

CIRJE-F-872

**A Variable Selection Criterion for Linear Discriminant
Rule and its Optimality in High Dimensional Setting**

Masashi Hyodo
Graduate School of Economics, University of Tokyo

Tatsuya Kubokawa
University of Tokyo

December 2012; Revised in February 2013

CIRJE Discussion Papers can be downloaded without charge from:

<http://www.cirje.e.u-tokyo.ac.jp/research/03research02dp.html>

Discussion Papers are a series of manuscripts in their draft form. They are not intended for circulation or distribution except as indicated by the author. For that reason Discussion Papers may not be reproduced or distributed without the written consent of the author.

A Variable Selection Criterion for Linear Discriminant Rule and its Optimality in High Dimensional Setting

Masashi Hyodo*and Tatsuya Kubokawa†

University of Tokyo

February 4, 2013

Abstract

In this paper, we suggest the new variable selection procedure, called MEC, for linear discriminant rule in the high-dimensional setup. MEC is derived as a second-order unbiased estimator of the misclassification error probability of the linear discriminant rule. It is shown that MEC not only decomposes into ‘fitting’ and ‘penalty’ terms like AIC and Mallows C_p , but also possesses an asymptotic optimality in the sense that MEC achieves the smallest possible conditional probability of misclassification in candidate variable sets. Through simulation studies, it is shown that MEC has good performances in the sense of selecting the true variable sets.

Key words and phrases: asymptotic optimality, high dimension, linear discriminant analysis, misclassification error, multivariate normal, second-order approximation, variable selection.

1 Introduction

In this paper, we consider the problem of classifying a future observation vector into one of the two population groups Π_1 and Π_2 . For each $i = 1, 2$, Π_i denotes a population from a multivariate normal distribution $\mathcal{N}_p(\boldsymbol{\mu}_i, \Sigma)$, and it is supposed that \boldsymbol{x}_{ij} , $j = 1, \dots, N_i$, are observed from the population Π_i . Here, $\boldsymbol{\mu}_i$, $i = 1, 2$, and Σ are unknown parameters, and they are estimated by the sample mean $\bar{\boldsymbol{x}}_i = N_i^{-1} \sum_{j=1}^{N_i} \boldsymbol{x}_{ij}$, $i = 1, 2$, and the pooled sample covariance matrix $S = n^{-1} \sum_{i=1}^2 \sum_{j=1}^{N_i} (\boldsymbol{x}_{ij} - \bar{\boldsymbol{x}}_i)(\boldsymbol{x}_{ij} - \bar{\boldsymbol{x}}_i)'$ for $n = N_1 + N_2 - 2$. When the dimension p is large, the model involves many unknown parameters, which causes a large misclassification error in the linear discriminant rule (LDR). Thus, it is desired to find an optimal subset of variables in LER. Variable selection methods for discriminant analysis have been studied by Fujikoshi (1985, 2002), Sakurai, Nakada and

*Faculty of Economics, University of Tokyo, 7-3-1 Hongo, Bunkyo-ku, Tokyo 113-0033, JAPAN, E-Mail: hyodohh@yahoo.co.jp

†Faculty of Economics, University of Tokyo, 7-3-1 Hongo, Bunkyo-ku, Tokyo 113-0033, JAPAN, E-Mail: tatsuya@e.u-tokyo.ac.jp

Fujikoshi (2012), Wilbur, Ghosh, Nakatsu and Doerge (2002) and others. Related to this issue, multiple testing problems for no additional information have been discussed by Rao (1948, 1970, 1973) and Kshirsagar (1972). In this paper, we suggest a new variable selection procedure based on error of misclassification and establish the optimality in high-dimensional situation.

To explain the new variable selection procedure, consider the following linear discriminant rule. Let $\mathbf{x} = (x_1, \dots, x_p)$ be a future observation in the full model. Let $\mathbf{j} = (j_1, \dots, j_{k(\mathbf{j})})$ be a subset of the set $\{1, 2, \dots, p\}$, and let $\mathbf{x}(\mathbf{j}) = (x_{j_1}, \dots, x_{j_{k(\mathbf{j})}})$ be the corresponding subvector of \mathbf{x} . The model based on the variable $\mathbf{x}(\mathbf{j})$ is denoted by \mathbf{j} . Let \mathcal{J} be a suitable family of subsets of $\{1, \dots, p\}$. The LDR for classifying \mathbf{x} based on the model \mathbf{j} is that \mathbf{x} is classified as coming from Π_1 , if $W(\mathbf{j}) > \alpha$, and from Π_2 , if $W(\mathbf{j}) < \alpha$, where α is cut off point for classification rule, and

$$W(\mathbf{j}) = (\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j}))' S(\mathbf{j})^{-1} \{ \mathbf{x}(\mathbf{j}) - \frac{1}{2}(\bar{\mathbf{x}}_1(\mathbf{j}) + \bar{\mathbf{x}}_2(\mathbf{j})) \}.$$

Here, $\bar{\mathbf{x}}_i(\mathbf{j})$, $i = 1, 2$, and $S(\mathbf{j})$ are the sample mean and the pooled sample covariance matrix in the model \mathbf{j} . Then the problem of variable selection in LDR is regarded as how to select the best subset \mathbf{j} from \mathcal{J} . To this end, we consider the conditional error probabilities of misallocation $L_1(\mathbf{j}) = P[W(\mathbf{j}) < \alpha | \mathbf{x}(\mathbf{j}) \in \Pi_1, \bar{\mathbf{x}}_1(\mathbf{j}), \bar{\mathbf{x}}_2(\mathbf{j}), S(\mathbf{j})]$ and $L_2(\mathbf{j}) = P[W(\mathbf{j}) > \alpha | \mathbf{x}(\mathbf{j}) \in \Pi_2, \bar{\mathbf{x}}_1(\mathbf{j}), \bar{\mathbf{x}}_2(\mathbf{j}), S(\mathbf{j})]$, which can be expressed as

$$L_g(\mathbf{j}) = \Phi \left((-1)^g \frac{(\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j}))' S(\mathbf{j})^{-1} \{ \boldsymbol{\mu}_g(\mathbf{j}) - (\bar{\mathbf{x}}_1(\mathbf{j}) + \bar{\mathbf{x}}_2(\mathbf{j}))/2 \}}{\sqrt{(\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j}))' S(\mathbf{j})^{-1} \Sigma(\mathbf{j}) S(\mathbf{j})^{-1} (\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j}))}} \right), \quad (1.1)$$

for $g = 1, 2$, where $\Phi(\cdot)$ is the standard normal distribution function, and $\boldsymbol{\mu}_g(\mathbf{j})$ and $\Sigma(\mathbf{j})$ denote the population mean and covariance matrix in the model \mathbf{j} . When π_i , $i = 1, 2$ is a prior probability of the group membership, the expected error rate is given by

$$R(\mathbf{j}) = \pi_1 R_1(\mathbf{j}) + \pi_2 R_2(\mathbf{j})$$

where $R_g(\mathbf{j})$ is the unconditional error of misallocation given by $R_g(\mathbf{j}) = E[L_g(\mathbf{j})]$ for $g = 1, 2$. The variable selection procedure proposed in this paper is an asymptotically unbiased estimator of the misclassification error $R(\mathbf{j})$ in the high-dimensional setting.

A naive procedure for selection of variables is to minimize

$$\Phi(-D(\mathbf{j})/2), \quad (1.2)$$

with respect to $\mathbf{j} \in \mathcal{J}$, where $D(\mathbf{j})$ is the sample Mahalanobis distance based on $\mathbf{x}(\mathbf{j})$, namely,

$$D(\mathbf{j}) = (\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j}))' S(\mathbf{j})^{-1} (\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j})). \quad (1.3)$$

However, $\Phi(-D(\mathbf{j})/2)$ has the bias $R(\mathbf{j}) - E[\Phi(-D(\mathbf{j})/2)]$ which is not negligible. MacLachlan (1976, 1980) derived a second order asymptotic unbiased estimator of $R(\mathbf{j})$ under the large sample framework, namely,

(A0): $n \rightarrow \infty$, but p is bounded.

Fujikoshi (1985) applied the estimator given by MacLachlan (1976,1980) to the variable selection problem, and investigated the asymptotic properties and the relationship with AIC. In the high-dimensional setting On the other hand, Raudys (1972) and Wyman, Young and Turner (1990) derived the asymptotic approximations of the error probability under the high-dimensional setting given by

$$(A1) : (n, p) \rightarrow \infty \text{ with } p/n \rightarrow c_0 \in [0, 1).$$

It is known that these approximations are also good in the large sample situation (A0) as seen from Fujikoshi, Ulyanov and Shimizu (2009).

In this paper, we derive a second-order unbiased estimator of $R(\mathbf{j})$ in the high-dimensional setting (A1). The unbiased estimator is here called the *Misclassification Error Criterion* (MEC), which is useful for selecting variables in linear discriminant rule. We show that MEC can be decomposed into the ‘fitting’ and ‘penalty’ terms, namely,

$$\text{MEC} = \Phi(-D(\mathbf{j})/2) + (\text{penalty}),$$

where the penalty term increases in the number of the unknown parameters. This is a desirable property that variable selection procedures like AIC and C_p should possess. We also show that MEC has an asymptotical optimality as a variable selection procedure in (A1). Such an optimality in the high-dimensional setting is not known as long as we know.

Recently, Kubokawa, Hyodo and Srivastava (2013) derived a second-order approximation of the error probability of misclassification (EPMC) for the ridge-type linear discriminant rule in high dimensional setting, and derived a second-order unbiased estimator of EPMC. Since the ridge-type linear discriminant rule is not invariant under scale transformations, their approach needs to calculate various kinds of fourth moments of the inverted Wishart matrix. It was hard to obtain such fourth moments, so that the approach used by Kubokawa, *et al.* (2013) cannot be used for developing an asymptotic optimality of MEC. Instead, the method used in this paper is to express $L_g(\mathbf{j})$ based on nine primitive random variables, namely four random variables having the standard normal distribution and five random variables having chi-square distributions. This approach not only makes it easier to derive the second-order approximation and the second-order unbiased estimator of $R(\mathbf{j})$, but also enables us to establish the asymptotic optimality of MEC as a variable selection procedure in both high-dimensional and large-sample situations.

The organization of this paper is as follows. In Section 2.1, we determine the asymptotically optimal cut off point α based on the expected error rate $R(\mathbf{j})$. In Section 2.2, we derive the second order unbiased estimator of $R(\mathbf{j})$ in high dimensional setting and propose the new variable selection procedure MEC based on this estimator. In Section 3.1, we show that MEC can be decomposed into the “fitting term” and “dimensionality penalty term” like Mallows C_p and AIC. In Section 3.2, we prove the asymptotic optimality of MEC. In Section 4, we investigate performances of MEC through numerical studies. The conclusion of our study is summarized in Section 5. Some preliminary results are given in Appendix.

2 MEC: A Variable Selection Procedure for High Dimensional Data

In this Section, we derive the second order asymptotic unbiased estimator of $R(\mathbf{j})$ under the following high dimensional framework (A1)-(A3).

2.1 The optimal cut off point based on expected probability of misclassification

An optimal rule of allocation can be defined by taking it to be one that minimizes the expected error rate $R(\mathbf{j})$. However, the expected error rate $R(\mathbf{j})$ is not explicit formula. So, we determine the optimal cut off point based on limiting value of $R(\mathbf{j})$. We assume the following conditions, in order to derive limiting value of $R(\mathbf{j})$.

$$(A1) : (n, p) \rightarrow \infty \text{ with } p/n \rightarrow c_0 \in [0, 1].$$

$$(A2) : (n, N_1, N_2) \rightarrow \infty \text{ with } N_1/n \rightarrow \gamma_1, N_2/n \rightarrow \gamma_2.$$

$$(A3) : \lim_{p \rightarrow \infty} \Delta^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) < \infty.$$

Let $\#\mathbf{j}$ be the cardinality of \mathbf{j} . In what follows, we treat the case that $(\#\mathbf{j}, n) \rightarrow \infty$ with $\#\mathbf{j}/n \rightarrow c \in [0, 1)$, but the derived results include the case of $\#\mathbf{j}$ bounded. We primarily consider the asymptotic approximation for the expected error rate of $E[L_1(\mathbf{j})]$, since the asymptotic approximations for the expected error rate of $E[L_2(\mathbf{j})]$ can be constructed similarly.

