

# High-dimensional consistency of rank estimation criteria in multivariate linear model



Yasunori Fujikoshi<sup>a</sup>, Tetsuro Sakurai<sup>b,\*</sup>

<sup>a</sup> Department of Mathematics, Graduate School of Science, Hiroshima University, 1-3-1 Kagamiyama, Higashi Hiroshima, Hiroshima 739-8626, Japan

<sup>b</sup> Center of General Education, Tokyo University of Science, Suwa 5000-1 Toyohira, Chino, Nagano 391-0292, Japan

## ARTICLE INFO

### Article history:

Received 13 April 2015

Available online 6 May 2016

### AMS 2000 subject classifications:

primary 62H12

secondary 62H10

### Keywords:

AIC

BIC

$C_p$

Consistency property

Dimensionality

Discriminant analysis

High-dimensional framework

Multivariate regression model

Multivariate linear model

Rank

Ridge-type criterion

Tuning parameter

## ABSTRACT

This paper is concerned with consistency properties of rank estimation criteria in a multivariate linear model, based on the model selection criteria AIC, BIC and  $C_p$ . The consistency properties of these criteria are studied under a high-dimensional framework with two different assumptions on the noncentrality matrix such that the number of response variables and the sample size tend to infinity. In general, it is known that under a large-sample asymptotic framework, the criteria based on AIC and  $C_p$  are not consistent, but the criterion based on BIC is consistent. However, we note that there are cases that the criteria based on AIC and  $C_p$  are consistent, but the criterion based on BIC is not consistent. Such consistency properties are also shown for the generalized criteria with a tuning parameter. Further, the modified criteria with a ridge-type estimator are also examined. Through a Monte Carlo simulation experiment, our results are checked numerically, and the estimation criteria are compared.

© 2016 Elsevier Inc. All rights reserved.

## 1. Introduction

This paper is concerned with the problem of estimating the rank (or dimensionality) of a regression coefficient matrix in a multivariate linear model. A multivariate regression model with a reduced rank is called a multivariate reduced-rank regression model (see Izenman [15], Reisel and Velu [19]). Our problem includes also that of estimating the number of meaningful discriminant functions in discriminant analysis.

One commonly used rank estimation method is based on sequential test procedures. In this method, the test in each step is a likelihood ratio test which was first proposed by Anderson [2], [3]. Recently, related to sparse reduced-rank regression, some penalized methods have been studied. Yuan et al. [23] used the Ky-Fan norm penalty for factor selection and shrinkage. Chen and Huang [8] proposed a simultaneous method for selecting the rank and the explanatory variables by using penalized regression with a group lasso penalty. For related works, see Bunea et al. [6], [7]. These penalized methods all involve finding the multivariate least squares under a restriction.

\* Corresponding author.

E-mail addresses: [fujikoshi\\_y@yahoo.co.jp](mailto:fujikoshi_y@yahoo.co.jp) (Y. Fujikoshi), [sakurai@rs.tus.ac.jp](mailto:sakurai@rs.tus.ac.jp) (T. Sakurai).

There are other methods based on the use of model selection criteria as follows. In a general multivariate regression, we have  $n$  observations  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)^\top$  and  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$  on  $p$  response variables  $\mathbf{y} = (y_1, \dots, y_p)^\top$  and  $q$  explanatory variables  $\mathbf{x} = (x_1, \dots, x_k)^\top$ , and

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\Theta} + \boldsymbol{\varepsilon}, \quad (1)$$

where  $\boldsymbol{\Theta}$  is a  $k \times p$  regression coefficient matrix, and  $\boldsymbol{\varepsilon} = (\epsilon_1, \dots, \epsilon_n)^\top$  is the error matrix, with the  $\epsilon_i$ 's independently distributed as  $\mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$ .

Our interest is to estimate the rank of  $\mathbf{C}\boldsymbol{\Theta}$ , where  $\mathbf{C}$  is a given  $q \times k$  matrix with rank  $q$ , and hence  $q \leq k$ . As the model selection approach, we consider the models  $M_j$ ,  $j = 0, \dots, m = \min(p, q)$ , where  $M_j$  is the model (1) with  $\text{rank}(\mathbf{C}\boldsymbol{\Theta}) = j$ . If  $M_j$  is selected by a model selection criterion, we estimate the rank as  $j$ . We are concerned with rank estimation methods that use the following model selection criteria: AIC (Akaike [1]),  $C_p$  (Mallows [18]) and BIC (Schwarz [20]). The rank estimation criteria based on AIC and  $C_p$  were proposed by Fujikoshi and Veitch [13] in multivariate linear model and canonical correlation analysis. In this context, Gunderson and Muirhead [14] consider a rank estimation criterion based on BIC. These rank estimation methods are denoted herein as AIC, BIC and  $C_p$ , for simplicity.

In general, a selection criterion is said to be consistent if the probability of selecting the true model tends to 1 under an asymptotic framework. Our rank estimation criteria depend on  $p, q, k, n$  and  $\boldsymbol{\Omega} = \text{diag}(\omega_1, \dots, \omega_q)$ , where  $\omega_1 \geq \dots \geq \omega_q \geq 0$  are the possible nonzero characteristic roots of the (standardized) noncentrality matrix  $\tilde{\boldsymbol{\Omega}}$  in (7). It is known (Fujikoshi [9]) that, under the large-sample asymptotic framework

$$n \rightarrow \infty; \quad p, q, k: \text{fixed}, \quad (2)$$

with  $\boldsymbol{\Omega} = O(n)$ , AIC and  $C_p$  are not consistent. Here, for a matrix  $\mathbf{A}$ ,  $\mathbf{A} = O(n)$  means that each element of  $\mathbf{A}$  is order  $n$ . Gunderson and Muirhead [14] showed the consistency of BIC in canonical correlation analysis.

One of our purposes is to study consistency properties of AIC, BIC and  $C_p$  when  $p/n$  tends to  $c$ . More precisely, we consider the high-dimensional asymptotic framework given by

$$q: \text{fixed}, \quad k \rightarrow \infty, \quad k/n \rightarrow 0, \quad p \rightarrow \infty, \quad n \rightarrow \infty, \quad p/n \rightarrow c \in [0, 1), \quad (3)$$

with each for (i)  $\boldsymbol{\Omega} = O(n)$  and (ii)  $\boldsymbol{\Omega} = O(np)$ . In the present study, it is shown that the criteria AIC and  $C_p$  are consistent under some additional assumptions depending on (i) and (ii). Further, we find that the criterion BIC is consistent under  $\boldsymbol{\Omega} = O(np)$ , but that it is not consistent under  $\boldsymbol{\Omega} = O(n)$ .

Instead of the criteria AIC, BIC and  $C_p$ , we may use the criteria A, B and C defined by

$$A_j = \text{AIC}_j - \text{AIC}_p, \quad B_j = \text{BIC}_j - \text{BIC}_m, \quad C_j = C_{p,j} - C_{p,m}.$$

For the details of these criteria, refer to (12), (13) and (14) in Section 3. In this paper, we propose generalizations and modifications of A, B and C. One such generalization is the criteria  $\text{IC}_\nu$  and  $C_{p,\nu}$  with a tuning parameter  $\nu$ , where  $\text{IC}_2 = A$ ,  $\text{IC}_{\ln n} = B$  and  $C_{p,2} = C$ . The other is the ridge-type criteria of A, B and C denoted by  $A_\lambda$ ,  $B_\lambda$  and  $C_\lambda$ , respectively. The ridge-type criteria are defined by using the estimator of the covariance matrix with a ridge parameter. The criteria are defined also for the case  $p > n - k$ . Our notation, for example,  $A_j$ ,  $A_\nu$ ,  $A_\lambda$ , might cause confusion. However, note that we use  $i$  and  $j$  for possible ranks,  $\nu$  for a tuning parameter, and  $\lambda$  for the ridge parameter. We derive sufficient conditions for the generalized criteria  $\text{IC}_\nu$  and  $C_{p,\nu}$  to be consistent under the high-dimensional asymptotic framework (3). Such sufficient conditions will be useful in selecting the tuning parameter  $\nu$ . For  $n - k > p$ , we note that the ridge-type criteria  $A_\lambda$ ,  $B_\lambda$  and  $C_\lambda$  have the same consistency properties as those of A, B and C under a high-dimensional asymptotic framework. We check our theoretical results through a Monte Carlo simulation experiment and compare numerically the selection probabilities of the criteria A, B and C with their modifications. For  $p > n - k$ , some tendencies in the ridge-type criteria are clarified through numerical experiments.

Our results may be applied to data sets in multivariate regression model and discriminant analyses with a relatively large number of response variables, which are carried out in many fields, including medicine. Examples of data with a very large number of response variables include stock data and genomic data. We may also apply our results to the analysis of data obtained by choosing fewer response variables less than the sample size. Similar high-dimensional consistency properties have been derived by Fujikoshi et al. [11] and Yanagihara et al. [22] for AIC, BIC and  $C_p$  in selection of the explanatory variables in a multivariate linear model.

The present paper is organized as follows. In Section 2, we present a reduced-rank multivariate linear model, and two important special cases are explained. Then we prepare three criteria and their modifications with a tuning parameter. In Section 3, we give sufficient conditions for the criteria  $\text{IC}_\nu$  and  $C_{p,\nu}$  to be consistent. As a special case, we give consistency properties of the criteria A, B and C. We check our theoretical results by conducting a Monte Carlo simulation experiment, and compare with the selection probabilities of the two criteria. In Section 4, we propose ridge-type criteria, whose consistency properties are then theoretically and numerically examined. In Section 5, we discuss our conclusions. The proofs of our results are given in the Appendix.

