

# A consistency property of the AIC for multivariate linear models when the dimension and the sample size are large\*

Hirokazu Yanagihara<sup>†</sup>, Hirofumi Wakaki<sup>‡</sup> and Yasunori Fujikoshi<sup>§</sup>

*Department of Mathematics, Graduate School of Science, Hiroshima University, 1-3-1*

*Kagamiyama, Higashi-Hiroshima, Hiroshima 739-8626, Japan*

*e-mail: [yanagi@math.sci.hiroshima-u.ac.jp](mailto:yanagi@math.sci.hiroshima-u.ac.jp); [wakaki@math.sci.hiroshima-u.ac.jp](mailto:wakaki@math.sci.hiroshima-u.ac.jp);*

*[fujikoshi\\_y@yahoo.co.jp](mailto:fujikoshi_y@yahoo.co.jp)*

**Abstract:** It is common knowledge that Akaike's information criterion (AIC) is not a consistent model selection criterion, and Bayesian information criterion (BIC) is. These have been confirmed from an asymptotic selection probability evaluated from a large-sample framework. However, when a high-dimensional asymptotic framework, such that the dimension of the response variables and the sample size are approaching  $\infty$ , is used for evaluating the selection probability, there are cases that the AIC for selecting variables in multivariate linear models is consistent, but the BIC is not. The AIC and BIC are included in a family of information criteria defined by adding a penalty term expressing the complexity of the model to a negative twofold maximum log-likelihood. By clarifying the condition of the penalty term to ensure the consistency, we derive conditions for consistency of the AIC, BIC and other information criteria under the high-dimensional asymptotic framework.

**MSC 2010 subject classifications:** Primary 62J05; secondary 62E20.

**Keywords and phrases:** AIC, bias-corrected AIC, BIC, consistent AIC, high-dimensional asymptotic framework, multivariate linear model, selection probability, variable selection.

Received September 2013.

## 1. Introduction

Let  $\mathbf{Y}$  be an  $n \times p$  observation matrix of  $p$  response variables, and let  $\mathbf{X}$  be an  $n \times k$  observation matrix of  $k$  nonstochastic explanatory variables, where  $n$  is the sample size, and it is assumed that  $n - p - k - 1 > 0$ . In order to ensure the

---

\*The authors wish to thank the referees, the associate editor and the editor for their helpful suggestions.

<sup>†</sup>The first author's research is partially supported by the Ministry of Education, Science, Sports, and Culture, a Grant-in-Aid for Challenging Exploratory Research, #25540012, 2013–2015.

<sup>‡</sup>The second author's research is partially supported by the Ministry of Education, Science, Sports, and Culture, a Grant-in-Aid for Scientific Research (C), #24500343, 2012–2014.

<sup>§</sup>The third author's research is partially supported by the Ministry of Education, Science, Sports, and Culture, a Grant-in-Aid for Scientific Research (C), #25330038, 2013–2015.

possibility of estimating the model, we also assume that  $\text{rank}(\mathbf{X}) = k (< n)$ . Suppose that  $j$  denotes a subset of  $\omega = \{1, \dots, k\}$  containing  $k_j$  elements, and  $\mathbf{X}_j$  denotes the  $n \times k_j$  matrix consisting of the columns of  $\mathbf{X}$  indexed by the elements of  $j$ . For example, if  $j = \{1, 2, 4\}$ , then  $\mathbf{X}_j$  consists of the first, second, and fourth columns of  $\mathbf{X}$ . Of course, it holds that  $\mathbf{X}_\omega = \mathbf{X}$  and  $k_\omega = k$ . Also, we let  $k_A$  denote the number of elements of a set  $A$ , i.e.,  $k_A = \#(A)$ . Then the following multivariate linear regression model with  $k_j$  explanatory variables is considered as the candidate model:

$$\mathbf{Y} \sim N_{n \times p}(\mathbf{X}_j \boldsymbol{\Theta}_j, \boldsymbol{\Sigma}_j \otimes \mathbf{I}_n), \quad (1.1)$$

where  $\boldsymbol{\Theta}_j$  is a  $k_j \times p$  unknown matrix of regression coefficients, and  $\boldsymbol{\Sigma}_j$  is a  $p \times p$  unknown covariance matrix. Here,  $\mathbf{A} \otimes \mathbf{B}$  denotes the Kronecker product of an  $m \times n$  matrix  $\mathbf{A}$  and a  $p \times q$  matrix  $\mathbf{B}$ , which is an  $mp \times nq$  matrix defined by

$$\mathbf{A} \otimes \mathbf{B} = \begin{pmatrix} a_{11}\mathbf{B} & \cdots & a_{1n}\mathbf{B} \\ \vdots & \ddots & \vdots \\ a_{m1}\mathbf{B} & \cdots & a_{mn}\mathbf{B} \end{pmatrix},$$

where  $a_{ij}$  is the  $(i, j)$ th element of  $\mathbf{A}$  (see, e.g., [18, chap. 16]). In particular, the model with  $\mathbf{X}_\omega$  (namely  $\mathbf{X}$ ) is called the full model. We will assume that the data are generated from the following true model:

$$\mathbf{Y} \sim N_{n \times p}(\mathbf{X}_{j_*} \boldsymbol{\Theta}_*, \boldsymbol{\Sigma}_* \otimes \mathbf{I}_n), \quad (1.2)$$

where  $j_*$  is a set of integers indicating the subset of explanatory variables in the true model. Henceforth, for simplicity, we represent  $\mathbf{X}_{j_*}$  and  $k_{j_*}$  as  $\mathbf{X}_*$  and  $k_*$ , respectively.

The multivariate linear regression model of (1.1) is one of basic models of multivariate analysis. This model is introduced in many multivariate statistical textbooks (see, e.g., [31, chap. 9], [33, chap. 4]), and even now is widely used in chemometrics, engineering, econometrics, psychometrics, and many other fields, for the prediction of multiple responses from a set of explanatory variables (see, e.g., [10, 24, 25, 40]). Since it is important to specify factors affecting response variables in regression analysis, searching for the optimal subset  $j$  is essential.

Akaike's information criterion (AIC), proposed by [1, 2], is widely used for selecting the best model. The AIC was proposed as an asymptotic unbiased estimator of the risk function assessed by the expected Kullback-Leibler (KL) loss [20] under the assumption that the candidate model includes the true model. One purpose of a model selection method based on the AIC is to choose a model that makes the risk function small. For that purpose, using the AIC for model selection will be asymptotically efficient when the true model is infinite (see, e.g., [27, 29, 39]). A Bayesian information criterion (BIC) proposed by [26] and a consistent AIC (CAIC) proposed by [6] are also widely used for model selection purposes. It is a well-known fact that, when the true model is included in a set of the candidate models, these two criteria are consistent in model selection, i.e., the probability of selecting the true model goes to 1 asymptotically, although

for the AIC is not. When using the AIC for model selection, this inconsistency property sometimes becomes a target for criticism, although the purpose of the AIC is not to choose the true model. The inconsistency property of the AIC is confirmed from the asymptotic probability of selecting the model, which is evaluated from the following asymptotic framework that represents an ordinary asymptotic procedure [11, 12, 22, 28]:

- A large-sample (LS) asymptotic framework: the sample size is approaching  $\infty$  under a fixed number of parameters. In this paper,  $\lim_{n \rightarrow \infty}$  means a limit as  $n \rightarrow \infty$  under the condition that the number of parameters is fixed.

In the case of multivariate linear models, although there are many bias-corrected AICs for the risk function (see, e.g., [3, 14, 17, 37, 38]), such a bias-corrected AIC is still not consistent for model selection.

In recent years, high-dimensional data analysis has been attracting the attention of many researchers. It is known that the LS asymptotic framework gives a poor approximation when the dimension is large. However, the following asymptotic framework gives a better approximation than the LS asymptotic framework when the dimension and the sample size are large, and sometimes even when the dimension is not so large [13, 15, 16]:

- A high-dimensional (HD) asymptotic framework: the sample size and the dimension of the response variables simultaneously approach  $\infty$  under the condition that  $c_{n,p} = p/n \rightarrow c_0 \in [0, 1)$ . For simplicity, we will write “ $(n, p) \rightarrow \infty$  simultaneously under the condition that  $c_{n,p} \rightarrow c_0$ ” as “ $c_{n,p} \rightarrow c_0$ ”, and  $\lim_{c_{n,p} \rightarrow c_0}$  means a limit under the HD asymptotic framework. It should be emphasized that we assume that  $p$  always goes to  $\infty$  in the HD asymptotic framework. Hence, the notation  $c_{n,p} \rightarrow 0$  does not mean the LS asymptotic framework.

When the HD asymptotic framework is used for evaluating the asymptotic probability of selecting the true model, there is a possibility that the AIC can become consistent. In fact, in this paper, we will prove that a variable selection method based on the AIC becomes consistent in multivariate linear models under a HD asymptotic framework. The AIC is included in a family of information criteria defined by adding a penalty term expressing the complexity of the model to a negative twofold maximum likelihood. By clarifying the condition of the penalty term to satisfy the consistency property, we will also prove that a variable selection method based on the bias-corrected AIC (AIC<sub>c</sub>), as proposed by [3], becomes consistent under more non-restrictive situation than that based on the AIC, and those based on the BIC and the CAIC are not necessarily consistent when  $c_0 \in (0, 1)$ . Additionally, we derive a sufficient condition to satisfy the consistency of the family of information criteria under an asymptotic framework such that the number of candidate models may approach  $\infty$ .

In this paper,  $o(x)$ ,  $O(x)$ ,  $o_p(x)$ , and  $O_p(x)$  used in a vector or matrix having finite dimension or size mean that the orders of all the elements in that vector or matrix are  $o(x)$ ,  $O(x)$ ,  $o_p(x)$ , and  $O_p(x)$ , respectively. Furthermore, the Landau notations indicate the orders as  $n \rightarrow \infty$  under a fixed number of parameters

when the LS asymptotic framework is considered. Meanwhile, those Landau notations are also used for the orders as  $c_{n,p} \rightarrow c_0$  when the HD asymptotic framework is considered. As stated already, we deal with not a strong consistency but a weak consistency. Hence, throughout the paper, the word “consistency” means weak consistency.

The remainder of the paper is organized as follows: In Section 2, we present the necessary notation for evaluating an asymptotic selection probability. In Section 3, the asymptotic probability of selecting the true model is calculated under the HD asymptotic framework. In Section 4, we compare with variable selection methods based on the AIC, AIC<sub>c</sub>, BIC and CAIC by conducting numerical experiments. In Section 5, we discuss our conclusions. Technical details are provided in [Appendix](#).

## 2. Preliminaries

In this section, we present and discuss the notation that we used for evaluating the asymptotic selection probability. First, we describe several classes of the set  $j$ . Let  $\mathcal{J}$  be a set of candidate models denoted by  $\mathcal{J} = \{j_1, \dots, j_K\}$ , where  $K$  is the number of candidate models. We then separate  $\mathcal{J}$  into two sets, one of which is a set of overspecified models, candidate models that include the true model, i.e.,  $\mathcal{J}_+ = \{j \in \mathcal{J} | j_* \subseteq j\}$ , and the other is a set of underspecified models that are not the overspecified models, i.e.,  $\mathcal{J}_- = \mathcal{J}_+^c \cap \mathcal{J}$ . We use the same terminology, “overspecified model” and “underspecified model”, as was used by [\[14\]](#).

Estimations for the unknown parameters  $\Theta_j$  and  $\Sigma_j$  in the model [\(1.1\)](#) are carried out by the maximum likelihood method, i.e.,  $\Theta_j$  and  $\Sigma_j$  are estimated by

$$\hat{\Theta}_j = (\mathbf{X}'_j \mathbf{X}_j)^{-1} \mathbf{X}'_j \mathbf{Y}, \quad \hat{\Sigma}_j = \frac{1}{n} \mathbf{Y}' (\mathbf{I}_n - \mathbf{P}_j) \mathbf{Y},$$

where  $\mathbf{P}_j$  is the projection matrix to the subspace spanned by the columns of  $\mathbf{X}_j$ , i.e.,  $\mathbf{P}_j = \mathbf{X}_j (\mathbf{X}'_j \mathbf{X}_j)^{-1} \mathbf{X}'_j$ . A family of information criteria in the model [\(1.1\)](#) is

$$\text{IC}_m(j) = \mathcal{L}(j) + np(\log 2\pi + 1) + m(j), \quad (2.1)$$

where  $\mathcal{L}(j) = n \log \det(\hat{\Sigma}_j)$  and  $m(j)$  is a positive constant expressing a penalty for the complexity of the model [\(1.1\)](#). An information criterion included in this family is specified by an individual penalty term  $m(j)$ . This family contains AIC, AIC<sub>c</sub>, BIC and CAIC as a special case.

$$m(j) = \begin{cases} 2\{pk_j + p(p+1)/2\} & \text{(AIC)} \\ 2n\{pk_j + p(p+1)/2\}/(n - k_j - p - 1) & \text{(AIC}_c\text{)} \\ \{pk_j + p(p+1)/2\} \log n & \text{(BIC)} \\ \{pk_j + p(p+1)/2\}(1 + \log n) & \text{(CAIC)} \end{cases} . \quad (2.2)$$

When  $p = 1$ , the AIC<sub>c</sub> coincides with the bias-corrected AIC proposed by [\[32\]](#). [\[9\]](#) showed that Sugiura’s bias-corrected AIC is a uniformly minimum-variance

unbiased estimator (UMVUE) of the risk function consisting of the expected KL loss when the candidate model includes the true model. By extending the result to the multivariate case, this property can be proved even when  $p > 1$ . The detailed proof is omitted because it can be obtained from the Lehman-Scheffé theorem and the fact that  $\hat{\Theta}_j$  and  $\hat{\Sigma}_j$  are complete sufficient statistics. Complete efficiencies of  $\hat{\Theta}_j$  and  $\hat{\Sigma}_j$  can be derived by slightly modifying the results of [30, pp. 18–20]. This property indicates that, for all the overspecified models, the AIC<sub>c</sub> is better than the AIC at estimating the risk function. The best subsets of  $\omega$  is chosen by minimizing  $IC_m(j)$ , i.e., it is presented as

$$\hat{j}_m = \arg \min_{j \in \mathcal{J}} IC_m(j).$$

Next, we describe a noncentrality matrix that plays a critical role for proving consistency. In fact, asymptotic behaviors of elements or eigen values of a noncentrality matrix are one of important factors that determines whether an information criterion is consistent or not. The noncentrality matrix is defined by

$$\Sigma_*^{-1/2} \Theta_*' \mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j) \mathbf{X}_* \Theta_* \Sigma_*^{-1/2}.$$