Suppose that $\mathbf{x} \in \Pi_1$. Then, a conditional distribution given $(\bar{\mathbf{x}}_1(\mathbf{j}), \bar{\mathbf{x}}_2(\mathbf{j}), S(\mathbf{j}))$ is written as $W(\mathbf{j}) | (\bar{\mathbf{x}}_1(\mathbf{j}), \bar{\mathbf{x}}_2(\mathbf{j}), S(\mathbf{j})) \sim \mathcal{N}_p(-U, V)$, where

$$\begin{aligned} U &= (\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j}))' S(\mathbf{j})^{-1} (\bar{\mathbf{x}}_1(\mathbf{j}) - \boldsymbol{\mu}_1(\mathbf{j})) - 2^{-1} D(\mathbf{j}), \\ V &= (\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j}))' S(\mathbf{j})^{-1} \boldsymbol{\Sigma}(\mathbf{j}) S(\mathbf{j})^{-1} (\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j})), \end{aligned}$$

where $D(\mathbf{j})$ is given in (1.3). Then, the expected error rate of $W(\mathbf{j})$ can be expressed as

$$E[L_1(\mathbf{j})] = E \left[\Phi \left(\frac{U + \alpha}{\sqrt{V}} \right) \right],$$

where $\Phi(\cdot)$ denotes the distribution function of a standard normal random variable.

As given in (6.1) and (6.2), U and V can be expanded as $U = U_0 + n^{-1/2}U_1 + n^{-1}U_2$ and $V = V_0 + n^{-1/2}V_1 + n^{-1}V_2$ with

$$U_0 = -\frac{1}{2(1-c)} \left(\Delta(\mathbf{j})^2 + \frac{c(\gamma_1 - \gamma_2)}{\gamma_1 \gamma_2} \right), \quad V_0 = \frac{1}{(1-c)^3} \left(\frac{c}{\gamma_1 \gamma_2} + \Delta(\mathbf{j})^2 \right),$$

$U_1 = O_p(1)$, $U_2 = O_p(1)$, $V_1 = O_p(1)$ and $V_2 = O_p(1)$. Then we observe that

$$\frac{U + \alpha}{V^{1/2}} = \left\{ U_0 + \frac{U_1}{\sqrt{n}} + \frac{U_2}{n} + \alpha \right\} \frac{1}{V_0^{1/2}} \left\{ 1 - \frac{1}{2V_0} \left(\frac{V_1}{\sqrt{n}} + \frac{V_2}{n} \right) + \frac{3V_1^2}{8nV_0^2} \right\} + O_p(n^{-3/2}), \quad (2.1)$$

which gives the expansion

$$(U + \alpha)/V^{1/2} = w_0 + w_1 + w_2 + O_p(n^{-3/2}),$$

where

$$\begin{aligned} w_0 &= V_0^{-1/2}(U_0 + \alpha), \\ w_1 &= \frac{1}{\sqrt{n}V_0^{1/2}} \left\{ U_1 - \frac{U_0 + \alpha}{2V_0} V_1 \right\}, \\ w_2 &= \frac{1}{nV_0^{1/2}} \left\{ U_2 - \frac{U_0 + \alpha}{2V_0} V_2 + \frac{3(U_0 + \alpha)}{8V_0^2} V_1^2 - \frac{1}{2V_0} U_1 V_1 \right\}. \end{aligned}$$

Using the Taylor series expansion again, we can approximate $E[L_1(\mathbf{j})]$ as

$$\begin{aligned} E[L_1(\mathbf{j})] &= E[\Phi(w_0 + (w_1 + w_2))] \\ &= \Phi(w_0) + \phi(w_0)E[w_1 + w_2 - \frac{1}{2}w_0w_1^2] + O(n^{-3/2}), \end{aligned}$$

where $\phi(\cdot)$ is the pdf of the standard normal distribution. Letting $H = E[w_1 + w_2 - \frac{1}{2}w_0w_1^2]$, we can write H as

$$\begin{aligned} H &= \frac{1}{V_0^{1/2}} \left\{ \frac{E[U_1]}{\sqrt{n}} + \frac{E[U_2]}{n} \right\} - \frac{U_0 + \alpha}{2V_0^{3/2}} \left\{ \frac{E[V_1]}{\sqrt{n}} + \frac{E[V_2]}{n} \right\} - \frac{U_0 + \alpha}{2V_0^{3/2}} \frac{E[U_1^2]}{n} \\ &\quad + \frac{U_0 + \alpha}{8V_0^{5/2}} \left(3 - \frac{(U_0 + \alpha)^2}{V_0} \right) \frac{E[V_1^2]}{n} - \frac{1}{2V_0^{3/2}} \left(1 - \frac{(U_0 + \alpha)^2}{V_0} \right) \frac{E[U_1 V_1]}{n}. \end{aligned} \quad (2.2)$$

Since U_1, U_2, V_1 and V_2 are given around (6.1) and (6.2), we can calculate the moments in (2.2), which yields the following theorem. Define H_U, H_V, H_1, H_2 and H_{12} by

$$\begin{aligned} H_U(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) &= -\frac{1}{2(1-c)^2} \left(\Delta(\mathbf{j})^2 + \frac{c(\gamma_1 - \gamma_2)}{\gamma_1 \gamma_2} \right), \\ H_V(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) &= \frac{(c+3)\Delta(\mathbf{j})^2 \gamma_1 \gamma_2 + (5-c)c}{(c-1)^4 \gamma_1 \gamma_2}, \\ H_1(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) &= \frac{\Delta(\mathbf{j})^4}{2(1-c)^3} + \frac{1}{(1-c)^3 \gamma_2} \left(\frac{c}{\gamma_1} + \Delta(\mathbf{j})^2 \right) + \frac{c(\gamma_1 - \gamma_2)^2}{2(1-c)^3 \gamma_1^2 \gamma_2^2}, \\ H_2(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) &= \frac{2(c+4)\Delta(\mathbf{j})^4}{(1-c)^7} + \frac{4\{(c+1)^2 + c\} \Delta(\mathbf{j})^2}{(1-c)^7 \gamma_1 \gamma_2} + \frac{2c\{(c+1)^2 + c\}}{(1-c)^7 \gamma_1^2 \gamma_2^2}, \\ H_{12}(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) &= -\frac{2\Delta(\mathbf{j})^4}{(1-c)^5} - \frac{2(c+1)\Delta(\mathbf{j})^2}{(1-c)^5 \gamma_2} - \frac{c(c+1)(\gamma_1^2 - \gamma_2^2)}{(1-c)^5 \gamma_1^2 \gamma_2^2}, \end{aligned}$$

for $\Delta(\mathbf{j}) = \{\boldsymbol{\mu}_1(\mathbf{j}) - \boldsymbol{\mu}_2(\mathbf{j})\}' \boldsymbol{\Sigma}(\mathbf{j})^{-1} \{\boldsymbol{\mu}_1(\mathbf{j}) - \boldsymbol{\mu}_2(\mathbf{j})\}$.

Theorem 2.1 *Assume the conditions (A1)-(A3). The second order approximation of $E[L_1(\mathbf{j})]$ is given by*

$$E[L_1(\mathbf{j})] = \Phi \left(\frac{U_0 + \alpha}{\sqrt{V_0}} \right) + \frac{1}{n} \phi \left(\frac{U_0 + \alpha}{\sqrt{V_0}} \right) H(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) + O(n^{-3/2}), \quad (2.3)$$

where, $H(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c)$ is given by

$$\begin{aligned}
H(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) = & \frac{1}{V_0^{1/2}} \left\{ H_U(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) - \frac{U_0 + \alpha}{2V_0} H_V(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) \right. \\
& - \frac{U_0 + \alpha}{2V_0} H_1(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) \\
& + \frac{U_0 + \alpha}{8V_0^2} \left(3 - \frac{(U_0 + \alpha)^2}{V_0} \right) H_2(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) \\
& \left. - \frac{1}{2V_0} \left(1 - \frac{(U_0 + \alpha)^2}{V_0} \right) H_{12}(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) \right\}. \quad (2.4)
\end{aligned}$$

It is noted that the second order approximation of $E[L_2(\mathbf{j})]$ can be also obtained by replacing $(\gamma_1, \gamma_2, \alpha)$ in (2.3) and (2.4) with $(\gamma_2, \gamma_1, -\alpha)$. Thus, it follows from Theorem 2.1 that the limiting value of $R(\mathbf{j})$ is given by

$$\pi_1 \Phi \left(-\frac{\Delta(\mathbf{j})^2 + \frac{c(\gamma_1 - \gamma_2)}{\gamma_1 \gamma_2} - (1 - c)\alpha}{2\sqrt{1 - c}\sqrt{\Delta^2(\mathbf{j}) + c/(\gamma_1 \gamma_2)}} \right) + \pi_2 \Phi \left(-\frac{\Delta(\mathbf{j})^2 - \frac{c(\gamma_1 - \gamma_2)}{\gamma_1 \gamma_2} + (1 - c)\alpha}{2\sqrt{1 - c}\sqrt{\Delta^2(\mathbf{j}) + c/(\gamma_1 \gamma_2)}} \right),$$

which can be minimized at

$$\alpha = \frac{c(\gamma_1 - \gamma_2)}{2\gamma_1 \gamma_2 (1 - c)} + \left(\frac{1}{(1 - c)^2} + \frac{c}{(1 - c)^2 \Delta(\mathbf{j})^2 \gamma_1 \gamma_2} \right) \log \left(\frac{\pi_2}{\pi_1} \right). \quad (2.5)$$

It is hereafter assumed that $\pi_1 = \pi_2 = 1/2$. Then, the asymptotically optimal classification rule is given by

$$\begin{aligned}
W(\mathbf{j}) = & (\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j}))' S(\mathbf{j})^{-1} \{ \mathbf{x}(\mathbf{j}) - \frac{1}{2}(\bar{\mathbf{x}}_1(\mathbf{j}) + \bar{\mathbf{x}}_2(\mathbf{j})) \} > (\text{resp. } <) \frac{c(\gamma_1 - \gamma_2)}{2\gamma_1 \gamma_2 (1 - c)} \\
\implies & \mathbf{x} \in \Pi_1 (\text{resp. } \Pi_2). \quad (2.6)
\end{aligned}$$

2.2 Derivation of MEC

We now derive a new variable selection procedure, called MEC, based on the misclassification error probability of the linear discriminant rule. As variable selection methods, Mallows's C_p and AIC are well known. Mallows's C_p is an unbiased estimator of the prediction error for a future observation relative to a quadratic loss, and AIC is motivated from minimization of prediction error relative to the Kullback-Leibler information. Recall that $L_1(\mathbf{j})$ given in (1.1) can be interpreted as a conditional prediction error that a future observation \mathbf{x} from Π_1 is misclassified into Π_2 . Thus, the expected error rate $R(\mathbf{j}) = (R_1(\mathbf{j}) + R_2(\mathbf{j}))/2$ for $\pi_1 = \pi_2 = 1/2$ is regarded as a predictive probability of misclassification, which suggests that an unbiased estimator of $R(\mathbf{j})$ possesses the ability to work as a variable selection procedure. In this subsection, we derive the second order unbiased estimator of $R(\mathbf{j})$ in the high dimensional setting (A1) or the large sample setting (A0).