## 2. Rank estimation criteria

### 2.1. Reduced-rank multivariate linear model

The multivariate regression model in (1) may be called a multivariate linear model when the  $x_i$ 's are any fixed variables, including explanatory variables or dummy variables. The model may be expressed as

$$\mathbf{Y} \sim \mathcal{N}_{n \times p}(\mathbf{X}\boldsymbol{\Theta}, \boldsymbol{\Sigma} \otimes \mathbf{I}_n), \quad (4)$$

where  $\boldsymbol{\Theta}$  is a  $k \times p$  unknown matrix of coefficients,  $\boldsymbol{\Sigma}$  is a  $p \times p$  unknown covariance matrix, and  $\mathbf{I}_n$  is the identity matrix of order  $n$ . The notation  $\mathcal{N}_{n \times p}(\cdot, \cdot)$  means the matrix normal distribution such that the mean of  $\mathbf{Y}$  is  $\mathbf{X}\boldsymbol{\Theta}$  and the covariance matrix of  $\text{vec}(\mathbf{Y})$  is  $\boldsymbol{\Sigma} \otimes \mathbf{I}_n$ , where  $\text{vec}(\mathbf{Y})$  is the  $np \times 1$  vector formed by stacking the columns of  $\mathbf{Y}$  under each other. We assume that  $\text{rank}(\mathbf{X}) = k$  and in most case,  $n - k \geq p$ . However, the restriction  $n - k \geq p$  is not necessarily assumed when we treat the ridge-type criterion.

Let  $\mathbf{C}$  be a given  $q \times k$  matrix with  $\text{rank}(\mathbf{C}) = q$ , and consider reduced-rank models  $M_j$ ; (1) or (4) with  $\text{rank}(\mathbf{C}\boldsymbol{\Theta}) = j$ ,  $j = 0, \dots, m = \min(p, q)$ . For simplicity, we write  $M_j$  as

$$M_j : \text{rank}(\mathbf{C}\boldsymbol{\Theta}) = j, \quad j = 0, \dots, m. \quad (5)$$

Our interest is to estimate  $\text{rank}(\mathbf{C}\boldsymbol{\Theta})$  by the criteria based on model selection criteria AIC, BIC and  $C_p$ . An LR statistic for  $M_j : \text{rank}(\mathbf{C}\boldsymbol{\Theta}) = j$  is given (see, e.g., Anderson [3]) by

$$\Lambda_{(j)} = \{(1 + \ell_{j+1}) \cdots (1 + \ell_m)\}^{-1}, \quad (6)$$

where  $\ell_1 \geq \cdots \geq \ell_m \geq 0$  are the possible non-zero characteristic roots of  $\mathbf{S}_h \mathbf{S}_e^{-1}$ ,

$$\mathbf{S}_e = \mathbf{Y}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{Y}, \quad \mathbf{S}_h = (\mathbf{C}\hat{\boldsymbol{\Theta}})^\top \{\mathbf{C}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{C}^\top\}^{-1} \mathbf{C}\hat{\boldsymbol{\Theta}},$$

and  $\hat{\boldsymbol{\Theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}$ . Here, without loss of generality we may assume that  $\ell_1 > \cdots > \ell_m > 0$ , since the probability of any two  $\ell_i$ s being equal is 0. It is well known (see, e.g., Anderson [3]) that  $\mathbf{S}_e$  and  $\mathbf{S}_h$  are independently distributed as a Wishart distribution  $\mathcal{N}_p(n - k, \boldsymbol{\Sigma})$  and a noncentral Wishart distribution  $\mathcal{N}_p(q, \boldsymbol{\Sigma}; \boldsymbol{\Sigma}^{1/2} \tilde{\boldsymbol{\Omega}} \boldsymbol{\Sigma}^{1/2})$ , respectively, where

$$\tilde{\boldsymbol{\Omega}} = \boldsymbol{\Sigma}^{-1/2} (\mathbf{C}\boldsymbol{\Theta})^\top \{\mathbf{C}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{C}^\top\}^{-1} \mathbf{C}\boldsymbol{\Theta} \boldsymbol{\Sigma}^{-1/2}. \quad (7)$$

The reduced-rank model includes the following two important special cases. One is a multivariate reduced-rank regression model which is given by (4) and (5) with  $\mathbf{C} = \mathbf{I}_k$ . From the rank constraint (5), the regression matrix  $\boldsymbol{\Theta}$  can be expressed as a product of two rank  $j$  matrices as follows:

$$\boldsymbol{\Theta} = \mathbf{G}\boldsymbol{\Xi},$$

where  $\mathbf{G}$  is of dimension  $k \times j$  and  $\boldsymbol{\Xi}$  is of dimension  $j \times p$ . Then

$$\mathbf{E}[\mathbf{Y}] = (\mathbf{X}\mathbf{G}) \cdot \boldsymbol{\Xi} = \mathbf{Z}\boldsymbol{\Xi}, \quad \mathbf{Z} = \mathbf{X}\mathbf{G}. \quad (8)$$

The model means that the  $j$  linear combinations  $\mathbf{z} = \mathbf{G}^\top \mathbf{x}$  of the  $k$  explanatory variables  $\mathbf{x}$  are sufficient to model the variation in the  $p$  response variables  $\mathbf{y}$ . The  $j$ -variate  $\mathbf{z}$  may be regarded as a factor or latent variate. In practice, the dimension  $j$  is unknown, and we need to estimate it.

The other special case is a reduced-rank problem in discriminant analysis, based on  $(q + 1)$   $p$ -variate normal populations with common covariance matrix  $\boldsymbol{\Sigma}$ . Let  $\boldsymbol{\mu}_i$  be the mean vector of the  $i$ th population. Suppose that a sample of size  $n_i$  is available from the  $i$ th population, and let  $\mathbf{y}_{ij}$  be the  $j$ th observation from the  $i$ th population. Then, the model of  $\mathbf{Y} = (\mathbf{y}_{11}, \dots, \mathbf{y}_{1n_1}, \dots, \mathbf{y}_{q+1,1}, \dots, \mathbf{y}_{q+1,n_{q+1}})^\top$  is a model (4) with  $k = q + 1$  and

$$\mathbf{X} = \begin{pmatrix} \mathbf{1}_{n_1} & \cdots & \mathbf{0} \\ \vdots & \ddots & \vdots \\ \mathbf{0} & \cdots & \mathbf{1}_{n_{q+1}} \end{pmatrix}, \quad \boldsymbol{\Theta} = \begin{pmatrix} \boldsymbol{\mu}_1^\top \\ \vdots \\ \boldsymbol{\mu}_{q+1}^\top \end{pmatrix}.$$

Now we consider the reduced-rank model with  $\mathbf{C} = (\mathbf{I}_q, -\mathbf{1}_q)$ , and  $\mathbf{C}\boldsymbol{\Theta} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_{q+1}, \dots, \boldsymbol{\mu}_q - \boldsymbol{\mu}_{q+1})$ . Then the matrices  $\mathbf{S}_e$  and  $\mathbf{S}_h$  are expressed as

$$\mathbf{S}_b = \sum_{i=1}^{q+1} n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}})(\bar{\mathbf{y}}_i - \bar{\mathbf{y}})^\top, \quad \mathbf{S}_w = \sum_{i=1}^{q+1} (n_i - 1) \mathbf{S}_i,$$

respectively, where  $\bar{\mathbf{y}}_i$  and  $\mathbf{S}_i$  are the mean vector and sample covariance matrix of the  $i$ th population, and  $\bar{\mathbf{y}}$  is the total mean vector defined by  $(1/n) \sum_{i=1}^{q+1} n_i \bar{\mathbf{y}}_i$ , in which  $n = \sum_{i=1}^{q+1} n_i$ . Then,  $\mathbf{S}_b$  and  $\mathbf{S}_w$  are called the between-group and the within-group sums of squares and products matrices. In general,  $\mathbf{S}_w$  and  $\mathbf{S}_b$  are independently distributed as a Wishart distribution  $\mathcal{N}_p(n - q - 1, \boldsymbol{\Sigma})$  and a noncentral Wishart distribution  $\mathcal{N}_p(q, \boldsymbol{\Sigma}; \boldsymbol{\Sigma}^{1/2} \tilde{\boldsymbol{\Omega}} \boldsymbol{\Sigma}^{1/2})$ , respectively, where

$$\tilde{\boldsymbol{\Omega}} = \boldsymbol{\Sigma}^{-1/2} \sum_{i=1}^{q+1} n_i (\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})^\top \boldsymbol{\Sigma}^{-1/2}, \quad \bar{\boldsymbol{\mu}} = (1/n) \sum_{i=1}^{q+1} n_i \boldsymbol{\mu}_i. \quad (9)$$

It is known (see, e.g., Kshirsagar [16]) that

$$M_j : \text{rank}(\mathbf{C}\boldsymbol{\Theta}) = j \Leftrightarrow \omega_1 \geq \cdots \geq \omega_j > \omega_{j+1} = \cdots = \omega_q = 0, \quad (10)$$

where  $\omega_1 \geq \omega_2 \geq \cdots \geq \omega_m \geq 0$  are the possible non-zero characteristic roots of  $\tilde{\boldsymbol{\Omega}}$ .

## 2.2. Criteria based on AIC, BIC, $C_p$ and their modifications

In general, AIC for a model  $M$  is defined (Akaike [1]) as

$$\text{AIC} = -2 \ln \hat{L} + 2d,$$

where  $\hat{L}$  is the maximum likelihood under  $M$ , and  $d$  is the number of independent parameters under  $M$ . The AIC for  $M_j$  is expressed as

$$\begin{aligned} \text{AIC}_j &= n \ln(1 + \ell_{j+1}) \cdots (1 + \ell_m) + n \ln |(1/n)\mathbf{S}_e| + np(\ln 2\pi + 1) \\ &\quad + 2 \left\{ j(p + q - j) + (k - q)p + \frac{1}{2}p(p + 1) \right\}. \end{aligned} \quad (11)$$

The expression (11) was obtained by Fujikoshi and Veitch [13] as an asymptotic unbiased estimator of the risk function based on Kullback–Leibler distance.

Based on  $\text{AIC}_j$ , if  $\min\{\text{AIC}_0, \dots, \text{AIC}_m\} = \text{AIC}_j$ , we estimate the rank as  $j$ . Instead of  $\text{AIC}_j$ , we may use

$$\begin{aligned} A_j &= \text{AIC}_j - \text{AIC}_m \\ &= n \ln \prod_{i=j+1}^m (1 + \ell_i) - 2(p - j)(q - j), \quad j = 0, \dots, m. \end{aligned} \quad (12)$$

Here,  $A_m = 0$ . Then the rank estimation method is equivalent to estimating the rank as  $j$  if  $\min\{A_0, \dots, A_m\} = A_j$ .

