# NONPARAMETRIC TEST FOR EIGENVALUES OF COVARIANCE MATRIX IN MULTIPOPULATION

Hidetoshi Murakami*, Emiko Hino* and Shin-ichi Tsukada**

We propose a nonparametric procedure to test the hypothesis that the $j$-th largest eigenvalues of a covariance matrix are equal in multipopulation. We apply the Mood test by using the principal component scores and deal the equality of eigenvalues with the equality of variance. We investigate the significance level and the power of test by simulation and show that this nonparametric test is useful for symmetric populations.

*Key words and phrases*: Eigenvalues, $k$-sample Mood test, nonparametric test, principal component score.

## 1. Introduction

Principal component analysis (PCA) is one of the most common and important methods in multivariate analysis, and many books on PCA have been published (Anderson (2003), Jackson (2003) and Jolliffe (2002)). Since it is difficult to obtain the exact distribution of eigenvalues of a covariance matrix under the nonnormal population, we have not seen testing of the hypothesis that the $j$-th largest eigenvalues are equal under multipopulation. For two populations, Sugiyama and Ushizawa (1998) proposed the nonparametric procedure which is the Ansari-Bradley test by using the principal component scores. In this paper, we extend the testing procedure under the multipopulation.

Suppose that $\boldsymbol{x}_1^{(i)}, \ldots, \boldsymbol{x}_{N_i}^{(i)}$ are the random observations from a $p$-dimensional distribution $\Lambda_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ with mean $\boldsymbol{\mu}_i$ and covariance matrix $\boldsymbol{\Sigma}_i$, where $i = 1, \ldots, k$. Let $\lambda_j^{(i)}$ be the $j$-th largest eigenvalue of covariance matrix $\boldsymbol{\Sigma}_i$ in the $i$-th population. For fixed $j$, we consider testing the hypothesis

$$H_0 : \lambda_j^{(1)} = \cdots = \lambda_j^{(k)}$$
$$H_1 : \text{not } H_0.$$

Let $\boldsymbol{h}_j^{(i)}$ be the sample eigenvector corresponding to the $j$-th largest sample eigenvalue $l_j^{(i)}$ of the $i$-th sample covariance matrix given by

$$\boldsymbol{S}^{(i)} = \frac{1}{N_i - 1} \sum_{\alpha=1}^{N_i} (\boldsymbol{x}_\alpha^{(i)} - \bar{\boldsymbol{x}}^{(i)})(\boldsymbol{x}_\alpha^{(i)} - \bar{\boldsymbol{x}}^{(i)})',$$

where $\bar{\boldsymbol{x}}^{(i)}$ is the $i$-th sample mean vector. As $\lambda_j^{(i)}$ is the asymptotic variance of principal components

$$y_{j\alpha}^{(i)} = \boldsymbol{h}'_j{}^{(i)}(\boldsymbol{x}_\alpha^{(i)} - \bar{\boldsymbol{x}}^{(i)}), \qquad \alpha = 1, \ldots, N_i,$$

we then apply the Mood test (1954) for equality of variance to test the hypothesis in Section 2. In Section 3, we investigate the significance level and the power of test by simulation.

## 2. Testing procedure

We deal with testing the equality of the $i$-th largest eigenvalues in the $k$-population using the principal component scores

$$\begin{aligned}
Y_1 &= \{y_{j1}^{(1)}, y_{j2}^{(1)}, \ldots, y_{jN_1}^{(1)}\}, \\
Y_2 &= \{y_{j1}^{(2)}, y_{j2}^{(2)}, \ldots, y_{jN_2}^{(2)}\}, \\
&\vdots \\
Y_k &= \{y_{j1}^{(k)}, y_{j2}^{(k)}, \ldots, y_{jN_k}^{(k)}\}.
\end{aligned}$$

The variance of principal component $y_{j\alpha}^{(i)}$ is as follows:

$$\mathrm{Var}[y_{j\alpha}^{(i)}] = \lambda_j^{(i)} - \frac{2}{N_i - 1} \sum_{q \neq j}^{p} \frac{m_{qj}^{(i)22}}{\lambda_q^{(i)} - \lambda_j^{(i)}} + O(N_i^{-2}),$$

where $m_{qj}^{(i)22} = E[x_q^{(i)2}x_j^{(i)2}] = E[x_{q\alpha}^{(i)2}x_{j\alpha}^{(i)2}]$. Therefore, the null hypothesis is equivalent with the equality for variance of the principal component when all eigenvalues are equal for $j = 1, \ldots, p$. If all eigenvalues except the eigenvalue of null hypothesis are not equal, the equality for variance of the principal component and the null hypothesis are not accurately equivalent, but are asymptotically equivalent. We may treat testing the null hypothesis as the equality for variances of the principal component in the case that the sample sizes $N_i$ are sufficiently large. In addition, we also need even larger sample sizes when the eigenvalues are close. Takeda (2001) treated this methodology under the multivariate contaminated normal distribution.

