

A fast and consistent variable selection method for high-dimensional multivariate linear regression with a large number of explanatory variables

Ryoya Oda* and Hirokazu Yanagihara

*Department of Mathematics, Graduate School of Science, Hiroshima University,
1-3-1 Kagamiyama, Higashi-Hiroshima, Hiroshima 739-8526, Japan
e-mail: oda.stat@gmail.com; yanagi-hiro@hiroshima-u.ac.jp*

Abstract: We put forward a variable selection method for selecting explanatory variables in a normality-assumed multivariate linear regression. It is cumbersome to calculate variable selection criteria for all subsets of explanatory variables when the number of explanatory variables is large. Therefore, we propose a fast and consistent variable selection method based on a generalized C_p criterion. The consistency of the method is provided by a high-dimensional asymptotic framework such that the sample size and the sum of the dimensions of response vectors and explanatory vectors divided by the sample size tend to infinity and some positive constant which are less than one, respectively. Through numerical simulations, it is shown that the proposed method has a high probability of selecting the true subset of explanatory variables and is fast under a moderate sample size even when the number of dimensions is large.

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

Keywords and phrases: Consistency, high-dimensional asymptotic framework, multivariate linear regression, variable selection.

Received March 2019.

Contents

1	Introduction	1387
2	Preliminaries	1390
3	Main results	1392
	3.1 Proposed selection method	1392
	3.2 Consistency of proposed selection method	1394
	3.3 Extension of the ZKB selection method	1394
4	Numerical studies	1396
A	Appendix A	1399
	A.1 Proof of equations (2.4) and (2.5)	1399
	A.2 Proof of equation (3.4)	1401
	A.3 Proof of Lemma 3.1	1401

*Corresponding author.

A.4 Proof of Theorem 3.1 1402
 A.5 Proof of Lemma 3.2 1404
 A.6 Proof of Theorem 3.2 1404
 A.7 Proof of Lemma A.1 1406
 A.8 Proof of Lemma A.2 1408
 A.9 Proof of Lemma A.3 1409
 A.10 R file related to this article 1410
 Acknowledgments 1410
 References 1411

1. Introduction

Multivariate linear regression is a widely known method of inferential analysis. It features in many theoretical and applied textbooks (see, e.g., [21, chap 9], [24, chap 4]) and it is used by researchers in many fields. Let \mathbf{Y} be an $n \times p$ observation matrix of p response variables and \mathbf{X} be an $n \times k$ observation matrix of k non-stochastic explanatory variables, where n is the sample size, and p and k are the numbers of response variables and explanatory variables, respectively. Let $N = n - p - k + 1$, and we assume that $\text{rank}(\mathbf{X}) = k < n$ and (n, p, k) satisfies $N - 4 > 0$ in proposing our method.

In actual empirical contexts, it is important to specify the factors affecting response variables. In multivariate linear regression, this is regarded as the problem of selecting a subset of explanatory variables. Suppose that j denotes a subset of the full set $\omega = \{1, \dots, k\}$ containing k_j elements, and \mathbf{X}_j denotes the $n \times k_j$ matrix consisting of columns of \mathbf{X} indexed by the elements of j , where k_A denotes the number of elements in a set A , i.e., $k_A = \#(A)$. Next, j expresses the subset of explanatory variables. For example, if $j = \{1, 2, 4\}$, then \mathbf{X}_j consists of the first, second and fourth column vectors of \mathbf{X} . Using the notation j , the candidate model with k_j explanatory variables is expressed as follows:

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

where $\boldsymbol{\Theta}_j$ is a $k_j \times p$ unknown matrix of regression coefficients and $\boldsymbol{\Sigma}_j$ is a $p \times p$ unknown covariance matrix. In particular, the total number of explanatory variables k_ω and the explanatory matrix \mathbf{X}_ω in the full model ω express k and \mathbf{X} , respectively. Herein, we assume that the data are generated from the following true model with k_{j^*} explanatory variables:

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

where $\boldsymbol{\Theta}_*$ is a $k_{j^*} \times p$ true unknown matrix of regression coefficients and $\boldsymbol{\Sigma}_*$ is a $p \times p$ true unknown covariance matrix assuming that $\boldsymbol{\Sigma}_*$ is positive definite. 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 complement of set A . For expository purposes, we represent k_{j^*} and \mathbf{X}_{j^*} as k_* and \mathbf{X}_* , respectively.

To systematize and optimize the configuration of models, variable selection criteria have been widely used. The C_p criterion was proposed by [13, 14]. In this paper, we focus on a generalized variable selection criterion based on the C_p criterion, termed the Generalized C_p (GC_p) criterion. The GC_p criterion for a linear regression with a single response was proposed by [1], and the counterpart for a multivariate linear regression with multiple responses was proposed by [15]. The GC_p criterion can express a wide variety of variable selection criteria, e.g., the C_p criterion for multivariate contexts proposed by [20], and the modified C_p (MC_p) criterion proposed by [3].

The best subset chosen by a variable selection criterion is usually defined as the subset of explanatory variables which minimizes the value of that criterion among all candidate subsets. The basic approach to identifying the best subset involves searching over all candidate subsets. We call this method the “full search method”. To elaborate, assuming a full search method is used, variable selection criteria for $2^k - 1$ subsets need to be calculated. Recently, increasing attention has been paid to investigating statistical methods for high-dimensional data, in which the dimension of response vectors p or the number of explanatory variables k is large. However, in high-dimensional data contexts, particularly where k is large, it may be impossible to apply the full search method because the total number of subsets of explanatory variables exponentially increases when k becomes large. For example, if $k = 40$ and the time taken to calculate a variable selection criterion for a subset is 0.01 seconds, then the time required to implement the full search method will be $(2^{40} - 1) \times 0.01$ seconds, i.e., about 35 years. Thus, for practical reasons, we need another search method when k is large. A practicable selection method was proposed by [17, 31] when k is large. This method is based on the behavior of variable selection criteria for the subset where a variable is removed from the full set ω . In that selection method, the best subset \hat{j} is determined as follows. For each explanatory variable, if the criterion for the subset where a variable is removed from ω is greater than the criterion for the full set ω , then the removed variable is regarded as the element of the best subset. Since this method is needed to calculate variable selection criteria for only k subsets and ω for searching the best subset \hat{j} , we expect that the method is faster than the full search method, and it is practical for high-dimensional data contexts. We call this method the “ZKB selection method” and consider it using a class of the GC_p criterion, where “ZKB” is formed from the initial letters of the authors in [31].

An important property of a variable selection criterion is its consistency. Consistency is achieved where the probability of selecting the true subset j_* converges to 1, i.e., $P(\hat{j} = j_*) \rightarrow 1$. However, since we do not know the true subset j_* , we often hope to specify j_* by variable selection. Then, we should use a variable selection criterion that maximizes the probability of selecting the true subset. It is expected that a consistent variable selection criterion has a high-probability of selecting the true subset j_* because in general the probability of selecting the true subset is approximated by the asymptotic probability. To this end, let LS, LR, LE and LTE be the large-sample (LS), large-response vector (LR), large-explanatory vector (LE) and large-true explanatory vector

(LTE) asymptotic frameworks such that only n , p , k and k_* tend to infinity, respectively. Further, they are denoted by LS: $n \rightarrow \infty$, LR: $p \rightarrow \infty$, LE: $k \rightarrow \infty$ and LTE: $k_* \rightarrow \infty$. LS was used by [16, 17, 19, 31] under the ZKB selection method. However, it is not appropriate to use LS for high-dimensional data because approximate accuracy using LS deteriorates as p or k become large. Hence, criteria used by [16, 17, 19, 31] may not have consistency under the ZKB selection method when p or k tend to infinity. In the context for the consistency of variable selection criteria under the full search method, [4, 27, 28] used the following asymptotic frameworks as $(p + k)/n \rightarrow c \in [0, 1)$:

- [4]: LS and LR,
- [27]: LS or (LS and LR),
- [28]: (LS and LR) or (LS and LR and LE) or (LS and LR and LTE) as $k/n \rightarrow 0$.

Since as described above [4, 27, 28] used asymptotic frameworks such that not only n but also p , k or k_* tend to ∞ , the probabilities of selecting the true subset will be high for high-dimensional data suited to the used asymptotic frameworks. However, the probabilities may become low for high-dimensional data not suited to the used asymptotic frameworks. Moreover, it is hard for us to judge whether p , k and k_* are large or not, and so we do not know which asymptotic framework is suitable to given data. Hence, to ensure the consistency, it is more desirable to use an asymptotic framework regardless of sizes of p , k and k_* .

In this paper, we consider the consistency of the GC_p criterion under the ZKB selection method and propose the new consistent ZKB selection method even in high-dimensional contexts. Moreover, we also propose the selection method which can perform group selections. To achieve this, we use the following high-dimensional (HD) asymptotic framework:

$$\text{HD} : n \rightarrow \infty, \frac{p + k}{n} \rightarrow c \in [0, 1).$$

Importantly, the HD asymptotic framework can be rewritten as

$$\text{HD: LS or (LS and LR) or (LS and LE) or (LS and LTE) or (LS and LR and LE) or (LS and LR and LTE) as } (p + k)/n \rightarrow c \in [0, 1).$$

This means that n always tends to infinity, but p , k and k_* may tend to infinity as $(p + k)/n \rightarrow c \in [0, 1)$. Hence, it is expected that our proposed method will have a high probability of selecting the true subset where n is large regardless of the sizes of p , k and k_* . Moreover, even when k is large under $N - 4 > 0$, our proposed method will be very fast although the full search methods used in like [4, 27, 28] cannot be calculable. In recent years, regularization methods are often used for estimating the regression coefficients. The lasso is famous as one of methods estimating the regression coefficients and selecting explanatory variables simultaneously. In multivariate linear regression, it is possible to select explanatory variables by the group lasso proposed by [29], and several papers (e.g., [11, 18, 26, 30]) proposed regularization methods by extending the group lasso for multivariate linear regression case. Moreover, a generalized adaptive elastic-net was proposed and the consistency properties of the method were

obtained by [26]. The consistency properties were provided by using asymptotic frameworks such that $\log k / \log n \rightarrow \nu \in [0, 1)$ or $\log k = o(n^{1-2\kappa})$ for some $\kappa \in (0, 1/2)$. However, the properties are not ensured as p tends to infinity. Our method is consistent even when p tends to infinity as long as $N \rightarrow \infty$. Further, our method is faster than an adaptive group lasso even when p or k are large.

The remainder of the paper is organized as follows. In section 2, we present the necessary notation and assumptions to ensure consistency of our method. In section 3, we put forward the proposed method, explicate its consistency, and present a fast algorithm. We also propose an extended ZKB selection method. In section 4, we conduct numerical experiments for verification purposes. Technical details are relegated to the Appendix.

2. Preliminaries

First, we present the GC_p criterion. Let \mathbf{S}_j be the unbiased estimator of $\boldsymbol{\Sigma}_j$ in model (1.1), which is defined by

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

where \mathbf{P}_j is the projection matrix to the subspace spanned by the columns of \mathbf{X}_j , i.e., $\mathbf{P}_j = \mathbf{X}_j(\mathbf{X}_j'\mathbf{X}_j)^{-1}\mathbf{X}_j'$. Then, the GC_p criterion in model (1.1) is defined by

$$GC_p(j) = (n - k_j)\text{tr}(\mathbf{S}_j\mathbf{S}_\omega^{-1}) + \alpha p k_j, \quad (2.1)$$

where α is a positive constant. The first and second terms in (2.1) express the residual sum of squares with the weighted matrix \mathbf{S}_ω^{-1} and α times the strength of the penalty for the number of elements of $\boldsymbol{\Theta}_j$ in model (1.1), respectively.

Next, we present notation and assumptions to ensure consistency of our method. For a subset $j \subset \omega$, let a $p \times p$ non-centrality matrix and parameter be denoted by

$$\boldsymbol{\Delta}_j = \boldsymbol{\Sigma}_*^{-1/2} \boldsymbol{\Theta}_*' \mathbf{X}_*' (\mathbf{I}_n - \mathbf{P}_{\omega_j}) \mathbf{X}_* \boldsymbol{\Theta}_* \boldsymbol{\Sigma}_*^{-1/2}, \quad \delta_j = \text{tr}(\boldsymbol{\Delta}_j). \quad (2.2)$$

where $\omega_j = j^c$ and j^c denotes as $\omega \setminus j$. It should be emphasized that $\boldsymbol{\Delta}_j = \mathbf{O}_{p,p}$ and $\delta_j = 0$ hold if and only if $j \subset j_*^c$, where $\mathbf{O}_{p,p}$ is a $p \times p$ matrix of zeros. To ensure the consistency of our method, the following three assumptions are prepared:

Assumption A1. The true subset j_* is included in the full set ω , i.e., $j_* \subset \omega$.

