



One-sided Multivariate Tests for High Dimensional Data

โดย

รศ.ดร.สำรวม จงเจริญ

งานวิจัยสมบูรณ์แบบได้รับทุนสนับสนุนจากคณะกรรมการส่งเสริมการวิจัย
สถาบันบัณฑิตพัฒนบริหารศาสตร์

EXECUTIVE SUMMARY

At present people involve a lot of data in various areas such as genetic microarrays, medical imaging, pharmaceutical science, finance, chemometrics, climatology or astronomy etc. Each of them consists of large amount of data. There might be many variables and/or many observations on each variables. One can think of each variable as an additional dimension, and so many variables corresponds to data setting in a high dimensional space. To analysis this kind of data ,in statistics areas, we called “ multivariate analysis”. Multivariate analysis deals with observations on more than one variable when there is or may be some dependence between the vaiables. The most basic phenomenon is that of correlation- the tendency of quantities to vary together. All classical computation statistical technique available such as the Hotelling’s T^2 approach, deal with the number of observations n is larger than the number of dimensions p , $n > p$. In some situations, people are exposed to situations where the number of dimensions increases dramatically, but, the number of subjects increases much more slowly. Especially, in most medical and pharmaceutical studies, due to budgetary constraints and other experimental restrictions, microarray also plays a key role in molecular biology, in pharmacy and in medicine for discovering certain diseases and developing new drugs. It is very common that microarray data containing gene expression values measured on thousands of genes from much fewer biological subjects. Actually, analysis of microarray data gave urge to search for a test statistic which works for high dimensionality. When one have to analysis the data with the number of subjects relatively smaller than the number of dimension, i.e. $n < p$, called high dimensional data, those classical statistical technique available are not powerful or cannot to be used because the estimated variance-covariance estimator is no longer non-singular.

Therefore, the high dimensionality problem was very challenging or interesting and has drawn a lot of attention for several decades since it is no exact solution satisfying the classical criteria for good tests. Much current research in statistics, both in statistical theory, and in many areas of application, such as genomics, climatology or astronomy, focuses on the problems and opportunities posed by availability of large amount of data. Modern computation techniques make it possible to deal with high dimensional data. Originally the statistical technique for testing $H_0 : \theta_1 = \theta_2 = \dots = \theta_p = 0$ agaist $H_1 : \text{at least } \theta_i \neq 0$ when the population has normal distribution with mean θ and covariance matrix V for $n \leq p$ was proposed by Dempster (1958, 1960), later more than 30 years, Bai and Saranadasa(1996) proposed another test statistic; in 2008 Srivastava and Du proposed another test statistic; and Ahmad, Werner and Brunner proposed another test statistic which all these tests we will give more detail later. For others statistical techniques related with high dimensional data appeared recently, see Yasunori Fujikoshi et al.(2004), Srivastava and Fujikoshi(2006), Siegfried Kropf et al. (2008) and Srivastava and Yanagihara (2010).

With high dimensional data, one may consider testing the null hypothesis $H_0: \theta = \underline{0}$ versus $H_1 - H_0$ where $H_1: \theta \in \Omega_p$ and $\Omega_p = \{x: x_i \geq 0 \text{ for } i = 1, 2, \dots, p\}$ is the p -dimensional nonnegative orthant or called multivariate one-sided testing. We applied Follmann's test, a test which Follmann (1996) proposed for one-sided modifications of the usual omnibus chi-squared test and Hotelling's T^2 test for $n > p$, with the tests of Dempster (1958, 1960), Bai and Saranadasa (1996), Srivastava and Du (2008) and Ahmad, Werner and Brunner (2008) for $n \leq p$. We proposed one-sided multivariate tests for $n < p$ from the combination of Follmann's test with the mentioned high dimensional tests above and found that only DF test (the combination of Dempster's and Follmann's test) and BSF^* (the combination of Bai and Saranadasa and Follmann's test) provided reasonable type I error rate for one-sided covariance structures. We compared the powers of these two tests and found that over all both tests, DF test and BSF^* test, gave almost the same powers in every p and n and every covariance matrices structure considered. Then we recommend DF test and BSF^* test for testing one-sided hypothesis of high dimensional data. we also recommend for data with $p \geq 20$ and $n > 10$ for any covariance matrices structure. We gave the example for using these proposed tests on DNA micro arrays which both tests gave the same result. Therefore, for one-sided multivariate tests that one believes that for each coordinate, the mean responses for treatment one are at least as large as those for treatment two and the data has the number n of available observations is smaller than the dimension p ($n \leq p$), the proposed tests, DF test and BSF^* test, perform best and are recommended for $p \geq 20$ and $n > 10$ under the circumstances considered in this paper.

ABSTRACT

For a multivariate normal population with size larger than dimension , $n > p$, Kudo (1963), Shorack (1967) and Perlman (1969) derived the likelihood ratio tests of the null hypothesis that the mean vector is zero with a one-sided alternative for a known covariance matrix, a partially known covariance matrix and a completely unknown covariance matrix, respectively. Because these tests may be tedious to use, Tang et al.(1989) developed approximate likelihood ratio tests and Follmann (1996) proposed one-sided modifications of the usual omnibus chi-squared test and Hotelling's T^2 test. Chongcharoen et al.(2002) considered a modification of Follmann's test (the new test) to include information of off diagonal of covariance matrix , which adjusts for possibly unequal variances. For the non-normal population, Boyett and Shuster (1977) proposed a nonparametric one-sided test and Chongcharoen et al. (2002) used their technique to develop nonparametric versions of Perlman's test, Follmann's test, the new test and the Tang-Gnecco-Geller test. Also Chongcharoen et al. (2002) considered known and partially known covariance matrices. Chongcharoen (2009) studied the powers of these one-sided tests for an unknown covariance matrices. In some situations, there are no longer data for $n > p$. That is, when the number n of available observations is smaller than the dimension p of the observed vectors. For example, the data comes from DNA micro arrays where thousands of gene expression levels are measured on relatively few subjects. The one-sided multivariate tests as above are no longer valid for this kind of data. The proposed tests are tests for one-sided multivariate tests with $n < p$ provided reasonable type I error rate for one-sided covariance structures. Their powers are compared for alternatives. An example for using these proposed tests on DNA micro arrays is given. However, the methodology is valid for any application which involves high-dimensional data.