The Ansari-Bradley test is known as a method of testing the variance. One of the assumptions for the Ansari-Bradley test is that the sample values are independent. However, there exist weak correlations between each principal component scores. Sugiyama and Ushizawa (1998) proved that the degree of dependence between each principal component score was weak when the sample size was sufficiently large under the multivariate normal distribution. Then they showed that the Ansari-Bradley test could be applicable to test the equality for variance of $Y_1$ and $Y_2$ (cf. Ansari and Bradley (1960)). It is well known that the asymptotic relative efficiency of the Mood test is higher than that of the Ansari-Bradley test (Gibbons and Chakraborti (2003)). Therefore, we apply the Mood test for a $k$-population.

Let $R_{jm}^{(i)}$ be the increasing order rank of $y_{jm}^{(i)}$ in the combined $N = N_1 + \cdots + N_k$ observations. The statistic of the Mood test is as follows:

$$M_k = \frac{180}{N(N+1)(N^2-4)} \sum_{i=1}^{k} N_i \left( \bar{M}_j^{(i)} - \frac{N^2-1}{12} \right)^2,$$

where

$$\bar{M}_j^{(i)} = \frac{1}{N_i} \sum_{m=1}^{N_i} \left( R_{jm}^{(i)} - \frac{N+1}{2} \right)^2.$$

The limiting distribution of the Mood statistic, named $M_k$, for $k$-population is a $\chi^2$ distribution with $k-1$ degrees of freedom under the null hypothesis (Tsai *et al.* (1975)).

## 3.  Simulation study

In this section, we examine the power of tests for equality of the $j$-th eigenvalues, using a significance level of 5%. To compare the power of tests, we carry out simulations for multivariate normal populations and multivariate contaminated normal populations. We assume that the number of population is three and investigate the behavior of the $M_k$ statistic under the trivariate distribution. The simulation is repeated a million times in each case.

When $N$ is even, we give the Ansari-Bradley statistic, namely $AB_{ke}$, for $k$-population as follows:

$$AB_{ke} = \frac{48(N-1)}{N(N^2-4)} \sum_{i=1}^{k} N_i \left( \bar{A}_j^{(i)} - \frac{N+2}{4} \right)^2.$$

If $N$ is odd, we give the Ansari-Bradley statistic, namely $AB_{ko}$, as follows:

$$AB_{ko} = \frac{48N^2}{N(N+1)(N^2+3)} \sum_{i=1}^{k} N_i \left( \bar{A}_j^{(i)} - \frac{(N+1)^2}{4N} \right)^2.$$

Here, $\bar{A}_j^{(i)}$ denotes

$$\bar{A}_j^{(i)} = \frac{1}{N_i} \sum_{m=1}^{N_i} \left( \frac{N+1}{2} - \left| R_{jm}^{(i)} - \frac{N+1}{2} \right| \right).$$

The limiting distribution of the Ansari-Bradley statistic for the $k$-population is also the $\chi^2$ distribution with $k-1$ degrees of freedom (Tsai *et al.* (1975)). Therefore we set the critical value of the Ansari-Bradley statistic and the Mood statistic as 5.991 for $k = 3$. We simulate under the normal populations $N(\mathbf{0}, \boldsymbol{\Sigma}_i)$ and the contaminated normal populations $0.95 \times N(\mathbf{0}, \boldsymbol{\Sigma}_i) + 0.05 \times N(\mathbf{0}, 3\boldsymbol{\Sigma}_i)$ in the following cases. Cases 1 and 2 are the cases under the null hypothesis. Under the alternative hypothesis; Cases 3, 4 and 5, the variance of $Y_i$ is different to each other.