In order to decompose the noncentrality matrix, the minimum overspecified model including  $j$  is prepared as

$$j_+ = j \cup j_*, \quad (j \in \mathcal{J}). \tag{2.3}$$

If  $j_*$  is arranged as  $j_* = \{\{j_* \cap j\}, \{j_* \cap j^c\}\}$ ,  $(\mathbf{I}_n - \mathbf{P}_j) \mathbf{X}_* = (\mathbf{O}_{n, k_{j_* \cap j}}, (\mathbf{I}_n - \mathbf{P}_j) \mathbf{X}_{j_* \cap j^c})$  is satisfied, where  $\mathbf{O}_{k,p}$  is a  $k \times p$  zero matrix. It is easy to see that  $\mathbf{X}_{j_* \cap j^c}$  is a full column rank matrix because it is assumed that  $\mathbf{X}$  is the full column rank matrix. Hence, the rank of  $\mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j) \mathbf{X}_*$  is calculated as

$$\text{rank}(\mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j) \mathbf{X}_*) = k_{j_* \cap j^c} = k_{j_+} - k_j \leq k_*, \quad (j \in \mathcal{J}_-).$$

This indicates that the rank of  $\mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j) \mathbf{X}_*$  is independent of  $p$  if  $k_*$  is an independent of  $p$ . Let the rank of the noncentrality matrix be denoted by  $\gamma_j$ . It follows from the inequality  $\text{rank}(\Theta_* \Sigma_*^{-1} \Theta_*') \leq \min\{p, k_*\}$  and a knowledge of an elementary linear algebra that

$$\gamma_j \leq \min\{\text{rank}(\mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j) \mathbf{X}_*), \text{rank}(\Theta_* \Sigma_*^{-1} \Theta_*')\} \leq \min\{p, k_{j_+} - k_j\}.$$

It notes that  $\gamma_j = k_{j_+} - k_j$  if  $\Theta_*$  is a full row rank matrix. Since the noncentrality matrix is a positive semidefinite matrix, and its rank is  $\gamma_j$ , it is decomposed as

$$\Sigma_*^{-1/2} \Theta_*' \mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j) \mathbf{X}_* \Theta_* \Sigma_*^{-1/2} = \Gamma_j \Gamma_j', \tag{2.4}$$

where  $\Gamma_j$  is a  $p \times \gamma_j$  matrix.  $\Gamma_j$  is a full column rank matrix in the case of large  $p$ , at least  $p \geq k_*$ . If we assume that the orders of elements of  $\mathbf{X}' \mathbf{X}$  are  $O(n)$  and elements of  $\Theta_*$  and  $\Sigma_*$  are independent of  $n$ , which are common assumptions in papers dealing with an asymptotic theory on the regression model [14, 17],

the orders of elements of  $\mathbf{\Gamma}_j \mathbf{\Gamma}'_j$  are  $O(n)$ . Let  $(\mathbf{A})_{ab}$  denote the  $(a, b)$ th element of a matrix  $\mathbf{A}$ . Notice that

$$\sum_{a=1}^p (\mathbf{\Gamma}_j \mathbf{\Gamma}'_j)_{aa} = \text{tr}(\mathbf{\Gamma}_j \mathbf{\Gamma}'_j) = \text{tr}(\mathbf{\Gamma}'_j \mathbf{\Gamma}_j) = \sum_{a=1}^{\gamma_j} (\mathbf{\Gamma}'_j \mathbf{\Gamma}_j)_{aa}.$$

Hence, if we assume that all the orders of the elements of  $\mathbf{\Gamma}_j \mathbf{\Gamma}'_j$  are  $O(n)$ , all the orders of the elements of  $\mathbf{\Gamma}'_j \mathbf{\Gamma}_j$  are uniformly equal, and  $\gamma_j$  is constant, then all the orders of the elements of  $\mathbf{\Gamma}'_j \mathbf{\Gamma}_j$  are  $O(np)$ . From this fact and the inequality  $\{(\mathbf{\Gamma}'_j \mathbf{\Gamma}_j)_{ab}\}^2 \leq (\mathbf{\Gamma}'_j \mathbf{\Gamma}_j)_{aa} (\mathbf{\Gamma}'_j \mathbf{\Gamma}_j)_{bb}$ ,  $(\mathbf{\Gamma}'_j \mathbf{\Gamma}_j)_{ab} = O(np)$  is obtained. Consequently, it is natural to assume that  $\mathbf{\Gamma}'_j \mathbf{\Gamma}_j = O(np)$  when  $\mathbf{X}' \mathbf{X} = O(n)$  is assumed.

Let  $\lambda_{j,1} \geq \dots \geq \lambda_{j,\gamma_j} > 0$  be eigen values of  $\mathbf{\Gamma}'_j \mathbf{\Gamma}_j$ . In order to evaluate the probability of selecting the model  $j$  by the  $\text{IC}_m$ , we introduce the following assumptions:

Assumption 1. The true model is included in the set of candidate models, i.e.,

$$j_* \in \mathcal{J}.$$

Assumption 2.  $\lim_{n \rightarrow \infty} n^{-1} \mathbf{X}' \mathbf{X} = \mathbf{R}_0$  exists and is positive definite, and

$$\lim_{n \rightarrow \infty} n^{-1} \mathbf{\Gamma}_j \mathbf{\Gamma}'_j = \mathbf{\Psi}_{j,0}$$
 exists and is not the zero matrix for all  $j \in \mathcal{J}_-$ .

Assumption 3. For all  $j \in \mathcal{J}_-$ ,  $\gamma_j$  is constant, and  $\limsup_{c_n, p \rightarrow c_0} (np)^{-1} \lambda_{j,1} < \infty$  and  $\liminf_{c_n, p \rightarrow c_0} (np)^{-1} \lambda_{j,\gamma_j} > 0$ .

For  $\mathbf{R}_0$  in assumption 2, we write a limiting value of  $n^{-1} \mathbf{X}'_j \mathbf{X}_\ell$  as  $\mathbf{R}_{j,\ell,0}$  for  $j, \ell \in \mathcal{J}$ . It is clear that  $\mathbf{R}_{j,\ell,0}$  is a submatrix of  $\mathbf{R}_0$ , and  $\mathbf{R}_{j,\ell,0}$  also exists if  $\mathbf{R}_0$  exists. Moreover, it notes that  $\mathbf{\Psi}_{j,0}$  still depends on  $p$  because  $\mathbf{\Psi}_{j,0}$  is the convergent value under the LS asymptotic framework.

### 3. Main results

In this section, we evaluate an asymptotic probability of selecting a model by the  $\text{IC}_m$  in (2.1). First, we describe the asymptotic selection probabilities of selecting the true model  $j_*$  under the ordinary asymptotic framework, i.e., the LS asymptotic framework. Using the ideas of [11, 12, 22, 28], we obtain the following Theorem 3.1 (the proof is given in Appendix A.1):

**Theorem 3.1.** *Suppose that assumptions 1 and 2 hold. A variable selection method based on the  $\text{IC}_m$  is consistent when  $n \rightarrow \infty$  if the following conditions are satisfied simultaneously:*

C1-1. For all  $j \in \mathcal{J}_-$ ,

$$\lim_{n \rightarrow \infty} \frac{m(j) - m(j_*)}{n} = 0.$$

C1-2. For all  $j \in \mathcal{J}_+ \setminus \{j_*\}$ ,

$$\lim_{n \rightarrow \infty} \{m(j) - m(j_*)\} = \infty.$$

If one of the above two conditions is not satisfied, a variable selection method based on the  $IC_m$  is not consistent when  $n \rightarrow \infty$ . Additionally, when  $m(j) = O(1)$  as  $n \rightarrow \infty$  and  $\lim_{n \rightarrow \infty} \{m(j) - m(\ell)\} = pm_0(k_j - k_\ell)$  for all  $j, \ell \in \mathcal{J}_+$ , the asymptotic probability of selecting the model  $j$  by the  $IC_m$  is

$$\begin{aligned} & \lim_{n \rightarrow \infty} P(\hat{j}_m = j) \\ &= \begin{cases} 0 & (j \in \mathcal{J}_-) \\ P(\cap_{\ell \in \mathcal{J}_+ \setminus \{j\}} (\mathbf{z}'_\ell \mathbf{z}_\ell - \mathbf{z}'_j \mathbf{z}_j) < m_0 p (k_\ell - k_j)) & (j \in \mathcal{J}_+) \end{cases}, \end{aligned} \tag{3.1}$$

where  $\mathbf{z}_j \sim N_{k_j p}(\mathbf{0}_{k_j p}, \mathbf{I}_{k_j p})$ ,  $Cov[\mathbf{z}_j, \mathbf{z}_\ell] = \mathbf{I}_p \otimes \mathbf{R}_{j,j,0}^{-1/2} \mathbf{R}_{j,\ell,0} \mathbf{R}_{\ell,\ell,0}^{-1/2}$ , and  $\mathbf{0}_p$  is the  $p$ -dimensional zero vector.

These results include the results of [22, 35] etc. as a special case for  $p = 1$ . Theorem 3.1 points out a well-known fact that, when  $n \rightarrow \infty$ , the AIC and the  $AIC_c$  are not consistent and the BIC and the CAIC are consistent in model selection. However, when behaviors of the information criteria are evaluated under the HD framework, we obtain new properties, as in Theorem 3.2 (the proof is given in Appendix A.2).

**Theorem 3.2.** *Suppose that assumptions 1 and 3 are satisfied. Then, a variable selection method based on the  $IC_m$  is consistent when  $c_{n,p} \rightarrow c_0$  if the following conditions are satisfied simultaneously:*

C2-1. For all  $j \in \mathcal{J}_-$ ,

$$\lim_{c_{n,p} \rightarrow c_0} \frac{m(j) - m(j_*)}{n \log p} > -\gamma_j.$$

C2-2. For all  $j \in \mathcal{J}_+ \setminus \{j_*\}$ ,

$$\lim_{c_{n,p} \rightarrow c_0} \frac{m(j) - m(j_*)}{p} > -\frac{1}{c_0} (k_j - k_*) \log(1 - c_0).$$

If the sign “ $>$ ” becomes “ $<$ ” in one of the above two conditions, a variable selection method based on the  $IC_m$  is not consistent when  $c_{n,p} \rightarrow c_0$ .

It notes that  $\lim_{c \rightarrow 0} c^{-1} \log(1 - c) = -1$  and  $c^{-1} \log(1 - c)$  is a monotonically decreasing function in  $0 \leq c < 1$ . From Theorem 3.2, consistency properties of specific criteria are clarified as the following corollary (the proof is given in Appendix A.4):

**Corollary 3.1.** *Suppose that assumptions 1 and 3 are satisfied.*

(i) *When  $c_{n,p} \rightarrow c_0$ , a variable selection method based on the AIC is consistent if  $c_0 \in [0, c_a)$ , and is not consistent if  $c_0 \in (c_a, 1)$ , where  $c_a$  ( $\approx 0.797$ ) is a constant satisfying*

$$\log(1 - c_a) + 2c_a = 0. \tag{3.2}$$

(ii) *When  $c_{n,p} \rightarrow c_0$ , a variable selection method based on the  $AIC_c$  is consistent*

- (iii) When  $c_{n,p} \rightarrow c_0$ , variable selection methods based on the BIC and the CAIC are consistent if  $c_0 \in [0, c_b)$ , and are not consistent if  $c_0 \in (c_b, 1)$ , where  $c_b = \min\{1, \min_{j \in \mathcal{S}_-} \gamma_j / (k_* - k_j)\}$  and  $\mathcal{S}_- = \{j \in \mathcal{J}_- | k_* - k_j > 0\}$ .

Corollary 3.1 shows that, there is no restriction of  $c_0$  in the condition for consistency of the  $\text{AIC}_c$  although it is restricted in the AIC. This indicates that it is possible that the bias correction to the risk function has a positive effect on selection of the true model. Moreover, Corollary 3.1 indicates that the BIC and the CAIC are not always consistent in variable selection when  $c_{n,p} \rightarrow c_0$ . If  $\Theta_*$  is the full row rank matrix,  $\gamma_j$  becomes  $k_{j+} - k_j$ . Since  $c_0 < 1$  and  $k_{j+} - k_j > k_* - k_j$  for all  $j \in \mathcal{S}_-$ ,  $\gamma_j > c_0(k_* - k_j)$  is satisfied if  $\Theta_*$  is the full row rank matrix. In contrast, if  $c_0 = 0$  then  $\gamma_j > c_0(k_* - k_j)$  is satisfied. Therefore, we can see that variable selection methods based on the BIC and the CAIC are consistent as  $c_{n,p} \rightarrow c_0$  if  $\Theta_*$  is the full row rank matrix, or  $c_{n,p}$  converges to 0. However, if  $\Theta_*$  is not the full row rank matrix and  $c_0 \in (0, 1)$ , we cannot determine as if variable selection methods based on the BIC and the CAIC are consistent as  $c_{n,p} \rightarrow c_0$ .

In order to clarify the condition to ensure inconsistency, assumption 3 is assumed in Theorem 3.2, i.e., we assume that the orders of eigen values of  $\Gamma_j' \Gamma_j$  are uniformly the same and  $\gamma_j$  is independent of  $n$  and  $p$ . If the aim is only to derive a sufficient condition for consistency, such a strong assumption like assumption 3 is unnecessary. In fact, for evaluating consistency, there is no need to assume the same orders for all the eigen values of  $\Gamma_j' \Gamma_j$ . Whether an information criterion is consistent strongly depends on the orders of divergence speeds of several eigen values. Hence, we clarify condition of the orders of divergence speeds of eigen values of  $\Gamma_j' \Gamma_j$  to ensure a consistency. In addition, most recently, many researchers pay close attention to “big data analysis”, and thus study on a theory of a variable selection when the number of candidate models approaches  $\infty$  (see, e.g., [19]). Hence, we derive a sufficient condition for the consistency by using the following asymptotic framework:

- A high-dimensional and large-model (HD-LM) asymptotic framework: the HD asymptotic framework under the condition that the following equations are satisfied:

$$\max_{j \in \mathcal{J}} \frac{k_j}{n} \rightarrow 0, \quad \exists l > 0 \text{ s.t. } K = o(p^l), \quad (3.3)$$

where  $K$  is the number of candidate models. In the HD-LM asymptotic framework,  $p$  always goes to  $\infty$ , and it makes no difference whether  $K$  is constant or  $K$  goes to  $\infty$ . This indicates that the HD asymptotic framework is a special case of the HD-LM asymptotic framework. For simplicity, we will write “ $(n, p) \rightarrow \infty$  simultaneously under the HD-LM asymptotic framework” as “ $c_{n,p} \rightarrow c_0$  under LM”, and  $\lim_{c_{n,p} \rightarrow c_0, \text{LM}}$  and  $\liminf_{c_{n,p} \rightarrow c_0, \text{LM}}$  mean a limit and a limit inferior under the HD-LM asymptotic framework, respectively.

