

## Consistent variable selection criteria in multivariate linear regression even when dimension exceeds sample size

Ryoya ODA

(Received Xxx 00, 0000)

ABSTRACT. This paper is concerned with the selection of explanatory variables in multivariate linear regression. The Akaike's information criterion and the  $C_p$  criterion cannot perform in high-dimensional situations such that the dimension of a vector stacked with response variables exceeds the sample size. To overcome this, we consider two variable selection criteria based on an  $L_2$  squared distance with a weighted matrix, namely the scalar-type generalized  $C_p$  criterion and the ridge-type generalized  $C_p$  criterion. We clarify conditions for their consistency under a hybrid-ultra-high-dimensional asymptotic framework such that the sample size always goes to infinity but the number of response variables may go to infinity. Numerical experiments show that the probabilities of selecting the true subset by criteria satisfying consistency conditions are high even when the dimension is larger than the sample size. Finally, we illuminate the practical utility of these criteria using empirical data.

### 1. Introduction

Multivariate linear regression is an important and very widely used inferential statistical methodology. It is the cornerstone of many theoretical and applied statistics textbooks (see, e.g., Srivastava, 2002, chap 9; Timm, 2002, chap 4) and it has widespread applications in many fields. Let  $\mathbf{Y} = (\mathbf{y}_{(1)}, \dots, \mathbf{y}_{(n)})'$  be an  $n \times p$  observation matrix stacking individual  $p$  response variables, and  $\mathbf{X} = (\mathbf{x}_{(1)}, \dots, \mathbf{x}_{(n)})'$  be an  $n \times k$  observation matrix stacking individual non-stochastic  $k$  explanatory variables, where  $n$  is the sample size. Note that  $\mathbf{X}$  may include the intercept term that the column vector is  $\mathbf{1}_n$ , where  $\mathbf{1}_n$  is an  $n$ -dimensional vector of ones. Assume that  $\text{rank}(\mathbf{X}) = k < n$  to ensure the existence of variable selection criteria used in this paper. We consider linear regression for  $n$  samples of a vector of individual  $p$  response variables and  $k$  explanatory variables on  $\{(\mathbf{y}'_{(i)}, \mathbf{x}'_{(i)})' \mid i = 1, \dots, n\}$ . Then, the multivariate linear regression is written as

$$\mathbf{Y} = \mathbf{X}\Theta + \mathcal{E},$$

---

The author is supported financially by Research Fellowships of the Japan Society for the Promotion of Science for Young Scientists.

2010 *Mathematics Subject Classification*. Primary 62J05; Secondary 62H12.

*Key words and phrases*. Hybrid-ultra-high-dimensional asymptotic framework, Multivariate linear regression, Non-normality, Selection consistency, Variable selection criterion.

where  $\Theta$  is a  $k \times p$  unknown matrix of regression coefficients, and each row of an  $n \times p$  error matrix  $\mathcal{E}$  is identically distributed with a mean vector  $\mathbf{0}_p$ , which is a  $p$ -dimensional vector of zeros, and a covariance matrix  $\Sigma$ .

In actual data analysis contexts, it is important to specify salient explanatory variables affecting response variables. In multivariate linear regression, this is regarded as the problem of selecting the best subset of explanatory variables. Variable selection criteria are widely used in empirical contexts to choose the best subset of explanatory variables. The Akaike's information criterion (AIC) (Akaike, 1973; 1974) and the  $C_p$  criterion (Sparks *et al.*, 1983) which is a multivariate version of Mallows'  $C_p$  criterion (Mallows, 1973; 1995) are well-known examples in this respect. The AIC and  $C_p$  criterion are estimators of risk functions, the Kullback-Leibler loss function and the mean squared prediction error standardized by the true covariance matrix, respectively. Further, as extensions of the AIC and  $C_p$  criterion, the generalized information criterion (GIC) and the generalized  $C_p$  ( $GC_p$ ) criterion were proposed by Nishii *et al.* (1988) and Nagai *et al.* (2012), respectively. The GIC and  $GC_p$  criterion were generalized from the AIC and  $C_p$  criterion by replacing "2" (the penalty term for model complexity) with any positive number. Note that the GIC includes the AIC, the Bayesian information criterion (BIC) proposed by Schwarz (1978), a consistent AIC (CAIC) proposed by Bozdogan (1987), and the Hannan-Quinn information criterion (HQC) proposed by Hannan and Quinn (1979). Further, the  $GC_p$  criterion includes the  $C_p$  criterion and the modified  $C_p$  ( $MC_p$ ) criterion proposed by Fujikoshi and Satoh (1997).

Importantly, there are increasing demands in recent years vis-a-vis analyzing high-dimensional data such that  $p$  exceeds  $n$  (for an example, see Wille *et al.*, 2004). For high-dimensional cases, we need a variable selection criterion which can be operationalized even when  $p > n$ . However, note that the GIC consists of the logarithm of the determinant of the sample covariance matrix, and the  $GC_p$  criterion consists of the inverse matrix of the sample covariance matrix. Therefore, since the sample covariance matrix becomes singular when  $p$  is larger than  $n$ , more precisely  $n - k < p$ , the GIC always gives  $-\infty$  and the  $GC_p$  criterion cannot be defined when  $p > n$ . However, criteria proposed by Fujikoshi *et al.* (2011), Yamamura *et al.* (2010), and Kubokawa and Srivastava (2012) are calculable even when  $p > n$ . Fujikoshi *et al.* (2011) proposed the prediction error (PE) criterion based on the mean squared prediction error. Yamamura *et al.* (2010) and Kubokawa and Srivastava (2012) proposed criteria using a ridge-type sample covariance matrix as an estimator of the true covariance matrix. Moreover, their criteria are exact or asymptotically unbiased estimators of risk functions under some conditions.

In this paper, we consider consistency as one of the asymptotic properties of variable selection criteria. In a given variable selection context, the desired outcome is to specify explanatory variables which substantively affect the response variable according to the nature and extent of available empirical data. In other words, it is hoped that the true subset of variables is identified as the

best subset by variable selection. Since we do not know the true subset, we use a variable selection criterion to maximize the probability of selecting the true subset. When the probability that the chosen subset is the true subset approaches 1, consistency is assured, i.e., the following equation holds:

$$P(\hat{j} = j_*) \rightarrow 1,$$

where  $\hat{j}$  is the best subset according to a variable selection criterion and  $j_*$  is the true subset. It is expected that a consistent variable selection criterion has a high probability of selecting the true subset when the amount of data is sufficient. Therefore, consistency is an important property of a variable selection criterion. In the context of  $n > p$ , assuming that the true distribution of the error vector is the multivariate normal distribution, Fujikoshi *et al.* (2014) and Yanagihara *et al.* (2015) obtained the consistency properties of criteria such as the AIC and  $C_p$  criterion. They used a moderate-high-dimensional asymptotic framework such that both  $n$  and  $p$  go to  $\infty$  but  $p$  does not exceed  $n$ . Moreover, Yanagihara *et al.* (2015) also used an asymptotic framework defined by adding  $k/n \rightarrow 0$  to the moderate-high-dimensional asymptotic framework. Relaxing the normality assumption, Yanagihara (2015) dealt with conditions for consistency of the GIC under the moderate-high-dimensional asymptotic framework. Under the normality assumption, Yanagihara (2016) obtained conditions for consistency of the  $GC_p$  criterion under a hybrid-moderate-high-dimensional asymptotic framework such that  $n$  goes to  $\infty$  and  $p$  may go to  $\infty$  but  $p/n$  converges to some positive constant included in  $[0, 1)$ . Relaxing the normality assumption, Yanagihara (2019) focused on conditions for consistency of the GIC and  $GC_p$  criterion under the hybrid-moderate-high-dimensional asymptotic framework. As such, therein,  $p$  does not exceed  $n$ . On the other hand, in the context where  $p > n$ , Katayama and Imori (2014) considered variable selection criteria based on a lasso-type estimation for the inverse of the covariance matrix. Under the normality assumption, they showed that the criteria are consistent in a restricted-ultra-high-dimensional asymptotic framework such that both  $n$  and  $p$  go to infinity but  $p$  may exceed  $n$  and  $\log p/n \rightarrow 0$  while  $k/n \rightarrow 0$ .

The aim of this paper is to obtain conditions for consistency of variable selection criteria (which are introduced in subsection 2.1) under non-normality and a high-dimensional asymptotic framework such that  $n$  goes to infinity but  $p$  may exceed  $n$ . To obtain conditions for consistency, the following hybrid-ultra-high-dimensional (HUHD) asymptotic framework is mainly used:

$$\text{HUHD} : n \rightarrow \infty, p/n \rightarrow c \in [0, \infty], k: \text{fixed},$$

where  $c = \infty$  means that  $p/n$  goes to  $\infty$ . The HUHD asymptotic framework has two key characteristics. First, the divergence speed of  $p$  is not restricted, hence this asymptotic framework incorporates an asymptotic framework such that  $n$  and  $p$  go to  $\infty$  but  $p$  may be larger than  $n$ , namely the ultra-high-dimensional (UHD) asymptotic framework, which is written as

$$\text{UHD} : (n, p) \rightarrow \infty, p/n \rightarrow c \in [0, \infty], k: \text{fixed}.$$

Second, the HUHD asymptotic framework also includes the large-sample asymptotic framework such that only  $n$  tends to  $\infty$ . From this, it is expected that consistent variable selection criteria under the HUHD asymptotic framework select the true subset with high probability regardless of the size of  $p$ .

The remainder of the paper is organized as follows. In section 2, we present the necessary notation and assumptions to clarify conditions for consistency. In section 3, we obtain conditions for consistency. In section 4, for the purposes of verification, we conduct numerical experiments and illuminate the practical utility of consistent criteria by using real data examples. Technical details are provided in the Appendix.

## 2. Preliminaries

**2.1. Models and Criteria.** Suppose that  $j$  denotes a subset of  $\omega = \{1, \dots, k\}$  containing  $k_j$  elements, and  $\mathbf{X}_j$  denotes an  $n \times k_j$  matrix consisting of columns of  $\mathbf{X}$  indexed by elements of  $j$ , where  $k_A$  is the number of elements in a set  $A$  denoted by  $k_A = \#(A)$ . For example, if  $j = \{1, 2, 4\}$ , then  $\mathbf{X}_j$  consists of the first, second, and fourth column vectors of  $\mathbf{X}$ . Then, the candidate model  $M_j$  with  $k_j$  explanatory variables from subset  $j$  is expressed as follows:

$$M_j : \mathbf{Y} = \mathbf{X}_j \boldsymbol{\Theta}_j + \boldsymbol{\mathcal{E}}_j, \quad (1)$$

where  $\boldsymbol{\Theta}_j$  is a  $k_j \times p$  unknown matrix of regression coefficients, and each row of  $\boldsymbol{\mathcal{E}}_j$  is identically distributed with a mean of  $\mathbf{0}_p$  and a covariance matrix  $\boldsymbol{\Sigma}_j$ . Let  $j_*$  ( $\subset \omega$ ) be a true subset, and assume that the data are generated from the following true model  $M_{j_*}$  with  $k_{j_*}$  true explanatory variables:

$$M_{j_*} : \mathbf{Y} = \mathbf{X}_{j_*} \boldsymbol{\Theta}_* + \boldsymbol{\mathcal{E}}_*,$$

where  $\boldsymbol{\Theta}_*$  is a  $k_{j_*} \times p$  unknown matrix of true regression coefficients and  $\boldsymbol{\mathcal{E}}_* = (\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_n)'$  is an  $n \times p$  true error matrix. Assume that  $\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_n$  are identically distributed according to a distribution of  $\boldsymbol{\varepsilon}$  with

$$E[\boldsymbol{\varepsilon}] = \mathbf{0}_p, \quad Cov[\boldsymbol{\varepsilon}] = \boldsymbol{\Sigma}_*, \quad E[||\boldsymbol{\varepsilon}||^4] < \infty,$$

where  $||\boldsymbol{\varepsilon}||^2 = \boldsymbol{\varepsilon}'\boldsymbol{\varepsilon}$  and  $\boldsymbol{\Sigma}_*$  is a  $p \times p$  true unknown covariance matrix. Although it is typical to assume independence of  $\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_n$ , here we assume a moment condition which relaxes independence; specifically, we assume that for any  $i \neq j$ ,  $\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_n$  are satisfied with the following moment condition:

$$\begin{aligned} E[\boldsymbol{\varepsilon}_i \boldsymbol{\varepsilon}_j'] &= E[\boldsymbol{\varepsilon}_i] E[\boldsymbol{\varepsilon}_j'], \quad E[||\boldsymbol{\varepsilon}_i||^2 ||\boldsymbol{\varepsilon}_j||^2] = E[||\boldsymbol{\varepsilon}_i||^2] E[||\boldsymbol{\varepsilon}_j||^2], \\ E[\boldsymbol{\varepsilon}_i \boldsymbol{\varepsilon}_i' \boldsymbol{\varepsilon}_j \boldsymbol{\varepsilon}_j'] &= E[\boldsymbol{\varepsilon}_i \boldsymbol{\varepsilon}_i'] E[\boldsymbol{\varepsilon}_j \boldsymbol{\varepsilon}_j']. \end{aligned}$$

Note that the above moment condition is similar to assuming independence. Without loss of generality, we sort column vectors of  $\mathbf{X}$  as  $\mathbf{X} = (\mathbf{X}_{j_*}, \mathbf{X}_{j_*^c})$ , where set  $A^c$  denotes the compliment of set  $A$ . Moreover, for expository purposes, we represent  $\mathbf{X}_{j_*}$ ,  $\mathbf{X}_\omega$ ,  $k_{j_*}$  and  $k_\omega$  as  $\mathbf{X}_*$ ,  $\mathbf{X}$ ,  $k_*$ , and  $k$ , respectively.