Case 1

$$\lambda_1^{(1)} = 6,\ \lambda_2^{(1)} = 3,\ \lambda_3^{(1)} = 1 \qquad \lambda_1^{(2)} = 6,\ \lambda_2^{(2)} = 3,\ \lambda_3^{(2)} = 1 \qquad \lambda_1^{(3)} = 6,\ \lambda_2^{(3)} = 3,\ \lambda_3^{(3)} = 1$$

$$\mathbf{\Sigma}_1 \qquad\qquad\qquad \mathbf{\Sigma}_2 \qquad\qquad\qquad \mathbf{\Sigma}_3$$

$$\begin{pmatrix} 6 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 1 \end{pmatrix} \qquad \begin{pmatrix} 6 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 1 \end{pmatrix} \qquad \begin{pmatrix} 6 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 1 \end{pmatrix}$$

Case 2

$$\lambda_1^{(1)} = 6,\ \lambda_2^{(1)} = 3,\ \lambda_3^{(1)} = 1 \qquad \lambda_1^{(2)} = 6,\ \lambda_2^{(2)} = 3,\ \lambda_3^{(2)} = 1 \qquad \lambda_1^{(3)} = 6,\ \lambda_2^{(3)} = 3,\ \lambda_3^{(3)} = 1$$

$$\mathbf{\Sigma}_1 \qquad\qquad\qquad \mathbf{\Sigma}_2 \qquad\qquad\qquad \mathbf{\Sigma}_3$$

$$\begin{pmatrix} 6 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 1 \end{pmatrix} \qquad \begin{pmatrix} 5.4776 & -0.7244 & -1.1998 \\ -0.7244 & 3.2010 & 0.1941 \\ -1.1998 & 0.1941 & 1.3215 \end{pmatrix} \qquad \begin{pmatrix} 2.75 & -1.9874 & -0.4874 \\ -1.9874 & 4.6856 & -0.875 \\ -0.4874 & -0.875 & 2.5643 \end{pmatrix}$$

Case 3

$$\lambda_1^{(1)} = 10,\ \lambda_2^{(1)} = 3,\ \lambda_3^{(1)} = 1 \qquad \lambda_1^{(2)} = 8,\ \lambda_2^{(2)} = 3,\ \lambda_3^{(2)} = 1 \qquad \lambda_1^{(3)} = 6,\ \lambda_2^{(3)} = 3,\ \lambda_3^{(3)} = 1$$

$$\mathbf{\Sigma}_1 \qquad\qquad\qquad \mathbf{\Sigma}_2 \qquad\qquad\qquad \mathbf{\Sigma}_3$$

$$\begin{pmatrix} 10 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 1 \end{pmatrix} \qquad \begin{pmatrix} 7.7359 & -0.8718 & -0.8953 \\ -0.8718 & 3.1000 & -0.1836 \\ -0.8953 & -0.1836 & 1.1642 \end{pmatrix} \qquad \begin{pmatrix} 2.75 & -1.9874 & -0.4874 \\ -1.9874 & 4.6856 & -0.875 \\ -0.4874 & -0.875 & 2.5643 \end{pmatrix}$$

Case 4

$$\lambda_1^{(1)} = 9,\ \lambda_2^{(1)} = 5,\ \lambda_3^{(1)} = 2 \qquad \lambda_1^{(2)} = 6,\ \lambda_2^{(2)} = 3,\ \lambda_3^{(2)} = 1 \qquad \lambda_1^{(3)} = 6,\ \lambda_2^{(3)} = 3,\ \lambda_3^{(3)} = 1$$

$$\mathbf{\Sigma}_1 \qquad\qquad\qquad \mathbf{\Sigma}_2 \qquad\qquad\qquad \mathbf{\Sigma}_3$$

$$\begin{pmatrix} 9 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 2 \end{pmatrix} \qquad \begin{pmatrix} 5.4776 & -0.7244 & -1.1998 \\ -0.7244 & 3.2010 & 0.1941 \\ -1.1998 & 0.1941 & 1.3215 \end{pmatrix} \qquad \begin{pmatrix} 2.75 & -1.9874 & -0.4874 \\ -1.9874 & 4.6856 & -0.875 \\ -0.4874 & -0.875 & 2.5643 \end{pmatrix}$$