**Theorem 3.3.** *Suppose that assumption 1 holds. A variable selection method based on the  $\text{IC}_m$  is consistent under the HD-LM asymptotic framework if the following conditions are satisfied:*



C3-1. For sufficiently large  $n$ , there exist positive constants  $\delta_1, \delta_2$  such that for all  $j \in \mathcal{J}_-$  there exists an integer  $q \in [1, \gamma_j]$  such that  $\lambda_{j,q}/q^2 > n^{\delta_1}$  and

$$\log \frac{\beta_{j,q}}{n - p - k_{j_+} + (q + 1)/2} - \frac{1}{qn} \{m(j_+) - m(j)\} > \delta_2,$$

where  $j_+$  is given by (2.3),  $\gamma_j = \text{rank}(\mathbf{\Gamma}_j)$ , and  $\beta_{j,q}$  is the geometric mean of the largest  $q$  eigen values of  $\mathbf{\Gamma}_j \mathbf{\Gamma}'_j$ , i.e.,

$$\beta_{j,q} = (\prod_{i=1}^q \lambda_{j,i})^{1/q}. \tag{3.4}$$

C3-2. For sufficiently large  $n$ , there exists a positive constant  $\delta$  such that for all  $j \in \mathcal{J}_+ \setminus \{j_*\}$ ,

$$\frac{m(j) - m(j_*)}{p(k_j - k_*)} + \frac{1}{c_{n,p}} \log(1 - c_{n,p}) > \delta.$$

The proof is given in Appendix A.5. Roughly speaking, the existence of  $\delta_2$  in condition C3-1 is related to the order of noncentrality matrix. Assumption 3 is equivalent to the condition that the orders of all the eigen values of  $\mathbf{\Gamma}'_j \mathbf{\Gamma}_j$  are  $O(np)$ . However, the condition  $\lambda_{j,q}/q^2 > n^{\delta_1}$  indicates that the orders of the eigen values of  $\mathbf{\Gamma}'_j \mathbf{\Gamma}_j$  do not need to be the same orders uniformly. Moreover, in Theorem 3.3,  $k_*$ ,  $\gamma_j$  and  $K$  do not have to be bounded. Hence, Theorem 3.3 can be applied to more non-restrictive situations than Theorem 3.2.

Although we have derived sufficient conditions for consistency in Theorem 3.3, it is hard to check from the conditions whether an information criterion considered is consistent. Hence, in order to establish an easy-to-understand formula, we rewrite the conditions by using a limit inferior. Besides, by using 1 as  $q$ , we simplify condition C3-1 although the sufficient conditions become restrictive.

**Corollary 3.2.** *Suppose that assumption 1 holds. A variable selection method based on the  $IC_m$  is consistent under the HD-LM asymptotic framework if the following conditions are satisfied:*

C3-1'.  $\inf_{j \in \mathcal{J}_-} \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} \frac{\log \lambda_{j,1}}{\log n} > 0$ , and

$$\inf_{j \in \mathcal{J}_-} \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} \left\{ \log \frac{\lambda_{j,1}}{n} - \frac{m(j_+) - m(j)}{n} \right\} > \log(1 - c_0).$$

C3-2.  $\inf_{j \in \mathcal{J}_+ \setminus \{j_*\}} \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} \frac{m(j) - m(j_*)}{p(k_j - k_*)} > -\frac{1}{c_0} \log(1 - c_0)$ .

The proof is given in Appendix A.9. It notes that not “min” but “inf” is used for conditions C3-1' and -2 because the number of candidate models may go to  $\infty$ . Although we cannot check whether condition C3-1' is satisfied from an actual data, we can derive the order of the divergence speed of the maximum eigen value of  $\mathbf{\Gamma}'_j \mathbf{\Gamma}_j$  to ensure a consistency. If the size of the order is small, we can consider that a possibility that an information criterion is consistent is high.

Hence, the size of the order helps to assess a quality of an information criterion in the sense of a possibility to have consistency.

Even if  $k_*$  is not bounded, Theorem 3.3 and Corollary 3.2 hold. However, in order to clarify the sufficient conditions to ensure consistency, we consider the simple case that  $k_*$  is bounded. Then, conditions to satisfy consistency properties of specific criteria are simplified as the following corollary (the proof is given in Appendix A.10):

**Corollary 3.3.** *Suppose that assumption 1 is satisfied, and  $k_*$  is bounded.*

- (i) *A variable selection method based on the AIC is consistent when  $c_{n,p} \rightarrow c_0$  under LM if  $c_0 \in [0, c_a)$  and*

$$\inf_{j \in \mathcal{J}_-} \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} \log \frac{\lambda_{j,1}}{n} > \log(1 - c_0) + 2k_*c_0, \quad (3.5)$$

where  $c_a$  is the constant given by (3.2).

- (ii) *A variable selection method based on the  $\text{AIC}_c$  is consistent when  $c_{n,p} \rightarrow c_0$  under LM if*

$$\inf_{j \in \mathcal{J}_-} \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} \log \frac{\lambda_{j,1}}{n} > \log(1 - c_0) + k_*c_0 \left\{ \frac{1}{1 - c_0} + \frac{1}{(1 - c_0)^2} \right\}. \quad (3.6)$$

- (iii) *Variable selection methods based on the BIC and the CAIC are consistent when  $c_{n,p} \rightarrow c_0$  under LM if*

$$\inf_{j \in \mathcal{J}_-} \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} \frac{\log(\lambda_{j,1}/n)}{\log n} > k_*c_0. \quad (3.7)$$

An example of the noncentrality matrix is shown in Appendix A.11. From Corollary 3.3, we can see that the  $\text{AIC}_c$  is consistent if  $\lim_{c_{n,p} \rightarrow c_0, \text{LM}} \log(\lambda_{j,1}/n) = \infty$ , and the BIC and the CAIC are consistent if  $\lim_{c_{n,p} \rightarrow c_0, \text{LM}} \log(\lambda_{j,1}/n)/\log n = \infty$ . Moreover, the AIC is consistent if  $c_0 < c_a$  and  $\lim_{c_{n,p} \rightarrow c_0, \text{LM}} \log(\lambda_{j,1}/n) = \infty$ . Hence, the AIC and the  $\text{AIC}_c$  has a superiority over the BIC and the CAIC in the sense of a possibility to have a consistency. Moreover, although the AIC is consistent under the restriction  $c_0 < c_a$ , there is no such a restriction in  $\text{AIC}_c$ . Consequently, we can judge that the  $\text{AIC}_c$  has a superiority over the AIC, BIC and CAIC in the sense of a possibility to have a consistency.

#### 4. Numerical study

In this section, we compare with the probabilities of selecting the true model by AIC,  $\text{AIC}_c$ , BIC and CAIC in (2.2), which were evaluated by Monte Carlo simulations based on 10,000 replications under several different values of  $n$  and  $p$ . A set of candidate models was  $\mathcal{J} = \{j_1, \dots, j_k\}$ , where  $j_\alpha = \{1, \dots, \alpha\}$  ( $\alpha = 1, \dots, k$ ). A  $1000 \times 156$  matrix  $\mathbf{M}\Phi(156)^{1/2}$  was generated, where an each element of  $\mathbf{M}$  was independent and identically chosen from  $U(-1, 1)$ , and  $\Phi(q)$  is a  $q \times q$  symmetric matrix whose the  $(a, b)$ th element was defined by  $(0.8)^{|a-b|}$ . Using this matrix, we constructed an  $n \times k$  matrix of explanatory variables  $\mathbf{X}$  as

TABLE 1  
Selection probabilities of the true model (%)

Case 1						Case 2 ( $c_0 = 0.02$ )						
$n$	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC
100	2	10	75.8	84.5	98.7	99.5	2	10	75.8	84.5	98.7	99.5
200	2	10	79.0	83.2	99.5	99.9	4	10	86.0	90.3	99.9	100.0
500	2	10	79.4	81.2	99.8	99.9	10	10	96.3	97.4	100.0	100.0
1000	2	10	79.1	79.9	99.9	100.0	20	10	99.3	99.6	100.0	100.0
$\infty$	2	10	80.2	80.2	100.0	100.0	$\infty$	10	100.0	100.0	100.0	100.0
Case 3						Case 4 ( $c_0 = 0.1$ )						
$n$	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC
100	10	10	92.9	99.3	100.0	100.0	10	10	92.9	99.3	100.0	100.0
200	10	10	95.1	98.3	100.0	100.0	20	10	98.6	99.9	100.0	100.0
500	10	10	96.3	97.5	100.0	100.0	50	10	100.0	100.0	100.0	100.0
1000	10	10	96.5	97.1	100.0	100.0	100	10	100.0	100.0	100.0	100.0
$\infty$	10	10	96.8	96.8	100.0	100.0	$\infty$	10	100.0	100.0	100.0	100.0
Case 5						Case 6 ( $c_0 = 0.3$ )						
$n$	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC
100	30	10	97.0	40.1	0.0	0.0	30	10	97.0	40.1	0.0	0.0
200	30	10	99.4	100.0	100.0	55.8	60	10	99.8	100.0	0.0	0.0
500	30	10	99.7	100.0	100.0	100.0	150	10	100.0	100.0	0.0	0.0
1000	30	10	99.8	99.9	100.0	100.0	300	10	100.0	100.0	0.0	0.0
$\infty$	30	10	99.9	99.9	100.0	100.0	$\infty$	10	100.0	100.0	0.0	0.0
Case 7 ( $c_0 = 0.0$ )						Case 8 ( $c_0 = 0.0$ )						
$n$	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC
100	30	10	97.0	40.1	0.0	0.0	30	10	97.0	40.1	0.0	0.0
200	32	10	99.4	100.0	100.0	15.8	40	10	99.8	100.0	59.5	0.0
500	35	10	99.9	100.0	100.0	100.0	50	10	100.0	100.0	100.0	100.0
1000	40	10	100.0	100.0	100.0	100.0	60	10	100.0	100.0	100.0	100.0
$\infty$	$\infty$	10	100.0	100.0	100.0	100.0	$\infty$	10	100.0	100.0	100.0	100.0
Case 9 ( $c_0 = 0.1$ )						Case 10 ( $c_0 = 0.1$ )						
$n$	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC
100	10	10	92.9	99.3	100.0	100.0	10	10	92.9	99.3	100.0	100.0
200	20	14	98.8	99.9	100.0	100.0	20	31	98.7	100.0	100.0	100.0
500	50	22	100.0	100.0	100.0	100.0	50	84	100.0	100.0	100.0	100.0
1000	100	31	100.0	100.0	100.0	100.0	100	156	100.0	100.0	100.0	100.0
$\infty$	$\infty$	$\infty$	100.0	100.0	—	—	$\infty$	$\infty$	100.0	100.0	—	—
Case 11 ( $c_0 = 0.3$ )						Case 12 ( $c_0 = 0.3$ )						
$n$	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC	$p$	$k$	AIC	AIC <sub>c</sub>	BIC	CAIC
100	30	10	97.0	40.1	0.0	0.0	30	10	97.0	40.1	0.0	0.0
200	60	14	99.8	100.0	0.0	0.0	60	31	99.7	100.0	0.0	0.0
500	150	22	100.0	100.0	0.0	0.0	150	84	100.0	100.0	0.0	0.0
1000	300	31	100.0	100.0	0.0	0.0	300	156	100.0	100.0	0.0	0.0
$\infty$	$\infty$	$\infty$	100.0	100.0	—	—	$\infty$	$\infty$	100.0	100.0	—	—

a submatrix of  $\mathbf{M}\Phi(156)^{1/2}$  from rows 1 to  $n$  and columns 1 to  $k$ . The true model was determined by  $\Theta_* = (1, 1, 3, -4, 5)' \mathbf{1}'_p$ ,  $j_* = \{1, 2, 3, 4, 5\}$ , and  $\Sigma_* = \Phi(p)$ , where  $\mathbf{1}_p$  was the  $p$ -dimensional vector of ones. Thus,  $j_\alpha$  with  $\alpha = 1, \dots, 4$  was the underspecified model, and  $j_\alpha$  with  $\alpha \geq 5$  was the overspecified model.

In our numerical study,  $\gamma_j = 1$  and  $\max_{j \in \mathcal{S}_-} (k_* - k_j) = 4$  hold. This implies that when  $c_0 > 1/4$ , the inequality  $\gamma_j > c_0(k_* - k_j)$  was not always satisfied for all  $j \in \mathcal{S}_-$ . Thus, the BIC and the CAIC were not consistent in variable selection when  $c_0 > 1/4$  under the fixed  $k$ .