Similarly, the rank estimation criteria  $B_j$  and  $C_j$  based on BIC and  $C_p$  are given as follows.

$$\begin{aligned} B_j &= \text{BIC}_j - \text{BIC}_m \\ &= n \ln \prod_{i=j+1}^m (1 + \ell_i) - (\ln n)(p - j)(q - j), \quad j = 0, \dots, m. \end{aligned} \quad (13)$$

$$\begin{aligned} C_j &= C_{p,j} - C_{p,m} \\ &= n \sum_{i=j+1}^m \ell_i - 2(p - j)(q - j), \quad j = 0, \dots, m. \end{aligned} \quad (14)$$

Here,  $B_m = 0$  and  $C_m = 0$ .

Using a tuning parameter  $v$ , we consider the following two modified criteria:

$$\text{IC}_{v,j} = n \ln \prod_{i=j+1}^m (1 + \ell_i) - v(p - j)(q - j), \quad j = 0, \dots, m. \quad (15)$$

$$C_{p,v,j} = n \sum_{i=j+1}^m \ell_i - v(p - j)(q - j), \quad j = 0, \dots, m. \quad (16)$$

Here,  $\text{IC}_{v,m} = 0$  and  $C_{p,v,m} = 0$ . Then,

$$\text{IC}_{2,j} = A_j, \quad \text{IC}_{\ln n,j} = B_j, \quad C_{p,2,j} = C_j.$$

## 3. High-dimensional consistency

### 3.1. Theoretical results

In this subsection, we are concerned with asymptotic behaviors of the criteria when  $p$  and  $n$  are large and  $q$  is fixed. Without loss of generality, we assume that  $p \geq q$ , so that  $m = \min(p, q) = q$ . Let the set of all possible ranks be denoted by  $\mathcal{F} = \{0, \dots, q\}$ . It is assumed that the true model is the model (4) with  $\boldsymbol{\Theta} = \boldsymbol{\Theta}_*$  and  $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_*$ , which is denoted by  $M_*$ . The true rank is defined by  $\text{rank}(\mathbf{C}\boldsymbol{\Theta}_*) = j_*$ . Then, the model  $M_{j_*}$  is the minimum reduced-rank model including  $M_*$ . For notational simplicity, we often write  $\boldsymbol{\Theta}_*$  and  $\boldsymbol{\Sigma}_*$  as simply  $\boldsymbol{\Theta}$  and  $\boldsymbol{\Sigma}$ .

We separate  $\mathcal{F}$  into two sets,  $\mathcal{F}_+ = \{j_*, \dots, q\}$ , the set of ranks of overspecified models, and  $\mathcal{F}_- = \mathcal{F}_+^c \cap \mathcal{F} = \{0, \dots, j_* - 1\}$ , the set of ranks of underspecified models. Further, we denote the set of ranks after deleting the true model from  $\mathcal{F}_+$  by  $\mathcal{F}_+ \setminus \{j_*\}$ , i.e.,  $\mathcal{F}_+ \setminus \{j_*\} = \{j_* + 1, \dots, q\}$ .

The estimation methods based on  $A_j$ ,  $B_j$  and  $C_j$  are expressed as

$$\hat{j}_A = \arg \min_{j \in \mathcal{F}} A_j, \quad \hat{j}_B = \arg \min_{j \in \mathcal{F}} B_j, \quad \text{and} \quad \hat{j}_C = \arg \min_{j \in \mathcal{F}} C_j,$$

respectively. Similarly, the estimation methods based on  $IC_{v,j}$  and  $C_{p,v,j}$  are expressed as

$$\hat{j}_{IC_v} = \arg \min_{j \in \mathcal{F}} IC_{v,j}, \quad \text{and} \quad \hat{j}_{C_{p,v}} = \arg \min_{j \in \mathcal{F}} C_{p,v,j},$$

respectively. In the following, we summarize some assumptions.

A1 (The true model  $M_*$  and rank  $j_*$ ): The true model is (1) or (4) with  $\Theta = \Theta_*$  and  $\Sigma = \Sigma_*$ . The true rank is  $\text{rank}(\mathbf{C}\Theta_*) = j_*$ .

A2 (The asymptotic framework):  $q$  is fixed,  $p \rightarrow \infty$ ,  $n \rightarrow \infty$ ,  $k \rightarrow \infty$ ,  $p/n \rightarrow c \in [0, 1)$ ,  $k/n \rightarrow 0$ .

Further, we make two types of assumptions on the order of the noncentrality matrix  $\tilde{\Omega}$  in (7). Since  $\text{rank}(\tilde{\Omega}) \leq q$ , we can write  $\tilde{\Omega} = \Gamma\Gamma^\top$ , where  $\Gamma$  is a  $p \times q$  matrix. Let

$$\Omega = \Gamma^\top \Gamma, \quad (17)$$

which is a  $q \times q$  matrix. In discriminant analysis with  $q = 1$ ,

$$\tilde{\Omega} = \frac{n_1 n_2}{n} \Sigma^{-1/2} (\mu_1 - \mu_2) (\mu_1 - \mu_2)^\top \Sigma^{-1/2},$$

and

$$\Omega = \omega = \frac{n_1 n_2}{n} (\mu_1 - \mu_2)^\top \Sigma^{-1} (\mu_1 - \mu_2),$$

which is  $(n_1 n_2 / n)$  times the squared Mahalanobis distance between two normal populations  $\mathcal{N}_p(\mu_1, \Sigma)$  and  $\mathcal{N}_p(\mu_2, \Sigma)$ . When  $n_i/n \rightarrow d_i > 0$ ,  $\omega = O(n)$  and also  $\omega = O(np)$ , depending on whether the squared Mahalanobis distance is  $O(1)$  or  $O(p)$ , where  $O(\cdot)$  is the usual order under a high-dimensional framework (3). Note that, when we consider the distributions of our criteria, without loss of generality we may assume

$$\Omega = \text{diag}(\omega_1, \dots, \omega_q), \quad (18)$$

where  $\omega_1 \geq \dots \geq \omega_q$  are the characteristic roots of  $\Omega$  or the non-zero characteristic roots of  $\tilde{\Omega}$ . Based on these considerations, we take up the following two types of assumptions on the noncentrality matrix.

A3 (The noncentrality matrix-1): For any  $i(1 \leq i \leq j_*)$ ,

$$\omega_i = n\delta_i = O(n), \quad \lim_{p/n \rightarrow c} \delta_i = \delta_i^* > 0.$$

A4 (The noncentrality matrix-2): For any  $i(1 \leq i \leq j_*)$ ,

$$\omega_i = np\xi_i = O(np), \quad \lim_{p/n \rightarrow c} \xi_i = \xi_i^* > 0.$$

In A3 and A4, it is assumed that the multiplicities of the  $\omega_i$ 's do not depend on  $p$  and  $n$ . In the following, we give sufficient conditions for  $IC_v$  and  $C_{p,v}$  to be consistent. Here, the consistency of, e.g.,  $IC_v$  means that the probability that  $IC_v$  selects the true rank  $j_*$  tends asymptotically to 1, i.e.,

$$\lim_{p/n \rightarrow c} \Pr(\hat{j}_{IC_v} = j_*) = 1.$$

Here, the notation  $\lim_{p/n \rightarrow c}$  is used as an abbreviation for the asymptotic framework A2 or (3).

**Theorem 3.1.** Suppose that assumption A1 is satisfied.

(1)  $IC_v$  is consistent if assumptions A2 and A3 and the inequality  $-c^{-1} \ln(1 - c) < v < -c^{-1} \ln(1 - c) + c^{-1} \ln(1 + \delta_{j_*}^*)$  are satisfied.

(2)  $IC_v$  is consistent if assumptions A2 and A4 and the inequality  $-c^{-1} \ln(1 - c) < v$  are satisfied.

When  $c = 0$  in Theorem 3.1, “ $c^{-1} \ln(1 - c)$ ” should be read as  $\lim_{c \rightarrow 0+} c^{-1} \ln(1 - c) = -1$ , and “ $c^{-1} \ln(1 + \delta_{j_*}^*)$ ” should be read as “ $\infty$ ”. Thus, when  $c = 0$ , “ $-c^{-1} \ln(1 - c) < v < -c^{-1} \ln(1 - c) + c^{-1} \ln(1 + \delta_{j_*}^*)$ ” should be read as “ $-c^{-1} \ln(1 - c) < v$ ”. These conventions will be used elsewhere herein. Theorem 3.1 implies the following Corollaries 3.1 and 3.2 except for Corollary 3.2(1).

**Table 3.1**  
Boundary values of  $\nu$  for consistency of  $IC_\nu$  and  $C_{p,\nu}$ .

$c$	$b_1(c) = -c^{-1} \ln(1 - c)$	$b_2(c) = (1 - c)^{-1}$
0.1	1.05	1.11
0.2	1.12	1.25
0.3	1.19	1.43
0.4	1.28	1.67
0.5	1.39	2.00
0.6	1.53	2.50
0.7	1.72	3.33
0.8	2.01	5.00
0.9	2.56	10.00

**Corollary 3.1.** Suppose that assumption A1 is satisfied. Further, assume that  $c \in [0, c_a)$ , where  $c_a$  ( $\approx 0.797$ ) is the larger constant satisfying  $\ln(1 - c_a) + 2c_a = 0$ .

- (1) AIC is consistent if assumptions A2 and A3 and the inequality  $\ln(1 + \delta_{j_*}^*) > (j_* - j)\{2c + \ln(1 - c)\}$  are satisfied.
- (2) AIC is consistent if assumptions A2 and A4 are satisfied.

**Corollary 3.2.** Suppose that assumption A1 is satisfied.

- (1) BIC is not consistent if assumptions A2 with  $c > 0$  and A3 are satisfied.
- (2) BIC is consistent if assumptions A2 and A4 are satisfied.

Similar results are obtained for the criteria  $C_{p,\nu}$  and  $C_p$ .

**Theorem 3.2.** Suppose that assumption A1 is satisfied.

- (1)  $C_{p,\nu}$  is consistent if assumptions A2 and A3 and the inequality  $(1 - c)^{-1} < \nu < (1 - c)^{-1} + \{c(1 - c)\}^{-1}\delta_{j_*}^*$  are satisfied.
- (2)  $C_{p,\nu}$  is consistent if assumptions A2 and A4 and the inequality  $(1 - c)^{-1} < \nu$  are satisfied.

