

# Consistency of Distance-based Criterion for Selection of Variables in High-dimensional Two-Group Discriminant Analysis

Tetsuro Sakurai\* and Yasunori Fujikoshi\*\*

*\*Center of General Education, Tokyo University of Science, Suwa,  
5000-1 Toyohira, Chino, Nagano 391-0292, Japan*

*\*\*Department of Mathematics, Graduate School of Science,*

*Hiroshima University, 1-3-1 Kagamiyama, Higashi Hiroshima, Hiroshima  
739-8626, Japan*

## Abstract

This paper is concerned with selection of variables in two-group discriminant analysis with the same covariance matrix. We propose a distance-based criterion (DC) drawing on the distance of each variables. The selection method can be applied for high-dimensional data. Sufficient conditions for the distance-based criterion to be consistent are provided when the dimension and the sample are large. Our results, and tendencies therein are explored numerically through a Monte Carlo simulation.

*AMS 2000 subject classification:* primary 62H12; secondary 62H30

*Key Words and Phrases:* Consistency, Discriminant analysis, High-dimensional framework, Selection of variables, Distance-based criterion.

# 1. Introduction

This paper is concerned with the variable selection problem in two-group discriminant analysis with the same covariance matrix. Some methods have been proposed for high-dimensional data as well as large-sample data. One of the methods are based on model selection criteria such as AIC and BIC, see, e.g., Fujikoshi (1984), Nishii et al. (1988). There are methods based on misclassification errors by MaLachlan (1976), Fujikoshi (1985), Hyodo and Kubokawa (2014), Yamada et al. (2017). For high-dimensional data, there are Lasso and other regularization methods by Clemmensen et al. (2011), Witten and Tibshirani (2011), etc.

In this paper we propose a distance-based criterion based on the distance of each variables, which is useful for high-dimensional data as well as large-sample data. Here, the distance of each variables is measured as the one except for the effect of other variables. The criterion involves a constant term which should be determined through some optimality. We propose a class of constants satisfying a consistency when the dimension and the sample size are large. The class depends on the parameters. With respect to the actual use, we also propose a class of estimators for the constant term which satisfies a consistency. Our results, and tendencies therein are explored numerically through a Monte Carlo simulation.

The remainder of the present paper is organized as follows. In Section 2, we present the relevant notation and the distance-based criterion. In Sections 3 we derive sufficient conditions for the distance-based criterion to be consistent under a high-dimensional case. In Section 4 we also propose a class of estimators for the constant term. The proposed distance-based criterion is also examined through a Monte Carlo simulation in Section 5. In Section 6, conclusions are offered. All proofs of our results are provided in the Appendix.

## 2. Distance-based Criterion

In two-group discriminant analysis, suppose that we have independent samples  $\mathbf{y}_1^{(i)}, \dots, \mathbf{y}_{n_i}^{(i)}$  from  $p$ -dimensional normal distributions  $N_p(\boldsymbol{\mu}^{(i)}, \boldsymbol{\Sigma})$ ,  $i = 1, 2$ . Let  $\mathbf{Y}$  be the total sample matrix defined by

$$\mathbf{Y} = (\mathbf{y}_1^{(1)}, \dots, \mathbf{y}_{n_1}^{(1)}, \mathbf{y}_1^{(2)}, \dots, \mathbf{y}_{n_2}^{(2)})'.$$

Let  $\Delta$  and  $D$  be the population and the sample Mahalanobis distances defined by  $\Delta = \{(\boldsymbol{\mu}^{(1)} - \boldsymbol{\mu}^{(2)})' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}^{(1)} - \boldsymbol{\mu}^{(2)})\}^{1/2}$ , and

$$D = \{(\bar{\mathbf{x}}^{(1)} - \bar{\mathbf{x}}^{(2)})' \mathbf{S}^{-1} (\bar{\mathbf{x}}^{(1)} - \bar{\mathbf{x}}^{(2)})\}^{1/2},$$

respectively. Here  $\bar{\mathbf{x}}^{(1)}$  and  $\bar{\mathbf{x}}^{(2)}$  are the sample mean vectors, and  $\mathbf{S}$  be the pooled sample covariance matrix based on  $n = n_1 + n_2$  samples. Let  $\boldsymbol{\beta}$  be the coefficients of the population discriminant function given by

$$\boldsymbol{\beta} = \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}^{(1)} - \boldsymbol{\mu}^{(2)}) = (\beta_1, \dots, \beta_p)'. \quad (2.1)$$

Suppose that  $j$  denotes a subset of  $\omega = \{1, \dots, p\}$  containing  $p_j$  elements, and  $\mathbf{y}_j$  denotes the  $p_j$  vector consisting of the elements of  $\mathbf{y}$  indexed by the elements of  $j$ . We use the notation  $D_j$  and  $D_\omega$  for  $D$  based on  $\mathbf{y}_j$  and  $\mathbf{y}_\omega (= \mathbf{y})$ , respectively.