Table 1 shows the probability of selecting the true model by the AIC, AIC<sub>c</sub>, BIC, and CAIC. For  $n = \infty$  or  $p = \infty$ , we list the theoretical values obtained

from Theorems 3.1, 3.2 and 3.3. A symbol “—” means that theoretical values are unclear because the sufficient condition for consistency does not hold. In the table, Cases 1, 3, and 5 are the results when  $n \rightarrow \infty$  under fixed  $p$  and  $k = 10$ , and Cases 2, 4, 6, 7, and 8 are the results when  $(n, p) \rightarrow \infty$  under a fixed  $k = 10$  and with  $c_0 = 0.02, 0.1, 0.3, 0.0$ , and  $0.0$ . Moreover, Cases 9, 10, 11 and 12 are the results when  $(n, p) \rightarrow \infty$  and with  $c_0 = 0.1$  and  $0.3$ , and  $k = 10 + [n^{1/2} - 10^{1/2}]$  and  $k = 10 + [n^{3/4} - 10^{3/4}]$ , where  $[ \ ]$  is the Gauss' symbol. From the table, we can see that in the cases of the AIC and the AIC<sub>c</sub>, the greater the dimension and sample size considered, the greater the probabilities became. Compared with the results obtained from the AIC and the AIC<sub>c</sub>, probabilities by the AIC<sub>c</sub> tended to be higher than those by the AIC when  $n$  was not small. In the cases of the BIC and the CAIC, the greater the dimension and sample size considered, the higher the selection probabilities became, with the exception of Case 6. This was because variable selection methods based on the BIC and the CAIC were not consistent in Case 6. Additionally, when  $n$  was small and  $p$  was large, the selection probabilities of the BIC and the CAIC were both very low. However, if the BIC and the CAIC were consistent in variable selection, these probabilities became high as  $n$  and  $p$  increased. Moreover, we can see that above tendencies were satisfied even if the number of explanatory variables becomes large.

We simulated several other models and obtained similar results. Since the theoretical difference between using the AIC and the AIC<sub>c</sub> occurs when  $c_{n,p} > 0.8$ , we should list the numerical results for such a case. However, when  $c_{n,p}$  is close to 1, the convergence of selection probabilities was extremely slow. Thus, we do not show simulation results for dimensions close to the sample size.

## 5. Conclusion and discussion

In this paper, we demonstrated that there is the case that the AIC for the multivariate linear regression model is consistent in variable selection when we approximate the probability of selecting the true model using the HD asymptotic framework. Although the AIC becomes consistent under the restriction  $c_0 < c_a$ , the AIC<sub>c</sub> becomes consistent without the restriction of  $c_0$ . This indicates that it is possible that correcting the bias to the risk function may have a positive effect on the selection of the true model. It is a well-known fact that variable selection methods based on the BIC and the CAIC are consistent if we approximate the probability of selecting the true model using the LS asymptotic framework. However, we found that there is a possibility that the BIC and the CAIC become inconsistent if we approximate the probability of selecting the true model using the HD asymptotic framework.

It is known that the LS asymptotic theory gives a poor approximation when the dimension is large. The HD asymptotic theory gives a better approximation than the LS asymptotic theory when the sample size and the dimension are large, and sometimes even when the dimension is not so large. Hence, the consistency property of the AIC that we demonstrated will be useful for high-dimensional data analysis. Usually, the HD asymptotic theory is used to improve

the approximations of the distributions of statistics. However, the results in this paper suggest a possibility that new insight can be provided by applying the HD asymptotic theory to high-dimensional data.

From the simulation study, we found that, the larger the dimension and sample size considered, the higher the selection probabilities became. This numerical result naturally implies that using multiple response variables at the same time as the model selection can increase the probability of selecting the true model. In other words, we should not select variables using only each response variable. That is a strong reason to apply the model selection procedure based on the multivariate linear regression model to high-dimensional data.

In this paper, we considered the case of  $n > p$  because  $\hat{\Sigma}_j$  becomes singular when  $p > n$ . Unfortunately,  $n > p$  is not always satisfied in the actual data. If our results can be extended to the case of  $n \leq p$ , we clarify the conditions to satisfy consistency property in many infinite-dimensional statistics, e.g., the time series analysis (see [4, 5]), spatiotemporal geostatistical analysis (see [7, 8]) and functional data analysis (see [23]). The singularity of  $\hat{\Sigma}_j$  can be avoided by using a ridge-type estimator of the covariance matrix, as demonstrated by [36]. We can expect that an AIC consisting of such a ridge-type estimator will be consistent in model selection.

## Appendix

### A.1. The proof of Theorem 3.1

Recall that the LS asymptotic framework is used for proving Theorem 3.1. We can see that  $\hat{\Sigma}_j \xrightarrow{P} \Sigma_*$  as  $n \rightarrow \infty$  holds when  $j \in \mathcal{J}_+$  and  $\hat{\Sigma}_j \xrightarrow{P} \Sigma_*^{1/2} \Psi_{j,0} \Sigma_*^{1/2} + \Sigma_*$  as  $n \rightarrow \infty$  holds when  $j \in \mathcal{J}_-$ , where  $\Psi_{j,0} = \lim_{n \rightarrow \infty} n^{-1} \Gamma_j \Gamma_j'$  and  $\Gamma_j$  is given by (2.4). Notice that  $\Psi_{j,0}$  is a positive semidefinite matrix. When  $\lim_{n \rightarrow \infty} \{m(j) - m(j_*)\}/n = 0$  for all  $j \in \mathcal{J}_-$ , we have

$$\begin{aligned} \frac{1}{n} \{IC_m(j) - IC_m(j_*)\} &\xrightarrow{P} \log \det(\Sigma_*^{1/2} \Psi_{j,0} \Sigma_*^{1/2} + \Sigma_*) - \log \det(\Sigma_*) \\ &= \log \det(\mathbf{I}_p + \Psi_{j,0}) > 0. \end{aligned}$$

This result implies that  $\lim_{n \rightarrow \infty} P(IC_m(j_*) > IC_m(j)) = 0$  for any  $j \in \mathcal{J}_-$ . Thus, we obtain

$$\begin{aligned} &\lim_{n \rightarrow \infty} P(\hat{j}_m = j) \\ &= \begin{cases} 0 & (j \in \mathcal{J}_-) \\ \lim_{n \rightarrow \infty} P(\cap_{\ell \in \mathcal{J}_+ \setminus \{j\}} \{IC_m(\ell) > IC_m(j)\}) & (j \in \mathcal{J}_+) \end{cases} \quad (\text{A.1}) \end{aligned}$$

From here to the end of proof, we assume  $j, \ell \in \mathcal{J}_+$ . Let  $\mathbf{V}$  and  $\mathbf{Z}_j$  be the  $p \times p$  and the  $k_j \times p$  matrices defined by

$$\mathbf{V} = \frac{1}{\sqrt{n}}(\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon} - n\mathbf{I}_p), \quad \mathbf{Z}_j = (\mathbf{X}'_j \mathbf{X}_j)^{-1/2} \mathbf{X}'_j \boldsymbol{\varepsilon},$$

where

$$\mathbf{E} = (\mathbf{Y} - \mathbf{X}_* \boldsymbol{\Theta}_*) \boldsymbol{\Sigma}_*^{-1/2}. \quad (\text{A.2})$$

It is well known that  $\mathbf{V}$  has an asymptotic normality as  $n \rightarrow \infty$ , and  $\mathbf{Z}_j \sim N_{k_j \times p}(\mathbf{O}_{k_j, p}, \mathbf{I}_{k_j p})$ . Furthermore, using

$$\boldsymbol{\Sigma}_*^{-1/2} \hat{\boldsymbol{\Sigma}}_j \boldsymbol{\Sigma}_*^{-1/2} = \frac{1}{n} \mathbf{E}' (\mathbf{I}_n - \mathbf{P}_j) \mathbf{E} = \frac{1}{n} (\mathbf{E}' \mathbf{E} - \mathbf{Z}_j' \mathbf{Z}_j),$$

we have

$$\boldsymbol{\Sigma}_*^{-1/2} \hat{\boldsymbol{\Sigma}}_j \boldsymbol{\Sigma}_*^{-1/2} = \mathbf{I}_p + \frac{1}{\sqrt{n}} \mathbf{V} - \frac{1}{n} \mathbf{Z}_j' \mathbf{Z}_j.$$

From the above expression, the first term of the  $\text{IC}_m(j)$  can be expanded as

$$\mathcal{L}(j) = n \log \det(\boldsymbol{\Sigma}_*) + \sqrt{n} \text{tr}(\mathbf{V}) - \{\text{tr}(\mathbf{V}^2)/2 + \text{tr}(\mathbf{Z}_j' \mathbf{Z}_j)\} + O_p(n^{-1/2}).$$

Let  $\mathbf{z}_j$  be a  $k_j p$ -dimensional random vector defined by  $\mathbf{z}_j = \text{vec}(\mathbf{Z}_j)$ , where  $\text{vec}(\mathbf{A})$  is an operator that transforms a matrix to a vector by stacking the first to the last columns of  $\mathbf{A}$ , i.e.,  $\text{vec}(\mathbf{A}) = (\mathbf{a}'_1, \dots, \mathbf{a}'_m)'$  when  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_m)$  (see, e.g., [18, chap. 16.2]). Then, it follows from the expansion and the equality  $\text{tr}(\mathbf{Z}_j' \mathbf{Z}_j) = \mathbf{z}_j' \mathbf{z}_j$  that

$$\text{IC}_m(\ell) - \text{IC}_m(j) = -(\mathbf{z}'_\ell \mathbf{z}_\ell - \mathbf{z}'_j \mathbf{z}_j) + m(\ell) - m(j) + O_p(n^{-1/2}). \quad (\text{A.3})$$

Hence, when  $\lim_{n \rightarrow \infty} \{m(j) - m(j_*)\} = \infty$  holds for all  $j \in \mathcal{J}_+ \setminus \{j_*\}$  we derive

$$\frac{1}{m(j) - m(j_*)} \{\text{IC}_m(j) - \text{IC}_m(j_*)\} \xrightarrow{p} 1 > 0. \quad (\text{A.4})$$

On the other hand, when  $\lim_{n \rightarrow \infty} m(j) = m_0 \{pk_j + p(p+1)/2\} < \infty$  holds,  $\lim_{n \rightarrow \infty} \{m(\ell) - m(j)\} = m_0 p(k_\ell - k_j)$  is satisfied. Consequently, by combining this result, and (A.3) and (A.4) with (A.1), Theorem 3.1 is proved.

## A.2. The proof of Theorem 3.2

At first, we describe the lemma which is used for proving Theorems 3.2 and 3.3 (the proof of lemma is given after this subsection).

**Lemma A.1.** *Let  $T = -l(pq)^{-1} \log \Lambda$  where  $\Lambda$  is distributed according to the Wilks' lambda distribution  $\Lambda_q(p, l+q)$ , and let  $\kappa_T^{(s)}$  be the  $s$ th order cumulant of  $T$ . Suppose that  $p/l \rightarrow \alpha$  (constant) and  $q/l \rightarrow 0$ . If  $\alpha > 0$  then*

$$\begin{aligned} \kappa_T^{(1)} &\rightarrow \frac{1}{\alpha} \log(1 + \alpha), \\ (ql)^{s-1} \kappa_T^{(s)} &\rightarrow \frac{2^{s-1}(s-2)!}{c} \left\{ 1 - \left( \frac{1}{1+c} \right)^{s-1} \right\} \quad (s \geq 2). \end{aligned}$$

*If  $\alpha = 0$  then  $\kappa_T^{(1)} \rightarrow 1$  and  $(pq)^{s-1} \kappa_T^{(s)} \rightarrow 2^{s-1}(s-1)!$ . Hence whether  $\alpha = 0$  or  $\alpha > 0$  for any positive integer  $m$  and any positive value  $\delta$ ,*

$$P(|T - \kappa_T^{(1)}| > \delta) \leq \frac{1}{\delta^{2m}} E[(T - \kappa_T^{(1)})^{2m}] = O((pq)^{-m}).$$

Recall that the HD asymptotic framework is used for proving Theorem 3.2. First, we consider the case of  $j \in \mathcal{J}_-$ . Let

$$\mathcal{A}_j = (\mathbf{I}_n - \mathbf{P}_j)\mathbf{X}_*\Theta_*\Sigma_*^{-1/2}.$$

Notice that  $\text{rank}(\mathcal{A}_j) = \gamma_j$  because  $\mathcal{A}'_j\mathcal{A}_j = \Gamma_j\Gamma'_j$  and  $\text{rank}(\Gamma_j) = \gamma_j$ , where  $\Gamma_j$  is given by (2.4). By using a singular value decomposition,  $\mathcal{A}_j$  can be rewritten as

$$\mathcal{A}_j = \mathbf{H}_j\mathbf{L}_j^{1/2}\mathbf{G}'_j,$$

where  $\mathbf{H}_j$  and  $\mathbf{G}_j$  are  $n \times \gamma_j$  and  $p \times \gamma_j$  matrices satisfying  $\mathbf{H}'_j\mathbf{H}_j = \mathbf{I}_{\gamma_j}$  and  $\mathbf{G}'_j\mathbf{G}_j = \mathbf{I}_{\gamma_j}$ , respectively, and  $\mathbf{L}_j$  is a  $\gamma_j \times \gamma_j$  diagonal matrix whose diagonal elements are squared singular values of  $\mathcal{A}_j$ . By using  $\mathcal{A}_j$  and  $\mathcal{E}$  given by (A.2), we have

$$n\Sigma_*^{-1/2}\hat{\Sigma}_j\Sigma_*^{-1/2} = \mathbf{W}_1 + \mathbf{W}_2 + \mathbf{W}_3, \quad n\Sigma_*^{-1/2}\hat{\Sigma}_{j+}\Sigma_*^{-1/2} = \mathbf{W}_1, \quad (\text{A.5})$$

where

$$\begin{aligned} \mathbf{W}_1 &= \mathcal{E}'(\mathbf{I}_n - \mathbf{P}_{j+})\mathcal{E}, & \mathbf{W}_2 &= \mathcal{E}'(\mathbf{P}_{j+} - \mathbf{H}_j\mathbf{H}'_j - \mathbf{P}_j)\mathcal{E}, \\ \mathbf{W}_3 &= (\mathcal{A}_j + \mathcal{E})'\mathbf{H}_j\mathbf{H}'_j(\mathcal{A}_j + \mathcal{E}). \end{aligned}$$

It follows from the equations  $\mathbf{P}_{j*}\mathbf{X}_* = \mathbf{X}_*$  and  $\mathbf{P}_j\mathbf{P}_{j+} = \mathbf{P}_j$  that

$$\begin{aligned} \mathcal{A}'_j\mathbf{P}_{j+} &= \Sigma_*^{-1/2}\Theta'_*\mathbf{X}'_*(\mathbf{P}_{j+} - \mathbf{P}_j) = \Sigma_*^{-1/2}\Theta'_*\mathbf{X}'_*\mathbf{P}_{j+}(\mathbf{I}_n - \mathbf{P}_j) \\ &= \Sigma_*^{-1/2}\Theta'_*\mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j) = \mathcal{A}'_j. \end{aligned}$$

