1.230	0.755	0.936	0.535
			0.95	1.411	0.879	1.115	0.657	0.122	0.043
		0.95	0.2	1.307	0.815	1.259	0.777	1.059	0.624
			0.95	1.399	0.872	1.159	0.693	0.233	0.094
Average				0.686	0.395	0.615	0.403	0.255	0.246

particular, it strongly improved the PMSE when δ and κ are small. Based on these results, we recommend using Method 6 to optimize the multivariate ridge regression.

A. Appendix

A.1. Proof of THEOREM 1. In this subsection, we prove THEOREM 1, which shows the properties of $\hat{\theta}(\lambda)$. Using d_j and \mathbf{z}_j in (4), we can write $\text{tr}(\mathbf{W}_\theta \mathbf{S}^{-1})$ and $\text{tr}(\mathbf{M}_\theta^{-1} \mathbf{M}_0)$ in (1) as

$$g(\theta) = \text{tr}(\mathbf{W}_\theta \mathbf{S}^{-1}) = \text{tr}(\mathbf{Y}' \mathbf{Y} \mathbf{S}^{-1}) - 2 \sum_{j=1}^k \frac{\|\mathbf{z}_j\|^2}{d_j + \theta} + \sum_{j=1}^k \frac{\|\mathbf{z}_j\|^2 d_j}{(d_j + \theta)^2} - n \hat{\boldsymbol{\mu}}' \mathbf{S}^{-1} \hat{\boldsymbol{\mu}},$$

$$h(\theta) = \text{tr}(\mathbf{M}_\theta^{-1} \mathbf{M}_0) = \sum_{j=1}^k \frac{d_j}{d_j + \theta}.$$

Table 4. MSE of $\hat{\lambda}$ in which $(k, n) = (10, 50)$

κ	δ	ρ_y	ρ_x	Method					
				1	2	3	4	5	6
0	0	0.2	0.2	1.173	0.940	0.977	0.862	0.016	0.007
			0.95	1.282	1.028	1.562	1.409	0.008	0.004
		0.95	0.2	1.129	0.902	0.890	0.751	0.014	0.007
			0.95	1.285	1.031	1.598	1.403	0.007	0.003
5	1	0.2	0.2	0.012	0.000	0.005	0.004	0.200	0.287
			0.95	0.203	0.111	0.183	0.126	0.407	0.437
		0.95	0.2	0.007	0.001	0.003	0.004	0.058	0.115
			0.95	0.182	0.096	0.120	0.067	0.372	0.414
	3	0.2	0.2	0.005	0.002	0.003	0.003	0.001	0.015
			0.95	0.014	0.000	0.013	0.024	0.492	0.601
		0.95	0.2	0.006	0.001	0.004	0.002	0.001	0.006
			0.95	1.261	1.016	1.138	0.917	0.455	0.328
10	1	0.2	0.2	0.004	0.002	0.001	0.008	0.158	0.248
			0.95	0.171	0.087	0.102	0.054	0.426	0.465
		0.95	0.2	0.007	0.001	0.004	0.003	0.047	0.104
			0.95	0.123	0.054	0.063	0.022	0.430	0.487
	3	0.2	0.2	1.224	0.986	1.217	0.979	1.066	0.8444
			0.95	1.325	1.067	1.161	0.912	0.229	0.147
		0.95	0.2	1.211	0.976	1.206	0.971	1.124	0.897
			0.95	1.304	1.050	1.194	0.951	0.459	0.325
Average				0.596	0.468	0.572	0.474	0.298	0.287

Since $d_j > 0$ and $\theta \geq 0$, we have

$$\dot{g}(\theta) = \frac{\partial g(\theta)}{\partial \theta} = 2\theta \sum_{j=1}^k \frac{\|z_j\|^2}{(d_j + \theta)^3} \geq 0, \quad (11)$$

with equality if and only if $\theta = 0$ or $\theta \rightarrow \infty$, and

$$\dot{h}(\theta) = \frac{\partial h(\theta)}{\partial \theta} = - \sum_{j=1}^k \frac{d_j}{(d_j + \theta)^2} \leq 0, \quad (12)$$

with equality if and only if $\theta \rightarrow \infty$. Therefore, $g(\theta)$ and $h(\theta)$ are strictly monotonic increasing and decreasing functions of $\theta \in [0, \infty]$, respectively. Since $GC_p(\theta, \lambda) = \lambda g(\theta) + 2ph(\theta)$, these results imply that

Table 5. MSE of $\hat{\lambda}$ in which $(k, n) = (15, 30)$

κ	δ	ρ_y	ρ_x	Method						
				1	2	3	4	5	6	
0	0	0.2	0.2	1.311	0.660	1.046	0.716	0.014	0.002	
			0.95	1.514	0.762	1.364	0.940	0.011	0.001	
	0.95	0.2	1.217	0.594	0.922	0.540	0.012	0.004		
		0.95	1.462	0.726	1.326	0.861	0.013	0.003		
5	1	0.2	0.2	0.285	0.042	0.163	0.036	0.242	0.328	
			0.95	0.864	0.328	0.717	0.356	0.080	0.096	
		0.95	0.2	0.030	0.024	0.013	0.048	0.207	0.471	
			0.95	0.316	0.044	0.201	0.049	0.338	0.412	
	3	0.2	0.2	0.012	0.047	0.005	0.069	0.181	0.473	
			0.95	0.232	0.017	0.141	0.048	0.427	0.522	
		0.95	0.2	0.004	0.063	0.003	0.068	0.001	0.111	
			0.95	0.139	0.002	0.043	0.018	0.224	0.390	
	10	1	0.2	0.2	0.099	0.001	0.038	0.020	0.336	0.506
				0.95	0.475	0.111	0.340	0.099	0.253	0.296
			0.95	0.2	0.044	0.014	0.018	0.033	0.097	0.301
				0.95	0.315	0.044	0.169	0.038	0.315	0.398
3		0.2	0.2	1.308	0.667	1.139	0.562	0.621	0.249	
			0.95	1.542	0.787	1.145	0.547	0.057	0.013	
		0.95	0.2	0.002	0.073	0.001	0.084	0.002	0.118	
			0.95	0.114	0.002	0.026	0.024	0.210	0.376	
15	1	0.2	0.2	0.071	0.005	0.025	0.023	0.260	0.462	
			0.95	0.468	0.109	0.307	0.074	0.244	0.292	
		0.95	0.2	0.033	0.020	0.013	0.040	0.074	0.267	
			0.95	1.540	0.786	1.142	0.540	0.050	0.011	
	3	0.2	0.2	1.259	0.642	1.146	0.571	0.841	0.374	
			0.95	1.527	0.779	1.121	0.556	0.104	0.029	
		0.95	0.2	0.006	0.056	0.004	0.063	0.001	0.083	
			0.95	0.121	0.002	0.033	0.018	0.149	0.305	
Average				0.583	0.264	0.450	0.251	0.192	0.246	

$$\hat{\theta}(\infty) = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, \infty) = \arg \min_{\theta \in [0, \infty]} \text{tr}(W_\theta S^{-1}) = 0,$$

$$\hat{\theta}(0) = \arg \min_{\theta \in [0, \infty]} GC_p(\theta, 0) = \arg \min_{\theta \in [0, \infty]} \text{tr}(M_\theta^{-1} M_0) = \infty.$$