When  $c = 0$ , “ $(1 - c)^{-1} < \nu < (1 - c)^{-1} + \{c(1 - c)\}^{-1}\delta_{j_*}^*$ ” should be read as “ $1 < \nu$ ”.

**Corollary 3.3.** Suppose that assumption A1 is satisfied. Further, assume that  $c \in [0, 0.5)$ .

- (1)  $C_p$  is consistent if assumptions A2 and A3 and the inequality  $\delta_{j_*}^* > (j_* - j)c(1 - 2c)$  are satisfied.
- (2)  $C_p$  is consistent if assumptions A2 and A4 are satisfied.

Theorems 3.1 and 3.2 are helpful for the selection of tuning parameters. For consistency of  $IC_\nu$ , it is necessary that  $b_1(c) = -c^{-1} \ln(1 - c) < \nu$ . Similarly, for consistency of  $C_{p,\nu}$ , it is necessary that  $b_2(c) = (1 - c)^{-1} < \nu$ . The boundary values for  $c = 0.1, 0.2, \dots, 0.9$  are given as in Table 3.1.

The boundary values of  $b_1(c)$  are relatively steady for a change of  $c$ . However, those of  $b_2(c)$  become large as  $c$  is increased beyond 0.5.

Our sufficient conditions for consistency are derived as follows. Let  $T_j$  be a general criterion for  $M_j$ ,  $j \in \mathcal{F}$ . Then, we attempt to show that

$$\forall j \neq j_* \in \mathcal{F}, \quad \frac{1}{h_{j,j_*}}(T_j - T_{j_*}) \geq D_{j,j_*} \xrightarrow{p} \alpha_{j,j_*} > 0 \quad (19)$$

where  $D_{j,j_*}$  is some quantity, and  $h_{j,j_*}$  is some positive constant depending on the models. It is easy to see that (19) implies  $P(\hat{j}_T = j_*) \rightarrow 1$ . For example, under assumptions A1, A2 and A3, we shall show in the Appendix that for  $j > j_*$ ;

$$\frac{1}{n} \{IC_{\nu;j} - IC_{\nu;j_*}\} \xrightarrow{p} (j - j_*)\{\ln(1 - c) + \nu c\}$$

for  $j < j_*$ ;

$$\begin{aligned} \frac{1}{n} \{IC_{\nu;j} - IC_{\nu;j_*}\} &\xrightarrow{p} \ln(1 + \delta_{j+1}^*) \cdots (1 + \delta_{j_*}^*) - (j_* - j)\{\ln(1 - c) + \nu c\} \\ &\geq (j_* - j)[\ln(1 + \delta_{j_*}^*) - \{\ln(1 - c) + \nu c\}]. \end{aligned}$$

Therefore, by obtaining the ranges of  $\nu$  such that the above right-hand sides are positive, we can obtain Theorem 3.1(1). Generally, we can say that for the consistency, it is necessary that the penalty term tends to infinity. When  $\nu = \ln n$ , we have that for  $j < j_*$ ,

$$\frac{1}{n \ln n} (B_j - B_{j_*}) \xrightarrow{p} -(j_* - j)c,$$

which implies  $\Pr(B_j > B_{j_*}) \rightarrow 0$  for  $c > 0$ . Therefore,

$$\begin{aligned}\Pr(\hat{J}_B = j_*) &= \Pr(B_j > B_{j_*}, \text{ for all } j \neq j_*) \\ &\leq \Pr(B_j > B_{j_*}) \rightarrow 0, \quad \text{for } j < j_*,\end{aligned}$$

which implies [Corollary 3.2\(1\)](#).

We note that the consistency properties in [Theorems 3.1](#) and [3.2](#) hold with slight modifications for the case where the set of candidate ranks is a subfamily  $\mathcal{G} \subset \mathcal{F}$ . In fact, for example, the estimation method based on AIC can be expressed as

$$\hat{J}_{A;\mathcal{G}} = \arg \min_{j \in \mathcal{G}} A_j.$$

Then, the consistency of  $\hat{J}_{A;\mathcal{G}}$  is given as [Theorem 3.1](#) with the following modifications: add “ $j_* \in \mathcal{G}$ ” to A1, and replace “any  $i$  ( $1 \leq i \leq j_*$ )” in A2 and A3 by “any  $i$  ( $1 \leq i \leq j_*$ )  $\in \mathcal{G}$ ”.

### 3.2. Differences between high-dimensional and large-sample results

From the proofs of [Theorems 3.1](#) and [3.2](#), we can see that the AIC and  $C_p$  criteria on the dimensionality in a multivariate linear model satisfy the following:

$$(i) \text{ if } c \in [0, c_a), \quad \lim_{p/n \rightarrow c} P(\hat{J}_A \in \mathcal{F}_+ \setminus \{j_*\}) = 0, \quad (20)$$

$$(ii) \text{ if } c \in [0, 0.5), \quad \lim_{p/n \rightarrow c} P(\hat{J}_C \in \mathcal{F}_+ \setminus \{j_*\}) = 0, \quad (21)$$

under  $\Omega = O(n)$  and  $\Omega = O(np)$ , where  $c_a$  is the constant given [Theorem 3.1](#). Note that the properties (20) and (21) are different from those under a large-sample framework. In fact, under the large-sample framework (2) and assumption A4, it is known (Fujikoshi [9]) that

$$\lim_{n \rightarrow \infty} \Pr(\hat{J}_A = j) = \Pr(\hat{J}_C = j) = h(j|j_0), \quad (22)$$

where for  $j = 0, \dots, j_* - 1$ ,  $h(j|j_*) = 0$ , and for  $j = j_*, \dots, q$ , the  $h(j|j_*)$ 's are positive and are expressed in terms of the characteristic roots  $z_1 > \dots > z_s$  of a  $s \times s$  Wishart matrix  $\mathbf{W}$  distributed as  $\mathcal{N}_s(t, \mathbf{I}_s)$  as

$$h(j|j_*) = \Pr \left( \sum_{i=k+1}^{j-j_0} z_i > 2p_{k-j_*}; k = 0, \dots, j-j_*-1, \text{ and, } \sum_{i=j-j_*+1}^k z_i > 2p_{j-j_*,k}; k = j-j_*+1, \dots, q \right). \quad (23)$$

Here,  $p_{ij} = (p - j_* - i)(q - j_* - i) - (p - j_* - j)(q - j_* - j)$ , and the density function of  $z_1, \dots, z_s$  is expressed as

$$\frac{\pi^{s^2/2}}{2^{st/2} \Gamma_s(\frac{1}{2}s) \Gamma_s(\frac{1}{2}t)} \exp \left( -\frac{1}{2} \sum_{i=1}^s z_i \right) \prod_{i=1}^s z_i^{(t-s-1)/2} \prod_{i < j}^s (z_i - z_j),$$

where  $\Gamma_a(b) = \pi^{a(a-1)/4} \prod_{i=1}^a \Gamma[b - (i-1)/2]$ .

For BIC, we have seen that under a high-dimensional asymptotic framework it is not consistent under  $\Omega = O(n)$ , but it is consistent under  $\Omega = O(np)$ . Under a large-sample framework, it is consistent when  $\Omega = O(n)$ .

### 3.3. Numerical study

In this subsection, we numerically examine the validity of our claims and tendencies for the ranks estimated by the criteria A, B, C, or equivalently the criteria AIC, BIC,  $C_p$ , and their extensions. Our numerical results are given for the rank or dimensionality estimation criteria in discriminant analysis with  $q+1$  groups based on the total sample size  $n$  of  $p$  response variables. Assume that  $p \geq q$ . The criteria are based on the nonzero characteristic roots  $\ell_1 > \dots > \ell_q$  of  $\mathbf{S}_b \mathbf{S}_w^{-1}$ . Without loss of generality, we may assume that  $\mathbf{S}_w$  and  $\mathbf{S}_b$  are independently distributed as  $\mathcal{N}_p(n-q-1, \mathbf{I}_p)$  and  $\mathcal{N}_p(q, \mathbf{I}_p; \Omega_p)$ , respectively. Here,  $\Omega_p = \text{diag}(\omega_1, \dots, \omega_q, 0, \dots, 0)$  and  $\omega_1, \dots, \omega_q$  are the possible nonzero characteristic roots of the noncentrality matrix  $\Omega$  defined by (9). Further, the sample roots  $\ell_1 > \dots > \ell_q$  may be regarded (see [Lemma A.1](#)) as those of  $\mathbf{B}\mathbf{W}^{-1}$ , where  $\mathbf{W}$  and  $\mathbf{B}$  are independently distributed as  $\mathcal{N}_q(n-p-1, \mathbf{I}_q)$  and  $\mathcal{N}_q(p, \mathbf{I}_q; \Omega)$ , respectively, and  $\Omega = \text{diag}(\omega_1, \dots, \omega_q)$ .

If we suppose that  $q = 5$ , then we have six candidate models:  $M_0, M_1, \dots, M_5$ . It is assumed that the minimum model including the true model is  $M_3$ , and so  $j_* = 3$ . The two types of characteristic roots  $\omega_i$ ,  $i = 1, \dots, 5$ , are defined as follows:

$$\begin{aligned}(a): \quad & \omega_1 = 2\omega_3, \quad \omega_2 = 1.5\omega_3, \quad \omega_3 = n, \quad \omega_4 = \omega_5 = 0, \\ (b): \quad & \omega_1 = 2\omega_3, \quad \omega_2 = 1.5\omega_3, \quad \omega_3 = np, \quad \omega_4 = \omega_5 = 0.\end{aligned}$$



**Table 3.2**Selection probabilities of the true rank by A, B and C when  $p/n = 1/6$ .

$(n, p)$	Case (a)			Case (b)		
	A	B	C	A	B	C
(30, 5)	0.76	0.82	0.72	0.76	0.95	0.70
(120, 20)	0.95	0.65	0.86	0.93	1.00	0.83
(300, 50)	1.00	0.13	0.97	1.00	1.00	0.98
(480, 80)	1.00	0.00	0.99	1.00	1.00	1.00
(600, 100)	1.00	0.00	1.00	1.00	1.00	1.00

**Table 3.3**Selection probabilities of the true rank by A, B and C when  $n = 100$ .

$p$	Case (a)			Case (b)		
	A	B	C	A	B	C
10	0.87	1.00	0.82	0.87	1.00	0.82
30	0.95	0.00	0.71	0.94	1.00	0.66
50	0.90	0.00	0.21	0.83	1.00	0.16
70	0.52	0.00	0.00	0.40	1.00	0.00
90	0.00	0.01	0.00	0.00	0.98	0.00

These (a) and (b) correspond to the noncentrality matrix-1 and -2 under assumptions A3 and A4. Several different values of  $n$  and  $p = cn$  were prepared for Monte Carlo simulations with  $10^4$  repetitions. Table 3.2 shows simulation results for

$$(n, p) = (30, 5), (120, 20), (300, 50), (480, 80), (600, 100).$$

In these cases, the values of  $p/n$  are all  $1/6$ , and assumptions A3 and A4 are satisfied. We use the selection probabilities as the relative frequencies. From Table 3.2, we can identify the following tendencies.

