**Nonparametric tests for testing equality of location parameters of two multivariate distributions**
By Chavan, Shirke

# Nonparametric tests for testing equality of location parameters of two multivariate distributions

Chavan A.R.*and Shirke D. T.

*Department of Statistics, Shivaji University, Kolhapur, India - 416004*

In this paper, we have proposed two nonparametric tests for testing the equality of location parameters of two multivariate distributions based on the notion of data depth. The proposed tests are extensions of the $M$-based test due to Li and Liu (2004). The performance of proposed tests has been assessed for symmetric as well as skewed multivariate distributions by simulation experiments. The tests have better performance in terms of power as compared to the $M$-based test and some of their competitors. The use of tests is illustrated with real life data.

**keywords:** Data depth, DD plot, Multivariate nonparametric tests, Location parameter, Permutation test.

## 1 Introduction

In several situations comparison between two data sets is required for number of reasons. The comparison can be based on the locations of these data sets. If multivariate data follow multivariate normal distribution then the task is easy as well known tests are available in the literature. However, if data do not follow multivariate normal distribution or we have no information about underlying distribution, nonparametric multivariate statistical methods are used to analyze data. One of the multivariate nonparametric statistical methods is based on the notion of the statistical data depth function, which was first introduced by Tukey (1975).

---

*Corresponding author: chavanatul2190@gmail.com

A data depth is a device for finding the location of multivariate data point with respect to a given data cloud. Larger depths are associated with more central points. Data depth gives a natural center-outward ranking to a multivariate data points with respect to data cloud. With the help of such rankings, Li and Liu (2004) proposed two depth-based nonparametric tests for multivariate location difference viz. $T$-based test and $M$-based test. These tests are developed using the Depth Depth (DD) plot (Liu et al., 1999). Dovoedo and Chakraborti (2015) have reported an extensive simulation study to evaluate the performance of these two tests for well known family of multivariate skewed distributions as well as multivariate symmetric distributions and compared performance of these tests for four popular affine-invariant depth functions, namely Mahalanobis depth, Spatial depth, Halfspace depth and Simplicial depth. We briefly discuss few of these in this article.

Several nonparametric tests have been proposed to deal with the multivariate two sample location problems as well as multi-sample location problems based on the concept of data depth. See Rousson (2002), Li et al. (2011), Chenouri and Small (2012) among others. Many of these methods are use permutation test to calculate the p-value.

In this paper we have proposed two nonparametric tests for testing equality of location parameters of two multivariate distributions based on the data depth, which are purely nonparametric. These tests are extensions of the $M$-based test introduced by Li and Liu (2004). Li and Liu (2004) use the most deepest point of two data clouds. We instead, consider some pre-specified number of most deepest points of the data clouds under comparison and construct tests based on these points. The performance of the proposed tests has been assessed by simulation experiments. The proposed tests give better performance in terms of power as compared to the $M$-based test and $T$-based test for symmetric as well as skewed multivariate distributions.

The rest of the paper is organized as follows. In section 2, we briefly discuss the notion of data depth, various data depth functions with their properties and DD plot. In section 3, we review the existing $T$-based and $M$-based tests of multivariate locations proposed by Li and Liu (2004). We describe the two new proposed nonparametric tests for testing the equality of locations using data depth in section 4. In section 5, we report simulation studies to compare performance of proposed tests with existing tests. In section 6, we apply the proposed tests to real life data. Section 7 contains some concluding remarks.

## 2 Statistical Data Depth Functions, Its Properties and DD Plot

### 2.1 Data Depth

Let $(X_1, X_2, ..., X_m)$ be a data set (cloud), where each $X_i \in \mathbb{R}^p$ is assumed to follow a continuous distribution with cumulative distribution function (CDF) $F(.)$, $i = 1, 2, ..., m$. Let $D(x, F)$ be the depth of a point $x$ with respect to $F$. A data depth is a function defined from $\mathbb{R}^p$ to $[0, \infty)$. Notion of data depth can be used to obtain the location of a given data points with respect to a data cloud. It measures the centrality of a given

data point with respect to a given data cloud. The deepest point using notion of data depth has the largest depth. Data depth gives a natural center-outward ranking to a data points with respect to data cloud. Such rankings were used for testing difference in location or scale parameters of two or more multivariate distributions, constructing nonparametric control charts, outlier detection and classification problem etc.

Tukey (1975) has first invented the word depth for picturing data. In literature, many different notions of data depth functions were proposed for capturing different probabilistic properties of multivariate data. Among them, the most popular choices of data depth functions are Mahalanobis depth (Mahalanobis, 1936), Simplicial depth (Liu, 1990), majority depth (Singh, 1991), half-space depth (Tukey, 1975), projection depth (Donoho and Gasko, 1992) etc. Some of these depth functions are reviewed in the following.

- **Mahalanobis Depth**

The Mahalanobis depth of a point $x \in \mathbb{R}^p$ with respect to $F$ on $\mathbb{R}^p$ is defined as,

$$MHD(x, F) = \frac{1}{(x-\mu_F)' \Sigma_F^{-1} (x-\mu_F)},$$

where $\mu_F$ is a location parameter or center and $\Sigma_F$ is the variance covariance matrix or dispersion matrix of $F$. The sample version of Mahalanobis depth can be obtained by replacing $\mu_F$ by $\bar{X}$ (sample mean) and $\Sigma_F$ by $S$ (sample variance covariance matrix).