One of the approaches for variable selection is to consider a set of models as follows. For each subset  $j$  of  $\omega = \{1, \dots, p\}$ , consider a variable selection model  $M_j$  defined by

$$M_j : \beta_i \neq 0 \text{ if } i \in j, \text{ and } \beta_i = 0 \text{ if } i \notin j. \quad (2.2)$$

We identify the selection of  $M_j$  with the selection of  $\mathbf{y}_j$ . Let  $\text{AIC}_j$  be the AIC for  $M_j$ . Then, it is known (see, e.g., Fujikoshi 1985) that

$$\begin{aligned} A_j &= \text{AIC}_j - \text{AIC}_\omega \\ &= n \log \left\{ 1 + \frac{g^2 (D_\omega^2 - D_j^2)}{n - 2 + g^2 D_j^2} \right\} - 2(p - p_j), \end{aligned} \quad (2.3)$$

where  $g = \sqrt{(n_1 n_2)/n}$ . Similarly, let  $\text{BIC}_j$  be the BIC for  $M_j$ , and we have

$$\begin{aligned} B_j &= \text{BIC}_j - \text{BIC}_\omega \\ &= n \log \left\{ 1 + \frac{g^2(D_\omega^2 - D_j^2)}{n - 2 + g^2 D_j^2} \right\} - (\log n)(p - p_j). \end{aligned} \quad (2.4)$$

In a large sample framework, it is known (see, Fujikoshi 1985, Nishii et al. 1988) that AIC is not consistent, but BIC is consistent.

The variable selection methods based on AIC and BIC are given as  $\min_j \text{AIC}_j$  and  $\min_j \text{BIC}_j$ , respectively. Therefore, such model selection methods have a computationally onerous problem when  $p$  is large, since the methods involve minimizing with respect to the subsets of  $2^p$ . To circumvent this issue, we consider a distance-based criterion (DC) drawing on the distance of each variables. A contribution of  $y_i$  in the distance between two groups may be measured by

$$D^2 - D_{(-i)}^2 \quad (2.5)$$

where  $(-i)$  is the subset of  $\omega = \{1, \dots, p\}$  obtained by omitting  $i$  from  $\omega$ . If " $D^2 - D_{(-i)}^2$ " is large, it may be considered that  $y_i$  has a large contribution in discriminant analysis. Conversely, if " $D^2 - D_{(-i)}^2$ " is small, it may be considered that  $y_i$  has a small contribution in discriminant analysis. Let us define

$$D_{d,i} = D^2 - D_{(-i)}^2 - d, \quad i = 1, \dots, p, \quad (2.6)$$

where  $d$  is a positive constant. We propose a distance-based criterion for the selection of variables defined by selecting the set of suffixes or the set of variables given by

$$\text{DC}_d = \{i \in \omega \mid D_{d,i} > 0\}, \quad (2.7)$$

or  $\{y_i \in \{y_1, \dots, y_p\} \mid D_{d,i} > 0\}$ . The notation  $\hat{j}_{\text{DC}_d}$  is also used for  $\text{DC}_d$ .

For a distance-based criterion, there is an important problem how to decide the constant term  $d$ . In Section 3 we derive a range of  $d$  which has a high-dimensional consistency. The range depends on the population parameters. For a practical usefulness, we also propose a class of estimators

for  $d$  given by

$$\hat{d} = \frac{n-2+g^2D^2}{(n-p-1)g^2} \left\{ a(np)^{1/3} + \frac{n-p-1}{n-p-3} \right\}, \quad (2.8)$$

where  $a$  is a constant satisfying  $a = O(1)$ . It is shown that  $\text{DC}_{\hat{d}}$  is consistent under a high-dimensional framework.

### 3. High-dimensional Consistency

For studying consistency properties of the variable selection criterion  $\text{DC}_d$ , it is assumed that the true model  $M_*$  is specified through the values of  $\boldsymbol{\mu}^{(1)}$ ,  $\boldsymbol{\mu}^{(2)}$  and  $\Sigma$ . For simplicity, we write the variable selection models in terms of  $\boldsymbol{\beta}$ , or  $M_j$ . Assume that the true model  $M_*$  is included in the full model  $M_\omega$ . Let the minimum model including  $M_*$  be  $M_{j_*}$ . For a notational simplicity, we regard the true model  $M_*$  as  $M_{j_*}$ . Let  $\mathcal{F}$  be the entire suite of candidate models, defined by

$$\mathcal{F} = \{\{1\}, \dots, \{k\}, \{1, 2\}, \dots, \{1, \dots, k\}\},$$

Let  $\mathcal{F}$  separate into two sets, one is a set of overspecified models, i.e.,  $\mathcal{F}_+ = \{j \in \mathcal{F} \mid j_* \subseteq j\}$  and the other is a set of underspecified models, i.e.,  $\mathcal{F}_- = \mathcal{F}_+^c \cap \mathcal{F}$ .

Here we list some of our main assumptions:

A1 (The true model):  $M_{j_*} \in \mathcal{F}$ .

A2 (The high-dimensional asymptotic framework):

$$p \rightarrow \infty, \quad n \rightarrow \infty, \quad p/n \rightarrow c \in (0, \infty).$$

A3 :  $p_*$  is finite, and  $\Delta^2 = O(1)$ .

A4 :  $n_i/n \rightarrow k_i > 0$ ,  $i = 1, 2$ .