On the other hand, from (11) and (12), we derive

Table 6. MSE of $\hat{\lambda}$ in which $(k, n) = (15, 50)$

κ	δ	ρ_y	ρ_x	Method					
				1	2	3	4	5	6
0	0	0.2	0.2	1.166	0.903	0.867	0.753	0.006	0.002
			0.95	1.256	0.973	1.247	1.059	0.004	0.001
		0.95	0.2	1.167	0.904	0.816	0.701	0.007	0.003
			0.95	1.254	0.972	1.163	0.997	0.004	0.002
5	1	0.2	0.2	0.023	0.001	0.023	0.025	0.520	0.625
			0.95	0.289	0.163	0.210	0.141	0.318	0.340
		0.95	0.2	0.002	0.007	0.001	0.009	0.073	0.154
			0.95	0.044	0.006	0.010	0.006	0.555	0.630
	3	0.2	0.2	0.000	0.015	0.000	0.017	0.053	0.126
			0.95	0.035	0.003	0.011	0.007	0.557	0.645
		0.95	0.2	0.006	0.041	0.006	0.042	0.013	0.059
			0.95	0.004	0.005	0.000	0.015	0.117	0.205
10	1	0.2	0.2	0.012	0.000	0.004	0.005	0.278	0.392
			0.95	0.174	0.080	0.078	0.036	0.441	0.472
		0.95	0.2	0.004	0.004	0.003	0.006	0.034	0.095
			0.95	0.063	0.014	0.021	0.005	0.452	0.532
	3	0.2	0.2	1.171	0.912	1.158	0.900	0.946	0.714
			0.95	1.261	0.982	1.069	0.817	0.126	0.073
		0.95	0.2	0.003	0.034	0.003	0.034	0.008	0.046
			0.95	0.004	0.005	0.000	0.014	0.088	0.168
15	1	0.2	0.2	0.006	0.003	0.001	0.011	0.232	0.345
			0.95	0.142	0.059	0.060	0.020	0.449	0.494
		0.95	0.2	0.002	0.007	0.001	0.009	0.030	0.088
			0.95	0.053	0.010	0.015	0.002	0.442	0.523
	3	0.2	0.2	1.197	0.932	1.184	0.920	1.002	0.761
			0.95	1.270	0.989	1.056	0.812	0.192	0.122
		0.95	0.2	0.007	0.002	0.007	0.002	0.003	0.005
			0.95	0.016	0.000	0.006	0.003	0.051	0.115
Average				0.380	0.287	0.322	0.263	0.250	0.276

$$\frac{\partial GC_p(\theta, \lambda)}{\partial \theta} = GC_p(\theta, \lambda) = 2 \sum_{j=1}^k \frac{pd_j^2(\theta r_{\lambda,j} - 1)}{(d_j + \theta)^3} = \sum_{j=1}^k \phi_j(\theta|\lambda),$$

where $r_{\lambda,j}$ is given by (5). Note that $\phi_j(\theta|\lambda) \leq 0$ when $\theta \in [0, (r_{\lambda,1}^+)^{-1}]$. Therefore, $GC_p(\theta, \lambda) \leq 0$ when $\theta \in [0, (r_{\lambda,1}^+)^{-1}]$. These imply that $GC_p(\theta, \lambda)$ is a

Table 7. Simulation results for the case in which $(k, n) = (5, 30)$

κ	δ	ρ_y	ρ_x	Method							
				1	2	3	4	5	6	7	8
0	0	0.2	0.2	88.34	87.53	88.00	87.16	87.30	86.76	87.44	86.88
			0.95	90.24	89.05	90.52	89.32	88.34	87.46	88.63	87.73
		0.95	0.2	88.49	87.64	88.15	87.28	87.39	86.82	87.53	86.94
			0.95	90.25	89.03	90.37	89.15	88.28	87.41	88.56	87.67
5	1	0.2	0.2	<i>94.80</i>	94.72	94.81	94.90	94.87	95.10	94.72	94.91
			0.95	92.17	91.54	92.09	91.48	91.25	90.96	91.34	<i>91.02</i>
		0.95	0.2	97.34	<i>97.31</i>	97.34	97.36	97.37	97.46	97.29	97.41
			0.95	93.47	92.94	93.47	93.00	92.74	92.43	92.76	<i>92.45</i>
	3	0.2	0.2	98.92	98.93	98.92	98.94	98.92	98.95	98.91	98.98
			0.95	95.44	95.23	95.69	95.75	95.40	95.43	95.19	<i>95.21</i>
		0.95	0.2	99.37	99.37	99.37	99.38	99.37	99.38	99.37	99.40
			0.95	<i>97.61</i>	97.54	97.66	97.82	97.72	97.86	97.54	97.67
Average				93.87	93.40	93.87	93.46	93.25	93.00	93.27	<i>93.02</i>

Table 8. Simulation results for the case in which $(k, n) = (5, 50)$

κ	δ	ρ_y	ρ_x	Method							
				1	2	3	4	5	6	7	8
0	0	0.2	0.2	92.22	91.94	92.06	91.70	91.63	91.45	91.73	<i>91.54</i>
			0.95	93.29	92.91	93.52	93.08	92.10	91.80	92.27	<i>91.98</i>
		0.95	0.2	92.22	91.96	92.08	91.72	91.67	91.49	91.76	<i>91.58</i>
			0.95	93.24	92.83	93.42	93.01	92.05	91.78	92.21	<i>91.93</i>
5	1	0.2	0.2	<i>97.46</i>	97.45	<i>97.46</i>	97.48	97.52	97.57	<i>97.46</i>	97.52
			0.95	94.96	94.79	95.00	94.82	94.54	94.46	94.57	<i>94.48</i>
		0.95	0.2	98.75	98.76	98.76	98.77	98.77	98.79	98.77	98.81
			0.95	96.19	96.06	96.28	96.20	95.96	95.92	<i>95.91</i>	95.87
3	0.2	0.2	0.2	99.57	99.57	99.57	99.57	99.57	99.57	99.57	<i>99.58</i>
			0.95	97.70	<i>97.67</i>	97.80	97.84	97.78	97.83	97.66	97.70
		0.95	0.2	99.79	<i>99.80</i>	99.79	<i>99.80</i>	99.79	<i>99.80</i>	<i>99.80</i>	99.81
			0.95	98.73	98.73	<i>98.74</i>	98.75	98.77	98.80	<i>98.74</i>	98.78
Average				96.18	96.04	96.21	96.06	95.84	95.77	95.87	<i>95.80</i>

Table 9. Simulation results for the case in which $(k, n) = (10, 30)$

κ	δ	ρ_y	ρ_x	Method							
				1	2	3	4	5	6	7	8
0	0	0.2	0.2	79.07	77.26	78.10	76.70	77.69	76.47	77.86	76.56
			0.95	81.19	78.61	80.54	78.24	78.86	77.14	79.13	77.28
		0.95	0.2	78.78	77.07	77.90	76.56	77.48	76.33	77.64	76.41
			0.95	81.01	78.44	80.29	77.72	78.72	76.89	79.03	77.08
5	1	0.2	0.2	87.14	86.60	87.12	87.07	86.75	86.88	86.72	86.76
			0.95	83.49	81.63	82.99	81.41	81.83	80.74	82.01	80.86
		0.95	0.2	91.60	91.36	91.45	91.52	91.46	91.66	91.38	91.56
			0.95	85.21	83.47	84.66	83.12	83.70	82.59	83.87	82.68
	3	0.2	0.2	95.74	95.71	95.71	95.80	95.72	95.85	95.67	95.83
			0.95	88.41	87.52	88.34	87.83	87.67	87.21	87.73	87.23
		0.95	0.2	97.50	97.51	97.50	97.56	97.50	97.56	97.47	97.58
			0.95	91.85	91.67	91.91	92.63	91.77	92.18	91.66	92.02
10	1	0.2	0.2	92.02	91.74	91.93	91.93	91.85	92.12	91.76	91.98
			0.95	84.62	82.96	84.27	82.91	83.14	82.18	83.29	82.28
		0.95	0.2	94.53	94.37	94.50	94.57	94.47	94.68	94.36	94.57
			0.95	86.19	84.74	85.81	84.66	84.91	84.02	85.03	84.10
	3	0.2	0.2	98.57	98.64	98.59	98.71	98.59	98.72	98.57	98.77
			0.95	91.50	90.99	91.59	91.67	91.12	91.23	91.06	91.12
		0.95	0.2	99.69	99.73	99.70	99.76	99.70	99.76	99.70	99.81
			0.95	94.15	93.83	94.20	94.37	94.00	94.19	93.87	94.02
Average				89.11	88.19	88.85	88.24	88.35	87.92	88.39	87.93

monotonic decreasing function with respect to $\theta \in [0, (r_{\lambda,1}^+)^{-1}]$. Thus, $\hat{\theta}(\lambda) > (r_{\lambda,1}^+)^{-1}$. On the other hand, if $\max_{j=1,\dots,k} r_{\lambda,j} \leq 0$ is satisfied, $\phi_j(\theta|\lambda) \leq 0$ holds for any θ . This fact means $GC_p(\theta, \lambda)$ is a monotonic decreasing function with respect to $\theta \in [0, \infty]$. Hence, we can see that $\hat{\theta}(\lambda) = \infty$ when $\max_{j=1,\dots,k} r_{\lambda,j} \leq 0$. Since $\max_{j=1,\dots,k} r_{\lambda,j} \leq 0$ holds when $\lambda < \min_{j=1,\dots,k} pd_j / \|\mathbf{z}_j\|^2$, $\hat{\theta}(\lambda) = \infty$ when $\lambda < \min_{j=1,\dots,k} pd_j / \|\mathbf{z}_j\|^2$ is satisfied.