Consider the linear discriminant rule (2.6) with the asymptotically optimal cut-off point α given in (2.5) for $\pi_1 = \pi_2 = 1/2$. For notational convenience, let

$$U_\alpha = U_0 + \alpha = -\frac{1}{2(1 - c)} \Delta(\mathbf{j})^2.$$

Note that $V_0 = (1 - c)^{-3}\{\Delta(\mathbf{j})^2 + c/(\gamma_1\gamma_2)\}$. Since the second-order expansion given in Theorem 2.1 is a function of $\Delta(\mathbf{j})^2$, we begin by obtaining a consistent estimator of $\Delta(\mathbf{j})^2$. Define $\widehat{\Delta}(\mathbf{j})^2$ by

$$\widehat{\Delta}(\mathbf{j})^2 = (1 - c)D(\mathbf{j}) - \frac{c}{\gamma_1\gamma_2}, \quad (2.7)$$

for $D(\mathbf{j})$ given in (1.3). Then, the estimator $\widehat{\Delta}(\mathbf{j})^2$ is expanded as

$$\widehat{\Delta}(\mathbf{j})^2 = \Delta(\mathbf{j})^2 + \frac{D_1}{\sqrt{n}} + \frac{D_2}{n} + O_p(n^{-3/2}),$$

where

$$\begin{aligned} D_1 &= -\frac{\sqrt{2}v_2(c + \Delta(\mathbf{j})^2\gamma_1\gamma_2)}{\sqrt{1 - c\gamma_1\gamma_2}} + \frac{\sqrt{2}\sqrt{c}v_1}{\gamma_1\gamma_2} + \frac{2\Delta(\mathbf{j})u_1}{\sqrt{\gamma_1\gamma_2}}, \\ D_2 &= -\frac{1}{(1 - c)\gamma_1\gamma_2} - \frac{\Delta(\mathbf{j})^2}{1 - c} + \frac{2\left(\frac{c}{\gamma_1\gamma_2} + \Delta(\mathbf{j})^2\right)}{1 - c}v_2^2 + \frac{1}{\gamma_1\gamma_2}u_1^2 \\ &\quad - \frac{2\sqrt{2}\Delta(\mathbf{j})}{\sqrt{1 - c}\sqrt{\gamma_1\gamma_2}}v_2u_1 - \frac{2\sqrt{c}}{\sqrt{1 - c}\gamma_1\gamma_2}v_1v_2 \end{aligned}$$

for u_1 and u_2 given in Lemma 6.1 and v_1 and v_2 defined above (6.1).

We consider to substitute the consistent estimator into the limiting term $\Phi((U_\alpha)V_0^{-1/2})$ or $\Phi(U_\alpha V_0^{-1/2})$ in Theorem 2.1. Let $\widehat{U}_\alpha = -2^{-1}(1 - c)^{-1}\widehat{\Delta}(\mathbf{j})^2$ and $\widehat{V}_0 = (1 - c)^{-3}\{\widehat{\Delta}(\mathbf{j})^2 + c/(\gamma_1\gamma_2)\}$. For the term $\Phi(U_\alpha V_0^{-1/2})$, however, the estimator $\Phi(\widehat{U}_\alpha \widehat{V}_0^{-1/2})$ is not a second order unbiased estimator of $\Phi(U_\alpha V_0^{-1/2})$, since $\Phi(U_\alpha V_0^{-1/2}) = O(1)$. Since $\widehat{\Delta}(\mathbf{j})^2 = \Delta(\mathbf{j})^2 + D_1/\sqrt{n} + D_2/n + O_p(n^{-3/2})$, it is noted that

$$\begin{aligned} \widehat{U}_\alpha &= U_\alpha + c_1 \left(\frac{D_1}{\sqrt{n}} + \frac{D_2}{n} \right) + O_p(n^{-3/2}), \\ \widehat{V}_0 &= V_0 + c_2 \left(\frac{D_1}{\sqrt{n}} + \frac{D_2}{n} \right) + O_p(n^{-3/2}), \end{aligned}$$

for $c_1 = -\{2(1 - c)\}^{-1}$ and $c_2 = (1 - c)^3$. Then, it follows from (2.1) that

$$\begin{aligned} \widehat{U}_\alpha \widehat{V}_0^{-1/2} &= U_\alpha V_0^{-1/2} + V_0^{-1/2} \left(c_1 \frac{D_1}{\sqrt{n}} - \frac{U_\alpha}{2V_0} c_2 \frac{D_1}{\sqrt{n}} \right) \\ &\quad + V_0^{-1/2} \left\{ c_1 \frac{D_2}{n} - \frac{U_\alpha}{2V_0} c_2 \frac{D_2}{n} + \frac{3U_\alpha}{8V_0^2} c_2^2 \frac{D_1^2}{n} - \frac{1}{2V_0} c_1 c_2 \frac{D_1^2}{n} \right\} + O_p(n^{-3/2}), \end{aligned}$$

which implies that

$$\mathbb{E}[\Phi(\widehat{U}_\alpha \widehat{V}_0^{-1/2})] = \Phi(U_\alpha V_0^{-1/2}) + \frac{1}{n} \phi(U_\alpha V_0^{-1/2}) K(\Delta(\mathbf{j})) + O(n^{-3/2}), \quad (2.8)$$

where

$$\begin{aligned} K(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) &= \frac{1}{V_0^{1/2}} c_1 \mathbb{E}[\sqrt{n}D_1 + D_2] - \frac{U_0}{2V_0^{3/2}} c_2 \mathbb{E}[\sqrt{n}D_1 + D_2] - \frac{U_\alpha}{2V_0^{3/2}} c_1^2 \mathbb{E}[D_1^2] \\ &\quad + \frac{U_\alpha}{8V_0^{5/2}} \left(3 - \frac{U_\alpha^2}{V_0} \right) c_2^2 \mathbb{E}[D_1^2] - \frac{1}{2V_0^{3/2}} \left(1 - \frac{U_\alpha^2}{V_0} \right) c_1 c_2 \mathbb{E}[D_1^2]. \end{aligned}$$

Combining (2.3) and (2.8), we can see that the approximation of $E[L_1(\mathbf{j})]$ is expressed as

$$\begin{aligned} E[L_1(\mathbf{j})] &= \Phi(U_\alpha V_0^{-1/2}) + \frac{1}{n} \phi(U_\alpha V_0^{-1/2}) H(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) + O(n^{-3/2}) \\ &= E[\Phi(\widehat{U}_\alpha \widehat{V}_0^{-1/2})] + \frac{1}{n} \phi(U_\alpha V_0^{-1/2}) \{H(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) - K(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c)\} + O(n^{-3/2}). \end{aligned}$$

To calculate the moments in $K(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c)$, note that $E[D_1] = 0$,

$$\begin{aligned} E[D_2] &= \frac{\Delta(\mathbf{j})^2}{1-c} + \frac{c}{(1-c)\gamma_1\gamma_2}, \\ E[D_1^2] &= \frac{2\Delta(\mathbf{j})^4}{1-c} + \frac{4\Delta(\mathbf{j})^2}{(1-c)\gamma_1\gamma_2} + \frac{2c}{(1-c)\gamma_1^2\gamma_2^2}. \end{aligned}$$

Replacing the unknown parameter $\Delta(\mathbf{j})$ in $H(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) - K(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c)$ with the consistent estimator $\widehat{\Delta}(\mathbf{j})$, we obtain the second order asymptotically unbiased estimator of $\gamma_1(\mathbf{j})$ given by

$$\widehat{R}_1(\mathbf{j}) = \Phi(\widehat{U}_\alpha \widehat{V}_0^{-1/2}) + \frac{1}{n} \phi(\widehat{U}_\alpha \widehat{V}_0^{-1/2}) \{H(\widehat{\Delta}(\mathbf{j}), \gamma_1, \gamma_2, c) - K(\widehat{\Delta}(\mathbf{j}), \gamma_1, \gamma_2, c)\}. \quad (2.9)$$

Since the second order asymptotically unbiased estimator of $R_2(\mathbf{j})$ can be obtained by interchanging γ_1 with γ_2 in (2.9), we get

$$\widehat{R}_2(\mathbf{j}) = \Phi(\widehat{U}_\alpha \widehat{V}_0^{-1/2}) + \frac{1}{n} \phi(\widehat{U}_\alpha \widehat{V}_0^{-1/2}) \{H(\widehat{\Delta}(\mathbf{j}), \gamma_2, \gamma_1, c) - K(\widehat{\Delta}(\mathbf{j}), \gamma_1, \gamma_2, c)\}. \quad (2.10)$$

Combining (2.9) and (2.10), we obtain the second order asymptotically unbiased estimator of $R(\mathbf{j})$ given by

$$MEC = \widehat{R}(\mathbf{j}) = \frac{1}{2} \left(\widehat{R}_1(\mathbf{j}) + \widehat{R}_2(\mathbf{j}) \right), \quad (2.11)$$

which can be used for variable selection. We here call $\widehat{R}(\mathbf{j})$ the *Misclassification Error Criterion* (MEC). The best subset of variables among \mathcal{J} is suggested as \mathbf{j} which minimizes $\widehat{R}(\mathbf{j})$ among \mathcal{J} .

3 Asymptotic Properties of MEC

3.1 Decomposition into fitting and penalty terms

A feature of variable selection procedures like AIC and C_p is that they are decomposed into the two terms: ‘fitting’ term and ‘penalty’ term for model complexity. It is interesting to investigate whether MEC has a similar feature.

It is noted that $\Phi(-D(\mathbf{j})/2)$ decreases as the cardinality of variable set $\#\mathbf{j}$ increases, namely, $\Phi(-D(\mathbf{j}_1)/2) > \Phi(-D(\mathbf{j}_2)/2)$ if $\#\mathbf{j}_2 > \#\mathbf{j}_1$. This means that $\Phi(-D(\mathbf{j})/2)$ is regarded as a fitting term. Thus, MEC is decomposed as

$$MEC(\mathbf{j}) = \Phi(-D(\mathbf{j})/2) + \widehat{b}(\mathbf{j}), \quad (3.1)$$

where $\hat{b}(\mathbf{j}) = \hat{R}(\mathbf{j}) - \Phi(-D(\mathbf{j})/2)$. Let $b(\mathbf{j}) = R(\mathbf{j}) - \mathbb{E}[\Phi(-D(\mathbf{j})/2)]$. Then it is noted that $\hat{b}(\mathbf{j})$ is a second order asymptotically unbiased estimator of $b(\mathbf{j})$, namely $\mathbb{E}[\hat{b}(\mathbf{j})] = R(\mathbf{j}) - \mathbb{E}[\Phi(-D(\mathbf{j})/2)] + O(n^{-3/2})$.