Assumption A2. There exists $c_1 > 0$ such that

$$\min_{\ell \in j_*} n^{-1} \mathbf{x}'_{\{\ell\}} (\mathbf{I}_n - \mathbf{P}_{\omega_{\{\ell\}}}) \mathbf{x}_{\{\ell\}} \geq c_1,$$

where $\mathbf{x}_{\{\ell\}}$ is the ℓ -th column vectors of \mathbf{X} .

Assumption A3. There exist $1/2 < c_A \leq 1$ and $c_2 > 0$ such that

$$n^{1-c_A} \min_{\ell \in j_*} \boldsymbol{\theta}'_{\{\ell\}} \boldsymbol{\Sigma}_*^{-1} \boldsymbol{\theta}_{\{\ell\}} \geq c_2, \quad (2.3)$$

where $\boldsymbol{\theta}_{\{\ell\}}$ is the ℓ -th column vectors of $\boldsymbol{\Theta}'_*$.

Assumption A1 is needed to consider consistency because the probability of selecting the true subset becomes 0 if it does not hold. Assumption A2 means that the minimum value among the sample variances of residuals resulting from the linear regression of $\mathbf{x}_{\{\ell\}}$ with the remaining $\mathbf{X}_{\omega_{\{\ell\}}}$ for $\ell \in j_*$ is always positive and does not converge to 0. We often see an assumption for explanatory variables such that the inequality $n^{-1} \lambda_{\min}(\mathbf{X}'\mathbf{X}) \geq c_1$, where $\lambda_{\min}(\mathbf{A})$ is the minimum eigenvalue of a square matrix \mathbf{A} . Assumption A2 is weaker than this assumption because the inequality $\min_{\ell \in j_*} \mathbf{x}'_{\{\ell\}} (\mathbf{I}_n - \mathbf{P}_{\omega_{\{\ell\}}}) \mathbf{x}_{\{\ell\}} \geq \lambda_{\min}(\mathbf{X}'\mathbf{X})$ holds. Assumption A3 is a weak assumption for the true regression coefficients and the true covariance matrix. If $c_A < 1$, Assumption A3 allows $\min_{\ell \in j_*} \boldsymbol{\theta}'_{\{\ell\}} \boldsymbol{\Sigma}_*^{-1} \boldsymbol{\theta}_{\{\ell\}}$ to converge to 0. Moreover, for all $\ell = 1, \dots, k_*$, the following inequality holds (the proof is given in Appendix A.1):

$$\boldsymbol{\theta}'_{\{\ell\}} \boldsymbol{\Sigma}_*^{-1} \boldsymbol{\theta}_{\{\ell\}} \geq \max_{a=1, \dots, p} \frac{\theta_{* \ell a}^2}{\sigma_{*a}^2}, \quad (2.4)$$

where $\theta_{* \ell a}$ is the (ℓ, a) -th element of $\boldsymbol{\Theta}_*$ and σ_{*a}^2 is the a -th diagonal element of $\boldsymbol{\Sigma}_*$. From (2.4), Assumption A3 can be rewritten as the assumption which does not rely on the correlations of response variables by replacing $\boldsymbol{\theta}'_{\{\ell\}} \boldsymbol{\Sigma}_*^{-1} \boldsymbol{\theta}_{\{\ell\}}$ with $\max_{a=1, \dots, p} \theta_{* \ell a}^2 / \sigma_{*a}^2$ in (2.3). If Assumptions A1–A3 are supported, the following inequality holds (the proof is given in Appendix A.1):

$$n^{-c_A} \delta_{\min} \geq c_1 c_2, \quad (2.5)$$

where $\delta_{\min} = \min_{\ell \in j_*} \delta_{\{\ell\}}$. The above equation restricts the divergence order of the non-centrality parameter $\delta_{\{\ell\}}$. If k is fixed and $c_A = 1$, (2.5) is as per what was put forward in [27].

Finally, we identify the upper bound of the rank of the non-centrality matrix $\boldsymbol{\Delta}_j$, which is used to ensure consistency. For a subset $j \subset \omega$ ($j \neq \omega$), let m_j and d_j be the number of elements of j and the rank of $\boldsymbol{\Delta}_j$ as follows:

$$m_j = \#(j), \quad d_j = \text{rank}(\boldsymbol{\Delta}_j). \quad (2.6)$$

In accordance with [28], it follows from Assumption A1 that the rank of $\mathbf{X}'_*(\mathbf{P}_\omega - \mathbf{P}_{\omega_j})\mathbf{X}_*$ is calculated as

$$\text{rank}(\mathbf{X}'_*(\mathbf{P}_\omega - \mathbf{P}_{\omega_j})\mathbf{X}_*) = \begin{cases} 0 & (j \subset j_*^c) \\ m_j & (j \subset j_*) \end{cases}.$$

It is straightforward that $\text{rank}(\boldsymbol{\Theta}_* \boldsymbol{\Sigma}_*^{-1} \boldsymbol{\Theta}'_*) \leq \min\{p, k_*\}$. Since $m_j \leq k_*$ holds when $j \subset j_*$, the following inequality can be derived:

$$\begin{aligned} d_j &\leq \min\{\text{rank}(\mathbf{X}'_*(\mathbf{P}_\omega - \mathbf{P}_{\omega_j})\mathbf{X}_*), \text{rank}(\boldsymbol{\Theta}_* \boldsymbol{\Sigma}_*^{-1} \boldsymbol{\Theta}'_*)\} \\ &\leq \begin{cases} 0 & (j \subset j_*^c) \\ \min\{m_j, p\} & (j \subset j_*) \end{cases}. \end{aligned} \quad (2.7)$$

3. Main results

3.1. Proposed selection method

We define a class of the GC_p criterion, denoted as the high-dimensionality-adjusted consistent generalized C_p ($HCGC_p$) criterion:

Definition 3.1. *The $HCGC_p$ criterion is defined by the GC_p criterion (2.1) satisfying*

$$\alpha = \frac{n-k}{N-2} + \beta, \quad \beta > 0 \text{ s.t. } \frac{\sqrt{p}}{k^{1/2r}}\beta \rightarrow \infty, \quad \frac{p}{n^{c_A}}\beta \rightarrow 0, \quad (3.1)$$

as $n \rightarrow \infty$, $(p+k)/n \rightarrow c \in [0, 1)$, for some $r \in \mathbb{N}$, where c_A is defined in Assumption A3.

We now introduce the ZKB selection method using a variable selection criterion (SC). The best subset chosen by the ZKB selection method using an SC is written as

$$\{\ell \in \omega \mid \text{SC}(\omega_{\{\ell\}}) > \text{SC}(\omega)\},$$

where $\omega_{\{\ell\}}$ expresses $\{\ell\}^c$ or $\omega \setminus \{\ell\}$. The ZKB selection method is based on the idea that the value of the SC for the subset where a true variable is removed from ω will be greater than that for ω asymptotically. We define the following best subset chosen by the ZKB selection method using the $HCGC_p$ criterion:

Definition 3.2. *The best subset chosen by the ZKB selection method using the $HCGC_p$ criterion is defined by*

$$\hat{j} = \{\ell \in \omega \mid HCGC_p(\omega_{\{\ell\}}) > HCGC_p(\omega)\}. \quad (3.2)$$

Next, to use this method in actual empirical contexts we have to decide the value of α because the $HCGC_p$ criterion is expressed as the class of criteria. Hence, we show the following value of α :

$$\tilde{\alpha} = \frac{n-k}{N-2} + \tilde{\beta}, \quad \tilde{\beta} = \frac{(n-k)\sqrt{N+p-4}}{(N-2)\sqrt{N-4}} \cdot \frac{k^{1/4} \log n}{\sqrt{p}}. \quad (3.3)$$

This $\tilde{\alpha}$ is based on [27]. It is straightforward to observe that $\tilde{\beta}$ is satisfied with $(\sqrt{p}/k^{1/2r})\tilde{\beta} \rightarrow \infty$ and $(p/n^{c_A})\tilde{\beta} \rightarrow 0$ as $n \rightarrow \infty$, $(p+k)/n \rightarrow c \in [0, 1)$ for $r \geq 3$ and $3/4 < c_A \leq 1$. Therefore, the GC_p criterion with $\alpha = \tilde{\alpha}$ is included in the class of the $HCGC_p$ criterion with $3/4 < c_A \leq 1$. In practice, regardless of whether there is the constant value $\{(n-k)\sqrt{N+p-4}\}/\{(N-2)\sqrt{N-4}\}$ in $\tilde{\beta}$, the criterion belongs to the class of the $HCGC_p$ criterion. However, the constant value plays a role in terms of stabilizing the behavior of $p^{-1/2}\{HCGC_p(\omega_{\{\ell\}}) - HCGC_p(\omega)\}$ for $\ell \in j_*^c$.

Since the ZKB selection method using the GC_p criterion only necessitates calculating the differences $GC_p(\omega_{\{\ell\}}) - GC_p(\omega)$ for $\ell = 1, \dots, k$, it can be expected that the calculation time associated with this method will be shorter than that for the full search method. However, it is important that $GC_p(\omega_{\{\ell\}})$ consists of the projection matrix $\mathbf{P}_{\omega_{\{\ell\}}} = \mathbf{X}_{\omega_{\{\ell\}}} (\mathbf{X}'_{\omega_{\{\ell\}}} \mathbf{X}_{\omega_{\{\ell\}}})^{-1} \mathbf{X}'_{\omega_{\{\ell\}}}$ and the calculation time of an inverse matrix costs about the cube of the size of the matrix. Hence, it is not advisable to calculate $(\mathbf{X}'_{\omega_{\{\ell\}}} \mathbf{X}_{\omega_{\{\ell\}}})^{-1}$ for each ℓ when k is large. To overcome this problem, we offer an efficient calculation of $GC_p(\omega_{\{\ell\}}) - GC_p(\omega)$. Let r_ℓ and \mathbf{z}_ℓ be the (ℓ, ℓ) -th element of $(\mathbf{X}'\mathbf{X})^{-1}$ and the ℓ -th column vector of $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$, respectively. Then, using r_ℓ and \mathbf{z}_ℓ , we can express $\mathbf{P}_\omega - \mathbf{P}_{\omega_{\{\ell\}}}$ as follows (the proof of (3.4) is given in Appendix A.2):

$$\mathbf{P}_\omega - \mathbf{P}_{\omega_{\{\ell\}}} = \frac{1}{r_\ell} \mathbf{z}_\ell \mathbf{z}'_\ell. \quad (3.4)$$

Using the above equation, $GC_p(\omega_{\{\ell\}}) - GC_p(\omega)$ can be expressed as

$$GC_p(\omega_{\{\ell\}}) - GC_p(\omega) = \frac{1}{r_\ell} \mathbf{z}'_\ell \mathbf{Y} \mathbf{S}_\omega^{-1} \mathbf{Y}' \mathbf{z}_\ell - p\alpha. \quad (3.5)$$

Note that (3.5) does not need to calculate $(\mathbf{X}'_{\omega_{\{\ell\}}} \mathbf{X}_{\omega_{\{\ell\}}})^{-1}$ if only $(\mathbf{X}'\mathbf{X})^{-1}$ can be calculated. Moreover, the calculation cost of the product of each $\mathbf{Y}' \mathbf{z}_\ell$ relies on n . Hence, we also present an efficient calculation of $\mathbf{z}'_\ell \mathbf{Y} \mathbf{S}_\omega^{-1} \mathbf{Y}' \mathbf{z}_\ell$ when p is small. Let \mathbf{t}_ℓ be the ℓ -th column vector of $\mathbf{S}_\omega^{-1/2} \mathbf{Y}' \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1}$. Then, the following equation can be derived:

$$\mathbf{z}'_\ell \mathbf{Y} \mathbf{S}_\omega^{-1} \mathbf{Y}' \mathbf{z}_\ell = \mathbf{t}'_\ell \mathbf{t}_\ell. \quad (3.6)$$

Since \mathbf{t}_ℓ is a p -dimensional vector, the calculation cost of $\mathbf{t}'_\ell \mathbf{t}_\ell$ does not rely on n . Therefore, we propose to use (3.5) (and also use (3.6) when p is small) to perform the ZKB selection method using the GC_p criterion.

Note that the proposed selection method is calculable when $N - 4 > 0$. When $k > n$, we can formally combine the proposed selection method and screening methods by [8, 10, 12], which can apply to screening explanatory variables for a multivariate linear regression with multiple responses. However, we should pay attention to use their methods because the screening properties are ensured when p or k_* are fixed although the consistency of the proposed selection method is ensured even when p and k_* may diverge. On the other hand, a multivariate linear regression can be regarded as a perfunctory linear regression on a single response from the explanatory matrix $(\mathbf{I}_p \otimes \mathbf{X})$. However, notice that in generally we cannot directly apply several screening methods (e.g., [2]) for a linear regression with a single response to our variable selection problem. This is because our selection problem can be regarded as a group selection problem for explanatory variables corresponding to the p -dimensional regression coefficient vectors.

