

## On Effect of Nonnormality on Size of Test for Dimensionality in Discriminant Analysis

Changha Hwang<sup>1)</sup>

### Abstract

In discriminant analysis the procedures commonly used to estimate the dimensionality involve testing a sequence of dimensionality hypotheses. There is a problem with the size of the test since dimensionality hypotheses are tested sequentially and thus they are actually conditional tests. The focus of this paper is to investigate in asymptotic sense what happens to the sequential testing procedure if the assumption of normality does not hold.

### 1. Introduction

In discriminant analysis, the study of dimensionality is quite interesting since it determines the number of discriminant functions required to describe group differences. The procedures commonly used to estimate this dimensionality involve testing a sequence of dimensionality hypotheses. These hypotheses are tested sequentially and thus they are actually conditional tests; that is, we test  $H_k$  after we have tested and rejected the hypotheses  $H_0, H_1, \dots, H_{k-1}$  in sequence. There is a problem with the size of the test since successive tests are not independent. Hwang(1995b) showed that the size of test is not affected asymptotically under the normality. The focus of this paper is to investigate in asymptotic sense "How is the size of the test affected under nonnormality by viewing this sequence of tests as conditional tests?".

### 2. Main Result

Let  $\mathbf{y}_{i1}, \dots, \mathbf{y}_{iq_i}$  ( $i=1, \dots, p$ ) be i.i.d.  $m \times 1$  absolutely continuous random vectors with mean  $\boldsymbol{\mu}_i$ , covariance matrix  $\boldsymbol{\Sigma}$  and finite fourth moments. Suppose that the samples are independent across populations. Let  $\bar{\mathbf{y}}_i$  be the sample mean of the  $q_i$  observations in the  $i$

---

1) Assistant Professor, Dept. of Statistics, Catholic University of Taegu-Hyosung, Kyungbuk, 713-702, Korea.

th sample and  $\bar{\mathbf{y}}$  be the sample mean of all  $n$  observations, ( $n = \sum_{i=1}^p q_i$ ). Then matrices  $A$  and  $B$  are defined as

$$A = \sum_{i=1}^p q_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}})(\bar{\mathbf{y}}_i - \bar{\mathbf{y}})' \text{ and } B = \sum_{i=1}^p \sum_{j=1}^{q_i} (\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)(\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)'.$$

The matrix  $\mathbf{Q}$  is defined as  $\mathbf{Q} = \Sigma^{-1} \sum_{i=1}^p q_i (\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})'$ , where  $\bar{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^p q_i \boldsymbol{\mu}_i$ . From now on, we will assume that  $p \geq m + 1$  so that  $AB^{-1}$  has  $m$  nonzero eigenvalues  $f_1 > \dots > f_m > 0$ . For the asymptotic theory there is no loss of generality in assuming that  $\mathbf{Q}$  is the diagonal matrix defined by  $\mathbf{Q} = \text{diag}\{w_1, \dots, w_m\}$ ,  $\mathbf{Q} = n_2 \boldsymbol{\Theta}$ , and  $\Sigma = I_m$  where  $n_2 = n - p$  and  $\boldsymbol{\Theta}$  is the fixed matrix defined by  $\boldsymbol{\Theta} = \text{diag}\{\theta_1, \dots, \theta_m\}$ . This means that we consider the case where  $A$ ,  $B$ ,  $\mathbf{Q}$ , and  $\Sigma$  are already transformed to canonical form. Thus, the dimensionality is the rank of  $\mathbf{Q}$ .

In practice, to determine the number of useful discriminant functions we test the sequence of dimensionality hypotheses,

$$H_k: \theta_{k+1} = \dots = \theta_m = 0 (\theta_k > 0), \quad k = 0, 1, \dots, m - 1.$$

By testing these hypotheses sequentially they are actually conditional tests. We test  $H_k$  given we have tested and rejected  $H_0, H_1, \dots, H_{k-1}$ , keeping in mind the effect on the significance level (the size of test). The likelihood ratio test statistic for  $H_k$  is given by

$$T_k = n_2 \sum_{i=k+1}^m \log(1 + f_i)$$

where  $f_1 > \dots > f_m > 0$  are the eigenvalues of  $AB^{-1}$ . For nonnormal populations, the asymptotic distribution of  $T_k$  is  $\chi^2_{(m-k)(n_1-k)}$  when  $H_k$  is true. See for details Hwang(1994).