In this subsection, we express the bias $b(\mathbf{j})$ with an explicit formula, and show that $b(\mathbf{j})$ works as dimensionality penalty under the high dimensional or large sample situation. We evaluate $\mathbb{E}[\Phi(-D(\mathbf{j})/2)]$ since the explicit formula of $R(\mathbf{j})$ is provided in subsection 2.1. Using Lemma 5.1, we can approximate $D(\mathbf{j})$ stochastically as

$$D(\mathbf{j})^2 = E_0 + \frac{E_1}{\sqrt{n}} + \frac{E_2}{n}, \quad (3.2)$$

where

$$\begin{aligned} E_0 &= \frac{1}{(1-c)} \left(\Delta(\mathbf{j})^2 + \frac{c}{\gamma_1 \gamma_2} \right), \\ E_1 &= -\frac{\sqrt{2}v_2(c + \Delta(\mathbf{j})^2\gamma_1\gamma_2)}{(1-c)^{3/2}\gamma_1\gamma_2} + \frac{\sqrt{2}\sqrt{c}v_1}{\gamma_1\gamma_2(1-c)} + \frac{2\Delta(\mathbf{j})u_1}{\sqrt{\gamma_1\gamma_2}(1-c)}, \\ E_2 &= -\frac{1}{(1-c)^2\gamma_1\gamma_2} - \frac{\Delta(\mathbf{j})^2}{(1-c)^2} + \frac{2\left(\frac{c}{\gamma_1\gamma_2} + \Delta(\mathbf{j})^2\right)}{(1-c)^2}v_2^2 + \frac{1}{\gamma_1\gamma_2(1-c)}u_1^2 \\ &\quad - \frac{2\sqrt{2}\Delta(\mathbf{j})}{(1-c)^{3/2}\sqrt{\gamma_1\gamma_2}}v_2u_1 - \frac{2\sqrt{c}}{(1-c)^{3/2}\gamma_1\gamma_2}v_1v_2 \end{aligned}$$

Using the stochastic expansion (3.2), we can see that

$$\mathbb{E} \left[\Phi \left(-\frac{D(\mathbf{j})}{2} \right) \right] = \Phi(m_0) + \frac{1}{n}\phi(m_0)\mathbb{E} \left[m_1 + m_2 - \frac{1}{2}m_0m_1^2 \right] + O(n^{-3/2}),$$

where

$$m_0 = -\frac{1}{2}E_0^{1/2}, \quad m_1 = -\frac{1}{4}E_0^{-1/2}\frac{E_1}{\sqrt{n}}, \quad m_2 = \frac{1}{16}E_0^{-3/2}\frac{E_1^2}{n} - \frac{1}{4}E_0^{-1/2}\frac{E_2}{n}.$$

We can obtain the expectations of m_1 , m_2 and m_1^2 using the moments $\mathbb{E}[E_1] = 0$,

$$\begin{aligned} \mathbb{E}[E_2] &= \frac{(3-2c)c + \Delta(\mathbf{j})^2\gamma_1\gamma_2}{(1-c)^2\gamma_1\gamma_2}, \\ \mathbb{E}[E_1^2] &= \frac{2(c + \Delta(\mathbf{j})^2\gamma_1\gamma_2)(\Delta(\mathbf{j})^2\gamma_1\gamma_2 + 2)}{(1-c)^3\gamma_1^2\gamma_2^2}. \end{aligned}$$

Thus, it is observed that

$$\mathbb{E} \left[\Phi \left(-\frac{D(\mathbf{j})}{2} \right) \right] = \Phi(m_0) + \frac{1}{n}\phi(m_0)M(\Delta(\mathbf{j})^2, \gamma_1, \gamma_2, c) + O(n^{-3/2}),$$

where

$$M(\Delta(\mathbf{j})^2, \gamma_1, \gamma_2, c) = \frac{1}{16}E_0^{-3/2}\mathbb{E}[E_1^2] - \frac{1}{4}E_0^{-1/2}\mathbb{E}[E_2] + \frac{1}{64}E_0^{-1/2}\mathbb{E}[E_1^2].$$

Then, we can get a second-order approximation of the bias term $b(\mathbf{j})$. Define $b^{(g_1, g_2)}(\mathbf{j})$ by

$$b^{(g_1, g_2)}(\mathbf{j}) = \Phi(\Delta_1(\mathbf{j})) - \Phi(\Delta_2(\mathbf{j})) + \frac{1}{n} \{ \phi(\Delta_1(\mathbf{j})) H(\Delta(\mathbf{j}), r_{g_1}, r_{g_2}, c) - \phi(\Delta_2(\mathbf{j})) M(\Delta(\mathbf{j}), \gamma_1, \gamma_2, c) \},$$

where

$$\Delta_1(\mathbf{j}) = -\frac{\sqrt{1-c}\Delta^2(\mathbf{j})}{2\sqrt{\Delta^2(\mathbf{j}) + c/(\gamma_1\gamma_2)}},$$

$$\Delta_2(\mathbf{j}) = -\frac{1}{2\sqrt{1-c}}\sqrt{\Delta^2(\mathbf{j}) + c/(\gamma_1\gamma_2)}.$$

Theorem 3.1 *Under the assumptions (A1)-(A3),*

$$b(\mathbf{j}) = \frac{1}{2} (b^{(1,2)}(\mathbf{j}) + b^{(2,1)}(\mathbf{j})) + O(n^{-3/2}).$$

From Theorem 3.1, an expression of the bias $b(\mathbf{j})$ implies that $b(\mathbf{j})$ depends on $\sharp(\mathbf{j})$ through $c = \lim_{n,p \rightarrow \infty} \sharp(\mathbf{j})/n$. Conceding the limiting term of $b(\mathbf{j})$, namely, $\Phi(\Delta_1(\mathbf{j})) - \Phi(\Delta_2(\mathbf{j}))$, it can be seen that it increases in $\sharp(\mathbf{j})$ through c for $\mathbf{j} \in \{\mathbf{j} | \Delta^2(\mathbf{j}) = \Delta^2\}$, since $\Delta_1(\mathbf{j})$ is increases in c , $\Delta_2(\mathbf{j})$ is decreases in c .

To make it clear that $b(\mathbf{j})$ works as a penalty of the cardinality $\sharp(\mathbf{j})$, we consider the large sample framework (A0). Then, it is observed

$$\begin{aligned} E[L_1(\mathbf{j})] &= \Phi\left(-\frac{\Delta(\mathbf{j})}{2}\right) + \frac{1}{n}\phi\left(-\frac{\Delta(\mathbf{j})}{2}\right) \left\{ \frac{1}{4\Delta(\mathbf{j})\gamma_1} \left(\frac{\Delta(\mathbf{j})^2}{4} + 3(\sharp(\mathbf{j}) - 1) \right) \right. \\ &\quad \left. + \frac{1}{4\Delta(\mathbf{j})\gamma_2} \left(\frac{\Delta(\mathbf{j})^2}{4} - \sharp(\mathbf{j}) + 1 \right) + \frac{1}{4}\Delta(\mathbf{j})(\sharp(\mathbf{j}) - 1) \right\} + o(n^{-1}), \\ E\left[\Phi\left(-\frac{D(\mathbf{j})}{2}\right)\right] &= \Phi\left(-\frac{\Delta(\mathbf{j})}{2}\right) + \frac{1}{n}\phi\left(-\frac{\Delta(\mathbf{j})}{2}\right) \left(\frac{\Delta(\mathbf{j})^3}{32} - \frac{\sharp(\mathbf{j})}{4\Delta(\mathbf{j})\gamma_1\gamma_2} - \frac{\Delta(\mathbf{j})\sharp(\mathbf{j})}{4} \right. \\ &\quad \left. + \frac{\Delta(\mathbf{j})}{16\gamma_1\gamma_2} + \frac{1}{4\Delta(\mathbf{j})\gamma_1\gamma_2} - \frac{\Delta(\mathbf{j})}{8} \right) + o(n^{-1}). \end{aligned}$$

Then, we get the following proposition.

Proposition 3.1 *Under the assumptions (A0), (A2) and (A3),*

$$b(\mathbf{j}) = \frac{1}{n}\phi\left(-\frac{\Delta(\mathbf{j})}{2}\right) \left\{ \frac{\sharp(\mathbf{j})}{\Delta(\mathbf{j})r_1} + \frac{\Delta(\mathbf{j})\sharp(\mathbf{j})}{2} - \frac{1}{\Delta(\mathbf{j})r_1} - \left(\frac{\Delta^3(\mathbf{j})}{32} + \frac{\Delta(\mathbf{j})}{8} \right) \right\} + o(n^{-1}).$$

Thus, the limiting term of $b(\mathbf{j})$, $\mathbf{j} \in \{\mathbf{j} | \Delta^2(\mathbf{j}) = \Delta^2\}$ increases as $\sharp(\mathbf{j})$ increases.

The arguments given in this subsection shows that MEC $\widehat{R}(\mathbf{j})$ is decomposed into the "fitting term" and the "penalty term", which is the same feature as in variable selection procedures like AIC and C_p .

3.2 Asymptotic optimality of MEC

In this subsection, we show that MEC has an asymptotic optimality in the high-dimensional setting (A1). The optimality is related to Li (1987), who showed that Mallows' C_p is asymptotically equivalent to the squared error loss in a linear regression model, and that the estimator selected by C_p asymptotically achieves the smallest possible squared error loss in the class of model average estimators. The squared error loss corresponds to the conditional misclassification error rate $L(\mathbf{j})$ in the classification problem, where

$$L(\mathbf{j}) = 2^{-1}\{L_1(\mathbf{j}) + L_2(\mathbf{j})\}.$$

Thus, we shall verify that $\widehat{R}(\mathbf{j})$ is asymptotically equivalent to the conditional misclassification error rate $L(\mathbf{j})$, and that MEC asymptotically achieves the smallest possible conditional misclassification error rate.

The primary goal of this section is to demonstrate that under reasonable conditions, MEC is asymptotically optimal in the sense that

$$\left| \frac{L(\widehat{\mathbf{j}})}{\inf_{\mathbf{j} \in \mathcal{J}} L(\mathbf{j})} \right| \xrightarrow{p} 1, \quad (3.3)$$

where $\widehat{\mathbf{j}}$ is the best selection satisfying $MEC(\widehat{\mathbf{j}}) = \min_{\mathbf{j} \in \mathcal{J}} MEC(\mathbf{j})$. It can be seen that sufficient conditions for (3.3) are given by

$$\begin{aligned} \text{(i)} \quad & \sup_{\mathbf{j} \in \mathcal{J}} \left| \frac{L(\mathbf{j})}{R(\mathbf{j})} - 1 \right| \xrightarrow{p} 0, \\ \text{(ii)} \quad & \sup_{\mathbf{j} \in \mathcal{J}} \left| \frac{\widehat{R}(\mathbf{j}) - L(\mathbf{j})}{R(\mathbf{j})} \right| \xrightarrow{p} 0. \end{aligned} \quad (3.4)$$

We shall check conditions (i) and (ii) using chebyshev's inequality. For any $\varepsilon > 0$,

$$\begin{aligned} & \Pr \left(\sup_{\mathbf{j} \in \mathcal{J}} \left| \frac{L(\mathbf{j})}{R(\mathbf{j})} - 1 \right| > \varepsilon \right) \\ & \leq \sum_{\mathbf{j} \in \mathcal{J}} \frac{\mathbb{E}[|L(\mathbf{j}) - R(\mathbf{j})|^{2m}]}{R(\mathbf{j})^{2m} \varepsilon^{2m}} \\ & \leq \sum_{\mathbf{j} \in \mathcal{J}} \left(\mathbb{E}[|L_1(\mathbf{j}) - R_1(\mathbf{j})|^{2m}]^{\frac{1}{2m}} + \mathbb{E}[|L_2(\mathbf{j}) - R_2(\mathbf{j})|^{2m}]^{\frac{1}{2m}} \right)^{2m} \left(\frac{1}{2R(\mathbf{j})\varepsilon} \right)^{2m} \\ & \leq \sup_{\mathbf{j} \in \mathcal{J}} \left(\mathbb{E}[|L_1(\mathbf{j}) - R_1(\mathbf{j})|^{2m}]^{\frac{1}{2m}} + \mathbb{E}[|L_2(\mathbf{j}) - R_2(\mathbf{j})|^{2m}]^{\frac{1}{2m}} \right)^{2m} \sum_{\mathbf{j} \in \mathcal{J}} \left(\frac{1}{2R(\mathbf{j})\varepsilon} \right)^{2m} \end{aligned}$$

and

$$\begin{aligned}
& \Pr \left(\sup_{\mathbf{j} \in \mathcal{J}} \left| \frac{\widehat{R}(\mathbf{j}) - L(\mathbf{j})}{R(\mathbf{j})} - 1 \right| > \varepsilon \right) \\
& \leq \sum_{\mathbf{j} \in \mathcal{J}} \frac{\mathbb{E}[|\widehat{R}(\mathbf{j}) - L(\mathbf{j})|^{2m}]}{R(\mathbf{j})^{2m} \varepsilon^{2m}} \\
& \leq \sum_{\mathbf{j} \in \mathcal{J}} \left(\mathbb{E}[|\widehat{R}_1(\mathbf{j}) - L_1(\mathbf{j})|^{2m}]^{\frac{1}{2m}} + \mathbb{E}[|\widehat{R}_2(\mathbf{j}) - L_2(\mathbf{j})|^{2m}]^{\frac{1}{2m}} \right)^{2m} \left(\frac{1}{2R(\mathbf{j})\varepsilon} \right)^{2m} \\
& \leq \sup_{\mathbf{j} \in \mathcal{J}} \left(\mathbb{E}[|\widehat{R}_1(\mathbf{j}) - L_1(\mathbf{j})|^{2m}]^{\frac{1}{2m}} + \mathbb{E}[|\widehat{R}_2(\mathbf{j}) - L_2(\mathbf{j})|^{2m}]^{\frac{1}{2m}} \right)^{2m} \sum_{\mathbf{j} \in \mathcal{J}} \left(\frac{1}{2R(\mathbf{j})\varepsilon} \right)^{2m}.
\end{aligned}$$

Thus, we need to evaluate $\mathbb{E}[|L_i(\mathbf{j}) - R_i(\mathbf{j})|^{2m}]$ and $\mathbb{E}[|\widehat{R}_i(\mathbf{j}) - L_i(\mathbf{j})|^{2m}]$ for $i = 1, 2$.