3.2. Consistency of proposed selection method

We ensure the consistency of the ZKB selection method using the $HCGC_p$ criterion (3.2). To do so, we present a lemma about sufficient conditions for consistency (the proof is given in Appendix A.3). Importantly, Lemma 3.1 does not rely on a specific asymptotic framework, indeed any such framework could be applied here.

Lemma 3.1. *Suppose that Assumption A1 and the following equations hold:*

$$\sum_{\ell \notin j_*} P(HCGC_p(\omega_{\{\ell\}}) > HCGC_p(\omega)) \rightarrow 0, \quad (3.7)$$

$$\sum_{\ell \in j_*} P(HCGC_p(\omega_{\{\ell\}}) < HCGC_p(\omega)) \rightarrow 0. \quad (3.8)$$

Then, the ZKB selection method using the $HCGC_p$ criterion (3.2) is consistent, that is $P(\hat{j} = j_*) \rightarrow 1$ holds.

By showing that the sufficient conditions (3.7) and (3.8) in Lemma 3.1 hold, the consistency of the ZKB selection method using the $HCGC_p$ criterion (3.2) can be obtained as follows (the proof is given in Appendix A.4):

Theorem 3.1. *Suppose that Assumptions A1–A3 hold. Then, the ZKB selection method using the $HCGC_p$ criterion (3.2) is consistent as $n \rightarrow \infty$, $(p+k)/n \rightarrow c \in [0, 1)$.*

From Theorem 3.1, the ZKB selection method using the $HCGC_p$ criterion with $\alpha = \tilde{\alpha}$ given by (3.3) is also consistent under Assumptions A1, A2 and Assumption A3 with $3/4 < c_A \leq 1$.

3.3. Extension of the ZKB selection method

In the previous subsections, we proposed the ZKB selection method using the $HCGC_p$ criterion (3.2). However, when the full model ω includes several explanatory variables such as multinomial variables, it will be not appropriate to use the ZKB selection method because whether such explanatory variables should be chosen or not should be decided simultaneously. To overcome this problem, we extend the ZKB selection method. Let \mathcal{J} be a family of sets of some explanatory variables denoted by $\mathcal{J} = \{j_1, \dots, j_q\}$, where q is the number of these sets. Since we suppose dummy variables or non-dummy variables as explanatory variables, we assume that m_{j_a} is finite, j_a is satisfied with $j_a \subset j_*$ or $j_a \subset j_*^c$ and $j_a \cap j_b = \emptyset$ ($a \neq b$) for $j_a, j_b \in \mathcal{J}$, where m_{j_a} is defined by (2.6). Then, it is clear that $\cup_{a=1}^q j_a = \omega$ holds. For example, if $k = 7$ and the sets of explanatory variables are $\{1\}$, $\{2\}$, $\{3, 5\}$ and $\{4, 6, 7\}$ then $\mathcal{J} = \{\{1\}, \{2\}, \{3, 5\}, \{4, 6, 7\}\}$, $q = 4$, and the subsets $\{3, 5\}$ and $\{4, 6, 7\}$ express the subsets of binomial and trinomial dummy variables, respectively. Using

this notation, we consider the following best subset chosen by the extended ZKB (EZKB) selection method using an SC:

$$\{j \in \mathcal{J} \mid \text{SC}(\omega_j) > \text{SC}(\omega)\}.$$

We observe that the EZKB selection method is equivalent to the ZKB selection method (3.2) when $m_j = 1 (\forall j \in \mathcal{J})$ or $q = k$. The EZKB selection method can accommodate the selection of grouped explanatory variables. We define the following best subset chosen by the EZKB selection method using the $HCGC_p$ criterion:

Definition 3.3. *The best subset chosen by the EZKB selection method using the $HCGC_p$ criterion is defined by*

$$\hat{j}_{\mathcal{J}} = \{j \in \mathcal{J} \mid HCGC_p(\omega_j) > HCGC_p(\omega)\}. \tag{3.9}$$

Next, we ensure the consistency of the EZKB selection method using the $HCGC_p$ criterion (3.9). Let $\mathcal{J}_+ = \{j \in \mathcal{J} \mid j \subset j_*\}$ and $\mathcal{J}_- = \{j \in \mathcal{J} \mid j \subset j_*^c\}$. Then, as with Lemma 3.1, we present the following lemma about sufficient conditions for consistency (the proof is given in Appendix A.5).

Lemma 3.2. *Suppose that Assumption A1 and the following equations hold:*

$$\begin{aligned} \sum_{j \in \mathcal{J}_+} P(HCGC_p(\omega_j) < HCGC_p(\omega)) &\rightarrow 0, \\ \sum_{j \in \mathcal{J}_-} P(HCGC_p(\omega_j) > HCGC_p(\omega)) &\rightarrow 0. \end{aligned}$$

Then, the EZKB selection method using the $HCGC_p$ criterion (3.9) is consistent.

Using Lemma 3.2, the consistency of the EZKB selection method using the $HCGC_p$ criterion (3.9) can be obtained as follows (the proof is given in Appendix A.6):

Theorem 3.2. *Suppose that Assumptions A1–A3 hold. Then, the EZKB selection method using the $HCGC_p$ criterion (3.9) is consistent as $n \rightarrow \infty$, $(p + k)/n \rightarrow c \in [0, 1)$.*

From Theorem 3.2, we can observe that the EZKB selection method using the $HCGC_p$ criterion is also consistent as with the ZKB selection method (3.2). Hence, as an example of the consistent EZKB selection method under $3/4 < c_A \leq 1$ in Assumption A3, we can use the method using the $HCGC_p$ criterion with $\alpha = \tilde{\alpha}$ in (3.3).

Finally, we provide an efficient calculation of $GC_p(\omega_j) - GC_p(\omega)$. Let \mathbf{R}_j and \mathbf{Z}_j be the $m_j \times m_j$ and $n \times m_j$ matrices consisting of the row and column elements of $(\mathbf{X}'\mathbf{X})^{-1}$ and the column vectors of $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$ indexed by the elements of j , respectively. For example, if $j = \{2, 5\}$, then \mathbf{R}_j and \mathbf{Z}_j are expressed as

$$\mathbf{R}_j = \begin{pmatrix} \tilde{x}_{22} & \tilde{x}_{25} \\ \tilde{x}_{52} & \tilde{x}_{55} \end{pmatrix}, \quad \mathbf{Z}_j = (\tilde{z}_2, \tilde{z}_5),$$

where \tilde{x}_{ab} is the (a, b) -element of $(\mathbf{X}'\mathbf{X})^{-1}$ and $\tilde{\mathbf{z}}_a$ is the a -th column vector of $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$. Then, using \mathbf{R}_j and \mathbf{Z}_j , $GC_p(\omega_j) - GC_p(\omega)$ can be expressed as

$$GC_p(\omega_j) - GC_p(\omega) = \text{tr}(\mathbf{R}_j^{-1} \mathbf{Z}_j' \mathbf{Y} \mathbf{S}_\omega^{-1} \mathbf{Y}' \mathbf{Z}_j) - m_j p \alpha. \quad (3.10)$$

The proof of the above equation is omitted because it essentially mimics (3.4). Although (3.10) requires the calculation of the inverse matrix of \mathbf{R}_j , it will not be computationally onerous because the size is finite.

4. Numerical studies

We present numerical results to explore the validity of our claim based on Monte Carlo simulations with 1,000 iterations executed in MATLAB 9.3.0 on a Panasonic CF-SV7U7FKVS with an Intel(R) Core(TM) i7-8650U CPU @ 1.90GHz 2.11 GHz and 16 GB of RAM. The probabilities of selecting the true subset and the CPU times are presented for the ZKB selection methods using the $HCGC_p$ criterion with $\alpha = \tilde{\alpha}$ given in (3.3) and the three GC_p criteria with $\alpha = 2, 2 \log \log n$ and $\log n$ (named GC_p^1, GC_p^2 and GC_p^3). The calculations were performed using (3.5) (and (3.6) if $p < 100$ and $k \geq p$). The explanatory matrix \mathbf{X} , the true coefficient matrix Θ_* and the true covariance matrix Σ_* were determined as follows:

$$\begin{aligned} \mathbf{X} &\sim N_{n \times k}(\mathbf{O}_{n,k}, \Psi \otimes \mathbf{I}_n), \quad \Theta_* \sim N_{k_* \times p}(\mathbf{O}_{k_*,p}, \mathbf{I}_p \otimes \mathbf{I}_{k_*}), \\ \Sigma_* &= \xi_1 \{(1 - \xi_2) \mathbf{I}_p + \xi_2 \mathbf{1}_p \mathbf{1}_p'\}, \end{aligned}$$

where Ψ is the $k \times k$ autoregressive matrix with the correlation ψ , i.e., $(\Psi)_{ab} = \psi^{|a-b|}$, and $\mathbf{1}_p$ is a p -dimensional vector of ones. Further, we set $\psi = 0.5$, $\xi_1 = 0.4$ and $\xi_2 = 0.8$. Although Theorems 3.1 and 3.2 were obtained by assuming that \mathbf{Y} is distributed according to the multivariate normal distribution under the true model, we also examine the probabilities under the non-normality in this simulation. Let $\mathcal{E} = (\varepsilon_1, \dots, \varepsilon_n)'$ be a $n \times p$ random matrix, where $\varepsilon_1, \dots, \varepsilon_n$ are independent and identically distributed according to the multivariate t -distribution with 10 degrees of freedom, mean $\mathbf{0}_p$ and covariance matrix $(5/4)\Sigma_*$. Then, we constructed the following two true models:

- Multivariate normal distribution: $\mathbf{Y} \sim N_{n \times p}(\mathbf{X}(\Theta_*', \mathbf{O}'_{k-k_*,p})', \Sigma_* \otimes \mathbf{I}_n)$.
- Multivariate t -distribution: $\mathbf{Y} = \mathbf{X}(\Theta_*', \mathbf{O}'_{k-k_*,p})' + (4/5)^{1/2} \mathcal{E}$.

For comparison, we also calculated the probabilities of selecting the true subset and the CPU times using the adaptive group lasso (AGL) proposed by [25]. The estimator of Θ by the AGL is written as

$$\begin{aligned} \hat{\Theta}_\tau &= \arg \min_{\Theta} f(\Theta|\tau), \\ f(\Theta|\tau) &= \text{tr}\{(\mathbf{Y} - \mathbf{X}\Theta)(\mathbf{Y} - \mathbf{X}\Theta)'\} + 2\tau \sum_{a=1}^k w_a \|\theta_a\|, \end{aligned} \quad (4.1)$$

where τ is a turning parameter, w_a is the weight for the norm $\|\boldsymbol{\theta}_a\| = (\boldsymbol{\theta}'_a \boldsymbol{\theta}_a)^{1/2}$, and $\boldsymbol{\theta}_a$ is the a -th column vector of $\boldsymbol{\Theta}'$. Each column vector of \mathbf{Y} and \mathbf{X} in (4.1) is centralized and standardized. To optimize (4.1), we used a coordinate descent algorithm based on [6]. The algorithm is given as follows. Let 100 candidates of τ be $\tau_t = \exp\{t \log(\tau_{\max} + 1)/(100 - 1)\} - 1$ ($t \in \{0, 1, 2, \dots, 99\}$), where $\tau_{\max} = \max_{a=1, \dots, k} w_a^{-1} \|\mathbf{Y}' \mathbf{x}_{\{a\}}\|$. Initialize $\hat{\boldsymbol{\Theta}}_{\tau_0} = \hat{\boldsymbol{\Theta}}_{\tau_0}^{\text{aft}} = (\hat{\boldsymbol{\theta}}_1^{(0)}, \dots, \hat{\boldsymbol{\theta}}_k^{(0)})' = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}$. For $t = 1, \dots, 99$,

1. Update $\hat{\boldsymbol{\Theta}}_{\tau_t}^{\text{bef}} \leftarrow \hat{\boldsymbol{\Theta}}_{\tau_{t-1}}^{\text{aft}}$ and $(\hat{\boldsymbol{\theta}}_1^{(t)}, \dots, \hat{\boldsymbol{\theta}}_k^{(t)})' \leftarrow \hat{\boldsymbol{\Theta}}_{\tau_{t-1}}^{\text{aft}}$. For each $a = 1, \dots, k$,
 - (1) Calculate $\mathbf{c}_a = \mathbf{Y}' \mathbf{x}_{\{a\}} - \sum_{i \neq a}^k \mathbf{x}'_{\{a\}} \mathbf{x}_{\{i\}} \hat{\boldsymbol{\theta}}_i^{(t)}$.
 - (2) If $\tau_t w_a \leq \|\mathbf{c}_a\|$, then update $\hat{\boldsymbol{\theta}}_a^{(t)} \leftarrow \{(\|\mathbf{c}_a\| - \tau_t w_a) / (\mathbf{x}'_{\{a\}} \mathbf{x}_{\{a\}} \|\mathbf{c}_a\|)\} \mathbf{c}_a$, otherwise $\hat{\boldsymbol{\theta}}_a^{(t)} \leftarrow \mathbf{0}_p$.
2. Update $\hat{\boldsymbol{\Theta}}_{\tau_t}^{\text{aft}} \leftarrow (\hat{\boldsymbol{\theta}}_1^{(t)}, \dots, \hat{\boldsymbol{\theta}}_k^{(t)})'$. If

$$\left| 1 - \frac{f(\hat{\boldsymbol{\Theta}}_{\tau_t}^{\text{aft}} | \tau_t)}{f(\hat{\boldsymbol{\Theta}}_{\tau_t}^{\text{bef}} | \tau_t)} \right| < \varepsilon_{\text{AGL}},$$

then define $\hat{\boldsymbol{\Theta}}_{\tau_t} = \hat{\boldsymbol{\Theta}}_{\tau_t}^{\text{aft}}$, otherwise go back to step 1.