Case 5

$$\lambda_1^{(1)} = 9,\ \lambda_2^{(1)} = 5,\ \lambda_3^{(1)} = 2 \qquad \lambda_1^{(2)} = 7.5,\ \lambda_2^{(2)} = 4,\ \lambda_3^{(2)} = 1.5 \qquad \lambda_1^{(3)} = 6,\ \lambda_2^{(3)} = 3,\ \lambda_3^{(3)} = 1$$

$$\mathbf{\Sigma}_1 \qquad\qquad\qquad \mathbf{\Sigma}_2 \qquad\qquad\qquad \mathbf{\Sigma}_3$$

$$\begin{pmatrix} 9 & 0 & 0 \\ 0 & 5 & 0 \\ 0 & 0 & 2 \end{pmatrix} \qquad \begin{pmatrix} 3.625 & -2.3776 & -0.6276 \\ -2.3776 & 5.9249 & -1.0625 \\ -0.6276 & -1.0625 & 3.4501 \end{pmatrix} \qquad \begin{pmatrix} 2.75 & -1.9874 & -0.4874 \\ -1.9874 & 4.6856 & -0.875 \\ -0.4874 & -0.875 & 2.5643 \end{pmatrix}$$

The following tables present the power of the Ansari-Bradley test and the Mood test. Tables 1(a)–(e) show the results of the normal population, and Tables 2(a)–(e) present the results of the contaminated normal population. We set $N_1 = N_2 = N_3 = 50$ for Tables 1(a) and 2(a), $N_1 = N_2 = N_3 = 100$ for Tables 1(b) and 2(b), $N_1 = N_2 = N_3 = 200$ for Tables 1(c) and 2(c), and $N_1 = 50$,

Table 1(a).  Normal population ($N_1 = N_2 = N_3 = 50$).

|  |  | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|---|---|---|---|---|---|---|
| $j = 1$ | $M_k$ | 0.041 | 0.041 | 0.242 | 0.224 | 0.169 |
|  | $AB_{ke}$ | 0.042 | 0.042 | 0.201 | 0.184 | 0.144 |
| $j = 2$ | $M_k$ | 0.048 | 0.047 | 0.049 | 0.335 | 0.253 |
|  | $AB_{ke}$ | 0.048 | 0.047 | 0.049 | 0.271 | 0.211 |
| $j = 3$ | $M_k$ | 0.056 | 0.056 | 0.056 | 0.564 | 0.445 |
|  | $AB_{ke}$ | 0.055 | 0.055 | 0.055 | 0.465 | 0.370 |

Table 1(b).  Normal population ($N_1 = N_2 = N_3 = 100$).

|  |  | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|---|---|---|---|---|---|---|
| $j = 1$ | $M_k$ | 0.046 | 0.046 | 0.478 | 0.428 | 0.327 |
|  | $AB_{ke}$ | 0.047 | 0.046 | 0.395 | 0.348 | 0.269 |
| $j = 2$ | $M_k$ | 0.049 | 0.049 | 0.050 | 0.615 | 0.488 |
|  | $AB_{ke}$ | 0.049 | 0.049 | 0.050 | 0.511 | 0.403 |
| $j = 3$ | $M_k$ | 0.053 | 0.052 | 0.053 | 0.867 | 0.763 |
|  | $AB_{ke}$ | 0.052 | 0.052 | 0.052 | 0.776 | 0.664 |

Table 1(c).  Normal population ($N_1 = N_2 = N_3 = 200$).

|  |  | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|---|---|---|---|---|---|---|
| $j = 1$ | $M_k$ | 0.048 | 0.048 | 0.802 | 0.732 | 0.600 |
|  | $AB_{ke}$ | 0.048 | 0.048 | 0.704 | 0.625 | 0.502 |
| $j = 2$ | $M_k$ | 0.049 | 0.049 | 0.050 | 0.903 | 0.807 |
|  | $AB_{ke}$ | 0.049 | 0.049 | 0.049 | 0.823 | 0.710 |
| $j = 3$ | $M_k$ | 0.051 | 0.051 | 0.051 | 0.993 | 0.974 |
|  | $AB_{ke}$ | 0.051 | 0.050 | 0.051 | 0.976 | 0.936 |

Table 1(d).  Normal population ($N_1 = 50, N_2 = 40, N_3 = 30$).

|  |  | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|---|---|---|---|---|---|---|
| $j = 1$ | $M_k$ | 0.038 | 0.038 | 0.168 | 0.181 | 0.116 |
|  | $AB_{ke}$ | 0.039 | 0.039 | 0.147 | 0.153 | 0.105 |
| $j = 2$ | $M_k$ | 0.046 | 0.046 | 0.048 | 0.309 | 0.206 |
|  | $AB_{ke}$ | 0.047 | 0.046 | 0.048 | 0.255 | 0.179 |
| $j = 3$ | $M_k$ | 0.059 | 0.059 | 0.059 | 0.540 | 0.379 |
|  | $AB_{ke}$ | 0.058 | 0.058 | 0.058 | 0.450 | 0.322 |

$N_2 = 40$ and $N_3 = 30$ for Tables 1(d) and 2(d), $N_1 = 200$, $N_2 = 150$ and $N_3 = 100$ for Tables 1(e) and 2(e).