Lemma 3.1 *Assume either condition (A0) or condition (A1). Under conditions (A2) and (A3), it holds that*

$$\mathbb{E}[|\sqrt{n}(L_i(\mathbf{j}) - R_i(\mathbf{j}))|^{2m}] < \infty, \quad \mathbb{E}[|\sqrt{n}(\widehat{R}_i(\mathbf{j}) - L_i(\mathbf{j}))|^{2m}] < \infty.$$

(Proof) Using stochastic expansions of U , V and $\widehat{\Delta}(\mathbf{j})$, it follows that

$$\begin{aligned}
\sqrt{n}(L_1(\mathbf{j}) - R_1(\mathbf{j})) &= \phi(\Delta_1(\mathbf{j})) \left(\frac{U_1}{\sqrt{V_0}} - \frac{U_\alpha V_1}{2V_0^{3/2}} \right) + o_p(1), \\
\sqrt{n}(\widehat{R}_i(\mathbf{j}) - R_i(\mathbf{j})) &= -\phi(\Delta_1(\mathbf{j})) \frac{1}{V_0} \left(\frac{1}{2(1-c)} + \frac{U_\alpha}{V_0(1-c)^3} \right) D_1 + o_p(1).
\end{aligned}$$

We begin by expanding the statistic $\sqrt{n}(L_1(\mathbf{j}) - R_1(\mathbf{j}))$ stochastically. Using (5.1) and (5.2), we can expand the statistic $\sqrt{n}(L_1(\mathbf{j}) - R_1(\mathbf{j}))$ as

$$\begin{aligned}
\sqrt{n}(L_1(\mathbf{j}) - R_1(\mathbf{j})) &= \phi(\Delta_1(\mathbf{j})) \left(\sqrt{\frac{(1-c)c \{c(4\gamma_2 - 2) + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2 (4\gamma_2 - 1)\}}{\gamma_1 \gamma_2}} \frac{v_1}{2\sqrt{2}(c + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2)^{3/2}} \right. \\
&\quad + \sqrt{\frac{1}{\gamma_1 \gamma_2}} \frac{c(1 - 2\gamma_2)}{\sqrt{2}\sqrt{c + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2}} v_2 + \frac{\Delta(\mathbf{j})^2 \sqrt{(1-c)c\gamma_1 \gamma_2}}{2\sqrt{2}\sqrt{c + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2}} v_3 \\
&\quad - \frac{c\Delta(\mathbf{j})^2 \sqrt{\gamma_1 \gamma_2}}{2\sqrt{2}\sqrt{c + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2}} v_4 + \frac{\Delta(\mathbf{j}) \{2c(\gamma_2 - 1) + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2 (2\gamma_2 - 1)\}}{2(1-c)^{-1/2} (c + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2)^{3/2}} u_1 \\
&\quad \left. - \frac{\sqrt{c}\Delta(\mathbf{j})\gamma_2}{\sqrt{c + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2}} u_2 + \sqrt{1 - cu_3} - \sqrt{cu_4} \right) + o_p(1). \tag{3.5}
\end{aligned}$$

It should be noted that the above statistic is a linear combination of independent random variables u_i and v_i with finite $2m$ -th moment under the assumptions (A1)-(A3). Thus, for any $m \in \mathbb{R}$, $\mathbb{E}[|\sqrt{n}(L_i(\mathbf{j}) - R_i(\mathbf{j}))|^{2m}] < \infty$ under the assumptions (A1)-(A3).

We here treat the case that p is a fixed constant and p/n tends to zero, i.e. (A0). Substitute $c = 0$ in (3.5), we obtain a stochastic expansion of the statistics $\sqrt{n}(L_1(\mathbf{j}) - R_1(\mathbf{j}))$ given by

$$\sqrt{n}(L_1(\mathbf{j}) - R_1(\mathbf{j})) = \frac{\gamma_2 - \gamma_1}{2\sqrt{\gamma_1\gamma_2}}u_1 + u_3 + o_p(1).$$

This shows that the above statistic is a linear combination of independent random variables u_1 and u_3 with finite $2m$ -th moment under the large sample framework (A0). Thus, for any $m \in \mathbb{R}$, $E[|\sqrt{n}(L_i(\mathbf{j}) - R_i(\mathbf{j}))|^{2m}] < \infty$ under the large sample framework (A0).

We next evaluate the statistic $\sqrt{n}(\widehat{R}_i(\mathbf{j}) - R_i(\mathbf{j}))$, which can be carried out similarly. Using (5.1) and (5.2), we can expand the statistic $\sqrt{n}(\widehat{R}_i(\mathbf{j}) - R_i(\mathbf{j}))$ as

$$\begin{aligned} \sqrt{n}(\widehat{R}_i(\mathbf{j}) - R_i(\mathbf{j})) &= -\phi(\Delta_1(\mathbf{j})) \frac{1}{V_0} \left(\frac{1}{2(1-c)} + \frac{U_\alpha}{V_0(1-c)^3} \right) \left(-\frac{\sqrt{2}(c + \Delta(\mathbf{j})^2\gamma_1\gamma_2)}{\sqrt{1-c\gamma_1\gamma_2}}v_2 \right. \\ &\quad \left. + \frac{\sqrt{2}\sqrt{c}}{\gamma_1\gamma_2}v_1 + \frac{2\Delta(\mathbf{j})}{\sqrt{\gamma_1\gamma_2}}u_1 \right) + o_p(1). \end{aligned} \quad (3.6)$$

It should be noted that above statistic is linear combination of independent random variables u_1 , v_1 and v_2 with finite $2m$ -th moment under the assumptions (A1)-(A3). Thus, for any $m \in \mathbb{R}$, $E[|\sqrt{n}(\widehat{R}_i(\mathbf{j}) - R_i(\mathbf{j}))|^{2m}] < \infty$ under the assumptions (A1)-(A3). In addition, consider the case of the large sample framework (A0). Substituting $c = 0$ in (3.6), we obtain the stochastic expansion of the statistic $\sqrt{n}(\widehat{R}_i(\mathbf{j}) - R_i(\mathbf{j}))$ given by

$$\sqrt{n}(\widehat{R}_i(\mathbf{j}) - R_i(\mathbf{j})) = \frac{\gamma_2 - \gamma_1}{2\sqrt{\gamma_1\gamma_2}}u_1 + u_3 + o_p(1).$$

the above statistic is a linear combination of independent random variables u_1 and u_3 with finite $2m$ -th moment under the large sample framework (A0). Thus, for any $m \in \mathbb{R}$, $E[|\sqrt{n}(\widehat{R}_i(\mathbf{j}) - R_i(\mathbf{j}))|^{2m}] < \infty$ under the large sample framework (A0). Therefore, the proof of Lemma 3.1 is complete. ■

Combining Lemma 3.1 and (3.4) gives the sufficient condition for the asymptotic optimality (3.3), given by

$$(C1) \quad \sum_{\mathbf{j} \in \mathcal{J}} \left(\frac{1}{2\sqrt{n}R(\mathbf{j})} \right)^{2m} \rightarrow 0.$$

Hence, we obtain the following theorem.

Theorem 3.2 *Assume either condition (A0) or condition (A1). Under conditions (A2), (A3) and (C1), MEC is asymptotically optimal in the sense that*

$$\left| \frac{L(\widehat{\mathbf{j}})}{\inf_{\mathbf{j} \in \mathcal{J}} L(\mathbf{j})} \right| \xrightarrow{p} 1.$$

It is noted that the condition (C1) holds if the cardinality of the family \mathcal{J} satisfies $\#\mathcal{J} = o(n^m)$. For example, for the family of the subsets $\mathcal{J} = \{\{1\}, \{1, 2\}, \{1, 2, 3\}, \dots, \{1, 2, 3, \dots, p\}\}$, we have $\#\mathcal{J} = p$, so that the condition $\#\mathcal{J} = o(n^m)$ is satisfied under (A1). However, when \mathcal{J} consists of all the subsets of $\{1, \dots, p\}$, we have $\#\mathcal{J} = 2^p$, which does not satisfy the condition. Thus, Theorem 3.2 implies that we should pay heed to cardinality of the family \mathcal{J} of candidate variable sets in the high-dimensional setting (A1).

4 Numerical Results

In this section, we investigate numerical properties of MEC by Monte Carlo simulations. The frequencies of selecting the true variables by selection procedures are calculated by simulation with 10,000 iterations.

The mean vector of Π_1 is set up by $\boldsymbol{\mu} = ((1/2)\mathbf{1}'_{10}, \mathbf{0}'_{p-10})'$, the mean of Π_2 is set up by zero vector and the covariance matrix of both groups is set up by I_p since the expected error rate $R(\mathbf{j})$ only depends on Mahalanobis distance $\Delta(\mathbf{j})$ and (n, p) . The data sets are generated as follows:

$$\mathbf{x}_{11}, \dots, \mathbf{x}_{1N_1} \sim \mathcal{N}(\boldsymbol{\mu}, I_p), \quad \mathbf{x}_{21}, \dots, \mathbf{x}_{2N_2} \sim \mathcal{N}(\mathbf{0}, I_p)$$

in each step. In our numerical study, we consider the ten candidate variable sets \mathbf{j}_i given by

$$\begin{aligned} \mathbf{j}_i &= \{1, 2, \dots, (10 - i + 1)\} \text{ for } i = 1, \dots, 5, \\ \mathbf{j}_i &= \{1, 2, \dots, (p - i + 6)\} \text{ for } i = 6, \dots, 10. \end{aligned}$$

In this setting, the variable set which minimizes $R(\mathbf{j})$ is \mathbf{j}_1 in all the cases as seen from Table 1.

Table 1. Values of $R(\mathbf{j})$.

