

CIRJE-F-709

**Selection of Variables in Multivariate Regression
Models for Large Dimensions**

Muni S. Srivastava
University of Toronto

Tatsuya Kubokawa
University of Tokyo

January 2010

CIRJE Discussion Papers can be downloaded without charge from:

<http://www.e.u-tokyo.ac.jp/cirje/research/03research02dp.html>

Discussion Papers are a series of manuscripts in their draft form. They are not intended for circulation or distribution except as indicated by the author. For that reason Discussion Papers may not be reproduced or distributed without the written consent of the author.

Selection of Variables in Multivariate Regression Models for Large Dimensions

Muni S. Srivastava* and Tatsuya Kubokawa†
University of Toronto and University of Tokyo

January 13, 2010

Abstract

The Akaike information criterion, AIC, and Mallows' C_p statistic have been proposed for selecting a smaller number of regressor variables in the multivariate regression models with fully unknown covariance matrix. All these criteria are, however, based on the implicit assumption that the sample size is substantially larger than the dimension of the covariance matrix. To obtain a stable estimator of the covariance matrix, it is required that the dimension of the covariance matrix be much smaller than the sample size. When the dimension is close to the sample size, it is necessary to use ridge type of estimators for the covariance matrix. In this paper, we use a ridge type of estimators for the covariance matrix and obtain the modified AIC and modified C_p statistic under the asymptotic theory that both the sample size and the dimension go to infinity. It is numerically shown that these modified procedures perform very well in the sense of selecting the true model in large dimensional cases.

Key words and phrases: Akaike information criterion, Mallows' C_p , large dimension, multivariate linear regression model, selection of variables.

1 Introduction

Consider a multivariate linear regression model in which p response variables y_1, \dots, y_p are regressed on k explanatory variables $x_{(1)}, \dots, x_{(K)}$, when n observations are available on y_1, \dots, y_p and $x_{(1)}, \dots, x_{(K)}$. Let \mathbf{Y} denotes the $n \times p$ observation matrix on the response variable, and $\widetilde{\mathbf{X}}$ denotes the $n \times K$ observation matrix on the K explanatory variables. Then the multivariate regression model is given by

$$\text{Full Model : } \mathbf{Y} = \widetilde{\mathbf{X}}\boldsymbol{\beta}_F + \mathbf{E}, \quad (1)$$

*Department of Statistics, University of Toronto, 100 St George Street, Toronto, Ontario, CANADA M5S 3G3, E-Mail: srivasta@utstat.toronto.edu

†Faculty of Economics, University of Tokyo, Hongo, Bunkyo-ku, Tokyo 113-0033, JAPAN, E-Mail: tatsuya@e.u-tokyo.ac.jp

where the n rows of \mathbf{E} are independent and identically distributed (iid) as multivariate normal with mean vector zero and the $p \times p$ covariance matrix $\mathbf{\Sigma}$, that is, $\mathbf{e}_i \sim \mathcal{N}_p(\mathbf{0}, \mathbf{\Sigma})$, where $\mathbf{E} = (\mathbf{e}_1, \dots, \mathbf{e}_n)'$ and $\mathbf{\beta}_F$ is a $K \times p$ matrix of unknown parameters. The $n \times p$ matrix \mathbf{Y} is given by $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)'$ and the $n \times K$ matrix $\widetilde{\mathbf{X}}$ is given by $\widetilde{\mathbf{X}} = (\mathbf{x}_1, \dots, \mathbf{x}_n)' = (\mathbf{x}_{(1)}, \dots, \mathbf{x}_{(K)})$ where \mathbf{y}_i 's are random p -vector and \mathbf{x}_i 's and $\mathbf{x}_{(i)}$ are, respectively, K and n -vectors considered known or fixed.

In this paper, we consider the model (1) as a full model and we want to address the problem of selecting the explanatory variables $x_{(1)}, \dots, x_{(K)}$ when n and p are large. When k variables $x_{(\gamma_1)}, \dots, x_{(\gamma_k)}$ are selected from $\{x_{(1)}, \dots, x_{(K)}\}$, the candidate model is written as

$$\text{Candidate Model : } \quad \mathbf{Y} = \mathbf{X}\mathbf{\beta} + \mathbf{E}, \quad (2)$$

where $\mathbf{X} = (\mathbf{x}_{(\gamma_1)}, \dots, \mathbf{x}_{(\gamma_k)})$, and $\mathbf{\beta}$ is a $k \times p$ matrix of unknown parameters. For simplicity, we hereafter write $\mathbf{X} = (\mathbf{x}_{(1)}, \dots, \mathbf{x}_{(k)})$ without any loss of generality. The above model is written as

$$\mathbf{Y} \sim \mathcal{N}_{n,p}(\mathbf{X}\mathbf{\beta}, \mathbf{I}_n, \mathbf{\Sigma}). \quad (3)$$

The *Akaike Information Criterion* (AIC) proposed by Akaike (1973, 1974) is recognized as a useful tool for selecting variables in linear regression models. For obtaining an expression for the AIC, we shall assume that the model given in (2) is an overspecified model, and the true model is given by

$$\text{True Model : } \quad \mathbf{Y} \sim \mathcal{N}_{n,p}(\mathbf{X}\mathbf{\beta}^*, \mathbf{I}_n, \mathbf{\Sigma}^*). \quad (4)$$

It will be assumed that the true model belongs to the overspecified model (2). Let $f(\mathbf{Y}; \mathbf{X}\mathbf{\beta}^*, \mathbf{\Sigma}^*)$ denote the pdf of the true model, namely,

$$f(\mathbf{Y}|\mathbf{X}\mathbf{\beta}^*, \mathbf{\Sigma}^*) = (2\pi)^{-pn/2} |\mathbf{\Sigma}^*|^{-n/2} \text{etr} \left[-\frac{1}{2} \mathbf{\Sigma}^{*-1} (\mathbf{Y} - \mathbf{X}\mathbf{\beta}^*)' (\mathbf{Y} - \mathbf{X}\mathbf{\beta}^*) \right].$$

Let $\widehat{\mathbf{\beta}}(\mathbf{Y})$ and $\widehat{\mathbf{\Sigma}}(\mathbf{Y})$ be estimators of $\mathbf{\beta}$ and $\mathbf{\Sigma}$ based on the candidate model. When the true model is predicted based on the candidate model, the prediction error relative to the the Kullback-Leibler information is given by

$$R_{KL}(\mathbf{\beta}, \mathbf{\Sigma}; \widehat{\mathbf{\beta}}, \widehat{\mathbf{\Sigma}}) = E_{\mathbf{Y}}^* [E_{\mathbf{Z}}^* [\log \{f(\mathbf{Z}|\mathbf{X}\mathbf{\beta}^*, \mathbf{\Sigma}^*) / f(\mathbf{Z}|\mathbf{X}\widehat{\mathbf{\beta}}(\mathbf{Y}), \widehat{\mathbf{\Sigma}}(\mathbf{Y}))\}]], \quad (5)$$

where \mathbf{Y} and \mathbf{Z} are independently distributed but having the same distribution as $f(\mathbf{Y}|\mathbf{X}\mathbf{\beta}^*, \mathbf{\Sigma}^*)$ and $f(\mathbf{Z}|\mathbf{X}\mathbf{\beta}^*, \mathbf{\Sigma}^*)$. Let us define the *Akaike Information* (AI) by

$$AI = -2E_{\mathbf{Y}}^* \left[E_{\mathbf{Z}}^* [\log f(\mathbf{Z}|\mathbf{X}\widehat{\mathbf{\beta}}(\mathbf{Y}), \widehat{\mathbf{\Sigma}}(\mathbf{Y}))] \right], \quad (6)$$

which is a model-related part of the prediction error $R_{KL}(\mathbf{\beta}, \mathbf{\Sigma}; \widehat{\mathbf{\beta}}, \widehat{\mathbf{\Sigma}})$. Then the AIC is generally defined as an asymptotically unbiased estimator of AI , where $\mathbf{\beta}$ and $\mathbf{\Sigma}$ are estimated by the maximum likelihood estimators (MLE), given by

$$\begin{aligned} \widehat{\mathbf{\beta}} &= (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}, \\ \widehat{\mathbf{\Sigma}}_0 &= \mathbf{S}/n = (\mathbf{Y} - \mathbf{X}\widehat{\mathbf{\beta}})(\mathbf{Y} - \mathbf{X}\widehat{\mathbf{\beta}})' / n, \end{aligned}$$

where $\mathbf{S} = (\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})(\mathbf{Y} - \mathbf{X}\hat{\boldsymbol{\beta}})'$. For more accounts on the AIC and the related selection criteria, see Sugiura (1978) and Konishi and Kitagawa (2007).

The AIC and the modified criterion in the multivariate linear regression model were derived by Fujikoshi and Satoh (1997) when $n \rightarrow \infty$ and p is bounded. In the large dimensional case, namely the case that $p \rightarrow \infty$, the MLE $\hat{\boldsymbol{\Sigma}}_0$ and the inverse matrix $\hat{\boldsymbol{\Sigma}}_0^{-1}$ must be instable or nonexistent, which means that the AIC based on the MLE $\hat{\boldsymbol{\Sigma}}_0$ is not appropriate. Srivastava and Kubokawa (2008) considered the ridge type estimator $\hat{\boldsymbol{\Sigma}}_\lambda = (\mathbf{S} + \hat{\lambda}\mathbf{I}_p)/n$ for a function $\hat{\lambda} = \hat{\lambda}(\mathbf{S})$ instead of the MLE, and derived the AIC when $p > n$, $p \rightarrow \infty$ and n is bounded based on the theory given in Srivastava (2007). Recently, Yamamura, Yanagihara and Srivastava (2009) obtained the AIC when $p > n$, $n - k = O(p^\delta)$ and $(n, p) \rightarrow \infty$ for $0 < \delta < 1/3$.

In this paper, we consider the case that

$$\nu_k \equiv n - k - p - 3 > 0 \text{ and } (n, p) \rightarrow \infty \text{ such that } p/n = c \text{ for } 0 < c < 1, \quad (7)$$

where the condition of $n - k - p - 3 > 0$ is required for the existence of the moment $E[\text{tr}[\mathbf{S}^{-2}]]$. Since $n - k > p + 1$, there exists the inverse matrix of the MLE $\hat{\boldsymbol{\Sigma}}_0$. Thus, the AIC based on the MLE are available, but not appropriate for large dimension p , because the MLE is very unstable when p is large, see, e.g., Johnston (2001). In this case, the ridge-type estimator $\hat{\boldsymbol{\Sigma}}_\lambda$ should be used instead of the MLE, and we obtain the AIC based on $\hat{\boldsymbol{\Sigma}}_\lambda$.

When a squared error loss function is employed instead of the Kullback-Leibler information, the prediction error is given by

$$R_{PE}(\boldsymbol{\beta}, \boldsymbol{\Sigma}; \hat{\boldsymbol{\beta}}) = E_{\mathbf{Y}}^*[E_{\mathbf{Z}}^*[\text{tr}[(\mathbf{Z} - \mathbf{X}\hat{\boldsymbol{\beta}}(\mathbf{Y}))\boldsymbol{\Sigma}^{-1}(\mathbf{Z} - \mathbf{X}\hat{\boldsymbol{\beta}}(\mathbf{Y}))']]]. \quad (8)$$

Corresponding to the derivation of the AIC, we can suggest an unbiased estimator of $R_{PE}(\boldsymbol{\beta}, \boldsymbol{\Sigma}; \hat{\boldsymbol{\beta}})$ for the model selection. The unbiased estimator is related to the *Mallows' C_p statistic* proposed by Mallows (1973), and we here call it the C_p -type statistic. In this paper, we also obtain the C_p -type statistic based on the ridge-type estimator $\hat{\boldsymbol{\Sigma}}_\lambda$.

The AIC and C_p -type statistics based on the ridge-type estimator $\hat{\boldsymbol{\Sigma}}_\lambda$ are given in Section 2. We also propose the double ridge AIC and C_p -type statistics based on $\hat{\boldsymbol{\Sigma}}_\lambda$ and the ridge regression estimator of $\boldsymbol{\beta}$, which can be expected to work well in the multicollinearity case of \mathbf{X} . The proofs of their derivation is given in Section 3. A simulation experiment is carried out in Section 4 to compare the AIC and C_p criteria for different value of λ including $\lambda = 0$, and it is shown that the usual AIC and C_p based on the MLE do not work in the large dimensional case, but the AIC and C_p statistics based on the ridge-type estimator perform very well in all the cases. We conclude in Section 5.