Using Equation (3) and $\dot{GC}_p(\theta, \lambda) = \lambda \dot{g}(\theta) + 2p\dot{h}(\theta)$, we have

$$\frac{\partial}{\partial \lambda} \dot{GC}_p(\hat{\theta}(\lambda), \lambda) = \dot{g}(\hat{\theta}(\lambda)) + \frac{\partial \hat{\theta}(\lambda)}{\partial \lambda} \ddot{GC}_p(\hat{\theta}(\lambda), \lambda) = 0, \quad (13)$$

where $\ddot{GC}_p(\theta, \lambda) = \partial^2 GC_p(\theta, \lambda) / (\partial \theta^2)$. Since $\hat{\theta}(\lambda)$ satisfies (2), i.e., $\hat{\theta}(\lambda)$ is the minimizer of $GC_p(\theta, \lambda)$, $GC_p(\theta, \lambda)$ is a convex function around the neigh-

Table 10. Simulation results for the case in which $(k, n) = (10, 50)$

κ	δ	ρ_y	ρ_x	Method							
				1	2	3	4	5	6	7	8
0	0	0.2	0.2	84.93	84.48	84.60	84.09	84.41	84.09	84.51	84.17
			0.95	86.22	85.43	85.59	84.96	84.99	84.48	85.19	84.63
	0.95	0.2	84.98	84.50	84.62	84.13	84.44	84.10	84.54	84.18	
		0.95	86.20	85.42	85.58	85.00	85.04	84.60	85.20	84.68	
5	1	0.2	0.2	92.65	92.60	92.67	92.72	92.66	92.74	92.60	92.66
			0.95	88.43	87.94	88.14	87.77	87.70	87.40	87.82	87.48
		0.95	0.2	95.43	95.43	95.44	95.47	95.45	95.51	95.42	95.50
			0.95	90.03	89.69	89.94	89.68	89.55	89.31	89.61	89.37
	3	0.2	0.2	97.90	97.92	97.90	97.93	97.90	97.93	97.90	97.95
			0.95	93.17	93.09	93.36	93.52	93.15	93.22	93.09	93.14
		0.95	0.2	99.27	99.28	99.27	99.28	99.27	99.28	99.27	99.30
			0.95	95.51	95.47	95.49	95.53	95.54	95.61	95.48	95.54
10	1	0.2	0.2	96.50	96.47	96.53	96.55	96.52	96.57	96.47	96.53
			0.95	89.71	89.33	89.58	89.32	89.16	88.94	89.23	89.00
		0.95	0.2	97.64	97.64	97.66	97.68	97.66	97.70	97.64	97.70
			0.95	91.07	90.75	91.09	90.91	90.63	90.48	90.67	90.51
	3	0.2	0.2	99.44	99.44	99.44	99.44	99.44	99.44	99.44	99.45
			0.95	95.66	95.64	95.85	95.97	95.75	95.89	95.65	95.77
		0.95	0.2	99.46	99.46	99.46	99.47	99.46	99.47	99.46	99.47
			0.95	97.31	97.29	97.35	97.38	97.36	97.43	97.30	97.37
Average				93.08	92.86	92.98	92.84	92.80	92.71	92.83	92.72

borhood of $\hat{\theta}(\lambda)$. Hence, we have $\ddot{GC}_p(\hat{\theta}(\lambda), \lambda) > 0$. Using this result and Equation (13), we obtain

$$\frac{\partial \hat{\theta}(\lambda)}{\partial \lambda} = - \frac{\dot{g}(\hat{\theta}(\lambda))}{\ddot{GC}_p(\hat{\theta}(\lambda), \lambda)}.$$

We derive $\dot{g}(\hat{\theta}(\lambda)) \geq 0$ because $g(\theta)$ is a strictly monotonic increasing function of $\theta \in [0, \infty]$. Hence, $\partial \hat{\theta}(\lambda) / (\partial \lambda) \leq 0$ is obtained. This implies that $\hat{\theta}(\lambda)$ is a monotonic decreasing function with respect to λ .

A.2. Proof of THEOREM 2. In this subsection, we present the proof of THEOREM 2, which describes the expansion of $\hat{\theta}(\lambda)$. In order to prove this theorem, we expand the GC_p criterion in (1) under fixed λ . Recall that $X' \mathbf{1}_n = \mathbf{0}_k$, $(n - k - 1)S = Y'(I_n - \mathbf{1}_n \mathbf{1}'_n / n - X M_0^{-1} X') Y$, and $Q' X' X Q = D$. Hence, we derive

Table 11. Simulation results for the case in which $(k, n) = (15, 30)$

κ	δ	ρ_y	ρ_x	Method							
				1	2	3	4	5	6	7	8
0	0	0.2	0.2	72.83	69.39	71.66	68.77	70.91	68.51	71.05	68.56
			0.95	76.39	71.49	74.39	70.03	72.67	69.23	72.91	69.43
		0.95	0.2	73.86	69.79	72.30	69.03	71.62	68.75	71.78	68.81
			0.95	76.96	71.53	74.71	70.21	72.86	69.12	73.14	69.26
5	1	0.2	0.2	81.70	79.12	81.18	79.31	80.22	78.70	80.33	78.73
			0.95	78.31	73.62	76.23	72.84	74.79	71.99	75.01	72.06
		0.95	0.2	93.65	93.73	93.78	94.48	93.61	94.47	93.50	94.31
			0.95	83.27	80.58	82.40	80.94	81.18	80.12	81.29	80.16
	3	0.2	0.2	95.37	95.79	95.46	96.48	95.45	96.66	95.33	96.48
			0.95	83.89	81.25	83.09	81.71	81.83	80.83	81.96	80.84
		0.95	0.2	99.19	99.30	99.19	99.32	99.20	99.34	99.19	99.40
			0.95	91.36	90.77	91.08	91.36	90.87	91.24	90.84	91.15
10	1	0.2	0.2	86.74	85.67	86.32	86.26	86.01	86.17	86.02	86.09
			0.95	80.23	76.29	78.81	75.65	77.28	74.85	77.44	74.89
		0.95	0.2	94.47	94.45	94.46	94.87	94.43	94.95	94.33	94.84
			0.95	84.20	81.78	83.37	82.11	82.29	81.24	82.40	81.27
	3	0.2	0.2	97.61	97.51	97.59	97.73	97.57	97.76	97.48	97.72
			0.95	86.41	84.68	85.85	85.24	85.06	84.55	85.13	84.55
		0.95	0.2	99.24	99.32	99.25	99.35	99.25	99.35	99.24	99.40
			0.95	92.31	91.76	91.97	92.10	91.90	92.16	91.86	92.08
15	1	0.2	0.2	89.95	89.49	89.75	90.12	89.59	90.23	89.54	90.11
			0.95	81.17	77.65	79.74	77.00	78.54	76.39	78.69	76.45
		0.95	0.2	95.13	95.08	95.11	95.38	95.09	95.47	95.00	95.39
			0.95	84.74	82.63	84.10	83.32	83.04	82.47	83.13	82.49
	3	0.2	0.2	97.99	98.09	98.00	98.23	97.99	98.25	97.95	98.28
			0.95	88.99	87.89	88.33	88.61	88.07	88.28	88.09	88.23
		0.95	0.2	99.16	99.21	99.16	99.23	99.16	99.23	99.15	99.27
			0.95	92.82	92.25	92.46	92.53	92.40	92.56	92.36	92.48
Average				87.78	86.08	87.13	86.15	86.53	85.82	86.58	85.81

$$\begin{aligned}
W_\theta \mathbf{S}^{-1} &= \mathbf{Y}'(\mathbf{I}_n - \mathbf{1}_n \mathbf{1}_n' / n - \mathbf{X} \mathbf{M}_\theta^{-1} \mathbf{X}')^2 \mathbf{Y} \mathbf{S}^{-1} \\
&= (n - k - 1) \mathbf{I}_p + \mathbf{Y}' \mathbf{X} (\mathbf{M}_0^{-1} - 2\mathbf{M}_\theta^{-1} + \mathbf{M}_\theta^{-1} \mathbf{M}_0 \mathbf{M}_\theta^{-1}) \mathbf{X}' \mathbf{Y} \mathbf{S}^{-1} \\
&= (n - k - 1) \mathbf{I}_p + \mathbf{Y}' \mathbf{X} \mathbf{Q} \mathbf{D}^{-1/2} \{ \mathbf{I}_k - \mathbf{D}(\mathbf{D} + \theta \mathbf{I}_k)^{-1} \}^2 \mathbf{D}^{-1/2} \mathbf{Q}' \mathbf{X}' \mathbf{Y} \mathbf{S}^{-1}.
\end{aligned}$$

