

# On model selection consistency using a kick-one-out method for selecting response variables in high-dimensional multivariate linear regression

Ryoya Oda\*<sup>†</sup> Hirokazu Yanagihara<sup>‡</sup> and Yasunori Fujikoshi<sup>§</sup>

(Last Modified: October 8, 2021)

## Abstract

This paper deals with the selection of non-redundant response variables in normality-assumed multivariate linear regression, where the redundancy of the response variables is defined by conditional expectation. A sufficient condition for model selection consistency is obtained using a kick-one-out method based on the generalized information criterion under a high-dimensional asymptotic framework such that the sample size tends to infinity, and the number of response variables and explanatory variables does not exceed the sample size but may tend to infinity. A consistent kick-one-out method using the obtained condition is proposed. Simulation results show that the proposed method has a high probability of selecting true non-redundant variables.

## 1 Introduction

The multivariate linear regression model is fundamental in inferential statistical analysis and is introduced in many statistics textbooks (e.g., [20, 22]). Suppose that  $\mathbf{y}$  is a  $p$ -dimensional response vector and  $\mathbf{x}$  is a  $k$ -dimensional explanatory vector. Then, the normality-assumed multivariate linear regression model with  $\mathbf{y}$  and  $\mathbf{x}$  is given by

$$\mathbf{y} \sim \mathcal{N}_p(\Theta' \mathbf{x}, \Sigma), \quad (1)$$

where  $\Theta$  is a  $k \times p$  unknown regression coefficients matrix and  $\Sigma$  is a  $p \times p$  unknown covariance matrix that is positive definite. In actual empirical contexts, selecting variables for the model is of key interest. It is generally expected that the accuracy of prediction is improved and that interpretation of the model is made easier by proper variable selection. In the literature on multivariate linear regression, numerous papers have dealt with the variable selection problem as it relates to selecting explanatory variables. On the other hand, it is also important to consider

---

\*Corresponding author. Email: ryoya-oda@hiroshima-u.ac.jp

<sup>†</sup>School of Informatics and Data Science, Hiroshima University, 1-3-1 Kagamiyama, Higashi-Hiroshima, Hiroshima, Japan

<sup>‡</sup>Graduate School of Advanced Science and Engineering, Hiroshima University, 1-3-1 Kagamiyama, Higashi-Hiroshima, Hiroshima, Japan

<sup>§</sup>Department of Mathematics, Graduate School of Science, Hiroshima University, 1-3-1 Kagamiyama, Higashi-Hiroshima, Hiroshima, Japan

the problem of selecting appropriate response variables, as including redundant response variables in the model will almost certainly lead to an improper understanding of the essential structure of the model. [8] defined redundancy using conditional expectation. Specifically, a response variable is determined to be redundant when the conditional expectation of the response variable given the non-redundant variables does not rely on the explanatory variables. This implies that the relationship between the redundant variables and the explanatory variables can be expressed by the relationship between the non-redundant variables and the explanatory variables. Thus, removing such redundant variables will promote a better understanding of the essential structure of the model.

As a result of the rapid advances in information technology and database systems in recent years, efforts to analyze high-dimensional data where not only the sample size but also the number of variables is large have become commonplace. Although [8] derived a formal Akaike's information criterion (AIC) [1, 2] for selecting response variables, the behavior of the formal AIC is not apparent for high-dimensional data. Moreover, variable selection using criteria such as the AIC is often performed by minimizing the criterion for all combinations of response variables, i.e., all subsets of  $\omega$ , where  $\omega = \{1, \dots, p\}$  is defined as the full set of response variable suffixes. This means that it may be impossible to apply the approach due to its high computational cost when the number of response variables is large. For this reason, we focus on the practicable selection method proposed by [12, 24], using the kick-one-out (KOO) method described in [3]. With the KOO method, if the criterion value for the subset with one variable removed from the full set  $\omega$  is greater than the criterion value for  $\omega$ , then the removed variable is selected. This makes it possible to apply the KOO method for high-dimensional data without the risk of incurring excessive computational costs.

In this paper, we propose the use of the generalized information criterion (GIC) developed in [11] as a variable selection criterion to address the redundancy issue described in [8]. Specifically, we propose a KOO method based on the GIC which has a model selection consistency property under a high-dimensional asymptotic framework such that the sample size always tends to infinity and the number of response variables and explanatory variables does not exceed the sample size but may tend to infinity. Consistency here means that the probability that the selected variables are identical to the true non-redundant variables tends to 1 under a high-dimensional asymptotic framework. Thus, it can be expected that the proposed method has a high probability of selecting true non-redundant variables when the sample size is large, regardless of the number of variables. In several previous works, a consistency property of the KOO method has been obtained for selection problems under a high-dimensional asymptotic framework in several multivariate models, e.g., selecting the explanatory variables in a multivariate linear regression model [3, 15, 16, 18] and selecting redundant variables in discriminant analysis [7, 10, 14, 17]. However, there is no consistent KOO method for selecting response variables in the sense of the redundancy described in [8].

The remainder of the paper is organized as follows: In section 2, we define the redundancy of response variables and introduce the KOO method based on the GIC. In section 3, we describe the model selection consistency property of the KOO method based on the GIC under a high-dimensional asymptotic framework. In section 4, we conduct numerical experiments for verification purposes. Technical details are relegated to the Appendix.

## 2 Framework for selecting response variables

### 2.1 Redundancy of response variables

We define the redundancy of response variables as in [8]. Let  $j$  denote a subset of  $\omega = \{1, \dots, p\}$  containing  $p_j$  elements and  $\mathbf{y}_j$  denote the  $p_j$ -dimensional vector consisting of the components of  $\mathbf{y}$  indexed by the elements of  $j$ . For example, if  $j = \{1, 2, 4\}$ , then  $\mathbf{y}_j$  consists of the first, second, and fourth elements of  $\mathbf{y}$ . Without loss of generality, we sort  $\mathbf{y}$  into  $\mathbf{y} = (\mathbf{y}'_j, \mathbf{y}'_{\bar{j}})'$  for a subset  $j$ , where  $\mathbf{y}_{\bar{j}}$  is the  $(p - p_j)$ -dimensional vector and  $\bar{A}$  denotes the complement of a set  $A$ . Similar to the division of  $\mathbf{y}$ , we can express the division of  $\Theta$  and  $\Sigma$  as follows:

$$\Theta = (\Theta_j, \Theta_{\bar{j}}), \quad \Sigma = \begin{pmatrix} \Sigma_{jj} & \Sigma_{j\bar{j}} \\ \Sigma'_{j\bar{j}} & \Sigma_{\bar{j}\bar{j}} \end{pmatrix}, \quad (2)$$

where  $\Theta_j$  and  $\Theta_{\bar{j}}$  are  $k \times p_j$  and  $k \times (p - p_j)$ ,  $\Sigma_{jj}$  and  $\Sigma_{\bar{j}\bar{j}}$  are  $p_j \times p_j$  and  $(p - p_j) \times (p - p_j)$ , and  $\Sigma_{j\bar{j}}$  is  $p_j \times (p - p_j)$ . Then, from a property of a multivariate normal distribution (e.g., [21]), the conditional distribution of  $\mathbf{y}_{\bar{j}}$  given  $\mathbf{y}_j$  can be written as

$$\mathbf{y}_{\bar{j}} | \mathbf{y}_j \sim \mathcal{N}_{p-p_j} \left( (\Theta_{\bar{j}} - \Theta_j \Sigma_{jj}^{-1} \Sigma_{j\bar{j}})' \mathbf{x} + \Sigma'_{j\bar{j}} \Sigma_{jj}^{-1} \mathbf{y}_j, \Sigma_{\bar{j}\bar{j}.j} \right), \quad (3)$$