From the results of the simulation, we can see the validity of the proposed nonparametric methods whether the covariance matrices are diagonal or not. When $j = 1, 2$, the Mood test $M_k$ and the Ansari-Bradley test $AB_{ke}$ may be

Table 1(e). Normal population ($N_1 = 200, N_2 = 150, N_3 = 100$).

|       |          | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|-------|----------|--------|--------|--------|--------|--------|
| $j = 1$ | $M_k$    | 0.047  | 0.047  | 0.613  | 0.641  | 0.420  |
|       | $AB_{ke}$ | 0.048  | 0.047  | 0.518  | 0.538  | 0.349  |
| $j = 2$ | $M_k$    | 0.049  | 0.049  | 0.050  | 0.852  | 0.641  |
|       | $AB_{ke}$ | 0.049  | 0.049  | 0.050  | 0.761  | 0.545  |
| $j = 3$ | $M_k$    | 0.052  | 0.052  | 0.052  | 0.985  | 0.898  |
|       | $AB_{ke}$ | 0.051  | 0.051  | 0.052  | 0.956  | 0.823  |

Table 2(a). Contaminated normal population ($N_1 = N_2 = N_3 = 50$).

|       |          | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|-------|----------|--------|--------|--------|--------|--------|
| $j = 1$ | $M_k$    | 0.040  | 0.040  | 0.233  | 0.215  | 0.163  |
|       | $AB_{ke}$ | 0.042  | 0.042  | 0.198  | 0.178  | 0.141  |
| $j = 2$ | $M_k$    | 0.048  | 0.047  | 0.048  | 0.329  | 0.248  |
|       | $AB_{ke}$ | 0.048  | 0.047  | 0.048  | 0.268  | 0.208  |
| $j = 3$ | $M_k$    | 0.055  | 0.055  | 0.055  | 0.555  | 0.438  |
|       | $AB_{ke}$ | 0.054  | 0.054  | 0.054  | 0.460  | 0.366  |

Table 2(b). Contaminated normal population ($N_1 = N_2 = N_3 = 100$).

|       |          | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|-------|----------|--------|--------|--------|--------|--------|
| $j = 1$ | $M_k$    | 0.044  | 0.044  | 0.462  | 0.411  | 0.311  |
|       | $AB_{ke}$ | 0.045  | 0.045  | 0.386  | 0.337  | 0.260  |
| $j = 2$ | $M_k$    | 0.047  | 0.047  | 0.048  | 0.601  | 0.473  |
|       | $AB_{ke}$ | 0.048  | 0.048  | 0.048  | 0.501  | 0.395  |
| $j = 3$ | $M_k$    | 0.051  | 0.050  | 0.050  | 0.858  | 0.751  |
|       | $AB_{ke}$ | 0.051  | 0.050  | 0.050  | 0.769  | 0.655  |

Table 2(c). Contaminated normal population ($N_1 = N_2 = N_3 = 200$).

|       |          | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|-------|----------|--------|--------|--------|--------|--------|
| $j = 1$ | $M_k$    | 0.045  | 0.045  | 0.789  | 0.715  | 0.582  |
|       | $AB_{ke}$ | 0.047  | 0.047  | 0.693  | 0.612  | 0.489  |
| $j = 2$ | $M_k$    | 0.046  | 0.047  | 0.047  | 0.895  | 0.794  |
|       | $AB_{ke}$ | 0.047  | 0.048  | 0.048  | 0.815  | 0.699  |
| $j = 3$ | $M_k$    | 0.048  | 0.048  | 0.048  | 0.992  | 0.971  |
|       | $AB_{ke}$ | 0.049  | 0.048  | 0.048  | 0.974  | 0.931  |

Table 2(d). Contaminated normal population ($N_1 = 50, N_2 = 40, N_3 = 30$).