Based on this result, the GC_p criterion is expressed as

Table 12. Simulation results for the case in which $(k, n) = (15, 50)$

κ	δ	ρ_y	ρ_x	Method							
				1	2	3	4	5	6	7	8
0	0	0.2	0.2	79.00	78.27	78.34	77.85	78.34	77.87	78.42	77.93
			0.95	80.06	78.92	78.88	78.24	78.74	78.07	78.90	78.17
		0.95	0.2	78.96	78.24	78.35	77.78	78.31	77.80	78.39	77.86
			0.95	80.47	79.24	79.11	78.33	79.01	78.24	79.18	78.34
5	1	0.2	0.2	89.06	88.80	89.42	89.64	88.86	88.89	88.84	88.83
			0.95	82.51	81.71	81.91	81.31	81.60	81.14	81.70	81.20
		0.95	0.2	97.59	97.65	97.60	97.69	97.63	97.76	97.61	97.77
			0.95	88.51	88.15	88.41	88.35	88.14	88.05	88.16	88.04
	3	0.2	0.2	98.10	98.12	98.10	98.13	98.11	98.16	98.10	98.18
			0.95	89.07	88.81	89.06	89.18	88.81	88.87	88.82	88.84
		0.95	0.2	99.76	99.77	99.76	99.77	99.76	99.77	99.76	99.77
			0.95	95.77	95.74	95.77	95.81	95.77	95.83	95.73	95.81
10	1	0.2	0.2	93.46	93.41	93.49	93.55	93.46	93.60	93.40	93.53
			0.95	84.38	83.67	83.97	83.45	83.59	83.15	83.68	83.22
		0.95	0.2	98.01	98.03	98.02	98.05	98.02	98.07	98.01	98.08
			0.95	89.21	88.93	89.27	89.22	88.92	88.88	88.93	88.87
	3	0.2	0.2	98.84	98.86	98.84	98.86	98.85	98.87	98.85	98.89
			0.95	92.05	91.85	92.06	92.11	91.89	91.92	91.86	91.87
		0.95	0.2	99.68	99.69	99.68	99.69	99.68	99.69	99.69	99.70
			0.95	96.56	96.56	96.57	96.62	96.58	96.65	96.55	96.64
15	1	0.2	0.2	95.10	95.05	95.13	95.17	95.10	95.19	95.05	95.14
			0.95	85.64	85.11	85.46	85.14	85.04	84.76	85.11	84.80
		0.95	0.2	98.37	98.39	98.38	98.41	98.39	98.43	98.38	98.44
			0.95	89.86	89.58	89.72	89.75	89.58	89.58	89.59	89.56
	3	0.2	0.2	99.46	99.46	99.46	99.46	99.46	99.46	99.45	99.48
			0.95	93.82	93.65	93.77	93.73	93.69	93.71	93.65	93.66
		0.95	0.2	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
			0.95	96.76	96.75	96.78	96.81	96.77	96.82	96.74	96.80
Average				91.79	91.52	91.62	91.50	91.50	91.40	91.52	91.41

$$GC_p(\theta, \lambda) = \lambda(n - k - 1)p + \lambda\theta^2 \operatorname{tr}\{\mathbf{D}^{-1/2}\mathbf{Q}'\mathbf{V}\mathbf{Q}\mathbf{D}^{-1/2}(\mathbf{D} + \theta\mathbf{I}_k)^{-2}\} + 2p\{\mathbf{D}(\mathbf{D} + \theta\mathbf{I}_k)^{-1}\}.$$

Letting $t_j = (\mathbf{D}^{-1/2}\mathbf{Q}'\mathbf{V}\mathbf{Q}\mathbf{D}^{-1/2})_{jj}$, we obtain

$$GC_p(\theta, \lambda) = \lambda(n - k - 1)p + \sum_{i=1}^k \left\{ \lambda \left(1 + \frac{\theta}{d_i} \right)^{-2} \frac{\theta^2 t_i}{d_i^2} + 2p \left(1 + \frac{\theta}{d_i} \right)^{-1} \right\}.$$

By Taylor expansion around $\theta = 0$, we have

$$\begin{aligned}
GC_p(\theta, \lambda) &= \lambda(n-k-1)p + \lambda \sum_{i=1}^k \sum_{\ell=1}^{\infty} \frac{(-1)^{\ell+1} \ell \theta^{\ell+1}}{d_i^{\ell+1}} t_i \\
&\quad + 2p \sum_{i=1}^k \left(1 - \sum_{\ell=1}^{\infty} \frac{(-1)^{\ell+1}}{d_i^{\ell}} \theta^{\ell} \right) \\
&= (\lambda(n-k-1) + 2k)p \\
&\quad + \sum_{\ell=1}^{\infty} \left\{ \lambda (-1)^{\ell+1} \ell \theta^{\ell+1} \sum_{i=1}^k \frac{t_i}{d_i^{\ell+1}} - 2p (-1)^{\ell+1} \theta^{\ell} \sum_{i=1}^k \frac{1}{d_i^{\ell}} \right\} \\
&= (\lambda(n-k-1) + 2k)p \\
&\quad + \sum_{\ell=1}^{\infty} (-1)^{\ell+1} \theta^{\ell} \{ \lambda \ell \theta \operatorname{tr}(\mathbf{V} \mathbf{M}_0^{-(\ell+2)}) - 2p \operatorname{tr}(\mathbf{M}_0^{-\ell}) \}.
\end{aligned}$$

Recall that $a_j = n^j \operatorname{tr}(\mathbf{V} \mathbf{M}_0^{-(j+2)})$ and $b_j = n^j \operatorname{tr}(\mathbf{M}_0^{-j})$. It follows that

$$GC_p(\theta, \lambda) = (\lambda(n-k-1) + 2k)p + \lim_{L \rightarrow \infty} \sum_{\ell=1}^L \frac{(-1)^{\ell+1} \theta^{\ell}}{n^{\ell}} \{ \lambda \ell \theta a_{\ell} - 2p b_{\ell} \}.$$

Then, the following equation is derived:

$$\frac{\partial}{\partial \theta} GC_p(\theta, \lambda) = \lim_{L \rightarrow \infty} \sum_{\ell=1}^L \frac{(-1)^{\ell+1} \ell \theta^{\ell-1}}{n^{\ell}} \{ \lambda(\ell+1) a_{\ell} \theta - 2p b_{\ell} \}.$$

Using the above equation and $\hat{\theta}(\lambda)$ satisfies (3), we obtain the equation in THEOREM 2.

A.3. Proof of THEOREM 3. In this subsection, we prove THEOREM 3, which shows the risk function with respect to λ . Recall that $g(\theta) = \operatorname{tr}(\mathbf{W}_{\theta} \mathbf{S}^{-1})$, $h(\theta) = \operatorname{tr}(\mathbf{M}_{\theta}^{-1} \mathbf{M}_0)$ and $GC_p(\theta, \lambda) = \lambda g(\theta) + 2p h(\theta)$. Since $\hat{\theta}(\lambda)$ satisfies (3), we obtain

$$\begin{aligned}
0 &= \frac{\partial}{\partial(\mathbf{Y})_{ij}} \left(\lambda \frac{\partial g(\theta)}{\partial \theta} \Big|_{\theta=\hat{\theta}(\lambda)} + 2p \frac{\partial h(\theta)}{\partial \theta} \Big|_{\theta=\hat{\theta}(\lambda)} \right) \\
&= \frac{\partial \hat{\theta}(\lambda)}{\partial(\mathbf{Y})_{ij}} \{ \lambda \ddot{g}(\hat{\theta}(\lambda)) + 2p \ddot{h}(\hat{\theta}(\lambda)) \} + \lambda \frac{\partial \dot{g}(\hat{\theta}(\lambda))}{\partial(\mathbf{Y})_{ij}},
\end{aligned}$$