where  $\Sigma_{\bar{j}\bar{j}.j} = \Sigma_{\bar{j}\bar{j}} - \Sigma'_{j\bar{j}} \Sigma_{jj}^{-1} \Sigma_{j\bar{j}}$ . From (3), if the equation  $\Theta_{\bar{j}} - \Theta_j \Sigma_{jj}^{-1} \Sigma_{j\bar{j}} = \mathbf{O}_{k, p-p_j}$  holds, then the conditional distribution of  $\mathbf{y}_{\bar{j}}$  given  $\mathbf{y}_j$  does not depend on the explanatory vector  $\mathbf{x}$ . In other words, the relationship between  $\mathbf{y}_{\bar{j}}$  and  $\mathbf{x}$  can be described by the relationship between  $\mathbf{y}_j$  and  $\mathbf{x}$ . In this sense, the model such that  $\mathbf{y}_{\bar{j}}$  is redundant in the relationship between  $\mathbf{y}$  and  $\mathbf{x}$  may be expressed as

$$M_j : (3) \text{ s.t. } \Theta_{\bar{j}} - \Theta_j \Sigma_{jj}^{-1} \Sigma_{j\bar{j}} = \mathbf{O}_{k, p-p_j}. \quad (4)$$

We note that (4) is also related to selecting response variables in a multivariate inverse regression (e.g., [5, 6, 13]). In fact, one of the purposes in a multivariate inverse regression is to estimate an unknown explanatory vector  $\mathbf{x}_0$  corresponding to an observed response vector  $\mathbf{y}_0$ . Then, if  $\Theta$  is full row rank, the classical estimator  $\hat{\mathbf{x}}_0$  of  $\mathbf{x}_0$  can be expressed as

$$\hat{\mathbf{x}}_0 = \arg \min_{\mathbf{x}_0} (\mathbf{y}_0 - \Theta' \mathbf{x}_0)' \Sigma^{-1} (\mathbf{y}_0 - \Theta' \mathbf{x}_0) = (\Theta \Sigma^{-1} \Theta')^{-1} \Theta \Sigma^{-1} \mathbf{y}_0 = \Xi_j \mathbf{y}_{0,j} + \Xi_{\bar{j}} \mathbf{y}_{0,\bar{j}},$$

where  $(\Xi_j, \Xi_{\bar{j}}) = (\Theta \Sigma^{-1} \Theta')^{-1} \Theta \Sigma^{-1}$  and  $\mathbf{y}_0 = (\mathbf{y}'_{0,j}, \mathbf{y}'_{0,\bar{j}})'$ . Moreover, it holds that  $\Xi_{\bar{j}} = \Theta_{\bar{j}} - \Theta_j \Sigma_{jj}^{-1} \Sigma_{j\bar{j}}$ . Hence, the redundancy model (4) can be also regarded as the redundancy of response variables for estimating the classical estimator  $\hat{\mathbf{x}}_0$  in a multivariate inverse regression.

### 2.2 Kick-one-out method based on GIC

We first introduce the generalized information criterion (GIC) for the redundancy model (4). Let  $\{(\mathbf{y}_{(i)}, \mathbf{x}_{(i)})\}$  ( $i = 1, \dots, n$ ) be a set of *i.i.d.* copies from  $(\mathbf{y}, \mathbf{x})$  and let  $\mathbf{Y} = (\mathbf{y}_{(1)}, \dots, \mathbf{y}_{(n)})'$  and  $\mathbf{X} = (\mathbf{x}_{(1)}, \dots, \mathbf{x}_{(n)})'$ , where  $n$  is the sample size. We assume  $\text{rank}(\mathbf{X}) = k < n$  and  $n - p - k - 1 > 0$  for the existence of the GIC and our proposed method. Suppose  $\mathbf{Y}_j$  and  $\mathbf{Y}_{\bar{j}}$  are the  $n \times p_j$  and  $n \times (p - p_j)$  partitioned matrices of  $\mathbf{Y} = (\mathbf{Y}_j, \mathbf{Y}_{\bar{j}})$  for the redundancy

model (4). The model (4) can be regarded as a redundancy model as in estimating an unknown explanatory vector  $\boldsymbol{x}_0$  by using the classical estimator  $\hat{\boldsymbol{x}}_0$  in a multivariate inverse regression. Hence, in accordance with [13], the maximum log-likelihood of  $\mathcal{N}_{n \times p}(\boldsymbol{X}\boldsymbol{\Theta}, \boldsymbol{\Sigma} \otimes \boldsymbol{I}_n)$  under the redundant model (4) is given by

$$np(\log 2\pi + 1) + n \log \left| \frac{1}{n} \boldsymbol{Y}'_j (\boldsymbol{I}_n - \boldsymbol{P}_X) \boldsymbol{Y}_j \right| + n \log \left| \frac{1}{n} \boldsymbol{Y}'_{\bar{j}} (\boldsymbol{I}_n - \boldsymbol{P}_{Y_j}) \boldsymbol{Y}_{\bar{j}} \right|,$$

where  $\boldsymbol{P}_A = \boldsymbol{A}(\boldsymbol{A}'\boldsymbol{A})^{-1}\boldsymbol{A}'$  for a square matrix  $\boldsymbol{A}$ . Then, the GIC for the model (4) is written as

$$\begin{aligned} \text{GIC}(j) = & np(\log 2\pi + 1) + n \log \left| \frac{1}{n} \boldsymbol{Y}'_j (\boldsymbol{I}_n - \boldsymbol{P}_X) \boldsymbol{Y}_j \right| + n \log \left| \frac{1}{n} \boldsymbol{Y}'_{\bar{j}} (\boldsymbol{I}_n - \boldsymbol{P}_{Y_j}) \boldsymbol{Y}_{\bar{j}} \right| \\ & + \alpha \left\{ \frac{1}{2} p(p+1) + kp_j \right\}, \end{aligned} \quad (5)$$

where  $\alpha$  is a positive constant expressing a penalty for the complexity of the model. Various variable selection criteria are included in the GIC by specifying  $\alpha$ . For example, the AIC [1, 2], HQC [9], BIC [19], and CAIC [4] are expressed by the GIC for the following  $\alpha$ :

$$\alpha = \begin{cases} 2 & \text{(AIC)} \\ 2 \log \log n & \text{(HQC)} \\ \log n & \text{(BIC)} \\ 1 + \log n & \text{(CAIC)} \end{cases}.$$

Next, we introduce the kick-one-out method based on the GIC. Denote by  $\ell$  the subset of  $\omega$  satisfying  $\#\ell = p-1$  and, moreover, denote by  $\bar{\ell}$  the compliment of  $\ell$ . Then, the best subset decided by the kick-one-out method based on the GIC can be written as

$$\hat{j} = \{\bar{\ell} \in \omega \mid \text{GIC}(\ell) > \text{GIC}(\omega)\}, \quad (6)$$

where  $\text{GIC}(\omega)$  is the GIC for the full set  $\omega$  and is defined for the model such that all the response variables are non-redundant, i.e.,  $\text{GIC}(\omega)$  is written as

$$\text{GIC}(\omega) = np(\log 2\pi + 1) + n \log \left| \frac{1}{n} \boldsymbol{Y}' (\boldsymbol{I}_n - \boldsymbol{P}_X) \boldsymbol{Y} \right| + \alpha \left\{ \frac{1}{2} p(p+1) + kp \right\}.$$