For our purpose, we need the following asymptotic expansion of test statistic  $T_k$ :

$$T_k = n_2 \sum_{i=k+1}^m \log(1 + \theta_i) + \sqrt{n_2} C + D + O_p(n_2^{-\frac{1}{2}}),$$

where

$$\begin{aligned} C &= \sum_{i=k+1}^m \frac{E_{\hat{u}}(n) - \theta_i U_{\hat{u}}(n)}{1 + \theta_i}, \\ D &= \sum_{i=k+1}^m \frac{F_{\hat{u}}(n)}{(1 + \theta_i)} - \sum_{i=1}^k \sum_{j=k+1}^m \frac{E_{\hat{y}}(n)^2}{(1 + \theta_j)(\theta_i - \theta_j)} \\ &+ \sum_{i=1}^k \sum_{j=k+1}^m \frac{4\theta_j E_{\hat{y}}(n) U_{\hat{y}}(n)}{(1 + \theta_i)(1 + \theta_j)(\theta_i - \theta_j)} - \sum_{i=k+1}^m \sum_{j=k+1}^m \frac{E_{\hat{y}}(n) U_{\hat{y}}(n)}{(1 + \theta_i)(1 + \theta_j)} \end{aligned}$$

$$\begin{aligned}
& + \sum_{i=1}^k \sum_{j=k+1}^m \frac{\theta_j(\theta_j - 2\theta_i - \theta_i^2) U_{ij}(n)^2}{(1 + \theta_i)(1 + \theta_j)(\theta_i - \theta_j)} + \sum_{i=k+1}^m \sum_{j=k+1}^m \frac{\theta_i(\theta_j + 2) U_{ij}(n)^2}{2(1 + \theta_i)(1 + \theta_j)} \\
& - \sum_{i=k+1}^m \sum_{j=k+1}^m \frac{E_{ij}(n)^2}{2(1 + \theta_i)(1 + \theta_j)}.
\end{aligned}$$

Furthermore,  $E_{ij}(n)$ ,  $F_{ij}(n)$ , and  $U_{ij}(n)$  are the  $ij$ th element of matrices  $E(n)$ ,  $F(n)$ , and  $U(n)$  defined as follows:

$$\begin{aligned}
E(n) &= \frac{1}{\sqrt{n_2}} \sum_{i=1}^p q_i [(\bar{\boldsymbol{\varepsilon}}_i - \bar{\boldsymbol{\varepsilon}})(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})' + (\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})(\bar{\boldsymbol{\varepsilon}}_i - \bar{\boldsymbol{\varepsilon}})'], \\
F(n) &= \sum_{i=1}^p q_i (\bar{\boldsymbol{\varepsilon}}_i - \bar{\boldsymbol{\varepsilon}})(\bar{\boldsymbol{\varepsilon}}_i - \bar{\boldsymbol{\varepsilon}})', \\
U(n) &= \frac{1}{\sqrt{n_2}} \sum_{i=1}^p \sum_{j=1}^{q_i} [(\boldsymbol{\varepsilon}_{ij} - \bar{\boldsymbol{\varepsilon}}_i)(\boldsymbol{\varepsilon}_{ij} - \bar{\boldsymbol{\varepsilon}}_i)' - I_m],
\end{aligned}$$

where  $\bar{\boldsymbol{y}}_i = \boldsymbol{\mu}_i + \bar{\boldsymbol{\varepsilon}}_i$ ,  $\bar{\boldsymbol{y}} = \bar{\boldsymbol{\mu}} + \bar{\boldsymbol{\varepsilon}}$ ,  $\bar{\boldsymbol{\varepsilon}}_i = \frac{1}{q_i} \sum_{j=1}^{q_i} \boldsymbol{\varepsilon}_{ij}$ , and  $\bar{\boldsymbol{\varepsilon}} = \frac{1}{n} \sum_{i=1}^p q_i \bar{\boldsymbol{\varepsilon}}_i$ . See for details Hwang(1995a).