- For case (a), the selection probabilities of the true rank by A and C are increasing when  $(n, p)$  with  $p/n = 1/6$  is increasing, and tend to 1.
- For case (a), the selection probabilities of the true rank by B do not increase even when  $(n, p)$  with  $p/n = 1/6$  is increasing
- For case (b), the selection probabilities of the true rank by A, B and C are increasing when  $(n, p)$  with  $p/n = 1/6$  is increasing, and tend to 1.

Next we examined the selection probabilities when  $n$  is fixed and  $p$  increases as follows:

$$n = 100, \quad p = 10, 30, 50, 70, 90.$$

In this case, the  $p/n$  values are 0.1, 0.3, 0.5, 0.7, 0.9. Table 3.3 gives the selection probabilities of the true model by A, B and C.

From Table 3.3, we can identify the following tendencies in the selection probabilities:

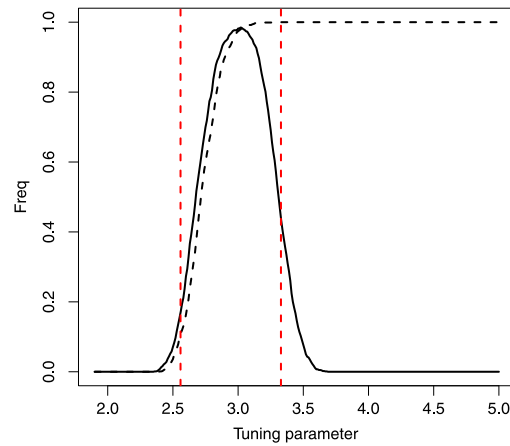
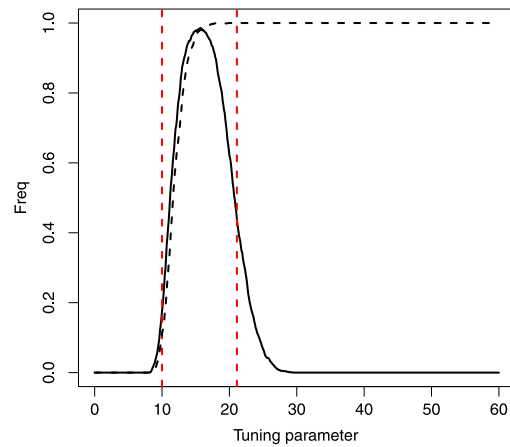
- The selection probabilities of the true rank by A for both case (a) and case (b) are increasing for  $10 \leq p \leq 30$  and decreasing for  $30 \leq p \leq 90$ , taking the maximum at  $p = 30$ .
- The selection probabilities of the true rank by C for both case (a) and case (b) take the maximum at  $p = 10$  and then are decreasing.
- For case (a) and case (b), the selection probabilities of the true rank by A and C are near 0 for  $p > 0.797 (\approx c_a)$  and  $p > 0.5$ , respectively.
- For case (a), the selection probabilities of the true rank by B are 1 when  $p = 10$ , but zero when  $30 \leq p \leq 90$ .
- For case (b), the selection probabilities of the true rank by B are 1 for all  $p (10 \leq p \leq 90)$ .

Next we examined a range of tuning parameters  $\nu$  such that  $IC_\nu$  and  $C_{p,\nu}$  are consistent. The experiment was performed for  $p/n = 0.9$ ,  $n = 1000$  and  $p = 900$ . Note that in this case we cannot expect consistency of the usual AIC and  $C_p$ . Further, the usual BIC is not consistent for case (a).

The numerical results are given in Figs. 3.1 and 3.2 whose horizontal axis and vertical axis show the values of  $\nu$  and the selection probabilities of the true rank, respectively. The solid lines and the dotted lines correspond to case (a) and case (b). In Fig. 3.1, the left dotted vertical line denotes  $\nu = -\frac{1}{c} \ln(1 - c)$ , and the right dotted vertical line denotes  $-\frac{1}{c} \ln(1 - c) + \frac{1}{c} \ln(1 + \delta_{j_*}^*)$ . In Fig. 3.2, the left dotted vertical line denotes  $\nu = \frac{1}{1-c}$ , and the right dotted vertical line denotes  $\frac{1}{1-c} + \frac{1}{c(1-c)} \delta_{j_*}^*$ . From Figs. 3.1 and 3.2, we can identify the following tendencies:

- For case (a),  $IC_\nu$  shall be consistent when  $-c^{-1} \ln(1 - c) < \nu < -c^{-1} \ln(1 - c) + \ln(1 + \delta_{j_*}^*)$ .
- For case (b),  $IC_\nu$  shall be consistent when  $-c^{-1} \ln(1 - c) < \nu$ .
- For case (a),  $C_{p,\nu}$  shall be consistent when  $(1 - c)^{-1} < \nu < (1 - c)^{-1} + c(1 - c)^{-1} \delta_{j_*}^*$ .
- For case (b),  $C_{p,\nu}$  shall be consistent when  $(1 - c)^{-1} < \nu$ .



Fig. 3.1. Selection probabilities by  $IC_v$ .Fig. 3.2. Selection probabilities by  $C_{p,v}$ .

## 4. Ridge-type criteria and their properties

### 4.1. Ridge-type criteria

When  $p > n - k$ ,  $\mathbf{S}_e$  becomes singular, and so we cannot use the criteria AIC, BIC and  $C_p$ . One way to overcome this problem is to use the ridge-type estimator of  $\Sigma$  defined by

$$\hat{\Sigma}_\lambda = \frac{1}{n}(\mathbf{S}_e + \lambda \mathbf{I}_p) = \frac{1}{n}\mathbf{S}_{e,\lambda}, \quad (24)$$

or its modifications. Here,  $\lambda = \{1/(np)\}\text{tr}\mathbf{S}_e$ . For a discussion on the use of  $\lambda$ , see Kubokawa and Srivastava [17]. Let  $\ell_{\lambda,1} > \dots > \ell_{\lambda,q}$  be the non-zero characteristic roots of  $\mathbf{S}_h\mathbf{S}_{e,\lambda}^{-1}$ . Then, we propose the following modifications of A, B and C:

$$\begin{aligned} A_{\lambda,j} &= n \ln \prod_{i=j+1}^q (1 + \ell_{\lambda,i}) - 2(p-j)(q-j), \quad j = 0, \dots, q, \\ B_{\lambda,j} &= n \ln \prod_{i=j+1}^q (1 + \ell_{\lambda,i}) - (\ln n)(p-j)(q-j), \quad j = 0, \dots, q, \\ C_{\lambda,j} &= n \sum_{i=j+1}^q \ell_{\lambda,i} - 2(p-j)(q-j), \quad j = 0, \dots, q. \end{aligned} \quad (25)$$

Here,  $A_{\lambda,q} = 0$ ,  $B_{\lambda,q} = 0$  and  $C_{\lambda,q} = 0$ . The criteria  $A_\lambda$ ,  $B_\lambda$  and  $C_\lambda$  are obtained from the criteria A, B and C by substituting  $\ell_{\lambda,i}$  to  $\ell_i$ .

The criteria  $A_\lambda$  and  $C_\lambda$  can be justified by deriving asymptotic unbiased estimators of the AIC-type or the  $C_p$ -type risks based on ridge-type estimators. In the following, we shall give an outline of the derivations. Let  $f(\mathbf{Y}; \boldsymbol{\Theta}, \boldsymbol{\Sigma})$  be the density function of  $\mathbf{Y}$ , and let  $\hat{\boldsymbol{\Theta}}$  be the maximum likelihood estimator of  $\boldsymbol{\Theta}$  when  $n - k > p$ . We consider the ridge-type estimator  $\hat{\boldsymbol{\Theta}}_\lambda$  defined from  $\hat{\boldsymbol{\Theta}}$  by substituting  $\mathbf{S}_{e,\lambda}$  for  $\mathbf{S}_e$ . The AIC-type risk of a candidate model  $M_j$  based on the ridge-type estimator is

$$R_{A,\lambda} = E_{\mathbf{Y}}^* E_{\mathbf{Z}}^* \left\{ -2 \ln f(\mathbf{Z}; \hat{\boldsymbol{\Theta}}_\lambda, \hat{\boldsymbol{\Sigma}}_\lambda) \right\}, \quad (26)$$

which is based on Kullback–Leibler information. Here,  $\mathbf{Z}; n \times p$  is a random matrix that has the same distribution as  $\mathbf{Y}$  and is independent of  $\mathbf{Y}$ , and  $E^*$  denotes the expectation of the true model  $M_*$ . The random matrix  $\mathbf{Z}$  may be regarded as a future observation matrix of  $\mathbf{Y}$ . When we estimate  $R_{A,\lambda}$  by

$$-2 \ln f(\mathbf{Y}; \hat{\boldsymbol{\Theta}}_\lambda, \hat{\boldsymbol{\Sigma}}_\lambda) = n \ln \prod_{i=j+1}^q (1 + \ell_{\lambda,i}) + n \ln |(1/n)\mathbf{S}_{e,\lambda}| + np\{1 + \ln(2\pi)\}, \quad (27)$$

where  $\ell_{\lambda,1} > \dots > \ell_{\lambda,q} > 0$  are the non-zero characteristic roots of  $\mathbf{S}_h \mathbf{S}_{e,\lambda}^{-1}$ . Then, the bias term is expressed as  $-b_{A,\lambda}$ , where

$$b_{A,\lambda} = E_{\mathbf{Y}}^* E_{\mathbf{Z}}^* \left\{ \text{tr} \hat{\boldsymbol{\Sigma}}_\lambda^{-1} (\mathbf{Z} - \mathbf{X}\boldsymbol{\Theta})^\top (\mathbf{Z} - \mathbf{X}\boldsymbol{\Theta}) \right\} - np. \quad (28)$$