### 3 Model selection consistency of the KOO method based on GIC

In this section, we demonstrate the model selection consistency property of the KOO method based on the GIC defined by (6). First, we define the notation used in deriving the consistency property. Let  $\mathcal{J}_+$  be the family of sets consisting of the overspecified subsets; that is,  $\mathcal{J}_+$  is given by

$$\mathcal{J}_+ = \{j \subset \omega \mid \boldsymbol{\Theta}_{\bar{j}} - \boldsymbol{\Theta}_j \boldsymbol{\Sigma}_{j_j}^{-1} \boldsymbol{\Sigma}_{j_{\bar{j}}} = \boldsymbol{O}_{k, p-p_j}\}.$$

Then, the true model  $j_*$  is defined as the subset such that the number of elements is smallest for all the subsets in  $\mathcal{J}_+$ ; that is,  $j_*$  is given by  $j_* = \arg \min_{j \in \mathcal{J}_+} p_j$ . The following assumptions are made:

**Assumption A1.**  $n \rightarrow \infty$ ,  $p/n \rightarrow c_1 \in [0, 1)$ ,  $k/n \rightarrow c_2 \in [0, 1 - c_1)$ .

**Assumption A2.** The true subset  $j_*$  is included in the full set  $\omega$ , i.e.,  $j_* \subset \omega$ .

Assumption A1 is our high-dimensional asymptotic framework. This assumption means that  $n$  always tends to infinity, but  $p$  and  $k$  can be fixed or can tend to infinity. Note that  $p_{j_*}$  is finite because we suppose that the number of non-redundant variables is small. Assumption A2 is needed to consider the model's selection consistency  $P(\hat{j} = j_*) \rightarrow 1$ .

Next, we re-form our expression of (6) to obtain the consistency property. For  $j = \ell$  ( $\#(\ell) = p - 1$ ), the divisions of  $\mathbf{Y}$ ,  $\Theta$  and  $\Sigma$  are given by

$$\mathbf{Y} = (\mathbf{Y}_\ell, \mathbf{y}_{\bar{\ell}}), \quad \Theta = (\Theta_\ell, \theta_{\bar{\ell}}), \quad \Sigma = \begin{pmatrix} \Sigma_{\ell\ell} & \sigma_{\ell\bar{\ell}} \\ \sigma'_{\ell\bar{\ell}} & \sigma_{\bar{\ell}\bar{\ell}} \end{pmatrix}.$$

By considering a distributional reduction of  $\text{GIC}(\ell) - \text{GIC}(\omega)$ , the best subset  $\hat{j}$  can be rewritten as follows (the proof is given in Appendix 1).

**Lemma 1.** For  $\ell \subset \omega$  ( $\#(\ell) = p - 1$ ), let  $s_\ell$  be a random variable distributed according to the chi-squared distribution  $\chi^2(n - p - k + 1)$ . Further, let  $h_\ell$  be a random variable such that it is conditionally independent of  $s_\ell$  given  $\mathbf{Y}_\ell$ , and the conditional distribution given  $\mathbf{Y}_\ell$  is the non-central chi-squared distribution  $\chi^2(k; \delta_\ell)$ . Here, the non-centrality parameter  $\delta_\ell$  is defined by

$$\delta_\ell = \sigma_{\bar{\ell}\bar{\ell}, \ell}^{-1} \boldsymbol{\xi}'_\ell \mathbf{X}' (\mathbf{I}_n - \mathbf{P}_{\mathbf{Y}_\ell}) \mathbf{X} \boldsymbol{\xi}_\ell,$$

where  $\sigma_{\bar{\ell}\bar{\ell}, \ell} = \sigma_{\bar{\ell}\bar{\ell}} - \sigma'_{\ell\bar{\ell}} \Sigma_{\ell\ell}^{-1} \sigma_{\ell\bar{\ell}}$  and  $\boldsymbol{\xi}_\ell = \theta_{\bar{\ell}} - \Theta_\ell \Sigma_{\ell\ell}^{-1} \sigma_{\ell\bar{\ell}}$ . Then, the best subset  $\hat{j}$  can be expressed as

$$\hat{j} = \left\{ \bar{\ell} \in \omega \mid \frac{h_\ell}{s_\ell} > \exp\left(\frac{k}{n}\alpha\right) - 1 \right\}.$$

From Lemma 1, we see that the behavior of the non-centrality parameter  $\delta_\ell$  plays an important role in variable selection. Some of the properties of  $\delta_\ell$  are described below (the proof is given in Appendix 4).

**Proposition 1.** For  $\ell \subset \omega$  ( $\#(\ell) = p - 1$ ), the following properties of  $\delta_\ell$  hold.

- (i) If  $\bar{\ell} \notin j_*$ , then  $\boldsymbol{\xi}_\ell = \mathbf{0}_k$ , and hence  $\delta_\ell = 0$ .
- (ii) If  $\bar{\ell} \in j_*$ , then we can express  $\delta_\ell$  under Assumption A1 as follows:

$$\delta_\ell = O_p\left(\text{tr}(\Theta_{j_*} \Sigma_{j_* j_*}^{-1} \Theta'_{j_*}) \lambda_{\max}(\mathbf{X}' \mathbf{X})\right),$$

where  $\lambda_{\max}(\mathbf{A})$  is the maximum eigenvalue of a square matrix  $\mathbf{A}$ , and  $\Theta_{j_*}$  and  $\Sigma_{j_* j_*}$  are the partition matrices  $\Theta_j$  and  $\Sigma_{jj}$  in (2) corresponding to  $j = j_*$ .

From (i) in Proposition 1, the non-centrality parameter  $\delta_\ell$  is 0 when  $\bar{\ell} \notin j_*$ . Moreover, from (ii) it may be expected that  $\delta_\ell$  diverges at about the same speed as the maximum eigenvalue of  $\mathbf{X}' \mathbf{X}$  when  $\bar{\ell} \in j_*$ . In fact, if  $\lambda_{\max}(\mathbf{X}' \mathbf{X}) = O(n)$  and  $\lambda_{\max}(\Theta'_{j_*} \Theta_{j_*}) = O(1)$ , then  $\delta_\ell = O_p(n)$ . However, we cannot determine the exact divergence speed of  $\delta_\ell$  because  $\delta_\ell$  includes the inverse of a non-central Wishart matrix. Hence, we require the following additional assumption.

**Assumption A3.** *There exists  $\tau > 0$  such that  $P(\min_{\bar{\ell} \in j_*} \delta_\ell > \tau n) \rightarrow 1$ .*

Assumption A3 means that  $n^{-1}\delta_\ell$  does not tend to 0 asymptotically in the sense of the convergence in probability, and this will be natural from (ii) in Proposition 1. In the selection contexts for explanatory variables (e.g., [3, 23]), it is assumed that the corresponding non-central parameters have  $O(n)$  and do not converge to 0. Even from that point of view, we consider that Assumption A3 will be appropriate. Moreover, it is possible to relax Assumption A3 and express it as  $P(\min_{\bar{\ell} \in j_*} \delta_\ell > \tau n^\kappa) \rightarrow 1$ , where  $\kappa$  is a constant satisfying  $1/2 < \kappa < 1$ . However, we do not consider such a relaxation, as it makes conditions for the consistency of the KOO method based on the GIC stricter; that is, the relaxation and the conditions are related to the transactions.