**Theorem 1** For each  $i = 1, \dots, p$ , let  $\boldsymbol{y}_{ij}: m \times 1, j = 1, \dots, q_i$  be a sequence of i.i.d. random vectors drawn from  $m$  multivariate elliptical populations with parameter  $\boldsymbol{\mu}_i$ , covariance matrix  $I_m$ , and finite fourth moments. Suppose that the  $p$  sequences are independent and put

$$\begin{aligned}
T_k &= n_2 \sum_{i=k+1}^m \log(1 + f_i) \\
V_k &= \frac{1}{\sqrt{n_2}} [T_k - n_2 \sum_{i=k+1}^m \log(1 + \theta_i)]
\end{aligned}$$

Then under  $H_k$ ,  $T_k$  is asymptotically independent of  $V_j, j = 0, 1, \dots, k-1$ .

**Proof** From the expansions of  $T_0, T_1, \dots, T_k$  under  $H_k$  we form two subvectors  $\boldsymbol{z}_1$  and  $\boldsymbol{z}_2$ , where  $\boldsymbol{z}_1$  contains the  $\sqrt{q_i} \bar{\boldsymbol{\varepsilon}}_{i \cdot s}$  variables which make up  $T_k$  and  $\boldsymbol{z}_2$  contains the  $\sqrt{q_i} \bar{\boldsymbol{\varepsilon}}_{i \cdot t}$  and  $\frac{1}{\sqrt{q_i}} \sum_{j=1}^{q_i} (\boldsymbol{\varepsilon}_{ij}^2 - 1)$  variables which make up  $V_0, V_1, \dots, V_{k-1}$ . Specifically,

$$\begin{aligned}
\boldsymbol{z}_1 &= (\sqrt{q_i} \bar{\boldsymbol{\varepsilon}}_{i \cdot k+1}, \dots, \sqrt{q_i} \bar{\boldsymbol{\varepsilon}}_{i \cdot m})' \\
\boldsymbol{z}_2 &= (\sqrt{q_i} \bar{\boldsymbol{\varepsilon}}_{i \cdot 1}, \dots, \sqrt{q_i} \bar{\boldsymbol{\varepsilon}}_{i \cdot k}; \frac{1}{\sqrt{q_i}} \sum_{j=1}^{q_i} (\boldsymbol{\varepsilon}_{ij}^2 - 1), \dots, \frac{1}{\sqrt{q_i}} \sum_{j=1}^{q_i} (\boldsymbol{\varepsilon}_{ij}^2 - 1))'
\end{aligned}$$

and define  $\boldsymbol{z} = (\boldsymbol{z}_1', \boldsymbol{z}_2')'$ . Here,  $\boldsymbol{\varepsilon}_{ij}$ ,  $\bar{\boldsymbol{\varepsilon}}_i$  and  $\bar{\boldsymbol{\varepsilon}}$  are denoted by

$$\boldsymbol{\varepsilon}_{ij} = (\boldsymbol{\varepsilon}_{ij1}, \dots, \boldsymbol{\varepsilon}_{ijm})', \quad \bar{\boldsymbol{\varepsilon}}_i = (\bar{\boldsymbol{\varepsilon}}_{i,1}, \dots, \bar{\boldsymbol{\varepsilon}}_{i,m})' \quad \text{and} \quad \bar{\boldsymbol{\varepsilon}} = (\bar{\boldsymbol{\varepsilon}}_{..1}, \dots, \bar{\boldsymbol{\varepsilon}}_{..m})',$$