TABLE OF CONTENTS

ABSTRACT

CHAPTER 1 INTRODUCTION.....	1
1.1 Introduction.....	1
1.2 The Objective of study	3
1.3 Scope of study	3
1.4 Research methodology	3
CHAPTER 2 LITERATURE REVIEW AND PRELIMINRIES.....	4
2.1 Dempster test	4
2.2 Bai and Saranadasa test	5
2.3 Srivastava and Du test	5
2.4 Ahmad, Werner and Brunner test	6
CHAPTER 3 TEST STATISTIC DEVELOPMENT	7
3.1 The estimated significant level for high dimensional multivariate tests	7
3.2 The proposed tests for high dimensional data.....	7
3.3 Simulation study.....	9
3.4 An example	9
3.5 Conclusion	10
TABLES	10-17
REFERENCE	18-19

Chapter 1

Introduction

1.1 Introduction

Suppose one uses a matched-pair design to compare the multivariate responses of two treatments. If the responses are p dimensional and $\theta = (\theta_1, \theta_2, \dots, \theta_p)'$ is the difference, treatment one minus treatment two, of the mean responses, then one may test the null hypothesis, $H_0 : \theta_1 = \theta_2 = \dots = \theta_p = 0$, to determine if there is a difference in the two treatments. Furthermore, if one believes that for each coordinate, the mean responses for treatment one are at least as large as those for treatment two, then the alternative can be constrained by $H_1 : \theta_i \geq 0$ for $i = 1, 2, \dots, p$.

Based on a random sample with $n > p$ from the normal distribution with mean θ and covariance matrix V , Kudo (1963), Shorack (1967) and Perlman (1969) derived the likelihood ratio test of H_0 versus $H_1 - H_0$ for the cases in which V is known, known up to a multiplicative constant and completely unknown, respectively. Because the likelihood ratio tests with restricted alternatives are complicated to use, Tang et al.(1989) proposed an approximate likelihood ratio test, and Follmann (1996) proposed one-sided modifications of the usual χ^2 and Hotelling's T^2 tests of H_0 versus $\sim H_0$ that are easier to implement. Using exact computations and Monte Carlo methods, Chongcharoen et al.(2002) compared the performance of Kudo's test, Follmann's test, a new test, which is a modification of Follmann's test, the permutation test of Boyett and Shuster(1977) and the Tang-Gnecco-Geller test for a known covariance matrix, and for a partially known covariance matrix, they compared the powers of these tests with Kudo's test replaced by Shorack's test. For a completely unknown covariance matrix, Chongcharoen (2009) studied the power of these one-sided tests for unknown covariance matrices with equal variances and unequal variances as well as tests obtained by combining the Boyett-Shuster technique (1977) with Follmann's test, the new test, Perlman's test and the Tang-Gnecco-Geller test.

In some situations, there are no longer data for $n > p$. That is, when the number n of available observations is smaller than the dimension p of the observed vectors. For example, the data comes from DNA micro arrays where thousands of gene expression levels are measured on relatively few subjects. The one-sided multivariate tests as above are no longer valid for this kind of data because the $p \times p$ sample covariance matrix S is singular with rank $n < p$, S^{-1} does not exist. Since now there have no one-sided multivariate tests available for the data which has the number n of available observations is smaller than the dimension p yet, therefore We interested in developing the one-sided multivariate tests for the data with $n < p$.

Throughout this paper, suppose X_1, X_2, \dots, X_n is a random sample from a p -dimensional multivariate normal distribution with unknown mean $\theta = (\theta_1, \theta_2, \dots, \theta_p)'$

and unknown positive definite covariance matrix V with $n \leq p$. One may consider testing the null hypothesis $H_0: \theta = 0$ versus $H_1 - H_0$ where $H_1: \theta \in \Omega_p$ and $\Omega_p = \{x : x_i \geq 0 \text{ for } i=1,2, \dots, p\}$ is the p -dimensional nonnegative orthant. The sample mean and covariance are

$$\bar{X} = \sum_{i=1}^n \frac{X_i}{n} \quad \text{and} \quad S = \sum_{i=1}^n \frac{(X_i - \bar{X})(X_i - \bar{X})'}{n-1}, \quad (1.1)$$

when $n < p$, S is a singular matrix.

The hypotheses H_0 and H_1 also arise in the one-way analysis of variance when the means are known to satisfy an order restriction. For observations which come from k normal populations whose means are known to satisfy a simple ordering, i.e. $H_s: \mu_1 \leq \mu_2 \leq \dots \leq \mu_k$, Bartholomew (1959a, 1959b, 1961) derived the likelihood ratio test of $\mu_1 = \mu_2 = \dots = \mu_k$ with the alternative restricted by H_s for the cases of known variances and variances known up to a multiplicative constant. Suppose the observations are Y_{ij} for $j = 1, 2, \dots, n_i$ and $i = 1, 2, \dots, k$, and the sample means are $\bar{Y}_1, \bar{Y}_2, \dots, \bar{Y}_k$. With known variances, $\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2$, Kudo (1963) noted that for $p = k - 1$, $X_i = \bar{Y}_{i+1} - \bar{Y}_i$ for $i = 1, 2, \dots, p$, $X = (X_1, X_2, \dots, X_p)'$ and $\theta = E(X)$, the hypotheses on μ are equivalent to H_0 and H_1 above, Bartholomew's and Kudo's tests are equivalent. With $w_i = n_i / \sigma_i^2$ for $i = 1, 2, \dots, k$, the correlation matrix for X satisfies

$$\rho_{i,i+1} = - \sqrt{\frac{w_i w_{i+2}}{(w_i + w_{i+1})(w_{i+1} + w_{i+2})}} \text{ for } i = 1, 2, \dots, p - 1 \quad (1.2)$$

and $\rho_{ij} = 0$ for $|i - j| \geq 2$.