2 Ridge-type variable selection procedures

2.1 Ridge-type AIC

In this section, we derive ridge-type AIC and C_p statistic based on the ridge-type estimator of Σ . The ridge-type estimator we want to propose for Σ is

$$\widehat{\Sigma}_\lambda = n^{-1}(\mathbf{S} + \widehat{\lambda}\mathbf{I}_p), \quad (9)$$

where

$$\widehat{\lambda} = c_n(\text{tr } \mathbf{S}/np), \quad \text{for } c_n = O(n^{-\delta}), \delta \geq 0. \quad (10)$$

Let us define the Akaike information based on the ridge-type estimator by

$$AI_\lambda = -2E_{\mathbf{Y}}^* \left[E_{\mathbf{Z}}^* [\log f(\mathbf{Z} | \mathbf{X}\widehat{\beta}(\mathbf{Y}), \widehat{\Sigma}_\lambda(\mathbf{Y}))] \right].$$

The Akaike information criterion is an asymptotically unbiased estimator of AI_λ based on $-2\log f(\mathbf{Y}; \mathbf{X}\widehat{\beta}, \widehat{\Sigma}_\lambda)$, where the bias is given by

$$\Delta_\lambda = \Delta_\lambda(\beta^*, \Sigma^*, \widehat{\beta}, \widehat{\Sigma}_\lambda) = AI_\lambda - E_{\mathbf{Y}}^* [-2\log f(\mathbf{Y} | \mathbf{X}\widehat{\beta}, \widehat{\Sigma}_\lambda)]. \quad (11)$$

When Δ_λ is estimated by Δ_λ^* , the AIC is provided by

$$AIC_\lambda = -2\log f(\mathbf{Y} | \mathbf{X}\widehat{\beta}, \widehat{\Sigma}_\lambda) + \Delta_\lambda^*, \quad (12)$$

We shall assume that $\lim_{p \rightarrow \infty} \text{tr } \Sigma/p \in (0, \infty)$. Under this assumption, it follows from Srivastava (2005) that $\text{tr } \mathbf{S}/np \rightarrow \text{tr } \Sigma/p$ in probability as $(n, p) \rightarrow \infty$. We shall consider the case when $c_n = n/p$, other choices of c_n can also be considered. We obtain an asymptotic expression for the bias in calculating the AIC under (7) and the assumption

$$\lim_{p \rightarrow \infty} \text{tr } [\Sigma]/p < \infty. \quad (13)$$

Theorem 2.1 *Assume the conditions (7) and (13). Then, Δ_λ given in (11) is approximated as*

$$\begin{aligned} \Delta_\lambda &= \frac{np(p+1+2k)}{n-k-p-1} \\ &+ \frac{c_n(n-k)}{p(n-k-p-1)} \left\{ \frac{(n+k)(n-k)}{(n-p)^2} - 1 \right\} \text{tr } [\Sigma^*] \text{tr } [\Sigma^{*-1}] + O(n^{-\delta}). \end{aligned} \quad (14)$$

The unknown quantity $\text{tr } [\Sigma^*] \text{tr } [\Sigma^{*-1}]$ is estimated based on the equality

$$E_{\mathbf{Y}}^* [\widehat{\lambda} \text{tr } [\widehat{\Sigma}_\lambda^{-1}]] = \frac{c_n(n-k)}{p(n-k-p-1)} \text{tr } [\Sigma^*] \text{tr } [\Sigma^{*-1}] + O(n^{-\delta}). \quad (15)$$

Combining these approximations yields AIC_λ given by

$$\begin{aligned} AIC_\lambda &= np \log 2\pi + n \log |\widehat{\Sigma}_\lambda| + \text{tr } [\widehat{\Sigma}_\lambda^{-1} \mathbf{S}] \\ &+ \frac{np(p+1+2k)}{n-k-p-1} + \left\{ \frac{(n+k)(n-k)}{(n-p)^2} - 1 \right\} \widehat{\lambda} \text{tr } [\widehat{\Sigma}_\lambda^{-1}], \end{aligned} \quad (16)$$

From Theorem 2.1, the AIC based on the MLE $\widehat{\boldsymbol{\Sigma}}_0$ is derived by putting $\widehat{\lambda} = 0$ in (16) as

$$AIC_0 = np \log 2\pi + n \log |\widehat{\boldsymbol{\Sigma}}_0| + np + \frac{np(p+1+2k)}{n-k-p-1}. \quad (17)$$

2.2 Ridge-type C_p

As explained above, the AIC is an asymptotically unbiased estimator of a part of the prediction error based on the Kullback-Leiber information. We here use the squared error loss function instead of the Kullback-Leibler information, and consider to estimate the prediction error given by

$$R_{PE}(\boldsymbol{\beta}, \boldsymbol{\Sigma}; \widehat{\boldsymbol{\beta}}) = E_{\mathbf{Y}}^* [E_{\mathbf{Z}}^* [\text{tr} [(\mathbf{Z} - \mathbf{X}\widehat{\boldsymbol{\beta}})\boldsymbol{\Sigma}^{-1}(\mathbf{Z} - \mathbf{X}\widehat{\boldsymbol{\beta}})']]].$$

This is rewritten as $R_{PE}(\boldsymbol{\beta}, \boldsymbol{\Sigma}; \widehat{\boldsymbol{\beta}}) = np + PE$, where

$$PE = E_{\mathbf{Y}}^* [\text{tr} [\boldsymbol{\Sigma}^{-1}(\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})' \mathbf{X}' \mathbf{X} (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta})]]. \quad (18)$$

Since an unbiased estimator of PE is related to the C_p statistic, we here call it the C_p -type statistic. According to the same arguments as in the derivation of the Mallows' C_p statistic, we estimate the covariance matrix $\boldsymbol{\Sigma}$ based on the full model (1). Let $\widetilde{\mathbf{S}} = n^{-1}(\mathbf{Y} - \widetilde{\mathbf{X}}\widetilde{\boldsymbol{\beta}})'(\mathbf{Y} - \widetilde{\mathbf{X}}\widetilde{\boldsymbol{\beta}})$ for $\widetilde{\boldsymbol{\beta}} = (\widetilde{\mathbf{X}}'\widetilde{\mathbf{X}})^{-1}\widetilde{\mathbf{X}}'\mathbf{Y}$. Then $\boldsymbol{\Sigma}$ is estimated by

$$\widetilde{\boldsymbol{\Sigma}}_{\lambda} = n^{-1}(\widetilde{\mathbf{S}} + \widetilde{\lambda}\mathbf{I}_p),$$

where

$$\widetilde{\lambda} = c_n(\text{tr} \widetilde{\mathbf{S}}/np), \quad \text{for } c_n = O(n^{-\delta}), \delta \geq 0.$$

When PE is estimated based on the statistic $\text{tr} [\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})]$, the bias is $\Delta_{PE} = PE - E_{\mathbf{Y}}^* [\text{tr} [\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})]]$. Then, the C_p statistic based on the ridge-type estimator $\widetilde{\boldsymbol{\Sigma}}_{\lambda}$ is given by

$$C_{p,\lambda} = \text{tr} [\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})] + \Delta_{PE}^*,$$

where Δ_{PE}^* is an estimator of Δ_{PE} .

Theorem 2.2 *Assume the conditions (7) and (13). Then, Δ_{PE} is evaluated as*

$$\Delta_{PE} = \frac{np(n-k-p-1)}{n-K-p-1} - \frac{c_n(n-K)}{p(n-K-p-1)} \text{tr} [\boldsymbol{\Sigma}^*] \text{tr} [\boldsymbol{\Sigma}^{*-1}] + O(n^{-\delta}). \quad (19)$$

Let us define C_{λ} by

$$C_{\lambda} = \text{tr} [\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1} \mathbf{S}] - \frac{np(n-k-p-1)}{n-K-p-1} + pk + \widetilde{\lambda} \text{tr} [\widetilde{\boldsymbol{\Sigma}}^{-1}]. \quad (20)$$

Then, C_{λ} is an asymptotically unbiased estimator of PE given in (18), namely, $E[C_{\lambda}] = PE + O(n^{-\delta})$.

When $\widetilde{\lambda} = 0$, from Theorem 2.2, we get Mallows' C_p statistic based on the MLE, given by

$$C_0 = n \text{tr} [\widetilde{\mathbf{S}}^{-1} \mathbf{S}] - \frac{np(n-k-p-1)}{n-K-p-1} + pk. \quad (21)$$

2.3 An extension to double ridge criteria for selection

We are often faced with the multicollinearity cases, where the variables $\mathbf{x}_1, \dots, \mathbf{x}_K$ are highly correlated for $\widetilde{\mathbf{X}} = (\mathbf{x}_1, \dots, \mathbf{x}_K)$. In this case, the inverse matrix of $\mathbf{X}'\mathbf{X}$ is not stable, and it is known that the least squares estimator $\widehat{\boldsymbol{\beta}}$ of $\boldsymbol{\beta}$ does not behave well. An alternative procedure is the ridge regression estimator

$$\widehat{\boldsymbol{\beta}}_\tau = (\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\mathbf{X}'\mathbf{Y},$$

where τ is a nonnegative constant. It is certain that $\widehat{\boldsymbol{\beta}}_\tau$ must be stable for an appropriate constant τ , which results in a good predictor based on $\widehat{\boldsymbol{\beta}}_\tau$. However, it may be important how to determine τ . A possible method in the framework of variable selection is that τ and the variable in \mathbf{X} can be chosen based on AIC or C_p . We thus extend the results given in the previous subsections to the criteria based on the ridge regression estimator $\widehat{\boldsymbol{\beta}}_\tau$ instead of $\widehat{\boldsymbol{\beta}}$, which we call here the double ridge criteria.

Let us define the Akaike information based on the double ridge-type estimators by

$$AI_{\lambda,\tau} = -2E_{\mathbf{Y}}^* \left[E_{\mathbf{Z}}^* [\log f(\mathbf{Z}|\mathbf{X}\widehat{\boldsymbol{\beta}}_\tau(\mathbf{Y}), \widehat{\boldsymbol{\Sigma}}_\lambda(\mathbf{Y}))] \right].$$

The Akaike information criterion is an asymptotically unbiased estimator of $AI_{\lambda,\tau}$ based on $-2\log f(\mathbf{Y}; \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau, \widehat{\boldsymbol{\Sigma}}_\lambda)$, where the bias is given by

$$\Delta_{\lambda,\tau} = \Delta_{\lambda,\tau}(\boldsymbol{\beta}^*, \boldsymbol{\Sigma}^*, \widehat{\boldsymbol{\beta}}_\tau, \widehat{\boldsymbol{\Sigma}}_\lambda) = AI_{\lambda,\tau} - E_{\mathbf{Y}}^* [-2\log f(\mathbf{Y}|\mathbf{X}\widehat{\boldsymbol{\beta}}_\tau, \widehat{\boldsymbol{\Sigma}}_\lambda)]. \quad (22)$$

Theorem 2.3 *Assume the conditions (7) and (13). Then, $\Delta_{\lambda,\tau}$ given in (22) is approximated as*

$$\begin{aligned} \Delta_{\lambda,\tau} = & \frac{np\{p+1+k+(1-\tau^2)\rho_\tau\}}{n-k-p-1} \\ & + \frac{c_n(n-k)}{p(n-k-p-1)} \left\{ \frac{\{n+(1-\tau^2)\rho_\tau\}(n-k)}{(n-p)^2} - 1 \right\} \text{tr}[\boldsymbol{\Sigma}^*] \text{tr}[\boldsymbol{\Sigma}^{*-1}] + O(n^{-\delta}), \end{aligned} \quad (23)$$

where $\rho_\tau = \text{tr}[\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}]^2$. The double ridge Akaike information criterion is given by

$$\begin{aligned} AIC_{\lambda,\tau} = & np \log 2\pi + n \log |\widehat{\boldsymbol{\Sigma}}_\lambda| + \text{tr}[\widehat{\boldsymbol{\Sigma}}_\lambda^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)] \\ & + \frac{np\{p+1+k+(1-\tau^2)\rho_\tau\}}{n-k-p-1} + \left\{ \frac{\{n+(1-\tau^2)\rho_\tau\}(n-k)}{(n-p)^2} - 1 \right\} \widehat{\lambda} \text{tr}[\widehat{\boldsymbol{\Sigma}}_\lambda^{-1}], \end{aligned} \quad (24)$$

When τ takes a value in the range of $[0, \tau_0]$ for a fixed τ_0 , the optimal ridge parameter τ and the optimal variables can be simultaneously and numerically selected so as to minimize the double ridge criterion $AIC_{\lambda,\tau}$.