|  |  | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|---|---|---|---|---|---|---|
| $j = 1$ | $M_k$ | 0.038 | 0.038 | 0.143 | 0.158 | 0.098 |
|  | $AB_{ke}$ | 0.040 | 0.040 | 0.131 | 0.139 | 0.094 |
| $j = 2$ | $M_k$ | 0.047 | 0.047 | 0.047 | 0.279 | 0.180 |
|  | $AB_{ke}$ | 0.047 | 0.047 | 0.047 | 0.233 | 0.159 |
| $j = 3$ | $M_k$ | 0.058 | 0.058 | 0.058 | 0.504 | 0.342 |
|  | $AB_{ke}$ | 0.056 | 0.056 | 0.056 | 0.421 | 0.292 |

Table 2(e). Contaminated normal population ($N_1 = 200, N_2 = 150, N_3 = 100$).

|  |  | Case 1 | Case 2 | Case 3 | Case 4 | Case 5 |
|---|---|---|---|---|---|---|
| $j = 1$ | $M_k$ | 0.045 | 0.045 | 0.593 | 0.617 | 0.401 |
|  | $AB_{ke}$ | 0.046 | 0.046 | 0.505 | 0.521 | 0.338 |
| $j = 2$ | $M_k$ | 0.047 | 0.047 | 0.047 | 0.836 | 0.621 |
|  | $AB_{ke}$ | 0.048 | 0.048 | 0.048 | 0.745 | 0.531 |
| $j = 3$ | $M_k$ | 0.050 | 0.050 | 0.050 | 0.982 | 0.887 |
|  | $AB_{ke}$ | 0.050 | 0.050 | 0.050 | 0.951 | 0.812 |

conservative under the null hypothesis. The power of the Mood test $M_k$ is greater than the power of the Ansari-Bradley test $AB_{ke}$ for every $j$-th eigenvalue in both the case that the sample sizes are equal or unequal. We have expected these results from the asymptotic relative efficiency of two tests. However, the power of both tests didn't depend on the distribution which was either a normal or a contaminated normal distribution.

Additionally, the simulation results indicate that it is difficult to keep the significance level when the sample sizes are small. Therefore, the sample size $N_i$ should be greater than 50 for the case of $k = 3$ and $p = 3$. In Case 3, the difference of eigenvalue is only $j = 1$. Therefore the power of tests increases only the case for the largest eigenvalue. It might require sufficiently large sample size, larger than 100 from tables. The difference of eigenvalues on Case 5 is greater than the difference on Case 4. Then the powers of tests on Case 5 are higher on Case 4.

## 4. Conclusion and discussion

In this paper, we propose the nonparametric test by using principal component scores under the multipopulation and apply the testing procedure under the normal population and the contaminated normal population when the population eigenvalues are separated and the sample sizes are large. Though the convergence for significance level of the procedure using Ansari-Bradley test and the procedure using Mood test is almost the same tendency, the power of Mood test is greater than the power of Ansari-Bradley test.

From the asymptotic relative efficiency of nonparametric test, we have expected that using the Mood test for $k$-population is more suitable than using the

Ansari-Bradley test for $k$-population. This result is showed by simulation.

It will also be important to develop where the principal component scores are evaluated from correlation matrices.

## Acknowledgements

### REFERENCES

Anderson, T. W. (2003). *An Introduction to Multivariate Statistical Analysis*, 3rd ed., John Wiley, New York.

Ansari, A. R. and Bradley, R. A. (1960). Rank sum tests for dispersion, *Ann. Math. Statist.*, **31**, 1174–1189.

Gibbons, J. D. and Chakraborti, S. (2003). *Nonparametric Statistical Inference*, 4th ed., Dekker, New York.

Jackson, J. E. (2003). *A User's Guide to Principal Components*, John Wiley, New York.

Jolliffe, I. T. (2002). *Principal Component Analysis*, 2nd ed., Springer, New York.

Mood, A. M. (1954). On the asymptotic efficiency of certain nonparametric two-sample tests, *Ann. Math. Statist.*, **25**, 514–522.

Sugiyama, T. and Ushizawa, K. (1998). A non-parametric method to test equality of intermediate latent roots of two populations in a principal component analysis, *J. Japan Statist. Soc.*, **28**, 227–235.

Takeda, Y. (2001). Permutation test for equality of each characteristic root in two populations, *J. Jpn. Soc. Comp. Statist.*, **14**, 1–10.

Tsai, W. S., Duran, B. S. and Lewis, T. O. (1975). Small-sample behavior of some multisample nonparametric tests for scale, *J. Amer. Statist. Assoc.*, **70**, 791–796.