In our setting, we used $\varepsilon_{\text{AGL}} = 0.01$, and w_a was given by $\|\hat{\boldsymbol{\theta}}_a^{\text{LSE}}\|^{-1}$, where $\hat{\boldsymbol{\theta}}_a^{\text{LSE}}$ is the least square estimator (LSE) of $\boldsymbol{\theta}_a$, i.e., $(\hat{\boldsymbol{\theta}}_1^{\text{LSE}}, \dots, \hat{\boldsymbol{\theta}}_k^{\text{LSE}})' = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}$. To choose the best turning parameter, we used three criteria as follows:

$$\hat{\tau}(\alpha_i) = \arg \min_{\tau_0, \dots, \tau_{99}} \text{IC}(\tau_t | \alpha_i),$$

$$\text{IC}(\tau_t | \alpha_i) = p^{-1} \text{tr}\{(\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\Theta}}_{\tau_t})' (\mathbf{Y} - \mathbf{X} \hat{\boldsymbol{\Theta}}_{\tau_t}) \mathbf{S}_{\omega}^{-1}\} + |\mathcal{A}_t| \alpha_i \quad (i = 1, 2, 3),$$

where $|\mathcal{A}_t|$ is the number of non-zero row vectors of $\hat{\boldsymbol{\Theta}}_{\tau_t}$, and $\alpha_1 = 2$, $\alpha_2 = 2 \log n$ and $\alpha_3 = \log n$. We name the AGL using $\text{IC}(\tau_t | \alpha_i)$ ($i = 1, 2, 3$) as AGL^1 , AGL^2 and AGL^3 , respectively. Table 1 shows the probabilities of selecting the true subset by the ZKB selection methods using the $HCGC_p$, GC_p^i ($i = 1, 2, 3$) denoted by $HCGC_p$, GC_p^i ($i = 1, 2, 3$) and AGL^i ($i = 1, 2, 3$) when \mathbf{Y} is distributed according to the multivariate normal distribution under the true model. From Table 1, we observe that the selection method using the $HCGC_p$ criterion always exhibits high probabilities of selecting the true subset for all combinations of n , p , k and k_* in Table 1. Although the probabilities by the method using the GC_p^3 criterion also achieve 100%, the performance by the method using the $HCGC_p$ criterion is better than those when the GC_p^3 criterion is used when the sample size is moderate. On the other hand, the probabilities by AGL^1 are low as the sample size increases in many cases. The probabilities by AGL^2 reach 100% only when the sample size is large and the dimensions are small. The probabilities by AGL^3 seem to increase slowly in some cases, but are low when k_* is large. Table 2 shows the probabilities when \mathbf{Y} is dis-

TABLE 1
 True subset selection probabilities (%) when \mathbf{Y} is distributed according to the multivariate normal distribution under the true model

n	p	k	k_*	$HCGC_p$	GC_p^1	GC_p^2	GC_p^3	AGL ¹	AGL ²	AGL ³
200	10	10	5	100.0	80.4	99.8	100.0	40.6	69.4	81.7
500	10	10	5	100.0	85.0	100.0	100.0	60.7	69.0	95.2
1000	10	10	5	100.0	85.5	100.0	100.0	74.4	76.9	97.7
2000	10	10	5	100.0	85.2	100.0	100.0	87.0	99.5	99.5
3000	10	10	5	100.0	83.9	100.0	100.0	11.0	100.0	100.0
200	160	10	5	99.9	0.0	0.0	0.3	0.0	0.0	0.3
500	400	10	5	100.0	0.0	0.0	37.4	0.0	0.0	30.5
1000	800	10	5	100.0	0.0	0.0	96.1	0.0	0.0	62.9
2000	1600	10	5	100.0	0.0	0.0	100.0	0.0	0.0	85.7
3000	2400	10	5	100.0	0.0	0.0	100.0	0.0	0.0	92.1
200	10	160	5	100.0	0.0	16.5	88.3	0.2	1.5	4.3
500	10	400	5	100.0	0.0	70.4	100.0	1.0	5.9	15.5
1000	10	800	5	100.0	0.0	90.5	100.0	15.7	26.0	36.6
2000	10	1600	5	100.0	0.0	93.7	100.0	44.0	84.1	86.7
3000	10	2400	5	100.0	0.0	95.4	100.0	24.1	24.3	35.1
200	10	160	80	100.0	0.2	34.0	93.5	0.0	0.0	0.1
500	10	400	200	100.0	0.1	83.1	99.9	0.0	0.0	3.1
1000	10	800	400	100.0	0.0	93.2	100.0	0.0	0.0	0.0
2000	10	1600	800	100.0	0.0	97.1	100.0	0.0	0.0	0.0
3000	10	2400	1200	100.0	0.0	97.5	100.0	0.0	0.0	0.0
200	80	80	5	100.0	0.0	0.0	31.5	0.0	0.0	2.5
500	200	200	5	100.0	0.0	0.0	99.4	0.0	5.3	19.4
1000	400	400	5	100.0	0.0	0.3	100.0	0.0	20.7	40.8
2000	800	800	5	100.0	0.0	78.2	100.0	0.0	41.3	66.1
3000	1200	1200	5	100.0	0.0	99.6	100.0	0.0	51.4	78.9
200	80	80	40	100.0	0.0	0.0	52.7	0.0	0.0	0.2
500	200	200	100	100.0	0.0	0.0	99.7	0.0	0.0	0.0
1000	400	400	200	100.0	0.0	3.7	100.0	0.0	0.0	1.3
2000	800	800	400	100.0	0.0	89.9	100.0	0.0	0.0	68.4
3000	1200	1200	600	100.0	0.0	100.0	100.0	0.0	0.0	96.0

tributed according to the multivariate t -distribution under the true model. We observe that the results in Table 2 are about the same as those in Table 1. Hence, it is expected that our results may hold even under the non-normality. The proofs of Theorems 3.1 and 3.2 in this paper are needed to calculate the moments of $GC_p(\omega_{\{\ell\}}) - GC_p(\omega)$ and $GC_p(\omega_j) - GC_p(\omega)$, and we calculated them by assuming that \mathbf{Y} is distributed according to the multivariate normal distribution under the true model. However, our results will be shown even under non-normality if another approach to the evaluation of the moments exists although we need to calculate the moments consisting of the inverse matrix \mathbf{S}_ω^{-1} . Table 3 shows the CPU times by the ZKB selection method using the $HCGC_p$ criterion denoted by $HCGC_p$ and AGL^3 when \mathbf{Y} is distributed according to the multivariate normal distribution under the true model, and the former is faster than the latter. The difference is particularly clear when the dimensions are large. In sum, the ZKB selection method using the $HCGC_p$ criterion with $\alpha = \tilde{\alpha}$ exhibits the highest probabilities of selecting the true subset and is faster than the AGLs.

TABLE 2
 True subset selection probabilities (%) when \mathbf{Y} is distributed according to the multivariate t -distribution under the true model

n	p	k	k_*	$HCGC_p$	GC_p^1	GC_p^2	GC_p^3	AGL ¹	AGL ²	AGL ³
200	10	10	5	100.0	80.6	99.5	100.0	51.0	67.8	81.8
500	10	10	5	100.0	85.2	99.9	100.0	76.5	92.2	96.9
1000	10	10	5	100.0	86.9	100.0	100.0	78.9	95.6	97.9
2000	10	10	5	100.0	87.4	100.0	100.0	83.3	99.5	99.5
3000	10	10	5	100.0	84.0	99.9	100.0	0.0	99.8	100.0
200	160	10	5	99.8	0.0	0.0	0.2	0.0	0.0	0.2
500	400	10	5	100.0	0.0	0.0	33.8	0.0	0.0	23.8
1000	800	10	5	100.0	0.0	0.0	96.2	0.0	0.0	65.7
2000	1600	10	5	100.0	0.0	0.0	100.0	0.0	0.0	81.9
3000	2400	10	5	100.0	0.0	0.0	100.0	0.0	0.0	89.8
200	10	160	5	100.0	0.0	19.7	87.9	0.2	0.8	3.9
500	10	400	5	100.0	0.0	72.7	100.0	3.0	9.5	20.2
1000	10	800	5	100.0	0.0	88.4	100.0	4.8	22.2	49.5
2000	10	1600	5	100.0	0.0	94.9	100.0	12.4	50.2	66.9
3000	10	2400	5	100.0	0.0	96.0	100.0	33.3	36.0	65.4
200	10	160	80	100.0	0.3	33.1	93.1	0.0	0.0	0.0
500	10	400	200	100.0	0.0	84.7	100.0	0.0	0.0	0.0
1000	10	800	400	100.0	0.0	93.0	100.0	0.0	0.0	0.0
2000	10	1600	800	100.0	0.0	96.9	100.0	0.0	0.0	0.0
3000	10	2400	1200	100.0	0.0	97.8	100.0	0.0	0.0	0.0
200	80	80	5	100.0	0.0	0.0	34.4	0.0	0.2	3.4
500	200	200	5	100.0	0.0	0.0	99.6	0.0	4.8	21.9
1000	400	400	5	100.0	0.0	0.1	100.0	0.0	14.7	37.3
2000	800	800	5	100.0	0.0	79.0	100.0	0.0	43.6	66.7
3000	1200	1200	5	100.0	0.0	99.8	100.0	0.0	48.5	76.2
200	80	80	40	100.0	0.0	0.0	52.9	0.0	0.0	0.3
500	200	200	100	100.0	0.0	0.0	99.8	0.0	0.0	0.0
1000	400	400	200	100.0	0.0	4.4	100.0	0.0	0.0	2.2
2000	800	800	400	100.0	0.0	90.5	100.0	0.0	0.0	68.1
3000	1200	1200	600	100.0	0.0	100.0	100.0	0.0	0.0	94.0

Appendix A

A.1. Proof of equations (2.4) and (2.5)

First, we show (2.5). For an arbitrary $\ell \in j_*$, we have the following equation:

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

From the above equation, $\Delta_{\{\ell\}}$ which is defined in (2.2) can be expressed as follows:

$$\begin{aligned} \Delta_{\{\ell\}} &= \Sigma_*^{-1/2} \left(\sum_{\ell \in j_*} \boldsymbol{\theta}_{\{\ell\}} \mathbf{x}'_{\{\ell\}} \right) (\mathbf{I}_n - \mathbf{P}_{\omega_{\{\ell\}}}) \left(\sum_{\ell \in j_*} \mathbf{x}_{\{\ell\}} \boldsymbol{\theta}'_{\{\ell\}} \right) \Sigma_*^{-1/2} \\ &= \Sigma_*^{-1/2} \boldsymbol{\theta}_{\{\ell\}} \mathbf{x}'_{\{\ell\}} (\mathbf{I}_n - \mathbf{P}_{\omega_{\{\ell\}}}) \mathbf{x}_{\{\ell\}} \boldsymbol{\theta}'_{\{\ell\}} \Sigma_*^{-1/2} \\ &= \mathbf{x}'_{\{\ell\}} (\mathbf{I}_n - \mathbf{P}_{\omega_{\{\ell\}}}) \mathbf{x}_{\{\ell\}} \Sigma_*^{-1/2} \boldsymbol{\theta}_{\{\ell\}} \boldsymbol{\theta}'_{\{\ell\}} \Sigma_*^{-1/2} . \end{aligned}$$