If the weights are equal, i.e. $w_1 = w_2 = \dots = w_k$, the correlation matrix in (2) is denoted by R_s . Also, Bartholomew (1959a, 1959b, 1961) considered an arbitrary partial order restriction, which includes the simple tree order, i.e. $H_T: \mu_1 \leq \mu_i$ for $i = 1, 2, \dots, k$. For this ordering, one takes differences, $X_i = \bar{Y}_{i+1} - \bar{Y}_1$ for $i = 1, 2, \dots, p$, and with $p = k - 1$ and w_i as above, the correlation matrix of $X = (X_1, X_2, \dots, X_p)'$ satisfies

$$\rho_{i,j} = \sqrt{\frac{w_{i+1} w_{j+1}}{(w_{i+1} + w_1)(w_{j+1} + w_1)}} \text{ for } 1 \leq i \neq j \leq p. \quad (1.3)$$

If the weights are equal, i.e. $w_1 = w_2 = \dots = w_k$, the correlation matrix in (3) is denoted by R_T . The powers of the proposed tests are compared for several correlation matrices including R_S and R_T .

1.2 The objective of study

The objective of this proposal is to develop one-sided multivariate tests for the data which has the number n of available observations is smaller than the dimension p .

1.3 Scope of study

When $n \leq p$, some one-sided multivariate statistic tests will be developed under the unrestricted alternatives tests as in chapter 2 or other reasonable test may be proposed. The developed tests are tested for reasonably estimated significance level and compare them with empirical powers under various covariance structures such as the matrices from the simple order correlation, the simple tree order correlation as well as some others form of the correlation matrices.

1.4 Research methodology

- (1) Study unrestricted alternative multivariate tests available for $n \leq p$ and then check their estimated significance level under covariance structure which it is possible to develop to be one-sided multivariate tests. This computations will be done by Monte Carlo technique.
- (2) Develop one-sided multivariate tests for $n \leq p$ from possible multivariate tests with unrestricted multivariate alternative for $n \leq p$.
- (3) Compute the estimated significance level of the proposed tests under covariance structure considered. This computations will be done by Monte Carlo technique.
- (4) Investigate the reasonable of the proposed tests under the difference covariance structure considered
- (5) Compute the power of the proposed tests under the difference covariance structure considered. This computations also will be done by Monte Carlo technique.
- (6) Compare the power of each one-sided multivariate tests for $n \leq p$ developed for each particular form of the covariance matrices.

Chapter 2

Literature review and preliminaries

Much current research in statistics, both in statistical theory, and in many areas of application, such as genomics, climatology or astronomy, focuses on the problems and opportunities posed by availability of large amount of data. There might be many variables and/or many observations on each variables. One can think of each variable as an additional dimension, and so many variables corresponds to data sitting in a high dimensional space. In mathematical themes one could follow Banach space theory, convex geometry, even topology and in statistics areas we called “multivariate analysis”. Multivariate analysis deals with observations on more than one variable when there is or may be some dependence between the variables. The most basic phenomenon is that of correlation- the tendency of quantities to vary together. All classical computation statistical technique available deal with the number of observations n is larger than the number of dimensions p , $n > p$ which are not possible for computing with the data that the number of observations n is less than the number of dimensions p , $n \leq p$. Modern computation techniques make it possible to deal with high dimensional data, the data that the number of observations n is less than the number of dimensions p , $n \leq p$. Originally the statistical technique for testing $H_0 : \theta_1 = \theta_2 = \dots = \theta_p = 0$ against $H_1 : \text{at least } \theta_i \neq 0$ when the population has normal distribution with mean θ and covariance matrix V for $n \leq p$ was proposed by Dempster (1958, 1960), later more than 30 years, Bai and Saranadasa(1996) proposed another test statistic; in 2008 Srivastava and Du proposed another test statistic; and Ahmad, Werner and Brunner proposed another test statistic which all these tests we will give more detail later. For others statistical techniques related with high dimensional data appeared recently, see Yasunori Fujikoshi et al.(2004), Srivastava and Fujikoshi(2006), Siegfried Kropf et al. (2008) and Srivastava and Yanagihara (2010).

Again the tests with unrestricted alternative test when $n \leq p$, that is, the tests with the hypothesis

$$\begin{aligned} H_0 : \theta_1 = \theta_2 = \dots = \theta_p = 0 \\ H_1 : \text{at least } \theta_i \neq 0 \end{aligned}$$

are proposed by several researchers recently related with the tests we proposed such as

2.1 Dempster test Dempster (1958, 1960) proposed the test statistic

$$D = \frac{n\bar{X}\bar{X}}{tr(S)} \quad (2.1)$$

where \bar{X} is the sample mean vector and S is the sample covariance defined as in (1.1). Under null hypothesis, Dempster showed that D has an approximate F-distribution with $[\hat{r}]$ and $[(n-1)\hat{r}]$ degrees of freedom, where $[a]$ denotes the largest integer less than or equal to a and

$$\hat{r} = p \frac{\hat{a}_1^2}{\hat{a}_2},$$

$$\hat{a}_1 = \frac{tr(S)}{p} \quad \text{and} \quad \hat{a}_2 = \frac{(n-1)^2}{(n-2)(n+1)} \frac{1}{p} \left[tr(S^2) - \frac{(tr(S))^2}{n-1} \right]$$

under condition $0 < \lim_{p \rightarrow \infty} a_i = \lim_{p \rightarrow \infty} \frac{(tr(V^i))}{p} = a_{i0} < \infty; i = 1, 2, 3, 4.$