Relating to our proof of a consistency of  $DC_d$  in (2.7), note that

$$DC_d = j_* \Leftrightarrow "D_{d,i} > 0 \text{ for } i \in j_*, \text{ and } "D_{d,i} \leq 0 \text{ for } i \notin j_*$$

Therefore,

$$\begin{aligned} P(DC_d = j_*) &= P\left(\bigcap_{i \in j_*} "D_{d,i} > 0" \bigcap_{i \notin j_*} "D_{d,i} < 0"\right) \\ &= 1 - P\left(\bigcup_{i \in j_*} "D_{d,i} \leq 0" \bigcup_{i \notin j_*} "D_{d,i} \geq 0"\right) \\ &\geq 1 - \sum_{i \in j_*} P(D_{d,i} \leq 0) - \sum_{i \notin j_*} P(D_{d,i} \geq 0). \end{aligned}$$

Our result will be obtained by proving the followings:

$$[F1] \equiv \sum_{i \in j_*} P(D_{d,i} \leq 0) \rightarrow 0. \quad (3.1)$$

$$[F2] \equiv \sum_{i \notin j_*} P(D_{d,i} \geq 0) \rightarrow 0, \quad (3.2)$$

[F1] is the probability that  $DC_d$  does not select the set of true variables. [F2] is the probability that  $DC_d$  select the set of no true variables. Our consistency depends on

$$\tau_{\min} = \min_{i \in j_*} \tau_i, \quad (3.3)$$

where  $\tau_i = \Delta^2 - \Delta_{(-i)}^2$ ,  $i \in \omega$ .

**Theorem 3.1.** *Suppose that assumptions A1, A2, A3 and A4 are satisfied. Let  $b = \tau_{\min}/(1 - c)$ , where  $\tau_{\min}$  is defined by (3.3). Then, if  $0 < d < b$ , the distance-based criterion  $DC_d$  satisfies [F1], i.e.,*

$$\sum_{i \in j_*} P(D_{d,i} \leq 0) \rightarrow 0. \quad (3.4)$$

Related to a sufficient condition for (3.2), we use the following result: for  $i \notin j_*$ ,

$$\begin{aligned} \mathbb{E} \{D^2 - D_{(-i)}^2\} &= \frac{n-2}{(n-p-3)(n-p-2)} \{g^{-2}(n-3) + \Delta^2\} \\ &\equiv q(p, n, \Delta^2). \end{aligned} \quad (3.5)$$

**Theorem 3.2.** *Suppose that assumptions A1, A2, A3 and A4 are satisfied. Then, if  $d = O(1)$  and  $d > q(p, n, \Delta^2)$ , the distance-based criterion  $DC_d$  satisfies*

$$\sum_{i \notin j_*} P(D_{d,i} \geq 0) \rightarrow 0, \quad (3.6)$$

Combining Theorems 3.1 and 3.2, a sufficiency condition for  $DC_d$  to be consistent is given as follows:.

**Theorem 3.3.** *Suppose that assumptions A1, A2, A3 and A4 are satisfied. Then, if  $d = O(1)$  and*

$$q(p, n, \Delta^2) < d < b = \tau_{\min}/(1-c), \quad (3.7)$$

*the distance-based criterion  $DC_d$  is consistent.*

For an actual use, it is necessary to get an estimator  $\hat{d}$  for  $d$ . This problem is discussed in Section 4.

## 4. A Method of Determining the Constant Term

In the previous section we have seen that the distance-based criterion  $DC_d$  is consistent if we use a constant term  $d$  such that  $d = O(1)$  and  $q(p, n, \Delta) < d < \tau_{\min}$ , where  $q(p, n, \Delta)$  is given by (3.5),  $\alpha_{\min} = \min_{i \in j_*} \alpha_i$  and  $\alpha_i =$

$\Delta^2 - \Delta_{(-i)}^2$ . However, such a  $d$  involves unknown parameters. We need to estimate  $d$  by using the estimators of  $\Delta^2$  and  $\tau_{min}$ . Here, we consider to construct an estimator by using the fact that  $DC_d$  is based on the test of  $\alpha_i = 0$  or  $\beta_i = 0$ . The likelihood test is based on (see, e.g., Rao 1973, Fujikoshi et al. 2010)

$$\frac{g^2 D_{\{i\},(-i)}^2}{n - 2 + g^2 D_{(-i)}^2}$$

whose null distribution is a ratio  $\chi_1^2/\chi_{n-p-1}^2$  of independent chi-squared variates  $\chi_1^2$  and  $\chi_{n-p-1}^2$ , where  $D_{\{i\},(-i)}^2 = D^2 - D_{(-i)}^2$ . Let us define  $\hat{d}$  as follows:

$$\hat{d} = \frac{n - 2 + g^2 D^2}{(n - p - 1)g^2} \left\{ a(np)^{1/3} + \frac{n - p - 1}{n - p - 3} \right\}. \quad (4.1)$$

Here  $a$  is a constant satisfying  $a = O(1)$ . Then, it is shown that  $DC_{\hat{d}}$  has a consistency.

**Theorem 4.1.** *Suppose that assumptions A1, A2, A3 and A4 are satisfied. Then, the distance-based criterion  $DC_{\hat{d}}$  with  $\hat{d}$  in (4.1) is consistent.*

As a special  $\hat{d}$ , we consider the followings:

$$\hat{d}_1 = \frac{n - 2 + g^2 D^2}{(n - p - 1)g^2} \left\{ (np)^{1/3} + \frac{n - p - 1}{n - p - 3} \right\}, \quad (4.2)$$

$$\hat{d}_2 = \frac{n - 2 + g^2 D^2}{(n - p - 1)g^2} \left\{ \left(1 - \frac{p}{n}\right)^2 (np)^{1/3} + \frac{n - p - 1}{n - p - 3} \right\}. \quad (4.3)$$

The estimators  $\hat{d}_1$  and  $\hat{d}_2$  are defined from  $\hat{d}$  by choosing  $a$  as  $a = 1$  and  $a = (1 - p/n)^2$ , respectively.