Here, for a notational simplicity, we express the true parameters as those in (4). Under a large-sample asymptotic framework and  $\tilde{\boldsymbol{\Omega}} = O(n)$ , the expectation in (28) can be evaluated by the same method as in Fujikoshi and Veitch [13] as follows.

$$b_{A,\lambda} = 2 \left\{ j(p + q - j) + (k - q)p + \frac{1}{2}p(p + 1) - \frac{1}{2p} \text{tr} \boldsymbol{\Sigma} \right\} + o(1). \quad (29)$$

This suggests that  $A_\lambda$  in (25) is an asymptotic unbiased estimator for  $R_{A,\lambda}$ . Estimating  $\boldsymbol{\Sigma}$  by a ridge-type estimator  $\hat{\boldsymbol{\Sigma}}_\lambda = (1/n)\mathbf{S}_{e,\lambda}$ , we can justify  $A_\lambda$  as an estimator of  $R_{A,\lambda}$ . The justification  $C_\lambda$  can be obtained by considering a ridge-type  $C_p$  risk defined by

$$R_{C,\lambda} = E_{\mathbf{Y}}^* E_{\mathbf{Z}}^* \left\{ \text{tr} \hat{\boldsymbol{\Sigma}}_\lambda^{-1} (\mathbf{Z} - \mathbf{X}\hat{\boldsymbol{\Theta}}_\lambda)^\top (\mathbf{Z} - \mathbf{X}\hat{\boldsymbol{\Theta}}_\lambda) \right\}, \quad (30)$$

and by deriving its asymptotic unbiased estimator under a large-sample framework and  $\tilde{\boldsymbol{\Omega}} = O(n)$ .

#### 4.2. Consistency of ridge-type criteria

In this subsection, we examine the consistency of ridge-type criteria when  $n - k > p$ . More precisely, it is shown that the criteria  $A_\lambda$ ,  $B_\lambda$  and  $C_\lambda$  have the same consistency properties as the criteria A, B and C, respectively, under an additional assumption. The results are stated as follows:

**Theorem 4.1.** Suppose that assumption A1 is satisfied,  $n - k > p$  and  $(1/p)\text{tr} \boldsymbol{\Sigma} \rightarrow \alpha_0$ . Then, we have the following results.

- (1)  $A_\lambda$  is consistent if  $c \in [0, c_a)$  and assumptions A2 and A3 and the inequality  $\ln(1 + \delta_{j_*}^*) > (j_* - j)\{2c + \ln(1 - c)\}$  are satisfied.
- (2)  $A_\lambda$  is consistent if  $c \in [0, c_a)$  and assumptions A2 and A4 are satisfied.
- (3)  $B_\lambda$  is not consistent if assumptions A2 and A3 are satisfied.
- (4)  $B_\lambda$  is consistent if assumptions A2 and A4 are satisfied.
- (5)  $C_\lambda$  is consistent if  $c \in [0, 0.5)$  and assumptions A2 and A3 and the inequality  $\delta_{j_*}^* > (j_* - j)c(1 - 2c)$  are satisfied.
- (6)  $C_\lambda$  is consistent if  $c \in [0, 0.5)$  and assumptions A2 and A4 are satisfied.

Theorem 4.1 is shown by noting that the limiting values of  $\ell_{\lambda,i}$ ,  $i = 1, \dots, q$  are the same as those of  $\ell_i$ ,  $i = 1, \dots, q$ . For the proof, see Lemma A.2 in the Appendix. It is possible to generalize Theorem 4.1 for a generalized criterion with a tuning parameter. The consistency properties of ridge-type criteria for  $p > n - k$  are studied numerically in the next section.

#### 4.3. Numerical study

##### The case of $n - q - 1 \geq p$

In this section, we consider the selection probabilities of ridge-type criteria  $A_\lambda$ ,  $B_\lambda$  and  $C_\lambda$  for case (a) and case (b) considered in Section 3.2. First, simulations were performed for  $\boldsymbol{\Sigma} = (0.8^{|i-j|})$ ,  $n = 100$  and  $p = 10, 20, \dots, 90$ . The results are given in Figs. 4.1 and 4.2. The horizontal axis and the vertical axis show the values of  $p$  and the selection probabilities of the true model, respectively. Here,  $(\alpha)$ ,  $(\beta)$  and  $(\gamma)$  denote the probabilities of selecting the true rank by  $(A_\lambda, A)$ ,  $(B_\lambda, B)$  and  $(C_\lambda, C)$ , respectively. The selection probabilities by A, B and C are denoted by solid lines. Their ridge-type ones are denoted

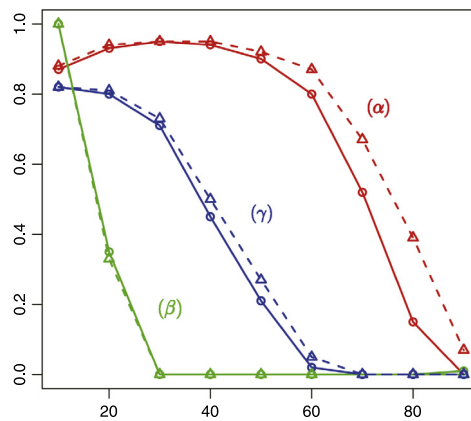


Fig. 4.1. Selection probabilities of the true rank by  $(A_\lambda, A)$ ,  $(B_\lambda, B)$  and  $(C_\lambda, C)$  for case (a).

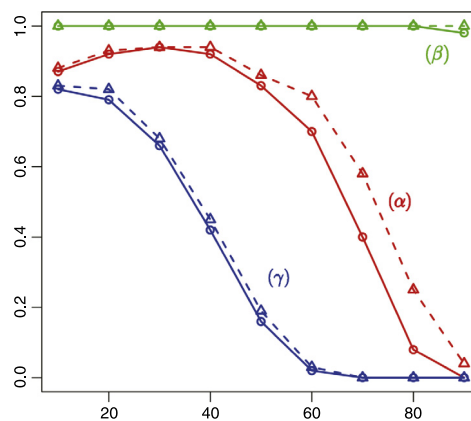


Fig. 4.2. Selection probabilities of the true rank by  $(A_\lambda, A)$ ,  $(B_\lambda, B)$  and  $(C_\lambda, C)$  for case (b).

by dotted lines. Based on Figs. 4.1 and 4.2, it is clear that when  $n - q - 1 \geq p$ , the selection probabilities of the true model by  $A_\lambda$ ,  $B_\lambda$  and  $C_\lambda$  are similar to the ones by their ridge-type criteria A, B and C, respectively.

#### The case of $n - q - 1 < p$

Next, in order to examine behaviors of ridge-type criteria  $A_\lambda$ ,  $B_\lambda$  and  $C_\lambda$  when  $n - q - 1 < p$ , we did simulation experiments for  $\Sigma = (0.8^{|i-j|})$ ,  $n = 30, 50, 100$  and  $c = p/n = 1, 1.5, 2, 2.5, 3$  for case (b). The numerical results are given in Table 4.1, where the values for “Under”, “True” and “Over” denote the probabilities of selecting the underspecified models, the true model and the overspecified models, respectively.

From Table 4.1, the following tendencies are clear for  $1 \leq c = p/n \leq 3$  and  $30 \leq n \leq 100$ .

- For each  $n$ , the selection probabilities of the true model by  $A_\lambda$  become large as  $c$  becomes large. The probabilities are 1 when  $c \geq 2$  and  $n = 30, 50$ , but for  $n = 100$ , the probabilities are 1 when  $c \geq 2.5$ .
- The selection probabilities of the true model by  $B_\lambda$  are 1 when  $c \leq 1.5$ . However, for  $c \geq 2$ ,  $B_\lambda$  underestimates the true model when  $p$  is near  $n$ , and chooses underspecified models when  $n$  is large.
- The criterion  $C_\lambda$  chooses overspecified ranks for all most cases. The cause seems to be that the smallest roots  $\ell_{\lambda,4}$  and  $\ell_{\lambda,5}$  are relatively large, and this gives large effects to  $C_p, \lambda$  in the comparison with  $A_\lambda$ .

## 5. Concluding remarks

In general, it is known that under the large-sample asymptotic framework (2), AIC and  $C_p$  are not consistent, but BIC is consistent. However, in this paper, we demonstrated that the criteria based on AIC and  $C_p$  for estimating the rank (dimensionality) in a multivariate linear model can have consistency, under the high-dimensional asymptotic framework (3). For consistency, some additional assumptions are required. For AIC, it is necessary that  $c \in [0, c_a)$ , where  $c_a \approx 0.797$ . For  $C_p$ , it is necessary that  $c \in [0, 0.5)$ . More precisely, consistency was considered under two types of assumptions on the largeness of the characteristic roots of the noncentrality matrix  $\Omega$  in (17). For the criterion based on BIC, we note that it is consistent when  $\Omega = O(np)$ , but it is not consistent when  $\Omega = O(n)$ . These results were extended for the general criteria  $IC_v$  and  $C_{p,v}$  with a tuning parameter in (15) and (16). We gave sufficient conditions for  $IC_v$  and  $C_{p,v}$  to be consistent. The

**Table 4.1**Selection probabilities by  $A_\lambda$ ,  $B_\lambda$  and  $C_\lambda$  when  $c = p/n > 1$ .

$c$	$n = 30$			$n = 50$			$n = 100$		
	Under	True	Over	Under	True	Over	Under	True	Over
Selection probabilities by $A_\lambda$									
1	0.00	0.81	0.19	0.00	0.79	0.21	0.00	0.55	0.45
1.5	0.00	0.93	0.07	0.00	0.87	0.13	0.00	0.24	0.76
2	0.00	1.00	0.00	0.00	1.00	0.00	0.00	0.89	0.11
2.5	0.00	1.00	0.00	0.00	1.00	0.00	0.00	1.00	0.00
3	0.00	1.00	0.00	0.00	1.00	0.00	0.00	1.00	0.00
Selection probabilities by $B_\lambda$									
1	0.00	1.00	0.00	0.00	1.00	0.00	0.00	1.00	0.00
1.5	0.13	0.87	0.00	0.06	0.94	0.00	0.00	1.00	0.00
2	1.00	0.00	0.00	1.00	0.00	0.00	1.00	0.00	0.00
2.5	1.00	0.00	0.00	1.00	0.00	0.00	1.00	0.00	0.00
3	1.00	0.00	0.00	1.00	0.00	0.00	1.00	0.00	0.00
Selection probabilities by $C_\lambda$									
1	0.00	0.01	0.99	0.00	0.00	1.00	0.00	0.00	1.00
1.5	0.00	0.00	1.00	0.00	0.00	1.00	0.00	0.00	1.00
2	0.00	0.00	1.00	0.00	0.00	1.00	0.00	0.00	1.00
2.5	0.00	0.00	1.00	0.00	0.00	1.00	0.00	0.00	1.00
3	0.00	0.00	1.00	0.00	0.00	1.00	0.00	0.00	1.00

sufficient conditions are useful in the selection of the tuning parameter  $\nu$ . When  $q$  is large, we showed that the probability of selecting ranks  $j > j_*$  by AIC tends to 0. It is left as a future problem to examine consistency properties of AIC and the other criteria.