2.2 Bai and Saranadasa test Bai and Saranadasa (1996) proposed test statistic as

$$\begin{aligned} BS &= \frac{n\bar{X}'\bar{X} - tr(S)}{\left[\frac{n}{n-1} \right]^{\frac{1}{2}} [2p\hat{a}_2]^{\frac{1}{2}}} \\ &= \frac{n\bar{X}'\bar{X} - tr(S)}{\left[\frac{2n(n-1)}{(n-2)(n+1)} \left(tr(S^2) - \frac{(tr(S))^2}{n-1} \right) \right]^{\frac{1}{2}}} \end{aligned} \quad (2.2)$$

under condition $0 < \lim_{p \rightarrow \infty} a_i = \lim_{p \rightarrow \infty} \frac{(tr(V^i))}{p} = a_{i0} < \infty, i = 1, 2, 3, 4$ and for

$\lambda_i = O(p^\gamma), 0 \leq \gamma \leq \frac{1}{2}$, where λ_i are the eigenvalues of V , then under the null hypothesis

$$\lim_{(n,p) \rightarrow \infty} P_0(BS \leq z) = \Phi(z).$$

That is, BS has the asymptotic distribution as $N(0,1)$.

2.3 Srivastava and Du test Srivastava and Du (2008) proposed test statistic as

$$SD = \frac{n\bar{X}'D_s^{-1}\bar{X} - \frac{(n-1)}{(n-3)}p}{\sqrt{2(tr(R^2) - \frac{p^2}{n-1})(1 + \frac{tr(R^2)}{p^{\frac{3}{2}})}}} \quad (2.3)$$

where $R = D_S^{-\frac{1}{2}} S D_S^{-\frac{1}{2}}$ is the sample correlation matrix and $D = \text{diag}(s_{11}, \dots, s_{pp})$ is the diagonal matrix with the diagonal elements of S defined in (1). Under conditions stated in their paper, when $\underline{\theta} = 0$, SD is asymptotically distributed as $N(0, 1)$.

2.4 Ahmad, Werner and Brunner test Ahmad, Werner and Brunner (2008) proposed another test statistic as

$$AWB = \frac{n\bar{X}\bar{X}}{B_0} \times \frac{B_1}{B_2} \quad (2.4)$$

where $X_k \sim N(\theta, V), V > 0, k = 1, 2, \dots, n$ and $A_k = X_k' X_k$ and $A_{kl} = X_k' X_l, k \neq l$ then

$$B_0 = \frac{1}{n} \sum_{k=1}^n A_k$$

$$B_1 = \frac{1}{n(n-1)} \sum_{k=1}^n \sum_{l=1}^n A_k A_l ; k \neq l$$

$$B_2 = \frac{1}{n(n-1)} \sum_{k=1}^n \sum_{l=1}^n A_{kl}^2 ; k \neq l$$

and it has an approximate chi-square distribution with degree of freedom $\tilde{f} = \frac{B_1}{B_2}$.

AWB 's test rejects H_0 if

$$AWB > \chi_{\alpha, \tilde{f}}^2$$

where $\chi_{\alpha, \tilde{f}}^2$ is the $(1-\alpha)^{th}$ quantile of the central chi-square distribution with \tilde{f} degrees of freedom.

Chapter 3

Test Statistic Development

3.1 The estimated significant level for high dimensional multivariate tests

In development of one-sided multivariate tests for high dimensional data, we consider from the available high dimensional unrestricted alternative multivariate tests such as Dempster's test, Bai and Saranadasa's test, Srivastava and Du's test and Ahmad, Werner and Brunner's test which we mentioned in chapter 2 . We study these tests by computing their estimated significant level under covariance structure which will use in one-sided tests, that is , R_S and R_T as well as the correlation matrix with all off-diagonal elements equals to -0.5 , called $\mathfrak{R} = R_1$, and the correlation matrix $\mathfrak{R} = R_2 = [\rho_{ij}]$, $\rho_{ij} = -0.5$ for $j = 2i$ and $\rho_{ij} = 0.7$ for elsewhere . To compute these estimated significant level, we use the Monte Carlo technique by considering on $p = 10, n = 5, 10$; $p = 20, n = 10, 20$; $p = 30, n = 10, 15, 20, 30$; $p = 40, n = 10, 20, 30, 40$; $p = 50, n = 10, 20, 25, 30, 40, 50$; $p = 60, n = 10, 20, 30, 40, 50, 60$; $p = 100, n = 10, 20, 50$; $p = 200, n = 10, 20, 50$; $p = 400, n = 10, 20, 50$. Each case is repeated 10,000 times and the proportion of rejections are recorded for each tests. All of these tests are conducted using the level of significance $\alpha = .05$. It were found that with their critical values, the estimated significance level of these tests, except D test statistic with R_S , under all covariance matrices considered, including R_S and R_T , are not consistency to $\alpha = .05$. That is, the estimated significance level of some tests are too large with some covariance matrices meanwhile the other tests give the estimated significance level too small with the other covariance structure considered. For instance, the estimated significance level for D 's test ranges from 0.048 to 0.057 for $p = 10, n \geq 10$ for the simple order and from 0.051 to 0.094 for the simple tree, the estimated significance level for BS 's test ranges from 0.055 to 0.78 for the simple order and from 0.070 to 0.111 for the simple tree, the estimated significance level for SD 's test ranges from 0.036 to 0.160 for the simple order and from 0.015 to 0.158 for the simple tree and the estimated significance level for AWB 's test ranges from 0.043 to 0.060 for the simple order and from 0.061 to 0.074 for the simple tree, detailed in Table 1. Since the performance of BS 's test, SD 's test and AWB 's test are extremely poor, but Bai and Saranadasa(1996) showed that their test, BS 's test, has asymptotic powers the same as those of D 's test, thus the only D 's test and BS 's test will be studied further for one-side alternatives.