Finally, we can demonstrate the consistency property of the KOO method  $\hat{j}$  based on the GIC under Assumptions A1-A3. To do so, we use the following expression for  $\alpha$ :

$$\alpha = \frac{n}{k} \log \left( 1 + \frac{k}{N-2} + \beta \right), \quad \beta > 0, \quad (7)$$

where  $N = n - p - k + 1$ . Then, from Lemma 1 and the definition of the method (6), the lower bound of the probability  $P(\hat{j} = j_*)$  can be derived as

$$\begin{aligned} P(\hat{j} = j_*) &= P \left( \left( \bigcap_{\bar{\ell} \notin j_*} \{ \text{GIC}(\ell) - \text{GIC}(\omega) \leq 0 \} \right) \cap \left( \bigcap_{\bar{\ell} \in j_*} \{ \text{GIC}(\ell) - \text{GIC}(\omega) > 0 \} \right) \right) \\ &\geq 1 - P \left( \bigcup_{\bar{\ell} \notin j_*} \left\{ \frac{h_\ell}{s_\ell} > \beta + \frac{k}{N-2} \right\} \right) - P \left( \bigcup_{\bar{\ell} \in j_*} \left\{ \frac{h_\ell}{s_\ell} \leq \beta + \frac{k}{N-2} \right\} \right), \end{aligned} \quad (8)$$

where  $h_\ell$  and  $s_\ell$  are defined in Lemma 1. Therefore, we look for the condition on  $\beta$  such that the two probabilities in the last expression of (8) tend to 0, in order to have  $\hat{j}$  manifest the consistency property. Such results are derived in Theorem 1 (the proof is given in Appendix 3).

**Theorem 1.** *Suppose that Assumptions A1-A3 hold. Then, the KOO method based on the GIC defined by (6) exhibits model selection consistency, i.e.,  $P(\hat{j} = j_*) \rightarrow 1$  holds, if for some  $r \in \mathbb{N}$  the following condition on  $\alpha$  is satisfied:*

$$\alpha = \frac{n}{k} \log \left( 1 + \frac{k}{N-2} + \beta \right), \quad \beta > 0 \text{ s.t. } \frac{n}{p^{1/2r} k^{1/2}} \beta \rightarrow \infty, \quad \beta \rightarrow 0. \quad (9)$$

Theorem 1 specifies the  $\alpha$  condition necessary for the KOO method based on the GIC to have the consistency property. Thus, we propose a KOO method based on the GIC with  $\alpha$  that satisfies (9). However, we need to decide on a value of  $\alpha$  in order to perform variable selection. An example of  $\alpha$  satisfying (9) is given as follows:

$$\alpha = \tilde{\alpha} = \frac{n}{k} \log \left( 1 + \frac{k}{N-2} + \beta \right), \quad \beta = \frac{p^{1/4} k^{1/2} (\log p) (\log n)}{n}. \quad (10)$$

Note that  $\tilde{\alpha}$  satisfies condition (9) for  $r \geq 2$ . The value  $\beta$  in (10) may satisfy both  $n\beta/(p^{1/4}k^{1/2}) \rightarrow \infty$  and  $\beta \rightarrow 0$  in a well-balanced manner. Although the consistency property in Theorem 1 is obtained by rewriting  $\alpha$  as  $\beta$ , we can derive the consistency property without  $\beta$  if  $k$  is fixed.

**Corollary 1.** *Suppose that Assumptions A1-A3 hold and assume that  $k$  is fixed. Then, the KOO method based on the GIC defined by (6) exhibits consistency, i.e.,  $P(\hat{j} = j_*) \rightarrow 1$  holds, if for some  $r \in \mathbb{N}$  the following condition of  $\alpha$  is satisfied:*

$$\frac{1}{p^{1/2r}}\alpha \rightarrow \infty, \quad \frac{1}{n}\alpha \rightarrow 0. \quad (11)$$

Moreover,  $\alpha$  satisfying condition (9) meets (11) if  $k$  is fixed.

The proof of Corollary 1 is omitted since it is similar to the proof of Theorem 1. From Corollary 1, we see that the AIC, HQC, BIC, and CAIC do not satisfy condition (11) unless  $p$  is fixed or  $p$  diverges at a very slow speed.

## 4 Numerical Studies

In this section, we report the numerical results of a simulation experiment in which the probabilities of selecting the true subset  $j_*$  using the KOO method based on the GIC were determined. Let  $\tilde{\alpha}$  be the  $\alpha$  given by (10), and note that  $\tilde{\alpha}$  satisfies the condition of consistency (9). The probabilities are calculated using Monte Carlo simulations with 10,000 iterations. For comparison, we calculated the results for several KOO methods based on GICs for specific  $\alpha$  values, including the AIC ( $\alpha = 2$ ), HQC ( $\alpha = 2 \log \log n$ ), BIC ( $\alpha = \log n$ ), and CAIC ( $\alpha = 1 + \log n$ ). Moreover, we implemented the KOO method based on the GIC for  $\alpha = n^{1/2}$ , which, from Corollary 1, is consistent when  $k$  is fixed. We set the true subset and the number of true response variables as  $j_* = \{1, 2, 3, 4, 5\}$  and  $p_{j_*} = 5$ , respectively. The response matrix  $\mathbf{Y}$  is generated from  $\mathcal{N}_{n \times p}(\mathbf{X}\Theta, \Sigma \otimes \mathbf{I}_n)$ , where  $\mathbf{X}$ ,  $\Sigma$  and  $\Theta$  are determined as shown below. The explanatory matrix  $\mathbf{X}$  is generated from  $\mathcal{N}_{n \times k}(\mathbf{O}_{n,k}, \Phi \otimes \mathbf{I}_n)$ , where the  $(a, b)$ -th element of the  $k \times k$  matrix  $\Psi$  is  $(0.5)^{|a-b|}$ . The covariance matrix  $\Sigma$  is given by  $\Sigma = 0.4\{(1 - 0.8)\mathbf{I}_p + 0.8\mathbf{1}_p\mathbf{1}_p'\}$ . Let the partition matrices of  $\Theta$  and  $\Sigma$  be as (2) corresponding to  $j = j_*$ ; then the partition matrices of  $\Theta$  are given by

$$\Theta_{j_*} \sim \mathcal{N}_{k \times p_{j_*}}(\mathbf{O}_{k,p_{j_*}}, \mathbf{I}_{p_{j_*}} \otimes \mathbf{I}_k), \quad \Theta_{\bar{j}_*} = \Theta_{j_*} \Sigma_{j_* j_*}^{-1} \Sigma_{j_* \bar{j}_*}.$$

Tables 1 and 2 show the probabilities for selecting the true subset  $j_*$  by KOO methods based on the six criteria, AIC, HQC, BIC, CAIC, GIC with  $\alpha = n^{1/2}$ , and GIC with  $\alpha = \tilde{\alpha}$ , where  $\tilde{\alpha}$  is given by (10). It is evident from the tables that the probabilities for the KOO method based on the GIC with  $\alpha = \tilde{\alpha}$  are high in all cases. This supports our assertion that the method with  $\alpha = \tilde{\alpha}$  is consistent under a high-dimensional asymptotic framework. The KOO methods based on the AIC, HQC, BIC, and CAIC, which do not satisfy (9), have low probabilities in at least some instances. The probabilities associated with the KOO method based on the GIC with  $\alpha = n^{1/2}$  are high when  $k$  is small, as  $\alpha = n^{1/2}$  satisfies condition (11). Based on these results, we recommend using  $\alpha = \tilde{\alpha}$  to find the true subset.

Table 1: True subset selection probabilities (%) by the KOO methods based on the six criteria.