We consider two variable selection criteria based on the following weighted  $L_2$  squared distance:

$$d(\mathbf{A}, \mathbf{B}|\mathbf{G}) = \text{tr}\{(\mathbf{A} - \mathbf{B})\mathbf{G}^{-1}(\mathbf{A} - \mathbf{B})'\},$$

where  $\mathbf{G}$  is a positive definite matrix. Let  $\mathbf{S}_j$  be an estimator of  $\boldsymbol{\Sigma}_j$  in the candidate model  $M_j$ , which is given by

$$\mathbf{S}_j = \frac{1}{n - k_j} \mathbf{Y}'(\mathbf{I}_n - \mathbf{P}_j)\mathbf{Y},$$

where  $\mathbf{I}_n$  is the  $n \times n$  identity matrix, and  $\mathbf{P}_j$  is the projection matrix to the subspace spanned by the columns of  $\mathbf{X}_j$ , i.e.,  $\mathbf{P}_j = \mathbf{X}_j(\mathbf{X}_j'\mathbf{X}_j)^{-1}\mathbf{X}_j'$ . Then, the minimum value of  $d(\mathbf{Y}, \mathbf{X}_j\boldsymbol{\Theta}_j|\mathbf{G})$  for  $\boldsymbol{\Theta}_j$  is expressed as

$$\min_{\boldsymbol{\Theta}_j} d(\mathbf{Y}, \mathbf{X}_j\boldsymbol{\Theta}_j|\mathbf{G}) = \text{tr}\{\mathbf{Y}'(\mathbf{I}_n - \mathbf{P}_j)\mathbf{Y}\mathbf{G}^{-1}\} = (n - k_j)\text{tr}(\mathbf{S}_j\mathbf{G}^{-1}). \quad (2)$$

The minimum value in (2) expresses a measurement about the goodness of fit for model  $M_j$ . Using (2) in the candidate model  $M_j$ , the following class of variable selection criteria is considered:

$$\mathcal{L}(j|\alpha, \mathbf{G}) = (n - k_j)\text{tr}(\mathbf{S}_j\mathbf{G}^{-1}) + \alpha pk_j, \quad (3)$$

where  $\alpha$  is a positive constant which expresses the complexity of the model  $M_j$ . It is straightforward that (3) with  $\alpha = 2$  and  $\mathbf{G} = \mathbf{S}_\omega$  is the  $C_p$  criterion proposed by Sparks *et al.* (1983) when  $n > p$ . Moreover, (3) with  $\mathbf{G} = \mathbf{S}_\omega$  is the  $GC_p$  criterion proposed by Nagai *et al.* (2012). However, the  $GC_p$  criterion cannot be defined when  $p > n$ . Therefore, we consider two criteria obtained by substituting one of two specific weighted matrices instead of  $\mathbf{S}_\omega$  into  $\mathbf{G}$  in (3). By substituting the scalar matrix  $p^{-1}\text{tr}(\mathbf{S}_\omega)\mathbf{I}_p$  into  $\mathbf{G}$ , we define the scalar-type generalized  $C_p$  ( $SGC_p$ ) criterion as follows:

$$SGC_p(j|\alpha) = p^{-1}\mathcal{L}(j|\alpha, p^{-1}\text{tr}(\mathbf{S}_\omega)\mathbf{I}_p) = (n - k_j)\frac{\text{tr}(\mathbf{S}_j)}{\text{tr}(\mathbf{S}_\omega)} + \alpha k_j. \quad (4)$$

Note that the  $SGC_p(j|\alpha)$  criterion is obtained by dividing  $\mathcal{L}(j|\alpha, p^{-1}\text{tr}(\mathbf{S}_\omega)\mathbf{I}_p)$  by  $p$  because the divided  $p$  is redundant for variable selection. The  $SGC_p$  criterion with  $\alpha = 2$  is essentially the same as the PE criterion proposed by Fujikoshi *et al.* (2011). Moreover, the value  $\text{tr}(\mathbf{S}_j)/\text{tr}(\mathbf{S}_\omega)$  in (4) corresponds to the MANOVA test statistic in Fujikoshi *et al.* (2004). They applied the Dempster trace criterion when  $p > n$  for tests about one and two sample mean vectors in Dempster (1958; 1960). Note that there is no inverse of the sample covariance matrix in the  $SGC_p$  criterion. Thus, this criterion is calculable even when  $p > n$ . Let  $\mathbf{S}_\lambda$  be the ridge-type sample covariance matrix, which is defined by

$$\mathbf{S}_\lambda = \mathbf{S}_\omega + \frac{\text{tr}(\mathbf{S}_\omega)}{\lambda}\mathbf{I}_p,$$

where  $\lambda$  is a positive ridge parameter. Then, by substituting  $\mathbf{S}_\lambda$  into  $\mathbf{G}$ , we define the ridge-type generalized  $C_p$  ( $RGC_p$ ) criterion as follows:

$$RGC_p(j|\alpha, \lambda) = \mathcal{L}(j|\alpha, \mathbf{S}_\lambda) = (n - k_j)\text{tr}(\mathbf{S}_j\mathbf{S}_\lambda^{-1}) + \alpha pk_j. \quad (5)$$

The first term in (5) is similar to that of the ridge-type  $C_p$  criterion used by Kubokawa and Srivastava (2012). If  $\mathbf{S}_\omega$  is invertible and  $\lambda = \infty$ , then (5) coincides with the  $GC_p$  criterion. However,  $\mathbf{S}_\omega$  is singular when  $p > n$ . The scalar matrix  $\lambda^{-1}\text{tr}(\mathbf{S}_\omega)\mathbf{I}_p$  keeps  $\mathbf{S}_\lambda$  invertible even in such case. The best subsets are given by minimizing the  $SGC_p$  criterion and  $RGC_p$  criterion, i.e., defined by

$$\hat{j}_S = \arg \min_{j \in \mathcal{J}} SGC_p(j|\alpha), \quad \hat{j}_R = \arg \min_{j \in \mathcal{J}} RGC_p(j|\alpha, \lambda), \quad (6)$$

where  $\mathcal{J}$  is a family of subsets of  $\omega$  denoted by  $\mathcal{J} = \{j_1, \dots, j_K\}$  and  $K$  is the number of candidate subsets.

**2.2. Assumptions for Consistency.** We prepare assumptions for consistency. To describe several classes of  $j$  that express the column indexes of  $\mathbf{X}$  in the candidate model (1), we separate  $\mathcal{J}$  into two sets, one is a family of over-specified subsets that includes the true subset, i.e.,  $\mathcal{J}_+ = \{j \in \mathcal{J} | j_* \subset j\}$ , and the other is a family of underspecified subsets that are not over-specified subsets, i.e.,  $\mathcal{J}_- = \mathcal{J}_+^c \cap \mathcal{J}$ . Let a  $p \times p$  non-centrality matrix and parameter be expressed by

$$\Delta_j = \Theta_*' \mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j) \mathbf{X}_* \Theta_*, \quad \delta_j^2 = \text{tr}(\Delta_j). \quad (7)$$

It should be noted that  $\Delta_j = \mathbf{O}_{p,p}$  and  $\delta_j^2 = 0$  hold from properties of projection matrices if and only if  $j \in \mathcal{J}_+$ , where  $\mathbf{O}_{p,p}$  is the  $p \times p$  matrix with all elements as zero. Then, we prepare the following assumptions for consistency:

- A1. The true subset  $j_*$  is included in  $\mathcal{J}$ , i.e.,  $j_* \in \mathcal{J}$ .
- A2.  $\limsup_{p \rightarrow \infty} \frac{1}{p} \text{tr}(\Sigma_*) < \infty$ .
- A3.  $\limsup_{p \rightarrow \infty} \frac{\kappa_4}{\text{tr}(\Sigma_*)^2} < \infty$ , where  $\kappa_4 = E[||\varepsilon||^4] - \text{tr}(\Sigma_*)^2 - 2\text{tr}(\Sigma_*^2)$ .
- A4. For every  $j \in \mathcal{J}_-$ , there exists  $\ell \in j_* \cap j^c$  such that

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \mathbf{x}'_\ell (\mathbf{I}_n - \mathbf{P}_{\omega_\ell}) \mathbf{x}_\ell > 0, \quad \liminf_{p \rightarrow \infty} \frac{1}{p} ||\boldsymbol{\theta}_\ell||^2 > 0,$$

where  $\omega_\ell = \{\ell\}^c$ , and  $\mathbf{x}_\ell$  and  $\boldsymbol{\theta}_\ell$  are the  $\ell$ -th column vectors of  $\mathbf{X}_*$  and  $\Theta'_*$ , respectively.

Assumption A1 is needed to consider consistency. From the definition of  $\mathcal{J}_+$ , the true subset  $j_*$  can be regarded as the smallest over-specified subset. Assumption A2 is a regularity assumption for the true covariance matrix  $\Sigma_*$ . If the number of response variables whose variances are  $O(p)$  is finite and the variances of the other response variables are  $O(1)$ , assumption A2 holds. Assumption A3 is the restriction for the fourth moment of  $\varepsilon$ . From properties of the multivariate normal distribution (e.g., Magnus and Neudecker, 1979; Himeno and Yamada, 2014),  $\kappa_4 = 0$  when  $\varepsilon$  is distributed according to the multivariate normal distribution. Moreover, some specific multivariate distributions such as the multivariate  $t$ -distribution or the multivariate contaminated normal distribution are satisfied with assumption A3. Assumption A4 concerns explanatory variables and true regression coefficients. In terms of explanatory variables, this means that a sample

covariance of residuals in the linear regression of  $\mathbf{x}_\ell$  with the remaining  $\mathbf{X}_{\omega_\ell}$  does not converge to 0. It is straightforward to show that this is weaker than assuming  $\liminf_{n \rightarrow \infty} n^{-1} \lambda_{\min}(\mathbf{X}'\mathbf{X}) > 0$ , where  $\lambda_{\min}(\mathbf{A})$  is the minimum eigenvalue of a square matrix  $\mathbf{A}$ . The assumption for true regression coefficients is essentially used in Katayama and Imori (2014). For example, when all the elements of each  $\boldsymbol{\theta}_\ell$  are non-zero constants not converging to 0, the assumption for true regression coefficients holds. Moreover, even when half of the elements of  $\boldsymbol{\theta}_\ell$  are zeros and the remaining half are non-zero constants not converging to 0, the assumption is satisfied. Hence, the assumption for the true regression coefficients will be not unrealistic. Further, if  $p$  diverges as fast as  $n$ , i.e.,  $c \in [0, \infty)$  in the HUHD asymptotic framework, the assumption for true regression coefficients can become weaker such as  $\liminf_{p \rightarrow \infty} q_p^{-1} \|\boldsymbol{\theta}_\ell\|^2 > 0$  for some  $q_p \rightarrow \infty$  ( $p \rightarrow \infty$ ). Note that assumption A4 does not always have to hold for every  $\ell \in j_*$ . For example, if  $\mathcal{J}$  is a set of nested subsets, i.e.,  $\mathcal{J} = \{\{1\}, \dots, \{1, \dots, k\}\}$ , then assumption A4 needs to hold only for  $\ell = k_*$ . If assumption A4 is supported, for every  $j \in \mathcal{J}_-$ , the following inequality holds (the proof is given in Appendix A):

$$\inf_{n > k, p \geq 1} \frac{1}{np} \lambda_{\max}(\boldsymbol{\Delta}_j) > 0, \quad (8)$$

where  $\lambda_{\max}(\mathbf{A})$  is the maximum eigenvalue of a square matrix  $\mathbf{A}$ .

Furthermore, we consider the following assumption that is regarded as a special case of assumption A3:

$$\text{A3}'. \quad \lim_{p \rightarrow \infty} \frac{\xi^2}{\text{tr}(\boldsymbol{\Sigma}_*)^2} = 0, \quad \text{where } \xi^2 = \max\{\kappa_4, \text{tr}(\boldsymbol{\Sigma}_*^2)\}.$$

Assumption A3' is used under the UHD asymptotic framework, and this assumption is stronger than assumption A3. For example, assumption A3' is satisfied if the following conditions hold:

$$\begin{aligned} \lim_{p \rightarrow \infty} \frac{\text{tr}(\boldsymbol{\Sigma}_*^2)}{\text{tr}(\boldsymbol{\Sigma}_*)^2} &= 0, \quad \boldsymbol{\varepsilon} = \boldsymbol{\Sigma}_*^{1/2} \mathbf{u}, \quad \mathbf{u} = (u_1, \dots, u_p)', \\ E[u_a] &= 0, \quad E[u_a^4] \leq r_u \quad (a = 1, \dots, p), \\ E[u_a^2 u_b^2] &= 1 \quad (a \neq b), \quad E[u_a u_b u_c u_d] = 0 \quad (a \neq b, c, d), \end{aligned} \quad (9)$$

where  $r_u$  is a positive constant not dependent on  $p$ . When  $\boldsymbol{\varepsilon} = \boldsymbol{\Sigma}_*^{1/2} \mathbf{u}$ ,  $\kappa_4$  is calculated as follows:

$$\kappa_4 = \sum_{a=1}^p \{(\boldsymbol{\Sigma}_*)_{aa}\}^2 (E[u_a^4] - 3) \leq |r_u - 3| \text{tr}(\boldsymbol{\Sigma}_*^2),$$

where  $(\mathbf{A})_{ab}$  expresses the  $(a, b)$ -th element of a matrix  $\mathbf{A}$ . The condition about the true covariance matrix  $\lim_{p \rightarrow \infty} \text{tr}(\boldsymbol{\Sigma}_*^2)/\text{tr}(\boldsymbol{\Sigma}_*)^2 = 0$  is called the sphericity condition, and it is often used for  $p \gg n$  setting (e.g., Aoshima *et al.*, 2018).

### 3. Main Results

**3.1. Conditions for Consistency of the  $SGC_p$  Criterion.** We obtain conditions for consistency of the  $SGC_p$  criterion (4). Recall that the best subset chosen by minimizing the  $SGC_p$  criterion is defined by (6). Then, the  $SGC_p$  criterion is consistent if  $P(\hat{j}_S = j_*) \rightarrow 1$ . The probability  $P(\hat{j}_S = j_*)$  can be expressed as

$$P(\hat{j}_S = j_*) = P\left(\bigcap_{j \in \mathcal{J} \cap \{j_*\}^c} \{SGC_p(j|\alpha) > SGC_p(j_*|\alpha)\}\right).$$

We separate  $\mathcal{J} \cap \{j_*\}^c$  into  $\mathcal{J}_+ \cap \{j_*\}^c$  and  $\mathcal{J}_-$  because the non-centrality matrix  $\Delta_j$  in (7) behaves differently for each of the cases of  $j \in \mathcal{J}_+ \cap \{j_*\}^c$  and  $j \in \mathcal{J}_-$ . From this and the subadditivity of a measure, a lower bound of  $P(\hat{j}_S = j_*)$  is written as

$$P(\hat{j}_S = j_*) \geq 1 - \bar{P}_S - \underline{P}_S,$$

where  $\bar{P}_S$  and  $\underline{P}_S$  are defined by

$$\bar{P}_S = P\left(\bigcup_{j \in \mathcal{J}_+ \cap \{j_*\}^c} \{SGC_p(j|\alpha) \leq SGC_p(j_*|\alpha)\}\right), \quad (10)$$

$$\underline{P}_S = P\left(\bigcup_{j \in \mathcal{J}_-} \{SGC_p(j|\alpha) \leq SGC_p(j_*|\alpha)\}\right). \quad (11)$$

To obtain conditions for consistency of the  $SGC_p$  criterion, we consider conditions such that  $\bar{P}_S$  and  $\underline{P}_S$  converge to 0. First, we prepare the results about the orders of several probabilities. For subsets  $j, h \subset \omega$ , let  $\mathbf{W}$ ,  $\mathbf{U}_j$ , and  $\mathbf{V}_{j,h}$  be random matrices defined by

$$\mathbf{W} = \boldsymbol{\varepsilon}'_*(\mathbf{I}_n - \mathbf{P}_\omega)\boldsymbol{\varepsilon}_*, \quad \mathbf{U}_j = \boldsymbol{\Theta}'_*\mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j)\boldsymbol{\varepsilon}_*, \quad \mathbf{V}_{j,h} = \boldsymbol{\varepsilon}'_*(\mathbf{P}_j - \mathbf{P}_h)\boldsymbol{\varepsilon}_*. \quad (12)$$

Then, we derive the following lemma about the orders of the tail probabilities for functions of (12) (the proof is given in Appendix B).