3.2 The proposed tests for high dimensional data

For the tests with restricted alternatives ,that is, to test the null hypothesis

$H_0 : \theta = 0$ versus $H_1 - H_0$ where $H_1 : \theta \in \Omega_p$ and $\Omega_p = \{x : x_i \geq 0 \text{ for } i=1,2, \dots, p\}$, one may applied Follmann's (1996) test to both D 's test and BS 's test. When applied with D 's test which is denoted DF , it rejects H_0 at level α if

$$D > F_{2\alpha; [\hat{r}], [(n-1)\hat{r}]} \quad \text{and} \quad \sum_{j=1}^p \bar{X}_j > 0$$

where $F_{2\alpha; [\hat{r}], [(n-1)\hat{r}]}$ is the $(1-2\alpha)^{th}$ quantile of the central F-distribution with $[\hat{r}]$ and $[(n-1)\hat{r}]$ degrees of freedom. By Theorem 2.1 of Follmann (1996), one has $1'\theta = 0$, and the significance level is approximated by

$$\begin{aligned} & \Pr(D > F_{1-2\alpha, [\hat{r}], [(n-1)\hat{r}]} \cap 1'\bar{X} > 0) \\ &= \Pr(D > F_{1-2\alpha, [\hat{r}], [(n-1)\hat{r}]}) \Pr(1'\bar{X} > 0) \\ &= (2\alpha) \times \frac{1}{2} \\ &= \alpha \end{aligned}$$

Also, when applied Follmann's idea to BS 's test, called BSF , one may reject H_0 at level α if

$$BS > Z_{2\alpha} \quad \text{and} \quad \sum_{j=1}^p \bar{X}_j > 0$$

where $Z_{2\alpha}$ is the $(1-2\alpha)^{th}$ quantile of the standard normal distribution. It also noted that, after Theorem 2.1 of Follmann (1996), the significance level of this test is α .

In the Monte Carlo study conditioned as above, it was found that DF shows reasonably well for all conditions considered, that is, the estimated significance level for DF ranges from 0.048 to 0.057 for the simple order with $p \geq 10$ and $n \geq 5$ and from 0.040 to 0.050 for the simple tree with $p \geq 20$ and $n > 10$ and the estimated significance level for BSF ranges from 0.051 to 0.085 for the simple order and from 0.048 to 0.087 for the simple tree. Because of the poor performance of this version of BSF , it is modified as BSF^* which rejects H_0 at level α if

$$BS^* > Z_{2\alpha} \quad \text{and} \quad \sum_{j=1}^p \bar{X}_j > 0$$

where

$$BS^* = \frac{n\bar{X}'\bar{X} - trS}{\left[\frac{2n}{(n-2)} \left(trS^2 - \frac{(trS)^2}{n-1} \right) \right]^{\frac{1}{2}}}$$

Then, BSF^* gives the Monte Carlo estimated significance level ranges from 0.046 to 0.057 for the simple order for every p and n considered and ranges from 0.045 to 0.056 for the simple tree with $p \geq 20$ and $n > 10$.

For the correlation matrix $\mathfrak{R} = R_1$ and $\mathfrak{R} = R_2$, DF gives the estimated significance level range from 0.046 to 0.058 for every p and n considered and range from 0.045 to 0.053 respectively when $p \geq 10$ and $n \geq 10$ respectively and BSF^* gives the estimated significance level range from 0.048 to 0.056 for $p \geq 40$ and $n \geq 20$ when $\mathfrak{R} = R_1$ and range from 0.049 to 0.056 for $p \geq 30$ and $n \geq 15$ when $\mathfrak{R} = R_2$. All estimated significance levels considered are given in Table 2.

3.3 Simulation Study

To compare these two tests, the performances of them are studied by Monte Carlo techniques for multivariate normal distributions with the correlation matrices R_S and R_T , that is for the simple order and the simple tree order correlations with equal weights as well as some other forms of correlation structures such as $\mathfrak{R} = R_1$ and $\mathfrak{R} = R_2$. Recall, R_S and R_T are given in (1.2) and (1.3), respectively. The mean vector for the alternative hypothesis is chosen in the non-negative orthant as $\theta = (v_1, v_2, \dots, v_p)'$; $v_{2k-1} = 0$ and $v_{2k} \sim \text{iid Unif}(0,1)$, $k = 1, 2, \dots, p/2$ so that the tests will be rejected. As before, 10,000 iterations are used. In each iteration, n multivariate normal X 's with the chosen mean vector and covariance of the form \mathfrak{R} are generated and the proportion of rejections for these tests were recorded. All of these tests are conducted using the level of significance $\alpha = 0.05$. Monte-Carlo estimates of power are given in Table 3.

From Table 3, it was shown that over all both tests, DF test and BSF^* test, gave almost the same powers in every p and n and every covariance matrices structure considered. Therefore, for protection some gains of using these two tests, we recommend these tests for high dimensional data when $p \geq 20$ and $n > 10$ for any covariance matrices structure.

3.4 An example

In this section, the proposed tests are applied to an example of DNA micro arrays data which the data are 8280 (p) genes expression information on 110 childhood suffering from acute lymphoblastic leukemia. To see the changes in gene expression after treatment, the data were cleaned and then obtained the difference of gene expression from before and after treatment of 50 children in 254 (p) genes expression (<http://www.ailab.si/supp/bi-cancer/projections/info>). The results of using these two tests is shown in Table 4. The p -values of DF test and BSF^* test equal to 0.0129 and 0.0000 respectively. Thus, all two tests lead to the rejection of the

hypothesis that the gene expression after treatment have the same level as before treatment.

3.5 Conclusion

For one-sided multivariate tests that one believes that for each coordinate, the mean responses for treatment one are at least as large as those for treatment two and the data has the number n of available observations is smaller than the dimension p ($n \leq p$), the proposed tests, DF test and BSF^* test, perform best and are recommended for $p \geq 20$ and $n > 10$ under the circumstances considered in this paper.