(p, N_1, N_2)	$R(\mathbf{j}_1)$	$R(\mathbf{j}_2)$	$R(\mathbf{j}_3)$	$R(\mathbf{j}_4)$	$R(\mathbf{j}_5)$	$R(\mathbf{j}_6)$	$R(\mathbf{j}_7)$	$R(\mathbf{j}_8)$	$R(\mathbf{j}_9)$	$R(\mathbf{j}_{10})$
(50, 50, 50)	0.15	0.17	0.18	0.19	0.21	0.25	0.25	0.25	0.24	0.24
(50, 100, 100)	0.14	0.15	0.17	0.18	0.20	0.19	0.19	0.19	0.18	0.18
(100, 100, 100)	0.14	0.15	0.17	0.18	0.20	0.25	0.25	0.25	0.25	0.25
(100, 200, 200)	0.14	0.15	0.163	0.18	0.20	0.19	0.19	0.19	0.19	0.19

We investigate the frequencies of selecting the true variable set minimizing $R(\mathbf{j})$ with several variable selection procedures. In this experiment, we compare the performances of AIC, BIC, Mc and MEC, where Mc is the large sample unbiased estimator suggested by MacLachlan (1976,1980). The AIC, BIC and Mc for the model \mathbf{j} are defined by

$$\begin{aligned} AIC(\mathbf{j}) &= -N \log \frac{1 + (N_1 N_2)/(Nn)D(\mathbf{j})^2}{1 + (N_1 N_2)/(Nn)D^2} + N \log |N^{-1}W| + Np(1 + \log 2\pi) + b_A, \\ BIC(\mathbf{j}) &= -N \log \frac{1 + (N_1 N_2)/(Nn)D(\mathbf{j})^2}{1 + (N_1 N_2)/(Nn)D^2} + N \log |N^{-1}W| + Np(1 + \log 2\pi) + b_B, \\ Mc(\mathbf{j}) &= \Phi \left(-\frac{D(\mathbf{j})}{2} \right) + \phi \left(-\frac{D(\mathbf{j})}{2} \right) \left(\frac{\#\mathcal{J}(\mathbf{j}) - 1}{N_1 D(\mathbf{j})} + 4(4\#\mathcal{J}(\mathbf{j}) - 1)D(\mathbf{j}) - \frac{D(\mathbf{j})^3}{32n} \right), \end{aligned}$$

where $D^2 = (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)'S^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$, $b_A = 2(2\sharp(\mathbf{j}) + p - \sharp(\mathbf{j}) + p(p+1)/2)$ and $b_B = \log N(2\sharp(\mathbf{j}) + p - \sharp(\mathbf{j}) + p(p+1)/2)$ for $N = N_1 + N_2$.

Tables 2-5 reports frequencies of selecting each variable sets \mathbf{j}_i by the variable selection procedures AIC, BIC, Mc and MEC. As seen from Tables 2-5, MEC is superior to the other criteria. Concerning AIC and MEC, their frequencies of selecting the true variable set get larger as the dimension p and sample size $N(= N_1 + N_2)$ are larger. Tables 2 and 4 treat the case of $p/N = 0.5$, which means that N is relatively small. In this case, the frequencies of selecting \mathbf{j}_1 by Mc is small. However, those values for Mc in Tables 3 and 5 are large. This arises from the reason that Mc is the variable selection criterion derived under the large sample framework (A0). However, MEC is excellent in all the cases, because it is derived under the high-dimensional setting.

Table 2. Comparison of frequencies of selecting \mathbf{j}_i for $p = 50$, $N_1 = N_2 = 50$.

\mathbf{j}_i	\mathbf{j}_1	\mathbf{j}_2	\mathbf{j}_3	\mathbf{j}_4	\mathbf{j}_5	\mathbf{j}_6	\mathbf{j}_7	\mathbf{j}_8	\mathbf{j}_9	\mathbf{j}_{10}
AIC	72.2	13.7	3.2	0.7	0.2	2.6	1.7	1.6	1.4	2.7
BIC	50.1	22.5	13.3	8.0	6.1	0.0	0.0	0.0	0.0	0.0
Mc	12.6	2.2	0.6	0.1	0.0	41.1	18.0	10.2	6.7	8.5
MEC	79.6	15.3	3.5	0.8	0.2	0.1	0.1	0.1	0.1	0.2

Table 3. Comparison of frequencies of selecting \mathbf{j}_i for $p = 50$, $N_1 = N_2 = 100$.

\mathbf{j}_i	\mathbf{j}_1	\mathbf{j}_2	\mathbf{j}_3	\mathbf{j}_4	\mathbf{j}_5	\mathbf{j}_6	\mathbf{j}_7	\mathbf{j}_8	\mathbf{j}_9	\mathbf{j}_{10}
AIC	96.6	2.3	0.1	0.0	0.1	0.2	0.0	0.1	0.2	0.4
BIC	85.1	12.1	2.3	0.4	0.1	0.0	0.0	0.0	0.0	0.0
Mc	84.0	1.9	0.0	0.0	0.0	3.4	2.5	2.3	2.0	3.9
MEC	97.1	2.5	0.1	0.0	0.0	0.0	0.0	0.0	0.1	0.2

Table 4. Comparison of frequencies of selecting \mathbf{j}_i for $p = 100$, $N_1 = N_2 = 100$.

\mathbf{j}_i	\mathbf{j}_1	\mathbf{j}_2	\mathbf{j}_3	\mathbf{j}_4	\mathbf{j}_5	\mathbf{j}_6	\mathbf{j}_7	\mathbf{j}_8	\mathbf{j}_9	\mathbf{j}_{10}
AIC	96.2	2.2	0.1	0.0	0.0	0.5	0.3	0.2	0.2	0.3
BIC	84.2	12.5	2.6	0.6	0.1	0.0	0.0	0.0	0.0	0.0
Mc	13.3	0.4	0.0	0.0	0.0	42.0	18.7	10.6	6.9	8.1
MEC	97.5	2.4	0.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Table 5. Comparison of frequencies of selecting \mathbf{j}_i for $p = 100$, $N_1 = N_2 = 200$.

\mathbf{j}_i	\mathbf{j}_1	\mathbf{j}_2	\mathbf{j}_3	\mathbf{j}_4	\mathbf{j}_5	\mathbf{j}_6	\mathbf{j}_7	\mathbf{j}_8	\mathbf{j}_9	\mathbf{j}_{10}
AIC	100	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
BIC	99.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Mc	96.9	0.1	0.0	0.0	0.0	0.6	0.6	0.5	0.4	0.9
MEC	100	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

We next investigate the relationship of $MEC(\mathbf{j})$ and the conditional error probability $L(\mathbf{j})$ by comparing the probabilities of selecting the true variable sets. Table 6 reports the frequencies of selecting the true variable sets with $MEC(\mathbf{j})$ and $L(\mathbf{j})$. This shows confirms that both $L(\mathbf{j})$ and $MEC(\mathbf{j})$ select the true variable sets with high frequencies, and this observation is related to the optimality (3.3).

Table 6. Comparison of frequencies of selecting the true variable sets

(p, N_1, N_2)	(50, 50, 50)	(50, 100, 100)	(100, 100, 100)	(100, 200, 200)
$L(\mathbf{j})$	88.8	97.7	97.7	100.0
$MEC(\mathbf{j})$	79.7	97.1	97.5	100.0

Finally, we check unbiasedness of MEC and Mc. For each subset \mathbf{j}_i , $i = 1, \dots, 10$, we compute the averages of $L(\mathbf{j})$, MEC and Mc by Monte Carlo simulations. In Figures 1-4, the averages of $L(\mathbf{j})$, $MEC(\mathbf{j})$ and $Mc(\mathbf{j})$ for each \mathbf{j} are plotted as “ \circ ”, “ \bullet ” and “ \times ”, respectively. As seen from the figures, MEC and Mc perform well for the subsets \mathbf{j}_i , $i = 1, \dots, 5$, while Mc is poor for the subsets \mathbf{j}_i , $i = 6, \dots, 10$. It is noted that the dimensions of the subsets \mathbf{j}_i , $i = 1, \dots, 5$ are not so large, but the dimensions of \mathbf{j}_i , $i = 6, \dots, 10$ are large. This is why Mc was the variable selection procedure derived in the large sample framework. MEC has good accuracies in both cases. Thus we conclude that MEC is more flexible than Mc concerning the sample size and dimension.

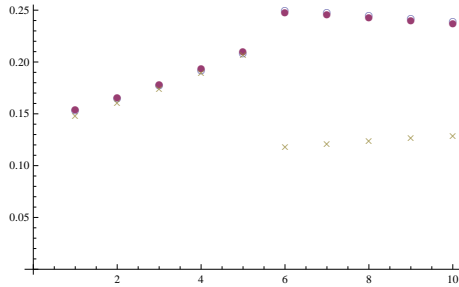


Fig.1. $(p, N_1, N_2) = (50, 50, 50)$.

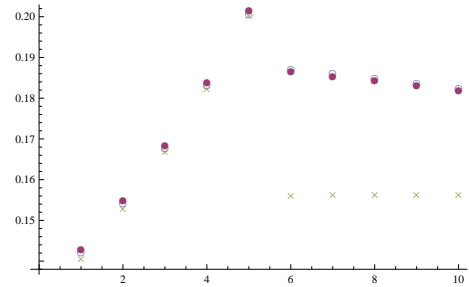


Fig.2. $(p, N_1, N_2) = (50, 100, 100)$.

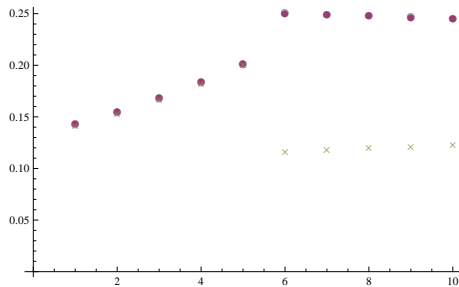


Fig.3. $(p, N_1, N_2) = (100, 100, 100)$.

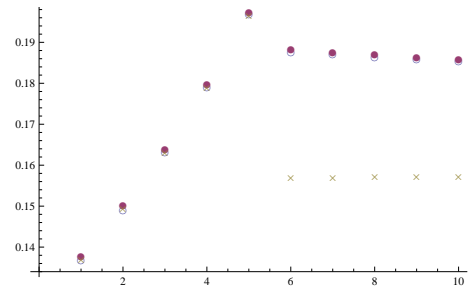


Fig.4. $(p, N_1, N_2) = (100, 200, 200)$.

5 Concluding Remarks

In this paper, we have suggested the new variable selection criterion MEC in linear discriminant analysis for high dimensional data. This is derived as a second-order unbiased estimator of the misclassification error probability. We have confirmed that MEC is decomposed into the “fitting term” and the “dimensionality penalty term” like AIC and Mallows’ C_p . Moreover, we have shown that MEC is asymptotically optimal in the sense of achieving the smallest possible conditional probability of misclassification in candidate variable sets. Also, the superiority of MEC has been verified in the sense of selecting the true variable sets by simulation.

It may be important to point out that the optimality (3.3) given in Theorem 3.2 is guaranteed by condition (C1), or $\sharp(\mathcal{J}) = o(n^m)$. This condition always holds in the large sample framework (A0), but it is not necessarily satisfied in the high-dimensional situation (A1). For example, for the family of the subsets $\mathcal{J} = \{\{1\}, \{1, 2\}, \{1, 2, 3\}, \dots, \{1, 2, 3, \dots, p\}\}$, we have $\sharp(\mathcal{J}) = p$, so that the condition $\sharp(\mathcal{J}) = o(n^m)$ is satisfied under (A1). However, when \mathcal{J} consists of all the subsets of $\{1, \dots, p\}$, we have $\sharp(\mathcal{J}) = 2^p$, which does not satisfy the condition. Thus, Theorem 3.2 implies that we should pay heed to cardinality of the family \mathcal{J} of candidate variable sets in the high-dimensional setting (A1).

6 Appendix

In this section, we derive the stochastic expansions for the statistics U and V . For the purpose, we prepare some stochastic expressions of required quadratic forms.