The C_p statistic can be similarly extended to the case of the double ridge criterion. Since the prediction error based on the ridge estimator $\widehat{\boldsymbol{\beta}}_\tau$ is written as

$$\begin{aligned} R_{PE}(\boldsymbol{\beta}, \boldsymbol{\Sigma}; \widehat{\boldsymbol{\beta}}_\tau) &= E_{\mathbf{Y}}^* [E_{\mathbf{Z}}^* [\text{tr} [(\mathbf{Z} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)\boldsymbol{\Sigma}^{-1}(\mathbf{Z} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)']]] \\ &= np + PE_\tau, \end{aligned}$$

where $PE_\tau = E_{\mathbf{Y}}^* [\text{tr} [\boldsymbol{\Sigma}^{-1}(\widehat{\boldsymbol{\beta}}_\tau - \boldsymbol{\beta})'\mathbf{X}'\mathbf{X}(\widehat{\boldsymbol{\beta}}_\tau - \boldsymbol{\beta})]]$. When PE_τ is estimated based on the statistic $\text{tr} [\widetilde{\boldsymbol{\Sigma}}_\lambda^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)]$, the bias is $\Delta_{PE,\tau} = PE_\tau - E_{\mathbf{Y}}^* [\text{tr} [\widetilde{\boldsymbol{\Sigma}}_\lambda^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)]]$. To evaluate $\Delta_{PE,\tau}$, we assume the condition

$$\lim_{p \rightarrow \infty} \boldsymbol{\beta}\boldsymbol{\Sigma}^{-1}\boldsymbol{\beta}'/p < \infty. \quad (25)$$

Theorem 2.4 *Assume the conditions (7), (13) and (25). Then, $\Delta_{PE,\tau}$ is evaluated as*

$$\begin{aligned} \Delta_{PE,\tau} &= p\rho_\tau - \frac{np(n-k-p-1+\tau^2\rho_\tau)}{n-K-p-1} \\ &\quad + \frac{c_n(n-K)}{p(n-K-p-1)} \text{tr} [\boldsymbol{\Sigma}^*] \text{tr} [\boldsymbol{\Sigma}^{*-1}] + O(1). \end{aligned} \quad (26)$$

Let us define $C_{\lambda,\tau}$ by

$$\begin{aligned} C_{\lambda,\tau} &= \text{tr} [\widetilde{\boldsymbol{\Sigma}}_\lambda^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)] \\ &\quad - \frac{np(n-k-p-1+\tau^2\rho_\tau)}{n-K-p-1} + p\rho_\tau + \widetilde{\lambda} \text{tr} [\widetilde{\boldsymbol{\Sigma}}^{-1}]. \end{aligned} \quad (27)$$

Then, $C_{\lambda,\tau}$ is an asymptotically unbiased estimator of PE_τ , namely, $E[C_{\lambda,\tau}] = PE_\tau + O(1)$.

Similarly to $AIC_{\lambda,\tau}$, the optimal ridge parameter τ and the optimal variables can be simultaneously and numerically selected so as to minimize $C_{\lambda,\tau}$ for $0 \leq \tau \leq \tau_0$.

3 Proofs of the main results

3.1 Proofs of Theorems 2.1 and 2.3.

Since Theorem 2.1 is a special case of Theorem 2.3, we here prove Theorem 2.3. For large p we consider the ridge-type estimator of the $p \times p$ covariance matrix $\boldsymbol{\Sigma}$ given in (9). In order to obtain AIC_λ defined in (17), we need to first evaluate Δ_λ under the true model and, if it depends on some of the unknown parameters, we may need to provide an estimated value of Δ_λ . To prove Theorem 2.1, we note that $-2 \log f(\mathbf{Y} | \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau, \widehat{\boldsymbol{\Sigma}}_\lambda)$ is given by

$$-2 \log f(\mathbf{Y} | \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau, \widehat{\boldsymbol{\Sigma}}_\lambda) = np \log(2\pi) + n \log |\widehat{\boldsymbol{\Sigma}}_\lambda| + \text{tr} [\widehat{\boldsymbol{\Sigma}}_\lambda^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)]$$

and $AI_{\lambda,\tau}$ is written as

$$AI_{\lambda,\tau} = E_{\mathbf{Y}}^* [E_{\mathbf{Z}}^* [np \log(2\pi) + n \log |\widehat{\boldsymbol{\Sigma}}_\lambda| + \text{tr} [\widehat{\boldsymbol{\Sigma}}_\lambda^{-1}(\mathbf{Z} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)'(\mathbf{Z} - \mathbf{X}\widehat{\boldsymbol{\beta}}_\tau)]]].$$

Taking the expectation with respect to \mathbf{Z} yields that

$$\begin{aligned} E_{\mathbf{Z}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}(\mathbf{Z} - \mathbf{X}\widehat{\beta}_\tau)'(\mathbf{Z} - \mathbf{X}\widehat{\beta}_\tau)]] \\ = n\text{tr}[\widehat{\Sigma}_\lambda^{-1}\Sigma] + \text{tr}[\widehat{\Sigma}_\lambda^{-1}(\widehat{\beta}_\tau - \beta)' \mathbf{X}' \mathbf{X}(\widehat{\beta}_\tau - \beta)], \end{aligned}$$

so that the bias is written as

$$\begin{aligned} \Delta_{\lambda,\tau} &= AI_{\lambda,\tau} - E_{\mathbf{Y}}^*[-2\log f(\mathbf{Y}|\mathbf{X}\widehat{\beta}_\tau, \widehat{\Sigma}_\lambda)] \\ &= E_{\mathbf{Y}}^*[n\text{tr}[\widehat{\Sigma}_\lambda^{-1}\Sigma] + \text{tr}[\widehat{\Sigma}_\lambda^{-1}(\widehat{\beta}_\tau - \beta)' \mathbf{X}' \mathbf{X}(\widehat{\beta}_\tau - \beta)] \\ &\quad - \text{tr}[\widehat{\Sigma}_\lambda^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\beta}_\tau)'(\mathbf{Y} - \mathbf{X}\widehat{\beta}_\tau)]]. \end{aligned} \quad (28)$$

It is here observed that

$$\begin{aligned} E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\beta}_\tau)'(\mathbf{Y} - \mathbf{X}\widehat{\beta}_\tau)]] \\ = E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}\mathbf{S}]] - 2E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\beta}_\tau)' \mathbf{X}(\widehat{\beta}_\tau - \beta)]] \\ + E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}(\widehat{\beta}_\tau - \widehat{\beta})' \mathbf{X}' \mathbf{X}(\widehat{\beta}_\tau - \widehat{\beta})]], \end{aligned} \quad (29)$$

where $\mathbf{S} = (\mathbf{Y} - \mathbf{X}\widehat{\beta})(\mathbf{Y} - \mathbf{X}\widehat{\beta})'$. Note that $\widehat{\beta}_\tau - \widehat{\beta} = -\tau(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = -\tau(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\widehat{\beta}$ and that $\mathbf{Y} - \mathbf{X}\widehat{\beta}$ is independent of β . Since \mathbf{S} is invariant under the sign change of $\mathbf{Y} - \mathbf{X}\widehat{\beta}$, it can be seen that

$$\begin{aligned} E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\beta}_\tau)' \mathbf{X}(\widehat{\beta}_\tau - \beta)]] \\ = \text{tr}[E_{\mathbf{Y}}^*[\widehat{\Sigma}_\lambda^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\beta}_\tau)'] E_{\mathbf{Y}}^*[\mathbf{X}(\widehat{\beta}_\tau - \beta)]] = 0. \end{aligned} \quad (30)$$

Also, note that $\widehat{\beta}_\tau - \beta \sim \mathcal{N}(-\tau(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\beta, (\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}, \Sigma)$ and $\widehat{\beta} - \beta \sim \mathcal{N}(\mathbf{0}, (\mathbf{X}'\mathbf{X})^{-1}, \Sigma)$. Then,

$$\begin{aligned} E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}(\widehat{\beta}_\tau - \beta)' \mathbf{X}' \mathbf{X}(\widehat{\beta}_\tau - \beta)]] \\ = E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}\Sigma]]\rho_\tau + \tau^2 E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}\beta'(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\beta]], \end{aligned} \quad (31)$$

$$\begin{aligned} E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}(\widehat{\beta}_\tau - \widehat{\beta})' \mathbf{X}' \mathbf{X}(\mathbf{X}\widehat{\beta}_\tau - \widehat{\beta})]] \\ = \tau^2 E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}\widehat{\beta}'(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\widehat{\beta}]] \\ = \tau^2 E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}\Sigma]]\rho_\tau + \tau^2 E_{\mathbf{Y}}^*[\text{tr}[\widehat{\Sigma}_\lambda^{-1}\beta'(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\beta]]. \end{aligned} \quad (32)$$

Combining these observations, from (28) we can express the bias as

$$\Delta_{\lambda,\tau} = E_{\mathbf{Y}}^*[\{n + (1 - \tau^2)\rho_\tau\}\text{tr}[\widehat{\Sigma}_\lambda^{-1}\Sigma] - \text{tr}[\widehat{\Sigma}_\lambda^{-1}\mathbf{S}]]. \quad (33)$$

From the equation given in problem 1.6 (i) of Srivastava and Khatri (1979, pp33), it is noted that

$$(\mathbf{I} + \widehat{\lambda}\mathbf{S}^{-1})^{-1} = \mathbf{I} - \widehat{\lambda}\mathbf{S}^{-1} + \widehat{\lambda}^2\mathbf{S}^{-2}(\mathbf{I} + \widehat{\lambda}\mathbf{S}^{-1})^{-1}. \quad (34)$$

Hence, Δ_λ is rewritten as

$$\begin{aligned}
\Delta_{\lambda,\tau} &= n\{n + (1 - \tau^2)\rho_\tau\}E_{\mathbf{Y}}^*[\text{tr}\{\{\mathbf{I} - \hat{\lambda}\mathbf{S}^{-1} + \hat{\lambda}^2\mathbf{S}^{-2}(\mathbf{I} + \hat{\lambda}\mathbf{S}^{-1})^{-1}\}\mathbf{S}^{-1}\boldsymbol{\Sigma}^*\}] \\
&\quad - nE_{\mathbf{Y}}^*[\text{tr}\{\mathbf{I} - \hat{\lambda}\mathbf{S}^{-1} + \hat{\lambda}^2\mathbf{S}^{-2}(\mathbf{I} + \hat{\lambda}\mathbf{S}^{-1})^{-1}\}] \\
&= nE_{\mathbf{Y}}^*[\{n + (1 - \tau^2)\rho_\tau\}\text{tr}\{\mathbf{S}^{-1}\boldsymbol{\Sigma}^*\} - p] - nE_{\mathbf{Y}}^*[\{n + (1 - \tau^2)\rho_\tau\}\hat{\lambda}\text{tr}\{\mathbf{S}^{-2}\boldsymbol{\Sigma}^*\} - \hat{\lambda}\text{tr}\{\mathbf{S}^{-1}\}] \\
&\quad + nE_{\mathbf{Y}}^*[\{n + (1 - \tau^2)\rho_\tau\}\hat{\lambda}^2\text{tr}\{\mathbf{S}^{-2}(\mathbf{I} + \hat{\lambda}\mathbf{S}^{-1})^{-1}\mathbf{S}^{-1}\boldsymbol{\Sigma}^*\} - \hat{\lambda}^2\text{tr}\{\mathbf{S}^{-2}(\mathbf{I} + \hat{\lambda}\mathbf{S}^{-1})^{-1}\}] \\
&= I_1 - I_2 + I_3. \quad (\text{say})
\end{aligned}$$