LEMMA 1. *Let  $\mathbf{W}$ ,  $\mathbf{U}_j$ , and  $\mathbf{V}_{j,h}$  be given by (12), and let  $r_1 > 0$ ,  $r_2 > 0$ ,  $r_3 < 0$ ,  $r_4 > 0$ ,  $r_5 > 0$ , and  $r_6 > 0$ . Then, under the HUHD asymptotic framework, the following results hold:*

(i) *If  $r_1 > \text{tr}(\boldsymbol{\Sigma}_*)$  and  $r_2 < \text{tr}(\boldsymbol{\Sigma}_*)$ , then we have*

$$P\left((n-k)^{-1}\text{tr}(\mathbf{W}) \geq r_1\right) = O\left(\xi^2 n^{-1}\{r_1 - \text{tr}(\boldsymbol{\Sigma}_*)\}^{-2}\right),$$

$$P\left((n-k)^{-1}\text{tr}(\mathbf{W}) \leq r_2\right) = O\left(\xi^2 n^{-1}\{\text{tr}(\boldsymbol{\Sigma}_*) - r_2\}^{-2}\right),$$

*where  $\xi^2$  is given in assumption A3'.*

(ii) *For  $j \not\supseteq j_*$ , we have*

$$P\left(\text{tr}(\mathbf{U}_j) \leq r_3\right) = O\left(\text{tr}(\boldsymbol{\Sigma}_*\Delta_j)|r_3|^{-2}\right),$$

*where  $\Delta_j$  is defined by (7).*

(iii) *For  $j \supseteq h$ , if  $r_4 > \text{tr}(\boldsymbol{\Sigma}_*)$ , then we have*

$$P\left(\text{tr}(\mathbf{V}_{j,h}) \geq (k_j - k_h)r_4\right) = O\left(\xi^2\{r_4 - \text{tr}(\boldsymbol{\Sigma}_*)\}^{-2}\right).$$

(iv) *For  $j \supseteq h$ , if  $r_6/r_5 \rightarrow 0$ , then we have*

$$P\left(\text{tr}(\mathbf{V}_{j,h}) - (k_j - k_h)\text{tr}(\boldsymbol{\Sigma}_*) + r_5 \leq r_6\right) = O\left(\xi^2 r_5^{-2}\right).$$

By using Lemma 1, we give the orders of  $\bar{P}_S$  and  $\underline{P}_S$  (the proof is given in Appendix C).

LEMMA 2. *Suppose that assumptions A1, A2, and A4 hold, and for some constants  $\tau_S$  satisfying  $0 < \tau_S < 1$ , the following equations hold:*

$$\lim_{n \rightarrow \infty, p/n \rightarrow c} \alpha \tau_S > 1, \quad \lim_{n \rightarrow \infty, p/n \rightarrow c} n^{-1} \alpha = 0, \quad (13)$$

*under the HUHD asymptotic framework. Then, the orders of  $\bar{P}_S$  and  $\underline{P}_S$  defined in (10) and (11) are given by*

$$\begin{aligned} \bar{P}_S &= O\left(\xi^2 \text{tr}(\Sigma_*)^{-2} \max\{(\alpha \tau_S - 1)^{-2}, n^{-1}(1 - \tau_S)^{-2}\}\right), \\ \underline{P}_S &= O\left(\xi^2 \text{tr}(\Sigma_*)^{-2} \max\{(\alpha \tau_S - 1)^{-2}, n^{-1}(1 - \tau_S)^{-2}\}\right) \\ &\quad + O\left(\max\{\xi^2 n^{-2} p^{-2}, \xi^2 \text{tr}(\Sigma_*)^{-2} n^{-1}, \lambda_{\max}(\Sigma_*) n^{-1} p^{-1}\}\right), \end{aligned}$$

where  $\xi^2$  is defined in assumption A3'.

Next, we obtain conditions for consistency of the  $SGC_p$  criterion (4). Note that the results in Lemma 2 are derived without assumptions A3 and A3'. We use assumption A3 or A3' to obtain consistency conditions, although the UHD asymptotic framework is used when assumption A3' is supported. It is straightforward that  $\limsup_{p \rightarrow \infty} \xi \text{tr}(\Sigma_*)^{-1} < \infty$  holds under assumption A3, but  $\lim_{p \rightarrow \infty} \xi \text{tr}(\Sigma_*)^{-1} = 0$  holds under assumption A3'. By using this fact and Lemma 2, we obtain consistency conditions about  $\alpha$  (the proof is given in Appendix D).

THEOREM 1. *Suppose that assumptions A1, A2, A3, and A4 hold. Then, the  $SGC_p$  criterion is consistent under the HUHD asymptotic framework if the following conditions are satisfied:*

$$\lim_{n \rightarrow \infty, p/n \rightarrow c} \alpha = \infty, \quad \lim_{n \rightarrow \infty, p/n \rightarrow c} \frac{\alpha}{n} = 0. \quad (14)$$

*Furthermore, when replacing assumption A3 with assumption A3', the  $SGC_p$  criterion is consistent under the UHD asymptotic framework if the following conditions are satisfied:*

$$\lim_{(n,p) \rightarrow \infty, p/n \rightarrow c} \alpha > 1, \quad \lim_{(n,p) \rightarrow \infty, p/n \rightarrow c} \frac{\alpha}{n} = 0. \quad (15)$$

From Theorem 1, if assumption A3' is supported, the  $SGC_p$  criterion is consistent under the UHD asymptotic framework even when  $\alpha$  is a constant not dependent on  $n$  and  $p$  such as  $\alpha = 2$ . When assumption A3' is not supported but assumption A3 is,  $\alpha$  should diverge to render the  $SGC_p$  criterion consistent. Moreover, if (14) holds, then (15) holds. It is difficult to verify whether assumption A3' holds using empirical data. Hence, we recommend that (14) be used to render the  $SGC_p$  criterion consistent by deciding  $\alpha$ . On the other hand, we also obtain conditions for inconsistency (the proof is given in Appendix E).

**THEOREM 2.** *Suppose that assumptions A1, A2, A3, and A4 hold. Let conditions of  $\alpha$  under the HUHD asymptotic framework be as follows:*

C1.  $\lim_{n \rightarrow \infty, p/n \rightarrow c} \alpha < 1$  and there exists  $j \in \mathcal{J}_+ \cap \{j_*\}^c$  such that

$$\lim_{n \rightarrow \infty, p/n \rightarrow c} \frac{\kappa_4 I(\kappa_4 > 0) + 2\text{tr}(\boldsymbol{\Sigma}_*^2)}{(1 - \alpha)^2 \text{tr}(\boldsymbol{\Sigma}_*)^2} < k_j - k_*, \quad (16)$$

where  $I(\kappa_4 > 0)$  is an indicator function, i.e., if  $\kappa_4 > 0$  then  $I(\kappa_4 > 0) = 1$ , otherwise  $I(\kappa_4 > 0) = 0$ .

C2. There exists  $j \subsetneq j_*$  such that

$$\lim_{n \rightarrow \infty, p/n \rightarrow c} \frac{\alpha \text{tr}(\boldsymbol{\Sigma}_*)}{\delta_j^2} > (k_* - k_j)^{-1}.$$

Then, if either of the conditions C1 or C2 is satisfied, the  $SGC_p$  criterion is inconsistent, i.e.,  $\lim_{n \rightarrow \infty, p/n \rightarrow c} P(\hat{j}_S = j_*) < 1$  holds under the HUHD asymptotic framework. Furthermore, when replacing assumption A3 with assumption A3', (16) and  $\lim_{(n,p) \rightarrow \infty, p/n \rightarrow c} P(\hat{j}_S = j_*) = 0$  always hold under the UHD asymptotic framework if  $\lim_{(n,p) \rightarrow \infty, p/n \rightarrow c} \alpha < 1$ .

We observe that the  $SGC_p$  criterion is inconsistent when  $\alpha$  is too small from condition C1 or too large from condition C2. Although we cannot cover all the consistency or inconsistency conditions of  $\alpha$  from only Theorems 1 and 2, these theorems nevertheless provide much information about the consistency or inconsistency of the  $SGC_p$  criterion.

**3.2. Conditions for Consistency of the  $RGC_p$  Criterion.** We obtain conditions for consistency of the  $RGC_p$  criterion (5). In the same way as subsection 3.1, a lower bound of  $P(\hat{j}_R = j_*)$  is written as

$$P(\hat{j}_R = j_*) \geq 1 - \bar{P}_R - \underline{P}_R,$$

where  $\bar{P}_R$  and  $\underline{P}_R$  are given by

$$\bar{P}_R = P(\cup_{j \in \mathcal{J}_+ \cap \{j_*\}^c} \{RGC_p(j|\alpha, \lambda) \leq RGC_p(j_*|\alpha, \lambda)\}), \quad (17)$$

$$\underline{P}_R = P(\cup_{j \in \mathcal{J}_-} \{RGC_p(j|\alpha, \lambda) \leq RGC_p(j_*|\alpha, \lambda)\}). \quad (18)$$

First, we obtain the orders of  $\bar{P}_R$  and  $\underline{P}_R$ . Then, we examine the orders by using moments of a statistic. It is difficult to calculate the moments of  $\mathbf{a}'\mathbf{S}_\lambda^{-1}\mathbf{a}$  because of the existence of the inverse matrix of  $\mathbf{S}_\lambda$ , where  $\mathbf{a}$  is a  $p$ -dimensional vector. Therefore, we do not evaluate  $\mathbf{a}'\mathbf{S}_\lambda^{-1}\mathbf{a}$  directly, but evaluate the following lower and upper bounds:

$$\|\mathbf{a}\|^2 \lambda_{\min}(\mathbf{S}_\lambda^{-1}) \leq \mathbf{a}'\mathbf{S}_\lambda^{-1}\mathbf{a} \leq \|\mathbf{a}\|^2 \lambda_{\max}(\mathbf{S}_\lambda^{-1}). \quad (19)$$

By using (19) and Lemma 1, we give the orders of  $\bar{P}_R$  and  $\underline{P}_R$  (the proof is given in Appendix F).

LEMMA 3. *Suppose that assumptions A1, A2, and A4 hold, and for some constants  $\tau_R$  satisfying  $0 < \tau_R < 1$  the following equations hold:*

$$\lim_{n \rightarrow \infty, p/n \rightarrow c} \lambda^{-1} p \alpha \tau_R > 1, \quad \lim_{n \rightarrow \infty, p/n \rightarrow c} n^{-1} (1 + \lambda^{-1}) p \alpha = 0,$$

*under the HUHD asymptotic framework. Then, the orders of  $\bar{P}_R$  and  $\underline{P}_R$  defined in (17) and (18) are given by*

$$\begin{aligned} \bar{P}_R &= O\left(\xi^2 \text{tr}(\Sigma_*)^{-2} \max\{(\lambda^{-1} p \alpha \tau_R - 1)^{-2}, n^{-1} (1 - \tau_R)^{-2}\}\right), \\ \underline{P}_R &= O\left(\xi^2 \text{tr}(\Sigma_*)^{-2} \max\{(\lambda^{-1} p \alpha \tau_R - 1)^{-2}, n^{-1} (1 - \tau_R)^{-2}\}\right) \\ &\quad + O\left(\max\{\xi^2 n^{-2} p^{-2}, \xi^2 \text{tr}(\Sigma_*)^{-2} n^{-1}, \lambda_{\max}(\Sigma_*) n^{-1} p^{-1}\}\right), \end{aligned}$$

*where  $\xi^2$  is defined in assumption A3'.*

By using Lemma 3, we obtain consistency conditions of the  $RGC_p$  criterion. Since the  $RGC_p$  criterion has the two parameters  $\alpha$  and  $\lambda$ , the conditions are connected with  $\alpha$  and  $\lambda$ .

THEOREM 3. *Suppose that assumptions A1, A2, A3, and A4 hold. Then, the  $RGC_p$  criterion is consistent under the HUHD asymptotic framework if the following conditions are satisfied:*

$$\lim_{n \rightarrow \infty, p/n \rightarrow c} \frac{p\alpha}{\lambda} = \infty, \quad \lim_{n \rightarrow \infty, p/n \rightarrow c} \frac{(1 + \lambda^{-1})p\alpha}{n} = 0. \quad (20)$$

*Furthermore, when replacing assumption A3 with assumption A3', the  $RGC_p$  criterion is consistent under the UHD asymptotic framework if the following conditions are satisfied:*

$$\lim_{(n,p) \rightarrow \infty, p/n \rightarrow c} \frac{p\alpha}{\lambda} > 1, \quad \lim_{(n,p) \rightarrow \infty, p/n \rightarrow c} \frac{(1 + \lambda^{-1})p\alpha}{n} = 0. \quad (21)$$

The proof of Theorem 3 is omitted because the theorem can be proved in the same way as Theorem 1. From Theorem 3, if we set  $\lambda = 1$  and  $\alpha = \tilde{\alpha}/p$  ( $\tilde{\alpha} > 0$ ), conditions (20) and (21) are the same as (14) and (15), respectively. Note that conditions (20) and (21) may be strong because they are derived using inequality (19). From Theorem 3, we observe that the larger  $\lambda$  becomes, the larger  $\alpha$  should be, to satisfy conditions (20) and (21). Furthermore, we also obtain conditions for inconsistency (the proof is given in Appendix G).

THEOREM 4. *Suppose that assumptions A1, A2, A3, and A4 hold. Let conditions of  $\alpha$  under the HUHD asymptotic framework be as follows:*

- C3.  $\lim_{n \rightarrow \infty, p/n \rightarrow c} (1 + \lambda^{-1}) p \alpha < 1$  and there exists  $j \in \mathcal{J}_+ \cap \{j_*\}^c$  such that

$$\lim_{n \rightarrow \infty, p/n \rightarrow c} \frac{\kappa_4 I(\kappa_4 > 0) + 2 \text{tr}(\Sigma_*^2)}{\{1 - (1 + \lambda^{-1}) p \alpha\}^2 \text{tr}(\Sigma_*)^2} < k_j - k_*. \quad (22)$$

Table 1. Assumptions and asymptotic behaviors of  $n$  and  $p$  to ensure consistency of six criteria.