TABLE 3
CPU times (s) when \mathbf{Y} is distributed according to the multivariate normal distribution under the true model

n	p	k	k_*	$HCGC_p$	AGL^3
200	10	10	5	0.0010	0.0091
500	10	10	5	0.0046	0.0108
1000	10	10	5	0.0124	0.0208
2000	10	10	5	0.0346	0.0428
3000	10	10	5	0.0606	0.0679
200	160	10	5	0.0036	0.0611
500	400	10	5	0.0438	0.8003
1000	800	10	5	0.2996	4.9315
2000	1600	10	5	2.1194	33.0587
3000	2400	10	5	6.6364	99.3778
200	10	160	5	0.0029	0.1048
500	10	400	5	0.0167	0.4377
1000	10	800	5	0.0734	1.3389
2000	10	1600	5	0.3844	4.7584
3000	10	2400	5	1.1095	11.4628
200	10	160	80	0.0024	0.0934
500	10	400	200	0.0146	0.4503
1000	10	800	400	0.0774	1.3550
2000	10	1600	800	0.4126	4.8661
3000	10	2400	1200	1.1370	11.8354
200	80	80	5	0.0045	0.1386
500	200	200	5	0.0327	1.5399
1000	400	400	5	0.5254	57.9087
2000	800	800	5	4.9127	492.7444
3000	1200	1200	5	15.5871	1724.8028
200	80	80	40	0.0059	0.1409
500	200	200	100	0.0312	1.5835
1000	400	400	200	0.5342	59.2730
2000	800	800	400	4.8972	516.9926
3000	1200	1200	600	15.7256	1730.3650

Hence, we have

$$\delta_{\{\ell\}} = \text{tr}(\Delta_{\{\ell\}}) = \mathbf{x}'_{\{\ell\}}(\mathbf{I}_n - \mathbf{P}_{\omega_{\{\ell\}}})\mathbf{x}_{\{\ell\}}\boldsymbol{\theta}'_{\{\ell\}}\boldsymbol{\Sigma}_*^{-1}\boldsymbol{\theta}_{\{\ell\}}.$$

The above equation leads to the following inequality:

$$\begin{aligned} n^{-c_A}\delta_{\min} &\geq \left\{ n^{-1} \min_{\ell \in J_*} \mathbf{x}'_{\{\ell\}}(\mathbf{I}_n - \mathbf{P}_{\omega_{\{\ell\}}})\mathbf{x}_{\{\ell\}} \right\} \left\{ n^{1-c_A} \min_{\ell \in J_*} \boldsymbol{\theta}'_{\{\ell\}}\boldsymbol{\Sigma}_*^{-1}\boldsymbol{\theta}_{\{\ell\}} \right\} \\ &\geq c_1 c_2. \end{aligned}$$

Next, we show (2.4). Let $\lambda_{\max}(\mathbf{A})$ be the maximum eigenvalue of the square matrix \mathbf{A} . Then, we have

$$\boldsymbol{\theta}'_{\{\ell\}}\boldsymbol{\Sigma}_*^{-1}\boldsymbol{\theta}_{\{\ell\}} = \lambda_{\max}(\boldsymbol{\Sigma}_*^{-1/2}\boldsymbol{\theta}_{\{\ell\}}\boldsymbol{\theta}'_{\{\ell\}}\boldsymbol{\Sigma}_*^{-1/2}) = \max_{\|e\|=1} e'\boldsymbol{\Sigma}_*^{-1/2}\boldsymbol{\theta}_{\{\ell\}}\boldsymbol{\theta}'_{\{\ell\}}\boldsymbol{\Sigma}_*^{-1/2}e.$$

By using $(e'_a\boldsymbol{\Sigma}_*e_a)^{-1/2}\boldsymbol{\Sigma}_*^{1/2}e_a$ as e , we have

$$\lambda_{\max}(\boldsymbol{\Sigma}_*^{-1/2}\boldsymbol{\theta}_{\{\ell\}}\boldsymbol{\theta}'_{\{\ell\}}\boldsymbol{\Sigma}_*^{-1/2}) \geq \max_{a=1,\dots,p} \frac{(e'_a\boldsymbol{\theta}_{\{\ell\}})^2}{\sigma_{*a}^2},$$

where e_a is the p -dimensional vector such that the a -th element is one and the other elements are zero. The above equation completes the proof of (2.4). \square

A.2. Proof of equation (3.4)

Without loss of generality, let $\mathbf{X} = (\mathbf{X}_{\omega_{\{\ell\}}}, \mathbf{X}_{\{\ell\}})$ for an $\ell \in \omega$. Further, let \mathbf{R}_ℓ , \mathbf{r}_ℓ and r_ℓ be satisfied with

$$\begin{pmatrix} \mathbf{R}_\ell & \mathbf{r}_\ell \\ \mathbf{r}'_\ell & r_\ell \end{pmatrix} = (\mathbf{X}'\mathbf{X})^{-1}.$$

Then, using the general formula for the inverse of a block matrix (e.g., [7, Theorem 8.5.11]), \mathbf{P}_ω and $\mathbf{P}_{\omega_{\{\ell\}}}$ can be expressed as follows:

$$\begin{aligned} \mathbf{P}_\omega &= \mathbf{X}_{\omega_{\{\ell\}}} \mathbf{R}_\ell \mathbf{X}'_{\omega_{\{\ell\}}} + \mathbf{X}_{\omega_{\{\ell\}}} \mathbf{r}_\ell \mathbf{X}'_{\{\ell\}} + \mathbf{X}_{\{\ell\}} \mathbf{r}'_\ell \mathbf{X}'_{\omega_{\{\ell\}}} + r_\ell \mathbf{X}_{\{\ell\}} \mathbf{X}'_{\{\ell\}}, \\ \mathbf{P}_{\omega_{\{\ell\}}} &= \mathbf{X}_{\omega_{\{\ell\}}} \mathbf{R}_\ell \mathbf{X}'_{\omega_{\{\ell\}}} + r_\ell^{-1} \mathbf{X}_{\omega_{\{\ell\}}} \mathbf{r}_\ell \mathbf{r}'_\ell \mathbf{X}'_{\omega_{\{\ell\}}}. \end{aligned}$$

From the above equations, we have

$$\mathbf{P}_\omega - \mathbf{P}_{\omega_{\{\ell\}}} = \frac{1}{r_\ell} \mathbf{X} \begin{pmatrix} \mathbf{r}_\ell \\ r_\ell \end{pmatrix} \begin{pmatrix} \mathbf{r}_\ell \\ r_\ell \end{pmatrix}' \mathbf{X}'.$$

Note that $\mathbf{X}(\mathbf{r}'_\ell, r_\ell)'$ is the k -th column vector of $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$. Therefore, (3.4) can be derived. \square

A.3. Proof of Lemma 3.1

We can express $P(\hat{j} = j_*)$ as follows:

$$\begin{aligned} P(\hat{j} = j_*) &= P\left(\left(\bigcap_{\ell \in j_*} \{HCGC_p(\omega_{\{\ell\}}) - HCGC_p(\omega) > 0\}\right) \right. \\ &\quad \left. \bigcap \left(\bigcap_{\ell \notin j_*} \{HCGC_p(\omega_{\{\ell\}}) - HCGC_p(\omega) \leq 0\}\right)\right). \end{aligned}$$

Then, the following lower bound of $P(\hat{j} = j_*)$ can be derived:

$$\begin{aligned} P(\hat{j} = j_*) &\geq 1 - \sum_{\ell \in j_*} P(HCGC_p(\omega_{\{\ell\}}) - HCGC_p(\omega) < 0) \\ &\quad - \sum_{\ell \notin j_*} P(HCGC_p(\omega_{\{\ell\}}) - HCGC_p(\omega) > 0). \end{aligned}$$

This completes the proof of Lemma 3.1. \square

A.4. Proof of Theorem 3.1

We first describe two lemmas. The first lemma gives another expression of $GC_p(\omega_j) - GC_p(\omega)$ for $j \subset \omega$ ($j \neq \omega$) (the proof is given in Appendix A.7):

Lemma A.1. *For $j \subset \omega$ ($j \neq \omega$), suppose that $\delta_{j,i}$ ($i = 1, \dots, m_j$) are constants satisfying $\text{tr}(\mathbf{\Delta}_j) = \sum_{i=1}^{m_j} \delta_{j,i}$ and $\delta_{j,i} \geq m_j^{-1} \lambda_{\max}(\mathbf{\Delta}_j)$, where $\mathbf{\Delta}_j$ and m_j are defined by (2.2) and (2.6), and $\lambda_{\max}(\mathbf{\Delta}_j)$ is the maximum eigenvalue of $\mathbf{\Delta}_j$. Let u_i , $u_{j,i}$ and v_i be random variables distributed according to $u_i \sim \chi^2(p)$, $u_{j,i} \sim \chi^2(p; \delta_{j,i})$ and $v_i \sim \chi^2(n - p - k + 1)$ ($i = 1, \dots, m_j$), where u_i and $u_{j,i}$ are independent of v_i for each i . Then, under Assumption A1, we have*

$$GC_p(\omega_j) - GC_p(\omega) = \begin{cases} (n-k) \sum_{i=1}^{m_j} \frac{u_i}{v_i} - m_j p \alpha & (j \subset j_*^c) \\ (n-k) \sum_{i=1}^{m_j} \frac{u_{j,i}}{v_i} - m_j p \alpha & (j \subset j_*) \end{cases}. \quad (\text{A.1})$$

The following lemma is needed to evaluate the divergence orders of the moments of $GC_p(\omega_j) - GC_p(\omega)$ (the proof is given in Appendix A.8).

Lemma A.2. *Let δ be a positive constant. And let u_1 , u_2 and v be random variables distributed according to $\chi^2(p)$, $\chi^2(p; \delta)$ and $\chi^2(N)$, where u_1 and u_2 are independent of v , and $N = n - p - k + 1$. Then, for $N - 4r > 0$ ($r \in \mathbb{N}$), we have*

$$E \left[\left(\frac{u_1}{v} - \frac{p}{N-2} \right)^{2r} \right] = O(p^r n^{-2r}),$$

$$E \left[\left(\frac{u_2}{v} - \frac{p+\delta}{N-2} \right)^{2r} \right] = O(\max\{(p+\delta)^r n^{-2r}, (p+\delta)^{2r} n^{-3r}\}),$$

as $n - p - k \rightarrow \infty$.

Applying the results of Lemma A.1 for $m_j = 1$ to $HCGC_p(\omega_{\{\ell\}}) - HCGC_p(\omega)$, we have

$$HCGC_p(\omega_{\{\ell\}}) - HCGC_p(\omega) = \begin{cases} (n-k) \frac{u}{v} - p \alpha & (\ell \notin j_*) \\ (n-k) \frac{u_\ell}{v} - p \alpha & (\ell \in j_*) \end{cases}, \quad (\text{A.2})$$

where u and u_ℓ are independent of v , and $u \sim \chi^2(p)$, $u_\ell \sim \chi^2(p; \delta_{\{\ell\}})$ and $v \sim \chi^2(N)$. From (A.2), we have

$$\begin{aligned} & \sum_{\ell \notin j_*} P(HCGC_p(\omega_{\{\ell\}}) > HCGC_p(\omega)) \\ &= (k - k_*) P \left(\frac{u}{v} > \frac{p}{n-k} \alpha \right) = (k - k_*) P \left(\frac{u}{v} - \frac{p}{N-2} > \rho \right) \end{aligned}$$

$$\leq (k - k_*)P\left(\left|\frac{u}{v} - \frac{p}{N-2}\right| \geq \rho\right), \tag{A.3}$$

where $\rho = \{p/(n - k)\}\beta$. Moreover, since $(N - 2)\rho\delta_{\min}^{-1} = O(p\beta n^{-c_A}) = o(1)$ from (2.5) and (3.1), when N is sufficiently large, we have

$$\begin{aligned} & \sum_{\ell \in j_*} P(HCGC_p(\omega_{\{\ell\}}) < HCGC_p(\omega)) \\ &= \sum_{\ell \in j_*} P\left(\frac{u_\ell}{v} < \frac{p}{n - k}\alpha\right) = \sum_{\ell \in j_*} P\left(\frac{u_\ell}{v} - \frac{p + \delta_{\{\ell\}}}{N - 2} < \rho + -\frac{\delta_{\{\ell\}}}{N - 2}\right) \\ &\leq \sum_{\ell \in j_*} P\left(\left|\frac{u_\ell}{v} - \frac{p + \delta_{\{\ell\}}}{N - 2}\right| \geq \frac{\delta_{\{\ell\}}}{N - 2} - \rho\right). \end{aligned} \tag{A.4}$$