where $\dot{g}(\hat{\theta}(\lambda)) = \partial g(\theta)/(\partial \theta)|_{\theta=\hat{\theta}(\lambda)}$, $\ddot{g}(\hat{\theta}(\lambda)) = \partial^2 g(\theta)/(\partial \theta^2)|_{\theta=\hat{\theta}(\lambda)}$, and $\dot{h}(\hat{\theta}(\lambda)) = \partial h(\theta)/(\partial \theta)|_{\theta=\hat{\theta}(\lambda)}$. Thus, we obtain

$$\frac{\partial \hat{\theta}(\lambda)}{\partial (\mathbf{Y})_{ij}} = -\frac{\lambda(\partial \dot{g}(\hat{\theta}(\lambda)))/(\partial (\mathbf{Y})_{ij})}{\ddot{G}\ddot{C}_p(\hat{\theta}(\lambda), \lambda)},$$

because, from Appendix A.1, $\ddot{G}\ddot{C}_p(\hat{\theta}(\lambda), \lambda) = \lambda \ddot{g}(\hat{\theta}(\lambda)) + 2p \dot{h}(\hat{\theta}(\lambda)) > 0$. By simple calculation, we have $\dot{g}(\hat{\theta}(\lambda)) = 2\hat{\theta}(\lambda) \text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-3} \mathbf{V})$. As in the proof of LEMMA 1, which is given in Appendix A.4, we obtain $\partial \text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-3} \mathbf{V})/(\partial (\mathbf{Y})_{ij}) = 2(\mathbf{S}^{-1} \mathbf{Y}' \mathbf{X} \mathbf{M}_{\hat{\theta}(\lambda)}^{-3} \mathbf{X}' \{\mathbf{I}_n - \mathbf{Y} \mathbf{S}^{-1} \mathbf{Y}' \mathbf{H}/(n-k-1)\})_{ji}$. Hence, we derive

$$\sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\hat{\theta}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \frac{\partial \hat{\theta}(\lambda)}{\partial (\mathbf{Y})_{ij}} = -\frac{4\lambda \hat{\theta}(\lambda) \text{tr}(\mathbf{V} \mathbf{M}_{\hat{\theta}(\lambda)}^{-5} \mathbf{M}_0)}{\ddot{G}\ddot{C}_p(\hat{\theta}(\lambda), \lambda)},$$

because $\mathbf{X}' \mathbf{1}_n = \mathbf{0}_k$ and $\mathbf{M}_{\hat{\theta}(\lambda)}^{-3} \mathbf{M}_0 \mathbf{M}_{\hat{\theta}(\lambda)}^{-2} = \mathbf{M}_{\hat{\theta}(\lambda)}^{-5} \mathbf{M}_0$. By simple calculation, we have $\ddot{G}\ddot{C}_p(\theta, \lambda) = \lambda \dot{g}(\theta) + 2p \dot{h}(\theta) = 2\{\lambda \theta \text{tr}(\mathbf{M}_{\theta}^{-3} \mathbf{V}) - p \text{tr}(\mathbf{M}_{\theta}^{-2} \mathbf{M}_0)\}$ and

$$\ddot{G}\ddot{C}_p(\hat{\theta}(\lambda), \lambda) = 2\{\lambda \text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-3} \mathbf{V}) - 3\lambda \hat{\theta}(\lambda) \text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-4} \mathbf{V}) + 2p \text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-3} \mathbf{M}_0)\}.$$

Thus, the theorem is proved.

A.4. Proof of LEMMA 1. In this subsection, we prove LEMMA 1, which shows the derivative of a_ℓ for any ℓ . Since $a_\ell = n^\ell \text{tr}(\mathbf{V} \mathbf{M}_0^{-(\ell+2)})$, $\mathbf{M}_0 = \mathbf{X}' \mathbf{X}$, and $\mathbf{V} = \mathbf{X}' \mathbf{Y} \mathbf{S}^{-1} \mathbf{Y}' \mathbf{X}$, we need only obtain the derivative of \mathbf{V} . We can see that

$$\begin{aligned} \frac{\partial \mathbf{S}^{-1}}{\partial (\mathbf{Y})_{ij}} &= -\mathbf{S}^{-1} \frac{\partial \mathbf{S}}{\partial (\mathbf{Y})_{ij}} \mathbf{S}^{-1} \\ &= -\frac{1}{n-k-1} \mathbf{S}^{-1} (\mathbf{e}_{j \cdot p} \mathbf{e}'_{i \cdot n} \mathbf{H} \mathbf{Y} + \mathbf{Y}' \mathbf{H} \mathbf{e}_{i \cdot n} \mathbf{e}'_{j \cdot p}) \mathbf{S}^{-1}, \end{aligned}$$

where $\mathbf{e}_{i \cdot n}$ is an n -dimensional vector whose i th element is one and other elements are zeros. Thus, we obtain

$$\begin{aligned} \frac{\partial \mathbf{V}}{\partial (\mathbf{Y})_{ij}} &= \mathbf{X}' \mathbf{e}_{i \cdot n} \mathbf{e}'_{j \cdot p} \mathbf{S}^{-1} \mathbf{Y}' \mathbf{X} + \mathbf{X}' \mathbf{Y} \mathbf{S}^{-1} \mathbf{e}_{j \cdot p} \mathbf{e}'_{i \cdot n} \mathbf{X} \\ &\quad - \frac{1}{n-k-1} \mathbf{X}' \mathbf{Y} \mathbf{S}^{-1} (\mathbf{e}_{j \cdot p} \mathbf{e}'_{i \cdot n} \mathbf{H} \mathbf{Y} + \mathbf{Y}' \mathbf{H} \mathbf{e}_{i \cdot n} \mathbf{e}'_{j \cdot p}) \mathbf{S}^{-1} \mathbf{Y}' \mathbf{X}. \end{aligned}$$

From $\partial a_\ell/(\partial (\mathbf{Y})_{ij}) = n^\ell \text{tr}\{\mathbf{M}_0^{-(\ell+2)} (\partial \mathbf{V})/(\partial (\mathbf{Y})_{ij})\}$, we derive this lemma.

A.5. Proof of THEOREM 4. From (8), we obtain

$$\begin{aligned} \text{PMSE}[\hat{\mathbf{Y}}_{\tilde{\theta}_{(L)}(\lambda)}] &= E_{\mathbf{Y}}[\text{tr}(\mathbf{W}_{\tilde{\theta}_{(L)}(\lambda)} \boldsymbol{\Sigma}^{-1}) + 2p\{\text{tr}(\mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-1} \mathbf{M}_0) + 1\}] \\ &\quad - 2E_{\mathbf{Y}} \left[\sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \frac{\partial \tilde{\theta}_{(L)}(\lambda)}{\partial (\mathbf{Y})_{ij}} \right]. \end{aligned}$$

Hence, we need only calculate $\partial \tilde{\theta}_{(L)}(\lambda) / (\partial (\mathbf{Y})_{ij})$. Using THEOREM 2, we derive the derivative as follows:

$$\begin{aligned} \frac{\partial \tilde{\theta}_{(L)}(\lambda)}{\partial (\mathbf{Y})_{ij}} &= \frac{\partial \tilde{\theta}_{(1)}(\lambda)}{\partial (\mathbf{Y})_{ij}} + \frac{\partial}{\partial (\mathbf{Y})_{ij}} \left[\frac{1}{2\lambda a_1} \sum_{\ell=0}^{L-1} \frac{1}{n^\ell} (-1)^{\ell+1} (\ell+1) \right. \\ &\quad \left. \times \tilde{\theta}_{(L-1)}^\ell(\lambda) \{\lambda(\ell+2)a_{\ell+1} \tilde{\theta}_{(L-1)}(\lambda) - 2pb_{\ell+1}\} \right]. \end{aligned}$$