6.1 Stochastic expression of quadratic forms

Lemma 6.1 *Let $\mathbf{z} \sim \mathcal{N}_p(\tau \mathbf{e}_1, I_p)$, $\mathbf{g} \sim \mathcal{N}_p(\mathbf{0}, I_p)$ and $W \sim \mathcal{W}_p(n, I_p)$. Then, it holds that*

$$\begin{aligned} \text{(i)} \quad \mathbf{z}'W^{-1}\mathbf{z} &= \frac{(u_1 + \tau)^2 + \tilde{v}_1}{\tilde{v}_2}, \\ \text{(ii)} \quad \mathbf{z}'W^{-2}\mathbf{z} &= \frac{(u_1 + \tau)^2 + \tilde{v}_1}{\tilde{v}_2^2} \left(1 + \frac{\tilde{v}_3}{\tilde{v}_4} \right), \\ \text{(iii)} \quad \tau \mathbf{e}_1'W^{-1}\mathbf{z} &= \frac{\tau}{\tilde{v}_2} \left\{ u_1 + \tau + \left(\frac{\tilde{v}_1 \tilde{v}_3}{\tilde{v}_4} \right)^{\frac{1}{2}} \frac{u_2}{\sqrt{\tilde{v}_5 + u_2^2}} \right\}, \\ \text{(iv)} \quad \mathbf{z}'W^{-1}\mathbf{g} &= \frac{\{(u_1 + \tau)^2 + \tilde{v}_1\}^{\frac{1}{2}}}{\tilde{v}_2} \left\{ u_3 - u_4 \left(\frac{\tilde{v}_3}{\tilde{v}_4} \right)^{\frac{1}{2}} \right\}, \end{aligned}$$

where, $\mathbf{e}_1 = (1, 0, 0, \dots, 0)$. Here, $u_i \sim \mathcal{N}(0, 1)$, $i = 1, 2, 3, 4$, $\tilde{v}_1 \sim \chi_{p-1}^2$, $\tilde{v}_2 \sim \chi_{N-p-1}^2$, $\tilde{v}_3 \sim \chi_{p-1}^2$, $\tilde{v}_4 \sim \chi_{N-p}^2$ and $\tilde{v}_5 \sim \chi_{p-2}^2$ and these variables are mutually independent.

(Proof) Fujikoshi (2000) has derived the stochastic expressions (i)-(iii). In this proof, we prove the stochastic expression (iv). Let P be the random orthogonal matrix with its

first column proportional to \mathbf{z} i.e. $PP' = I_p$. Then $PWP' \sim \mathcal{W}_p(n, I_p)$ and

$$(PWP')^{-1} = \begin{pmatrix} \{\mathbf{y}'(I - \Pi_Y)\mathbf{y}\}^{-1} & * \\ -\{\mathbf{y}'(I - \Pi_Y)\mathbf{y}\}^{-1}(Y'Y)^{-1}Y'\mathbf{y} & Y'(I - \Pi_Y)Y \end{pmatrix},$$

where $\Pi_Y = Y(Y'Y)^{-1}Y'$ and $\Pi_y = \mathbf{y}(\mathbf{y}'\mathbf{y})^{-1}\mathbf{y}'$. At first, we obtain the stochastic expressions (i)-(iii) by using Fujikoshi (2000). The parts (i)-(iii) can be expressed as

$$\begin{aligned} \mathbf{z}'W^{-1}\mathbf{z} &= \frac{\mathbf{z}'\mathbf{z}}{\mathbf{y}'(I - \Pi_Y)\mathbf{y}}, \\ \mathbf{z}'W^{-2}\mathbf{z} &= \frac{\mathbf{z}'\mathbf{z}}{(\mathbf{y}'(I - \Pi_Y)\mathbf{y})^2} \left(1 + \frac{\mathbf{y}'Y(Y'Y)^{-2}Y'\mathbf{y}}{\mathbf{y}'\Pi_Y\mathbf{y}} \mathbf{y}'\Pi_Y\mathbf{y} \right), \\ \tau\mathbf{e}'_1W^{-1}\mathbf{z} &= \frac{\tau}{\mathbf{y}'(I - \Pi_Y)\mathbf{y}} \left\{ \tau + \mathbf{e}'_1\mathbf{u} + (-\mathbf{z}'_2\mathbf{z}_2)^{1/2}P'_2\mathbf{e}_1 \right\}' (Y'Y)^{-1}Y'\mathbf{y}. \end{aligned}$$

From Fujikoshi (2002), we have

$$\frac{(-\mathbf{z}'_2\mathbf{z}_2)^{1/2}P'_2\mathbf{e}_1 \right\}' (Y'Y)^{-1}Y'\mathbf{y}}{(\mathbf{z}'_2\mathbf{z}_2)\mathbf{e}'_1P_2P'_2\mathbf{e}_1\mathbf{y}'Y(Y'Y)^{-2}Y'\mathbf{y}} = \frac{u_2}{\sqrt{u_2^2 + v_5}},$$

where

$$\begin{aligned} u_1 &= \mathbf{e}'_1\mathbf{u} \sim \mathcal{N}(0, 1), \quad u_2 \sim \mathcal{N}(0, 1), \quad \tilde{v}_1 = \mathbf{u}'(I_{p-1} - (\mathbf{w}'\mathbf{w})^{-1}\mathbf{w}\mathbf{w}')\mathbf{u} \sim \chi_{p-1}^2, \\ \tilde{v}_2 &= \mathbf{y}'(I - \Pi_Y)\mathbf{y} \sim \chi_{N-p-1}^2, \quad \tilde{v}_3 = \mathbf{y}'\Pi_Y\mathbf{y} \sim \chi_{p-1}^2, \\ \tilde{v}_4 &= \mathbf{y}'\Pi_Y\mathbf{y}(\mathbf{y}'Y(Y'Y)^{-2}Y'\mathbf{y})^{-1} \sim \chi_{N-p}^2, \quad \tilde{v}_5 = \mathbf{y}'\Pi_Y\mathbf{y}(\mathbf{y}'Y(Y'Y)^{-2}Y'\mathbf{y})^{-1} \sim \chi_{p-2}^2. \end{aligned}$$

Next, we prove (iv). The part (iv) denotes

$$\begin{aligned} \mathbf{z}'W^{-1}\mathbf{g} &= \frac{(\mathbf{z}'\mathbf{z})^{1/2}}{\mathbf{y}'(I - \Pi_Y)\mathbf{y}} \left(u_3 - \frac{\mathbf{y}'Y(Y'Y)^{-1}\mathbf{g}_2}{(\mathbf{y}'Y(Y'Y)^{-2}Y'\mathbf{y})^{1/2}} \right. \\ &\quad \left. \times \left(\frac{\mathbf{y}'Y(Y'Y)^{-2}Y'\mathbf{y}}{\mathbf{y}'Y(Y'Y)^{-1}Y'\mathbf{y}} \right)^{1/2} (\mathbf{y}'Y(Y'Y)^{-1}Y'\mathbf{y})^{1/2} \right), \end{aligned}$$

where $\mathbf{g} = (u_3, \mathbf{g}'_2)'$. On the other hand, we note that

$$u_3 = \mathbf{e}'_1\mathbf{g} \sim \mathcal{N}(0, 1), \quad u_4 = (\mathbf{y}'Y(Y'Y)^{-2}Y'\mathbf{y})^{-1/2}\mathbf{y}'Y(Y'Y)^{-1}\mathbf{g}_2 \sim \mathcal{N}(0, 1).$$

Then, the variables u_i , $i = 1, 2, 3, 4$ and v_i , $i = 1, 2, 3, 4, 5$ are mutually independent. ■

6.2 Stochastic expansions of U and V

To expand U and V stochastically, define random variables \mathbf{z}_1 and \mathbf{z}_2 by

$$\begin{aligned} \mathbf{z}_1 &= \sqrt{\frac{N_1N_2}{N}}\Sigma(\mathbf{j})^{-1/2}(\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j})), \\ \mathbf{z}_2 &= \frac{1}{\sqrt{N}}\Sigma(\mathbf{j})^{-1/2}(N_1\bar{\mathbf{x}}_1(\mathbf{j}) + N_2\bar{\mathbf{x}}_2(\mathbf{j}) - N_1\boldsymbol{\mu}_1(\mathbf{j}) - N_2\boldsymbol{\mu}_2(\mathbf{j})), \\ W &= n\Sigma(\mathbf{j})^{-1/2}S(\mathbf{j})\Sigma(\mathbf{j})^{-1/2}. \end{aligned}$$

It is seen that \mathbf{z}_1 , \mathbf{z}_2 and W are mutually independently distributed as $\mathbf{z}_1 \sim \mathcal{N}_p(\boldsymbol{\tau}(\mathbf{j}), I_{\#(\mathbf{j})})$, $\mathbf{z}_2 \sim \mathcal{N}_p(\mathbf{0}, I_{\#(\mathbf{j})})$ and $W \sim \mathcal{W}_p(n, I_{\#(\mathbf{j})})$, respectively, where $\boldsymbol{\tau}(\mathbf{j}) = \sqrt{(N_1 N_2)/N} \boldsymbol{\Sigma}(\mathbf{j})^{-1/2} (\boldsymbol{\mu}_1(\mathbf{j}) - \boldsymbol{\mu}_2(\mathbf{j}))$. Using these variables, we can rewrite U and V as

$$U = -\frac{(N_1 - N_2)n}{2N_1 N_2} \mathbf{z}'_1 W^{-1} \mathbf{z}_1 + \frac{n}{\sqrt{N_1 N_2}} \mathbf{z}'_1 W^{-1} \mathbf{z}_2 - \frac{n}{N_1} \boldsymbol{\tau}(\mathbf{j})' W^{-1} \mathbf{z}_1,$$

$$V = \frac{n^2 N}{N_1 N_2} \mathbf{z}'_1 W^{-2} \mathbf{z}_1.$$

Let Γ be an orthogonal matrix with its first column proportional to $\boldsymbol{\tau}(\mathbf{j})$ i.e. $\Gamma \Gamma' = I_p$. Then $\Gamma W \Gamma' \sim \mathcal{W}_p(n, I_p)$, $\Gamma \mathbf{z}_1 \sim \mathcal{N}_p(\tau \mathbf{e}_1, I_{\#(\mathbf{j})})$, $\Gamma \mathbf{z}_2 \sim \mathcal{N}_p(\mathbf{0}, I_{\#(\mathbf{j})})$ and $W \sim \mathcal{W}_p(n, I_{\#(\mathbf{j})})$, respectively, where $\tau^2 = (N_1 N_2)/N \Delta(\mathbf{j})^2$. For simplicity, we denote the statistics $\Gamma W \Gamma'$, $\Gamma \mathbf{z}_1$ and $\Gamma \mathbf{z}_2$ by W , \mathbf{z}_1 and \mathbf{z}_2 , respectively. We can rewrite the statistics U and V as

$$U = -\frac{(N_1 - N_2)n}{2N_1 N_2} \mathbf{z}'_1 W^{-1} \mathbf{z}_1 + \frac{n}{\sqrt{N_1 N_2}} \mathbf{z}'_1 W^{-1} \mathbf{z}_2 - \frac{n}{N_1} \tau \mathbf{e}'_1 W^{-1} \mathbf{z}_1,$$

$$V = \frac{n^2 N}{N_1 N_2} \mathbf{z}'_1 W^{-2} \mathbf{z}_1.$$

Using Lemma 6.1, we express U and V as

$$U = -\frac{(N_1 - N_2)n}{2N_1 N_2} \frac{(u_1 + \tau)^2 + \tilde{v}_1}{\tilde{v}_2} + \frac{n}{\sqrt{N_1 N_2}} \frac{\{(u_1 + \tau)^2 + \tilde{v}_1\}^{\frac{1}{2}}}{\tilde{v}_2} \left\{ u_3 - u_4 \left(\frac{\tilde{v}_3}{\tilde{v}_4} \right)^{\frac{1}{2}} \right\}$$

$$- \frac{n}{N_1} \frac{\tau}{\tilde{v}_2} \left\{ u_1 + \tau + \left(\frac{\tilde{v}_1 \tilde{v}_3}{\tilde{v}_4} \right)^{\frac{1}{2}} \frac{u_2}{\sqrt{\tilde{v}_5 + u_2^2}} \right\},$$

$$V = \frac{n^2 N}{N_1 N_2} \frac{(u_1 + \tau)^2 + \tilde{v}_1}{\tilde{v}_2^2} \left(1 + \frac{\tilde{v}_3}{\tilde{v}_4} \right).$$