In Theorem 3.3, we have seen that  $DC_d$  is consistent under a high-dimensional framework if  $d$  satisfies  $q(p, n, \Delta^2) < d < \tau_{min}/(1 - c)$ . The expectation of  $\hat{d}$  is related to the range as follows:

- (1)  $E(\hat{d}) > q(p, n, \Delta^2)$ ,
- (2)  $\lim_{p/n \rightarrow c} E(\hat{d}) = 0 \leq \tau_{min}/(1 - c)$ .



The above result follows from that

$$\begin{aligned} \mathbb{E}[\hat{d}] &= \left\{ a(np)^{1/3} + \frac{n-p-1}{n-p-3} \right\} \frac{n-2 + g^2 \frac{p+g^2\Delta^2}{n-p-3}}{(n-p-1)g^2} \\ &= a(np)^{1/3} \left\{ \frac{(n-2)(n-3)}{g^2(n-p-3)(n-p-1)} + \frac{n-2}{(n-p-3)(n-p-1)} \Delta^2 \right\} \\ &\quad + \frac{(n-2)(n-3)}{g^2(n-p-3)^2} + \frac{n-2}{(n-p-3)^2} \Delta^2. \end{aligned}$$

## 5. Numerical Study

In this section we numerically explore the validity of our claims through three distance-based criteria,  $\text{DC}_{d_0}$ ,  $\text{DC}_{\hat{d}_1}$ , and  $\text{DC}_{\hat{d}_2}$ . Here,  $d_0$  is the midpoint in the interval  $[q(p, n, \Delta^2), \tau_{\min}/(1-c)]$  in (3.7).

The true model was assumed as follows: the true dimension is  $p_* = 4, 8$ , the true mean vectors;

$$\boldsymbol{\mu}_1 = \alpha(1, \dots, 1, 0, \dots, 0)', \quad \boldsymbol{\mu}_2 = \alpha(-1, \dots, -1, 0, \dots, 0)',$$

and the true covariance matrix is  $\boldsymbol{\Sigma}_* = \mathbf{I}_p$ .

The selection rates associated with these criteria are given in Tables 5.1 to 5.4. "Under", "True", and "Over" denote the underspecified models, the true model, and the overspecified models, respectively. We focused on selection rates for  $10^3$  replications in Tables 5.1 ~ 5.4. From these tables, we can identify the following tendencies.

Table 5.1. Selection rates of  $\text{DC}_{d_0}$ ,  $\text{DC}_{\hat{d}_1}$  and  $\text{DC}_{\hat{d}_2}$  for  $p_* = 4$  and  $\alpha = 1$

$p_* = 4, \alpha = 1$			$d_0$			$\widehat{d}_1$			$\widehat{d}_2$		
$n_1$	$n_2$	$p$	Under	True	Over	Under	True	Over	Under	True	Over
50	50	10	0.34	0.65	0.01	0.44	0.56	0.00	0.25	0.75	0.01
100	100	10	0.11	0.89	0.00	0.00	1.00	0.00	0.00	0.99	0.01
200	200	10	0.01	0.99	0.00	0.00	1.00	0.00	0.00	1.00	0.00
50	50	25	0.39	0.51	0.10	0.95	0.05	0.00	0.45	0.52	0.03
100	100	50	0.19	0.81	0.01	0.63	0.37	0.00	0.06	0.94	0.00
200	200	100	0.03	0.97	0.00	0.06	0.94	0.00	0.00	1.00	0.00
50	50	50	0.55	0.17	0.29	1.00	0.00	0.00	0.54	0.24	0.22
100	100	100	0.28	0.59	0.13	1.00	0.00	0.00	0.10	0.61	0.29
200	200	200	0.09	0.91	0.00	0.99	0.01	0.00	0.00	0.92	0.08

Table 5.2. Selection rates of  $DC_{d_0}$ ,  $DC_{\widehat{d}_1}$  and  $DC_{\widehat{d}_2}$  for  $p_* = 4$  and  $\alpha = 10$

$p_* = 4, \alpha = 10$			$d_0$			$\widehat{d}_1$			$\widehat{d}_2$		
$n_1$	$n_2$	$p$	Under	True	Over	Under	True	Over	Under	True	Over
50	50	10	0.28	0.71	0.01	0.18	0.82	0.00	0.10	0.90	0.01
100	100	10	0.07	0.93	0.00	0.00	1.00	0.00	0.00	1.00	0.00
200	200	10	0.00	1.00	0.00	0.00	1.00	0.00	0.00	1.00	0.00
50	50	25	0.34	0.62	0.04	0.76	0.24	0.00	0.19	0.78	0.03
100	100	50	0.10	0.90	0.00	0.18	0.82	0.00	0.01	0.98	0.01
200	200	100	0.01	0.99	0.00	0.00	1.00	0.00	0.00	1.00	0.00
50	50	50	0.49	0.29	0.22	1.00	0.00	0.00	0.30	0.36	0.34
100	100	100	0.22	0.72	0.06	1.00	0.01	0.00	0.02	0.69	0.29
200	200	200	0.05	0.95	0.00	0.81	0.19	0.00	0.00	0.90	0.10

Table 5.3. Selection rates of  $DC_{d_0}$ ,  $DC_{\widehat{d}_1}$  and  $DC_{\widehat{d}_2}$  for  $p_* = 8$  and  $\alpha = 1$