$n$	$p$	$k$	AIC	HQC	BIC	CAIC	$\alpha = n^{1/2}$	$\alpha = \tilde{\alpha}$
100	10	10	68.78	98.25	99.97	100.00	100.00	100.00
300	10	10	81.54	99.86	100.00	100.00	100.00	100.00
500	10	10	84.22	99.93	100.00	100.00	100.00	100.00
1000	10	10	85.15	100.00	100.00	100.00	100.00	100.00
3000	10	10	85.84	100.00	100.00	100.00	100.00	100.00
5000	10	10	85.56	100.00	100.00	100.00	100.00	100.00
10000	10	10	86.18	100.00	100.00	100.00	100.00	100.00
100	40	10	0.00	9.98	84.49	97.62	99.79	99.92
300	120	10	0.00	7.19	96.98	99.70	100.00	100.00
500	200	10	0.00	4.61	98.51	99.88	100.00	100.00
1000	400	10	0.00	2.01	99.63	99.94	100.00	100.00
3000	1200	10	0.00	0.13	99.91	100.00	100.00	100.00
5000	2000	10	0.00	0.03	99.93	100.00	100.00	100.00
10000	4000	10	0.00	0.00	99.99	100.00	100.00	100.00
100	10	40	71.86	99.96	100.00	97.78	0.00	99.54
300	10	120	99.80	100.00	100.00	47.78	0.00	100.00
500	10	200	100.00	100.00	99.48	0.82	0.00	100.00
1000	10	400	100.00	100.00	23.00	0.00	0.00	100.00
3000	10	1200	100.00	100.00	0.00	0.00	0.00	100.00
5000	10	2000	100.00	100.00	0.00	0.00	0.00	100.00
10000	10	4000	100.00	100.00	0.00	0.00	0.00	100.00
100	20	20	19.46	96.12	99.99	100.00	37.67	100.00
300	60	60	25.69	100.00	100.00	100.00	0.00	100.00
500	100	100	56.37	100.00	100.00	100.00	0.00	100.00
1000	200	200	95.81	100.00	100.00	100.00	0.00	100.00
3000	600	600	100.00	100.00	100.00	99.98	0.00	100.00
5000	1000	1000	100.00	100.00	100.00	75.95	0.00	100.00
10000	2000	2000	100.00	100.00	100.00	0.00	0.00	100.00



Table 2: True subset selection probabilities (%) by the KOO methods based on the six criteria.

$n$	$p$	$k$	AIC	HQC	BIC	CAIC	$\alpha = n^{1/2}$	$\alpha = \tilde{\alpha}$
100	80	10	0.00	0.00	0.00	0.00	2.21	34.80
300	240	10	0.00	0.00	0.00	0.00	91.30	99.94
500	400	10	0.00	0.00	0.00	0.00	99.74	100.00
1000	800	10	0.00	0.00	0.00	0.00	100.00	100.00
3000	2400	10	0.00	0.00	0.00	0.00	100.00	100.00
5000	4000	10	0.00	0.00	0.00	0.00	100.00	100.00
10000	8000	10	0.00	0.00	0.00	0.00	100.00	100.00
100	10	80	0.00	27.91	83.72	4.83	0.00	28.62
300	10	240	0.07	100.00	0.00	0.00	0.00	50.04
500	10	400	0.08	100.00	0.00	0.00	0.00	61.44
1000	10	800	0.28	100.00	0.00	0.00	0.00	73.37
3000	10	2400	0.52	100.00	0.00	0.00	0.00	87.57
5000	10	4000	0.47	100.00	0.00	0.00	0.00	91.32
10000	10	8000	0.38	100.00	0.00	0.00	0.00	94.46
100	40	40	0.00	0.10	81.39	85.35	0.00	70.28
300	120	120	0.00	5.61	98.03	24.26	0.00	99.82
500	200	200	0.00	67.13	87.81	0.52	0.00	100.00
1000	400	400	0.00	99.98	7.78	0.00	0.00	100.00
3000	1200	1200	0.00	100.00	0.00	0.00	0.00	100.00
5000	2000	2000	0.00	100.00	0.00	0.00	0.00	100.00
10000	4000	4000	0.00	100.00	0.00	0.00	0.00	100.00

## Appendix 1: Proof of Lemma 1

From (5), the GIC for  $j = \ell$  ( $\#(\ell) = p - 1$ ) is written as

$$\begin{aligned} \text{GIC}(\ell) = & np(\log 2\pi + 1) + n \log \left| \frac{1}{n} \mathbf{Y}'_{\ell} (\mathbf{I}_n - \mathbf{P}_{\mathbf{X}}) \mathbf{Y}_{\ell} \right| + n \log \left\{ \frac{1}{n} \mathbf{y}'_{\bar{\ell}} (\mathbf{I}_n - \mathbf{P}_{\mathbf{Y}_{\ell}}) \mathbf{y}_{\bar{\ell}} \right\} \\ & + \alpha \left\{ \frac{1}{2} p(p+1) + k(p-1) \right\}. \end{aligned}$$

On the other hand, the GIC for the full set  $\omega$  can also be expressed as

$$\begin{aligned} \text{GIC}(\omega) = & np(\log 2\pi + 1) + n \log \left| \frac{1}{n} \mathbf{Y}'_{\ell} (\mathbf{I}_n - \mathbf{P}_{\mathbf{X}}) \mathbf{Y}_{\ell} \right| + n \log \left\{ \frac{1}{n} \mathbf{y}'_{\bar{\ell}} (\mathbf{I}_n - \mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})}) \mathbf{y}_{\bar{\ell}} \right\} \\ & + \alpha \left\{ \frac{1}{2} p(p+1) + kp \right\}. \end{aligned}$$

Hence,  $\text{GIC}(\ell) - \text{GIC}(\omega)$  is expressed by

$$\begin{aligned} \text{GIC}(\ell) - \text{GIC}(\omega) &= n \log \frac{\mathbf{y}'_{\bar{\ell}} (\mathbf{I}_n - \mathbf{P}_{\mathbf{Y}_{\ell}}) \mathbf{y}_{\bar{\ell}}}{\mathbf{y}'_{\bar{\ell}} (\mathbf{I}_n - \mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})}) \mathbf{y}_{\bar{\ell}}} - k\alpha \\ &= n \log \left\{ 1 + \frac{\mathbf{y}'_{\bar{\ell}} (\mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})} - \mathbf{P}_{\mathbf{Y}_{\ell}}) \mathbf{y}_{\bar{\ell}}}{\mathbf{y}'_{\bar{\ell}} (\mathbf{I}_n - \mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})}) \mathbf{y}_{\bar{\ell}}} \right\} - k\alpha. \end{aligned}$$

Let  $\tilde{s}_{\ell} = \sigma_{\bar{\ell}\bar{\ell},\ell}^{-1} \mathbf{y}'_{\bar{\ell}} (\mathbf{I}_n - \mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})}) \mathbf{y}_{\bar{\ell}}$  and  $\tilde{h}_{\ell} = \sigma_{\bar{\ell}\bar{\ell},\ell}^{-1} \mathbf{y}'_{\bar{\ell}} (\mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})} - \mathbf{P}_{\mathbf{Y}_{\ell}}) \mathbf{y}_{\bar{\ell}}$ . We can show that  $\tilde{s}_{\ell}$  and  $\tilde{h}_{\ell}$  satisfy the requirements imposed on  $s_{\ell}$  and  $h_{\ell}$  in Lemma 1, respectively.