Table 1: Attained significance level of Dempster's test (D), Bai and Saranadasa's test (BS), Srivastava and DU's test(SD) and Ahmad-Werner-Brunner's test (AWN) under the null hypothesis for correlation matrix $\mathfrak{R} = R_s, \mathfrak{R} = R_T, \mathfrak{R} = R_1$ and $\mathfrak{R} = R_2$ at $\alpha = 0.05$

p	n	D	BS	SD	AWB
$\mathfrak{R} = R_s$					
10	5	0.066	0.078	0.160	0.044
	10	0.049	0.068	0.094	0.052
20	10	0.051	0.061	0.096	0.049
	20	0.053	0.067	0.072	0.060
30	10	0.054	0.059	0.097	0.051
	15	0.053	0.062	0.080	0.052
	20	0.055	0.067	0.070	0.059
	30	0.049	0.061	0.060	0.055
40	10	0.054	0.056	0.096	0.048
	20	0.056	0.067	0.071	0.059
	30	0.051	0.061	0.062	0.056
	40	0.050	0.062	0.054	0.057
50	10	0.057	0.058	0.096	0.050
	20	0.056	0.061	0.068	0.057
	25	0.049	0.058	0.058	0.051
	30	0.049	0.058	0.056	0.053
	40	0.052	0.062	0.057	0.057
	50	0.048	0.058	0.050	0.054
60	10	0.054	0.054	0.095	0.048
	20	0.054	0.060	0.066	0.055
	30	0.050	0.058	0.057	0.053
	40	0.052	0.059	0.055	0.055
	50	0.048	0.058	0.050	0.053
	60	0.054	0.066	0.056	0.060
100	10	0.056	0.052	0.088	0.048
	20	0.053	0.055	0.056	0.052
	50	0.049	0.054	0.049	0.051
200	10	0.055	0.049	0.066	0.045
	20	0.051	0.051	0.048	0.049
	50	0.049	0.053	0.045	0.051
400	10	0.054	0.044	0.050	0.043
	20	0.053	0.050	0.036	0.051
	50	0.052	0.053	0.040	0.052

Table 1: (continue.)

p	n	D	BS	SD	AWB	
$\mathfrak{R} = R_T$						
10	5	0.094	0.111	0.158	0.061	
	10	0.069	0.093	0.090	0.074	
20	10	0.075	0.096	0.086	0.074	
	20	0.061	0.080	0.060	0.072	
30	10	0.071	0.092	0.079	0.068	
	15	0.062	0.083	0.062	0.069	
	20	0.059	0.080	0.053	0.070	
30	30	0.054	0.078	0.046	0.071	
	40	10	0.076	0.092	0.080	0.069
		20	0.058	0.079	0.051	0.069
30		0.054	0.076	0.043	0.070	
40		0.055	0.074	0.042	0.071	
50	10	0.079	0.096	0.080	0.074	
	20	0.065	0.084	0.054	0.075	
	25	0.057	0.079	0.045	0.069	
	30	0.058	0.078	0.044	0.072	
	40	0.054	0.075	0.039	0.068	
50	50	0.051	0.070	0.034	0.068	
	60	10	0.075	0.091	0.075	0.069
		20	0.061	0.081	0.048	0.071
		30	0.054	0.076	0.038	0.070
		40	0.056	0.077	0.036	0.072
50		0.056	0.076	0.036	0.072	
60	60	0.055	0.076	0.035	0.074	
	100	10	0.079	0.094	0.072	0.071
		20	0.064	0.081	0.046	0.072
50		0.049	0.070	0.029	0.066	
200	10	0.081	0.096	0.066	0.073	
	20	0.062	0.083	0.036	0.073	
	50	0.050	0.072	0.022	0.068	
400	10	0.082	0.099	0.055	0.074	
	20	0.063	0.081	0.027	0.070	
	50	0.053	0.072	0.015	0.070	

Table 1: (continue.)

p	n	D	BS	SD	AWB
$\mathfrak{R} = R_1$					
10	5	0.073	0.121	0.081	0.069
	10	0.053	0.097	0.050	0.082
20	10	0.055	0.101	0.038	0.070
	20	0.045	0.079	0.025	0.062
30	10	0.053	0.095	0.030	0.058
	15	0.047	0.084	0.021	0.056
	20	0.046	0.079	0.018	0.056
40	30	0.044	0.077	0.015	0.055
	10	0.053	0.100	0.027	0.057
	20	0.047	0.079	0.016	0.054
50	30	0.047	0.077	0.012	0.056
	40	0.048	0.074	0.011	0.056
	10	0.049	0.096	0.023	0.048
	20	0.049	0.080	0.013	0.053
60	25	0.048	0.077	0.013	0.052
	30	0.048	0.079	0.012	0.054
	40	0.049	0.079	0.010	0.056
	50	0.047	0.074	0.009	0.054
	10	0.053	0.094	0.022	0.051
100	20	0.048	0.081	0.013	0.051
	30	0.047	0.076	0.009	0.052
	40	0.045	0.070	0.008	0.051
	50	0.049	0.075	0.009	0.055
	60	0.050	0.075	0.008	0.054
200	10	0.050	0.097	0.016	0.044
	20	0.051	0.080	0.009	0.051
	50	0.047	0.073	0.005	0.050
400	10	0.056	0.101	0.013	0.045
	20	0.051	0.079	0.006	0.049
	50	0.048	0.073	0.002	0.049
400	10	0.054	0.095	0.008	0.041
	20	0.054	0.082	0.004	0.050
	50	0.049	0.072	0.001	0.048

Table 1: (continue.)