Using this result and  $\mathcal{A}'_j\mathbf{P}_j = \mathbf{O}_{p,n}$  yields  $\mathbf{H}'_j\mathbf{P}_j = \mathbf{O}_{\gamma_j,n}$  and  $\mathbf{H}'_j\mathbf{P}_{j+} = \mathbf{H}'_j$ . These imply that

$$\begin{aligned} \mathbf{H}_j\mathbf{H}'_j(\mathbf{I}_n - \mathbf{P}_{j+}) &= \mathbf{O}_{n,n}, & \mathbf{H}_j\mathbf{H}'_j(\mathbf{P}_{j+} - \mathbf{H}_j\mathbf{H}'_j - \mathbf{P}_j) &= \mathbf{O}_{n,n}, \\ (\mathbf{I}_n - \mathbf{P}_{j+})(\mathbf{P}_{j+} - \mathbf{H}_j\mathbf{H}'_j - \mathbf{P}_j) &= \mathbf{O}_{n,n}, \\ (\mathbf{P}_{j+} - \mathbf{H}_j\mathbf{H}'_j - \mathbf{P}_j)^2 &= \mathbf{P}_{j+} - \mathbf{H}_j\mathbf{H}'_j - \mathbf{P}_j. \end{aligned}$$

From the above results and the multivariate version of the Cochran theorem (see, e.g., [30, chap. 2.8]), we can see that  $\mathbf{W}_1$ ,  $\mathbf{W}_2$ , and  $\mathbf{W}_3$  are  $p \times p$  mutually independent random matrices distributed according to the Wishart or the noncentral Wishart distributions;

$$\mathbf{W}_1 \sim W_p(n - k_{j+}, \mathbf{I}_p), \quad \mathbf{W}_2 \sim W_p(d_j, \mathbf{I}_p), \quad \mathbf{W}_3 \sim W_p(\gamma_j, \mathbf{I}_p; \Gamma_j\Gamma'_j), \quad (\text{A.6})$$

where  $d_j = k_{j+} - k_j - \gamma_j$ . It follows from (A.5) and (A.6), and the property of the Wishart distributions (see [16, p. 57 th. 3.2.4]) that

$$\begin{aligned} \frac{1}{n}\{\mathcal{L}(j) - \mathcal{L}(j+)\} &= \log \frac{\det(\mathbf{W}_1 + \mathbf{W}_2 + \mathbf{W}_3)}{\det(\mathbf{W}_1 + \mathbf{W}_2)} + \log \frac{\det(\mathbf{W}_1 + \mathbf{W}_2)}{\det(\mathbf{W}_1)} \\ &= -\log \frac{\det(\mathbf{U}_1)}{\det(\mathbf{U}_1 + \mathbf{U}_2)} - \log \frac{\det(\mathbf{U}_3)}{\det(\mathbf{U}_3 + \mathbf{U}_4)}, \end{aligned} \quad (\text{A.7})$$

where  $\mathbf{U}_1$ ,  $\mathbf{U}_2$ ,  $\mathbf{U}_3$ , and  $\mathbf{U}_4$  are random matrices distributed according to the Wishart or the noncentral Wishart distributions;

$$\begin{aligned} \mathbf{U}_1 &\sim W_{\gamma_j}(n - k_j - p, \mathbf{I}_{\gamma_j}), & \mathbf{U}_2 &\sim W_{\gamma_j}(p, \mathbf{I}_{\gamma_j}; \mathbf{\Gamma}'_j \mathbf{\Gamma}_j), \\ \mathbf{U}_3 &\sim W_{d_j}(n - k_j - \gamma_j - p, \mathbf{I}_{d_j}), & \mathbf{U}_4 &\sim W_{d_j}(p, \mathbf{I}_{d_j}). \end{aligned} \quad (\text{A.8})$$

Here,  $\mathbf{U}_1$  and  $\mathbf{U}_2$  are mutually independent, and  $\mathbf{U}_3$  and  $\mathbf{U}_4$  are also mutually independent. When  $c_{n,p} \rightarrow c_0 \in [0, 1)$ , we have

$$\frac{1}{n - k_j - p} \mathbf{U}_1 \xrightarrow{p} \mathbf{I}_{\gamma_j}, \quad \frac{1}{n - k_j - \gamma_j - p} \mathbf{U}_3 \xrightarrow{p} \mathbf{I}_{d_j}, \quad \frac{1}{p} \mathbf{U}_4 \xrightarrow{p} \mathbf{I}_{d_j}. \quad (\text{A.9})$$

From the definition of the noncentral Wishart distribution, a different expression of  $\mathbf{U}_2$  is given as  $\mathbf{U}_2 = (\mathbf{Z} + \mathbf{\Gamma}_j)'(\mathbf{Z} + \mathbf{\Gamma}_j)$ , where  $\mathbf{Z} \sim N_{p \times \gamma_j}(\mathbf{O}_{p, \gamma_j}, \mathbf{I}_{p \gamma_j})$ . Let  $\mathbf{\Gamma}'_j \mathbf{\Gamma}_j)^{1/2} = \mathbf{\Delta}_j$ . Then, we have

$$\mathbf{\Delta}_j^{-1} \mathbf{U}_2 \mathbf{\Delta}_j^{-1} = \mathbf{\Delta}_j^{-1} (\mathbf{Z}' \mathbf{Z} + \mathbf{\Gamma}'_j \mathbf{Z} + \mathbf{Z}' \mathbf{\Gamma}_j + \mathbf{\Delta}_j^2) \mathbf{\Delta}_j^{-1}. \quad (\text{A.10})$$

Recall that that  $\limsup_{c_{n,p} \rightarrow c_0} np\lambda_{j,\gamma}^{-1} < \infty$  is derived from assumption 3. It follows from the above result and  $E[\mathbf{Z}' \mathbf{Z}] = p\mathbf{I}_{\gamma_j}$  that

$$E[\text{tr}(\mathbf{\Delta}_j^{-1} \mathbf{Z}' \mathbf{Z} \mathbf{\Delta}_j^{-1})] = p \text{tr}(\mathbf{\Delta}_j^{-2}) \leq \frac{p\gamma_j}{\lambda_{j,\gamma_j}} \rightarrow 0.$$

This equation implies that

$$\mathbf{\Delta}_j^{-1} \mathbf{Z}' \mathbf{Z} \mathbf{\Delta}_j^{-1} \xrightarrow{p} \mathbf{O}_{\gamma_j, \gamma_j}. \quad (\text{A.11})$$

Moreover, it is easy to see that  $E[\mathbf{\Gamma}'_j \mathbf{Z}] = \mathbf{O}_{\gamma_j, \gamma_j}$  and  $E[\mathbf{Z}' \mathbf{\Gamma}_j] = \mathbf{O}_{\gamma_j, \gamma_j}$ , and

$$E[\text{tr}\{(\mathbf{\Delta}_j^{-1} \mathbf{\Gamma}'_j \mathbf{Z} \mathbf{\Delta}_j^{-1})' (\mathbf{\Delta}_j^{-1} \mathbf{\Gamma}'_j \mathbf{Z} \mathbf{\Delta}_j^{-1})\}] = \text{tr}(\mathbf{I}_{\gamma_j}) \text{tr}(\mathbf{\Delta}_j^{-2}) \leq \frac{\gamma_j}{\lambda_{j,\gamma_j}} \rightarrow 0.$$

These equations imply that

$$\mathbf{\Delta}_j^{-1} \mathbf{\Gamma}'_j \mathbf{Z} \mathbf{\Delta}_j^{-1} \xrightarrow{p} \mathbf{O}_{\gamma_j, \gamma_j}, \quad \mathbf{\Delta}_j^{-1} \mathbf{Z}' \mathbf{\Gamma}_j \mathbf{\Delta}_j^{-1} \xrightarrow{p} \mathbf{O}_{\gamma_j, \gamma_j}. \quad (\text{A.12})$$

From (A.10), (A.11) and (A.12), we derive the convergence in probability of  $\mathbf{\Delta}_j^{-1} \mathbf{U}_2 \mathbf{\Delta}_j^{-1}$  as

$$\mathbf{\Delta}_j^{-1} \mathbf{U}_2 \mathbf{\Delta}_j^{-1} \xrightarrow{p} \mathbf{I}_{\gamma_j}. \quad (\text{A.13})$$

Notice that

$$E[\text{tr}(\mathbf{\Delta}_j^{-1} \mathbf{U}_1 \mathbf{\Delta}_j^{-1})] = (n - k_j - p) \text{tr}(\mathbf{\Delta}_j^{-2}) \leq \frac{(n - k_j - p)\gamma_j}{\lambda_{j,\gamma_j}} \rightarrow 0.$$

This equation implies that

$$\mathbf{\Delta}_j^{-1} \mathbf{U}_1 \mathbf{\Delta}_j^{-1} \xrightarrow{p} \mathbf{O}_{\gamma_j, \gamma_j}. \quad (\text{A.14})$$



Combining the equations (A.9), (A.13) and (A.14) yields

$$\Delta_j^{-1}(\mathbf{U}_1 + \mathbf{U}_2)\Delta_j^{-1} \xrightarrow{P} \mathbf{I}_{\gamma_j}, \quad \frac{1}{n - k_j - \gamma_j - p}(\mathbf{U}_3 + \mathbf{U}_4) \xrightarrow{P} \frac{1}{1 - c_0} \mathbf{I}_{d_j}.$$

Using the results of the convergence in probability, the first and second terms in (A.7) are expanded as

$$\begin{aligned} -\log \frac{\det(\mathbf{U}_1)}{\det(\mathbf{U}_1 + \mathbf{U}_2)} &= \log \left( \frac{p}{1 - c_{n,p} - k_j/n} \right)^{\gamma_j} + \log \det\{(np)^{-1}\Delta_j^2\} \\ &\quad - \log \frac{\det\{\mathbf{U}_1/(n - k_j - p)\}}{\det\{\Delta_j^{-1}(\mathbf{U}_1 + \mathbf{U}_2)\Delta_j^{-1}\}} \\ &= \gamma_j \log p - \gamma_j \log(1 - c_0) \\ &\quad + \log \det\{(np)^{-1}\Delta_j^2\} + o_p(1), \end{aligned} \tag{A.15}$$

and

$$\begin{aligned} -\log \frac{\det(\mathbf{U}_3)}{\det(\mathbf{U}_3 + \mathbf{U}_4)} &= -\log \frac{\det\{\mathbf{U}_3/(n - k_j - \gamma_j - p)\}}{\det\{(\mathbf{U}_3 + \mathbf{U}_4)/(n - k_j - \gamma_j - p)\}} \\ &= -d_j \log(1 - c_0) + o_p(1). \end{aligned} \tag{A.16}$$

Notice that

$$\gamma_j \log \left( \frac{\lambda_{j,\gamma_j}}{np} \right) \leq \log \det \left( \frac{1}{np} \Delta_j^2 \right) \leq \gamma_j \log \left( \frac{\lambda_{j,1}}{np} \right).$$

It follows from the above result and assumption 3 that

$$\liminf_{c_{n,p} \rightarrow c_0} \log \det \left( \frac{1}{np} \Delta_j^2 \right) > 0, \quad \limsup_{c_{n,p} \rightarrow c_0} \log \det \left( \frac{1}{np} \Delta_j^2 \right) < \infty.$$

Therefore, we have

$$\frac{\log \det \{(np)^{-1}\Delta_j^2\}}{\log p} \rightarrow 0.$$

Using the above equation after substituting the equations (A.15) and (A.16) into (A.7) yields

$$\frac{1}{n \log p} \{\mathcal{L}(j) - \mathcal{L}(j_+)\} \xrightarrow{P} \gamma_j > 0. \tag{A.17}$$

Using the same idea as in the derivation of (A.16), it can be shown that

$$\frac{1}{n \log p} \{\mathcal{L}(j_+) - \mathcal{L}(j_*)\} \xrightarrow{P} 0. \tag{A.18}$$

From the results (A.17) and (A.18), when condition C2-1 holds, the difference between the information criteria of the model  $j$  and the true model  $j_*$  is con-

vergent as

$$\begin{aligned} & \frac{1}{n \log p} \{\text{IC}_m(j) - \text{IC}_m(j_*)\} \\ &= \frac{1}{n \log p} \{\mathcal{L}(j) - \mathcal{L}(j_+) + \mathcal{L}(j_+) - \mathcal{L}(j_*) + m(j) - m(j_*)\} \quad (\text{A.19}) \\ &\xrightarrow{p} \gamma_j + \lim_{c_{n,p} \rightarrow c_0} \frac{m(j) - m(j_*)}{n \log p} > 0. \end{aligned}$$

Next, we consider the case of  $j \in \mathcal{J}_+$ . Notice that

$$n \Sigma_*^{-1/2} \hat{\Sigma}_j \Sigma_*^{-1/2} \sim W_p(n - k_j, I_p), \quad n \Sigma_*^{-1/2} \hat{\Sigma}_{j_*} \Sigma_*^{-1/2} \sim W_p(n - k_*, I_p).$$

It follows from the property of the Wishart distributions (see [16, p. 57 th. 3.2.4]) that

$$\mathcal{L}(j) - \mathcal{L}(j_*) = n \log \Lambda, \quad (\text{A.20})$$

where  $\Lambda$  is distributed according to the Wilks' lambda distribution  $\Lambda_{r_j}(p, n - k_* - p)$  and  $r_j = k_j - k_*$ . By using Lemma A.1, we have

$$\begin{aligned} \frac{1}{p} \{\mathcal{L}(j) - \mathcal{L}(j_*)\} &= \left( -\frac{nr_j}{n - k_* - p - r_j} \right) \left( -\frac{n - k_* - p - r_j}{pr_j} \log \Lambda \right) \\ &\xrightarrow{p} -\frac{r_j}{1 - c_0} \left( \frac{c_0}{1 - c_0} \right)^{-1} \log \left( 1 + \frac{c_0}{1 - c_0} \right) = \frac{r_j}{c_0} \log(1 - c_0). \end{aligned}$$

Therefore, when condition C2-2 holds, the difference between the information criteria of the model  $j$  and the true model  $j_*$  is convergent as

$$\begin{aligned} \frac{1}{p} \{\text{IC}_m(j) - \text{IC}_m(j_*)\} &= \frac{1}{p} \{\mathcal{L}(j) - \mathcal{L}(j_*) + m(j) - m(j_*)\} \\ &\xrightarrow{p} \frac{r_j}{c_0} \log(1 - c_0) + \lim_{c_{n,p} \rightarrow c_0} \frac{m(j) - m(j_*)}{p} > 0. \end{aligned} \quad (\text{A.21})$$

Consequently, from (A.19) and (A.21), Theorem 3.2 is proved.