Define variables v_1 , v_2 , v_3 , v_4 and v_5 by

$$v_1 = \frac{\tilde{v}_1 - (p-1)}{\sqrt{2(p-1)}}, \quad v_2 = \frac{\tilde{v}_2 - (N-p-1)}{\sqrt{2(N-p-1)}}, \quad v_3 = \frac{\tilde{v}_3 - (p-1)}{\sqrt{2(p-1)}}, \quad v_4 = \frac{\tilde{v}_4 - (N-p)}{\sqrt{2(N-p)}},$$

$$v_5 = \frac{\tilde{v}_5 - (p-2)}{\sqrt{2(p-2)}}.$$

Note that v_i is asymptotically distributed as $\mathcal{N}(0, 1)$ under condition (A1). Using Taylor series expansion based on these variables, we can expand U stochastically as

$$U = U_0 + \frac{1}{\sqrt{n}} U_1 + \frac{1}{n} U_2, \quad (6.1)$$

where

$$\begin{aligned}
U_0 &= -\frac{1}{2(1-c)} \left(\Delta(\mathbf{j})^2 + \frac{c(\gamma_1 - \gamma_2)}{\gamma_1 \gamma_2} \right), \\
U_1 &= \frac{\sqrt{\gamma_1 \gamma_2 \left(\frac{c}{\gamma_1 \gamma_2} + \Delta(\mathbf{j})^2 \right)}}{(1-c)\sqrt{\gamma_1 \gamma_2}} u_3 - \frac{\sqrt{c} \sqrt{\gamma_1 \gamma_2 \left(\frac{c}{\gamma_1 \gamma_2} + \Delta(\mathbf{j})^2 \right)}}{(1-c)^{3/2} \sqrt{\gamma_1 \gamma_2}} u_4 \\
&\quad + \left(\frac{\Delta(\mathbf{j})^2}{\sqrt{2}(1-c)^{3/2}} + \frac{c}{\sqrt{2}(1-c)^{3/2} \gamma_1 \gamma_2} - \frac{\sqrt{2}c}{(1-c)^{3/2} \gamma_1} \right) v_2 \\
&\quad + \left(-\frac{\Delta(\mathbf{j}) \gamma_2 \sqrt{\frac{1}{\gamma_1 \gamma_2}}}{1-c} + \frac{2\Delta(\mathbf{j}) \sqrt{\gamma_1 \gamma_2}}{(1-c) \gamma_1} - \frac{\Delta(\mathbf{j}) \sqrt{\gamma_1 \gamma_2}}{(1-c) \gamma_1 \gamma_2} \right) u_1 \\
&\quad - \frac{\sqrt{c} \Delta(\mathbf{j}) \gamma_2 \sqrt{\frac{1}{\gamma_1 \gamma_2}}}{(1-c)^{3/2}} u_2 + \left(\frac{\sqrt{2} \sqrt{c}}{(1-c) \gamma_1} - \frac{\sqrt{c}}{\sqrt{2}(1-c) \gamma_1 \gamma_2} \right) v_1, \\
U_2 &= \left(\frac{\sqrt{c} \left(u_3 - \frac{\sqrt{c} u_4}{\sqrt{1-c}} \right)}{\sqrt{2}(1-c) \gamma_1 \gamma_2 \sqrt{\frac{c}{\gamma_1 \gamma_2} + \Delta(\mathbf{j})^2}} - \frac{\Delta(\mathbf{j}) u_2 \sqrt{\frac{\gamma_2}{\gamma_1}}}{\sqrt{2}(1-c)^{3/2}} \right. \\
&\quad \left. + \frac{\sqrt{c} v_2}{(1-c)^{3/2} \gamma_1 \gamma_2} - \frac{2\sqrt{c} v_2}{(1-c)^{3/2} \gamma_1} \right) v_1 \\
&\quad + \left(-\frac{\sqrt{2} \left(u_3 - \frac{\sqrt{c} u_4}{\sqrt{1-c}} \right) \sqrt{\frac{c}{\gamma_1 \gamma_2} + \Delta(\mathbf{j})^2}}{(1-c)^{3/2}} - \frac{\sqrt{2} \Delta(\mathbf{j}) u_1 \sqrt{\frac{\gamma_2}{\gamma_1}}}{(1-c)^{3/2}} \right. \\
&\quad \left. + \frac{\sqrt{2} \Delta(\mathbf{j}) u_1}{(1-c)^{3/2} \sqrt{\gamma_1 \gamma_2}} + \frac{\sqrt{2} \sqrt{c} \Delta(\mathbf{j}) u_2 \sqrt{\frac{\gamma_2}{\gamma_1}}}{(1-c)^2} \right) v_2 \\
&\quad + \frac{\Delta(\mathbf{j}) \left(\sqrt{1-c} u_3 - \sqrt{c} u_4 \right) u_1}{(1-c)^{3/2} \sqrt{c + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2}} - \frac{\left(u_4 \sqrt{c + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2} + \Delta(\mathbf{j}) \gamma_2 u_2 \right) v_3}{\sqrt{2-2c(1-c)} \sqrt{\gamma_1 \gamma_2}} \\
&\quad + \frac{\sqrt{c} \left(u_4 \sqrt{c + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2} + \Delta(\mathbf{j}) \gamma_2 u_2 \right) v_4}{\sqrt{2}(1-c)^2 \sqrt{\gamma_1 \gamma_2}} - \frac{v_2^2 (-2c\gamma_2 + c + \Delta(\mathbf{j})^2 \gamma_1 \gamma_2)}{(c-1)^2 \gamma_1 \gamma_2} \\
&\quad + \frac{\Delta(\mathbf{j}) v_5 u_2 \sqrt{\frac{\gamma_2}{\gamma_1}}}{\sqrt{2}(1-c)^{3/2}} - \frac{(1-2\gamma_2) u_1^2}{2(1-c) \gamma_1 \gamma_2} + \frac{1-2\gamma_2}{2\gamma_1 \gamma_2 - 2c\gamma_1 \gamma_2},
\end{aligned}$$

Here, u_i , $i = 1, \dots, 4$ and v_j , $j = 1, \dots, 5$ are mutually independent and they are asymptotically distributed as $\mathcal{N}(0, 1)$.

Using similar arguments, we can expand V stochastically as

$$V = V_0 + \frac{V_1}{\sqrt{n}} + \frac{V_2}{n}, \quad (6.2)$$

where

$$\begin{aligned}
V_0 &= \frac{1}{(1-c)^3} \left(\frac{c}{\gamma_1 \gamma_2} + \Delta(\mathbf{j})^2 \right), \\
V_1 &= \frac{\frac{\sqrt{2}\sqrt{c}v_1}{\gamma_1 \gamma_2} + 2u_1 \sqrt{\frac{\Delta(\mathbf{j})^2}{\gamma_1 \gamma_2}}}{(1-c)^3} + \frac{c \left(\frac{c}{\gamma_1 \gamma_2} + \Delta(\mathbf{j})^2 \right) \left(-\frac{2\sqrt{2}v_2}{\sqrt{1-c}} + \frac{\sqrt{2}v_3}{\sqrt{c}} - \frac{\sqrt{2}v_4}{\sqrt{1-c}} \right)}{(1-c)^3}, \\
V_2 &= \frac{1}{(1-c)^3} v_2^2 \left(\frac{6\Delta(\mathbf{j})^2}{1-c} + \frac{6c}{(1-c)\gamma_1 \gamma_2} \right) \\
&\quad - \frac{1}{(1-c)^3} \left(\frac{4\sqrt{2}\Delta(\mathbf{j})v_2 u_1 \sqrt{\frac{1}{\gamma_1 \gamma_2}}}{\sqrt{1-c}} + \frac{4\sqrt{c}v_1 v_2}{\sqrt{1-c}\gamma_1 \gamma_2} - \frac{u_1^2}{\gamma_1 \gamma_2} + \frac{1}{\gamma_1 \gamma_2} \right) \\
&\quad + c \left(\frac{\sqrt{2}v_3}{(1-c)^3 \sqrt{c}} - \frac{\sqrt{2}v_4}{(1-c)^{7/2}} \right) \left(\frac{\sqrt{2}\sqrt{c}v_1}{\gamma_1 \gamma_2} - \frac{2\sqrt{2}v_2 \left(\frac{c}{\gamma_1 \gamma_2} + \Delta(\mathbf{j})^2 \right)}{\sqrt{1-c}} \right) \\
&\quad + 2u_1 \sqrt{\frac{\Delta(\mathbf{j})^2}{\gamma_1 \gamma_2}} + \frac{c}{(1-c)^3} \left(\frac{2v_4^2}{1-c} - \frac{2v_3 v_4}{\sqrt{1-c}\sqrt{c}} \right) \left(\frac{c}{\gamma_1 \gamma_2} + \Delta(\mathbf{j})^2 \right).
\end{aligned}$$

It can be seen that the stochastic expansions of U and V under the large sample framework (A0) is provided by replacing c in the limiting terms U_0 and V_0 with p/n , and by replacing c in first and second terms U_1 , U_2 , V_1 and V_2 with 0.

Acknowledgments.

Research of the second author was supported in part by Grant-in-Aid for Scientific Research (21540114 and 23243039) from Japan Society for the Promotion of Science.

References

- [1] Fujikoshi, Y. (1985). Selection of variables in two-group discriminant analysis by error rate and Akaike's information criterion. *J. Multivariate. Anal.*, **17**, 27-37.
- [2] Fujikoshi, Y. (2002). Selection of variables for discriminant analysis in a high-dimensional case. *Sankhya Ser. A*, **64**, 256-257.
- [3] Fujikoshi, Y., Ulyanov, and Shimizu, R. (2010). *Multivariate Statistics: High-Dimensional and Large-Sample Approximations*. Wiley, Hoboken, N. J.
- [4] Fujikoshi, Y. and Seo, T. (1998). Asymptotic approximations for EPMC's of the linear and the quadratic discriminant function when the sample size and the dimension are large. *Random Operators and Stochastic Equations*, **6**, 269-280.
- [5] Kshirsagar, A. M. (1972). *Multivariate Analysis*. New York Marcel Dekker.
- [6] Kubokawa, T, Hyodo, M., and Srivastava, M.S. (2013). Asymptotic expansion and estimation of EPMC for linear classification rules in high dimension. *J. Multivariate Analysis*, to appear.

- [7] McLachlan, G. J. (1976). A criterion for selecting variables for the linear discriminant function. *Biometrics*, **32**, 529-515.
- [8] McLachlan, G. J. (1980). Selection of variables in discriminant analysis (Letter to the Editor). *Biometrics*, **36**, 554.
- [9] Rao, C. R. (1948). Tests of significance in multivariate analysis. *Biometrika*, **35**, 58-79.
- [10] Rao, C. R. (1973). *Linear Statistical Inference and Its Applications*. Second Edition. New York Wiley.
- [11] Raudys, S. (1972). On the amount of a priori information in designing the classification algorithm. *Technical Cybernetics*, **4**, 168-174.
- [12] Sakurai, T., Nakada, T., and Fujikoshi, Y. (2012). High-Dimensional AICs for Selection of Redundancy Models in Discriminant Analysis. Technical Report 12-13, Hiroshima Statistical Research Group, Hiroshima University.
- [13] Wilbur, J.D., Ghosh, J.K., Nakatsu, C.H., and Doerge, R.W. (2002). Variable selection in high-dimensional multivariate binary data with application to analysis of microbial community DNA finger prints. *Biometrics*. **58**, 378-386.
- [14] Wyman, F. J., Young, D. M. and Turner, D. W. (1990). A comparison of asymptotic error rate expansions for sample linear discriminant function. *Pattern Recognition*. **23**, 775-783.