Criterion	Assumptions	Asymptotic behavior
1	A1, A2, A3', A4	$p \rightarrow \infty$
2	A1, A2, A3, A4	free
3	A1, A2, A3, A4	$\log \log p / \log n \rightarrow 0$
4	A1, A2, A3', A4	$p \rightarrow \infty$
5	A1, A2, A3, A4	free
6	A1, A2, A3, A4	$\log \log p / \log n \rightarrow 0$

C4. There exists  $j \subsetneq j_*$  such that

$$\lim_{n \rightarrow \infty, p/n \rightarrow c} \frac{p \alpha \text{tr}(\Sigma_*)}{\lambda \delta_j^2} > (k_* - k_j)^{-1}.$$

Then, if either of the conditions C3 or C4 is satisfied, the  $RGC_p$  criterion is inconsistent, i.e.,  $\lim_{n \rightarrow \infty, p/n \rightarrow c} P(\hat{j}_R = j_*) < 1$  holds under the HUHD asymptotic framework. Furthermore, when replacing assumption A3 with assumption A3', (22) and  $\lim_{(n,p) \rightarrow \infty, p/n \rightarrow c} P(\hat{j}_R = j_*) = 0$  always hold under the UHD asymptotic framework if  $\lim_{(n,p) \rightarrow \infty, p/n \rightarrow c} (1 + \lambda^{-1})p\alpha < 1$ .

From Theorem 4, we observe that  $\lambda$  should be large in order not to satisfy conditions C3 and C4. However, if  $\lambda$  is large,  $p\alpha\lambda^{-1}$  in (20) and (21) is small and then the condition of  $\alpha$  to have consistency becomes restricted.

## 4. Numerical Experiments

**4.1. Criteria for Numerical Experiments.** To conduct numerical experiments, we use the following six criteria:

Criterion 1: the  $SGC_p$  criterion with  $\alpha = 2$ .

Criterion 2: the  $SGC_p$  criterion with  $\alpha = \log n$ .

Criterion 3: the  $SGC_p$  criterion with  $\alpha = (\log n / \log \log p)^{1/2}$ .

Criterion 4: the  $RGC_p$  criterion with  $\alpha = 2p^{-1}$  and  $\lambda = 1$ .

Criterion 5: the  $RGC_p$  criterion with  $\alpha = p^{-1} \log n$  and  $\lambda = 1$ .

Criterion 6: the  $RGC_p$  criterion with  $\alpha = p^{-1}(n \log n / \log \log p)^{1/2}$  and  $\lambda = n^{1/2}$ .

Table 1 shows the assumptions and asymptotic behaviors of  $n$  and  $p$  to ensure the consistency of the above six criteria. We observe that to ensure consistency,  $p$  has to diverge for criteria 1 and 4, but  $p$  does not have to diverge for criteria 2, 3, 5, and 6. Further, criteria 3 and 6 are consistent when  $\log \log p / \log n \rightarrow 0$ . Since this slightly restricts the behavior of  $p$ , it may not be suitable where  $p$  increases dramatically. However, such a case is unrealistic, so this behavior is reasonable for empirical contexts. Note that the penalty terms  $k_j\alpha$  or  $k_j p\alpha$  in criteria 1, 2, 4, and 5 do not include  $p$ , but those in criteria 3 and 6 do.

For comparison, we also consider criteria in Katayama and Imori (2014) given by

$$\text{HGIC}(j) = p + \log |(1 - k_j/n) \mathbf{D}_{\mathbf{S}_j}| + \beta p k_j,$$

where  $\mathbf{D}_{\mathbf{S}_j} = \text{diag}\{(\mathbf{S}_j)_{11}, \dots, (\mathbf{S}_j)_{pp}\}$  and  $\text{diag}\{(\mathbf{A})_{11}, \dots, (\mathbf{A})_{pp}\}$  is the diagonal matrix with diagonal elements corresponding to those of a  $p \times p$  matrix  $\mathbf{A}$ . Especially, we use the following three HGICs from their paper:

Criterion 7: the HGIC with  $\beta = n^{-1}(\log p)(\log \log p)^{1/2}$ .

Criterion 8: the HGIC with  $\beta = n^{-1}(\log p)(\log \log p)$ .

Criterion 9: the HGIC with  $\beta = n^{-1}(\log p)(\log \log p)^{3/2}$ .

From Katayama and Imori (2014), criteria 7, 8, and 9 are consistent under several assumptions such as normality when  $p \rightarrow \infty$  and  $\log p/n \rightarrow 0$  for our numerical studies.

**4.2. Simulations.** We verify the foregoing exposition by simulations. The probabilities of selecting the true subset  $j_*$  were evaluated by Monte Carlo simulations with 10,000 iterations. Ten subsets  $j_m = \{1, \dots, m\}$  ( $m = 1, \dots, 10$ ), with several different values of  $n$  and  $p$ , were prepared for these simulations. We generated the explanatory matrix  $\mathbf{X}$  as follows. We independently generated  $s_1, \dots, s_n$  from  $U(-1, 1)$ , where  $U(a, b)$  denotes a uniform distribution on the range  $(a, b)$ . Using  $s_1, \dots, s_n$ , we constructed an  $n \times k$  matrix of explanatory variables  $\mathbf{X}$ , where the  $(a, b)$ -th element is defined by  $s_a^{b-1}$  ( $a = 1, \dots, n$ ;  $b = 1, \dots, k$ ). The true subset was determined by  $j_* = \{1, 2, 3, 4, 5\}$ . The true coefficient matrix  $\Theta_*$  adhered to the following structure:

$$\Theta_* = (\theta_1, \dots, \theta_{k_*})', \quad \theta_a = \begin{cases} \left( a(-1)^{a-1} \mathbf{1}'_{\lfloor p/2 \rfloor}, \mathbf{0}'_{\lceil p/2 \rceil} \right)' & (a : \text{odd}) \\ \left( \mathbf{0}'_{\lfloor p/2 \rfloor}, a(-1)^{a-1} \mathbf{1}'_{\lceil p/2 \rceil} \right)' & (a : \text{even}) \end{cases},$$

where  $\lfloor \cdot \rfloor$  and  $\lceil \cdot \rceil$  are the floor and ceiling functions, respectively. For these numerical simulations, we expressed  $\mathcal{E}_*$  as  $\mathbf{Z}_* \Sigma_*^{1/2}$ , where  $\mathbf{Z}_* = (\mathbf{z}_1, \dots, \mathbf{z}_n)'$  and  $\mathbf{z}_1, \dots, \mathbf{z}_n$  are independent and identically distributed from  $\mathbf{z} = (z_1, \dots, z_p)'$  with mean  $\mathbf{0}_p$  and covariance matrix  $\mathbf{I}_p$ . Let  $\boldsymbol{\nu} = (\nu_1, \dots, \nu_p)'$ ,  $\boldsymbol{\zeta} = (\zeta_1, \dots, \zeta_p)'$   $\sim$  *i.i.d.*  $N_p(\mathbf{0}_p, \mathbf{I}_p)$ , and  $\tau \sim \chi^2(10)$  be mutually independent random vectors and variable. Then,  $\mathbf{z}$  is generated from the following four distributions:

(D1) multivariate normal distribution:  $\mathbf{z} = \boldsymbol{\nu}$ .

(D2) multivariate  $t$ -distribution with 10 degrees of freedom:  $\mathbf{z} = (8/\tau)^{1/2} \boldsymbol{\nu}$ .

(D3) independent skew-normal distribution with shape parameter 10:

$$z_a = \left( 1 - \frac{2}{\pi} \eta^2 \right)^{-1/2} \left( \frac{\nu_a}{\sqrt{1 + 10^2}} + \eta |\zeta_a| - \sqrt{\frac{2}{\pi}} \eta \right) \quad (a = 1, \dots, p),$$

where  $\eta = 10/\sqrt{1 + 10^2}$ .

(D4) independent log-normal distribution:

$$z_a = \frac{\exp(\nu_a) - \sqrt{e}}{\sqrt{e(e-1)}} \quad (a = 1, \dots, p).$$

Note that distributions (D1)-(D4) are satisfied with  $\kappa_4 = O(\text{tr}(\boldsymbol{\Sigma}_*^2))$ . The true covariance matrix  $\boldsymbol{\Sigma}_*$  was set as the following two structures:

(S1) exchangeable structure with correlation 0.8:

$$\boldsymbol{\Sigma}_* = (1 - 0.8)\mathbf{I}_p + 0.8\mathbf{1}_p\mathbf{1}_p'.$$

(S2) autoregressive structure with correlation 0.8:  $(\boldsymbol{\Sigma}_*)_{ab} = (0.8)^{|a-b|}$ .

Note that assumption A3' is not satisfied when the true covariance matrix  $\boldsymbol{\Sigma}_*$  is (S1), but assumption A3' is satisfied when the true covariance matrix  $\boldsymbol{\Sigma}_*$  is (S2) under distributions (D1)-(D4). Under these settings, we used the 8 combinations of the four distributions and the two true covariance matrices (S1) and (S2). Tables 2-9 show the probabilities of selecting the true subset  $j_*$  using each of the nine criteria. In each table, the probabilities of selecting the true subset  $j_*$  were evaluated for distributions (D1)-(D4) and the two covariance matrices (S1) and (S2). When the true covariance matrix  $\boldsymbol{\Sigma}_*$  has an exchangeable structure, i.e., in Tables 2, 4, 6, and 8, it appears that criteria 2, 5, and 6 are consistent for both cases where only  $n$  is large and where  $n$  and  $p$  are large, but criteria 1 and 4 are not consistent. This is because assumption A3 is satisfied for the cases of (S1) and distributions (D1)-(D4), but assumption A3' is not satisfied for such cases. Moreover, although criterion 3 is consistent from Table 1, it looks inconsistent in Tables 2, 4, 6, and 8. This is because the penalty term in criterion 3 is smaller than that in criterion 1 for our numerical simulations. On the other hand, when the true covariance matrix  $\boldsymbol{\Sigma}_*$  has an autoregressive structure, i.e., in Tables 3, 5, 7, and 9, we observe that criteria 1 and 4 also are consistent except for the case that only  $n$  is large because (S2) is satisfied with  $\lim_{p \rightarrow \infty} \text{tr}(\boldsymbol{\Sigma}_*^2)/\text{tr}(\boldsymbol{\Sigma}_*)^2 = 0$ , so assumption A3' is satisfied for the cases of (S2) and distributions (D1)-(D4). This result accords with Theorem 1 and Theorem 3. In Tables 2-9, criteria 7, 8, and 9 are consistent when  $n$  and  $p$  are large, but they are not consistent when only  $n$  is large. Further, we observe that the probabilities by criteria 7, 8, and 9 are low when  $p/n = 10$  and  $n \leq 100$ . In sum, the probabilities by criterion 6 are the highest across Tables 2-9.

**4.3. Empirical Examples.** First, we verify the probabilities of selecting the true subsets by using real data. The dataset pertains to 8 groups ( $g = 1, \dots, 8$ ) of black cotton fibers dyed by Indigo and its derivative dyes. Each cotton fiber has 55 samples, and each sample has 541 variables, which are the absorbances for wavelengths from 240 nm to 780 nm in steps of 1 nm. Let the explanatory matrix be denoted as  $\mathbf{X} = (\mathbf{T}, \mathbf{1}_9) \otimes \mathbf{1}_{25}$ , where  $\mathbf{T} = (\mathbf{e}_1, \dots, \mathbf{e}_8)$  and  $\mathbf{e}_a$  ( $a = 1, \dots, 8$ ) is a 9-dimensional vector such that the  $(a + 1)$ -th element is one and the other elements are zeros, and the symbol  $\otimes$  denotes the Kronecker product (see, e.g., Harville, 1997). Here, the 9-th column vector of  $\mathbf{X}$  expresses the intercept term.









Moreover, let the family of candidate subsets be all of the subsets included in the intercept term, i.e.,  $\mathcal{J} = \{j \in \mathfrak{P}(\{1, \dots, 9\}) \mid j \cap \{9\} \neq \emptyset\}$ , where  $\mathfrak{P}(A)$  is the power set of a set  $A$ . Then, for each group  $b = 1, \dots, 8$ , we carried out the following two steps:

- Step 1. Let  $\mathbf{U}_g$  ( $g = 1, \dots, 8$ ) be the  $25 \times 541$  response matrices by random sampling without replacement from group  $g$ . Further, let  $\mathbf{U}_{9,b}$  be the  $25 \times 541$  response matrices by random sampling without replacement from the remaining samples in group  $b$ . Then, the response matrix is constructed as  $\mathbf{Y}_b = (\mathbf{U}'_1, \dots, \mathbf{U}'_8, \mathbf{U}'_{9,b})'$ .
- Step 2. Let the coefficient matrix  $\Theta_b$  given by  $\Theta_b = (\theta_{1,b}, \dots, \theta_{8,b}, \theta_{9,b})'$ . Then, apply multivariate linear regression with  $\mathbf{X}$  and  $\Theta_b$  to the response matrix  $\mathbf{Y}_b$ , and choose the best subset by performing variable selection from the explanatory variables excepting the intercept, i.e., from the elements of  $\mathcal{J}$ .

From steps 1 and 2, we have  $n = 225$ ,  $p = 541$ , and  $k = 9$  in this example. Note that  $\theta_{b,b}$  should be  $\mathbf{0}_p$  and the remainder should not be  $\mathbf{0}_p$ , because  $\mathbf{U}_{9,b}$  is extracted from the same group as  $\mathbf{U}_b$ . Hence, we know that the true subset is  $j_{*,b} = \{1, \dots, 9\} \cap \{b\}^c$  when  $\mathbf{Y}_b$  is used as the response matrix. Moreover, to increase calculation speed, instead of a variable selection method such as (6), we used the best subset  $\tilde{j}$  by the following method:

$$\tilde{j} = \{\ell \in \omega \mid \text{SC}(\omega_\ell) > \text{SC}(\omega)\}, \quad (23)$$

where  $\text{SC}(j)$  expresses the value of a variable selection criterion (SC) for model  $M_j$ , and  $\omega_\ell$  is defined in assumption A4. The selection method as per (23) was proposed by Zhao *et al.* (1986). From Nishii *et al.* (1988), it is known that when  $k$  is fixed, a criterion under (23) is consistent if the criterion under the selection method such as (6) is consistent. For these settings, we iterated steps 1 and 2 10,000 times for each group  $b = 1, \dots, 8$ . Table 10 shows the probabilities of selecting the true subset by the nine criteria for each group  $b = 1, \dots, 8$ . We observe that the probabilities by criterion 6 are highest except where  $b = 5, 6$ . However, all nine criteria have very low probabilities where  $b = 5, 6$ . This is because groups 5 and 6 are very similar. Actually, letting  $\bar{\mathbf{y}}_g$  be the sample mean vector of group  $g$ , we have  $\|\bar{\mathbf{y}}_5 - \bar{\mathbf{y}}_6\| \doteq 0.46$  but  $\|\bar{\mathbf{y}}_g - \bar{\mathbf{y}}_h\| \geq 1.60$  for the cases of  $g, h \neq 5, 6$  ( $g \neq h$ ). Hence, groups 5 and 6 will be very similar on average. Moreover, criterion 6 selected  $\{1, \dots, 9\} \cap \{5, 6\}^c$  as the best subset for many iterations when  $b = 5, 6$ .