First, we consider  $\tilde{s}_{\ell}$ . From a property of a conditional distribution of a multivariate normal distribution (e.g., [21]),  $\mathbf{y}_{\bar{\ell}}$  can be expressed as

$$\mathbf{y}_{\bar{\ell}} = \mathbf{Y}_{\ell} \boldsymbol{\Sigma}_{\ell\ell}^{-1} \boldsymbol{\sigma}_{\ell\bar{\ell}} + \mathbf{X} \boldsymbol{\xi}_{\ell} + \sigma_{\bar{\ell}\bar{\ell},\ell}^{1/2} \mathbf{u}_{\ell},$$

where  $\mathbf{u}_{\ell}$  is independent of  $\mathbf{Y}_{\ell}$  and  $\mathbf{u}_{\ell} \sim \mathcal{N}_n(\mathbf{0}_n, \mathbf{I}_n)$ . Since  $(\mathbf{I}_n - \mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})}) \mathbf{Y}_{\ell} \boldsymbol{\Sigma}_{\ell\ell}^{-1} \boldsymbol{\sigma}_{\ell\bar{\ell}} = \mathbf{0}_n$  and  $(\mathbf{I}_n - \mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})}) \mathbf{X} \boldsymbol{\xi}_{\ell} = \mathbf{0}_n$  hold from a property of the projection matrix, we have

$$\tilde{s}_{\ell} = \mathbf{u}'_{\ell} (\mathbf{I}_n - \mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})}) \mathbf{u}_{\ell}. \quad (12)$$

From Cochran's Theorem (e.g., [8]), we have  $\tilde{s}_{\ell} \sim \chi^2(n - p - k + 1)$ .

Next, we consider  $\tilde{h}_{\ell}$ . From a property of the projection matrix,  $(\mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})} - \mathbf{P}_{\mathbf{Y}_{\ell}}) \mathbf{Y}_{\ell} \boldsymbol{\Sigma}_{\ell\ell}^{-1} \boldsymbol{\sigma}_{\ell\bar{\ell}} = \mathbf{0}_n$  holds. Hence,  $\tilde{h}_{\ell}$  can be expressed as

$$\tilde{h}_{\ell} = \sigma_{\bar{\ell}\bar{\ell},\ell}^{-1} (\mathbf{X} \boldsymbol{\xi}_{\ell} + \sigma_{\bar{\ell}\bar{\ell},\ell}^{1/2} \mathbf{u}_{\ell})' (\mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})} - \mathbf{P}_{\mathbf{Y}_{\ell}}) (\mathbf{X} \boldsymbol{\xi}_{\ell} + \sigma_{\bar{\ell}\bar{\ell},\ell}^{1/2} \mathbf{u}_{\ell}). \quad (13)$$

Since  $\mathbf{Y}_{\ell}$  and  $\mathbf{u}_{\ell}$  are independent, the conditional distribution of  $\tilde{h}_{\ell}$  given  $\mathbf{Y}_{\ell}$  is the non-central chi-squared distribution  $\chi^2(k; \delta_{\ell})$  from Cochran's Theorem, where  $\delta_{\ell}$  is expressed as

$$\delta_{\ell} = \sigma_{\bar{\ell}\bar{\ell},\ell}^{-1} \boldsymbol{\xi}'_{\ell} \mathbf{X}' (\mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})} - \mathbf{P}_{\mathbf{Y}_{\ell}}) \mathbf{X} \boldsymbol{\xi}_{\ell} = \sigma_{\bar{\ell}\bar{\ell},\ell}^{-1} \boldsymbol{\xi}'_{\ell} \mathbf{X}' (\mathbf{I}_n - \mathbf{P}_{\mathbf{Y}_{\ell}}) \mathbf{X} \boldsymbol{\xi}_{\ell}.$$

Finally, we see that  $\tilde{s}_{\ell}$  and  $\tilde{h}_{\ell}$  are conditionally independent given  $\mathbf{Y}_{\ell}$  from (12), (13) and Cochran's Theorem because  $(\mathbf{I}_n - \mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})}) (\mathbf{P}_{(\mathbf{Y}_{\ell}, \mathbf{X})} - \mathbf{P}_{\mathbf{Y}_{\ell}}) = \mathbf{O}_{n,n}$  holds. Therefore, the proof of Lemma 1 is completed.  $\square$

## Appendix 2: Proof of Proposition 1

First, we show (i). Without loss of generality, we can rewrite  $\mathbf{y}$  by  $\mathbf{y} = (\mathbf{y}'_{j_*}, \mathbf{y}'_{\bar{j}_*})' = (\mathbf{y}'_{j_*}, \mathbf{y}'_{\bar{j}_* \cap \ell}, y_{\bar{\ell}})'$  for  $\bar{\ell} \notin j_*$ . Then, the proof of (i) is completed by letting  $\Theta_1 = \Theta_{j_*}$  in the following lemma (the proof is given in Appendix 4).

**Lemma A.1.** *Suppose that the divisions of  $\Theta$  and  $\Sigma$  are given by*

$$\begin{aligned} \Theta &= (\Theta_1, \Theta_2, \Theta_3) = (\Theta_1, \Theta_{(23)}) = (\Theta_{(12)}, \Theta_3), \\ \Sigma &= \begin{pmatrix} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} \\ \Sigma'_{12} & \Sigma_{22} & \Sigma_{23} \\ \Sigma'_{13} & \Sigma'_{23} & \Sigma_{33} \end{pmatrix} = \begin{pmatrix} \Sigma_{11} & \Sigma_{1(23)} \\ \Sigma'_{1(23)} & \Sigma_{(23)(23)} \end{pmatrix} = \begin{pmatrix} \Sigma_{(12)(12)} & \Sigma_{(12)3} \\ \Sigma'_{(12)3} & \Sigma_{33} \end{pmatrix}. \end{aligned}$$

If  $\Theta_{(23)} = \Theta_1 \Sigma_{11}^{-1} \Sigma_{1(23)}$ , then we have  $\Theta_3 = \Theta_{(12)} \Sigma_{(12)(12)}^{-1} \Sigma_{(12)3}$ .

Next, we show (ii). Without loss of generality, we can rewrite  $\mathbf{y}$  by  $\mathbf{y} = (\mathbf{y}'_{j_*}, \mathbf{y}'_{\bar{j}_*})' = (\mathbf{y}'_{j_*}, \mathbf{y}'_{j_* \cap \ell}, y_{\bar{\ell}})'$  for  $\bar{\ell} \in j_*$ . Let the partition matrix of  $\Theta$  and  $\Sigma$  corresponding to the true subset  $j_*$  be as follows:

$$\Theta = (\Theta_{\bar{j}_*}, \Theta_{j_*}), \quad \Sigma = \begin{pmatrix} \Sigma_{\bar{j}_* \bar{j}_*} & \Sigma'_{j_* \bar{j}_*} \\ \Sigma_{j_* \bar{j}_*} & \Sigma_{j_* j_*} \end{pmatrix}.$$