### A.3. The proof of Lemma A.1

The limiting values are easily obtained by using the following bounds for the cumulants of  $-\log \Lambda$ . Let  $\kappa^{(s)}$  be the  $s$ th cumulant of  $-\log \Lambda$  then

$$\begin{aligned} b^{(s)}(l+1, p, q) &< \kappa^{(s)} < b^{(s)}(l-1/2, p, q), \quad (s = 1, 2, \dots), \\ b^{(1)}(l, p, q) &= l \log \left( \frac{l}{l+q} \right) - (l+p) \log \left( \frac{l+p}{l+p+q} \right) - q \log \left( \frac{l+q}{l+p+q} \right), \\ b^{(2)}(l, p, q) &= 2 \log \left\{ 1 + \frac{pq}{(l+p+q)l} \right\}, \\ b^{(s)}(l, p, q) &= \frac{2^{s-1}(s-3)!}{l^{s-2}} \left\{ 1 - \left( \frac{l}{l+q} \right)^{s-2} - \left( \frac{l}{l+p} \right)^{s-2} + \left( \frac{l}{l+p+q} \right)^{s-2} \right\}. \end{aligned}$$

These bounds for  $s \geq 2$  are given by [34]. The bound for  $\kappa^{(1)}$  is also obtained by using the same method. The order of the  $2m$ th order central moment is given by

$$E[(T - \kappa_T^{(1)})^{2m}] = \sum_{l=1}^m \frac{(2m)!}{l!} \sum_{\substack{s_1 + \dots + s_l = 2m \\ s_1 \geq 2, \dots, s_l \geq 2}} \frac{\kappa_T^{(s_1)} \cdots \kappa_T^{(s_l)}}{s_1! \cdots s_l!} = \sum_{l=1}^m O((pq)^{-(2m-l)}).$$

**A.4. The proof of Corollary 3.1**

Recall that the HD asymptotic framework is used for proving Corollary 3.1. First, we consider the cases of the AIC and the AIC<sub>c</sub>. Notice that  $m(j) - m(j_*)$  in the AIC<sub>c</sub> can be expanded as

$$\begin{aligned} m(j) - m(j_*) &= \frac{p(k_j - k_*)(2 - c_{n,p} - 1/n)}{\{1 - c_{n,p} - (k_j + 1)/n\}\{1 - c_{n,p} - (k_* + 1)/n\}} \\ &= \frac{(k_j - k_*)(2 - c_{n,p})p}{(1 - c_{n,p})^2} + O(pn^{-1}). \end{aligned} \tag{A.22}$$

Hence, differences between the penalty terms of the AICs and the AIC<sub>c</sub>s are convergent as

$$\lim_{c_{n,p} \rightarrow c_0} \frac{1}{n \log p} \{m(j) - m(j_*)\} = 0.$$

This indicates that condition C2-1 holds in AIC and AIC<sub>c</sub>. Furthermore, it follows from the equality (A.22) that

$$\lim_{c_{n,p} \rightarrow c_0} \frac{1}{p} \{m(j) - m(j_*)\} = \begin{cases} 2(k_j - k_*) & \text{(AIC)} \\ (k_j - k_*)\{(1 - c_0)^{-1} + (1 - c_0)^{-2}\} & \text{(AIC}_c\text{)} \end{cases}.$$

Notice that  $c^{-1} \log(1 - c) + 2$  and  $c^{-1} \log(1 - c) + (1 - c)^{-1} + (1 - c)^{-2}$  are monotonically decreasing and increasing functions in  $0 \leq c < 1$ , respectively. Hence, when  $j \in \mathcal{J} \setminus \{j_*\}$ , the penalty terms in the AIC<sub>c</sub> always satisfy condition C2-2 and these in AIC satisfy condition C2-2 if  $c_0 \in [0, c_a)$ , where  $c_a$  is a constant satisfying  $\log(1 - c_a) + 2c_a = 0$ .

Next, we consider the cases of the BIC and the CAIC. When  $j \in \mathcal{J}_+ \setminus \{j_*\}$ , the differences between the penalty terms of the BICs and the CAICs are

$$\lim_{c_{n,p} \rightarrow c_0} \frac{1}{p \log n} \{m(j) - m(j_*)\} = k_j - k_* > 0.$$

Thus, condition C2-2 holds. Moreover, it is easy to obtain

$$\frac{1}{n \log p} \{m(j) - m(j_*)\} = \begin{cases} c_{n,p}(k_j - k_*) \left( -\frac{\log c_{n,p}}{\log p} + 1 \right) & \text{(BIC)} \\ c_{n,p}(k_j - k_*) \left( \frac{1 - \log c_{n,p}}{\log p} + 1 \right) & \text{(CAIC)}. \end{cases}$$

Since  $\lim_{c \rightarrow 0} c \log c = 0$  holds, we derive

$$\lim_{c_{n,p} \rightarrow c_0} \frac{1}{n \log p} \{m(j) - m(j_*)\} = c_0(k_j - k_*).$$

When  $j \in \mathcal{S}^c \cap \mathcal{J}_-$ , condition C2-1 is satisfied because  $c_0(k_j - k_*) \geq 0$  holds. When  $j \in \mathcal{S}_-$ , condition C2-1 is satisfied if  $c_0 < \gamma_j/(k_* - k_j)$  holds for all  $j \in \mathcal{S}_-$ .

**A.5. The proof of Theorem 3.3**

At first, we describe two lemmas which are used for proving Theorem 3.3 (the proofs of lemmas are given after this subsection).

**Lemma A.2.** Let  $\chi_n^2$  be a random variable distributed according to the chi-square distribution with  $n$  degrees of freedom. If  $z > n$ ,

$$P(\chi_n^2 > z) < \exp \left\{ -\frac{z}{2} \left( 1 - \frac{n}{z} - \frac{n}{z} \log \frac{n}{z} \right) \right\} < \exp \left[ -\frac{z}{2} \left\{ 1 - \left( \frac{n}{z} \right)^2 \right\} \right].$$

**Lemma A.3.** Let  $T = q^{-1} \log \det(\mathbf{V})$ , where  $\mathbf{V} \sim W_q(n, \mathbf{I}_q)$ , and let  $\kappa_T^{(s)}$  be the  $s$ th order cumulant of  $T$ . Then

$$\log(n - q) < \kappa_T^{(1)} < \log \left( n - \frac{q - 1}{2} \right). \tag{A.23}$$

Moreover if  $q/n \rightarrow 0$  as  $n \rightarrow \infty$

$$(qn)^{s-1} \kappa_T^{(s)} \rightarrow 2^{s-1} (s - 2)! \quad (s \geq 2), \tag{A.24}$$

and hence for any positive integer  $l$ ,

$$P(|T - \kappa_T^{(1)}| > \delta) \leq \frac{1}{\delta^{2l}} E[(T - \kappa_T^{(1)})^{2l}] = O((qn)^{-l}).$$

Recall that the HD-LM asymptotic framework is used for proving Theorem 3.3. First, we consider the case of  $j \in \mathcal{J}_-$ . Let  $d_j = k_{j_+} - k_j$ . As in the proof of Theorem 3.2, represent

$$\frac{1}{n} \{ \mathcal{L}(j) - \mathcal{L}(j_+) \} = \log \frac{\det(\mathbf{W}_1 + \mathbf{W}_2)}{\det(\mathbf{W}_1)} = \log \frac{\det(\mathbf{U}_1 + \mathbf{U}_2)}{\det(\mathbf{U}_1)}, \tag{A.25}$$

where  $\mathbf{W}_1$  and  $\mathbf{W}_2$  are independent, also  $\mathbf{U}_1$  and  $\mathbf{U}_2$  are independent, and

$$\begin{aligned} \mathbf{W}_1 &\sim W_p(n - k_{j_+}, \mathbf{I}_p), & \mathbf{W}_2 &\sim W_p(d_j, \mathbf{I}_p; \mathbf{\Gamma}_j \mathbf{\Gamma}'_j), \\ \mathbf{U}_1 &\sim W_{d_j}(n - p - k_j, \mathbf{I}_{d_j}), & \mathbf{U}_2 &\sim W_{d_j}(p, \mathbf{I}_{d_j}; \mathbf{\Omega}_j), \end{aligned}$$

with  $\mathbf{\Omega}_j = \text{diag}(\lambda_{j,1}, \dots, \lambda_{j,d_j})$ . It should be kept in mind that we recycle some notations to denote different random matrices from those in the proof of Theorem 3.2. Let  $q$  be the integer for  $j$  in condition C3-1. Express

$$\mathbf{U}_2 = (\mathbf{Z} + \mathbf{\Gamma}_{j,q})(\mathbf{Z} + \mathbf{\Gamma}_{j,q})' + \mathbf{U}_3,$$

where  $\mathbf{Z}$  and  $\mathbf{U}_3$  are mutually independent random matrices defined by

$$\mathbf{Z} \sim N_{d_j \times q}(\mathbf{O}_{d_j,q}, \mathbf{I}_{d_j} \otimes \mathbf{I}_q), \quad \mathbf{U}_3 \sim W_{d_j}(p - q, \mathbf{I}_{d_j}; \mathbf{\Omega}_j - \mathbf{\Gamma}_{j,q} \mathbf{\Gamma}'_{j,q}),$$

Here,  $\mathbf{\Gamma}_{j,q} = (\mathbf{\Omega}_{j,q}^{1/2}, \mathbf{O}_{q,d_j-q})'$  and  $\mathbf{\Omega}_{j,q} = \text{diag}(\lambda_{j,1}, \dots, \lambda_{j,q})$ . Then

$$\begin{aligned} \frac{\det(\mathbf{U}_1 + \mathbf{U}_2)}{\det(\mathbf{U}_1)} &> \frac{\det\{\mathbf{U}_1 + (\mathbf{Z} + \mathbf{\Gamma}_{j,q})(\mathbf{Z} + \mathbf{\Gamma}_{j,q})'\}}{\det(\mathbf{U}_1)} \\ &= \frac{\det\{\mathbf{V}_1 + (\mathbf{Z} + \mathbf{\Gamma}_{j,q})'(\mathbf{Z} + \mathbf{\Gamma}_{j,q})\}}{\det(\mathbf{V}_1)}, \end{aligned} \tag{A.26}$$

where

$$\begin{aligned} \mathbf{V}_1 &= \{(\mathbf{Z} + \mathbf{\Gamma}_{j,q})'(\mathbf{Z} + \mathbf{\Gamma}_{j,q})\}^{1/2} \{(\mathbf{Z} + \mathbf{\Gamma}_{j,q})'\mathbf{U}_1^{-1}(\mathbf{Z} + \mathbf{\Gamma}_{j,q})\}^{-1} \\ &\quad \times \{(\mathbf{Z} + \mathbf{\Gamma}_{j,q})'(\mathbf{Z} + \mathbf{\Gamma}_{j,q})\}^{1/2}. \end{aligned}$$

We can show that  $\mathbf{V}_1$  and  $\mathbf{Z}$  are independent, and  $\mathbf{V}_1 \sim W_q(n-p-k_{j_+}+q, \mathbf{I}_q)$  (see [16, p. 57 th. 3.2.4]). Let  $\delta_3 = \delta_2/2$  and  $h = 1 - \exp(-\delta_2/2)$ . Then  $0 < h < 1$  and

$$\log \frac{\beta_{q,j}(1-h)}{n-p-k_{j_+}+(q+1)/2} - \frac{1}{qn} \{m(j_+) - m(j)\} > \delta_3, \tag{A.27}$$

If an event

$$A_{j,h} : \text{tr}\{(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1/2}\mathbf{\Gamma}'_{j,q}\mathbf{Z}\mathbf{Z}'\mathbf{\Gamma}_{j,q}(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1/2}\} < h^2\lambda_{j,q}/4,$$

occurs, then for any unit vector  $\mathbf{b}$

$$\begin{aligned} &|\mathbf{b}'(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1/2}(\mathbf{Z}'\mathbf{\Gamma}_{j,q} + \mathbf{\Gamma}'_{j,q}\mathbf{Z})(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1/2}\mathbf{b}| \\ &\leq 2\{\mathbf{b}'(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1}\mathbf{b}\mathbf{b}'(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1/2}\mathbf{\Gamma}'_{j,q}\mathbf{Z}\mathbf{Z}'\mathbf{\Gamma}_{j,q}(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1/2}\mathbf{b}\}^{1/2} \\ &\leq 2\left[\lambda_{j,q}^{-1}\text{tr}\left\{(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1/2}\mathbf{\Gamma}'_{j,q}\mathbf{Z}\mathbf{Z}'\mathbf{\Gamma}_{j,q}(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1/2}\right\}\right]^{1/2} < h. \end{aligned}$$

Hence

$$\begin{aligned} &\det\{\mathbf{V}_1 + (\mathbf{Z} + \mathbf{\Gamma}_{j,q})'(\mathbf{Z} + \mathbf{\Gamma}_{j,q})\} \\ &> \det(\mathbf{Z}'\mathbf{\Gamma}_{j,q} + \mathbf{\Gamma}'_{j,q}\mathbf{Z} + \mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q}) \\ &> \det(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})\det(\mathbf{I}_q - h\mathbf{I}_q) = \beta_{j,q}^q(1-h)^q, \end{aligned} \tag{A.28}$$

where  $\beta_{j,q}$  is given by (3.4). Using (A.25), (A.26), (A.28) and (A.27), we obtain

$$\begin{aligned} &P(\hat{j}_m = j) \\ &\leq P(\text{IC}_m(j) - \text{IC}_m(j_+) < 0) \\ &< P(A_{j,h}^c) \\ &\quad + P(q \log \beta_{j,q} + q \log(1-h) + \{m(j) - m(j_+)\}/n < \log \det(\mathbf{V}_1)) \\ &< P(A_{j,h}^c) \\ &\quad + P(\log\{n-p-k_{j_+}+(q+1)/2\} + \delta_3 < q^{-1} \log \det(\mathbf{V}_1)). \end{aligned} \tag{A.29}$$