Recall that $\tilde{\theta}_{(0)}(\lambda) = 0$. From LEMMA 1, we have

$$\begin{aligned} \frac{\partial \tilde{\theta}_{(1)}(\lambda)}{\partial (\mathbf{Y})_{ij}} &= -\frac{pb_1}{\lambda a_1^2} \frac{\partial a_1}{\partial (\mathbf{Y})_{ij}} \\ &= -\frac{2n\tilde{\theta}_{(1)}(\lambda)}{a_1} \left(\mathbf{S}^{-1} \mathbf{Y}' \mathbf{X} \mathbf{M}_0^{-3} \mathbf{X}' \left(\mathbf{I}_n - \frac{1}{n-k-1} \mathbf{Y} \mathbf{S}^{-1} \mathbf{Y}' \mathbf{H} \right) \right)_{ji}, \end{aligned}$$

and

$$\begin{aligned} &\frac{\partial}{\partial (\mathbf{Y})_{ij}} \left[\frac{1}{2\lambda a_1} \sum_{\ell=0}^{L-1} \frac{1}{n^\ell} (-1)^{\ell+1} (\ell+1) \tilde{\theta}_{(L-1)}^\ell(\lambda) \{\lambda(\ell+2)a_{\ell+1} \tilde{\theta}_{(L-1)}(\lambda) - 2pb_{\ell+1}\} \right] \\ &= -\frac{n}{\lambda a_1^2} \left(\mathbf{S}^{-1} \mathbf{Y}' \mathbf{X} \mathbf{M}_0^{-3} \mathbf{X}' \left(\mathbf{I}_n - \frac{\mathbf{Y} \mathbf{S}^{-1} \mathbf{Y}' \mathbf{H}}{n-k-1} \right) \right)_{ji} \\ &\quad \times \sum_{\ell=0}^{L-1} \frac{1}{n^\ell} (-1)^{\ell+1} (\ell+1) \tilde{\theta}_{(L-1)}^\ell(\lambda) \{\lambda(\ell+2)a_{\ell+1} \tilde{\theta}_{(L-1)}(\lambda) - 2pb_{\ell+1}\} \\ &\quad + \frac{1}{2\lambda a_1} \sum_{\ell=0}^{L-1} \frac{1}{n^\ell} (-1)^{\ell+1} (\ell+1) \tilde{\theta}_{(L-1)}^{\ell-1}(\lambda) \\ &\quad \times \{\lambda(\ell+1)(\ell+2)a_{\ell+1} \tilde{\theta}_{(L-1)}(\lambda) - 2p\ell b_{\ell+1}\} \frac{\partial \tilde{\theta}_{(L-1)}(\lambda)}{\partial (\mathbf{Y})_{ij}} \end{aligned}$$

$$\begin{aligned}
& + \frac{n}{a_1} \sum_{\ell=0}^{L-1} (-1)^{\ell+1} (\ell+1)(\ell+2) \tilde{\theta}_{(L-1)}^{\ell+1}(\lambda) \\
& \times \left(\mathbf{S}^{-1} \mathbf{Y}' \mathbf{X} \mathbf{M}_0^{-(\ell+3)} \mathbf{X}' \left(\mathbf{I}_n - \frac{1}{n-k-1} \mathbf{Y} \mathbf{S}^{-1} \mathbf{Y}' \mathbf{H} \right) \right)_{ji}.
\end{aligned}$$

Thus, we obtain

$$\begin{aligned}
& \sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \\
& \times \frac{\partial}{\partial (\mathbf{Y})_{ij}} \left[\frac{1}{2\lambda a_1} \sum_{\ell=0}^{L-1} \frac{1}{n^\ell} (-1)^{\ell+1} (\ell+1) \right. \\
& \quad \left. \times \tilde{\theta}_{(L-1)}^\ell(\lambda) \{ \lambda(\ell+2) a_{\ell+1} \tilde{\theta}_{(L-1)}(\lambda) - 2pb_{\ell+1} \} \right] \\
& = - \frac{n}{\lambda a_1^2} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{V}) \\
& \quad \times \sum_{\ell=0}^{L-1} \frac{1}{n^\ell} (-1)^{\ell+1} (\ell+1) \tilde{\theta}_{(L-1)}^\ell(\lambda) \{ \lambda(\ell+2) a_{\ell+1} \tilde{\theta}_{(L-1)}(\lambda) - 2pb_{\ell+1} \} \\
& \quad + \frac{1}{2\lambda a_1} \sum_{\ell=0}^{L-1} \frac{1}{n^\ell} (-1)^{\ell+1} (\ell+1) \tilde{\theta}_{(L-1)}^{\ell-1}(\lambda) \\
& \quad \times \{ \lambda(\ell+1)(\ell+2) a_{\ell+1} \tilde{\theta}_{(L-1)}(\lambda) - 2p\ell b_{\ell+1} \} \\
& \quad \times \sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \frac{\partial \tilde{\theta}_{(L-1)}(\lambda)}{\partial (\mathbf{Y})_{ij}} \\
& \quad + \frac{n}{a_1} \sum_{\ell=0}^{L-1} (-1)^{\ell+1} (\ell+1)(\ell+2) \tilde{\theta}_{(L-1)}^{\ell+1}(\lambda) \text{tr}(\mathbf{M}_0^{-(\ell+2)} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{V}),
\end{aligned}$$

and

$$\sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \frac{\partial \tilde{\theta}_{(1)}(\lambda)}{\partial (\mathbf{Y})_{ij}} = - \frac{2n \tilde{\theta}_{(1)}(\lambda)}{a_1} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(L)}(\lambda)}^{-2} \mathbf{V}).$$

We derive this theorem by substituting these results into (8).

A.6. The criteria for optimizing λ when we use $\tilde{\theta}_{(2)}(\lambda)$. In this subsection, we calculate the criterion for optimizing λ when we use $\tilde{\theta}_{(2)}(\lambda)$. From (6), we obtain

$$\begin{aligned}\tilde{\theta}_{(2)}(\lambda) &= \tilde{\theta}_{(1)}(\lambda) + \frac{1}{n\lambda a_1} \tilde{\theta}_{(1)}(\lambda) \{3\lambda a_2 \tilde{\theta}_{(1)}(\lambda) - 2pb_2\} \\ &= \tilde{\theta}_{(1)}(\lambda) + \frac{1}{n} \tilde{\theta}_{(1)}^2(\lambda) \left\{ 3 \frac{a_2}{a_1} - 2 \frac{b_2}{b_1} \right\}.\end{aligned}$$

Since $b_\ell = n^\ell \text{tr}(\mathbf{M}_0^{-\ell})$ does not depend on $(\mathbf{Y})_{ij}$, we derive

$$\begin{aligned}\frac{\partial \tilde{\theta}_{(2)}(\lambda)}{\partial (\mathbf{Y})_{ij}} &= \frac{\partial \tilde{\theta}_{(1)}(\lambda)}{\partial (\mathbf{Y})_{ij}} + \frac{1}{n} \frac{\partial}{\partial (\mathbf{Y})_{ij}} \left\{ \tilde{\theta}_{(1)}^2(\lambda) \left(3 \frac{a_2}{a_1} - 2 \frac{b_2}{b_1} \right) \right\} \\ &= \frac{\partial \tilde{\theta}_{(1)}(\lambda)}{\partial (\mathbf{Y})_{ij}} + \frac{1}{n} \left\{ \frac{\partial \tilde{\theta}_{(1)}^2(\lambda)}{\partial (\mathbf{Y})_{ij}} \left(3 \frac{a_2}{a_1} - 2 \frac{b_2}{b_1} \right) + 3 \tilde{\theta}_{(1)}^2(\lambda) \frac{\partial a_2/a_1}{\partial (\mathbf{Y})_{ij}} \right\} \\ &= \frac{\partial \tilde{\theta}_{(1)}(\lambda)}{\partial (\mathbf{Y})_{ij}} \left\{ 1 + \frac{2\tilde{\theta}_{(1)}(\lambda)}{n} \left(3 \frac{a_2}{a_1} - 2 \frac{b_2}{b_1} \right) \right\} + \frac{3\tilde{\theta}_{(1)}^2(\lambda)}{n} \frac{\partial a_2/a_1}{\partial (\mathbf{Y})_{ij}}.\end{aligned}$$

Hence, the third term of (7) is obtained using the result of $\partial \tilde{\theta}_{(1)}(\lambda)/(\partial (\mathbf{Y})_{ij})$ as follows:

$$\begin{aligned}& \sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \frac{\partial \tilde{\theta}_{(2)}(\lambda)}{\partial (\mathbf{Y})_{ij}} \\ &= -\frac{2n\tilde{\theta}_{(1)}(\lambda)}{a_1} \left\{ 1 + \frac{2\tilde{\theta}_{(1)}(\lambda)}{n} \left(3 \frac{a_2}{a_1} - 2 \frac{b_2}{b_1} \right) \right\} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) \\ & \quad + \frac{3\tilde{\theta}_{(1)}^2(\lambda)}{na_1^2} \sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \left(a_1 \frac{\partial a_2}{\partial (\mathbf{Y})_{ij}} - a_2 \frac{\partial a_1}{\partial (\mathbf{Y})_{ij}} \right).\end{aligned}$$