We first evaluate I_3 . Since $p/n \rightarrow c$, $0 < c < 1$, it follows from Bai and Yin (1993) that \mathbf{S}/n is almost surely bounded by a constant matrix. Also, we have assumed that $\lim_{p \rightarrow \infty} \text{tr}[\boldsymbol{\Sigma}]/p < \infty$. Hence, it can be seen that

$$\begin{aligned}
&n\{n + (1 - \tau^2)\rho_\tau\}\text{tr}\{\mathbf{S}^{-2}(\mathbf{I} + \hat{\lambda}\mathbf{S}^{-1})^{-1}\mathbf{S}^{-1}\boldsymbol{\Sigma}^*\} \\
&\leq n\{n + (1 - \tau^2)\rho_\tau\}\text{tr}\{\mathbf{S}^{-3}\boldsymbol{\Sigma}^*\} = \frac{n+k}{n} \frac{p}{n} \frac{\text{tr}\{(\mathbf{S}/n)^{-3}\boldsymbol{\Sigma}^*\}}{p} = O_p(1),
\end{aligned}$$

and

$$n\text{tr}\{\mathbf{S}^{-2}(\mathbf{I} + \hat{\lambda}\mathbf{S}^{-1})^{-1}\} \leq n\text{tr}\{\mathbf{S}^{-2}\} = \frac{p}{n} \frac{\text{tr}\{(\mathbf{S}/n)^{-2}\}}{p} = O_p(1). \quad (35)$$

Also, note that $\hat{\lambda} = c_n \text{tr}[\mathbf{S}]/(np) = c_n \text{tr}[\mathbf{S}/n]/p = O_p(n^{-\delta})$ since $c_n = O(n^{-\delta})$, $\delta \geq 0$. These evaluations mean that $I_3 = O(n^{-2\delta})$. Since $E_{\mathbf{Y}}^*[\text{tr}\{\mathbf{S}^{-1}\boldsymbol{\Sigma}^*\}] = p/(n - k - p - 1)$, it is easy to see that

$$I_1 = \frac{np(p+1+2k)}{n-k-p-1},$$

which is of order $O(n^2)$. To estimate I_2 , we can express I_2 as

$$I_2 = \frac{c_n}{p} E_{\mathbf{Y}}^*[\{n + (1 - \tau^2)\rho_\tau\}\text{tr}[\mathbf{S}]\text{tr}[\mathbf{S}^{-2}\boldsymbol{\Sigma}^*] - \text{tr}[\mathbf{S}]\text{tr}[\mathbf{S}^{-1}]].$$

Thus, from Lemmas A.1 and A.2, it follows that

$$\begin{aligned}
I_2 &= \frac{c_n}{p(n-k-p-1)} \left\{ \left[\frac{\{n + (1 - \tau^2)\rho_\tau\}(n-k-1)(n-k+1)}{(n-k-p+1)(n-k-p-3)} - (n-k) \right] \text{tr}[\boldsymbol{\Sigma}^*]\text{tr}[\boldsymbol{\Sigma}^{*-1}] \right. \\
&\quad \left. + 2p - \frac{\{n + (1 - \tau^2)\rho_\tau\}p}{n-k-p-3} \left[\frac{(n-k)^2 - 1}{n-k-p+1} - \frac{(n-k)^2 - 5(n-k) + 2p + 2}{n-k-p} \right] \right\}.
\end{aligned}$$

Since $\{n + (1 - \tau^2)\rho_\tau\}(n-k-1)(n-k+1)/\{(n-k-p+1)(n-k-p-3)\} = \{n + (1 - \tau^2)\rho_\tau\}(n-k)(n-k)/(n-p)^2 + O(1)$, I_2 can be approximated as

$$I_2 = \frac{c_n(n-k)}{p(n-k-p-1)} \left\{ \frac{\{n + (1 - \tau^2)\rho_\tau\}(n-k)}{(n-p)^2} - 1 \right\} \text{tr}[\boldsymbol{\Sigma}^*]\text{tr}[\boldsymbol{\Sigma}^{*-1}] + O(n^{-\delta}).$$

Combining the above evaluations, we get

$$\begin{aligned}
\Delta_{\lambda,\tau} &= \frac{np\{p+1+k+(1-\tau^2)\rho_\tau\}}{n-k-p-1} \\
&\quad + \frac{c_n(n-k)}{p(n-k-p-1)} \left\{ \frac{\{n + (1 - \tau^2)\rho_\tau\}(n-k)}{(n-p)^2} - 1 \right\} \text{tr}[\boldsymbol{\Sigma}^*]\text{tr}[\boldsymbol{\Sigma}^{*-1}] + O(n^{-\delta}).
\end{aligned} \quad (36)$$

Since $\Delta_{\lambda,\tau}$ involves unknown quantity, we need to estimate $\text{tr}[\boldsymbol{\Sigma}^*]\text{tr}[\boldsymbol{\Sigma}^{*-1}]$. Using (34) and (35), we can observe that

$$\begin{aligned} E[\text{tr}[\widehat{\boldsymbol{\Sigma}}_{\lambda}^{-1}\mathbf{S}]] &= np - \frac{c_n}{p} E_{\mathbf{Y}}^*[\text{tr}[\mathbf{S}]\text{tr}[\mathbf{S}^{-1}]] + O(n^{-2\delta}) \\ &= np - \frac{c_n(n-k)}{p(n-k-p-1)} \text{tr}[\boldsymbol{\Sigma}^*]\text{tr}[\boldsymbol{\Sigma}^{*-1}] + O(n^{-\delta}), \end{aligned}$$

where Lemma A.1 is used to show the second equality. Since $np - \text{tr}[\widehat{\boldsymbol{\Sigma}}_{\lambda}^{-1}\mathbf{S}] = \hat{\lambda}\text{tr}[\widehat{\boldsymbol{\Sigma}}_{\lambda}^{-1}]$, it follows that

$$E_{\mathbf{Y}}^*[\hat{\lambda}\text{tr}[\widehat{\boldsymbol{\Sigma}}_{\lambda}^{-1}]] = \frac{c_n(n-k)}{p(n-k-p-1)} \text{tr}[\boldsymbol{\Sigma}^*]\text{tr}[\boldsymbol{\Sigma}^{*-1}] + O(n^{-\delta}), \quad (37)$$

which is substituted into (36) to get the expression

$$\begin{aligned} \Delta_{\lambda,\tau} &= \frac{np\{p+1+k+(1-\tau^2)\rho_{\tau}\}}{n-k-p-1} \\ &\quad + \left\{ \frac{\{n+(1-\tau^2)\rho_{\tau}\}(n-k)}{(n-p)^2} - 1 \right\} E_{\mathbf{Y}}^*[\hat{\lambda}\text{tr}[\widehat{\boldsymbol{\Sigma}}_{\lambda}^{-1}]] + O(n^{-\delta}). \end{aligned} \quad (38)$$

Hence, the approximated value of AIC_{λ} stated in Theorem 2.1 is obtained. \blacksquare

3.2 Proofs of Theorems 2.2 and 2.4

Since Theorem 2.2 is a special case of Theorem 2.4, we here prove Theorem 2.4. Letting $PE_{\tau} = E_{\mathbf{Y}}^*[\text{tr}[\boldsymbol{\Sigma}^{-1}(\widehat{\boldsymbol{\beta}}_{\tau} - \boldsymbol{\beta})'\mathbf{X}'\mathbf{X}(\widehat{\boldsymbol{\beta}}_{\tau} - \boldsymbol{\beta})]]$, we can see that

$$PE_{\tau} = p\rho_{\tau} + \tau^2 \text{tr}[\boldsymbol{\Sigma}^{-1}\boldsymbol{\beta}'(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X} + \tau\mathbf{I})^{-1}\boldsymbol{\beta}].$$

We shall obtain an asymptotic unbiased estimator of PE_{τ} based on $\text{tr}[\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_{\tau})'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_{\tau})]$. The expectation can be evaluated as

$$\begin{aligned} &E_{\mathbf{Y}}^*[\text{tr}[\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_{\tau})'(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_{\tau})]] \\ &= E_{\mathbf{Y}}^*[\text{tr}[\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}\mathbf{S}]] - 2E_{\mathbf{Y}}^*[\text{tr}[\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})'(\mathbf{X}\widehat{\boldsymbol{\beta}}_{\tau} - \widehat{\boldsymbol{\beta}})]] \\ &\quad + E_{\mathbf{Y}}^*[\text{tr}[\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}(\widehat{\boldsymbol{\beta}}_{\tau} - \widehat{\boldsymbol{\beta}})'(\mathbf{X}\widehat{\boldsymbol{\beta}}_{\tau} - \widehat{\boldsymbol{\beta}})]]. \end{aligned}$$

It is noted that $\mathbf{Y} - \widetilde{\mathbf{X}}\widetilde{\boldsymbol{\beta}}$, $\widetilde{\mathbf{X}}\widetilde{\boldsymbol{\beta}} - \mathbf{X}\widehat{\boldsymbol{\beta}}$ and $\widehat{\boldsymbol{\beta}}$ are mutually independent for $\widetilde{\boldsymbol{\beta}} = (\widetilde{\mathbf{X}}'\widetilde{\mathbf{X}})^{-1}\widetilde{\mathbf{X}}'\mathbf{Y}$. Since $\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}} = (\mathbf{Y} - \widetilde{\mathbf{X}}\widetilde{\boldsymbol{\beta}}) + (\widetilde{\mathbf{X}}\widetilde{\boldsymbol{\beta}} - \mathbf{X}\widehat{\boldsymbol{\beta}})$, the product term can be evaluated as

$$\begin{aligned} &E_{\mathbf{Y}}^*[\text{tr}[\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}(\mathbf{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}})'(\mathbf{X}\widehat{\boldsymbol{\beta}}_{\tau} - \widehat{\boldsymbol{\beta}})]] \\ &= E_{\mathbf{Y}}^*[\text{tr}[\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}\{(\mathbf{Y} - \widetilde{\mathbf{X}}\widetilde{\boldsymbol{\beta}}) + (\widetilde{\mathbf{X}}\widetilde{\boldsymbol{\beta}} - \mathbf{X}\widehat{\boldsymbol{\beta}})\}'(\mathbf{X}\widehat{\boldsymbol{\beta}}_{\tau} - \widehat{\boldsymbol{\beta}})]] \\ &= \text{tr}[E_{\mathbf{Y}}^*[\widetilde{\boldsymbol{\Sigma}}_{\lambda}^{-1}(\mathbf{Y} - \widetilde{\mathbf{X}}\widetilde{\boldsymbol{\beta}})']E_{\mathbf{Y}}^*[\mathbf{X}(\widehat{\boldsymbol{\beta}}_{\tau} - \widehat{\boldsymbol{\beta}})]] \\ &= 0, \end{aligned}$$

where the same arguments as in (30) have been used to show the last equality. Similar to (32), it can be observed that

$$\begin{aligned} E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1}(\hat{\beta}_\tau - \hat{\beta})' \mathbf{X}' \mathbf{X}(\hat{\beta}_\tau - \hat{\beta})]] \\ = \tau^2 E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \Sigma]] \rho_\tau + \tau^2 E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \beta' (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1} \mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1} \beta]]. \end{aligned}$$

Hence, the bias can be evaluated as

$$\begin{aligned} \Delta_{PE,\tau} &= PE_\tau - E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} (\mathbf{Y} - \mathbf{X} \hat{\beta}_\tau)' (\mathbf{Y} - \mathbf{X} \hat{\beta}_\tau)]] \\ &= p \rho_\tau + \tau^2 \text{tr}[\Sigma^{-1} \beta' (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1} \mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1} \beta] \\ &\quad - E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \mathbf{S}] - \tau^2 \rho_\tau E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \Sigma]]] \\ &\quad - \tau^2 \text{tr}[E_{\mathbf{Y}}^*[\tilde{\Sigma}_\lambda^{-1}] \beta' (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1} \mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1} \beta]. \end{aligned} \quad (39)$$