Next, we provide an example of variable selection using empirical data from Wille *et al.* (2004) as well as Yamamura *et al.* (2010). There are 795 genes which may exhibit associations with 39 genes from two biosynthesis pathways in *Arabidopsis thaliana*. All variables were logarithmically transformed. We configured the former 795 genes to response variables ( $p = 795$ ) with the latter 39 genes and an intercept as explanatory variables ( $k = 40$ ). The sample size is  $n = 118$ . We searched for the best subset of these models by using the selection method (23). Table 11 shows the explanatory variables selected by each criterion and the

Table 10. True subset selection probabilities (%) for each group  $b = 1, \dots, 8$  in the black cotton fibers dataset

$b$	Criterion								
	1	2	3	4	5	6	7	8	9
1	79.96	97.09	76.19	90.82	99.55	99.98	56.07	4.63	0.04
2	84.12	98.33	80.43	94.15	99.84	100.00	99.88	99.96	99.29
3	97.94	100.00	96.79	99.80	100.00	100.00	92.85	16.50	0.47
4	86.62	98.75	83.16	95.37	99.86	100.00	32.92	3.48	0.03
5	5.65	0.11	8.41	1.66	0.00	0.00	0.00	0.00	0.00
6	12.14	0.42	16.45	4.31	0.01	0.00	0.00	0.00	0.00
7	72.52	92.94	68.48	85.56	91.70	98.86	90.40	60.48	21.15
8	99.57	100.00	98.98	99.96	100.00	100.00	100.00	100.00	100.00

number of elements of the best subsets. From Table 11, we observe that criteria 7, 8, and 9 selected zero explanatory variables, and criteria 2 and 5 selected few variables. On the other hand, criteria 3 and 6 selected about half of the variables.

## 5. Conclusions and Discussions

We obtained the conditions for consistency of the  $SGC_p$  criterion and  $RGC_p$  criterion under the HUHD and UHD asymptotic frameworks. Importantly, consistency is established under non-normality and does not rely on the divergence speed of the dimension of the vector stacked with response variables  $p$ . Numerical studies suggest that criterion 6 has the highest probabilities of selecting the true subset, although consistency of criterion 6 holds when  $\log \log p / \log n \rightarrow 0$ .

Herein, the scalar matrix  $p^{-1}\text{tr}(\mathbf{S}_\omega)\mathbf{I}_p$  and the ridge-type sample covariance matrix  $\mathbf{S}_\lambda$  were used as  $\mathbf{G}$  in the weighted  $L_2$  squared distance  $d(\mathbf{A}, \mathbf{B}|\mathbf{G})$ . The  $SGC_p$  criterion and  $RGC_p$  criterion are invariant under transformation by a scalar times orthogonal matrices of  $\mathbf{Y}$ , i.e.,  $\mathbf{Y} : \mathbf{Y} \rightarrow a\mathbf{Y}\mathbf{F}$ , where  $\mathbf{F}$  satisfies  $\mathbf{F}\mathbf{F}' = \mathbf{F}'\mathbf{F} = \mathbf{I}_p$  and  $a \in \mathbb{R}$ . However, they are not invariant under transformation by nonsingular matrices of  $\mathbf{Y}$ , so their consistency is affected by the elements of  $\Sigma_*$  even for overspecified subsets. This is often the case in high-dimensional contexts such that  $p > n$ . On the other hand, using  $\text{diag}\{(\mathbf{S}_\omega)_{11}, \dots, (\mathbf{S}_\omega)_{pp}\}$  or  $\mathbf{S}_\omega + \lambda^{-1}\text{diag}\{(\mathbf{S}_\omega)_{11}, \dots, (\mathbf{S}_\omega)_{pp}\}$  as  $\mathbf{G}$  may eradicate the influence of the diagonal elements of  $\Sigma_*$ . Hence, it is also important to examine consistency in such cases. To do so would require assuming normality of the error vector and this represents fruitful terrain for future research.

Finally, we consider the influence of increasing  $p$  on consistency. To do so, another expression of multivariate linear regression is given by

$$\text{vec}(\mathbf{Y}) = (\mathbf{I}_p \otimes \mathbf{X})\text{vec}(\Theta) + \text{vec}(\mathcal{E}),$$

where  $\text{vec}(\mathbf{A})$  is the  $np$ -dimensional vector consisting of the columns of an  $n \times p$  matrix  $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n)$  and is defined by  $\text{vec}(\mathbf{A}) = (\mathbf{a}'_1, \dots, \mathbf{a}'_n)'$  (see, e.g., Harville, 1997). From the above expression, multivariate linear regression is regarded as univariate linear regression with the  $np$ -dimensional response vector  $\text{vec}(\mathbf{Y})$  and the explanatory matrix  $\mathbf{I}_p \otimes \mathbf{X}$  formally. From this, at first glance

Table 11. Selected explanatory variables based on the Arabidopsis thaliana dataset

Name	Criterion								
	1	2	3	4	5	6	7	8	9
Intercept	1	1	1	1	1	1	0	0	0
AACT1	1	0	1	1	0	1	0	0	0
AACT2	0	0	1	0	0	1	0	0	0
CMK	0	0	1	0	0	0	0	0	0
DPPS1	0	0	0	0	0	0	0	0	0
DPPS2	1	0	1	1	0	1	0	0	0
DPPS3	0	0	0	0	0	0	0	0	0
DXPS1	0	0	0	0	0	0	0	0	0
DXPS2(cla1)	1	0	1	1	0	1	0	0	0
DXPS3	0	0	1	0	0	0	0	0	0
DXR	1	0	1	1	0	1	0	0	0
FPPS1	0	0	0	0	0	0	0	0	0
FPPS2	0	0	0	0	0	0	0	0	0
GGPPS1mt	0	0	0	0	0	0	0	0	0
GGPPS2	0	0	0	0	0	0	0	0	0
GGPPS3	0	0	0	0	0	0	0	0	0
GGPPS4	0	0	0	0	0	0	0	0	0
GGPPS5	0	0	0	0	0	0	0	0	0
GGPPS6	1	0	1	1	0	1	0	0	0
GGPPS8	0	0	0	0	0	0	0	0	0
GGPPS9	0	0	0	0	0	0	0	0	0
GGPPS10	0	0	0	0	0	0	0	0	0
GGPPS11	0	0	1	0	0	0	0	0	0
GGPPS12	1	0	1	1	0	1	0	0	0
GPPS	1	0	1	1	0	1	0	0	0
HDR	1	0	1	1	0	1	0	0	0
HDS	1	0	1	1	0	1	0	0	0
HMGR1	1	0	1	1	0	1	0	0	0
HMGR2	0	0	1	0	0	1	0	0	0
HMGS	0	0	1	0	0	0	0	0	0
IPPI1	1	0	1	1	0	1	0	0	0
IPPI2	0	0	1	0	0	1	0	0	0
MCT	0	0	1	0	0	0	0	0	0
MECPS	0	0	1	0	0	1	0	0	0
MK	0	0	0	0	0	0	0	0	0
MPDC1	0	0	0	0	0	0	0	0	0
MPDC2	0	0	1	0	0	0	0	0	0
PPDS1	0	0	0	0	0	0	0	0	0
PPDS2mt	0	0	0	0	0	0	0	0	0
UPPS1	1	0	1	1	0	1	0	0	0
$\#(j)$	13	1	23	13	1	17	0	0	0

(1: selected variable, 0: non-selected variable. )

it seems that the dimension  $p$  has a role in increasing the sample size. However, from the results in Lemma 2 and Lemma 3, the probabilities of selecting  $j_*$  by the consistent criteria in this paper always approach 1 by diverging  $n$ , but do not always approach 1 by diverging only  $p$ . Moreover, increasing  $p$  leads to fast convergence of the probability of selecting the true subset under assumption A3', but this is not always the case under assumption A3. This difference depends on the assumption about  $\Sigma_*$  and  $\kappa_4$  since  $\xi \text{tr}(\Sigma_*)^{-1} = o(1)$  holds under assumption A3' not A3. This may also be verified from our simulations. Hence, to ensure fast convergence of the probability of selecting the true subset, a small sample size may be sufficient under assumption A3' when  $p$  is large. As per subsection 2.2, assumption A3' holds when (9) is supported. Since the sphericity condition  $\lim_{p \rightarrow \infty} \text{tr}(\Sigma_*^2)/\text{tr}(\Sigma_*)^2 = 0$  is equivalent to  $\lim_{p \rightarrow \infty} \lambda_{\max}(\Sigma_*)/\text{tr}(\Sigma_*) = 0$ , note that this condition implies that the maximum eigenvalue of  $\Sigma_*$  is not particularly large in the sense that  $\lambda_{\max}(\Sigma_*) = o(p)$  under assumption A2. However, in general  $\lambda_{\max}(\Sigma_*)$  tends to be very large for high-dimensional cases. Thus, it may not be suitable to assume the sphericity condition for high-dimensional cases. Aoshima and Yata (2018; 2019) considered methods to translate statistics under the strongly spiked model  $\liminf_{p \rightarrow \infty} \lambda_{\max}(\Sigma_*^2)/\text{tr}(\Sigma_*^2) > 0$  into those under the non-strongly spiked model  $\lim_{p \rightarrow \infty} \lambda_{\max}(\Sigma_*^2)/\text{tr}(\Sigma_*^2) = 0$ . By applying their idea to criteria for multivariate linear regression used in this paper, fast convergence of the probability of selecting the true subset can be ensured even under assumption A3, and, again, this should be explored in future research.

## Appendix

**A. Proof of equation (8).** Let  $j \in \mathcal{J}_-$ . From properties of projection matrices, for any  $\ell \in j_* \cap j^c$ , we have the following equation:

$$(\mathbf{I}_n - \mathbf{P}_{\omega_\ell})\mathbf{x}_{\ell_1} \begin{cases} = \mathbf{0}_n & (\ell_1 \in j_* \cap \{\ell\}^c) \\ \neq \mathbf{0}_n & (\ell_1 \in j_* \cap \{\ell\}) \end{cases}.$$

Using the above equation,  $\Theta_*' \mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_{\omega_\ell}) \mathbf{X}_* \Theta_*$  can be expressed as follows:

$$\begin{aligned} \Theta_*' \mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_{\omega_\ell}) \mathbf{X}_* \Theta_* &= \left( \sum_{\ell \in j_*} \boldsymbol{\theta}_\ell \mathbf{x}'_\ell \right) (\mathbf{I}_n - \mathbf{P}_{\omega_\ell}) \left( \sum_{\ell \in j_*} \mathbf{x}_\ell \boldsymbol{\theta}'_\ell \right) \\ &= \boldsymbol{\theta}_\ell \mathbf{x}'_\ell (\mathbf{I}_n - \mathbf{P}_{\omega_\ell}) \mathbf{x}_\ell \boldsymbol{\theta}'_\ell \\ &= \mathbf{x}'_\ell (\mathbf{I}_n - \mathbf{P}_{\omega_\ell}) \mathbf{x}_\ell \boldsymbol{\theta}_\ell \boldsymbol{\theta}'_\ell. \end{aligned}$$

Since we have

$$\mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j) \mathbf{X}_* - \mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_{\omega_\ell}) \mathbf{X}_* = \mathbf{X}'_*(\mathbf{P}_{\omega_\ell} - \mathbf{P}_j) \mathbf{X}_*,$$

and  $\mathbf{X}'_*(\mathbf{P}_{\omega_\ell} - \mathbf{P}_j) \mathbf{X}_*$  is positive-semidefinite, the following equation can be derived:

$$\lambda_{\max}(\Delta_j) \geq \lambda_{\max}(\Theta_*' \mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_{\omega_\ell}) \mathbf{X}_* \Theta_*) = \mathbf{x}'_\ell (\mathbf{I}_n - \mathbf{P}_{\omega_\ell}) \mathbf{x}_\ell \boldsymbol{\theta}_\ell \boldsymbol{\theta}'_\ell.$$

Hence, equation (8) can be derived from assumption A4.  $\square$

**B. Proof of Lemma 1.** We need a lemma to prove Lemma 1. To derive the upper bounds of probabilities, we use the variances of  $(n - k)^{-1}\text{tr}(\mathbf{W})$ ,  $\text{tr}(\mathbf{U}_j)$ , and  $\text{tr}(\mathbf{V}_{j,h})$ . The results for the variances are as follows (the proof is given in Appendix H):

LEMMA B.1. *Let  $\mathbf{A}$  be an  $n \times n$  symmetric matrix and  $\mathbf{B}$  be a  $p \times n$  matrix. Then, the following results hold:*

- (i)  $E[\text{tr}(\boldsymbol{\varepsilon}'_* \mathbf{A} \boldsymbol{\varepsilon}_*)] = \text{tr}(\mathbf{A})\text{tr}(\boldsymbol{\Sigma}_*)$ .
- (ii)  $E[\text{tr}(\mathbf{B} \boldsymbol{\varepsilon}_*^2)] = \text{tr}(\boldsymbol{\Sigma}_* \mathbf{B} \mathbf{B}')$ .
- (iii)  $E[\text{tr}(\boldsymbol{\varepsilon}'_* \mathbf{A} \boldsymbol{\varepsilon}_*^2)] = \left(\sum_{i=1}^n \{(\mathbf{A})_{ii}\}^2\right) \kappa_4 + \text{tr}(\mathbf{A})^2 \text{tr}(\boldsymbol{\Sigma}_*)^2 + 2\text{tr}(\mathbf{A}^2) \text{tr}(\boldsymbol{\Sigma}_*^2)$ , where  $\kappa_4 = E[|\boldsymbol{\varepsilon}|^4] - \text{tr}(\boldsymbol{\Sigma}_*)^2 - 2\text{tr}(\boldsymbol{\Sigma}_*^2)$ , which is defined in assumption A3.