Then, we set the matrices  $\Psi_\ell$ ,  $\Gamma_\ell$ ,  $\Psi$ , and  $\Gamma$  as follows:

$$\begin{aligned} \Psi_\ell &= \Gamma'_\ell \Theta' \Theta \Gamma_\ell, \quad \Gamma_\ell = \begin{pmatrix} \Sigma_{\ell\ell}^{-1/2} & -\sigma_{\bar{\ell}\bar{\ell},\ell}^{-1/2} \Sigma_{\ell\ell}^{-1} \sigma_{\ell\bar{\ell}} \\ \mathbf{0}'_{p-1} & \sigma_{\bar{\ell}\bar{\ell},\ell}^{-1/2} \end{pmatrix}, \\ \Psi &= \Gamma' \Theta' \Theta \Gamma, \quad \Gamma = \begin{pmatrix} \Sigma_{\bar{j}_* \bar{j}_*}^{-1/2} & \mathbf{O}_{p-p_{j_*}, p_{j_*}} \\ -\Sigma_{j_* \bar{j}_*}^{-1} \Sigma_{j_* \bar{j}_*} \Sigma_{\bar{j}_* \bar{j}_*}^{-1/2} & \Sigma_{j_* j_*}^{-1/2} \end{pmatrix}. \end{aligned}$$

Note that the  $(p, p)$ -th element of  $\Psi_\ell$  is  $\sigma_{\bar{\ell}\bar{\ell},\ell}^{-1} \xi'_\ell \xi_\ell$ , and  $\Sigma = (\Gamma'_\ell)^{-1} \Gamma_\ell^{-1} = (\Gamma')^{-1} \Gamma^{-1}$  holds. Moreover,  $\Psi$  can be calculated as

$$\Psi = \begin{pmatrix} \mathbf{O}_{p-p_{j_*}, p-p_{j_*}} & \mathbf{O}_{p-p_{j_*}, p_{j_*}} \\ \mathbf{O}_{p_{j_*}, p-p_{j_*}} & \Sigma_{j_* \bar{j}_*}^{-1/2} \Theta'_{j_*} \Theta_{j_*} \Sigma_{j_* j_*}^{-1/2} \end{pmatrix}.$$

Then, we have

$$\begin{aligned} \delta_\ell &\leq \sigma_{\bar{\ell}\bar{\ell},\ell}^{-1} \xi'_\ell \xi_\ell \lambda_{\max}(\mathbf{X}' \mathbf{X}) \\ &\leq \text{tr}(\Psi_\ell) \lambda_{\max}(\mathbf{X}' \mathbf{X}) \\ &= \text{tr}\{\Psi \Gamma^{-1} \Gamma_\ell \Gamma'_\ell (\Gamma')^{-1}\} \lambda_{\max}(\mathbf{X}' \mathbf{X}) \\ &= \text{tr}\{\Psi \Gamma^{-1} \Sigma^{-1} (\Gamma')^{-1}\} \lambda_{\max}(\mathbf{X}' \mathbf{X}) \\ &= \text{tr}(\Psi) \lambda_{\max}(\mathbf{X}' \mathbf{X}) \\ &= \text{tr}(\Theta_{j_*} \Sigma_{j_* j_*}^{-1} \Theta'_{j_*}) \lambda_{\max}(\mathbf{X}' \mathbf{X}). \end{aligned}$$

Therefore, the proof of (ii) in Proposition 1 is completed.  $\square$

### Appendix 3: Proof of Theorem 1

First, we show  $P(\cup_{\bar{\ell} \notin j_*} \{h_\ell s_\ell^{-1} > \beta + k/(N-2)\}) \rightarrow 0$  in (8). Since  $\delta_\ell = 0$  for  $\bar{\ell} \notin j_*$  from Proposition 1, we see that  $h_\ell \sim \chi^2(k)$  and, moreover,  $s_\ell$  is independent of  $h_\ell$ . Then, for any  $r \in \mathbb{N}$  we have

$$\begin{aligned} P\left(\bigcup_{\bar{\ell} \notin j_*} \left\{\frac{h_\ell}{s_\ell} > \beta + \frac{k}{N-2}\right\}\right) &\leq \sum_{\bar{\ell} \notin j_*} P\left(\frac{h_\ell}{s_\ell} - \frac{k}{N-2} > \beta\right) \\ &\leq (p - p_{j_*})\beta^{-2r} E\left[\left(\frac{h_\ell}{s_\ell} - \frac{k}{N-2}\right)^{2r}\right]. \end{aligned}$$

Then, the order of the moment of the above equation is  $O(k^r n^{-2r})$  from [15, Lemma A.2]. Hence, we have  $P(\cup_{\bar{\ell} \notin j_*} \{h_\ell s_\ell^{-1} > \beta + k/(N-2)\}) = O(pk^r \beta^{-2r} n^{-2r}) \rightarrow 0$ .

Next, we show  $P(\cup_{\bar{\ell} \in j_*} \{h_\ell s_\ell^{-1} \leq \beta + k/(N-2)\}) \rightarrow 0$  in (8). From a property of the non-central chi-squared distribution,  $h_\ell$  can be expressed as follows:

$$h_\ell = \delta_\ell + t_\ell + 2\delta_\ell^{1/2} v_\ell,$$

where  $t_\ell$  and  $v_\ell$  are random variables, and the conditional distributions of those given  $\mathbf{Y}_\ell$  are  $t_\ell | \mathbf{Y}_\ell \sim \chi^2(k)$  and  $v_\ell | \mathbf{Y}_\ell \sim \mathcal{N}(0, 1)$ . Let  $E = \{\min_{\bar{\ell} \in j_*} \delta_\ell > \tau n\}$ , where  $\tau$  is given by Assumption A3. Then, we have

$$\begin{aligned} P\left(\bigcup_{\bar{\ell} \in j_*} \left\{\frac{h_\ell}{s_\ell} \leq \beta + \frac{k}{N-2}\right\}\right) &= P\left(\left\{\bigcup_{\bar{\ell} \in j_*} \left\{\frac{h_\ell}{s_\ell} \leq \beta + \frac{k}{N-2}\right\}\right\} \cap (E \cup E^c)\right) \\ &\leq P\left(\left\{\bigcup_{\bar{\ell} \in j_*} \left\{\frac{h_\ell}{s_\ell} \leq \beta + \frac{k}{N-2}\right\}\right\} \cap E\right) + P(E^c) \\ &\leq \sum_{\bar{\ell} \in j_*} P\left(\tau + n^{-1}t_\ell + 2\tau^{1/2}n^{-1/2}v_\ell \leq n^{-1}s_\ell \left(\beta + \frac{k}{N-2}\right)\right) \\ &\quad + P(E^c). \end{aligned}$$

By calculating the variances of  $t_\ell$ ,  $v_\ell$  and  $(N-2)^{-1}s_\ell$ , we see that  $t_\ell = k + O_p(k^{1/2})$ ,  $v_\ell = O_p(1)$ , and  $(N-2)^{-1}s_\ell = 1 + O_p(n^{-1/2})$  hold. Moreover,  $P(E^c) = o(1)$  holds from Assumption A3. Hence, we have  $P(\cup_{\bar{\ell} \in j_*} \{h_\ell s_\ell^{-1} \leq \beta + k/(N-2)\}) \rightarrow 0$ . Therefore, the proof of Theorem 1 is completed.  $\square$

### Appendix 4: Proof of Lemma A.1

From the assumption of Lemma A.1, the following equations hold:

$$\Theta_2 = \Theta_1 \Sigma_{11}^{-1} \Sigma_{12}, \quad \Theta_3 = \Theta_1 \Sigma_{11}^{-1} \Sigma_{13}. \quad (14)$$