Since $\beta \Sigma^{-1} \beta' / p$ is bounded for large p from the condition (25), it is seen that

$$\begin{aligned} \text{tr}[\Sigma^{-1} \beta' (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1} \mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1} \beta] \\ \leq \text{tr}[\beta \Sigma^{-1} \beta' (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1}] = O(1). \end{aligned}$$

Similarly,

$$\begin{aligned} \text{tr}[E_{\mathbf{Y}}^*[\tilde{\Sigma}_\lambda^{-1}] \beta' (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1} \mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1} \beta] \\ \leq \text{tr}[\beta E_{\mathbf{Y}}^*[n \tilde{\mathbf{S}}^{-1}] \beta' (\mathbf{X}' \mathbf{X} + \tau \mathbf{I})^{-1}] = O(1). \end{aligned}$$

Thus,

$$\Delta_{PE,\tau} = p \rho_\tau - E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \mathbf{S}] - \tau^2 \rho_\tau E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \Sigma]]] + O(1). \quad (40)$$

We first evaluate $E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \mathbf{S}]]$. Noting that $\mathbf{S} = \mathbf{Y}'(\mathbf{I} - \mathbf{X}(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}') \mathbf{Y}$ and $\tilde{\mathbf{S}} = \mathbf{Y}'(\mathbf{I} - \tilde{\mathbf{X}}(\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}') \mathbf{Y}$, there exists a $p \times (K - k)$ random matrix \mathbf{U} such that $\mathbf{S} = \tilde{\mathbf{S}} + \mathbf{U} \mathbf{U}'$ and \mathbf{U} is distributed as $\mathbf{U}' \sim \mathcal{N}_{K-k}(\mathbf{0}, \mathbf{I}_N, \Sigma^*)$, independent of $\tilde{\mathbf{S}}$. Since $\tilde{\Sigma}_\lambda$ is a function of $\tilde{\mathbf{S}}$, it is seen that

$$\begin{aligned} E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \mathbf{S}]] &= E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} (\tilde{\mathbf{S}} + \mathbf{U} \mathbf{U}')]] \\ &= E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \tilde{\mathbf{S}}] + (K - k) \text{tr}[\tilde{\Sigma}_\lambda^{-1} \Sigma]] \\ &= E_{\mathbf{Y}}^*[n \text{tr}[(\tilde{\mathbf{S}} + \tilde{\lambda} \mathbf{I})^{-1} \tilde{\mathbf{S}}] + n(K - k) \text{tr}[(\tilde{\mathbf{S}} + \tilde{\lambda} \mathbf{I})^{-1} \Sigma]] \\ &= E_{\mathbf{Y}}^*[n \text{tr}[(\mathbf{I} + \tilde{\lambda} \tilde{\mathbf{S}}^{-1})^{-1}] + n(K - k) \text{tr}[(\mathbf{I} + \tilde{\lambda} \tilde{\mathbf{S}}^{-1})^{-1} \tilde{\mathbf{S}}^{-1} \Sigma]]. \end{aligned}$$

From (34) and the fact that $(\mathbf{I} + \hat{\lambda} \mathbf{S}^{-1})^{-1} = \mathbf{I} - \hat{\lambda} \mathbf{S}^{-1} (\mathbf{I} + \hat{\lambda} \mathbf{S}^{-1})^{-1}$, it follows that

$$\begin{aligned} E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \mathbf{S}]] &= E_{\mathbf{Y}}^*[n \text{tr}[\mathbf{I} - \tilde{\lambda} \tilde{\mathbf{S}}^{-1} + \tilde{\lambda}^2 \tilde{\Sigma}^{-2} (\mathbf{I} + \tilde{\lambda} \tilde{\mathbf{S}}^{-1})^{-1}] \\ &\quad + n(K - k) \text{tr}[\{\mathbf{I} - \tilde{\lambda} \tilde{\Sigma}^{-1} (\mathbf{I} + \tilde{\lambda} \tilde{\mathbf{S}}^{-1})^{-1}\} \tilde{\mathbf{S}}^{-1} \Sigma]]. \end{aligned}$$

Using the arguments as in (35), we can see that

$$\begin{aligned}
E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1} \mathbf{S}]] &= np - E_{\mathbf{Y}}^*[n\tilde{\lambda}\text{tr}[\tilde{\mathbf{S}}^{-1}] + n(K-k)\text{tr}[\tilde{\mathbf{S}}^{-1}\Sigma]] + O(n^{-2\delta}) \\
&= np - \frac{c_n}{p} E_{\mathbf{Y}}^*[\text{tr}[\tilde{\mathbf{S}}]\text{tr}[\tilde{\mathbf{S}}^{-1}]] + n(K-k)\frac{p}{n-K-p-1} + O(n^{-2\delta}) \\
&= \frac{np(n-k-p-1)}{n-K-p-1} - \frac{c_n(n-K)}{p(n-K-p-1)}\text{tr}[\Sigma^*]\text{tr}[\Sigma^{*-1}] + O(n^{-\delta}). \quad (41)
\end{aligned}$$

where Lemma A.1 is used at the third equality.

Using a similar argument, we next evaluate $E_{\mathbf{Y}}^*[\text{tr}[\tilde{\Sigma}_\lambda^{-1}\Sigma]]$. Since $(\mathbf{I} + \tilde{\lambda}\tilde{\mathbf{S}}^{-1})^{-1} = \mathbf{I} - \tilde{\lambda}\tilde{\mathbf{S}}^{-1}(\mathbf{I} + \tilde{\lambda}\tilde{\mathbf{S}}^{-1})^{-1}$, it can be seen that

$$\begin{aligned}
\text{tr}[\tilde{\Sigma}_\lambda^{-1}\Sigma] &= n\text{tr}[\tilde{\mathbf{S}}^{-1}\Sigma] - n\tilde{\lambda}\text{tr}[\tilde{\mathbf{S}}^{-1}\Sigma\tilde{\mathbf{S}}^{-1}(\mathbf{I} + \tilde{\lambda}\tilde{\mathbf{S}}^{-1})^{-1}] \\
&= \frac{np}{n-K-p-1} + O(n^{-\delta}).
\end{aligned}$$

Hence from (40) and (41), we can see that

$$\begin{aligned}
\Delta_{PE,\tau} &= p\rho_\tau - \frac{np(n-k-p-1)}{n-K-p-1} + \frac{c_n(n-K)}{p(n-K-p-1)}\text{tr}[\Sigma^*]\text{tr}[\Sigma^{*-1}] \\
&\quad - \tau^2\rho_\tau\frac{np}{n-K-p-1} + O(1) \\
&= p\rho_\tau - \frac{np(n-k-p-1 + \tau^2\rho_\tau)}{n-K-p-1} + \frac{c_n(n-K)}{p(n-K-p-1)}\text{tr}[\Sigma^*]\text{tr}[\Sigma^{*-1}] + O(1). \quad (42)
\end{aligned}$$

From (15), it follows that

$$E_{\mathbf{Y}}^*[\tilde{\lambda}\text{tr}[\tilde{\Sigma}_\lambda^{-1}]] = \frac{c_n(n-K)}{p(n-K-p-1)}\text{tr}[\Sigma^*]\text{tr}[\Sigma^{*-1}] + O(n^{-\delta}),$$

where K is the rank of $\tilde{\mathbf{X}}$. Hence, we get the $C_{p,\tau}$ type criterion given in (27). \blacksquare

4 Simulation and empirical studies

4.1 Simulation experiments

We now investigate the numerical performances of the ridge-type and double ridge-type AICs and C_p statistics derived in Section 2 through simulation and compare them in terms of the frequencies of selecting the true model.

As the true model, we consider the model that $\mathbf{Y} \sim \mathcal{N}_{n,p}(\tilde{\mathbf{X}}\boldsymbol{\beta}^*, \mathbf{I}_n, \Sigma^*)$, where $\tilde{\mathbf{X}} = (\mathbf{x}_{(1)}, \dots, \mathbf{x}_{(K)})$ is a matrix of regressor variables in a full model given in (1),

$$\boldsymbol{\beta}^* = ((\boldsymbol{\beta}_1^*)', \dots, (\boldsymbol{\beta}_{k^*}^*)', \mathbf{0}, \dots, \mathbf{0})', \quad \beta_{ij}^* = 2(-1)^i(u_{ij} + i), \quad i = 1, \dots, k^*, j = 1, \dots, p,$$

for random variable u_{ij} from a uniform distribution on the interval $[0, 1]$, and

$$\Sigma^* = \begin{pmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_p \end{pmatrix} \begin{pmatrix} \rho^{|1-1|} & \rho^{|1-2|} & \dots & \rho^{|1-p|} \\ \rho^{|2-1|} & \rho^{|2-2|} & \dots & \rho^{|2-p|} \\ \vdots & \vdots & \ddots & \vdots \\ \rho^{|p-1|} & \rho^{|p-2|} & \dots & \rho^{|p-p|} \end{pmatrix} \begin{pmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_p \end{pmatrix}.$$

for a constant ρ on the interval $(-1, 1)$ and $\sigma_i = 2 + (p - i + 1)/p$.

The simulation experiments have been carried out for $n = 76$, $K = 7$, $\rho = 0.7$, $p = 10, 20, 30, 40, 50, 60$. For the $n \times K$ matrix $\widetilde{\mathbf{X}}$ of the regressor variables in the full model (1), the row vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ for $\mathbf{X}' = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ are generated as mutually independent random variables distributed as $\mathcal{N}_k(\mathbf{0}, \Sigma_x)$ where $\Sigma_x = (1 - \rho_x)\mathbf{I}_K + \rho_x\mathbf{J}_K$ for $\rho_x = 0.7$, where $\mathbf{J}_K = \mathbf{j}_K\mathbf{j}'_K$ for $\mathbf{j}_K = (1, \dots, 1)'$, a K -vector of ones. The above true model is expressed as

$$M_{k^*} \quad \mathbf{Y} = \widetilde{\mathbf{X}}\boldsymbol{\beta}^* + \boldsymbol{\epsilon},$$

where $1 \leq k^* \leq 7$, $\boldsymbol{\beta}^* = ((\boldsymbol{\beta}_1)', \dots, (\boldsymbol{\beta}_{k^*})', \mathbf{0}, \dots, \mathbf{0})'$, and $\boldsymbol{\epsilon}$ is a random variable having $\boldsymbol{\epsilon} \sim \mathcal{N}_{n,p}(\widetilde{\mathbf{X}}\boldsymbol{\beta}^*, \mathbf{I}_n, \Sigma^*)$. Let us write the model using the first m regressor variables $\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m$ by M_m . Then, the full model is M_7 and the true model is M_{k^*} . As candidate models, we consider the nested subsets M_1, \dots, M_7 , namely,

$$M_m \quad \mathbf{y} = \widetilde{\mathbf{X}}\boldsymbol{\beta}^{(m)} + \boldsymbol{\epsilon},$$

where $\boldsymbol{\beta}^{(m)} = (\beta_1, \dots, \beta_m, 0, \dots, 0)'$.

In the simulation experiments, 20 observations of the regressor variables $\widetilde{\mathbf{X}}$ are generated, and for each observation of $\widetilde{\mathbf{X}}$, 50 observations of the response variable \mathbf{y} are generated from the true model M_{k^*} for $k^* = 4$. Thus, we have $20 \times 50 (= 1,000)$ total data sets. For each data set, we calculate the values of AIC_0 , AIC_λ , C_0 and C_λ with $c_n = n/p$ given in (17), (16), (21) and (20), respectively, for the seven candidate models M_1, \dots, M_7 , and we select the models minimizing the values of the selection procedures. For each criterion and each candidate model M_m , the number of selecting the model M_m is counted for 1,000 data set. We thus obtain the frequencies of the model M_m selected by the criteria by dividing the number by 1,000.

Table 1 reports the frequencies in the cases of $p = 10, 20, 30, 40, 50, 60$ under the true model M_4 , namely $k^* = 4$. From this table, it is seen that all the criteria perform well for small p in the sense of selecting the true model. For larger p , AIC_0 and C_0 based on the MLE of Σ perform much worse, while AIC_λ and C_λ based on the ridge-type estimator of Σ perform quite well. Table 2 handles the extreme cases of $p = 65$, namely $\nu_K = n - K - p - 3 = 1$ for $k^* = 2, 3, 4, 5, 6, 7$. For the extreme cases reported in this table, AIC_λ and C_λ work still well.