Let  $j \supseteq h$ . Since  $\mathbf{I}_n - \mathbf{P}_\omega$  and  $\mathbf{P}_j - \mathbf{P}_h$  are symmetric idempotent matrices, we can identify that

$$\begin{aligned} \sum_{i=1}^n \{(\mathbf{I}_n - \mathbf{P}_\omega)_{ii}\}^2 &\leq \sum_{i=1}^n (\mathbf{I}_n - \mathbf{P}_\omega)_{ii} = \text{tr}(\mathbf{I}_n - \mathbf{P}_\omega) = n - k, \\ \sum_{i=1}^n \{(\mathbf{P}_j - \mathbf{P}_h)_{ii}\}^2 &\leq \sum_{i=1}^n (\mathbf{P}_j - \mathbf{P}_h)_{ii} = \text{tr}(\mathbf{P}_j - \mathbf{P}_h) = k_j - k_h. \end{aligned}$$

From the above equations and Lemma B.1, we can evaluate the expectations and variances of  $(n - k)^{-1}\text{tr}(\mathbf{W})$ ,  $\text{tr}(\mathbf{U}_j)$ , and  $\text{tr}(\mathbf{V}_{j,h})$  as follows:

$$\begin{aligned} E[(n - k)^{-1}\text{tr}(\mathbf{W})] &= \text{tr}(\boldsymbol{\Sigma}_*), \quad \text{Var}[(n - k)^{-1}\text{tr}(\mathbf{W})] \leq 3(n - k)^{-1}\xi^2, \\ E[\text{tr}(\mathbf{U}_j)^2] &= \text{tr}(\boldsymbol{\Sigma}_* \boldsymbol{\Delta}_j), \\ E[\text{tr}(\mathbf{V}_{j,h})] &= (k_j - k_h)\text{tr}(\boldsymbol{\Sigma}_*), \quad \text{Var}[\text{tr}(\mathbf{V}_{j,h})] \leq 3(k_j - k_h)\xi^2. \end{aligned}$$

Then, we obtain the results of Lemma 1 by using Chebyshev's inequality. First, we derive the results of (i), (ii), and (iii) as follows:

$$\begin{aligned} &P((n - k)^{-1}\text{tr}(\mathbf{W}) \geq r_1) \\ &= P((n - k)^{-1}\text{tr}(\mathbf{W}) - \text{tr}(\boldsymbol{\Sigma}_*) \geq r_1 - \text{tr}(\boldsymbol{\Sigma}_*)) \\ &\leq P(|(n - k)^{-1}\text{tr}(\mathbf{W}) - \text{tr}(\boldsymbol{\Sigma}_*)| \geq r_1 - \text{tr}(\boldsymbol{\Sigma}_*)) \\ &\leq \text{Var}[(n - k)^{-1}\text{tr}(\mathbf{W})] \{r_1 - \text{tr}(\boldsymbol{\Sigma}_*)\}^{-2} = O(\xi^2 n^{-1} \{r_1 - \text{tr}(\boldsymbol{\Sigma}_*)\}^{-2}), \\ &P((n - k)^{-1}\text{tr}(\mathbf{W}) \leq r_2) \\ &= P((n - k)^{-1}\text{tr}(\mathbf{W}) - \text{tr}(\boldsymbol{\Sigma}_*) \leq r_2 - \text{tr}(\boldsymbol{\Sigma}_*)) \\ &\leq P(|(n - k)^{-1}\text{tr}(\mathbf{W}) - \text{tr}(\boldsymbol{\Sigma}_*)| \geq \text{tr}(\boldsymbol{\Sigma}_*) - r_2) \\ &\leq \text{Var}[(n - k)^{-1}\text{tr}(\mathbf{W})] \{\text{tr}(\boldsymbol{\Sigma}_*) - r_2\}^{-2} = O(\xi^2 n^{-1} \{\text{tr}(\boldsymbol{\Sigma}_*) - r_2\}^{-2}), \\ &P(\text{tr}(\mathbf{U}_j) \leq r_3) \\ &\leq P(|\text{tr}(\mathbf{U}_j)| \geq |r_3|) \\ &\leq E[\text{tr}(\mathbf{U}_j)^2] |r_3|^{-2} = O(\text{tr}(\boldsymbol{\Sigma}_* \boldsymbol{\Delta}_j) |r_3|^{-2}), \end{aligned}$$

$$\begin{aligned}
& P(\operatorname{tr}(\mathbf{V}_{j,h}) \geq (k_j - k_h)r_4) \\
&= P(\operatorname{tr}(\mathbf{V}_{j,h}) - (k_j - k_h)\operatorname{tr}(\boldsymbol{\Sigma}_*) \geq (k_j - k_h)\{r_4 - \operatorname{tr}(\boldsymbol{\Sigma}_*)\}) \\
&\leq \operatorname{Var}[\operatorname{tr}(\mathbf{V}_{j,h})](k_j - k_h)^{-2}\{r_4 - \operatorname{tr}(\boldsymbol{\Sigma}_*)\}^{-2} = O(\xi^2\{r_4 - \operatorname{tr}(\boldsymbol{\Sigma}_*)\}^{-2}).
\end{aligned}$$

Next, we obtain result (iv). When  $n$  is sufficiently large or both  $n$  and  $p$  are sufficiently large, we have

$$-r_5 + r_6 < 0, \quad (r_5 - r_6)^{-1} = O(r_5^{-1}).$$

Hence, result (iii) can be derived as follows:

$$\begin{aligned}
& P\left(\operatorname{tr}(\mathbf{V}_{j,h}) - (k_j - k_j)\operatorname{tr}(\boldsymbol{\Sigma}_*) + r_5 \leq r_6\right) \\
&\leq P(|\operatorname{tr}(\mathbf{V}_{j,h}) - (k_j - k_h)\operatorname{tr}(\boldsymbol{\Sigma}_*)| \geq r_5 - r_6) \\
&\leq \operatorname{Var}[\operatorname{tr}(\mathbf{V}_{j,h})](r_5 - r_6)^{-2} = O(\xi^2 r_5^{-2}).
\end{aligned}$$

□

**C. Proof of Lemma 2.** First, we obtain the order of  $\bar{P}_S$ . For  $j \in \mathcal{J}_+ \cap \{j_*\}^c$ , let  $\mathbf{W} = \boldsymbol{\mathcal{E}}'_*(\mathbf{I}_n - \mathbf{P}_\omega)\boldsymbol{\mathcal{E}}_*$  and  $\mathbf{V}_{j,j_*} = \boldsymbol{\mathcal{E}}'_*(\mathbf{P}_j - \mathbf{P}_{j_*})\boldsymbol{\mathcal{E}}_*$  defined by (12). It is straightforward that the equation  $(\mathbf{I}_n - \mathbf{P}_\omega)\mathbf{X}_* = (\mathbf{P}_j - \mathbf{P}_{j_*})\mathbf{X}_* = \mathbf{O}_{n,k_*}$  holds. Then, we have

$$\operatorname{tr}\{\mathbf{Y}'(\mathbf{I}_n - \mathbf{P}_\omega)\mathbf{Y}\} = \operatorname{tr}(\mathbf{W}), \quad \operatorname{tr}\{\mathbf{Y}'(\mathbf{P}_j - \mathbf{P}_{j_*})\mathbf{Y}\} = \operatorname{tr}(\mathbf{V}_{j,j_*}).$$

Using the above equations,  $SGC_p(j|\alpha) - SGC_p(j_*|\alpha)$  is calculated as

$$\begin{aligned}
SGC_p(j|\alpha) - SGC_p(j_*|\alpha) &= -(n-k) \frac{\operatorname{tr}\{\mathbf{Y}'(\mathbf{P}_j - \mathbf{P}_{j_*})\mathbf{Y}\}}{\operatorname{tr}(\mathbf{W})} + (k_j - k_*)\alpha \\
&= -(n-k) \frac{\operatorname{tr}(\mathbf{V}_{j,j_*})}{\operatorname{tr}(\mathbf{W})} + (k_j - k_*)\alpha. \tag{C.1}
\end{aligned}$$

Let  $E_S$  be an event defined by

$$E_S = \{(n-k)^{-1}\operatorname{tr}(\mathbf{W}) \geq \tau_S \operatorname{tr}(\boldsymbol{\Sigma}_*)\}. \tag{C.2}$$

Then, by using (C.1) and (C.2), we have

$$\begin{aligned}
\bar{P}_S &= P(\cup_{j \in \mathcal{J}_+ \cap \{j_*\}^c} \{\operatorname{tr}(\mathbf{V}_{j,j_*}) \geq (n-k)^{-1}\operatorname{tr}(\mathbf{W})(k_j - k_*)\alpha\}) \\
&= P(\{\cup_{j \in \mathcal{J}_+ \cap \{j_*\}^c} \{\operatorname{tr}(\mathbf{V}_{j,j_*}) \geq (n-k)^{-1}\operatorname{tr}(\mathbf{W})(k_j - k_*)\alpha\}\} \cap (E_S \cup E_S^c)) \\
&\leq P(\cup_{j \in \mathcal{J}_+ \cap \{j_*\}^c} \{\operatorname{tr}(\mathbf{V}_{j,j_*}) \geq (k_j - k_*)\operatorname{tr}(\boldsymbol{\Sigma}_*)\alpha\tau_S\}) + P(E_S^c) \\
&\leq \sum_{j \in \mathcal{J}_+ \cap \{j_*\}^c} P(\operatorname{tr}(\mathbf{V}_{j,j_*}) \geq (k_j - k_*)\operatorname{tr}(\boldsymbol{\Sigma}_*)\alpha\tau_S) + P(E_S^c). \tag{C.3}
\end{aligned}$$

From (i) and (iii) of Lemma 1, the orders of two terms in (C.3) are as follows:

$$\begin{aligned}
& \sum_{j \in \mathcal{J}_+ \cap \{j_*\}^c} P(\operatorname{tr}(\mathbf{V}_{j,j_*}) \geq (k_j - k_*)\operatorname{tr}(\boldsymbol{\Sigma}_*)\alpha\tau_S) \\
&= O(\xi^2 \operatorname{tr}(\boldsymbol{\Sigma}_*)^{-2} (\alpha\tau_S - 1)^{-2}), \\
P(E_S^c) &= O(\xi^2 \operatorname{tr}(\boldsymbol{\Sigma}_*)^{-2} n^{-1} (1 - \tau_S)^{-2}).
\end{aligned}$$

From the above equations and (C.3), we have

$$\bar{P}_S = O\left(\xi^2 \text{tr}(\boldsymbol{\Sigma}_*)^{-2} \max\{(\alpha\tau_S - 1)^{-2}, n^{-1}(1 - \tau_S)^{-2}\}\right). \quad (\text{C.4})$$

Next, we obtain the order of  $\underline{P}_S$ . For  $j \in \mathcal{J}_-$ , let

$$j_+ = j \cup j_*, \quad E_{S,j} = \{SGC_p(j_+|\alpha) - SGC_p(j_*|\alpha) \geq 0\}.$$

Using  $j_+$  and  $E_{S,j}$ , we have

$$\begin{aligned} \underline{P}_S &= P\left(\bigcup_{j \in \mathcal{J}_-} \{SGC_p(j|\alpha) - SGC_p(j_+|\alpha) + SGC_p(j_+|\alpha) - SGC_p(j_*|\alpha) \leq 0\}\right) \\ &= P\left(\bigcup_{j \in \mathcal{J}_-} \{SGC_p(j|\alpha) - SGC_p(j_+|\alpha) + SGC_p(j_+|\alpha) - SGC_p(j_*|\alpha) \leq 0\}\right. \\ &\quad \left. \cap (E_{S,j} \cup E_{S,j}^c)\right) \\ &\leq P\left(\bigcup_{j \in \mathcal{J}_-} \{SGC_p(j|\alpha) - SGC_p(j_+|\alpha) \leq 0\}\right) + P(\bigcup_{j \in \mathcal{J}_-} E_{S,j}^c). \end{aligned} \quad (\text{C.5})$$

Since  $j_+ \in \mathcal{J}_+$  is the same as (C.4), the order of  $P(\bigcup_{j \in \mathcal{J}_-} E_{S,j}^c)$  is calculated as

$$P(\bigcup_{j \in \mathcal{J}_-} E_{S,j}^c) = O\left(\xi^2 \text{tr}(\boldsymbol{\Sigma}_*)^{-2} \max\{(\alpha\tau_S - 1)^{-2}, n^{-1}(1 - \tau_S)^{-2}\}\right). \quad (\text{C.6})$$

Notice that

$$\text{tr}\{\mathbf{Y}'(\mathbf{P}_{j_+} - \mathbf{P}_j)\mathbf{Y}\} = \text{tr}(\mathbf{V}_{j_+,j}) + 2\text{tr}(\mathbf{U}_j) + \delta_j^2,$$

where  $\delta_j^2$  and  $\mathbf{U}_j = \boldsymbol{\Theta}'_* \mathbf{X}'_*(\mathbf{I}_n - \mathbf{P}_j)\boldsymbol{\mathcal{E}}_*$  are defined by (7) and (12), respectively. From this,  $SGC_p(j|\alpha) - SGC_p(j_+|\alpha)$  is calculated as

$$\begin{aligned} &SGC_p(j|\alpha) - SGC_p(j_+|\alpha) \\ &= (n - k) \frac{\text{tr}\{\mathbf{Y}'(\mathbf{P}_{j_+} - \mathbf{P}_j)\mathbf{Y}\}}{\text{tr}(\mathbf{W})} - (k_{j_+} - k_j)\alpha \\ &= (n - k)\text{tr}(\mathbf{W})^{-1} \{ \text{tr}(\mathbf{V}_{j_+,j}) + 2\text{tr}(\mathbf{U}_j) + \delta_j^2 \} - (k_{j_+} - k_j)\alpha. \end{aligned} \quad (\text{C.7})$$

Let  $E_1$  and  $E_{2,j}$  be events defined by

$$E_1 = \left\{ (n - k)^{-1} \text{tr}(\mathbf{W}) \leq \frac{3}{2} \text{tr}(\boldsymbol{\Sigma}_*) \right\}, \quad E_{2,j} = \left\{ \text{tr}(\mathbf{U}_j) \geq -\frac{1}{4} \delta_j^2 \right\}. \quad (\text{C.8})$$

Then, by using (C.7) and (C.8), we have

$$\begin{aligned} &P\left(\bigcup_{j \in \mathcal{J}_-} \{SGC_p(j|\alpha) - SGC_p(j_+|\alpha) \leq 0\}\right) \\ &= P\left(\bigcup_{j \in \mathcal{J}_-} \{ \text{tr}(\mathbf{V}_{j_+,j}) + 2\text{tr}(\mathbf{U}_j) + \delta_j^2 \leq (n - k)^{-1} \text{tr}(\mathbf{W})(k_{j_+} - k_j)\alpha \}\right) \\ &= P\left(\bigcup_{j \in \mathcal{J}_-} \{ \text{tr}(\mathbf{V}_{j_+,j}) + 2\text{tr}(\mathbf{U}_j) + \delta_j^2 \leq (n - k)^{-1} \text{tr}(\mathbf{W})(k_{j_+} - k_j)\alpha \}\right. \\ &\quad \left. \cap (E_1 \cup E_1^c)\right) \\ &\leq P\left(\bigcup_{j \in \mathcal{J}_-} \left\{ \text{tr}(\mathbf{V}_{j_+,j}) + 2\text{tr}(\mathbf{U}_j) + \delta_j^2 \leq \frac{3}{2}(k_{j_+} - k_j)\text{tr}(\boldsymbol{\Sigma}_*)\alpha \right\}\right) \\ &\quad + P(E_1^c) \end{aligned}$$

$$\begin{aligned}
&= P \left( \bigcup_{j \in \mathcal{J}_-} \left\{ \text{tr}(\mathbf{V}_{j_+,j}) + 2\text{tr}(\mathbf{U}_j) + \delta_j^2 \leq \frac{3}{2}(k_{j_+} - k_j)\text{tr}(\boldsymbol{\Sigma}_*)\alpha \right\} \cap (E_{2,j} \cup E_{2,j}^c) \right) \\
&\quad + P(E_1^c) \\
&\leq \sum_{j \in \mathcal{J}_-} P \left( \text{tr}(\mathbf{V}_{j_+,j}) + \frac{1}{2}\delta_j^2 \leq \frac{3}{2}(k_{j_+} - k_j)\text{tr}(\boldsymbol{\Sigma}_*)\alpha \right) \\
&\quad + P(E_1^c) + \sum_{j \in \mathcal{J}_-} P(E_{2,j}^c). \tag{C.9}
\end{aligned}$$