p	n	D	BS	SD	AWB
$\mathfrak{R} = R_2$					
10	5	0.078	0.118	0.103	0.074
	10	0.054	0.091	0.058	0.090
20	10	0.052	0.096	0.042	0.078
	20	0.044	0.083	0.029	0.077
30	10	0.049	0.094	0.035	0.064
	15	0.044	0.082	0.023	0.062
	20	0.047	0.086	0.022	0.068
40	30	0.045	0.077	0.018	0.063
	10	0.053	0.096	0.030	0.061
	20	0.046	0.083	0.018	0.060
50	30	0.048	0.079	0.015	0.061
	40	0.047	0.077	0.014	0.060
	10	0.052	0.097	0.024	0.057
	20	0.050	0.086	0.016	0.058
60	25	0.046	0.078	0.014	0.055
	30	0.049	0.079	0.014	0.058
	40	0.046	0.073	0.012	0.058
	50	0.048	0.076	0.008	0.060
	10	0.052	0.092	0.025	0.054
	20	0.052	0.084	0.013	0.058
100	30	0.046	0.074	0.011	0.054
	40	0.045	0.074	0.009	0.053
	50	0.049	0.072	0.009	0.056
	60	0.048	0.075	0.009	0.057
200	10	0.050	0.093	0.017	0.047
	20	0.051	0.083	0.010	0.054
	50	0.049	0.073	0.005	0.053
400	10	0.056	0.097	0.012	0.047
	20	0.052	0.083	0.007	0.050
	50	0.048	0.074	0.003	0.050
400	10	0.052	0.096	0.008	0.041
	20	0.051	0.081	0.004	0.048
	50	0.050	0.075	0.001	0.050

Table 2: Attained significance level of DF and BSF^* under the null hypothesis when the covariance matrices are $\mathfrak{R} = R_S, \mathfrak{R} = R_T, \mathfrak{R} = R_1$ and $\mathfrak{R} = R_2$, respectively

p	n	$\mathfrak{R} = R_S$		$\mathfrak{R} = R_T$		$\mathfrak{R} = R_1$		$\mathfrak{R} = R_2$	
		DF	BSF^*	DF	BSF^*	DF	BSF^*	DF	BSF^*
10	5	0.056	0.056	0.067	0.071	0.058	0.074	0.063	0.073
	10	0.051	0.054	0.056	0.061	0.047	0.059	0.053	0.062
20	10	0.051	0.052	0.056	0.059	0.052	0.063	0.049	0.060
	20	0.053	0.057	0.049	0.053	0.050	0.056	0.043	0.051
30	10	0.054	0.053	0.054	0.058	0.051	0.060	0.049	0.058
	15	0.052	0.052	0.047	0.052	0.051	0.058	0.045	0.052
	20	0.056	0.057	0.047	0.053	0.051	0.057	0.047	0.053
40	30	0.047	0.049	0.046	0.051	0.047	0.049	0.046	0.051
	10	0.054	0.051	0.055	0.059	0.052	0.062	0.051	0.062
	20	0.055	0.056	0.047	0.053	0.050	0.054	0.046	0.052
50	30	0.051	0.053	0.047	0.052	0.049	0.052	0.048	0.052
	40	0.048	0.050	0.041	0.046	0.046	0.048	0.051	0.055
	10	0.055	0.051	0.057	0.061	0.051	0.061	0.051	0.062
	20	0.051	0.051	0.049	0.054	0.050	0.053	0.049	0.055
	25	0.048	0.049	0.044	0.050	0.051	0.054	0.050	0.053
60	30	0.049	0.050	0.045	0.049	0.050	0.053	0.048	0.051
	40	0.050	0.052	0.044	0.049	0.050	0.051	0.048	0.049
	50	0.049	0.052	0.040	0.045	0.047	0.048	0.049	0.050
	10	0.057	0.052	0.057	0.060	0.048	0.056	0.051	0.058
	20	0.050	0.049	0.050	0.056	0.055	0.059	0.049	0.054
	30	0.046	0.047	0.045	0.050	0.050	0.052	0.047	0.050
100	40	0.051	0.052	0.045	0.049	0.048	0.050	0.049	0.051
	50	0.051	0.053	0.044	0.051	0.052	0.053	0.047	0.047
	60	0.054	0.056	0.045	0.050	0.051	0.051	0.048	0.049
	10	0.054	0.050	0.058	0.061	0.053	0.062	0.051	0.061
	20	0.049	0.048	0.047	0.053	0.051	0.053	0.050	0.054
200	50	0.047	0.049	0.043	0.048	0.051	0.051	0.049	0.049
	10	0.052	0.046	0.059	0.063	0.054	0.063	0.053	0.063
	20	0.054	0.052	0.049	0.055	0.052	0.056	0.052	0.055
400	50	0.049	0.049	0.045	0.049	0.050	0.050	0.050	0.050
	10	0.052	0.044	0.058	0.062	0.055	0.064	0.051	0.058
	20	0.057	0.053	0.050	0.055	0.052	0.055	0.053	0.056
	50	0.050	0.050	0.044	0.049	0.052	0.052	0.049	0.049

Table 3: Empirical powers of DF and BSF^* under the alternative hypothesis when the covariance matrices are $\mathfrak{R} = R_S, \mathfrak{R} = R_T, \mathfrak{R} = R_1$ and $\mathfrak{R} = R_2$, respectively