$p_* = 8, \alpha = 1$			$d_0$			$\widehat{d}_1$			$\widehat{d}_2$		
$n_1$	$n_2$	$p$	Under	True	Over	Under	True	Over	Under	True	Over
50	50	10	0.81	0.18	0.01	1.00	0.00	0.00	1.00	0.00	0.00
100	100	10	0.52	0.48	0.00	0.79	0.21	0.00	0.65	0.35	0.00
200	200	10	0.20	0.80	0.00	0.03	0.97	0.00	0.02	0.98	0.00
50	50	25	0.86	0.05	0.09	1.00	0.00	0.00	1.00	0.00	0.00
100	100	50	0.68	0.26	0.06	1.00	0.00	0.00	0.97	0.03	0.00
200	200	100	0.30	0.69	0.01	1.00	0.00	0.00	0.53	0.47	0.00
50	50	50	0.93	0.00	0.07	1.00	0.00	0.00	1.00	0.00	0.00
100	100	100	0.79	0.05	0.16	1.00	0.00	0.00	0.96	0.02	0.02
200	200	200	0.52	0.40	0.09	1.00	0.00	0.00	0.53	0.43	0.04

Table 5.4. Selection rates of  $DC_{d_0}$ ,  $DC_{\hat{d}_1}$  and  $DC_{\hat{d}_2}$  for  $p_* = 8$  and  $\alpha = 10$

$p_* = 8, \alpha = 1$			$d_0$			$\hat{d}_1$			$\hat{d}_2$		
$n_1$	$n_2$	$p$	Under	True	Over	Under	True	Over	Under	True	Over
50	50	10	0.79	0.19	0.02	1.00	0.00	0.00	0.99	0.01	0.00
100	100	10	0.48	0.52	0.00	0.61	0.40	0.00	0.48	0.52	0.00
200	200	10	0.14	0.86	0.00	0.00	1.00	0.00	0.00	1.00	0.00
50	50	25	0.84	0.08	0.08	1.00	0.00	0.00	1.00	0.00	0.00
100	100	50	0.60	0.33	0.06	1.00	0.00	0.00	0.91	0.09	0.00
200	200	100	0.25	0.75	0.00	1.00	0.00	0.00	0.27	0.73	0.00
50	50	50	0.91	0.00	0.09	1.00	0.00	0.00	0.99	0.00	0.01
100	100	100	0.74	0.08	0.18	1.00	0.00	0.00	0.89	0.06	0.04
200	200	200	0.45	0.49	0.06	1.00	0.00	0.00	0.31	0.63	0.06

From Tables 5.1 ~ 5.4, we can identify the following tendencies.

- In all the DC criteria with  $d_0$ ,  $\hat{d}_1$  and  $\hat{d}_1$ , there are cases which their selection rates are high. Especially, when the dimension  $p$  is small as in  $p = 10$ , their selection rates are near to 1. However, when  $n$  and  $p$  are near, their selection rates do not near to one when  $n$  and  $p$  are not large. On the other hand, for the case which a consistency can not be seen in Tables, it will be expected to be consistent as the values of  $n$  and  $p$  are large.
- In all the DC criteria with  $d_0$ ,  $\hat{d}_1$  and  $\hat{d}_1$ , their selection rates of the true model are increasing as  $n$  and  $p$  are large.
- In all the DC criteria with  $d_0$ ,  $\hat{d}_1$  and  $\hat{d}_1$ , their selection rates of the true model are decreasing as  $n$  and  $p$  are closer.
- In all the DC criteria with  $d_0$ ,  $\hat{d}_1$  and  $\hat{d}_1$ , their selection rates of the true model are decreasing as the dimension  $p_*$  is increasing.
- In all the DC criteria with  $d_0$ ,  $\hat{d}_1$  and  $\hat{d}_1$ , their selection rates of the true model are increasing as the distance between groups or  $\alpha$  is increasing.
- The DC criteria with  $d_0$  has a tendency of not selecting overspecified models as  $p$  is small in comparison with  $n$ , and has a tendency of selecting overspecified models as  $p$  and  $n$  are near.

- The DC criteria with  $\widehat{d}_1$  does not select overspecified models, and has a tendency of choosing the variables strictly.
- The DC criteria with  $\widehat{d}_2$  has a tendency of not selecting overspecified models when  $p$  is small, and has a tendency of selecting overspecified models as  $n$  and  $p$  are near, as in the DC criteria with  $d_0$ .

## 6. Concluding Remarks

In this paper we propose a distance-based criterion (DC) for the variable selection problem, based on drawing on the significance of each the variables except for the effect due to other variables. The criterion involves a constant term  $d$ , and is denoted by  $\text{DC}_d$ .  $\text{DC}_d$  need not to examine all the subsets, but need to examine only the  $p$  subsets  $\omega_{(-i)}$ . The circumvent computational complexities associated with all the subsets selection procedures been resolved. Further, it was identified that a range of  $d$  and a class of its estimators such that  $\text{DC}_d$  has a high-dimensional consistency property. Especially the criteria  $\text{DC}_d$  with  $d = d_0$ ,  $d = \widehat{d}_1$  and  $d = \widehat{d}_2$  were numerically examined. However, an optimality problem on  $\widehat{d}$  is left for future research. A study of high-dimensional consistency properties when  $p_*$  is infinite and  $\Delta^2 = O(p)$  are also left.