Since  $\text{tr}\{(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1/2}\mathbf{\Gamma}'_{j,q}\mathbf{Z}\mathbf{Z}'\mathbf{\Gamma}_{j,q}(\mathbf{\Gamma}'_{j,q}\mathbf{\Gamma}_{j,q})^{-1/2}\}$  is distributed according to the chi-square distribution with  $q^2$  degrees of freedom, using Lemma A.2 and

condition C3-1 we obtain

$$\begin{aligned} P(A_{j,h}^c) &< \exp \left[ -\frac{h^2 \lambda_{j,q}}{8} \left\{ 1 - \left( \frac{4q^2}{h^2 \lambda_{j,q}} \right)^2 \right\} \right] \\ &< \exp \left[ -\frac{h^2 q^2 n^{\delta_1}}{8} \left\{ 1 - \left( \frac{4}{h^2 n^{\delta_1}} \right)^2 \right\} \right], \end{aligned}$$

for sufficiently large  $n$ . Using Lemma A.3, we obtain

$$\begin{aligned} &P(\log\{n - p - k_{j_+} + (q + 1)/2\} + \delta_3 < q^{-1} \log \det(\mathbf{V}_1)) \\ &< P(q^{-1} |\log \det(\mathbf{V}_1) - E[\log \det(\mathbf{V}_1)]| > \delta_3) \\ &= O(\{q(n - p - k_{j_+} + q)\}^{-l}). \end{aligned} \tag{A.30}$$

Next, we consider the case of  $j \in \mathcal{J}_+ \setminus \{j_*\}$ . Let  $\Lambda$  be a random variable distributed according to  $\Lambda_{r_j}(p, n - k_* - p)$ , which is given in (A.20). From the equation in (3.3) and condition C3-2, by using Lemma A.1 we obtain that

$$\begin{aligned} &P(\hat{j}_m = j) \\ &\leq P(-(r_j c_{n,p})^{-1} \log \Lambda > (r_j p)^{-1} \{m(j) - m(j_*)\}) \\ &\leq P(-(r_j c_{n,p})^{-1} \log \Lambda + c_{n,p}^{-1} \log(1 - c_{n,p}) > \delta) \\ &\leq P((r_j c_{n,p})^{-1} |\log \Lambda - E[\log \Lambda]| > \delta/2) = O((r_j p)^{-l}), \end{aligned} \tag{A.31}$$

for sufficiently large  $n$ . In the last inequality, we used the fact that  $E[(r_j c_{n,p})^{-1} \log \Lambda] + c_{n,p}^{-1} \log(1 - c_{n,p}) \rightarrow 0$ .

From (A.29), (A.30), (A.31) and the equation in (3.3),

$$\max_{j \in \mathcal{J} \setminus \{j_*\}} P(\hat{j}_m = j) = O(p^{-l}).$$

Hence the equation in (3.3) leads that

$$P(\hat{j}_m = j_*) = 1 - \sum_{j \in \mathcal{J} \setminus \{j_*\}} P(\hat{j}_m = j) \rightarrow 1.$$

#### A.6. The proof of Lemma A.2

Notice that

$$\begin{aligned} P(\chi_n^2 > z) &= \frac{1}{2\Gamma(n/2)} \int_z^\infty e^{-x/2} (x/2)^{n/2-1} dx \\ &= \frac{1}{2\Gamma(n/2)} \int_z^\infty e^{-rx/2} e^{-(1-r)x/2} (x/2)^{n/2-1} dx \\ &< \frac{e^{-rz/2}}{2\Gamma(a/2)} \int_0^\infty e^{-(1-r)x/2} (x/2)^{n/2-1} dx = \frac{e^{-rz/2}}{(1-r)^{n/2}}, \end{aligned}$$

where  $\Gamma(x)$  is the gamma function. Taking the minimum with respect to  $r$  we get the first inequality. The second inequality can be obtained by the fact that  $1 + \log x < x$ .

**A.7. The proof of Lemma A.3**

At first, we describe the lemma which is used for proving Lemma A.3 (the proof of lemma is given after this subsection).

**Lemma A.4.** *Let  $\psi(a)$  be the digamma function defined by  $\psi(a) = d \log \Gamma(a)/da$ . Then,  $\psi(z) > \log(z - 1/2)$  holds if  $z > 1/2$  and  $\psi(z) < \log z$  if  $z > 0$ .*

Using the fact that  $\det(\mathbf{V}) \sim \prod_{i=1}^q \chi_{n-p+i}^2$ , where  $\chi_{n-p+1}^2, \dots, \chi_{n-p+q}^2$  are independent random variables and  $\chi_{n-p+i}^2$  is distributed according to the chi-square distribution with  $n - p + i$  degrees of freedom (see [21, p. 100 th. 3.2.15]), the moment generating function of  $\log \det(\mathbf{V})$  is given by

$$g(t) = \log E[\det(\mathbf{V})^t] = t(q \log 2) + \sum_{i=1}^q \sum_{s=1}^{\infty} \frac{t^s}{s!} \psi^{(s-1)}\left(\frac{n-q+i}{2}\right),$$

where  $\psi^{(s)}(a)$  is the  $s$ th order derivative of  $\psi(a)$ . Hence the first order cumulant of  $\log \det(\mathbf{V})$  is given by

$$\kappa^{(1)} = q \log 2 + \sum_{i=1}^q \psi\left(\frac{n-q+i}{2}\right).$$

Using Lemma A.4 and the fact that  $\log x$  is increasing and concave function of  $x$ ,

$$\begin{aligned} q \log(n-q) &< \sum_{i=1}^q \log(n-q+i-1) < \kappa^{(1)} \\ &< \sum_{i=1}^q \log(n-q+i) < q \log\left(n-q + \frac{q+1}{2}\right). \end{aligned}$$

The  $s$ th order cumulant of  $\log \det(\mathbf{V})$  can be expressed as

$$\kappa^{(s)} = \sum_{i=1}^q \sum_{k=0}^{\infty} \frac{(-1)^s (s-1)!}{\left(\frac{n-q+i}{2} + k\right)^s}, \quad (s \geq 2).$$

Since  $f(x, y) = 2^s(n - q + x + 2y)^{-s}$  is a decreasing and convex function of  $x$  and  $y$ ,

$$\begin{aligned} &\int_0^{\infty} \left\{ \int_1^{q+1} \frac{2^s (s-1)!}{(n-q+x+2y)^s} dx \right\} dy < (-1)^s \kappa^{(s)} \\ &< \int_{-1/2}^{\infty} \left\{ \int_{-1/2}^{q+1/2} \frac{2^s (s-1)!}{(n-q+x+2y)^s} dx \right\} dy. \end{aligned}$$

Calculating the integrals and taking the limits, we obtain (A.24). The  $2l$ th order central moment is the sum of the products of cumulants,  $\kappa_T^{(s_1)} \dots \kappa_T^{(s_k)}$  such that  $s_1 + \dots + s_k = 2l$ . Hence the order of the  $2l$ th order central moment is equal to the order of  $(\kappa_T^{(2)})^l$ .

### A.8. The proof of Lemma A.4

The digamma function can be expressed as

$$\psi(z) = -C + \sum_{k=0}^{\infty} \left( \frac{1}{1+k} - \frac{1}{z+k} \right),$$

where  $C$  is the Euler's constant defined by  $C = \lim_{n \rightarrow \infty} (\sum_{k=1}^n k^{-1} - \log n)$ . Hence

$$\lim_{n \rightarrow \infty} \{\log n - \psi(n)\} = \lim_{n \rightarrow \infty} \{\log(n-1/2) - \psi(n)\} = 0. \quad (\text{A.32})$$

Since  $f(x) = (z+x)^{-2}$  is convex and decreasing function of  $x$ ,

$$\begin{aligned} \frac{1}{z} &= \int_0^{\infty} \frac{1}{(z+x)^2} dx < \sum_{k=0}^{\infty} \frac{1}{(z+k)^2} = \frac{d}{dz} \psi(z) \\ &< \int_{-1/2}^{\infty} \frac{1}{(z+x)^2} dx = \frac{1}{(z-1/2)}. \end{aligned}$$

Hence  $\log z - \psi(z)$  is decreasing function, and  $\log(z-1/2) - \psi(z)$  is increasing function of  $z$ . Combining these properties with (A.32) we get the desired results.

### A.9. The proof of Corollary 3.2

Recall that the HD-LM asymptotic framework is used for proving Corollary 3.2. Notice that  $k_j/n \rightarrow 0$  holds when the equation in (3.3) holds. It is easy to see that condition C3-2 is rewritten as

$$\inf_{j \in \mathcal{J}_+ \setminus \{j_*\}} \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} \left\{ \frac{m(j) - m(j_*)}{p(k_j - k_*)} + \frac{1}{c_{n,p}} \log(1 - c_{n,p}) \right\} > 0.$$

Recall that  $\beta_{j,1} = \lambda_{j,1}$ . Hence, when  $q = 1$ , condition C3-1 is rewritten as

$$\begin{aligned} \inf_{j \in \mathcal{J}_-} \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} \frac{\log \lambda_{j,1}}{\log n} &> 0, \\ \inf_{j \in \mathcal{J}_-} \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} \left\{ \log \frac{\lambda_{j,1}}{n - p - k_j + 1} - \frac{m(j_+) - m(j)}{n} \right\} &> 0. \end{aligned}$$

Notice that

$$\log \frac{\lambda_{j,1}}{n - p - k_j + 1} = \log \frac{\lambda_{j,1}}{n} + \log \frac{n}{n - p - k_j + 1},$$

and

$$\begin{aligned} \lim_{c_{n,p} \rightarrow c_0, \text{LM}} \frac{1}{c_{n,p}} \log(1 - c_{n,p}) &= \frac{1}{c_0} \log(1 - c_0), \\ \lim_{c_{n,p} \rightarrow c_0, \text{LM}} \log \frac{n}{n - p - k_j + 1} &= -\log(1 - c_0). \end{aligned}$$

Hence, Corollary 3.2 is proved.



**A.10. The proof of Corollary 3.3**

Recall that the HD-LM asymptotic framework is used for proving Corollary 3.3. Notice that  $k_j/n \rightarrow 0$  holds for all  $j \in \mathcal{J}$  when the equation in (3.3) holds. From the above result and (A.22), we derive a limit of  $\{m(j) - m(j_*)\}/\{p(k_j - k_*)\}$  in each criterion as

$$\lim_{c_{n,p} \rightarrow c_0, \text{LM}} \frac{m(j) - m(j_*)}{p(k_j - k_*)} = \begin{cases} 2 & \text{(AIC)} \\ (1 - c_0)^{-1} + (1 - c_0)^{-2} & \text{(AIC}_c\text{)} \\ \infty & \text{(BIC, CAIC)} \end{cases} .$$

Therefore, condition C3-2 in the AIC<sub>c</sub>, BIC and CAIC holds, and that in the AIC holds if  $c_0 < c_a$ . Moreover,  $k_{j_+} - k_j$  is also bounded when  $k_*$  is bounded. It follows from  $k_{j_+} - k_j \leq k_*$  that

$$-\frac{m(j_+) - m(j)}{n} \geq \begin{cases} -2c_{n,p}k_* & \text{(AIC)} \\ -\frac{c_{n,p}(2 - c_{n,p} - 1/n)}{\{1 - c_{n,p} - (k_j + 1)/n\}^2} & \text{(AIC}_c\text{)} \\ -c_{n,p}k_* \log n & \text{(BIC)} \\ -c_{n,p}k_*(1 + \log n) & \text{(CAIC)} \end{cases} .$$

Hence, we have

$$\begin{aligned} \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} -\frac{m(j_+) - m(j)}{n} &= \begin{cases} -2k_*c_0 & \text{(AIC)} \\ -k_*c_0 \{(1 - c_0)^{-1} + (1 - c_0)^{-2}\} & \text{(AIC}_c\text{)} \end{cases} , \\ \liminf_{c_{n,p} \rightarrow c_0, \text{LM}} -\frac{m(j_+) - m(j)}{n \log n} &= -k_*c_0 \text{ (BIC, CAIC)}. \end{aligned}$$

Hence, condition C3-2 in the AIC and the AIC<sub>c</sub> holds if (3.5) and (3.6) are satisfied, respectively, and that in the BIC and the CAIC holds if (3.7) is satisfied. Consequently, Corollary 3.2 is proved.

**A.11. Example of the noncentrality matrix**

Theoretically, it is natural to describe conditions in terms of the eigen values of the noncentrality matrices, since the unknown parameters  $\Theta_*$  and  $\Sigma_*$  affect the distribution of each criterion only through the eigen values. However, the effect of each  $\Theta_*$  and  $\Sigma_*$  may be of interest. So we illustrate the conditions by a two-way MANOVA model with a certain structure for the regression matrix and the covariance matrix of the error term.