It is interesting to investigate how the double ridge criteria $AIC_{\lambda,\tau}$ and $C_{\lambda,\tau}$ work in multicollinearity cases, where $AIC_{\lambda,\tau}$ and $C_{\lambda,\tau}$ are given in (24) and (27). To clarify the difference between the ridge-type and the double ridge-type criteria, we consider the extreme case of $n = 22$, $K = 7$, $p = 10$ and $\nu_K = n - K - p - 3 = 2$. For the $n \times K$ matrix $\widetilde{\mathbf{X}} = (\mathbf{x}_{(1)}, \dots, \mathbf{x}_{(7)})$ of the regressor variables in the full model (1), it is supposed that $\mathbf{x}_{(3)}$, $\mathbf{x}_{(5)}$ and $\mathbf{x}_{(7)}$ are generated as

Table 1: Frequencies selected by the four criteria AIC_0 , AIC_λ , C_0 and C_λ in 1,000 replications for $n = 76$, $K = 7$, $p = 10, 20, 30, 40, 50, 60$ and $\nu_K = n - K - p - 3$ under the true model M_4 , namely $k^* = 4$

M_k	AIC_0	AIC_λ	C_0	C_λ	AIC_0	AIC_λ	C_0	C_λ	
$p = 10, \nu_K = 56$					$p = 20, \nu_K = 46$				
M_1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_4	99.7	99.6	90.3	99.0	100.0	100.0	91.6	99.2	
M_5	0.3	0.4	7.2	0.9	0.0	0.0	6.7	0.7	
M_6	0.0	0.0	2.0	0.1	0.0	0.0	1.4	0.0	
M_7	0.0	0.0	0.5	0.0	0.0	0.0	0.3	0.1	
$p = 30, \nu_K = 36$					$p = 40, \nu_K = 26$				
M_1	1.8	0.0	0.0	0.0	100.0	0.0	0.0	0.0	
M_2	0.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_4	97.6	100.0	89.4	99.3	0.0	100.0	83.3	99.2	
M_5	0.0	0.0	6.9	0.7	0.0	0.0	10.8	0.8	
M_6	0.0	0.0	2.5	0.0	0.0	0.0	3.4	0.0	
M_7	0.0	0.0	1.2	0.0	0.0	0.0	2.5	0.0	
$p = 50, \nu_K = 16$					$p = 60, \nu_K = 6$				
M_1	100.0	0.0	0.0	0.0	100.0	0.0	0.0	0.0	
M_2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_4	0.0	100.0	69.3	99.0	0.0	100.0	53.1	100.0	
M_5	0.0	0.0	14.4	1.0	0.0	0.0	17.9	0.0	
M_6	0.0	0.0	8.3	0.0	0.0	0.0	13.1	0.0	
M_7	0.0	0.0	8.0	0.0	0.0	0.0	15.9	0.0	

Table 2: Frequencies selected by the four criteria AIC_0 , AIC_λ , C_0 and C_λ in 1,000 replications for the extreme case of $n = 76$, $K = 7$, $p = 65$, namely $\nu_K = n - K - p - 3 = 1$

M_k	AIC_0	AIC_λ	C_0	C_λ	AIC_0	AIC_λ	C_0	C_λ	
$k^* = 2$					$k^* = 3$				
M_1	100.0	0.0	0.0	0.7	100.0	0.0	0.0	0.2	
M_2	0.0	100.0	50.7	99.3	0.0	0.0	0.0	0.7	
M_3	0.0	0.0	9.7	0.0	0.0	100.0	52.0	99.1	
M_4	0.0	0.0	8.3	0.0	0.0	0.0	14.7	0.0	
M_5	0.3	0.1	8.1	0.0	0.0	0.0	7.9	0.0	
M_6	0.0	0.0	8.8	0.0	0.0	0.0	9.3	0.0	
M_7	0.0	0.0	14.4	0.0	0.0	0.0	16.1	0.0	
$k^* = 4$					$k^* = 5$				
M_1	100.0	0.0	0.0	0.0	100.0	0.0	0.0	0.0	
M_2	0.0	0.0	0.0	0.1	0.0	0.0	0.0	0.0	
M_3	0.0	0.0	0.0	0.1	0.0	0.0	0.0	0.0	
M_4	0.0	100.0	54.8	99.8	0.0	0.0	0.0	0.0	
M_5	0.0	0.0	15.2	0.0	0.0	100.0	63.3	100.0	
M_6	0.0	0.0	14.3	0.0	0.0	0.0	17.3	0.0	
M_7	0.0	0.0	15.7	0.0	0.0	0.0	19.4	0.0	
$k^* = 6$					$k^* = 7$				
M_1	100.0	0.0	0.0	0.0	100.0	0.0	0.0	0.0	
M_2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
M_6	0.0	100.0	71.6	100.0	0.0	0.0	0.0	0.0	
M_7	0.0	0.0	28.4	0.0	0.0	100.0	100.0	100.0	

follows:

$$\begin{aligned}\mathbf{x}_{(3)} &= 0.3\mathbf{x}_{(1)} + 0.7\mathbf{x}_{(2)} + \varepsilon Z_1, \\ \mathbf{x}_{(5)} &= 0.5\mathbf{x}_{(3)} + 0.5\mathbf{x}_{(4)} + \varepsilon Z_2, \\ \mathbf{x}_{(7)} &= 0.7\mathbf{x}_{(5)} + 0.3\mathbf{x}_{(6)} + \varepsilon Z_3,\end{aligned}$$

where ε is a positive constant and Z_1 , Z_2 and Z_3 are mutually independently distributed as a standard normal distribution. For smaller ε , $\widetilde{\mathbf{X}}$ is closer to the multicollinearity case. In this experiment, we treat the two cases: $\varepsilon = 1$ and $\varepsilon = 0.0001$, which correspond to the non-multicollinearity and the multicollinearity cases, respectively. In the multicollinearity case, the ridge parameter τ in the ridge regression estimator $\widehat{\boldsymbol{\beta}}_\tau$ should be large since $(\mathbf{X}'\mathbf{X})^{-1}$ is instable. Define $L(\widetilde{\mathbf{X}})$ by

$$L(\widetilde{\mathbf{X}}) = \{|\widetilde{\mathbf{X}}'\widetilde{\mathbf{X}}/\text{tr}[\widetilde{\mathbf{X}}'\widetilde{\mathbf{X}}/K]| - \log(|\widetilde{\mathbf{X}}'\widetilde{\mathbf{X}}/\text{tr}[\widetilde{\mathbf{X}}'\widetilde{\mathbf{X}}/K]|) - 1\}/K,$$

which measures the discrepancy between the two matrices $\widetilde{\mathbf{X}}'\widetilde{\mathbf{X}}$ and $\text{tr}[\widetilde{\mathbf{X}}'\widetilde{\mathbf{X}}/K]\mathbf{I}_K$. $L(\widetilde{\mathbf{X}})$ takes a large value when $\widetilde{\mathbf{X}}$ is close to the multicollinearity. For the double ridge AIC, we select the regressor variables and the ridge parameter τ so as to minimize $AIC_{\lambda,\tau}$ for $0 \leq \tau \leq L(\widetilde{\mathbf{X}})/5$. For the double ridge C_p , the regressor variables and the ridge parameter τ are selected to minimize $C_{\lambda,\tau}$ for $0 \leq \tau \leq L(\widetilde{\mathbf{X}})/10$. The frequencies selected by AIC_λ , $AIC_{\lambda,\tau}$, C_λ and $C_{\lambda,\tau}$ in this experiment are reported in Table 3. For $\varepsilon = 1$, the non-multicollinearity case, there are little difference between $(AIC_{\lambda,\tau}, C_{\lambda,\tau})$ and (AIC_λ, C_λ) . For $\varepsilon = 0.0001$, which is close to the multicollinearity case, the double ridge criteria $AIC_{\lambda,\tau}$ and $C_{\lambda,\tau}$ are slightly better than AIC_λ and C_λ . When the true model is M_6 , we can observe that $AIC_{\lambda,\tau}$ performs well while AIC_λ does not work.

4.2 An application to posted land price data

We here treat the posted land price data along the Keikyu train line which connects the suburbs in Kanagawa prefecture to the Tokyo metropolitan area. Those who live in the suburbs take this line to work or study in Tokyo every weekday. Thus, it is expected that the land price depends on the distance from Tokyo. We use the selection procedures AIC_0 , AIC_λ , C_0 and C_λ to search for the covariates which affect the land price.

The posted land price data for fifteen years from 1987 to 2001 are available for 47 sites along the Keikyu train line. Each site is indexed by i , namely, $i = 1, \dots, n$ for $n = 47$. The values which are transformed by logarithm from the posted land price (Yen) per m^2 of the i -th site for the fifteen years are described by $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})$ for $T = 15$. For each y_{it} , we consider the following five explanatory variables: T_{1i} is the time to take from the nearby station to the Tokyo station around 8:30 in the morning, T_{2i} is the time to take on foot from the site i to the nearby station and FAR_i and ACR_i denote, respectively, the floor-area ratio and the acreage of the site i . Also, TKY_i is the dummy variable indicating whether the site i is in Tokyo or in Kanagawa prefecture, namely $TKY_i = 0$ if the site i is in Tokyo, otherwise $TKY_i = 1$. As the full model, we consider the mixed linear model

$$y_{it} = \beta_{0t} + T_{1i}\beta_{1t} + (T_{1i}^2)\beta_{2t} + T_{2i}\beta_{3t} + FAR_i\beta_{4t} + TKY_i\beta_{5t} + ACR_i\beta_{6t} + e_{it}.$$

For simplicity, the regressor variables 1, T_{1i} , T_{1i}^2 , T_{2i} , FAR_i , TKY_i and ACR_i are denoted by x_{0i} , x_{1i} , x_{2i} , x_{3i} , x_{4i} , x_{5i} and x_{6i} . Let $\mathbf{X} = (\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_6) = (\mathbf{x}_{(1)}, \dots, \mathbf{x}_{(N)})'$, which

Table 3: Frequencies selected by the four criteria AIC_λ , $AIC_{\lambda,\tau}$, C_λ and $C_{\lambda,\tau}$ in 1,000 replications for $n = 22$, $K = 7$, $p = 10$ and $\nu_K = n - K - p - 3 = 2$ in the case of multicollinearity under the true models M_2, M_4, M_6

	$\varepsilon = 1$, non-multicollinearity				$\varepsilon = 0.0001$, multicollinearity			
M_k	AIC_λ	$AIC_{\lambda,\tau}$	C_λ	$C_{\lambda,\tau}$	AIC_λ	$AIC_{\lambda,\tau}$	C_λ	$C_{\lambda,\tau}$
M_2 : the true model								
M_1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
M_2	100.0	100.0	95.9	95.9	100.0	100.0	95.9	99.8
M_3	0.0	0.0	3.7	3.7	0.0	0.0	3.7	0.0
M_4	0.0	0.0	0.3	0.3	0.0	0.0	0.3	0.2
M_5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
M_6	0.0	0.0	0.1	0.1	0.0	0.0	0.1	0.0
M_7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
M_4 : the true model								
M_1	0.0	0.0	0.0	0.0	0.1	0.0	0.0	0.0
M_2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
M_3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
M_4	100.0	100.0	97.0	96.9	99.9	100.0	97.0	99.3
M_5	0.0	0.0	2.4	2.5	0.0	0.0	2.4	0.0
M_6	0.0	0.0	0.3	0.3	0.0	0.0	0.3	0.7
M_7	0.0	0.0	0.3	0.3	0.0	0.0	0.3	0.0
M_6 : the true model								
M_1	4.8	4.2	0.0	0.0	83.5	0.0	0.0	0.0
M_2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
M_3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
M_4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
M_5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
M_6	95.2	95.8	97.0	97.0	16.5	100.0	97.0	100.0
M_7	0.0	0.0	3.0	3.0	0.0	0.0	3.0	0.0

is an $n \times 7$ matrix, for $\mathbf{x}_j = (x_{j1}, \dots, x_{jN})'$ and $\mathbf{x}'_{(i)} = (x_{0i}, x_{1i}, x_{2i}, x_{3i}, x_{4i}, x_{5i}, x_{6i})$. Also let $\mathbf{Y} = (\mathbf{y}'_1, \dots, \mathbf{y}'_n)'$ for $\mathbf{y}_i = (y_{i1}, \dots, y_{iT})$, and \mathbf{E} is similarly defined. Then, the model is expressed as $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{E}$, where $\boldsymbol{\beta} = (\boldsymbol{\beta}_{(1)}, \dots, \boldsymbol{\beta}_{(T)})$, which is a $7 \times T$ matrix, for $\boldsymbol{\beta}_{(t)} = (\beta_{0t}, \dots, \beta_{6t})'$.