By using (14) and the general formula for the inverse of a block matrix,  $\Theta_{(12)}\Sigma_{(12)(12)}^{-1}\Sigma_{(12)3}$  is expanded as

$$\begin{aligned}
\Theta_{(12)}\Sigma_{(12)(12)}^{-1}\Sigma_{(12)3} &= (\Theta_1, \Theta_2) \begin{pmatrix} \Sigma_{11}^{-1} + \Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22.1}^{-1}\Sigma'_{12}\Sigma_{11}^{-1} & -\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22.1}^{-1} \\ -\Sigma_{22.1}^{-1}\Sigma'_{12}\Sigma_{11}^{-1} & \Sigma_{22.1}^{-1} \end{pmatrix} \begin{pmatrix} \Sigma_{13} \\ \Sigma_{23} \end{pmatrix} \\
&= \Theta_1\Sigma_{11}^{-1}\Sigma_{13} + \Theta_1\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22.1}^{-1}\Sigma'_{12}\Sigma_{11}^{-1}\Sigma_{13} - \Theta_2\Sigma_{22.1}^{-1}\Sigma'_{12}\Sigma_{11}^{-1}\Sigma_{13} \\
&\quad - \Theta_1\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22.1}^{-1}\Sigma_{23} + \Theta_2\Sigma_{22.1}^{-1}\Sigma_{23} \\
&= \Theta_1\Sigma_{11}^{-1}\Sigma_{13} + \Theta_1\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22.1}^{-1}\Sigma'_{12}\Sigma_{11}^{-1}\Sigma_{13} - \Theta_1\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22.1}^{-1}\Sigma'_{12}\Sigma_{11}^{-1}\Sigma_{13} \\
&\quad - \Theta_1\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22.1}^{-1}\Sigma_{23} + \Theta_1\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22.1}^{-1}\Sigma_{23} \\
&= \Theta_1\Sigma_{11}^{-1}\Sigma_{13} \\
&= \Theta_3,
\end{aligned}$$

where  $\Sigma_{22.1} = \Sigma_{22} - \Sigma'_{12}\Sigma_{11}^{-1}\Sigma_{12}$ . Therefore, the proof of Lemma A.1 is completed.  $\square$

## Acknowledgments

This work was supported by funding from JSPS KAKENHI (grant numbers JP19K21672, JP20K14363 and JP20H04151 to Ryoya Oda; and JP16H03606, JP18K03415 and JP20H04151 to Hirokazu Yanagihara).

## 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 & F. Csáki), pp. 267–281. Akadémiai Kiadó, Budapest. doi:10.1007/978-1-4612-1694-0\_15
- [2] AKAIKE, H. (1974). A new look at the statistical model identification. *IEEE Trans. Automatic Control*, **AC-19**, 716–723. doi:10.1109/TAC.1974.1100705
- [3] BAI, Z. D., FUJIKOSHI, Y. & HU, J. (2018). Strong consistency of the AIC, BIC,  $C_p$  and KOO methods in high-dimensional multivariate linear regression. TR No. 18–09, *Statistical Research Group*, Hiroshima University.
- [4] BOZDOGAN, H. (1987). Model selection and Akaike’s information criterion (AIC): The general theory and its analytical extensions. *Psychometrika*, **52**, 345–370. doi:10.1007/BF02294361
- [5] BROWN, P. J. (1982). Multivariate calibration. *J. R. Statist. Soc. Ser. B*, **44**, 287–321. doi:10.1111/j.2517-6161.1982.tb01209.x
- [6] FUJIKOSHI, Y. & NISHII, R. (1986). Selection of variables in a multivariate inverse regression problem. *Hiroshima Math. J.*, **13**, 269–277. doi:10.32917/hmj/1206130428
- [7] FUJIKOSHI, Y. & SAKURAI, T. (2019). Consistency of test-based method for selection of variables in high-dimensional two-group discriminant analysis. *Jpn. J. Stat. Data Sci.*, **2**, 155–171. doi:10.1007/s42081-019-00032-4

- [8] FUJIKOSHI, Y., ULYANOV, V. V. & SHIMIZU, R. (2010). *Multivariate Statistics: High-Dimensional and Large-Sample Approximations*. John Wiley & Sons, Inc., Hoboken, New Jersey.
- [9] HANNAN, E. J. & QUINN, B. G (1979). The determination of the order of an autoregression. *J. Roy. Statist. Soc. Ser. B*, **26**, 270–273. doi:10.1111/j.2517-6161.1979.tb01072.x
- [10] NAKAGAWA, I., WATANABE, H. & HYODO, M. (2021). Kick-one-out-based variable selection method for Euclidean distance-based classifier in high-dimensional settings. *J. Multivariate Anal.*, **184**, 104756. doi:10.1016/j.jmva.2021.104756
- [11] NISHII, R. (1984). Asymptotic properties of criteria for selection of variables in multiple regression. *Ann. Statist.*, **12**, 758–765. doi:10.1214/aos/1176346522
- [12] NISHII, R., BAI, Z. D. & KRISHNAIAH, P. R. (1988). Strong consistency of the information criterion for model selection in multivariate analysis. *Hiroshima Math. J.*, **18**, 451–462. doi:10.32917/hmj/1206129611
- [13] ODA, R., MIMA, Y., YANAGIHARA, H. & FUJIKOSHI, Y. (2020). A high-dimensional bias-corrected AIC for selecting response variables in multivariate calibration. *Comm. Statist. Theory Methods*, **50**, 3453–3476. doi:10.1080/03610926.2019.1705978
- [14] ODA, R., SUZUKI, Y., YANAGIHARA, H. & FUJIKOSHI, Y. (2020). A consistent variable selection method in high-dimensional canonical discriminant analysis. *J. Multivariate Anal.*, **175**, 104561. doi:10.1016/j.jmva.2019.104561
- [15] ODA, R. & YANAGIHARA, H. (2020). A fast and consistent variable selection method for high-dimensional multivariate linear regression with a large number of explanatory variables. *Electron. J. Statist.*, **14**, 1386–1412. doi:10.1214/20-EJS1701
- [16] ODA, R. & YANAGIHARA, H. (2021). A consistent likelihood-based variable selection method in normal multivariate linear regression. *Smart Innov. Syst. Tec.*, **238**, 391–401. doi:10.1007/978-981-16-2765-1\_33
- [17] SAKURAI, T. & FUJIKOSHI, Y. (2018). Consistency of distance-based criterion for selection of variables in high-dimensional two-group discriminant analysis. TR No. 18–05, *Statistical Research Group*, Hiroshima University.
- [18] SAKURAI, T. & FUJIKOSHI, Y. (2020). Exploring consistencies of information criterion and test based criterion for high dimensional multivariate regression models under three covariance structures. *Recent Developments in Multivariate and Random Matrix Analysis*, 313–334. doi:10.1007/978-3-030-56773-6\_18
- [19] SCHWARZ, G. (1978). Estimating the dimension of a model. *Ann. Statist.*, **6**, 461–464. doi:10.1214/aos/1176344136
- [20] SRIVASTAVA, M. S. (2002). *Methods of Multivariate Statistics*. John Wiley & Sons, New York.

- [21] SRIVASTAVA, M. S. & KHATRI, C. G. (1979). *An Introduction to Multivariate Statistics*. North-Holland, New York.
- [22] TIMM, N. H. (2002). *Applied Multivariate Analysis*. Springer-Verlag, New York.
- [23] YANAGIHARA, H. (2016). A high-dimensionality-adjusted consistent  $C_p$ -type statistic for selecting variables in a normality-assumed linear regression with multiple responses. *Procedia Comput. Sci.*, **96**, 1096–1105. doi:10.1016/j.procs.2016.08.151
- [24] ZHAO, L. C., KRISHNAIAH, P. R. & BAI, Z. D. (1986). On detection of the number of signals in presence of white noise. *J. Multivariate Anal.*, **20**, 1–25. doi:10.1016/0047-259X(86)90017-5