Suppose there are two factors A and B with  $a$  levels and  $b$  levels, respectively. We consider the case that a characteristic is observed at time  $t$  ( $t = 1, \dots, p$ ) for  $m$  individuals in the cell  $(l, k)$  ( $l = 1, \dots, a; k = 1, \dots, b$ ). Let  $\mathbf{y}_{ilk} = (y_{ilk1}, \dots, y_{ilkp})'$  be  $p$ -dimensional observation vector for the  $i$ th individual from the cell  $(l, k)$  ( $i = 1, \dots, m$ ). Then two-way MANOVA model is represented as

$$y_{ilkt} = \theta_t + \theta_{lt}^{(A)} + \theta_{kt}^{(B)} + \theta_{lkt}^{(AB)} + \epsilon_{ilkt},$$

$$(i = 1, \dots, m; l = 1, \dots, a; k = 1, \dots, b; t = 1, \dots, p),$$

where  $\theta_{lt}^{(A)}$ 's and  $\theta_{kt}^{(B)}$ 's are the main effects of the factor A and B, respectively,  $\theta_{lkt}^{(AB)}$ 's are the effects of the interaction, and  $\epsilon = (\epsilon_{ilk1}, \dots, \epsilon_{ilkp})' \sim N_p(\mathbf{0}_p, \Sigma)$ . In order to assure the model identifiability, we assume that

$$\theta_{1t}^{(A)} = \theta_{1t}^{(B)} = \theta_{1kt}^{(AB)} = \theta_{l1t}^{(AB)} = 0, \quad (i = 1, \dots, p; l = 1, \dots, a; k = 1, \dots, b).$$

Let

$$\mathbf{Y} = (\mathbf{y}_{111}, \dots, \mathbf{y}_{11m}, \mathbf{y}_{121}, \dots, \mathbf{y}_{12m}, \dots, \mathbf{y}_{ab1}, \dots, \mathbf{y}_{abm})',$$

$$\mathbf{X} = \left[ \mathbf{1}_{abm} \left| \begin{array}{c} \left[ \mathbf{0}'_{a-1} \right] \\ \left[ \mathbf{I}_{a-1} \right] \end{array} \otimes \mathbf{1}_{bm} \right| \mathbf{1}_a \otimes \begin{array}{c} \left[ \mathbf{0}'_{b-1} \right] \\ \left[ \mathbf{I}_{b-1} \right] \end{array} \otimes \mathbf{1}_m \left| \begin{array}{c} \left[ \mathbf{0}'_{a-1} \right] \\ \left[ \mathbf{I}_{a-1} \right] \end{array} \otimes \begin{array}{c} \left[ \mathbf{0}'_{b-1} \right] \\ \left[ \mathbf{I}_{b-1} \right] \end{array} \otimes \mathbf{1}_m \right. \right].$$

main effect of A                      main effect of B                      interaction

Then  $\mathbf{Y} \sim N_{n \times p}(\mathbf{X}\Theta, \Sigma \otimes \mathbf{I}_n)$ , where  $n = abm$  and

$$\Theta = (\boldsymbol{\theta}, \boldsymbol{\theta}_2^{(A)}, \dots, \boldsymbol{\theta}_a^{(A)}, \boldsymbol{\theta}_2^{(B)}, \dots, \boldsymbol{\theta}_b^{(B)}, \boldsymbol{\theta}_{22}^{(AB)}, \boldsymbol{\theta}_{23}^{(AB)}, \dots, \boldsymbol{\theta}_{ab}^{(AB)})',$$

$$\boldsymbol{\theta}_l^{(A)} = (\theta_{l1}^{(A)}, \dots, \theta_{lp}^{(A)})',$$

$$\boldsymbol{\theta}_k^{(B)} = (\theta_{k1}^{(B)}, \dots, \theta_{kp}^{(B)})', \quad (l = 2, \dots, a; k = 2, \dots, b),$$

$$\boldsymbol{\theta}_{lk}^{(AB)} = (\theta_{lkt}^{(AB)}, \dots, \theta_{lkp}^{(AB)})',$$

Suppose that the true model has no interactions and the effects of the factor B are parallel, that is

$$\boldsymbol{\theta}_k^{(B)} = \eta_k \boldsymbol{\theta}_B \quad (k = 2, \dots, b)$$

with unknown parameters  $\eta_2, \dots, \eta_b$  and a  $p \times 1$  vector  $\boldsymbol{\theta}_B = (\theta_1^{(B)}, \dots, \theta_p^{(B)})'$ . Hence the true model is represented as  $j_* = \{1, 2, \dots, a + b - 1\}$  and

$$\mathbf{X}_* = \left[ \mathbf{1}_{abm} \left| \begin{array}{c} \left[ \mathbf{0}'_{a-1} \right] \\ \left[ \mathbf{I}_{a-1} \right] \end{array} \otimes \mathbf{1}_{bm} \right| \mathbf{1}_a \otimes \begin{array}{c} \left[ \mathbf{0}'_{b-1} \right] \\ \left[ \mathbf{I}_{b-1} \right] \end{array} \otimes \mathbf{1}_m \right],$$

$$\Theta_* = (\boldsymbol{\theta}, \boldsymbol{\theta}_2^{(A)}, \dots, \boldsymbol{\theta}_a^{(A)}, \boldsymbol{\theta}_B \boldsymbol{\eta}')',$$

with  $\boldsymbol{\eta} = (\eta_2, \dots, \eta_b)'$ . In order to distinguish the variations due to the difference among individuals and the observation error we consider the case that  $\epsilon_{ilk}^{(t)} = u_{ilk}^{(t)} + v_i$ , where  $v_i, u_{ilk}^{(t)}$  ( $l = 1, \dots, a; k = 1, \dots, b; t = 1, \dots, p$ ) are mutually independent,  $u_{ilk}^{(t)} \sim N(0, \sigma^2)$ , and  $v_i \sim N(0, \tau^2)$ . Then

$$\Sigma_* = \sigma^2 \mathbf{I}_p + \tau^2 \mathbf{1}_p \mathbf{1}_p'.$$

Under the above setup, the noncentrality matrix of an underspecified model  $j = \{1, 2, \dots, a\}$  is given by

$$\Gamma_j \Gamma_j' = am \Sigma_*^{-1/2} \boldsymbol{\theta}_B \boldsymbol{\eta}' (\mathbf{I}_{b-1} - b^{-1} \mathbf{1}_{b-1} \mathbf{1}_{b-1}') \boldsymbol{\eta} \boldsymbol{\theta}_B' \Sigma_*^{-1/2}.$$

Note that the structures of the regression matrix and the covariance matrix are just for the illustration, and the user of the models does not assume these structures. So the number of the unknown parameters of the model  $j$  is  $k_j p + p(p+1)/2$  with  $k_j = a$ .

Since  $\gamma_j = \text{rank}(\mathbf{\Gamma}_j) = 1$ ,  $q = 1$  in condition C3-1, and some algebraic calculation leads that

$$\beta_{j,q} = \lambda_{j,1} = \mathbf{\Gamma}'_j \mathbf{\Gamma}_j = (np) Q_B S_B,$$

where

$$Q_B = \frac{1}{\sigma^2} \frac{1}{p} \sum_{t=1}^p (\theta_t^{(B)} - \bar{\theta}_B)^2 + \frac{1}{\sigma^2 + p\tau^2} \bar{\theta}_B^2, \quad \bar{\theta}_B = \frac{1}{p} \sum_{t=1}^p \theta_t^{(B)},$$

$$S_B = \frac{1}{b} \sum_{k=1}^b (\eta_k - \bar{\eta})^2, \quad \eta_1 = 0, \quad \bar{\eta} = \frac{1}{b} \sum_{k=1}^b \eta_k.$$

If the factor B is a kind of medication for example, the effect is likely to diminish as time goes by. In this case  $Q_B \rightarrow 0$  as  $p \rightarrow \infty$ . However Theorem 3.3 says that even in this case  $IC_m$  has a possibility to have the consistency property. Actually, Corollary 3.3 says that it is sufficient for the AIC and the AIC<sub>c</sub> having consistency if  $pQ_B \rightarrow \infty$ .

## References

- [1] AKAIKE, H. (1973). Information theory and an extension of the maximum likelihood principle. In *2nd. International Symposium on Information Theory* (eds. B. N. Petrov and F. Csáki) 267–281. Akadémiai Kiadó, Budapest. [MR0483125](#)
- [2] AKAIKE, H. (1974). A new look at the statistical model identification. *IEEE Trans. Automatic Control* **AC-19** 716–723. [MR0423716](#)
- [3] BEDRICK, E. J. and TSAI, C.-L. (1994). Model selection for multivariate regression in small samples. *Biometrics* **50** 226–231.
- [4] BOSQ, D. (2000). *Linear Processes in Function Spaces. Theory and Applications*. Springer-Verlag, New York. [MR1783138](#)
- [5] BOSQ, D. and BLANKE, D. (2007). *Inference and Prediction in Large Dimensions*. John Wiley & Sons, Ltd., Paris. [MR2364006](#)
- [6] BOZDOGAN, H. (1987). Model selection and Akaike's information criterion (AIC): The general theory and its analytical extensions. *Psychometrika* **52** 345–370. [MR0914460](#)
- [7] CHRISTAKOS, G. (2000). *Modern Spatiotemporal Geostatistics*. Oxford University Press, New York.
- [8] CRESSIE, N. and WIKLE, C. K. (2011). *Statistics for Spatio-Temporal Data*. John Wiley & Sons, Inc., Hoboken. [MR2848400](#)
- [9] DAVIES, S. J., NEATH, A. A. and CAVANAUGH, J. E. (2006). Estimation optimality of corrected AIC and modified  $C_p$  in linear regression model. *International Statist. Review* **74** 161–168.

- [10] DIEN, S. J. V., IWATANI, S., USUDA, Y. and MATSUI, K. (2006). Theoretical analysis of amino acid-producing *Escherichia coli* using a stoichiometric model and multivariate linear regression. *J. Biosci. Bioeng.* **102** 34–40.
- [11] FUJIKOSHI, Y. (1983). A criterion for variable selection in multiple discriminant analysis. *Hiroshima Math. J.* **13** 203–214. [MR0693557](#)
- [12] FUJIKOSHI, Y. (1985). Selection of variables in two-group discriminant analysis by error rate and Akaike's information criteria. *J. Multivariate Anal.* **17** 27–37. [MR0797518](#)
- [13] FUJIKOSHI, Y. and SAKURAI, T. (2009). High-dimensional asymptotic expansions for the distributions of canonical correlations. *J. Multivariate Anal.* **100** 231–242. [MR2460489](#)
- [14] FUJIKOSHI, Y. and SATOH, K. (1997). Modified AIC and  $C_p$  in multivariate linear regression. *Biometrika* **84** 707–716. [MR1603952](#)
- [15] FUJIKOSHI, Y. and SEO, T. (1998). Asymptotic approximations for EPMC's of the linear and the quadratic discriminant functions when the sample sizes and the dimension are large. *Random Oper. Stochastic Equations* **6** 269–280. [MR1631003](#)
- [16] FUJIKOSHI, Y., SHIMIZU, R. and ULYANOV, V. V. (2010). *Multivariate Statistics: High-Dimensional and Large-Sample Approximations*. John Wiley & Sons, Inc., Hoboken, New Jersey. [MR2640807](#)
- [17] FUJIKOSHI, Y., YANAGIHARA, H. and WAKAKI, H. (2005). Bias corrections of some criteria for selection multivariate linear regression models in a general case. *Amer. J. Math. Management Sci.* **25** 221–258. [MR2255869](#)
- [18] HARVILLE, D. A. (1997). *Matrix Algebra from a Statistician's Perspective*. Springer-Verlag, New York. [MR1467237](#)
- [19] KIM, Y., KWON, S. and CHOI, H. (2012). Consistent model selection criteria on high dimensions. *J. Mach. Learn. Res.* **13** 1037–1057. [MR2930632](#)
- [20] KULLBACK, S. and LEIBLER, R. A. (1951). On information and sufficiency. *Ann. Math. Statist.* **22** 79–86. [MR0039968](#)
- [21] MUIRHEAD, R. J. (1982). *Aspects of Multivariate Statistical Theory*. John Wiley & Sons, Inc., New York. [MR0652932](#)
- [22] NISHII, R. (1984). Asymptotic properties of criteria for selection of variables in multiple regression. *Ann. Statist.* **12** 758–765. [MR0740928](#)
- [23] RAMSAY, J. O. and SILVERMAN, B. W. (2005). *Functional Data Analysis* (2nd. ed.). Springer, New York. [MR2168993](#)
- [24] SÂRBU, C., ONIŞOR, C., POSA, M., KEVRESAN, S. and KUHAJDA, K. (2008). Modeling and prediction (correction) of partition coefficients of bile acids and their derivatives by multivariate regression methods. *Talanta* **75** 651–657.
- [25] SAXÉN, R. and SUNDELL, J. (2006).  $^{137}\text{Cs}$  in freshwater fish in Finland since 1986 – a statistical analysis with multivariate linear regression models. *J. Environ. Radioactiv.* **87** 62–76.
- [26] SCHWARZ, G. (1978). Estimating the dimension of a model. *Ann. Statist.* **6** 461–464. [MR0468014](#)

- [27] SHAO, J. (1997). An asymptotic theory for linear model selection. *Statist. Sinica* **7** 221–264. [MR1466682](#)
- [28] SHIBATA, R. (1976). Selection of the order of an autoregressive model by Akaike's information criterion. *Biometrika* **63** 117–126. [MR0403130](#)
- [29] SHIBATA, R. (1980). Asymptotically efficient selection of the order of the model for estimating parameters of a linear process. *Ann. Statist.* **8** 147–164. [MR0557560](#)
- [30] SIOTANI, M., HAYAKAWA, T. and FUJIKOSHI, Y. (1985). *Modern Multivariate Statistical Analysis: A Graduate Course and Handbook*. American Sciences Press, Columbus, Ohio. [MR0832440](#)
- [31] SRIVASTAVA, M. S. (2002). *Methods of Multivariate Statistics*. John Wiley & Sons, New York. [MR1915968](#)
- [32] SUGIURA, N. (1978). Further analysis of the data by Akaike's information criterion and the finite corrections. *Commun. Statist. Theory Methods* **A7** 13–26.
- [33] TIMM, N. H. (2002). *Applied Multivariate Analysis*. Springer-Verlag, New York. [MR1908225](#)
- [34] WAKAKI, H. (2006). Edgeworth expansion of Wilks' lambda statistic. *J. Multivariate Anal.* **97** 1958–1964. [MR2301267](#)
- [35] WOODROOFE, M. (1982). On model selection and the arc sine laws. *Ann. Statist.* **10** 1182–1194. [MR0673653](#)
- [36] YAMAMURA, M., YANAGIHARA, H. and SRIVASTAVA, M. S. (2010). Variable selection in multivariate linear regression models with fewer observations than the dimension. *Japanese J. Appl. Statist.* **39** 1–19.
- [37] YANAGIHARA, H. (2006). Corrected version of AIC for selecting multivariate normal linear regression models in a general nonnormal case. *J. Multivariate Anal.* **97** 1070–1089. [MR2276149](#)
- [38] YANAGIHARA, H., KAMO, K. and TONDA, T. (2011). Second-order bias-corrected AIC in multivariate normal linear models under nonnormality. *Canad. J. Statist.* **39** 126–146. [MR2815342](#)
- [39] YANG, Y. (2005). Can the strengths of AIC and BIC be shared? A conflict between model identification and regression estimation. *Biometrika* **92** 937–950. [MR2234196](#)
- [40] YOSHIMOTO, A., YANAGIHARA, H. and NINOMIYA, Y. (2005). Finding factors affecting a forest stand growth through multivariate linear modeling. *J. Jpn. For. Soc.* **87** 504–512 (in Japanese).