Notice that

$$\frac{\text{tr}(\boldsymbol{\Sigma}_*)}{np} \left( \frac{3}{2}\alpha - 1 \right) \rightarrow 0, \quad \text{tr}(\boldsymbol{\Sigma}_* \boldsymbol{\Delta}_j) \leq \lambda_{\max}(\boldsymbol{\Sigma}_*)\delta_j^2.$$

Hence, by using (8) and (i), (ii), and (iii) of Lemma 1, the orders of three terms in (C.9) can be derived as follows:

$$\begin{aligned}
&\sum_{j \in \mathcal{J}_-} P \left( \text{tr}(\mathbf{V}_{j_+,j}) + \frac{1}{2}\delta_j^2 \leq \frac{3}{2}(k_{j_+} - k_j)\text{tr}(\boldsymbol{\Sigma}_*)\alpha \right) \\
&= \sum_{j \in \mathcal{J}_-} P \left( \text{tr}(\mathbf{V}_{j_+,j}) - (k_{j_+} - k_j)\text{tr}(\boldsymbol{\Sigma}_*) + \frac{1}{2}\delta_j^2 \leq (k_{j_+} - k_j)\text{tr}(\boldsymbol{\Sigma}_*) \left( \frac{3}{2}\alpha - 1 \right) \right) \\
&\leq \sum_{j \in \mathcal{J}_-} P \left( \frac{\text{tr}(\mathbf{V}_{j_+,j}) - (k_{j_+} - k_j)\text{tr}(\boldsymbol{\Sigma}_*)}{np} + \frac{1}{2}\tilde{\delta} \leq (k_{j_+} - k_j) \frac{\text{tr}(\boldsymbol{\Sigma}_*)}{np} \left( \frac{3}{2}\alpha - 1 \right) \right) \\
&= O(\xi^2 n^{-2} p^{-2}), \tag{C.10}
\end{aligned}$$

$$P(E_1^c) = O(\xi^2 \text{tr}(\boldsymbol{\Sigma}_*)^{-2} n^{-1}), \tag{C.11}$$

$$\sum_{j \in \mathcal{J}_-} P(E_{2,j}^c) = \sum_{j \in \mathcal{J}_-} O(\text{tr}(\boldsymbol{\Sigma}_* \boldsymbol{\Delta}_j) \delta_j^{-4}) = O(\lambda_{\max}(\boldsymbol{\Sigma}_*) n^{-1} p^{-1}), \tag{C.12}$$

where  $\tilde{\delta}$  is a positive constant satisfying  $0 < \tilde{\delta} < \min_{j \in \mathcal{J}_-} \inf_{n > k, p \geq 1} (np)^{-1} \delta_j^2$ . From (C.5), (C.6), (C.9), (C.10), (C.11), and (C.12), we have

$$\begin{aligned}
\underline{P}_S &= O \left( \xi^2 \text{tr}(\boldsymbol{\Sigma}_*)^{-2} \max \{ (\alpha \tau_S - 1)^{-2}, n^{-1} (1 - \tau_S)^{-2} \} \right) \\
&\quad + O \left( \max \{ \xi^2 n^{-2} p^{-2}, \xi^2 \text{tr}(\boldsymbol{\Sigma}_*)^{-2} n^{-1}, \lambda_{\max}(\boldsymbol{\Sigma}_*) n^{-1} p^{-1} \} \right). \tag{C.13}
\end{aligned}$$

Therefore, (C.4) and (C.13) complete the proof of Lemma 2.  $\square$

**D. Proof of Theorem 1.** First, we obtain the consistency conditions under assumptions A1, A2, A3, and A4. Note that under assumptions A2 and A3, the following equations hold:

$$\frac{\xi}{\text{tr}(\boldsymbol{\Sigma}_*)} = O(1), \quad \frac{\xi}{p} = O(1), \quad \frac{\lambda_{\max}(\boldsymbol{\Sigma}_*)}{p} = O(1).$$

Let us take  $\tau_S = 1/2$  in Lemma 2. By using Lemma 2 and the above equations, the orders of  $\overline{P}_S$  and  $\underline{P}_S$  are as follows:

$$\begin{aligned}\overline{P}_S &= O\left(\max\{(\alpha/2 - 1)^{-2}, n^{-1}\}\right), \\ \underline{P}_S &= O\left(\max\{(\alpha/2 - 1)^{-2}, n^{-1}\}\right) + O(n^{-1}).\end{aligned}$$

The above equations and (13) give the consistency conditions in (14).

Next, we obtain the consistency conditions under assumptions A1, A2, A3', and A4. Let us take  $\tau_S = 1 - n^{-1/2}$  in Lemma 2. Then, using (13), we have

$$\begin{aligned}(\alpha\tau_S - 1)^{-2} &= (\alpha - 1)^{-2} \left\{1 - \frac{\alpha}{\sqrt{n}(\alpha - 1)}\right\}^{-2} = O((\alpha - 1)^{-2}), \\ n^{-1}(1 - \tau_S)^{-2} &= 1.\end{aligned}$$

Note that under assumptions A2 and A3', the following equations hold:

$$\frac{\xi}{\text{tr}(\boldsymbol{\Sigma}_*)} = o(1), \quad \frac{\xi}{p} = o(1), \quad \frac{\lambda_{\max}(\boldsymbol{\Sigma}_*)}{p} = o(1).$$

Hence, the orders of  $\overline{P}_S$  and  $\underline{P}_S$  are as follows:

$$\overline{P}_S = o((\alpha - 1)^{-2}) + o(1), \quad \underline{P}_S = o((\alpha - 1)^{-2}) + o(1).$$

The above equations and (13) give the consistency conditions in (15).  $\square$

**E. Proof of Theorem 2.** First, we show the inconsistency under condition C1. Let  $\mathbf{W}$  and  $\mathbf{V}_{j,j_*}$  be defined by (12) and let  $E_3 = \{(n - k)^{-1}\text{tr}(\mathbf{W}) \leq (1 + n^{-1/4})\text{tr}(\boldsymbol{\Sigma}_*)\}$ . For any  $j \in \mathcal{J}_+ \cap \{j_*\}^c$ , we have

$$\begin{aligned}&P(\hat{j}_S = j_*) \\ &= P\left(\bigcap_{h \in \mathcal{J} \cap \{j_*\}^c} \{SGC_p(h|\alpha) > SGC_p(j_*|\alpha)\}\right) \\ &\leq P(SGC_p(j|\alpha) > SGC_p(j_*|\alpha)) \\ &= P\left(\text{tr}(\mathbf{V}_{j,j_*}) < \alpha(k_j - k_*)(n - k)^{-1}\text{tr}(\mathbf{W})\right) \\ &\leq P\left(\text{tr}(\mathbf{V}_{j,j_*}) - (k_j - k_*)\text{tr}(\boldsymbol{\Sigma}_*) < (k_j - k_*)\text{tr}(\boldsymbol{\Sigma}_*)\{(1 + n^{-1/4})\alpha - 1\}\right) \\ &\quad + P(E_3^c).\end{aligned}\tag{E.1}$$

Moreover, when  $n$  is sufficiently large or  $n$  and  $p$  are sufficiently large, we have

$$\begin{aligned}&P\left(\text{tr}(\mathbf{V}_{j,j_*}) - (k_j - k_*)\text{tr}(\boldsymbol{\Sigma}_*) < (k_j - k_*)\text{tr}(\boldsymbol{\Sigma}_*)\{(1 + n^{-1/4})\alpha - 1\}\right) \\ &\leq P\left(|\text{tr}(\mathbf{V}_{j,j_*}) - (k_j - k_*)\text{tr}(\boldsymbol{\Sigma}_*)| \geq (k_j - k_*)\text{tr}(\boldsymbol{\Sigma}_*)\{1 - (1 + n^{-1/4})\alpha\}\right) \\ &\leq \frac{\text{Var}[\text{tr}(\mathbf{V}_{j,j_*})]}{(k_j - k_*)^2\text{tr}(\boldsymbol{\Sigma}_*)^2\{1 - (1 + n^{-1/4})\alpha\}^2} \\ &\leq \frac{\kappa_4 I(\kappa_4 > 0) + 2\text{tr}(\boldsymbol{\Sigma}_*^2)}{(k_j - k_*)\text{tr}(\boldsymbol{\Sigma}_*)^2\{1 - (1 + n^{-1/4})\alpha\}^2}\end{aligned}$$

$$= (k_j - k_*)^{-1} (1 - \alpha)^{-2} \left( 1 - \frac{n^{-1/4} \alpha}{1 - \alpha} \right)^{-2} \left\{ \frac{\kappa_4 I(\kappa_4 > 0) + 2 \text{tr}(\boldsymbol{\Sigma}_*^2)}{\text{tr}(\boldsymbol{\Sigma}_*)^2} \right\}. \quad (\text{E.2})$$

Further, by using (i) in Lemma 1, the order of  $P(E_3^c)$  is as follows:

$$P(E_3^c) = O(\xi^2 \text{tr}(\boldsymbol{\Sigma}_*)^{-2} n^{-1/2}). \quad (\text{E.3})$$

From (E.1), (E.2), and (E.3), condition C1 gives the following inequality:

$$\begin{aligned} & \lim_{n \rightarrow \infty, p/n \rightarrow c} P(\hat{j}_S = j_*) \\ & \leq (k_j - k_*)^{-1} \left\{ \lim_{n \rightarrow \infty, p/n \rightarrow c} \frac{\kappa_4 I(\kappa_4 > 0) + 2 \text{tr}(\boldsymbol{\Sigma}_*^2)}{(1 - \alpha)^2 \text{tr}(\boldsymbol{\Sigma}_*)^2} \right\} < 1. \end{aligned}$$

Next, we show the inconsistency under condition C2. For  $j \not\subseteq j_*$ , let  $E_4 = \{(n - k)^{-1} \text{tr}(\mathbf{W}) \geq (1 - n^{-1/4}) \text{tr}(\boldsymbol{\Sigma}_*)\}$  and  $E_{5,j} = \{\text{tr}(\mathbf{U}_j) \leq n^{-1/4} \delta_j^2\}$ , where  $\mathbf{U}_j$  is defined by (12). Then, we have

$$\begin{aligned} & P(\hat{j}_S = j_*) \\ & \leq P(SGC_p(j|\alpha) > SGC_p(j_*|\alpha)) \\ & = P(\text{tr}(\mathbf{V}_{j_*,j}) + 2\text{tr}(\mathbf{U}_j) + \delta_j^2 > \alpha(k_* - k_j)(n - k)^{-1} \text{tr}(\mathbf{W})) \\ & \leq P\left(\text{tr}(\mathbf{V}_{j_*,j}) > (k_* - k_j) \text{tr}(\boldsymbol{\Sigma}_*) (1 - n^{-1/4}) \alpha - (1 + 2n^{-1/4}) \delta_j^2\right) \\ & \quad + P(E_4) + P(E_{5,j}^c). \end{aligned} \quad (\text{E.4})$$

From condition (C2), it is straightforward to identify that

$$\lim_{n \rightarrow \infty, p/n \rightarrow c} \frac{(k_* - k_j) \text{tr}(\boldsymbol{\Sigma}_*) \{(1 - n^{-1/4}) \alpha - 1\}}{(1 + 2n^{-1/4}) \delta_j^2} > 1.$$

Hence, when  $n$  is sufficiently large or  $n$  and  $p$  are sufficiently large, we have

$$\begin{aligned} & P\left(\text{tr}(\mathbf{V}_{j_*,j}) > (k_* - k_j) \text{tr}(\boldsymbol{\Sigma}_*) (1 - n^{-1/4}) \alpha - (1 + 2n^{-1/4}) \delta_j^2\right) \\ & \leq \frac{\text{Var}[\text{tr}(\mathbf{V}_{j_*,j})]}{[(k_* - k_j) \text{tr}(\boldsymbol{\Sigma}_*) \{(1 - n^{-1/4}) \alpha - 1\} - (1 + 2n^{-1/4}) \delta_j^2]^2} = O(n^{-2}). \end{aligned} \quad (\text{E.5})$$

Further, by using (i) and (ii) in Lemma 1, the orders of  $P(E_4^c)$  and  $P(E_{5,j}^c)$  are as follows:

$$P(E_4^c) = O(\xi^2 \text{tr}(\boldsymbol{\Sigma}_*)^{-2} n^{-1/2}), \quad P(E_{5,j}^c) = O(\lambda_{\max}(\boldsymbol{\Sigma}_*) p^{-1} n^{-1/2}). \quad (\text{E.6})$$

Equations (E.4), (E.5), and (E.6) give  $\lim_{n \rightarrow \infty, p/n \rightarrow c} P(\hat{j}_S = j_*) = 0$ .