We proposed the ridge-type criteria  $A_\lambda$ ,  $B_\lambda$  and  $C_\lambda$  in (25) which are also defined for the case where  $p > n - k$ . It was shown that  $A_\lambda$  and  $C_\lambda$  are asymptotic unbiased estimators of AIC-type and  $C_p$ -type, respectively, under a large-sample framework. Further, these ridge-type criteria have the same consistency properties as the criteria A, B and C, respectively, when  $n - k > p$ . For the case  $n - k < p$ , we studied the consistency properties numerically in the range  $1 \leq p/n \leq 3$ . It was found that  $A_\lambda$  becomes consistent as  $c = p/n$  becomes large,  $B_\lambda$  is consistent if  $n$  is close to  $p$ , and  $C_{p,\lambda}$  estimates overspecified ranks. The theoretical justifications for these observations are left as future research.

Recently, sparse estimation methods under a reduced rank based on penalization techniques have been proposed by Yuan et al. [23] and Bunea et al. [6]. Chen and Hung [8] and Bunea et al. [7] also considered simultaneous methods for dimension reduction and variable selection. These methods are based on least squares methods, and so it seems that the correlations of the response variables have been ignored. In our approach, the correlations of the response variables have been taken into consideration, but in most case, it is assumed that  $n - k > p$ , except for the ridge-type criteria. In both the penalized methods and the generalized criteria  $A_\nu$ ,  $B_\nu$  and  $C_{p,\nu}$ , there are problems in how to decide the tuning parameter.

## Acknowledgments

We wish to thank the editor and two reviewers for their critical and helpful comments on the first version, which allowed us to greatly improve our manuscript. The first author's research is partially supported by a Grant-in-Aid for Scientific Research (C), #25330038, 2013–2015, from the Ministry of Education, Science, Sports and Culture, Japan.

## Appendix. Proofs of theorems and lemmas

### Lemma A.1

First, we prepare a lemma on the limiting behavior of the characteristic roots of  $\mathbf{S}_h \mathbf{S}_e^{-1}$  in a high-dimensional case.

**Lemma A.1.** Let  $\mathbf{S}_e$  and  $\mathbf{S}_h$  be independently distributed as a Wishart distribution  $\mathcal{N}_p(n - k, \Sigma)$  and a noncentral Wishart distribution  $\mathcal{N}_p(q, \Sigma; \Sigma^{1/2} \tilde{\Omega} \Sigma^{1/2})$ , respectively. Here, it is assumed that  $n - k \geq p$ . Let  $\ell_1 > \dots > \ell_q$  and  $\omega_1 \geq \dots \geq \omega_q$  be the possible nonzero characteristic roots of  $\mathbf{S}_h \mathbf{S}_e^{-1}$  and  $\tilde{\Omega}$ , respectively. We assume that  $\text{rank}(\tilde{\Omega}) = a$ , and hence  $\omega_1 \geq \dots \geq \omega_a > \omega_{a+1} = \dots = \omega_q = 0$ . Further, we assume that the multiplicities of the  $\omega_i$ 's do not depend on  $p$  and  $n$ . For the limiting behavior of  $\ell_1 > \dots > \ell_q$  under a high-dimensional asymptotic framework

$$p \rightarrow \infty, \quad n \rightarrow \infty, \quad k \rightarrow \infty, \quad p/n \rightarrow c \in [0, 1), \quad k/n \rightarrow 0, \quad (\text{A.1})$$

we have the following results:

(1) Suppose that for any  $i$  ( $1 \leq i \leq a$ ),  $\omega_i = n\delta_i = O(n)$  and  $\lim_{p/n \rightarrow c} \delta_i = \delta_i^* > 0$ . Then

$$\ell_i \xrightarrow{p} \frac{c}{1-c} + \frac{1}{1-c} \delta_i^*, \quad i = 1, \dots, a, \quad \text{and} \quad \ell_i \xrightarrow{p} \frac{c}{1-c}, \quad i = a+1, \dots, q.$$

(2) Suppose that for any  $i$  ( $1 \leq i \leq a$ ),  $\omega_j = np\xi_i = O(np)$  and  $\lim_{p/n \rightarrow c} \xi_i = \xi_i^* > 0$ . Then

$$\frac{1}{p}\ell_i \xrightarrow{p} \frac{1}{1-c}\xi_i^*, \quad i = 1, \dots, a, \quad \text{and} \quad \ell_i \xrightarrow{p} \frac{c}{1-c}, \quad i = a+1, \dots, q.$$

It is known (see Fujikoshi, Ulyanov and Shimizu [12]) that the nonzero characteristic roots  $\ell_1 > \dots > \ell_q > 0$  of  $\mathbf{S}_h \mathbf{S}_e^{-1}$  may be regarded as those of  $\mathbf{B}\mathbf{W}^{-1}$ , where  $\mathbf{W}$  and  $\mathbf{B}$  are independently distributed as a central Wishart distribution  $\mathcal{N}_q(m, \mathbf{I}_q)$  and a noncentral Wishart distribution  $\mathcal{N}_q(p, \mathbf{I}_q; \mathbf{\Omega})$ , respectively. Here,  $m = n - k - p + q$ ,  $\mathbf{\Omega} = \text{diag}(\omega_1, \dots, \omega_q)$ , and

$$\omega_1 \geq \dots \geq \omega_a > \omega_{a+1} = \dots = \omega_q = 0.$$

Using this property, Fujikoshi [10] has derived limiting distributions of the  $\ell_i$ 's under a general multiplicity of the  $\delta_i$ 's and  $\xi_i$ 's. Lemma A.1 follows from Fujikoshi [10].

**The Proof of Theorem 3.1.** In the proof of Theorem 3.1, it is assumed that the true rank or dimensionality is  $j_*$ . Since the number of possible models is finite, it is sufficient to show that the values of  $\text{IC}_{v,j} - \text{IC}_{v,j_*}$  converge to positive values.

Note that for  $j > j_*$ ,

$$\text{IC}_{v,j} - \text{IC}_{v,j_*} = -n \ln\{(1 + \ell_{j_*+1}) \cdots (1 + \ell_j)\} + v(j - j_*)(p + q - j - j_*),$$

and for  $j < j_*$

$$\text{IC}_{v,j} - \text{IC}_{v,j_*} = n \ln\{(1 + \ell_{j+1}) \cdots (1 + \ell_{j_*})\} + v(j - j_*)(p + q - j - j_*).$$

Suppose that  $\mathbf{\Omega} = O(n)$ . Then, using Lemma A.1(1), we have that for  $j > j_*$ ,

$$\frac{1}{n} (\text{IC}_{v,j} - \text{IC}_{v,j_*}) \xrightarrow{p} (j - j_*)\{\ln(1 - c) + vc\}.$$

The limiting value is positive when  $-c^{-1} \ln(1 - c) < v$ . Next suppose that  $j < j_*$ . Then, using Lemma A.1(1), we have

$$\begin{aligned} \frac{1}{n} (\text{IC}_{v,j} - \text{IC}_{v,j_*}) &\xrightarrow{p} \ln(1 + \delta_{j+1}^*) \cdots (1 + \delta_{j_*}^*) - (j_* - j)\{\ln(1 - c) + vc\} \\ &\geq \ln(1 + \delta_{j_*}^*) \cdots (1 + \delta_{j_*}^*) - (j_* - j)\{\ln(1 - c) + vc\} \\ &= (j_* - j)[\ln(1 + \delta_{j_*}^*) - \{\ln(1 - c) + vc\}]. \end{aligned}$$

The limiting value is positive when

$$v < -\frac{1}{c} \ln(1 - c) + \frac{1}{c} \ln(1 + \delta_{j_*}^*).$$

For the proof for case  $c = 0$ , from the above result, it holds that for  $0 < c < 1$ ,

$$\lim \frac{1}{p} (\text{IC}_{v,j} - \text{IC}_{v,j_*}) \geq (j_* - j) \left[ \frac{1}{c} \ln(1 + \delta_{j_*}^*) - \left\{ \frac{1}{c} \ln(1 - c) + v \right\} \right].$$

The result is obtained by considering  $\lim_{c \rightarrow 0}$ .

Now we shall prove the result (2). For  $j > j_*$ , the limiting behavior of  $\ell_j$  under  $\omega_j = O(np)$  is the same as that under  $\omega_j = O(n)$ . Therefore, the limiting value of  $(1/n) \{\text{IC}_{v,j} - \text{IC}_{v,j_*}\}$  is positive when  $-c^{-1} \ln(1 - c) < v$ . From Lemma A.1(2), we have

$$\frac{1}{np} (\text{IC}_{v,j} - \text{IC}_{v,j_*}) \xrightarrow{p} j_* - j.$$

This proves Theorem 3.1(2).