Applying Markov’s inequality to (A.3) and (A.4), the following upper bounds can be derived:

$$\begin{aligned} (k - k_*)P\left(\left|\frac{u}{v} - \frac{p}{N - 2}\right| \geq \rho\right) &\leq (k - k_*)\rho^{-2r} E\left[\left(\frac{u}{v} - \frac{p}{N - 2}\right)^{2r}\right], \\ \sum_{\ell \in j_*} P\left(\left|\frac{u_\ell}{v} - \frac{p + \delta_{\{\ell\}}}{N - 2}\right| \geq \frac{\delta_{\{\ell\}}}{N - 2} - \rho\right) \\ &\leq \sum_{\ell \in j_*} \left(\frac{\delta_{\{\ell\}}}{N - 2} - \rho\right)^{-2\tilde{r}} E\left[\left(\frac{u_\ell}{v} - \frac{p + \delta_{\{\ell\}}}{N - 2}\right)^{2\tilde{r}}\right], \end{aligned}$$

where r is a natural number defined by (3.1) and \tilde{r} are any natural numbers. From the above equations and Lemma A.2, the following equation can be derived:

$$\sum_{\ell \notin j_*} P(HCGC_p(\omega_{\{\ell\}}) > HCGC_p(\omega)) = O(kp^{-r}\beta^{-2r}) = o(1).$$

Moreover, for sufficiently large \tilde{r} , since $k_*p^{\tilde{r}}n^{-2\tilde{r}c_A} = o(1)$, $k_*n^{-\tilde{r}c_A} = o(1)$, $k_*p^{2\tilde{r}}n^{-2\tilde{r}c_A - \tilde{r}} = o(1)$ and $k_*n^{-\tilde{r}} = o(1)$, we have

$$\begin{aligned} & \sum_{\ell \in j_*} P(HCGC_p(\omega_{\{\ell\}}) < HCGC_p(\omega)) \\ &= \sum_{\ell \in j_*} O\left((p + \delta_{\{\ell\}})^{\tilde{r}}\delta_{\{\ell\}}^{-2\tilde{r}} + (p + \delta_{\{\ell\}})^{2\tilde{r}}\delta_{\{\ell\}}^{-2\tilde{r}}n^{-\tilde{r}}\right) \\ &= O\left(k_*(p + \delta_{\min})^{\tilde{r}}\delta_{\min}^{-2\tilde{r}} + k_*(p + \delta_{\min})^{2\tilde{r}}\delta_{\min}^{-2\tilde{r}}n^{-\tilde{r}}\right) = o(1). \end{aligned}$$

These equations and Lemma 3.1 complete the proof of Theorem 3.1. □

A.5. Proof of Lemma 3.2

We can express $P(\hat{j}_{\mathcal{J}} = j_*)$ as follows:

$$P(\hat{j}_{\mathcal{J}} = j_*) = P\left(\left(\bigcap_{j \in \mathcal{J}_+} \{HCGC_p(\omega_j) - HCGC_p(\omega) > 0\}\right) \cap \left(\bigcap_{j \in \mathcal{J}_-} \{HCGC_p(\omega_j) - HCGC_p(\omega) \leq 0\}\right)\right).$$

Then, the following lower bound of $P(\hat{j}_{\mathcal{J}} = j_*)$ can be derived:

$$P(\hat{j}_{\mathcal{J}} = j_*) \geq 1 - \sum_{j \in \mathcal{J}_+} P(HCGC_p(\omega_j) - HCGC_p(\omega) < 0) - \sum_{j \in \mathcal{J}_-} P(HCGC_p(\omega_j) - HCGC_p(\omega) > 0).$$

Therefore, Lemma 3.2 can be derived. \square

A.6. Proof of Theorem 3.2

We can apply the results of Lemma A.1 to this proof, i.e., we can express the following distribution forms of $HCGC_p(\omega_j) - HCGC_p(\omega)$:

$$HCGC_p(\omega_j) - HCGC_p(\omega) = \begin{cases} (n-k) \sum_{i=1}^{m_j} \frac{u_i}{v_i} - m_j p \alpha & (j \in \mathcal{J}_-) \\ (n-k) \sum_{i=1}^{m_j} \frac{u_{j,i}}{v_i} - m_j p \alpha & (j \in \mathcal{J}_+) \end{cases}, \quad (\text{A.5})$$

where u_i and $u_{j,i}$ are independent of v_i , and

$$u_i \sim \chi^2(p), \quad u_{j,i} \sim \chi^2(p; \delta_{j,i}), \quad v_i \sim \chi^2(N) \quad (i = 1, \dots, m_j).$$

Here, $\delta_{j,i}$ ($i = 1, \dots, m_j$) are constants satisfying $\sum_{i=1}^{m_j} \delta_{j,i} = \text{tr}(\mathbf{\Delta}_j)$ and $\delta_{j,i} \geq m_j^{-1} \lambda_{\max}(\mathbf{\Delta}_j)$, where $\mathbf{\Delta}_j$ is given by (2.2). When $j \in \mathcal{J}_+$, let ℓ be an element of j , i.e., $\ell \in j$. Then, since $\mathbf{I}_n - \mathbf{P}_{\omega_{\{\ell\}}}$ and $\mathbf{P}_{\omega_{\{\ell\}}} - \mathbf{P}_{\omega_j}$ are semi-positive definite, the following equation can be derived:

$$\text{tr}(\mathbf{\Delta}_j) = \delta_{\{\ell\}} + \text{tr}\{\mathbf{\Sigma}_*^{-1/2} \mathbf{\Theta}'_* \mathbf{X}'_* (\mathbf{P}_{\omega_{\{\ell\}}} - \mathbf{P}_{\omega_j}) \mathbf{X}_* \mathbf{\Theta}_* \mathbf{\Sigma}_*^{-1/2}\} \geq \delta_{\{\ell\}}.$$

In addition, let $d_j = \text{rank}(\mathbf{\Delta}_j)$ which is defined by (2.6). From (2.7), we observe that d_j is bounded. Since $d_j \lambda_{\max}(\mathbf{\Delta}_j) \geq \text{tr}(\mathbf{\Delta}_j)$ holds, the following inequality is obtained:

$$\delta_{j,i} \geq m_j^{-1} \lambda_{\max}(\mathbf{\Delta}_j) \geq (m_j d_j)^{-1} \text{tr}(\mathbf{\Delta}_j) \geq (m_j d_j)^{-1} \delta_{\{\ell\}}. \quad (\text{A.6})$$

Now, we derive the divergence orders of $\sum_{j \in \mathcal{J}_-} P(HCGC_p(\omega_j) > HCGC_p(\omega))$ and $\sum_{j \in \mathcal{J}_+} P(HCGC_p(\omega_j) < HCGC_p(\omega))$. From (2.5), (3.1), (A.5) and (A.6), when N is sufficiently large, we have

$$\begin{aligned} & \sum_{j \in \mathcal{J}_-} P(HCGC_p(\omega_j) > HCGC_p(\omega)) \\ &= \sum_{j \in \mathcal{J}_-} P\left(\sum_{i=1}^{m_j} \frac{u_i}{v_i} > \frac{m_j p}{n-k} \alpha\right) \leq \sum_{j \in \mathcal{J}_-} \sum_{i=1}^{m_j} P\left(\frac{u_i}{v_i} > \frac{p}{n-k} \alpha\right) \\ &= \sum_{j \in \mathcal{J}_-} \sum_{i=1}^{m_j} P\left(\frac{u_i}{v_i} - \frac{p}{N-2} > \rho\right) \leq \sum_{j \in \mathcal{J}_-} \sum_{i=1}^{m_j} P\left(\left|\frac{u_i}{v_i} - \frac{p}{N-2}\right| \geq \rho\right), \end{aligned} \tag{A.7}$$

$$\begin{aligned} & \sum_{j \in \mathcal{J}_+} P(HCGC_p(\omega_j) < HCGC_p(\omega)) \\ &= \sum_{j \in \mathcal{J}_+} P\left(\sum_{i=1}^{m_j} \frac{u_{j,i}}{v_i} < \frac{m_j p}{n-k} \alpha\right) \leq \sum_{j \in \mathcal{J}_+} \sum_{i=1}^{m_j} P\left(\frac{u_{j,i}}{v_i} < \frac{p}{n-k} \alpha\right) \\ &= \sum_{j \in \mathcal{J}_+} \sum_{i=1}^{m_j} P\left(\frac{u_{j,i}}{v_i} - \frac{p + \delta_{j,i}}{N-2} < \rho - \frac{\delta_{j,i}}{N-2}\right) \\ &\leq \sum_{j \in \mathcal{J}_+} \sum_{i=1}^{m_j} P\left(\left|\frac{u_{j,i}}{v_i} - \frac{p + \delta_{j,i}}{N-2}\right| \geq \frac{\delta_{j,i}}{N-2} - \rho\right), \end{aligned} \tag{A.8}$$

where $\rho = \{p/(n-k)\}\beta$. Then, by applying Markov's inequality to (A.7) and (A.8), their following upper bounds can be derived:

$$\begin{aligned} & \sum_{j \in \mathcal{J}_-} \sum_{i=1}^{m_j} P\left(\left|\frac{u_i}{v_i} - \frac{p}{N-2}\right| \geq \rho\right) \leq \sum_{j \in \mathcal{J}_-} m_j \rho^{-2r} E\left[\left(\frac{u_1}{v_1} - \frac{p}{N-2}\right)^{2r}\right], \\ & \sum_{j \in \mathcal{J}_+} \sum_{i=1}^{m_j} P\left(\left|\frac{u_{j,i}}{v_i} - \frac{p + \delta_{j,i}}{N-2}\right| \geq \frac{\delta_{j,i}}{N-2}\right) \\ & \leq \sum_{j \in \mathcal{J}_+} \sum_{i=1}^{m_j} \left(\frac{\delta_{j,i}}{N-2} - \rho\right)^{-2\tilde{r}} E\left[\left(\frac{u_{j,i}}{v_i} - \frac{p + \delta_{j,i}}{N-2}\right)^{2\tilde{r}}\right], \end{aligned}$$

where \tilde{r} are any natural numbers. Hence, from the above equations and Lemma A.2, the following equations can be derived:

$$\sum_{j \in \mathcal{J}_-} m_j \rho^{-2r} E\left[\left(\frac{u_1}{v_1} - \frac{p}{N-2}\right)^{2r}\right] = O(kp^{-r}\beta^{-2r}) = o(1).$$

Note that m_j is bounded and $\#(\mathcal{J}_+) \leq k_*$, and it follows from (A.6) that $\delta_{j,i}^{-1} \leq m_j d_j \delta_{\{\ell\}}^{-1}$. Hence, for sufficiently large \tilde{r} , we have

$$\sum_{j \in \mathcal{J}_+} \sum_{i=1}^{m_j} \left(\frac{\delta_{j,i}}{N-2} - \rho\right)^{-2\tilde{r}} E\left[\left(\frac{u_{j,i}}{v_i} - \frac{p + \delta_{j,i}}{N-2}\right)^{2\tilde{r}}\right]$$

$$= \sum_{j \in \mathcal{J}_+} \sum_{i=1}^{m_j} O((p + \delta_{j,i})^{\bar{r}} \delta_{j,i}^{-2\bar{r}} + (p + \delta_{j,i})^{2\bar{r}} \delta_{j,i}^{-2\bar{r}} n^{-\bar{r}}) = o(1).$$

Therefore, from Lemma 3.2, Theorem 3.2 can be shown. \square

A.7. Proof of Lemma A.1

First, we derive results for the case of $j \subset j_*^c$. Denote the elements of j as a_1, \dots, a_{m_j} satisfying $a_s \neq a_t$ ($s \neq t$), i.e., $j = \{a_1, \dots, a_{m_j}\}$. Further, let $j_{-,0} = \omega_j$ and $j_{-,i} = j_{-,i-1} \cup \{a_i\}$ ($i = 1, \dots, m_j$). Then, it holds that $j_{-,m_j} = \omega$, and we can express $GC_p(\omega_j) - GC_p(\omega)$ as follows:

$$\begin{aligned} & GC_p(\omega_j) - GC_p(\omega) \\ &= \sum_{i=1}^{m_j} \{GC_p(j_{-,i-1}) - GC_p(j_{-,i})\} \\ &= (n-k) \sum_{i=1}^{m_j} \text{tr}[\mathbf{Y}'(\mathbf{P}_{j_{-,i}} - \mathbf{P}_{j_{-,i-1}})\mathbf{Y}\{\mathbf{Y}'(\mathbf{I}_n - \mathbf{P}_\omega)\mathbf{Y}\}^{-1}] - m_j p \alpha. \end{aligned} \quad (\text{A.9})$$

Let

$$\mathbf{W}_{j,i} = \Sigma_*^{-1/2} \mathbf{Y}'(\mathbf{P}_{j_{-,i}} - \mathbf{P}_{j_{-,i-1}})\mathbf{Y}\Sigma_*^{-1/2}, \quad \mathbf{W} = \Sigma_*^{-1/2} \mathbf{Y}'(\mathbf{I}_n - \mathbf{P}_\omega)\mathbf{Y}\Sigma_*^{-1/2}.$$