where  $\bar{\varepsilon}_{i,r} = \frac{1}{q_i} \sum_{j=1}^{q_i} \varepsilon_{ijr}$  and  $\bar{\varepsilon}_{..r} = \frac{1}{n} \sum_{i=1}^p \sum_{j=1}^{q_i} \varepsilon_{ijr}$ ,  $r = 1, \dots, m$ . By the Multivariate Central Limit Theorem  $\mathbf{z}$  converges in distribution to multivariate normal with mean  $\mathbf{0}$ . The elements of asymptotic covariance matrix can be computed as follows: For  $s = k+1, \dots, m$ ,  $t = 1, \dots, k$ , as  $n \rightarrow \infty$ ,

$$\begin{aligned} \text{Cov}(\sqrt{q_i} \bar{\varepsilon}_{i,s}, \sqrt{q_i} \bar{\varepsilon}_{i,t}) &= E[(\sqrt{q_i} \bar{\varepsilon}_{i,s})(\sqrt{q_i} \bar{\varepsilon}_{i,t})] \\ &\rightarrow 0, \end{aligned}$$

since  $\sqrt{q_i} \bar{\varepsilon}_{i,s}$  and  $\sqrt{q_i} \bar{\varepsilon}_{i,t}$  are asymptotically independent. Also,

$$\begin{aligned} \text{Cov}(\sqrt{q_i} \bar{\varepsilon}_{i,s}, \frac{1}{\sqrt{q_i}} \sum_{j=1}^{q_i} (\varepsilon_{ijt}^2 - 1)) &= E[(\sqrt{q_i} \bar{\varepsilon}_{i,s})(\frac{1}{\sqrt{q_i}} \sum_{j=1}^{q_i} (\varepsilon_{ijt}^2 - 1))] \\ &= E[\varepsilon_{is} \varepsilon_{it}^2] = 0 \end{aligned} \quad (1)$$

because for elliptical distributions all third moments are zero (see, for example, Gang(1987)). These covariance expressions show that the elements of  $\mathbf{z}_1$  are asymptotically independent of the elements of  $\mathbf{z}_2$  for elliptical distributions. Thus, under  $H_k$ , the test statistic  $T_k$  is asymptotically independent of the statistics  $V_j$ ,  $j=0, \dots, k-1$ . ■

Therefore, this result agrees with the result for multivariate normal distributions. From (1) we see that the elements of  $\mathbf{z}_1$  are generally not asymptotically independent of the elements of  $\mathbf{z}_2$ . Thus, under  $H_k$ , the test statistic  $T_k$  is generally not asymptotically independent of the statistics  $V_j$ ,  $j=0, \dots, k-1$ . It is shown that this result is sensitive to certain departures from normality.

### 3. Simulation Study and Conclusion

A Monte Carlo experiment was carried out to see how inferences regarding the sequential testing procedure based on the assumption of multivariate normality are affected if this assumption is violated. In particular, suppose we are sampling from an elliptical t-distribution on 5 degrees of freedom. Recall that the asymptotic distribution of  $T_k = n_2 \sum_{i=k+1}^m \log(1 + f_i)$  is  $\chi^2_{(n_1-k)(m-k)}$  for nonnormal populations. The statistic  $T_k$  was used in the study. The study consisted of generating 500 samples of size  $n_2 = 50, 100, 200$  of an 4-variate elliptical t-distribution on 5 degrees of freedom for 6 populations with parameters  $\mu_i (i=1, \dots, 6)$  and  $V = \frac{3}{5} I_4$ . These samples can be generated using the following

relationship:

$$y_{ij} = \mu_i + Z^{-\frac{1}{2}} (4V)x,$$

where  $x \sim N_4(0, I_m)$  and  $Z \sim \chi_5^2$ . For further details of simulations see Hwang(1994a). Generation of the samples, computation of the sample eigenvalue and the analysis were conducted using SAS and SAS/IML.

Table 1 and 2 present the observed unconditional and conditional significance levels, respectively. As  $n_2$  increases, the observed conditional and unconditional significance levels become closer to each other. This result agrees with the result for multivariate normal distributions. To conclude, we see for multivariate elliptical populations the size of the test is not affected asymptotically by viewing this sequence of tests as conditional tests but this result is sensitive to certain departures from normality.