p	n	$\mathfrak{R} = R_S$		$\mathfrak{R} = R_T$		$\mathfrak{R} = R_1$		$\mathfrak{R} = R_2$	
		DF	BSF^*	DF	BSF^*	DF	BSF^*	DF	BSF^*
10	5	0.721	0.738	0.475	0.499	0.328	0.394	0.227	0.261
	10	0.998	0.999	0.746	0.771	0.593	0.675	0.357	0.406
20	10	0.999	0.999	0.675	0.705	0.423	0.481	0.236	0.277
	20	1.000	1.000	0.994	0.996	0.942	0.954	0.603	0.643
30	10	1.000	1.000	0.852	0.872	0.625	0.681	0.337	0.384
	15	1.000	1.000	0.932	0.943	0.758	0.790	0.419	0.453
	20	1.000	1.000	0.991	0.994	0.939	0.952	0.509	0.543
	30	1.000	1.000	1.000	1.000	0.998	0.998	0.784	0.808
40	10	1.000	1.000	0.826	0.843	0.558	0.609	0.301	0.339
	20	1.000	1.000	0.983	0.988	0.859	0.877	0.471	0.498
	30	1.000	1.000	0.999	0.999	0.988	0.990	0.676	0.693
	40	1.000	1.000	1.000	1.000	0.997	0.997	0.776	0.786
50	10	1.000	1.000	0.833	0.851	0.544	0.589	0.303	0.339
	20	1.000	1.000	0.996	0.997	0.921	0.930	0.567	0.591
	25	1.000	1.000	1.000	1.000	0.988	0.990	0.702	0.720
	30	1.000	1.000	0.999	0.999	0.975	0.977	0.679	0.693
	40	1.000	1.000	1.000	1.000	1.000	1.000	0.868	0.874
	50	1.000	1.000	1.000	1.000	0.999	0.999	0.886	0.891
60	10	1.000	1.000	0.842	0.862	0.556	0.614	0.286	0.323
	20	1.000	1.000	0.996	0.998	0.904	0.916	0.549	0.570
	30	1.000	1.000	0.999	0.999	0.968	0.971	0.612	0.624
	40	1.000	1.000	1.000	1.000	0.999	0.999	0.864	0.870
	50	1.000	1.000	1.000	1.000	1.000	1.000	0.894	0.897
	60	1.000	1.000	1.000	1.000	1.000	1.000	0.969	0.969
100	10	1.000	1.000	0.871	0.890	0.589	0.632	0.328	0.362
	20	1.000	1.000	0.992	0.995	0.872	0.885	0.501	0.518
	50	1.000	1.000	1.000	1.000	1.000	1.000	0.886	0.887
200	10	1.000	1.000	0.844	0.863	0.545	0.583	0.303	0.335
	20	1.000	1.000	0.993	0.995	0.863	0.875	0.497	0.514
	50	1.000	1.000	1.000	1.000	1.000	1.000	0.872	0.871
400	10	1.000	1.000	0.849	0.869	0.560	0.599	0.299	0.329
	20	1.000	1.000	0.992	0.995	0.845	0.856	0.475	0.491
	50	1.000	1.000	1.000	1.000	1.000	1.000	0.904	0.902

Table 4: Observed p -values for testing the changes in gene expression after treatment for leukemia data.

	DF	BSF^*
Leukemia data:		
Statistic	115.733	10.7072
Average sum	438022	438022
p -values	0.0129	0.0000

References

1. Ahmad, M.R., Werner, C. and Brunner, E. 2008. Analysis of high-dimension repeated measures designs: The one sample case. *Computational Statistics and Data Analysis* 53 416-427.
2. Bartholomew, D.J. 1959a. A test of homogeneity for ordered alternatives. *Biometrika*, 46:36-48.
3. Bartholomew, D.J. 1959b. A test of homogeneity for ordered alternatives ii. *Biometrika*, 46:328-35.
4. Bartholomew, D.J. 1961. A test of homogeneity of means under restricted alternatives (with discussion) . *J.R. Statist. Soc(B)*, 23:239-81.
5. Bai, Z. and Saranadasa, H. 1996. Effect of high dimension: an example of a two sample problem, *Statist. Sinica* 6, 311-329.
6. Boyett, J.M. and Shuster, J.J. 1977. Nonparametric one-sided tests in multivariate analysis with medical applications. *Journal of the American Statistical Association*, 72:665-68.
7. Chongcharoen, S., Wright, F.T., and Singh, B. 2002. Powers of some one-sided multivariate tests with the population covariance matrix known up to a multiplicative constant. *Journal of Statistical Planning and Inference* 107, 103-121.
8. Chongcharoen, S. 2009. Powers of some one-sided multivariate tests with unknown population covariance matrix. *Songklanakarin Journal of Science and Technology*, 31(3), 351-359.
9. Dempster, A.P. 1958. A high dimensional two sample significance test, *Ann. Math. Statist.* 29, 995-1010.
10. Dempster, A.P. 1960. A significance test for the separation of two highly multivariate small samples, *Biometrics* 16, 41-50.
11. Follmann, D. 1996. A simple multivariate test for one-sided alternatives. *Journal of the American Statistical Association*, 91:854-61.
12. Kudo, A. 1963. A multivariate analogue of the one-sided test. *Biometrika*, 50:403-18.
13. Robertson, T., Wright, F.T., and Dykstra, R.L. 1988. *Ordered Restricted Statistical Inference*. John Wiley & Sons.
14. Perlman, M.D. 1969. One-sided problems in multivariate analysis. *Ann. Math. Statist.*, 40:549-67.
15. Shorack, G.R. 1967. Testing against ordered alternatives in model I analysis of variance: Normal theory and nonparametrics. *Ann. Math. Statist.*, 38:1740-53.
16. Siegfried Kropf , Jurgen Lauter, Daniela Kose and Dietrich von Rosen 2008. Comparison of exact parametric tests for high-dimensional data. *Computational Statistics and Data Analysis*.
17. Srivastava M.S. and Yasunori Fujikoshi 2006. Multivariate analysis of variance with fewer observations than the dimension. *The Journal of Multivariate Analysis* 97, 1927-1940.
18. Srivastava and Yanagihara 2010. Testing the equality of several covariance matrices with fewer observations than the dimension. *The Journal of Multivariate Analysis* 101, 1319-1329.

19. Srivastava, M.D. and Du, M. 2008. A test for mean vector with fewer observations than the dimension. *The Journal of Multivariate Analysis* 99, 386-402.
20. Tang, D.I., Gnecco, C. and Geller, N.L. 1989. An approximate likelihood ratio test for a normal mean vector with nonnegative components with application to clinical trials. *Biometrika*, 76:577-83.
21. Yasunori Fujikoshi, Tetsuto Himeno and Hirofumi Wakaki 2004. Asymptotic results of a high dimensional MANOVA test and power comparison when the dimension is larger compared to the sample size. *Journal of Japan Statistics Society*. Vol. 34 No.1, 19-26.