## Appendix: Proofs of Theorems 3.1 and 3.2

### A1 Preliminary Lemmas

First we summarize the distributional results related to the Mahalanobis distance and the statistics  $D_{d,i}$  in (2.5). For a notational simplicity, consider a decomposition of  $\mathbf{y} = (\mathbf{y}'_1, \mathbf{y}'_2)'$ ,  $\mathbf{y}$ ;  $p_1 \times 1$ ,  $\mathbf{p}_2$ ;  $p_2 \times 1$ . Let  $D_1$  and  $D$  be the

sample Mahalanobis distances based on  $\mathbf{y}_1$  and  $\mathbf{y}$ , respectively. Let  $D_{2,1}^2 = D^2 - D_1^2$ . Similarly, the corresponding population quantities are expressed as  $\Delta_1$ ,  $\Delta$  and  $\Delta_{2,1}^2$ . Let the coefficient vector  $\boldsymbol{\beta}$  of the linear discriminant function in (2.1) decompose as  $\boldsymbol{\beta} = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)'$ ,  $\boldsymbol{\beta}_1$ ;  $p_1 \times 1$ ,  $\boldsymbol{\beta}_2$ ;  $p_2 \times 1$ . The following lemmas (see, e.g., Fujikoshi et al. 2010) are used.

**Lemma A1.1.** *For the population Mahalanobis distances  $\Delta$  and  $\Delta_1$ , it holds that*

- (1)  $\Delta_{2,1}^2 = \Delta^2 - \Delta_1^2 \geq 0$ .
- (2)  $\Delta_{2,1}^2 = 0 \Leftrightarrow \boldsymbol{\beta}_2 = \mathbf{0}$ .

Note that

$$\Delta^2 - \Delta_{(-i)}^2 = 0 \Leftrightarrow \beta_i = 0. \quad (\text{A1.1})$$

**Lemma A1.2.** *For the sample Mahalanobis distances  $D_1$  and  $D$ , it holds that*

- (1)  $D_1^2 = (n-2)g^{-2}\chi_{p_1}^2(g^2\Delta_1^2)\{\chi_{n-p_1-1}^2\}^{-1}$ .
- (2) If  $\Delta_{2,1}^2 = 0$ ,  $g^2(D^2 - D_1^2)\{n-2+g^2D_1^2\}^{-1} = \chi_{p_2}^2\{\chi_{n-p-1}^2\}^{-1}$ .
- (3) If  $\Delta_{2,1}^2 = 0$ ,  $(D^2 - D_1^2)\{n-2+g^2D_1^2\}^{-1}$  and  $D_1^2$  are independent,
- (4)  $D_{2,1}^2 = (n-2)g^{-2}\chi_{p_2}^2\left(g^2\Delta_{2,1}^2 \cdot \frac{1}{1+R}\right)\{\chi_{n-p-1}^2\}^{-1}(1+R)$ ,

where  $R = \chi_{p_1}^2(g^2\Delta_1^2)\{\chi_{n-p_1-1}^2\}^{-1}$ . Here,  $\chi_{p_1}^2(\cdot)$ ,  $\chi_{n-p_1-1}^2$ ,  $\chi_{p_2}^2(\cdot)$ , and  $\chi_{n-p-1}^2$  are independent chi-square variates.

From Lemma A1.1,  $D^2$  and  $D_{(-i)}^2$  can be expressed as

$$D^2 = (n-2)g^{-2}\frac{\chi_p^2(g^2\Delta^2)}{\chi_{n-p-1}^2}, \quad D_{(-i)}^2 = (n-2)g^{-2}\frac{\chi_{p-1}^2(g^2\Delta_{(-i)}^2)}{\chi_{n-p-2}^2}, \quad (\text{A1.2})$$

where  $g = \sqrt{(n_1n_2)/n}$ . It is easy to see that

$$\frac{1}{m}\chi_m^2 \xrightarrow{P} 1, \quad \frac{1}{m}\chi_m^2(\delta^2) \xrightarrow{P} 1 + \eta^2, \quad (\text{A1.3})$$

if  $\eta^2 = \lim_{m \rightarrow \infty} \delta^2/m$ .

## A2 Proofs of Theorems 3.1 and 3.2

Without loss of generality, we may  $j_* = \{1, \dots, s\}$  and  $p_* = s$ . Then, [F1] and [F2] are expressed as

$$[\text{F1}] = \sum_{i=1}^s P(D_{d,i} \leq 0), \quad [\text{F2}] = \sum_{i=s+1}^p P(D_{d,i} \geq 0).$$

### Proof of Theorem 3.1

Using (A1.2) and (A1.3),

$$\begin{aligned} D^2 &= (n-2)g^{-2} \frac{\chi_p^2(g^2\Delta^2)}{\chi_{n-p-1}^2} \\ &= (n-2) \frac{n}{n_1 n_2} \cdot \frac{\chi_p^2(g^2\Delta^2)}{p} \cdot \frac{n-p-1}{\chi_{n-p-1}^2} \cdot \frac{p}{n-p-1} \\ &\xrightarrow{p} \frac{c}{k_1 k_2 (1-c)} + \frac{1}{1-c} \Delta^2. \end{aligned}$$

Similarly,

$$D_{(-i)}^2 \xrightarrow{p} \frac{c}{k_1 k_2 (1-c)} + \frac{1}{1-c} \Delta_{(-i)}^2.$$