Table 1: Unconditional Significance Levels for elliptical  $t(5)$ ,  $(m=4, p=6)$

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$	$V$	$T^{.05}$			$T^{.10}$		
					$n_2 = 50$	$n_2 = 100$	$n_2 = 200$	$n_2 = 50$	$n_2 = 100$	$n_2 = 200$
0.2	0	0	0	1	0.020	0.038	0.060	0.052	0.086	0.116
0.8	0	0	0	1	0.038	0.040	0.056	0.090	0.110	0.114
6	0	0	0	1	0.042	0.044	0.052	0.094	0.108	0.108
0.4	0.2	0	0	2	0.020	0.052	0.062	0.050	0.100	0.114
0.8	0.4	0	0	2	0.026	0.062	0.068	0.084	0.116	0.118
6	2	0	0	2	0.042	0.062	0.062	0.110	0.118	0.124
0.4	0.2	0.1	0	3	0.004	0.004	0.006	0.012	0.014	0.032
2	1	0.8	0	3	0.042	0.052	0.052	0.090	0.120	0.118
6	4	2	0	3	0.060	0.060	0.066	0.124	0.126	0.114

Table 2: Conditional Significance Levels for elliptical  $t(5)$ ,  $(m=4, p=6)$

$\theta_1$	$\theta_2$	$\theta_3$	$\theta_4$	$V$	$T^{.05}$			$T^{.10}$		
					$n_2 = 50$	$n_2 = 100$	$n_2 = 200$	$n_2 = 50$	$n_2 = 100$	$n_2 = 200$
0.2	0	0	0	1	0.062	0.053	0.061	0.106	0.103	0.117
0.8	0	0	0	1	0.041	0.040	0.056	0.092	0.110	0.114
6	0	0	0	1	0.042	0.044	0.052	0.094	0.108	0.108
0.4	0.2	0	0	2	0.058	0.066	0.063	0.150	0.113	0.115
0.8	0.4	0	0	2	0.036	0.063	0.068	0.099	0.117	0.118
6	2	0	0	2	0.042	0.062	0.062	0.110	0.118	0.124
0.4	0.2	0.1	0	3	0.250	0.222	0.059	0.250	0.214	0.174
2	1	0.8	0	3	0.067	0.054	0.052	0.118	0.122	0.118
6	4	2	0	3	0.060	0.060	0.066	0.124	0.126	0.114

## References

- [1] Gang, L. (1987). Moments of a random vector and its quadratic forms. *J. Statist. Appl. Prob.*, 2, 219-229. [Reprinted in *Statistical Inference in Elliptically Contoured and Related Distributions*(Fang, K. and Anderson T.W., ed.), Allerton Press, 1990, 433-440.]
- [2] Hwang, C. (1991). Model Selection Methods in Discriminant Analysis. Ph.D Thesis, Univ. of Michigan, Ann Arbor, Michigan.
- [3] Hwang, C. (1994a). On estimating the dimensionality in discriminant analysis. *Communications in Statistics: Theory and Methods*, 23, 2197-2215.
- [4] Hwang, C. (1994b). Characterization of the asymptotic distributions of certain eigenvalues in a general setting. *Journal of the Korean Statistical Society*, 23, 13-32.
- [5] Hwang, C. (1995a). A note on the asymptotic distributions of dimensionality Estimators in discriminant analysis. *The Korean Communications in Statistics*, 2, 320-329.
- [6] Hwang, C. (1995b). Size of test for dimensionality in discriminant analysis. *Journal of Statistical Theory and Methods*, 6, 9-15.
- [7] Seo, T., Kanda, T. and Fujikoshi, Y. (1993). The effects of nonnormality on tests for dimensionality in canonical correlation and MANOVA models. Technical Report No. 93-9, Hiroshima University.
- [8] Siotani, M., Hayakawa, T., and Fujikoshi, Y. (1985). *Modern Multivariate Statistical Analysis: A Graduate Course and Handbook*. American Sciences Press, Inc.