Table 4 reports values of AIC_0 , AIC_λ , C_0 and C_λ for several candidate models, where the regressor variable which minimizes AIC_λ is added to the model based on the forward selection rule. Among these candidate models, the minimum value of AIC_λ is -670 and attained by the model with the regressor variables $\{\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_4, \mathbf{x}_5\}$, while AIC_0 and C_λ select $\{\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_5\}$ or $\{\mathbf{x}_0, \mathbf{x}_1, \mathbf{x}_4, \mathbf{x}_5\}$. It is also observed that C_0 selects the full model, which shows that C_0 does not work well for this example. According to these observations based on AIC_0 , AIC_λ and C_λ , we can recommend the model given by

$$y_{it} = \beta_{0t} + T_{1i}\beta_{1t} + FAR_i\beta_{4t} + TKY_i\beta_{5t} + e_{it}.$$

Although values of $AIC_{\lambda,\tau}$ and $C_{\lambda,\tau}$ are not reported in Table 4, it is noted that their values are very close to those of AIC_λ and C_λ , respectively.

We here investigate whether the selected model is endorsed by a testing procedure. The general linear hypothesis is expressed as a testing of hypothesis

$$H : \mathbf{C}\boldsymbol{\beta} = \mathbf{0} \quad \text{vs} \quad A : \mathbf{C}\boldsymbol{\beta} \neq \mathbf{0}$$

where \mathbf{C} is a known $m \times 7$ matrix of rank $m \leq 7$. The error sum of squares and products is given by the matrix

$$\mathbf{V} = \mathbf{Y}'(\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{Y},$$

and the sum of squares and products due to regression under the hypotheses H is

$$\mathbf{W} = \hat{\boldsymbol{\beta}}'\mathbf{C}'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}\mathbf{C}\hat{\boldsymbol{\beta}}.$$

To test the hypothesis H we need to compare these matrices. The likelihood ratio test rejects the hypothesis H if

$$\frac{|\mathbf{V}|}{|\mathbf{V} + \mathbf{W}|} = U_{p,m,f} \leq U_{p,m,f,\alpha}$$

where $f = n - 7$ and $U_{p,m,f,\alpha}$ is the upper $100\alpha\%$ point of the distribution of $U_{p,m,f}$. The asymptotic approximation for $U_{p,m,f}$ is given by

$$\begin{aligned} &P[-\{f - (p - m + 1)/2\} \log U_{p,m,f} \geq z] \\ &= P[\chi_{pm}^2 \geq z] + f^{-2}\gamma_2 \{P[\chi_{pm+4}^2 \geq z] - P[\chi_{pm}^2 \geq z]\} \end{aligned} \quad (43)$$

where $\gamma_2 = pm(p^2 + p - 5)/48$. See Srivastava (2002, p.282).

Let $\boldsymbol{\beta}$ be decomposed into $\boldsymbol{\beta} = (\boldsymbol{\beta}'_0, \boldsymbol{\beta}'_1, \dots, \boldsymbol{\beta}'_k)'$. When the null hypothesis $H : \boldsymbol{\beta}_0 = \boldsymbol{\beta}_1 = \boldsymbol{\beta}_4 = \boldsymbol{\beta}_5 = \mathbf{0}$ is tested, the P-value given in (43) is 0.000, and the hypothesis is rejected strongly. When each hypothesis $H_i : \boldsymbol{\beta}_i = \mathbf{0}$ is tested for $i = 0, 1, \dots, 6$, the P-values are given by $P_0 = P_1 = P_2 = P_4 = P_5 = 0.000$, $P_3 = 0.022$ and $P_6 = 0.008$, where P_i is the P-value given in (43) for testing H_i . When the P-values can be also obtained numerically based on simulation experiments, those values for testing the hypotheses H_3 and H_6 are 0.027 and 0.007. Thus, it may be plausible that the variables \mathbf{x}_3 and \mathbf{x}_6 , namely T_{2i} and ACR_i are deleted from the regressor variables.

Table 4: Selection of regressor variables in the posted land price data

k	x_i	AIC_0	AIC_λ	C_0	C_λ
1	x_0	-3040	105	2051	286
2	x_0, x_1	-3106	-490	1279	68
3	x_0, x_1, x_2	-3108	-489	906	54
3	x_0, x_1, x_3	-3070	-438	1271	109
3	x_0, x_1, x_4	-3104	-619	1158	74
3	x_0, x_1, x_5	-3174	-521	353	-14
3	x_0, x_1, x_6	-3081	-445	1215	95
4	x_0, x_1, x_4, x_2	-3101	-594	777	58
4	x_0, x_1, x_4, x_3	-3064	-578	1150	113
4	x_0, x_1, x_4, x_5	-3167	-670	232	-9
4	x_0, x_1, x_4, x_6	-3067	-560	1136	108
5	x_0, x_1, x_4, x_5, x_2	-3145	-601	160	29
5	x_0, x_1, x_4, x_5, x_3	-3129	-634	211	30
5	x_0, x_1, x_4, x_5, x_6	-3128	-605	215	26
6	$x_0, x_1, x_4, x_5, x_3, x_2$	-3106	-566	127	68
6	$x_0, x_1, x_4, x_5, x_3, x_6$	-3085	-565	190	67
7	$x_0, x_1, x_4, x_5, x_3, x_2, x_6$	-3056	-490	105	105

5 Concluding remarks

The variable selection problem in the multivariate linear regression model is addressed under the asymptotic condition that both n and p tend to infinity subject to $n - k - p - 3 > 0$ and $\lim_{(n,p) \rightarrow \infty} p/n = c$ for $0 < c < 1$. In this paper, we have proposed the modified AIC and C_p statistic, denoted by AIC_λ and C_λ , based on the ridge-type estimator of Σ instead of the MLE, and proved their analytical justifications, namely, they are asymptotic unbiased estimators of the quantities related to the prediction errors. We also have extended the modified AIC and C_p statistic to the double ridge-type criteria which use the ridge regression estimator of β instead of the least squares estimator.

Through simulation studies reported in Tables 1 and 2, it is seen that AIC_0 and C_0 statistic, based on MLE of Σ , perform well for small p and large n , as it should be. The performances of AIC_λ and C_λ are, however, equally good and somewhat better. In contrast for large p , the performance of AIC_0 and C_0 are rather poor in comparison to the performance of AIC_λ and C_λ . In the case close to the multicollinearity, the double ridge-type criteria $AIC_{\lambda,\tau}$ and $C_{\lambda,\tau}$ have been shown to work well. Thus we recommend the use of AIC_λ and C_λ , or $AIC_{\lambda,\tau}$ and $C_{\lambda,\tau}$ for all p so long as $n - k - p > 3$.

A Appendix

Lemma A.1 *Let $\mathbf{S} \sim \mathcal{W}_p(\boldsymbol{\Sigma}, m)$. Then,*

$$E[(\text{tr } \mathbf{S})(\text{tr } \mathbf{S}^{-1})] = \frac{m}{m-p-1} \text{tr } \boldsymbol{\Sigma} \text{tr } \boldsymbol{\Sigma}^{-1} - \frac{2p}{m-p-1}.$$

Proof. Since $\text{tr } \mathbf{S}$ and $\text{tr } \mathbf{S}^{-1}$ are invariant under an orthogonal transformation, we may assume that $\boldsymbol{\Sigma}$ is a diagonal matrix with diagonal elements σ_i , $i = 1, \dots, p$, $\boldsymbol{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_p)$. Thus, with $\mathbf{W} \sim \mathcal{W}_p(\mathbf{I}, m)$, $\mathbf{W} = (w_{ij})$, $\mathbf{W}^{-1} = (w^{ij})$, we get

$$\text{tr } \mathbf{S} = \text{tr } \boldsymbol{\Sigma} \mathbf{W} = \sum_i \sigma_i w_{ii}.$$

Hence,

$$(\text{tr } \mathbf{S})(\text{tr } \mathbf{S}^{-1}) = \left(\sum_{i=1}^p \sigma_i w_{ii} \right) \left(\sum_{i=1}^p \sigma_i^{-1} w^{ii} \right) = \sum_{i=1}^p w_{ii} w^{ii} + \sum_{i \neq j} \sigma_i \sigma_j^{-1} w_{ii} w^{jj}.$$

Noting that $E[w_{ii} w^{ii}] = E[w_{pp} w^{pp}]$ for any i , and $E[w_{ii} w^{jj}] = E[w_{11} w^{pp}]$ for any $i \neq j$, we get

$$E[(\text{tr } \mathbf{S})(\text{tr } \mathbf{S}^{-1})] = pE[w_{pp} w^{pp}] + \sum_{i \neq j} \sigma_i \sigma_j^{-1} E[w_{11} w^{pp}].$$

Consider now the triangular factorization of $\mathbf{W} = \mathbf{T} \mathbf{T}'$, where

$$\mathbf{T} = \begin{pmatrix} \mathbf{T}_1 & \mathbf{0} \\ \mathbf{t}'_{12} & t_{pp} \end{pmatrix}.$$

Then,

$$\mathbf{W} = \begin{pmatrix} \mathbf{T}_1 \mathbf{T}'_1 & \mathbf{T}_1 \mathbf{t}_{12} \\ \mathbf{t}'_{12} \mathbf{T}'_1 & t_{pp}^2 + \mathbf{t}'_{12} \mathbf{t}_{12} \end{pmatrix} = \begin{pmatrix} \mathbf{W}_{11} & \mathbf{w}_{12} \\ \mathbf{w}'_{12} & w_{pp} \end{pmatrix}$$

and

$$w^{pp} = (w_{pp} - \mathbf{w}'_{12} \mathbf{W}_{11}^{-1} \mathbf{w}_{12})^{-1} = (t_{pp}^2)^{-1}.$$

Hence,

$$\begin{aligned} E[w_{pp} w^{pp}] &= E \left[\frac{t_{pp}^2 + \mathbf{t}'_{12} \mathbf{t}_{12}}{t_{pp}^2} \right] = 1 + E \left[\frac{\mathbf{t}'_{12} \mathbf{t}_{12}}{t_{pp}^2} \right] \\ &= 1 + E[\mathbf{t}'_{12} \mathbf{t}_{12}] E[t_{pp}^{-2}] = 1 + \frac{p-1}{m-p-1} = \frac{m-2}{m-p-1}, \end{aligned}$$

since $\mathbf{t}_{12} \sim \mathcal{N}_{p-1}(\mathbf{0}, \mathbf{I})$ is independently distributed of t_{pp}^2 , and t_{pp}^2 is distributed as chisquare with $m-p+1$ degrees of freedom, see Srivastava and Khatri (1979, Lemma 3.2.1, pp 74). Similarly,

$$E[w_{11} w^{pp}] = E[t_{11}^2 / t_{pp}^2] = E[t_{11}^2] E[1/t_{pp}^2] = \frac{m}{m-p-1}.$$

Hence,

$$\begin{aligned}
E[(\text{tr } \mathbf{S})(\text{tr } \mathbf{S}^{-1})] &= \frac{(m-2)p}{m-p-1} + \frac{m}{m-p-1} \sum_{i \neq j} \sigma_i \sigma_j^{-1} \\
&= \frac{m}{m-p-1} [(\text{tr } \boldsymbol{\Sigma})(\text{tr } \boldsymbol{\Sigma}^{-1}) - p] + \frac{(m-2)p}{m-p-1} \\
&= \frac{m}{m-p-1} (\text{tr } \boldsymbol{\Sigma})(\text{tr } \boldsymbol{\Sigma}^{-1}) - \frac{2p}{m-p-1}.
\end{aligned}$$

