Contents lists available at ScienceDirect

## Journal of Multivariate Analysis

journal homepage: www.elsevier.com/locate/jmva

## A variable selection criterion for linear discriminant rule and its optimality in high dimensional and large sample data

### Masashi Hyodo<sup>a,\*</sup>, Tatsuya Kubokawa<sup>b</sup>

<sup>a</sup> Department of Mathematical Information Science, Tokyo University of Science, 1-3, Kagurazaka, Shinjuku-ku, Tokyo 162-8601, Japan
<sup>b</sup> Faculty of Economics, University of Tokyo, 7-3-1 Hongo, Bunkyo-ku, Tokyo 113-0033, Japan

#### ARTICLE INFO

Article history: Received 11 December 2012 Available online 28 October 2013

AMS subject classifications: 62H30 62H12 62E20

Keywords: Asymptotic optimality High dimension and large sample Linear discriminant analysis Misclassification error Multivariate normal Second-order approximation Variable selection

#### ABSTRACT

In this paper, we suggest the new variable selection procedure, called MEC, for linear discriminant rule in the high dimensional and large sample setup. MEC is derived as a second-order unbiased estimator of the misclassification error probability of the linear discriminant rule (LDR). It is shown that MEC not only asymptotically 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.

© 2013 Elsevier Inc. All rights reserved.

#### 1. Introduction

In this paper, we consider the problem of classifying a future observation vector into one of the two population groups  $\Pi_1$  and  $\Pi_2$ . For each  $i = 1, 2, \Pi_i$  denotes a population from a multivariate normal distribution  $\mathcal{N}_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$ , and it is supposed that  $\mathbf{x}_{ij}$ ,  $j = 1, \ldots, N_i$ , are observed from the population  $\Pi_i$ . Here,  $\boldsymbol{\mu}_i$ , i = 1, 2, and  $\boldsymbol{\Sigma}$  are unknown parameters, and they are estimated by the sample mean  $\bar{\mathbf{x}}_i = N_i^{-1} \sum_{j=1}^{N_i} \mathbf{x}_{ij}$  i = 1, 2, and the pooled sample covariance matrix  $S = n^{-1} \sum_{i=1}^{2} \sum_{j=1}^{N_i} (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i) (\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)'$  for  $n = N - 2N = N_1 + N_2$ . Shao et al. [14] show that the misclassification error rate of the linear discriminant rule (LDR) is asymptotically close to ones of the optimal rule when p diverges to infinity at a rate slower than  $n^{1/2}$ . From this perspective, LDR works well for the situation where the number of variables used for classification is much smaller than the training sample size. Thus, it is desired to find an optimal subset of variables in the sense that misclassification error rate can be small. Variable selection methods for discriminant analysis have been studied by Fujikoshi [1,3], Sakurai et al. [13], Wilbur et al. [16] and others. Related to this issue, multiple testing problems for no additional information have been discussed by Rao [10,11] and Kshirsagar [5]. In this paper, we suggest a new variable selection procedure based on misclassification error rate and establish the optimality in high dimensional and large sample setting.

\* Corresponding author. E-mail addresses: hyodoh\_h@yahoo.co.jp (M. Hyodo), tatsuya@e.u-tokyo.ac.jp (T. Kubokawa).







<sup>0047-259</sup>X/\$ - see front matter © 2013 Elsevier Inc. All rights reserved. http://dx.doi.org/10.1016/j.jmva.2013.10.005

To explain the new variable selection procedure, consider the following linear discriminant rule. Let  $\mathbf{x} = (x_1, \ldots, x_p)'$  be a future observation in the full model. Let  $\mathbf{j} = (j_1, \ldots, j_{k(j)})'$  be a subset of the set  $\{1, 2, \ldots, p\}$ , and let  $\mathbf{x}(\mathbf{j}) = (x_{j_1}, \ldots, x_{j_{k(j)}})'$  be the corresponding sub vector of  $\mathbf{x}$ . The model based on the variable  $\mathbf{x}(\mathbf{j})$  is denoted by  $\mathbf{j}$ . Let  $\mathbf{\mathcal{J}}$  be a suitable family of subsets of  $\{1, \ldots, p\}$ . The LDR for classifying  $\mathbf{x}$  based on the model  $\mathbf{j}$  is that  $\mathbf{x}$  is classified as coming from  $\Pi_1$ , if  $W(\mathbf{j}) > \alpha$ , and from  $\Pi_2$ , if  $W(\mathbf{j}) < \alpha$ , where  $\alpha$  is a cut off point for classification rule, and