Therefore,

$$D^2 - D_{(-i)}^2 \xrightarrow{p} \frac{1}{1-c} (\Delta^2 - \Delta_{(-i)}^2)$$

If  $i \in j_* = \{1, \dots, s\}$ , then,  $\tau_i = \Delta^2 - \Delta_{(-i)}^2 > 0$ ,  $i = 1, \dots, s$ , and

$$D^2 - D_{(-i)}^2 \xrightarrow{p} \frac{1}{1-c} \tau_i.$$

Therefore, if  $d < \frac{1}{1-c} \tau_{\min}$ , where  $\tau_{\min} = \min_{i=1, \dots, s} \tau_i$ ,  $D_{d,i} = D^2 - D_{(-i)}^2 - d$  is positive, and

$$\begin{aligned} P(D_{d,i} \leq 0) &\rightarrow 0 \\ \Rightarrow [\text{F1}] &= \sum_{i=1}^s P(D_{d,i} \leq 0) \rightarrow 0, \end{aligned}$$

since  $s$  is finite.

### Proof of Theorem 3.2

Next we shall show [F2] under the assumptions in Theorem 3.2. Assume that  $i \in \{s + 1, \dots, p\}$ . Then, we have  $\Delta^2 = \Delta_{(-i)}^2$ . In general, note that  $D^2 - D_{(-i)}^2 > 0$ . Using Lemma A1.2 (1),

$$\mathbb{E}(D^2) = \frac{(n-2)np}{n_1 n_2 (n-p-3)} + \frac{n-2}{n-p-3} \Delta^2,$$

and hence

$$\mathbb{E} \{ D^2 - D_{(-i)}^2 \} = q(p, n, \Delta^2),$$

where  $q(p, n, \Delta^2)$  is given by (3.5). Suppose that

$$d > q(p, n, \Delta^2).$$

Then, we have

$$\begin{aligned} & P(D^2 - D_{(-i)}^2 \geq d) \\ &= P(D^2 - D_{(-i)}^2 - \mathbb{E}[D^2 - D_{(-i)}^2] \geq d - \mathbb{E}[D^2 - D_{(-i)}^2]) \\ &\leq P(|D^2 - D_{(-i)}^2 - \mathbb{E}[D^2 - D_{(-i)}^2]| \geq d - \mathbb{E}[D^2 - D_{(-i)}^2]) \\ &\leq \frac{1}{(d - \mathbb{E}[D^2 - D_{(-i)}^2])^2} \text{Var}(D^2 - D_{(-i)}^2). \end{aligned}$$

The last inequality follows from Chebyshev's inequality. It is easy to see that  $\mathbb{E}[D^2 - D_{(-i)}^2] = O(n^{-1})$ . Since  $d = O(1)$  from our assumption,

$$(d - \mathbb{E}[D^2 - D_{(-i)}^2])^2 = O(1).$$

Using Lemma A1.2 (4), the denominator is expressed as

$$\begin{aligned} \text{Var}(D^2 - D_{(-i)}^2) &= \mathbb{E}\{(D^2 - D_{(-i)}^2)^2\} - \{\mathbb{E}(D^2 - D_{(-i)}^2)\}^2 \\ &= (n-2)^2 g^{-4} \mathbb{E} \left\{ \left( \frac{\chi_1^2}{\chi_{n-p-1}^2} \right)^2 \left( 1 + \frac{\chi_{p-1}^2 (g^2 \Delta^2)}{\chi_{n-p}^2} \right)^2 \right\} \\ &\quad - (n-2)^2 g^{-4} \frac{1}{(n-p-3)^2} \left( 1 + \frac{p-1+g^2 \Delta^2}{n-p-2} \right)^2. \end{aligned}$$

The expectation of the first term can be computed as

$$\begin{aligned}
& \mathbb{E} \left\{ \left( \frac{\chi_1^2}{\chi_{n-p-1}^2} \right)^2 \left( 1 + \frac{\chi_{p-1}^2(g^2\Delta^2)}{\chi_{n-p}^2} \right)^2 \right\} \\
&= \frac{1}{(n-p-3)(n-p-5)} \mathbb{E} \left\{ 1 + 2 \frac{\chi_{p-1}^2(g^2\Delta^2)}{\chi_{n-p}^2} + \left( \frac{\chi_{p-1}^2(g^2\Delta^2)}{\chi_{n-p}^2} \right)^2 \right\} \\
&= \frac{1}{(n-p-3)(n-p-5)} \left\{ 1 + 2 \frac{p-1+g^2\Delta^2}{n-p-2} \right. \\
&\quad \left. + \frac{(p-1)^2 + 2(p-1) + 2(p-1)g^2\Delta^2 + 4g^2\Delta^2 + g^4\Delta^4}{(n-p-2)(n-p-4)} \right\} \\
&= \frac{1}{(n-p-3)(n-p-5)} \left\{ 1 + O(1) + O(1) \right\} = O(n^{-2}).
\end{aligned}$$

The second term is evaluated as

$$\begin{aligned}
& (n-2)^2 g^{-4} \frac{1}{(n-p-3)^2} \left( 1 + \frac{p-1+g^2\Delta^2}{n-p-2} \right)^2 \\
&= O(1) \times \frac{1}{(n-p-3)^2} (1 + O(1))^2 = O(n^{-2}).
\end{aligned}$$

Therefore, we have

$$\text{Var}(D^2 - D_{(-i)}^2) = (n-2)^2 g^{-4} O(n^{-2}) - O(n^{-2}) = O(n^{-2}).$$

These imply the followings:

$$\begin{aligned}
& P(D^2 - D_{(-i)}^2 \geq d) \leq O(n^{-2}), \text{ and hence} \\
& \sum_{i=s+1}^p P(D^2 - D_{(-i)}^2 \geq d) \rightarrow 0,
\end{aligned}$$

which proves "[F2]  $\rightarrow 0$ ".