■

Lemma A.2 *Let $\mathbf{S} \sim \mathcal{W}_p(\boldsymbol{\Sigma}, m)$. Then,*

$$\begin{aligned}
E[(\text{tr } \mathbf{S})\text{tr } (\boldsymbol{\Sigma} \mathbf{S}^{-2})] &= \frac{(m-1)(m+1)}{(m-p+1)(m-p-1)(m-p-3)} (\text{tr } \boldsymbol{\Sigma})(\text{tr } \boldsymbol{\Sigma}^{-1}) \\
&\quad - \frac{p}{(m-p-1)(m-p-3)} \left(\frac{m^2-1}{m-p+1} - \frac{m^2-5m+2p+2}{m-p} \right).
\end{aligned}$$

Proof. As explained above, we assume that $\boldsymbol{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_p)$, and $\mathbf{W} \sim \mathcal{W}_p(\mathbf{I}, m)$,

$$\begin{aligned}
E[(\text{tr } \mathbf{S})(\text{tr } \boldsymbol{\Sigma} \mathbf{S}^{-2})] &= E[(\text{tr } \boldsymbol{\Sigma} \mathbf{W})(\text{tr } \boldsymbol{\Sigma}^{-1} \mathbf{W}^{-2})] \\
&= E\left[\left(\sum_i \sigma_i w_{ii} \right) \left(\sum_i \sigma_i^{-1} (\mathbf{W}^{-2})_{ii} \right) \right] \\
&= E\left[\sum_{i=1}^p w_{ii} (\mathbf{W}^{-2})_{ii} + \sum_{i \neq j} \sigma_i \sigma_j^{-1} w_{ii} (\mathbf{W}^{-2})_{jj} \right] \\
&= p E[w_{pp} (\mathbf{W}^{-2})_{pp} + w_{11} (\mathbf{W}^{-2})_{pp} \sum_{i \neq j} \sigma_i \sigma_j^{-1}].
\end{aligned}$$

Note that

$$(\mathbf{W}^{-2})_{pp} = \frac{1}{t_{pp}^4} [1 + \mathbf{t}'_{12} (\mathbf{T}'_1 \mathbf{T}_1)^{-1} \mathbf{t}_{12}],$$

and $w_{pp} = t_{pp}^2 + \mathbf{t}'_{12} \mathbf{t}_{12}$. Thus,

$$\begin{aligned}
E[w_{pp} (\mathbf{W}^{-2})_{pp}] &= E[(t_{pp}^2 + \mathbf{t}'_{12} \mathbf{t}_{12})(1 + \mathbf{t}'_{12} (\mathbf{T}'_1 \mathbf{T}_1)^{-1} \mathbf{t}_{12} / t_{pp}^4)] \\
&= E[(t_{pp}^{-2} + \mathbf{t}'_{12} \mathbf{t}_{12} t_{pp}^{-4})(1 + \mathbf{t}'_{12} (\mathbf{T}'_1 \mathbf{T}_1)^{-1} \mathbf{t}_{12})] \\
&= E\left[\left(\frac{1}{m-p-1} + \frac{\mathbf{t}'_{12} \mathbf{t}_{12}}{(m-p-1)(m-p-3)} \right) (1 + \mathbf{t}'_{12} (\mathbf{T}'_1 \mathbf{T}_1)^{-1} \mathbf{t}_{12}) \right] \\
&= \frac{1}{m-p-1} E[1 + \mathbf{t}'_{12} (\mathbf{T}'_1 \mathbf{T}_1)^{-1} \mathbf{t}_{12}] + E\left[\frac{\mathbf{t}'_{12} \mathbf{t}_{12} \mathbf{t}'_{12} (\mathbf{T}'_1 \mathbf{T}_1)^{-1} \mathbf{t}_{12} + \mathbf{t}'_{12} \mathbf{t}_{12}}{(m-p-1)(m-p-3)} \right] \\
&= \frac{1}{m-p-1} \{1 + E[\text{tr } (\mathbf{T}'_1 \mathbf{T}_1)^{-1}]\} + E\left[\frac{\mathbf{t}'_{12} \mathbf{t}_{12} \mathbf{t}'_{12} (\mathbf{T}'_1 \mathbf{T}_1)^{-1} \mathbf{t}_{12} + \mathbf{t}'_{12} \mathbf{t}_{12}}{(m-p-1)(m-p-3)} \right].
\end{aligned}$$

Note that $\mathbf{t}'_{12}\mathbf{t}_{12}\mathbf{t}'_{12}(\mathbf{T}'_1\mathbf{T}_1)^{-1}\mathbf{t}_{12} = \text{tr}(\mathbf{t}_{12}\mathbf{t}'_{12})^2(\mathbf{T}'_1\mathbf{T}_1)^{-1}$, where $\mathbf{t}_{12} \sim \mathcal{N}_{p-1}(\mathbf{0}, \mathbf{I}_{p-1})$ and \mathbf{T}_1 are independently distributed. Hence, $\mathbf{t}_{12}\mathbf{t}'_{12} \sim \mathbf{W}_{p-1}(\mathbf{I}, 1)$. From Srivastava and Khatri (1979, Problem 3.2, pp 97),

$$E[(\mathbf{t}_{12}\mathbf{t}'_{12})^2] = 2\mathbf{I}_{p-1} + (p-1)\mathbf{I}_{p-1} = (p+1)\mathbf{I}_{p-1}.$$

Hence,

$$\begin{aligned} E[\mathbf{t}'_{12}\mathbf{t}_{12}\mathbf{t}'_{12}(\mathbf{T}'_1\mathbf{T}_1)^{-1}\mathbf{t}_{12}] &= (p+1)E[\text{tr}(\mathbf{T}'_1\mathbf{T}_1)^{-1}] \\ &= (p+1)E[\text{tr}(\mathbf{T}_1\mathbf{T}'_1)^{-1}] = \frac{(p+1)(p-1)}{m-p}, \end{aligned}$$

since $\mathbf{T}_1\mathbf{T}'_1 \sim \mathcal{W}_{p-1}(\mathbf{I}_{p-1}, m)$, and $E[(\mathbf{T}_1\mathbf{T}'_1)^{-1}] = (m-p)^{-1}\mathbf{I}_{p-1}$. Hence,

$$E[w_{pp}(\mathbf{W}^{-2})_{pp}] = \frac{1}{m-p-1} \left[1 + \frac{p-1}{m-p} + \frac{p-1}{m-p-3} + \frac{(p+1)(p-1)}{(m-p)(m-p-3)} \right].$$

We shall need to calculate

$$\begin{aligned} E[w_{11}(\mathbf{W}^{-2})_{pp}] &= E\left[\frac{t_{11}^2}{t_{pp}^4}(1 + \mathbf{t}'_{12}(\mathbf{T}'_1\mathbf{T}_1)^{-1}\mathbf{t}_{12})\right] \\ &= \frac{1}{(m-p-1)(m-p-3)} E[m + t_{11}^2\mathbf{t}'_{12}(\mathbf{T}'_1\mathbf{T}_1)^{-1}\mathbf{t}_{12}]. \end{aligned}$$

It may be noted that t_{11} is the (1, 1)st element of \mathbf{T}_1 , so we need to write \mathbf{T}_1 as

$$\mathbf{T}_1 = \begin{pmatrix} t_{11} & \mathbf{0} \\ \mathbf{t}_{31} & \mathbf{T}_3 \end{pmatrix}, \quad \mathbf{T}_1^{-1} = \begin{pmatrix} t_{11}^{-1} & \mathbf{0} \\ -\mathbf{T}_3^{-1}t_{11}^{-1}\mathbf{t}_{31} & \mathbf{T}_3^{-1} \end{pmatrix}.$$

Thus,

$$\begin{aligned} E[t_{11}^2\text{tr}(\mathbf{T}'_1\mathbf{T}_1)^{-1}\mathbf{t}_{12}\mathbf{t}'_{12}] &= E[t_{11}^2\text{tr}(\mathbf{T}'_1\mathbf{T}_1)^{-1}] \\ &= E\left[t_{11}^2 \left\{ t_{11}^{-2} + \frac{\text{tr}\mathbf{T}_3^{-1}\mathbf{t}_{31}\mathbf{t}'_{31}(\mathbf{T}_3^{-1})'}{t_{11}^2} + \text{tr}\mathbf{T}_3^{-1}(\mathbf{T}_3^{-1})' \right\}\right] \\ &= E\left[1 + \text{tr}\mathbf{t}_{31}\mathbf{t}'_{31}(\mathbf{T}_3\mathbf{T}'_3)^{-1} + t_{11}^2\text{tr}(\mathbf{T}_3\mathbf{T}'_3)^{-1}\right] \\ &= 1 + (m+1)E[\text{tr}(\mathbf{T}_3\mathbf{T}'_3)^{-1}] = 1 + \frac{(m+1)(p-2)}{m-p+1}. \end{aligned}$$

Combining all the above calculation, we get

$$E[w_{11}(\mathbf{W}^{-2})_{pp}] = \frac{(m-1)(m+1)}{(m-p+1)(m-p-1)(m-p-3)}.$$

Hence, after some simplification, we get

$$\begin{aligned} E[(\text{tr}\mathbf{S})(\text{tr}\mathbf{\Sigma}\mathbf{S}^{-2})] &= \frac{(m-1)(m+1)}{(m-p+1)(m-p-1)(m-p-3)}(\text{tr}\mathbf{\Sigma})(\text{tr}\mathbf{\Sigma}^{-1}) \\ &\quad - \frac{p}{(m-p-1)(m-p-3)} \left(\frac{m^2-1}{m-p-1} - \frac{m^2-5m+2p+2}{m-p} \right). \end{aligned}$$

■

Acknowledgments. The research was supported by NSERC. The research of the second author was supported in part by Grant-in-Aid for Scientific Research (19200020 and 21540114), Japan.

References

- [1] Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *2nd International Symposium on Information Theory*, (B.N. Petrov and Csaki, F, eds.), 267-281, Akademia Kiado, Budapest.
- [2] Akaike, H. (1974). A new look at the statistical model identification. System identification and time-series analysis. *IEEE Trans. Autom. Contr.*, **AC-19**, 716-723.
- [3] Bai, Z.D., and Yin, Y.Q. (1993). Limit of the smallest eigenvalue of a large dimensional sample covariance matrix. *Ann. Prob.*, **21**, 1275-1294.
- [4] Fujikoshi, Y. and Satoh, K. (1997). Modified AIC and C_p in multivariate linear regression. *Biometrika*, **84**, 707-716.
- [5] Johnston, I.M. (2001). On the distribution of the largest eigenvalue in principal component analysis. *Ann. Statist.*, **29**, 295-327.
- [6] Konishi, S. and Kitagawa, G. (2007). *Information Criteria and Statistical Modeling*. Springer.
- [7] Mallows, C.L. (1973). Some comments on C_p . *Technometrics*, **15**, 661-676.
- [8] Srivastava, M.S. (2002). *Methods of Multivariate Statistics*, Wiley, New York.
- [9] Srivastava, M.S. (2005). Some test concerning the covariance matrix in high dimensional data. *J. Japan Statist. Soc.*, **35**, 251-272.
- [10] Srivastava, M.S. (2007). Multivariate theory for analyzing high dimensional data. *J. Japan Statist. Soc.*, **37**, 53-86.
- [11] Srivastava, M.S., and Khatri, C.G. (1979). *An Introduction to Multivariate Statistics*. North-Holland, New York.
- [12] Srivastava, M.S. and Kubokawa, T. (2007). Comparison of discrimination methods for high dimensional data. *J. Japan Statist. Soc.*, **37**, 123-134.
- [13] Srivastava, M.S. and Kubokawa, T. (2008). Akaike information criterion for selecting components of the mean vector in high dimensional data with fewer observations. *J. Japan Statist. Soc.*, **38**, 259-283.
- [14] Sugiura, N. (1978). Further analysis of the data by Akaike's information criterion and the finite corrections. *Commun. Statist. - Theory Methods*, **1**, 13-26.
- [15] Yamamura, M., Yanagihara, H. and Srivastava, M.S. (2009). Variable selection in multivariate linear regression models with fewer observations than the dimension. Unpublished manuscript.