Note that $\mathbf{P}_{j_{-,i}} - \mathbf{P}_{j_{-,i-1}}$ and $\mathbf{I}_n - \mathbf{P}_\omega$ are symmetric idempotent matrices, and it holds that $(\mathbf{P}_{j_{-,i}} - \mathbf{P}_{j_{-,i-1}})(\mathbf{I}_n - \mathbf{P}_\omega) = \mathbf{O}_{n,n}$ and $(\mathbf{P}_{j_{-,i}} - \mathbf{P}_{j_{-,i-1}})\mathbf{X}_* = (\mathbf{I}_n - \mathbf{P}_\omega)\mathbf{X}_* = \mathbf{O}_{n,k_*}$. Then, from a property of the Wishart distribution and Cochran's Theorem (e.g., [5, chap 2, Theorem 2.4.2]), we can state that $\mathbf{W}_{j,i}$ and \mathbf{W} are independent, and $\mathbf{W}_{j,i} \sim W_p(1, \mathbf{I}_p)$ and $\mathbf{W} \sim W_p(n-k, \mathbf{I}_p)$. Thus, (A.9) is expressed as

$$GC_p(\omega_j) - GC_p(\omega) = (n-k) \sum_{i=1}^{m_j} \text{tr}(\mathbf{W}_{j,i} \mathbf{W}^{-1}) - m_j p \alpha. \quad (\text{A.10})$$

From a property of the Wishart distribution, $\mathbf{W}_{j,i}$ can be expressed as $\mathbf{W}_{j,i} = \mathbf{z}_i \mathbf{z}_i'$, where \mathbf{z}_i is independent of \mathbf{W} , and $\mathbf{z}_i \sim N_p(\mathbf{0}_p, \mathbf{I}_p)$. Then, we express $\mathbf{z}_i' \mathbf{W}^{-1} \mathbf{z}_i$ as

$$\mathbf{z}_i' \mathbf{W}^{-1} \mathbf{z}_i = \frac{\mathbf{z}_i' \mathbf{z}_i}{\{(\mathbf{z}_i' \mathbf{z}_i)^{-1/2} \mathbf{z}_i' \mathbf{W}^{-1} \mathbf{z}_i (\mathbf{z}_i' \mathbf{z}_i)^{-1/2}\}^{-1}}.$$

Let $u_i = \mathbf{z}_i' \mathbf{z}_i$ and $v_i = \{(\mathbf{z}_i' \mathbf{z}_i)^{-1/2} \mathbf{z}_i' \mathbf{W}^{-1} \mathbf{z}_i (\mathbf{z}_i' \mathbf{z}_i)^{-1/2}\}^{-1}$. Then, from a property of the Wishart distribution, we can state that u_i and v_i are independent, and $u_i \sim \chi^2(p)$ and $v_i \sim \chi^2(n-p-k+1)$. Therefore, $\text{tr}(\mathbf{W}_{j,i} \mathbf{W}^{-1})$ is expressed as

$$\text{tr}(\mathbf{W}_{j,i} \mathbf{W}^{-1}) = \frac{u_i}{v_i}.$$

From the above equation and (A.10), we can derive (A.1) for the case of $j \subset j_*^c$.

Next, we derive results for the case of $j \subset j_*$. Then, $GC_p(\omega_j) - GC_p(\omega)$ is expressed as

$$GC_p(\omega_j) - GC_p(\omega) = (n - k)\text{tr}[\mathbf{Y}'(\mathbf{P}_\omega - \mathbf{P}_{\omega_j})\mathbf{Y}\{\mathbf{Y}'(\mathbf{I}_n - \mathbf{P}_\omega)\mathbf{Y}\}^{-1}] - m_j p \alpha. \quad (\text{A.11})$$

Let $\mathbf{W}_j = \Sigma_*^{-1/2}\mathbf{Y}'(\mathbf{P}_\omega - \mathbf{P}_{\omega_j})\mathbf{Y}\Sigma_*^{-1/2}$. Note that $\mathbf{P}_\omega - \mathbf{P}_{\omega_j}$ is symmetric and idempotent, and it holds that $(\mathbf{P}_\omega - \mathbf{P}_{\omega_j})(\mathbf{I}_n - \mathbf{P}_\omega) = \mathbf{O}_{n,n}$. Then, from a property of the non-central Wishart distribution and Cochran's Theorem, we can state that \mathbf{W}_j and \mathbf{W} are independent, and $\mathbf{W}_j \sim W_p(m_j, \mathbf{I}_p; \Delta_j)$ and $\mathbf{W} \sim W_p(n - k, \mathbf{I}_p)$. Thus, (A.11) is expressed as

$$GC_p(\omega_j) - GC_p(\omega) = (n - k)\text{tr}(\mathbf{W}_j\mathbf{W}^{-1}) - m_j p \alpha. \quad (\text{A.12})$$

Let the spectral decomposition of Δ_j be $\Delta_j = \mathbf{Q}_j\mathbf{\Lambda}_j\mathbf{Q}_j'$, where \mathbf{Q}_j is the $p \times p$ orthogonal matrix and $\mathbf{\Lambda}_j$ is the $p \times p$ diagonal matrix whose a -th diagonal element is an eigenvalue $\lambda_{j,a}$, i.e., $\mathbf{\Lambda}_j = \text{diag}(\lambda_{j,1}, \dots, \lambda_{j,p})$ ($\lambda_{j,1} \geq \dots \geq \lambda_{j,p}$). Let $\mathbf{B}_{j,1} = \mathbf{Q}_j'\mathbf{W}_j\mathbf{Q}_j$ and $\mathbf{B}_{j,2} = \mathbf{Q}_j'\mathbf{W}\mathbf{Q}_j$. Then, from a property of the non-central Wishart distribution, we can state that $\mathbf{B}_{j,1}$ and $\mathbf{B}_{j,2}$ are independent and $\mathbf{B}_{j,1} \sim W_p(m_j, \mathbf{I}_p; \mathbf{\Lambda}_j)$ and $\mathbf{B}_{j,2} \sim W_p(n - k, \mathbf{I}_p)$. Let $d_j = \text{rank}(\Delta_j)$ be defined in (2.6). It is obvious that $\lambda_{j,d_j+1} = \dots = \lambda_{j,p} = 0$. Since it holds that $d_j \leq m_j$ from (2.7), let Γ_j be as follows:

$$\Gamma_j = \begin{pmatrix} \mathbf{\Lambda}_{j,0}^{1/2} & \mathbf{O}_{d_j, p-d_j} \\ \mathbf{O}_{m_j-d_j, d_j} & \mathbf{O}_{m_j-d_j, p-d_j} \end{pmatrix}, \quad \mathbf{\Lambda}_{j,0} = \text{diag}(\lambda_{j,1}, \dots, \lambda_{j,d_j}).$$

By using Γ_j , we can express $\mathbf{B}_{j,1}$ as $\mathbf{B}_{j,1} = (\mathbf{E}_j + \Gamma_j)'(\mathbf{E}_j + \Gamma_j)$, where $\mathbf{E}_j \sim N_{m_j \times p}(\mathbf{O}_{m_j, p}, \mathbf{I}_p \otimes \mathbf{I}_{m_j})$ and \mathbf{E}_j is independent of $\mathbf{B}_{j,2}$. Let $\mathbf{H} = (\mathbf{h}_1, \dots, \mathbf{h}_{m_j})$ be a $m_j \times m_j$ orthogonal matrix satisfying $\mathbf{h}_1 = m_j^{-1/2}\mathbf{1}_{m_j}$, and we express $\mathbf{H}\Gamma_j$ as $\mathbf{H}\Gamma_j = (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_{m_j})'$. Then, we have

$$\begin{aligned} (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_{m_j})' &= \mathbf{H} \begin{pmatrix} \mathbf{\Lambda}_{j,0}^{1/2} & \mathbf{O}_{d_j, p-d_j} \\ \mathbf{O}_{m_j-d_j, d_j} & \mathbf{O}_{m_j-d_j, p-d_j} \end{pmatrix} \\ &= (\sqrt{\lambda_{j,1}}\mathbf{h}_1, \dots, \sqrt{\lambda_{j,d_j}}\mathbf{h}_{d_j}, \mathbf{O}_{m_j, p-d_j}). \end{aligned}$$

Now, we put $\delta_{j,i} = \|\boldsymbol{\eta}_i\|^2$ ($i = 1, \dots, m_j$). Then, from the above equation, it is straightforward that $\delta_{j,i} \geq m_j^{-1}\lambda_{j,1}$ ($i = 1, \dots, m_j$) and $\text{tr}(\Delta_j) = \sum_{i=1}^{m_j} \delta_{j,i}$. Let $(\mathbf{z}_{j,1}, \dots, \mathbf{z}_{j,m_j})' = \mathbf{H}\mathbf{E}_j$. Since $\mathbf{z}_{j,1}, \dots, \mathbf{z}_{j,m_j} \sim N_p(\mathbf{0}_p, \mathbf{I}_p)$, $\mathbf{B}_{j,1}$ can be expressed as

$$\begin{aligned} \mathbf{B}_{j,1} &= (\mathbf{E}_j + \Gamma_j)' \mathbf{H}' \mathbf{H} (\mathbf{E}_j + \Gamma_j) = (\mathbf{H}\mathbf{E}_j + \mathbf{H}\Gamma_j)' (\mathbf{H}\mathbf{E}_j + \mathbf{H}\Gamma_j) \\ &= \sum_{i=1}^{m_j} (\mathbf{z}_{j,i} + \boldsymbol{\eta}_i)(\mathbf{z}_{j,i} + \boldsymbol{\eta}_i)'. \end{aligned}$$

Then, we can express $(\mathbf{z}_{j,i} + \boldsymbol{\eta}_i)' \mathbf{B}_{j,2}^{-1}(\mathbf{z}_{j,i} + \boldsymbol{\eta}_i)$ as

$$\begin{aligned} & (\mathbf{z}_{j,i} + \boldsymbol{\eta}_i)' \mathbf{B}_{j,2}^{-1}(\mathbf{z}_{j,i} + \boldsymbol{\eta}_i) \\ &= \frac{\|\mathbf{z}_{j,i} + \boldsymbol{\eta}_i\|^2}{\{ \|\mathbf{z}_{j,i} + \boldsymbol{\eta}_i\|^{-1} (\mathbf{z}_{j,i} + \boldsymbol{\eta}_i)' \mathbf{B}_{j,2}^{-1}(\mathbf{z}_{j,i} + \boldsymbol{\eta}_i) \|\mathbf{z}_{j,i} + \boldsymbol{\eta}_i\|^{-1} \}^{-1}}. \end{aligned}$$

Let $u_{j,i} = \|\mathbf{z}_{j,i} + \boldsymbol{\eta}_i\|^2$ and $v_i = \{ \|\mathbf{z}_{j,i} + \boldsymbol{\eta}_i\|^{-1} (\mathbf{z}_{j,i} + \boldsymbol{\eta}_i)' \mathbf{B}_{j,2}^{-1}(\mathbf{z}_{j,i} + \boldsymbol{\eta}_i) \|\mathbf{z}_{j,i} + \boldsymbol{\eta}_i\|^{-1} \}^{-1}$. Then, from a property of the Wishart distribution, we can state that $u_{j,i}$ and v_i are independent, and $u_{j,i} \sim \chi^2(p; \delta_{j,i})$ and $v_i \sim \chi^2(n - p - k + 1)$. Therefore, $\text{tr}(\mathbf{W}_j \mathbf{W}^{-1})$ is expressed as

$$\begin{aligned} \text{tr}(\mathbf{W}_j \mathbf{W}^{-1}) &= \text{tr}(\mathbf{Q}'_j \mathbf{W}_j \mathbf{Q}_j \mathbf{Q}'_j \mathbf{W}^{-1} \mathbf{Q}_j) = \text{tr}(\mathbf{B}_{j,1} \mathbf{B}_{j,2}^{-1}) \\ &= \sum_{i=1}^{m_j} (\mathbf{z}_{j,i} + \boldsymbol{\eta}_i)' \mathbf{B}_{j,2}^{-1}(\mathbf{z}_{j,i} + \boldsymbol{\eta}_i) = \sum_{i=1}^{m_j} \frac{u_{j,i}}{v_i}. \end{aligned}$$

From the above equation and (A.12), we can derive (A.1) for the case of $j \subset j_*$. \square

A.8. Proof of Lemma A.2

We first describe a lemma concerning the central moments of chi-square and non-central chi-square random variables; this is required for proving Lemma A.2 (the proof is given in Appendix A.9).

Lemma A.3. *Let $X_1 \sim \chi^2(t)$ and $X_2 \sim \chi^2(t; \psi)$, where ψ is a positive constant. Then, we have*

$$E[(X_1 - t)^h] = \begin{cases} 1 & (h = 0) \\ 0 & (h = 1) \\ O(t^{\lfloor h/2 \rfloor}) & (h \geq 2) \end{cases}, \quad (\text{A.13})$$

$$E[\{X_2 - (t + \psi)\}^h] = \begin{cases} 1 & (h = 0) \\ 0 & (h = 1) \\ O((t + \psi)^{\lfloor h/2 \rfloor}) & (h \geq 2) \end{cases}. \quad (\text{A.14})$$

Moreover, when $t - 2h > 0$, we have

$$E \left[\left(\frac{1}{X_1} - \frac{1}{t-2} \right)^h \right] = \begin{cases} 1 & (h = 0) \\ 0 & (h = 1) \\ O(t^{-2h + \lfloor h/2 \rfloor}) & (h \geq 2) \end{cases}, \quad (\text{A.15})$$

where $\lfloor h \rfloor$ is the floor function defined by $\lfloor h \rfloor = \max\{m \in \mathbb{Z} \mid m \leq h\}$.