### A3 Proof of Theorem 4.1

As in the proofs of Theorems 3.1 and 3.2, assume that  $j_* = \{1, \dots, s\}$ , and show that

$$\begin{aligned} [\text{F1}] &\equiv \sum_{i=1}^s P\left(D^2 - D_{(-i)}^2 \leq \widehat{d}\right) \rightarrow 0, \\ [\text{F2}] &\equiv \sum_{i=s+1}^p P\left(D^2 - D_{(-i)}^2 \geq \widehat{d}\right) \rightarrow 0. \end{aligned}$$

First consider [F1], and assume that  $i \in \{1, \dots, p_*\}$ . In the proof of Theorem 3.1, it has been shown that

$$D^2 \xrightarrow{p} \frac{c}{k_1 k_2 (1-c)} + \frac{1}{1-c} \Delta^2.$$

Using this result, we have

$$\widehat{d} \xrightarrow{p} 0.$$

These imply that

$$D^2 - D_{(-i)}^2 - \widehat{d} \xrightarrow{p} \frac{1}{1-c} \tau_i > 0,$$

and hence

$$P(D^2 - D_{(-i)}^2 \leq \widehat{d}) = 0,$$

which implies "[F1]  $\rightarrow 0$ ".

Next, consider to show "[F2]  $\rightarrow 0$ ". Assume that  $i \in \{s+1, \dots, p\}$ . Then,  $\Delta^2 = \Delta_{(-i)}^2$ . Denote that  $D^2 - D_{(-i)}^2 = D_{\{i\} \cdot (-i)}$ . Using Lemma A1.2(2) and (3), we can write as

$$\begin{aligned} P\left(D^2 - D_{(-i)}^2 \geq \widehat{d}\right) &= P\left(\frac{\chi_1^2}{\chi_{n-p-1}^2} \geq \frac{g^2 \widehat{d}}{n-2+g^2 D_{(i)}^2}\right) \\ &\leq P\left(\frac{\chi_1^2}{\chi_{n-p-1}^2} \geq \frac{g^2 \widehat{d}}{n-2+g^2 D^2}\right) \\ &= P\left(F_{1, n-p-1} - \frac{n-p-1}{n-p-3} \geq a(np)^{1/3}\right). \end{aligned}$$

Noting that the mean of  $F_{1,n-p-1}$  is  $(n-p-1)(n-p-3)^{-1}$ ,

$$\begin{aligned} P\left(D^2 - D_{(-i)}^2 \geq \hat{d}\right) &\leq P\left(|F_{1,n-p-1} - \mathbb{E}[F_{1,n-p-1}]| \geq a(np)^{1/3}\right) \\ &\leq \frac{1}{a^2(np)^{2/3}} \text{Var}(F_{1,n-p-1}) \\ &= \frac{1}{a^2(np)^{2/3}} \frac{2n^2(n-1)}{(n-2)^2(n-4)}. \end{aligned}$$

This implies that

$$\begin{aligned} [\text{F2}] &\leq \sum_{i=j_*+1}^p \frac{1}{a^2(np)^{2/3}} \frac{2n^2(n-1)}{(n-2)^2(n-4)} \\ &\leq \frac{1}{a^2} \frac{p}{(np)^{2/3}} \frac{2n^2(n-1)}{(n-2)^2(n-4)} = O(n^{-1/3}), \end{aligned}$$

which proves "[F2]  $\rightarrow 0$ ".

## Acknowledgements

The first author's research is partially supported by the Ministry of Education, Science, Sports, and Culture, a Grant-in-Aid for Scientific Research (C), 16K00047, 2016-2018.

## References

- [1] CLEMMENSEN, L., HASTIE, T., WITTEN, D. M. and ERBELL, B. (2011). Sparse discriminant analysis. *Technometrics*, **53**, 406–413.
- [2] FUJIKOSHI, Y. (1985). Selection of variables in two-group discriminant analysis by error rate and Akaike's information criteria. *J. Multivariate Anal.*, **17**, 27–37.
- [3] FUJIKOSHI, Y., ULYANOV, V. V. and SHIMIZU, R. (2010). *Multivariate Statistics: High-Dimensional and Large-Sample Approximations*. Wiley, Hoboken, N.J.

- [4] HYODO, M. and KUBOKAWA, T. (2014). A variable selection criterion for linear discriminant rule and its optimality in high dimensional and large sample data. *J. Multivariate Anal.*, **123**, 364–379.
- [5] NISHII, R. , BAI, Z. D. and KRISHNAIA, P. R. (1988). Strong consistency of the information criterion for model selection in multivariate analysis. *Hiroshima Math. J.*, **18**, 451–462.
- [6] RAO, C. R. (1973). *Linear Statistical Inference and Its Applications*, 2nd ed. Wiley, New-York.
- [7] YAMADA, T., SAKURAI, T. and FUJIKOSHI, Y. (2017). High-dimensional asymptotic results for EPMCs of W- and Z- rules. Hiroshima Statistical Research Group, TR 17-11.
- [8] WITTEN, D. W. and TIBSHIRANI, R. (2011). Penalized classification using Fisher’s linear discriminant. *J. R. Statist. Soc. B*, **73**, 753–772.