Finally, when we replace assumption A3 with assumption A3', the results in this case can be derived from (E.1), (E.2), and (E.3) because of  $\xi \text{tr}(\boldsymbol{\Sigma}_*)^{-1} = o(1)$ .  $\square$

**F. Proof of Lemma 3.** For  $j \in \mathcal{J}_+ \cap \{j_*\}^c$ , using (19), we have

$$\begin{aligned}
 & RGC_p(j|\alpha, \lambda) - RGC_p(j_*|\alpha, \lambda) \\
 &= -\text{tr}\{\mathbf{Y}'(\mathbf{P}_j - \mathbf{P}_{j_*})\mathbf{Y}\mathbf{S}_\lambda^{-1}\} + (k_j - k_*)p\alpha \\
 &\geq -\text{tr}(\mathbf{V}_{j,j_*})\lambda_{\max}(\mathbf{S}_\lambda^{-1}) + (k_j - k_*)p\alpha \\
 &\geq -\lambda(n - k)\frac{\text{tr}(\mathbf{V}_{j,j_*})}{\text{tr}(\mathbf{W})} + (k_j - k_*)p\alpha \\
 &= \lambda\{SGC_p(j|\alpha) - SGC_p(j_*|\alpha)\} + (k_j - k_*)(p - \lambda)\alpha, \tag{F.1}
 \end{aligned}$$

where  $\mathbf{V}_{j,j_*}$  and  $\mathbf{W}$  are given by (12). Moreover, for  $j \in \mathcal{J}_-$ , using (19), we have

$$\begin{aligned}
 & RGC_p(j|\alpha, \lambda) - RGC_p(j_+|\alpha, \lambda) \\
 &= \text{tr}\{\mathbf{Y}'(\mathbf{P}_{j_+} - \mathbf{P}_j)\mathbf{Y}\mathbf{S}_\lambda^{-1}\} - (k_{j_+} - k_j)p\alpha \\
 &\geq \lambda_{\min}(\mathbf{S}_\lambda^{-1})\text{tr}\{\mathbf{Y}'(\mathbf{P}_{j_+} - \mathbf{P}_j)\mathbf{Y}\} - (k_{j_+} - k_j)p\alpha \\
 &\geq (1 + \lambda^{-1})^{-1}(n - k)\text{tr}(\mathbf{W})^{-1}\text{tr}\{\mathbf{Y}'(\mathbf{P}_{j_+} - \mathbf{P}_j)\mathbf{Y}\} - (k_{j_+} - k_j)p\alpha \\
 &= (1 + \lambda^{-1})^{-1}\{SGC_p(j|\alpha) - SGC_p(j_+|\alpha)\} + (k_{j_+} - k_j)\{(1 + \lambda^{-1})^{-1} - p\}\alpha, \tag{F.2}
 \end{aligned}$$

where  $j_+ = j \cup j_*$ . From (F.1) and (F.2), we can replace  $RGC_p(j|\alpha, \lambda) - RGC_p(j_*|\alpha, \lambda)$  and  $RGC_p(j|\alpha, \lambda) - RGC_p(j_+|\alpha, \lambda)$  with  $SGC_p(j|\alpha) - SGC_p(j_*|\alpha)$  and  $SGC_p(j|\alpha) - SGC_p(j_+|\alpha)$ , respectively. Therefore, in the same way as the proof of Lemma 2, the results of Lemma 3 can be derived.  $\square$

**G. Proof of Theorem 4.** For  $j \in \mathcal{J}_+ \cap \{j_*\}^c$ , using (19), we have

$$\begin{aligned}
 & RGC_p(j|\alpha, \lambda) - RGC_p(j_*|\alpha, \lambda) \\
 &\leq -\text{tr}(\mathbf{V}_{j,j_*})\lambda_{\min}(\mathbf{S}_\lambda^{-1}) + (k_j - k_*)p\alpha \\
 &\leq -(1 + \lambda^{-1})^{-1}(n - k)\text{tr}(\mathbf{W})^{-1}\text{tr}(\mathbf{V}_{j,j_*}) + (k_j - k_*)p\alpha \\
 &= (1 + \lambda^{-1})^{-1}\{SGC_p(j|\alpha) - SGC_p(j_*|\alpha)\} + (k_j - k_*)\{p - (1 + \lambda^{-1})^{-1}\}\alpha. \tag{G.1}
 \end{aligned}$$

For  $j \subsetneq j_*$ , using (19), we have

$$\begin{aligned}
 & RGC_p(j|\alpha, \lambda) - RGC_p(j_*|\alpha, \lambda) \\
 &\leq \lambda_{\max}(\mathbf{S}_\lambda^{-1})\text{tr}\{\mathbf{Y}'(\mathbf{P}_{j_*} - \mathbf{P}_j)\mathbf{Y}\} - (k_* - k_j)p\alpha \\
 &\leq \lambda(n - k)\text{tr}(\mathbf{W})^{-1}\text{tr}\{\mathbf{Y}'(\mathbf{P}_{j_*} - \mathbf{P}_j)\mathbf{Y}\} - (k_* - k_j)p\alpha \\
 &= \lambda\{SGC_p(j|\alpha) - SGC_p(j_*|\alpha)\} - (k_* - k_j)(\lambda - p)\alpha. \tag{G.2}
 \end{aligned}$$

By using (G.1) and (G.2), in the same way as the proof of Theorem 2, the results of Theorem 4 can be derived.  $\square$

**H. Proof of Lemma B.1.** First, we calculate the expectation  $E[\text{tr}(\mathcal{E}'_* \mathbf{A} \mathcal{E}_*)]$  to prove (i). It is straightforward that

$$E[\text{tr}(\mathcal{E}'_* \mathbf{A} \mathcal{E}_*)] = \sum_{i,j} (\mathbf{A})_{ij} E[\varepsilon'_i \varepsilon_j] = \sum_{i=1}^n (\mathbf{A})_{ii} E[\varepsilon'_i \varepsilon_i] = \text{tr}(\mathbf{A}) \text{tr}(\boldsymbol{\Sigma}_*),$$

where the summation  $\sum_{i,j}$  is defined by  $\sum_{i=1}^n \sum_{j=1}^n$ .

Next, we calculate the expectation  $E[\text{tr}(\mathbf{B} \mathcal{E}_*)^2]$  in (ii). Let  $\mathbf{b}_i$  be the  $i$ -th column vector of  $\mathbf{B}$ . Then, we have

$$E[\text{tr}(\mathbf{B} \mathcal{E}_*)^2] = \sum_{i,j} \mathbf{b}'_i E[\varepsilon_i \varepsilon'_j] \mathbf{b}_j = \sum_{i=1}^n \mathbf{b}'_i E[\varepsilon_i \varepsilon'_i] \mathbf{b}_i = \text{tr}(\boldsymbol{\Sigma}_* \mathbf{B} \mathbf{B}').$$

Finally, we calculate the expectation  $E[\text{tr}(\mathcal{E}'_* \mathbf{A} \mathcal{E}_*)^2]$  in (ii). The expectation  $E[\text{tr}(\mathcal{E}'_* \mathbf{A} \mathcal{E}_*)^2]$  can be expressed as follows:

$$\begin{aligned} E[\text{tr}(\mathcal{E}'_* \mathbf{A} \mathcal{E}_*)^2] &= \sum_{i,j,k,\ell} (\mathbf{A})_{ij} (\mathbf{A})_{kl} E[(\varepsilon'_i \varepsilon_j)(\varepsilon'_k \varepsilon_\ell)] \\ &= \sum_{i=1}^n \{(\mathbf{A})_{ii}\}^2 E[(\varepsilon'_i \varepsilon_i)^2] + \sum_{i \neq j} (\mathbf{A})_{ii} (\mathbf{A})_{jj} E[(\varepsilon'_i \varepsilon_i)(\varepsilon'_j \varepsilon_j)] \\ &\quad + 2 \sum_{i \neq j} \{(\mathbf{A})_{ij}\}^2 E[(\varepsilon'_i \varepsilon_j)^2] \\ &= \left( \sum_{i=1}^n \{(\mathbf{A})_{ii}\}^2 \right) E[\|\varepsilon\|^4] + \left( \sum_{i \neq j} (\mathbf{A})_{ii} (\mathbf{A})_{jj} \right) \text{tr}(\boldsymbol{\Sigma}_*)^2 \\ &\quad + 2 \left( \sum_{i \neq j} \{(\mathbf{A})_{ij}\}^2 \right) \text{tr}(\boldsymbol{\Sigma}_*^2), \end{aligned}$$

where the summation  $\sum_{i \neq j}$  is defined by  $\sum_{j=1}^n \sum_{i:i \neq j}$ . Hence, given that

$$\sum_{i \neq j} (\mathbf{A})_{ii} (\mathbf{A})_{jj} = \text{tr}(\mathbf{A})^2 - \sum_{i=1}^n \{(\mathbf{A})_{ii}\}^2, \quad \sum_{i \neq j} \{(\mathbf{A})_{ij}\}^2 = \text{tr}(\mathbf{A}^2) - \sum_{i=1}^n \{(\mathbf{A})_{ii}\}^2,$$

we can calculate  $E[\text{tr}(\mathcal{E}'_* \mathbf{A} \mathcal{E}_*)^2]$  as follows:

$$E[\text{tr}(\mathcal{E}'_* \mathbf{A} \mathcal{E}_*)^2] = \left( \sum_{i=1}^n \{(\mathbf{A})_{ii}\}^2 \right) \kappa_4 + \text{tr}(\mathbf{A})^2 \text{tr}(\boldsymbol{\Sigma}_*)^2 + 2 \text{tr}(\mathbf{A}^2) \text{tr}(\boldsymbol{\Sigma}_*^2).$$

□

### Acknowledgement

I wish to express my deepest gratitude to Prof. Hirokazu Yanagihara at Hiroshima University for his valuable advice and encouragement and introducing

me to various fields of mathematical statistics during the academic years 2014–2020. I also got a lot of advices about not only the personal manners as a researcher but also my private life from him, so I could not have come this far without his helps. In addition, I would like to thank Prof. Yasunori Fujikoshi at Hiroshima University for many helpful comments and suggestions about new research themes, Prof. Hirofumi Wakaki at Hiroshima University for his advice and help and Dr. Mariko Yamamura at Radiation Effects Research Foundation for her encouragement. Also, I thank to Dr. Shinpei Imori, Dr. Shintaro Hashimoto and Dr. Heewon Park at Hiroshima University for their encouragements, especially, Dr. Shinpei Imori for his valuable comments for numerical studies in this paper. Moreover, I am also grateful to Dr. Hiromi Itamiya at National Research Institute of Police Science, for providing me with the black cotton fiber dataset used in one of the empirical examples. I thank to my colleagues, seniors and juniors for their helps. I would also like to thank the referee for valuable comments.

## References

- [1] Aoshima, M., Shen, D., Shen, H., Yata, K., Zhou, Y. H. & Marron, J. (2018). A survey of high dimension low sample size asymptotics. *Aust. Nz. J. Stat.*, **60**, 4–19.
- [2] Aoshima, M. & Yata, K. (2018). Two-sample tests for high-dimension, strongly spiked eigenvalue models. *Stat. Sinica*, **28**, 43–62.
- [3] Aoshima, M. & Yata, K. (2019). Distance-based classifier by data transformation for high-dimension, strongly spiked eigenvalue models. *Ann. I. Stat. Math.*, **71**, 473–503.
- [4] 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. 995–1010. Akadémiai Kiadó, Budapest.
- [5] Akaike, H. (1974). A new look at the statistical model identification. *Institute of Electrical and Electronics Engineers. Transactions on Automatic Control* **AC – 19**, 716–723.
- [6] Bozdogan, H. (1987). Model selection and Akaike’s information criterion (AIC): the general theory and its analytical extensions. *Psychometrika*, **52**, 345–370.
- [7] Dempster, A. P. (1958). A high dimensional two sample significance test. *Ann. Math. Statist.*, **29**, 995–1010.
- [8] Dempster, A. P. (1960). A significance test for the separation of two highly multivariate small samples. *Biometrics*, **16**, 41–50.
- [9] Fujikoshi, Y., Kan, T., Takahashi, S. & Sakurai, T. (2011). Prediction error criterion for selecting variables in a linear regression model. *Ann. I. Stat. Math.*, **63**, 387–403.
- [10] Fujikoshi, Y., Himeno, T. & Wakaki, H. (2004). Asymptotic results of a high dimensional MANOVA test and power comparison when the dimension is large compared to the sample size. *J. Japan Statist. Soc.*, **34**, 19–26.
- [11] Fujikoshi, Y., Sakurai, T. & Yanagihara, H. (2014). Consistency of high-dimensional AIC-type and  $C_p$ -type criteria in multivariate linear regression. *J. Multivariate Anal.*, **123**, 184–200.
- [12] Fujikoshi, Y. & Satoh, K. (1997). Modified AIC and  $C_p$  in multivariate linear regression. *Biometrika*, **84**, 707–716.
- [13] Hannan, E. J. & Quinn, B. G. (1979). The determination of the order of an autoregression. *J. Roy. Statist. Soc. Ser. B*, **26**, 270–273.
- [14] Harville, D. A. (1997). *Matrix Algebra from a Statistician’s Perspective*. Springer-Verlag, New York.

- [15] Himeno, T. & Yamada, T. (2014). Estimations for some functions of covariance matrix in high dimension under non-normality and its applications. *J. Multivariate Anal.*, **130**, 27–44.
- [16] Katayama, S. & Imori, S. (2014). Lasso penalized model selection criteria for high-dimensional multivariate linear regression analysis. *J. Multivariate Anal.*, **132**, 138–150.
- [17] Kubokawa, T. & Srivastava, M. S. (2012). Selection of variables in multivariate regression models for large dimensions. *Comm. Statist. A-Theory Methods*, **41**, 2465–2489.
- [18] Magnus, J. R. & Neudecker, H. (1979). The commutation matrix: some properties and applications. *Ann. Statist.*, **7**, 381–894.
- [19] Mallows, C. L. (1973). Some comments on  $C_p$ . *Technometrics*, **15**, 661–675.
- [20] Mallows, C. L. (1995). More comments on  $C_p$ . *Technometrics*, **37**, 362–372.
- [21] Nagai, I., Yanagihara, H. & Satoh, K. (2012). Optimization of ridge parameters in multivariate generalized ridge regression by plug-in methods. *Hiroshima Math. J.*, **42**, 301–324.
- [22] Nishii, R., Bai, Z. D. & Krishnaiah, P. R. (1988). Strong consistency information criterion for model selection in multivariate analysis. *Hiroshima Math. J.*, **18**, 451–462.
- [23] Sparks, R. S., Coutsourides, D. & Troskie, L. (1983). The multivariate  $C_p$ . *Comm. Statist. A-Theory Methods*, **12**, 1775–1793.
- [24] Schwarz, G. (1978). Estimating the dimension of a model. *Ann. Statist.*, **6**, 461–464.
- [25] Srivastava, M. S. (2002). *Methods of Multivariate Statistics*. John Wiley & Sons, New York.
- [26] Timm, N. H. (2002). *Applied Multivariate Analysis*. Springer-Verlag, New York.
- [27] Wille, A., Zimmermann, P., Vranova, E., Fürholz, A., Laule, O., Bleuler, S., Hennig, L., Prelic, A., von Rohr, P., Thiele, L., Zitzler, E., Gruissem, W. & Bühlmann, P. (2004). Sparse graphical Gaussian modeling of the isoprenoid gene network in *Arabidopsis thaliana*. *Genome Biol.*, **5**, 1–13.
- [28] Yamamura, M., Yanagihara, H. & Srivastava, M. S. (2010). Variable selection in multivariate linear regression models with fewer observations than the dimension. *Japan. J. Appl. Stat.*, **39**, 1–19.
- [29] Yanagihara, H. (2015). Conditions for consistency of a log-likelihood-based information criterion in normal multivariate linear regression models under the violation of the normality assumption. *J. Japan Statist. Soc.*, **45**, 21–56.
- [30] 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.
- [31] Yanagihara, H. (2019). Evaluation of consistency of model selection criteria in multivariate linear regression models by large-sample and high-dimensional asymptotic theory under nonnormality. *J. Jpn. Stat. Soc. Jpn. Issue*, **48**, 1–13.
- [32] Yanagihara, H., Wakaki, H. & Fujikoshi, Y. (2015). A consistency property of the AIC for multivariate linear models when the dimension and the sample size are large. *Electron. J. Statist.*, **9**, 869–897.
- [33] 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.

*Ryoya Oda*

*Department of Mathematics*

*Graduate School of Science*

*Hiroshima University*

*Higashi-Hiroshima 739-8526, JAPAN*

*E-mail: ryoya-oda@hiroshima-u.ac.jp*