Let $\xi = 1/(N - 2)$ and $\xi_\delta = p + \delta$. Then, we have

$$\frac{u_1}{v} - \frac{p}{N-2} = (u_1 - p)(v^{-1} - \xi) + p(v^{-1} - \xi) + \xi(u_1 - p),$$

$$\frac{u_2}{v} - \frac{p + \delta}{N - 2} = (u_2 - \xi_\delta)(v^{-1} - \xi) + \xi_\delta(v^{-1} - \xi) + \xi(u_2 - \xi_\delta).$$

Hence, from the multinomial theorem, we have

$$E \left[\left(\frac{u_1}{v} - \frac{p}{N - 2} \right)^{2r} \right] = \sum_{\substack{a+b+c=2r \\ 0 \leq a, b, c \leq 2r}} \frac{(2r)!}{a!b!c!} p^b \xi^c E[(u_1 - p)^{a+c}] E[(v^{-1} - \xi)^{a+b}], \tag{A.16}$$

$$E \left[\left(\frac{u_2}{v} - \frac{p + \delta}{N - 2} \right)^{2r} \right] = \sum_{\substack{a+b+c=2r \\ 0 \leq a, b, c \leq 2r}} \frac{(2r)!}{a!b!c!} \xi_\delta^b \xi^c E[(u_2 - \xi_\delta)^{a+c}] E[(v^{-1} - \xi)^{a+b}]. \tag{A.17}$$

From (A.13) and (A.15), the divergence order in (A.16) is maximized when $a = b = 0, c = 2r$ because of $pn^{-1} = O(1)$. Moreover, from (A.14) and (A.15), the divergence order in (A.17) is maximized when either $a = b = 0, c = 2r$ or $a = c = 0, b = 2r$. Therefore, we can derive the divergence orders as follows:

$$\begin{aligned} E \left[\left(\frac{u_1}{v} - \frac{p}{N - 2} \right)^{2r} \right] &= O(p^r n^{-2r}), \\ E \left[\left(\frac{u_2}{v} - \frac{p + \delta}{N - 2} \right)^{2r} \right] &= O(\max\{(p + \delta)^r n^{-2r}, (p + \delta)^{2r} n^{-3r}\}). \end{aligned} \quad \square$$

A.9. Proof of Lemma A.3

We elaborate only on the case of $h \geq 2$ because it is straightforward when $h = 0, 1$. First, we derive (A.13) and (A.14). Let h_1, \dots, h_d be natural numbers satisfying $\sum_{i=1}^d h_i = h$ and $2 \leq h_1, \dots, h_d$. From [22], we can state that h -th central moments can be expressed as the linear combination of the products of h_1, \dots, h_d -th cumulants. From [9, 23], h -th cumulants of $X_1 - t$ and $X_2 - (t + \psi)$ can, respectively, be expressed as follows:

$$\kappa_{h,1} = 2^{h-1}(h - 1)!t, \quad \kappa_{h,2} = 2^{h-1}(h - 1)!(t + h\psi).$$

Then, we observe that the maximum order term of each h -th central moment is $\kappa_{2,i}^{h/2}$ if h is even and $\kappa_{2,i}^{(h-1)/2-1} \kappa_{3,i}$ if h is odd ($i = 1, 2$). This completes (A.13) and (A.14).

Next, we derive (A.15). From the multinomial theorem, we have

$$\begin{aligned} &E \left[\left(\frac{1}{X_1} - \frac{1}{t - 2} \right)^h \right] \\ &= \sum_{i=0}^h \frac{h!}{i!(h - i)!} \left(-\frac{1}{t - 2} \right)^{h-i} E \left[\left(\frac{1}{X_1} \right)^i \right] \end{aligned}$$

$$\begin{aligned}
&= \left(-\frac{1}{t-2}\right)^h + \sum_{i=1}^h \frac{h!}{i!(h-i)!} \left(-\frac{1}{t-2}\right)^{h-i} \prod_{d=1}^i \frac{1}{t-2d} \\
&= \left(-\frac{1}{t-2}\right)^h \\
&\quad \cdot \prod_{d=1}^h \frac{1}{t-2d} \left[\{-(t-2)\}^h + \sum_{i=0}^{h-1} \frac{h!}{i!(h-i)!} \{-(t-2)\}^i \prod_{d=1}^{h-i} \{t-2h+2(d-1)\} \right].
\end{aligned}$$

Let $T \sim \chi^2(t-2h)$, then it is known that

$$E[T^{h-i}] = \begin{cases} 1 & (i = h) \\ \prod_{d=1}^{h-i} \{t-2h+2(d-1)\} & (i \leq h-1) \end{cases}.$$

Hence, by letting $s = \{-(t-2)\}^{-h} \prod_{d=1}^h (t-2d)^{-1}$, we have

$$\begin{aligned}
&E \left[\left(\frac{1}{X_1} - \frac{1}{t-2} \right)^h \right] \\
&= \left(-\frac{1}{t-2}\right)^h \prod_{d=1}^h \frac{1}{t-2d} \left\{ \sum_{i=0}^h \frac{h!}{i!(h-i)!} \{-(t-2)\}^i E[T^{h-i}] \right\} \\
&= sE[\{T - (t-2)\}^h] \\
&= s \sum_{i=0}^h \frac{h!}{i!(h-i)!} \{-2(h-1)\}^i E[\{T - (t-2h)\}^{h-i}]. \tag{A.18}
\end{aligned}$$

Note that $s = O(t^{-2h})$ and it follows from (A.13) that

$$E[\{T - (t-2h)\}^{h-i}] = \begin{cases} 1 & (i = h) \\ 0 & (i = h-1) \\ O(t^{\lfloor (h-i)/2 \rfloor}) & (i \leq h-2) \end{cases}. \tag{A.19}$$

The equations (A.18) and (A.19) complete (A.15). \square

A.10. R file related to this article

The R file ‘‘hcgcp.R’’ to perform the EZKB selection method using the $HCGC_p$ criterion can be found online at https://github.com/roda6288/hcgcp_function.

Acknowledgments

The authors would like to thank two reviewers and an Associate Editor for many constructive comments. The first author’s research is supported by a Research Fellowship for Young Scientists from the Japan Society for the Promotion of Science. The second author’s research is partially supported by a Grant-in-Aid for Scientific Research (C) from the Ministry of Education, Science, Sports, and Culture #18K03415.

References

- [1] ATKINSON, A. C. (1980). A note on the generalized information criterion for choice of a model. *Biometrika* **67** 413–418.
- [2] CHUDIK, A., KAPETANIOS, G. and PESARAN, M. H. (2018). A one covariate at a time, multiple testing approach to variable selection in high-dimensional linear regression models. *Econometrica* **86** 1479–1512. [MR3843496](#)
- [3] FUJIKOSHI, Y. and SATOH, K. (1997). Modified AIC and C_p in multivariate linear regression. *Biometrika* **84** 707–716. [MR1603952](#)
- [4] FUJIKOSHI, Y., SAKURAI, T. and YANAGIHARA, H. (2014). Consistency of high-dimensional AIC-type and C_p -type criteria in multivariate linear regression. *J. Multivariate Anal.* **123** 184–200. [MR3130429](#)
- [5] FUJIKOSHI, Y., ULYANOV, V. V. and SHIMIZU, R. (2010). *Multivariate Statistics: High-Dimensional and Large-Sample Approximations*. John Wiley & Sons, Inc., Hoboken, New Jersey. [MR2640807](#)
- [6] FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R. (2010). Regularization paths for generalized linear models via coordinate descent. *J. Stat. Softw.* **33(1)** 1–22. [MR1082147](#)
- [7] HARVILLE, D. A. (1997). *Matrix Algebra from a Statistician's Perspective*. Springer-Verlag, New York. [MR1467237](#)
- [8] HE, Y., JIANG, T., WEN, J. and XU, G. (2018). Likelihood ratio test in multivariate linear regression: from low to high dimension. [arXiv:1812.06894](#).
- [9] LANCASTER, H. O. (1982). Chi-square distribution. In *Encyclopedia of Statistical Sciences, Vol. 1* (eds. S. Kotz & N. L. Johnson), 439–442. John Wiley & Sons, New York. [MR0719028](#)
- [10] LI, R., ZHONG, W. and ZHU, L. (2012). Feature screening via distance correlation learning. *J. Amer. Statist. Assoc.* **107** 1129–1139. [MR3010900](#)
- [11] LI, Y., NAN, B. and ZHU, J. (2015). Multivariate sparse group lasso for the multivariate multiple linear regression with an arbitrary group structure. *Biometrics* **71** 354–363. [MR3366240](#)
- [12] LUO, S. (2018). Variable selection in high-dimensional sparse multi-response linear regression models. *Stat. Pap.* <https://doi.org/10.1007/s00362-018-0989-x>.
- [13] MALLOWS, C. L. (1973). Some comments on C_p . *Technometrics* **15** 661–675.
- [14] MALLOWS, C. L. (1995). More comments on C_p . *Technometrics* **37** 362–372. [MR1365719](#)
- [15] NAGAI, I., YANAGIHARA, H. and SATOH, K. (2012). Optimization of ridge parameters in multivariate generalized ridge regression by plug-in methods. *Hiroshima Math. J.* **42** 301–324. [MR3050124](#)
- [16] NISHII, R. (1988). Maximum likelihood principle and model selection when the true model is unspecified. *J. Multivariate Anal.* **27** 392–403. [MR0970962](#)
- [17] NISHII, R., BAI, Z. D. and KRISHNAIAH, P. R. (1988). Strong consistency

- of the information criterion for model selection in multivariate analysis. *Hiroshima Math. J.* **18** 451–462. [MR0991240](#)
- [18] PENG, J., ZHU, J., BERGAMASCHI, A., HAN, W., NOH, D. Y., POLLACK, J. R. and WANG, P. (2010). Regularized multivariate regression for identifying master predictors with application to integrative genomics study of breast cancer. *Ann. Appl. Stat.* **4**(1) 53–77. [MR2758084](#)
- [19] RAO, C. R. and WU, Y. (1989). A strongly consistent procedure for model selection in a regression problem. *Biometrika* **76** 369–374. [MR1016028](#)
- [20] SPARKS, R. S., COUTSOURIDES, D. and TROSKIE, L. (1983). The multivariate C_p . *Comm. Statist. A – Theory Methods* **12** 1775–1793. [MR0704853](#)
- [21] SRIVASTAVA, M. S. (2002). *Methods of Multivariate Statistics*. John Wiley & Sons, New York. [MR1915968](#)
- [22] STUART, A. and ORD, J. K. (1994). *Kendall’s Advanced Theory of Statistics. Vol. 1. Distribution Theory* (6th ed.). Edward Arnold, London; distributed in the United States of America by Oxford University Press, New York. [MR0246399](#)
- [23] TIKU, M. (1985). Noncentral chi-square distribution. In *Encyclopedia of Statistical Sciences, Vol. 6* (eds. S. Kotz & N. L. Johnson), 276–280, John Wiley & Sons, New York.
- [24] TIMM, N. H. (2002). *Applied Multivariate Analysis*. Springer-Verlag, New York. [MR1908225](#)
- [25] WANG, H. and LENG, C. (2008). A note on adaptive group lasso. *Comput. Stat. Data An.* **52** 5277–5286. [MR2526593](#)
- [26] XIN, X., HU, J. and LIU, L. (2017). On the oracle property of a generalized adaptive elastic-net for multivariate linear regression with a diverging number of parameters. *J. Multivariate Anal.* **162** 16–32. [MR3719332](#)
- [27] 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.
- [28] YANAGIHARA, H., WAKAKI, H. and 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. [MR3338666](#)
- [29] YUAN, M. and LIN, Y. (2006). Model selection and estimation in regression with grouped variables. *J. Roy. Stat. Soc. B* **68** 49–67. [MR2212574](#)
- [30] YUAN, M., EKICI, A., LU, Z. and MONTEIRO, R. (2007). Dimension reduction and coefficient estimation in multivariate linear regression. *J. Roy. Stat. Soc. B* **69** 329–346. [MR2323756](#)
- [31] ZHAO, L. C., KRISHNAIAH, P. R. and BAI, Z. D. (1986). On detection of the number of signals in presence of white noise. *J. Multivariate Anal.* **20** 1–25. [MR0862239](#)