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

Here,  $\bar{\mathbf{x}}_i(\mathbf{j}) i = 1, 2, \text{ and } S(\mathbf{j})$  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  $L_1(\mathbf{j}) = P[W(\mathbf{j}) < \alpha | \mathbf{x}(\mathbf{j}) \in \Pi_1, \overline{\mathbf{x}}_1(\mathbf{j}), \overline{\mathbf{x}}_2(\mathbf{j}), S(\mathbf{j})]$  and  $L_2(\mathbf{j}) = P[W(\mathbf{j}) \geq \alpha | \mathbf{x}(\mathbf{j}) \in \Pi_2, \overline{\mathbf{x}}_1(\mathbf{j}), \overline{\mathbf{x}}_2(\mathbf{j}), S(\mathbf{j})]$  and  $L_2(\mathbf{j}) = P[W(\mathbf{j}) \geq \alpha | \mathbf{x}(\mathbf{j}) \in \Pi_2, \overline{\mathbf{x}}_1(\mathbf{j}), \overline{\mathbf{x}}_2(\mathbf{j}), S(\mathbf{j})]$ , which can be expressed as

$$L_{g}(\mathbf{j}) = \Phi\left( (-1)^{g} \frac{(\bar{\mathbf{x}}_{1}(\mathbf{j}) - \bar{\mathbf{x}}_{2}(\mathbf{j}))'S(\mathbf{j})^{-1} \{\mu_{g}(\mathbf{j}) - (\bar{\mathbf{x}}_{1}(\mathbf{j}) + \bar{\mathbf{x}}_{2}(\mathbf{j}))/2\} - \alpha}{\sqrt{(\bar{\mathbf{x}}_{1}(\mathbf{j}) - \bar{\mathbf{x}}_{2}(\mathbf{j}))'S(\mathbf{j})^{-1}\Sigma(\mathbf{j})S(\mathbf{j})^{-1}(\bar{\mathbf{x}}_{1}(\mathbf{j}) - \bar{\mathbf{x}}_{2}(\mathbf{j}))}} \right)$$
(1.1)

for g = 1, 2, where  $\Phi(\cdot)$  is the standard normal distribution function, and  $\mu_g(\mathbf{j})$  and  $\Sigma(\mathbf{j})$  denote the population mean and covariance matrix in the model  $\mathbf{j}$ . When  $\pi_i$ , 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  $R_g(j)$  is the unconditional error of misallocation given by  $R_g(j) = E[L_g(j)]$  for g = 1, 2. The variable selection procedure proposed in this paper is an asymptotically unbiased estimator of the misclassification error R(j) in the high dimensional and large sample setting.

A naive procedure for selection of variables is the method of minimizing

$$\Phi\left(-D(\mathbf{j})/2\right),\tag{1.2}$$

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

$$D(\mathbf{j})^2 = (\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j}))'S(\mathbf{j})^{-1}(\bar{\mathbf{x}}_1(\mathbf{j}) - \bar{\mathbf{x}}_2(\mathbf{j})).$$
(1.3)

However,  $\Phi(-D(\mathbf{j})/2)$  has the bias  $R(\mathbf{j}) - E[\Phi(-D(\mathbf{j})/2)]$  which is not negligible. McLachlan [8,9] derived a second order asymptotic unbiased estimator of  $R(\mathbf{j})$  under the large sample framework, namely,

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

Fujikoshi [1] applied the estimator given by McLachlan [8,9] to the variable selection problem, and investigated the asymptotic properties and the relationship with AIC. On the other hand, Raudys [12] and Wyman et al. [17] derived the asymptotic approximations of the error probability under the high dimensional and large sample setting given by

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

This setup not only includes the large sample setting (A0) as  $c_0 = 0$ , but also covers the case of large dimension p subject to p < n. This is practically important, because, as pointed out by Siotani [15], it is known that large sample approximations under (A0) are not good when p is large. In fact, as seen from Fujikoshi et al. [4] and Kubokawa et al. [6], the approximations under the setting (A1) give good approximations for large p less than n as well as small p.

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

$$MEC = \Phi \left( -D(\mathbf{j})/2 \right) + (\text{penalty}) + o_p(1),$$

where the penalty term increases in the dimension of model j. 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 optimality in the high dimensional and large sample setting is not known as long as we know.

Recently, Kubokawa et al. [6] derived a second-order approximation of the error probability of misclassification (EPMC) for the ridge-type linear discriminant rule in the high dimensional and large sample 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. [6] 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. These stochastic expressions are essentially derived by Fujikoshi [3]. 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 settings.

Download English Version:

# https://daneshyari.com/en/article/1145670

Download Persian Version:

https://daneshyari.com/article/1145670

Daneshyari.com