From LEMMA 1, we obtain

$$\begin{aligned}& \frac{3\tilde{\theta}_{(1)}^2(\lambda)}{na_1^2} \sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \left(a_1 \frac{\partial a_2}{\partial (\mathbf{Y})_{ij}} - a_2 \frac{\partial a_1}{\partial (\mathbf{Y})_{ij}} \right) \\ &= \frac{6\tilde{\theta}_{(1)}^2(\lambda)}{a_1^2} \sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \\ & \quad \times \left(\mathbf{S}^{-1} \mathbf{Y}' \mathbf{X} \mathbf{M}_0^{-3} (na_1 \mathbf{M}_0^{-1} - a_2 \mathbf{I}_k) \mathbf{X}' \left(\mathbf{I}_n - \frac{1}{n-k-1} \mathbf{Y} \mathbf{S}^{-1} \mathbf{Y}' \mathbf{H} \right) \right)_{ji} \\ &= \frac{6\tilde{\theta}_{(1)}^2(\lambda)}{a_1^2} \text{tr} \{ \mathbf{X} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V} \mathbf{M}_0^{-3} (na_1 \mathbf{M}_0^{-1} - a_2 \mathbf{I}_k) \mathbf{X}' \} \\ &= \frac{6\tilde{\theta}_{(1)}^2(\lambda)}{a_1^2} \text{tr} \{ \mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V} (na_1 \mathbf{M}_0^{-1} - a_2 \mathbf{I}_k) \}.\end{aligned}$$

Thus, we have

$$\begin{aligned}
& - \sum_{i=1}^n \sum_{j=1}^p (\mathbf{X} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{X}' \mathbf{Y})_{ij} \frac{\partial \tilde{\theta}_{(2)}(\lambda)}{\partial (\mathbf{Y})_{ij}} \\
&= \frac{2n\tilde{\theta}_{(1)}(\lambda)}{a_1} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) + \frac{4\tilde{\theta}_{(1)}^2(\lambda)}{a_1} \left(3 \frac{a_2}{a_1} - 2 \frac{b_2}{b_1} \right) \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) \\
&\quad + \frac{6\tilde{\theta}_{(1)}^2(\lambda)}{a_1^2} \text{tr}\{\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V} (a_2 \mathbf{I}_k - na_1 \mathbf{M}_0^{-1})\} \\
&= \frac{2n\tilde{\theta}_{(1)}(\lambda)}{a_1} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) + \frac{18\tilde{\theta}_{(1)}^2(\lambda)a_2}{a_1^2} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) \\
&\quad - \frac{8\tilde{\theta}_{(1)}^2(\lambda)b_2}{a_1 b_1} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) - \frac{6n\tilde{\theta}_{(1)}^2(\lambda)}{a_1} \text{tr}(\mathbf{M}_0^{-3} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}).
\end{aligned}$$

When we use THEOREM 4, we obtain the same result. Using this result, we obtain the C_p type criteria for optimizing λ are defined as

$$\begin{aligned}
C_p^{(2)}(\lambda) &= \text{tr}(\mathbf{W}_{\tilde{\theta}_{(2)}(\lambda)} \mathbf{S}^{-1}) + 2p \text{tr}(\mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-1} \mathbf{M}_0) \\
&\quad + \frac{4n\tilde{\theta}_{(1)}(\lambda)}{a_1} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) + \frac{36\tilde{\theta}_{(1)}^2(\lambda)a_2}{a_1^2} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) \\
&\quad - \frac{16\tilde{\theta}_{(1)}^2(\lambda)b_2}{a_1 b_1} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) - \frac{12n\tilde{\theta}_{(1)}^2(\lambda)}{a_1} \text{tr}(\mathbf{M}_0^{-3} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}),
\end{aligned}$$

and

$$\begin{aligned}
MC_p^{(2)}(\lambda) &= c_M \text{tr}(\mathbf{W}_{\tilde{\theta}_{(2)}(\lambda)} \mathbf{S}^{-1}) + 2p \text{tr}(\mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-1} \mathbf{M}_0) \\
&\quad + \frac{4n\tilde{\theta}_{(1)}(\lambda)}{a_1} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) + \frac{36\tilde{\theta}_{(1)}^2(\lambda)a_2}{a_1^2} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) \\
&\quad - \frac{16\tilde{\theta}_{(1)}^2(\lambda)b_2}{a_1 b_1} \text{tr}(\mathbf{M}_0^{-2} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}) - \frac{12n\tilde{\theta}_{(1)}^2(\lambda)}{a_1} \text{tr}(\mathbf{M}_0^{-3} \mathbf{M}_{\tilde{\theta}_{(2)}(\lambda)}^{-2} \mathbf{V}).
\end{aligned}$$

A.7. Proof of THEOREM 5. In this subsection, we show an asymptotic expansion of the PMSE $[\hat{\mathbf{Y}}_{\tilde{\theta}(\lambda)}]$ for obtaining $\hat{\lambda}_0$. Since the PMSE $[\hat{\mathbf{Y}}_{\tilde{\theta}(\lambda)}]$ is obtained as (9), we consider expanding each term for obtaining $\hat{\lambda}_0$. We obtain

$$\begin{aligned}
\mathbf{W}_{\hat{\theta}(\lambda)} \boldsymbol{\Sigma}^{-1} &= \mathbf{Y}'(\mathbf{I}_n - \mathbf{1}_n \mathbf{1}_n' / n - \mathbf{X} \mathbf{M}_{\hat{\theta}(\lambda)}^{-1} \mathbf{X}')^2 \mathbf{Y} \boldsymbol{\Sigma}^{-1} \\
&= (n - k - 1) \mathbf{S} \boldsymbol{\Sigma}^{-1} \\
&\quad + \mathbf{Y}' \mathbf{X} (\mathbf{M}_0^{-1} - 2 \mathbf{M}_{\hat{\theta}(\lambda)}^{-1} + \mathbf{M}_{\hat{\theta}(\lambda)}^{-1} \mathbf{M}_0 \mathbf{M}_{\hat{\theta}(\lambda)}^{-1}) \mathbf{X}' \mathbf{Y} \boldsymbol{\Sigma}^{-1} \\
&= (n - k - 1) \mathbf{S} \boldsymbol{\Sigma}^{-1} \\
&\quad + \mathbf{Y}' \mathbf{X} \mathbf{Q} \mathbf{D}^{-1/2} (\mathbf{I}_k - \mathbf{D} \{ \mathbf{D} + \hat{\theta}(\lambda) \mathbf{I}_k \}^{-1})^2 \mathbf{D}^{-1/2} \mathbf{Q}' \mathbf{X}' \mathbf{Y} \boldsymbol{\Sigma}^{-1},
\end{aligned}$$

because $\mathbf{Y}'(\mathbf{I}_n - \mathbf{1}_n \mathbf{1}_n' / n - \mathbf{X} \mathbf{M}_0^{-1} \mathbf{X}') \mathbf{Y} = (n - k - 1) \mathbf{S}$, $\mathbf{X}' \mathbf{1}_n = \mathbf{0}_k$ and $\mathbf{Q}' \mathbf{X}' \mathbf{X} \mathbf{Q} = \mathbf{D}$. Hence, we obtain

$$\text{tr}(\mathbf{W}_{\hat{\theta}(\lambda)} \boldsymbol{\Sigma}^{-1}) = (n - k - 1) \text{tr}(\mathbf{S} \boldsymbol{\Sigma}^{-1}) + \theta^2 \text{tr}\{(\mathbf{D} + \theta \mathbf{I}_k)^{-2} \mathbf{D}^{-1} \mathbf{Q}' \mathbf{V}^* \mathbf{Q}\}.$$

Since \mathbf{S} is an unbiased estimator of $\boldsymbol{\Sigma}$, we have

$$E_{\mathbf{Y}}[\text{tr}(\mathbf{W}_{\hat{\theta}(\lambda)} \boldsymbol{\Sigma}^{-1})] = (n - k - 1)p + E_{\mathbf{Y}} \left[\sum_{j=1}^k \left(\frac{\hat{\theta}(\lambda)}{d_j + \hat{\theta}(\lambda)} \right)^2 \frac{(\mathbf{Q}' \mathbf{V}^* \mathbf{Q})_{jj}}{d_j} \right].$$