**The Proof of Theorem 3.2.** Note that for  $j > j_*$ ,

$$C_{p,v,j} - C_{p,v,j_*} = -n(\ell_{j_*+1} + \dots + \ell_j) + v(j - j_*)(p + q - j - j_*),$$

and for  $j < j_*$ ,

$$C_{p,v,j} - C_{p,v,j_*} = n(\ell_{j+1} + \dots + \ell_{j_*}) + v(j - j_*)(p + q - j - j_*).$$

Therefore, Theorem 3.2 can be proved in the same way as was Theorem 3.1, with the help of Lemma A.1.

**Lemma A.2.** Let  $\mathbf{S}_e$  and  $\mathbf{S}_h$  be independently distributed as a Wishart distribution  $\mathcal{N}_p(n - k, \mathbf{I}_p)$  and a noncentral Wishart distribution  $\mathcal{N}_p(q, \mathbf{I}_p; \mathbf{\Omega})$ , respectively. Here, we assume that  $n - k \geq p \geq q$  and  $\mathbf{\Omega} = \text{diag}(\omega_1, \dots, \omega_p)$ . Let  $\mathbf{S}_{e,\lambda} = \mathbf{S}_e + \lambda \mathbf{I}_p$ , where  $\lambda = (np)^{-1} \text{tr} \mathbf{\Sigma} \mathbf{S}_e$ . Let  $\ell_1 > \dots > \ell_q$  and  $\ell_{\lambda,1} > \dots > \ell_{\lambda,q}$  be the possible nonzero characteristic roots of  $\mathbf{S}_h \mathbf{S}_e^{-1}$  and  $\mathbf{S}_h \mathbf{S}_{e,\lambda}^{-1}$ , respectively. We assume that (i) we are under the high-dimensional asymptotic framework given by (A.1), (ii) when  $p/n \rightarrow c \in [0, 1)$ ,  $(1/p) \text{tr} \mathbf{\Sigma} \rightarrow \alpha_0$ , and (iii)  $\omega_1 \geq \dots \geq \omega_a > \omega_{a+1} = \dots = 0$ . Then, the characteristic roots  $\ell_{\lambda,1} > \dots > \ell_{\lambda,q}$  have the same limiting values as those of the characteristic roots  $\ell_1 > \dots > \ell_q$  given in Lemma A.1.

In general, it holds that  $\ell_i \geq \ell_{\lambda,i}$ ,  $i = 1, \dots, q$ . Noting that

$$\begin{aligned} (\mathbf{S}_e + \lambda \mathbf{I}_p)^{-1} &= \mathbf{S}_e^{-1} (\mathbf{I}_p + \lambda \mathbf{S}_e^{-1})^{-1} \\ &= \mathbf{S}_e^{-1} - \lambda \mathbf{S}_e^{-2} (\mathbf{I}_p + \lambda \mathbf{S}_e^{-1})^{-1}, \end{aligned}$$

the following decomposition is obtained:

$$\mathbf{S}_h \mathbf{S}_e^{-1} = \mathbf{S}_h \mathbf{S}_e^{-1} (\mathbf{I}_p + \lambda \mathbf{S}_e^{-1})^{-1} + \lambda \mathbf{S}_h \mathbf{S}_e^{-2} (\mathbf{I}_p + \lambda \mathbf{S}_e^{-1})^{-1}.$$

Using Weyl's Theorem (see Seber [21, p. 117]), we have

$$\ell_i \leq \ell_{\lambda,i} + \lambda \operatorname{tr} \mathbf{S}_h \mathbf{S}_e^{-2} (\mathbf{I}_p + \lambda \mathbf{S}_e^{-1})^{-1} \leq \ell_{\lambda,i} + \lambda \operatorname{tr} (\mathbf{S}_h \mathbf{S}_e^{-1}) (\mathbf{S}_e^{-1}).$$

Note that  $2\operatorname{tr} \mathbf{A} \mathbf{B} \leq \operatorname{tr} \mathbf{A}^2 + \operatorname{tr} \mathbf{B}^2$  for any square matrices  $\mathbf{A}$  and  $\mathbf{B}$ . Therefore,

$$\begin{aligned} 2\operatorname{tr} (\mathbf{S}_h \mathbf{S}_e^{-1}) (\mathbf{S}_e^{-1}) &\leq \operatorname{tr} (n^{-\gamma/2} \mathbf{S}_h \mathbf{S}_e^{-1})^2 + \operatorname{tr} (n^{\gamma/2} \mathbf{S}_e^{-1})^2 \\ &= n^{-\gamma} (\ell_1^2 + \dots + \ell_q^2) + n^{\gamma} (n-k)^{-2} \operatorname{tr} \{ (n-k)^{-1} \mathbf{S}_e \}^{-2} \end{aligned}$$

where  $\gamma$  is a positive constant. By the Marčenko–Pastur law (see Bai and Silverstein [5]), it is known (Bai et al. [4]) that

$$\lim_n \frac{1}{n} \operatorname{tr} \{ (n-k)^{-1} \mathbf{S}_e \}^{-2} = \frac{c}{(1-c)^3}.$$

When  $\omega_j = n\delta_j = O(n)$  and  $\lim \delta_j = \delta_j^* > 0$ , from Lemma A.1 we have

$$\lim (\ell_1^2 + \dots + \ell_q^2) = \sum_{i=1}^a \left( \frac{c}{1-c} + \frac{1}{1-c} \delta_i^* \right)^2 + (q-a) \left( \frac{c}{1-c} \right)^2.$$

Noting that  $\lambda \rightarrow \alpha$ , and taking  $0 < \gamma < 1$ , we get  $\lim \ell_i \leq \lim \ell_{\lambda,i}$ , and hence  $\lim \ell_{\lambda,i} = \lim \ell_i$ . When  $\omega_j = np\xi_j = O(np)$  and  $\lim \xi_j = \lim \xi_j^*$ , we have seen that  $(1/p)\ell_j \rightarrow (1-c)^{-1}\xi_j^*$ ,  $j = 1, \dots, a$  and  $\ell_j \rightarrow c(1-c)^{-1}$ ,  $j = a+1, \dots, q$ . By a similar argument as that in the case of  $\omega_j = O(n)$ , we can show that  $\tilde{\ell}_j$  has the same limiting value as  $\ell_j$ , for  $j = 1, \dots, q$ .

## References

- [1] H. Akaike, Information theory and an extension of the maximum likelihood principle, in: B.N. Petrov, F. Csáki (Eds.), 2nd. International Symposium on Information Theory, Akadémiai Kiadó, Budapest, 1973, pp. 267–281.
- [2] T.W. Anderson, Estimating linear restrictions on regression coefficients for multivariate normal distributions, *Ann. Math. Statist.* 22 (1951) 327–351.
- [3] T.W. Anderson, Introduction to Multivariate Statistical Analysis, third ed., Wiley, Hoboken, NJ, 2003.
- [4] Z.W. Bai, K.P. Choi, Y. Fujikoshi, Limiting behavior of eigenvalues in MANOVA with high-dimension by RMT, 2016 (submitted for publication).
- [5] Z.W. Bai, J.W. Silverstein, Spectral Analysis of Large Dimensional Random Matrices, second ed., Springer, 2010.
- [6] F. Bunea, Y. She, M.H. Wegkamp, Optimal selection of reduced rank estimators of high-dimensional matrices, *Ann. Statist.* 39 (2011) 1282–1309.
- [7] F. Bunea, Y. She, M.H. Wegkamp, Joint variable and rank selection for parsimonious estimation of high-dimensional matrices, *Ann. Statist.* 40 (2012) 2359–2388.
- [8] L. Chen, J.Z. Huang, Sparse reduced-rank regression for simultaneous dimension reduction and variable selection, *J. Amer. Statist. Assoc.* 107 (2012) 1533–1545.
- [9] Y. Fujikoshi, Two methods for estimation of dimensionality in canonical correlation analysis and the multivariate linear model, in: K. Matsushita (Ed.), Statistical Theory and Data Analysis, Elsevier Science, Amsterdam, 1985, pp. 233–240.
- [10] Y. Fujikoshi, High-dimensional asymptotic distributions of characteristic roots in multivariate linear model and canonical correlation analysis, 2016 (submitted for publication).
- [11] Y. Fujikoshi, T. Sakurai, H. Yanagihara, Consistency of high-dimensional AIC-type and  $C_p$ -type criteria in multivariate linear regression, *J. Multivariate Anal.* 123 (2014) 184–200.
- [12] Y. Fujikoshi, V.V. Ulyanov, R. Shimizu, Multivariate Statistics: High-Dimensional and Large-Sample Approximations, Wiley, Hoboken, NJ, 2010.
- [13] Y. Fujikoshi, L.G. Veitch, Estimation of dimensionality in canonical correlation analysis, *Biometrika* 66 (1979) 345–351.
- [14] B.K. Gunderson, R.J. Muirhead, On estimating the dimensionality in canonical correlation analysis, *J. Multivariate Anal.* 62 (1997) 121–136.
- [15] A.J. Izenman, Reduced-Rank Regression for the multivariate linear model, *J. Multivariate Anal.* 5 (1975) 248–262.
- [16] A.M. Kshirsagar, Multivariate Analysis, Marcel Dekker, New York, 1972.
- [17] T. Kubokawa, M.S. Srivastava, Selection of variables in multivariate regression models for large dimensions, *Commun. Stat. - Theory Methods* 41 (2012) 2465–2489.
- [18] C.L. Mallows, Some comments on  $C_p$ , *Technometrics* 15 (1973) 661–675.
- [19] G.C. Reisel, R.P. Velu, Multivariate Reduced-Rank Regression, in: Lecture Notes in Statistics, vol. 136, Springer, New York, 1998.
- [20] G. Schwarz, Estimating the dimension of a model, *Ann. Statist.* 6 (1978) 461–464.
- [21] G.A.F. Seber, A Matrix Handbook for Statisticians, Wiley, Hoboken, NJ, 2008.
- [22] H. Yanagihara, H. Wakaki, Y. Fujikoshi, A consistency property of AIC for multivariate linear model when the dimension and the sample size are large, *Electron. J. Stat.* 9 (2015) 59–81.
- [23] M. Yuan, A. Ekici, Z. Lu, R. Monteiro, Dimension reduction and coefficient estimation in multivariate linear model, *J. R. Stat. Soc. B* 69 (2007) 329–346.