Then, since $d_i = O(n)$ and $\mathbf{V}^* = O_p(n^2)$, we can expand the above equation as follows:

$$\begin{aligned}
E_{\mathbf{Y}}[\text{tr}(\mathbf{W}_{\hat{\theta}(\lambda)} \boldsymbol{\Sigma}^{-1})] &= (n - k - 1)p + E_{\mathbf{Y}} \left[\sum_{j=1}^k \frac{\hat{\theta}^2(\lambda)}{d_j^3} (\mathbf{Q}' \mathbf{V}^* \mathbf{Q})_{jj} + O_p(n^{-2}) \right] \\
&= (n - k - 1)p + E_{\mathbf{Y}} \left[\frac{a_1^* \hat{\theta}^2(\lambda)}{n} + O_p(n^{-2}) \right].
\end{aligned}$$

From simple calculation by Taylor expansion and noting that $a_j = O_p(1)$ and $b_j = O(1)$, we derive

$$\begin{aligned}
\text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-1} \mathbf{M}_0) &= k - \frac{b_1 \hat{\theta}(\lambda)}{n} + O_p(n^{-2}), \\
\lambda \hat{\theta}(\lambda) \text{tr}(\mathbf{V} \mathbf{M}_{\hat{\theta}(\lambda)}^{-5} \mathbf{M}_0) &= \frac{\lambda \hat{\theta}(\lambda) a_2}{n^2} + O_p(n^{-3}), \\
\lambda \text{tr}(\mathbf{V} \mathbf{M}_{\hat{\theta}(\lambda)}^{-3}) &= \frac{\lambda a_1}{n} + O_p(n^{-2}), \\
\lambda \theta \text{tr}(\mathbf{V} \mathbf{M}_{\hat{\theta}(\lambda)}^{-4}) &= \frac{\lambda \hat{\theta}(\lambda) a_2}{n^2} + O_p(n^{-3}), \\
\text{tr}(\mathbf{M}_{\hat{\theta}(\lambda)}^{-3} \mathbf{M}_0) &= \frac{b_2}{n^2} + O_p(n^{-3}).
\end{aligned}$$

By substituting these results into the $\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}]$ in (9), we obtain the asymptotic expansion of the $\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}]$, as follows:

$$\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}] = (n+k+1)p + E_Y \left[\frac{a_1^* \hat{\theta}^2(\lambda)}{n} - \frac{2pb_1 \hat{\theta}(\lambda)}{n} + \frac{4a_2 \hat{\theta}(\lambda)}{na_1} + O_p(n^{-2}) \right].$$

From (6), which is proved in Appendix A.2, we have $\hat{\theta}(\lambda) = pb_1/(\lambda a_1) + O_p(n^{-1})$. Hence, we consider minimizing the following approximated PMSE:

$$\text{PMSE}[\hat{\mathbf{Y}}_{\hat{\theta}(\lambda)}] = (n+k+1)p + E_Y \left[\frac{pb_1}{na_1} \left(\frac{pb_1 a_1^*}{\lambda^2 a_1} - \frac{2pb_1}{\lambda} + \frac{4a_2}{\lambda a_1} \right) \right] + O(n^{-2}).$$

Hence we obtain the asymptotic optimal λ^* , which minimizes the second term of the above equation as in THEOREM 5.

Acknowledgement

The author wishes to express my deepest gratitude to Dr. Hirokazu Yanagihara of Hiroshima University for his valuable advice and encouragements and introducing me to various fields of mathematical statistics. In addition, I would like to thank Professor Hirofumi Wakaki of Hiroshima University for previous comments and Professor Yasunori Fujikoshi of Hiroshima University for providing helpful comments and suggestions. Also, I thank to Dr. Kengo Kato of Hiroshima University for his encouragements and comments. Furthermore, I sincere thank Professor Megu Ohtaki, Dr. Kenichi Satoh, and Dr. Tetsuji Tonda of Hiroshima University for their encouragements and valuable comments and introducing me to various fields of applied statistics. I also would like to thank Dr. Tomoyuki Akita, Dr. Shinobu Fujii, and Dr. Kunihiko Taniguchi of Hiroshima University for their advice and encouragements. Special thanks are due to Professor Takashi Kanda of Hiroshima Institute of Technology for his encouragements and his recommendation to attend Hiroshima University.

References

- [1] A. C. Atkinson, A note on the generalized information criterion for choice of a model, *Biometrika*, **67** (1980), 413–418.
- [2] P. J. Brown and J. V. Zidek, Adaptive multivariate ridge regression, *Ann. Statist.*, **8** (1980), 64–74.
- [3] B. Efron, The estimation of prediction error: covariance penalties and cross-validation, *J. Amer. Statist. Assoc.*, **99** (2004), 619–632.
- [4] S. J. V. Dien, S. Iwatani, Y. Usuda and K. Matsui, Theoretical analysis of amino acid-producing *Escherichia coli* using a stoichiometric model and multivariate linear regression, *J. Biosci. Bioeng.*, **102** (2006), 34–40.

- [5] Y. Fujikoshi and K. Satoh, Modified AIC and C_p in multivariate linear regression, *Biometrika*, **84** (1997), 707–716.
- [6] Y. Fujikoshi, H. Yanagihara and H. Wakaki, Bias corrections of some criteria for selecting multivariate linear models in a general nonnormal case, *Amer. J. Math. Management Sci.*, **25** (2005), 221–258.
- [7] Y. Haitovsky, On multivariate ridge regression, *Biometrika*, **74** (1987), 563–570.
- [8] A. E. Hoerl and R. W. Kennard, Ridge regression: biased estimation for nonorthogonal problems, *Technometrics*, **12** (1970), 55–67.
- [9] J. C. Lagarias, J. A. Reeds, M. H. Wright and P. E. Wright, Convergence properties of the Nelder-Mead simplex method in low dimensions, *SIAM J. Optim.*, **9** (1998), 112–147.
- [10] J. F. Lawless, Mean squared error properties of generalized ridge estimators, *J. Amer. Statist. Assoc.*, **76** (1981), 462–466.
- [11] C. L. Mallows, Some comments on C_p , *Technometrics*, **15** (1973), 661–675.
- [12] C. L. Mallows, More comments on C_p , *Technometrics*, **37** (1995), 362–372.
- [13] C. Sárbu, C. Onisor, M. Posa, S. Kevresan and K. Kuhajda, Modeling and prediction (correction) of partition coefficients of bile acids and their derivatives by multivariate regression methods, *Talanta*, **75** (2008), 651–657.
- [14] R. Saxén and J. Sundell, ^{137}Cs in freshwater fish in Finland since 1986—a statistical analysis with multivariate linear regression models, *J. Environ. Radioactiv.*, **87** (2006), 62–76.
- [15] X. Shen and J. Ye, Adaptive model selection, *J. Amer. Statist. Assoc.*, **97** (2002), 210–221.
- [16] M. Siotani, T. Hayakawa and Y. Fujikoshi, *Modern Multivariate Statistical Analysis: A Graduate Course and Handbook*, American Sciences Press, Columbus, Ohio, 1985.
- [17] B. Skagerberg, J. MacGregor and C. Kiparissides, Multivariate data analysis applied to low-density polyethylene reactors, *Chemometr. Intell. Lab. Syst.*, **14** (1992), 341–356.
- [18] R. S. Sparks, D. Coutsourides and L. Troskie, The multivariate C_p , *Comm. Statist. A—Theory Methods*, **12** (1983), 1775–1793.
- [19] M. S. Srivastava, *Methods of Multivariate Statistics*, John Wiley & Sons, New York, 2002.
- [20] N. H. Timm, *Applied Multivariate Analysis*, Springer-Verlag, New York, 2002.
- [21] H. Yanagihara and K. Satoh, An unbiased C_p criterion for multivariate ridge regression, *J. Multivariate Anal.*, **101** (2010), 1226–1238.
- [22] J. Ye, On measuring and correcting the effects of data mining and model selection, *J. Amer. Statist. Assoc.*, **93** (1998), 120–131.
- [23] A. Yoshimoto, H. Yanagihara and Y. Ninomiya, Finding factors affecting a forest stand growth through multivariate linear modeling, *J. Jpn. For. Res.*, **87** (2005), 504–512 (in Japanese).

Isamu Nagai
Department of Mathematics
Graduate School of Science
Hiroshima University
Higashi-Hiroshima 739-8526, Japan
E-mail: inagai@hiroshima-u.ac.jp