- **Simplicial Depth**

The simplicial depth of a point $x \in \mathbb{R}^p$ with respect to $F$ on $\mathbb{R}^p$ is defined as,

$$SD(x, F) = Pr_F(s[X_1, X_2, ..., X_{p+1}] \ni x),$$

where $X_1, X_2, ..., X_{p+1}$ are independent and identically distributed observations from $F$ and $s[X_1, X_2, ..., X_{p+1}]$ is a closed simplex whose vertices are $X_1, X_2, ..., X_{p+1}$. The Sample version of simplicial depth can be obtained by replacing $F$ by $F_m$ in this expression. That is,

$$SD(x, F_m) = \binom{m}{p+1}^{-1} \sum_* I(x \epsilon s[X_{i1}, X_{i2}, ..., X_{ip+1}]),$$

where $(*)$ runs over all possible subsets of $X_1, X_2, ..., X_m$ of size $(p+1)$. Larger the depth $SD(x, F_m)$ indicates $x$ is contained in more simplices generated from the sample.

- **Tukey's Halfspace Depth**

Tukey's halfspace depth of a point $x \in \mathbb{R}^p$ with respect to probability measure $P$ on $\mathbb{R}^p$ is defined as the minimum probability mass carried by any closed half space containing $x$, that is,

$$HSD(x, F) = \inf_H \{P(H) : H \text{ is a closed halfspace containing x } \},$$

The sample version of $HSD(x, F)$ is obtained by replacing $F$ by $F_m$. If $k = 1$ then $HSD(x, F) = min\{F(x), 1 - F(x^-)\}$.

## 2.2 Properties of Depth Function

A depth function $D(x, F)$ is a non-negative function lies between $[0, \infty)$. According to Zuo and Serfling (2000), the depth function should satisfy the following four properties.

1. **Affine-invariance**: Suppose $x \in \mathbb{R}^p$ be a any given data point. Let A be any invertible matrix and $b \in \mathbb{R}^p$, then depth of a point $Ax + b$ with respect to $F$ is equal to the depth of a point with respect to $F$. That is, $D(Ax + b, F) = D(x, F)$.

2. **Maximality at a center**: If $F$ is centrally symmetric about $x_0 \in \mathbb{R}^p$, then depth of $x_0$ is the largest depth among all data points. That is,

$$D(x_0, F) \geq D(x, F) \text{ for any } x \in \mathbb{R}^p$$

3. **Monotonicity relative to any deepest point**: If $D(x_0, F) \geq D(x, F)$ for any $x \in \mathbb{R}^p$, then $D(x_0 + \lambda(x - x_0), F)$ is monotone non-increasing over $[0, \infty)$ for $\lambda \in [0, 1]$.

4. **Vanishing at infinity**: If $||x|| \longrightarrow \infty$ then $D(x, F) \longrightarrow 0$, where $||x||$ is the Euclidean norm in $\mathbb{R}^p$.

In the following section, we describe DD plot.

## 2.3 Depth-Depth Plot (DD Plot)

Let $(X_1, X_2, ..., X_m)$ and $(Y_1, Y_2, ..., Y_n)$ be two random samples from two continuous distributions $F$ and $G$ respectively, where $X_i, Y_j \in \mathbb{R}^p$, $i = 1, 2, ..., m$ and $j = 1, 2, ..., n$. Let $D(x, F)$ and $D(x, G)$ be the depths of a point $x \in Z$ with respect to $F$ and $G$ respectively, where $Z = X \cup Y$. Let

$$DD(F, G) = \{(D(x, F), D(x, G)), \quad \forall x \in Z\}.$$

The empirical version of $DD(F, G)$ based on the above described random samples is given by,

$$DD(F_m, G_n) = \{(D(x, F_m), D(x, G_n)), \quad \forall x \in Z\}.$$

DD plot is a two-dimensional graph, which is the plot of points in the set $DD(F_m, G_n)$. The DD plot can be used as a convenient diagnostic tool for graphical comparison of two multivariate samples. Difference in locations or scales or skewness or kurtosis are associated with different patterns observed on the DD plots. If $F = G$ then the points on the empirical DD Plot should fall on a $45^0$ line segment. This is illustrated in Figure 1(a), which is the DD plot of two multivariate samples drawn from the biariate normal distribution with mean vector $\underline{\mu} = \underline{0}$ and dispersion matrix $I_2$, where $I_2$ is the identity matrix of order two. That is $N_2(\underline{0}, I_2)$. The departure of $F$ from $G$ will indicate departure of points from $45^0$ line segment and Figure 1(b), Figure 2(a), Figure 2(b) and

Figure 3 reveal different patterns of DD plot that indicate the location differences, large location differences, scale differences and skewness differences (both location and scale differences) respectively. From Figure 1(b), the DD plot has a leaf-shaped figure with the cusp lying on the diagonal line towards the upper right corner and the leaf steam at the lower left corner point (0,0) when there is a shift in location parameters of two multivariate samples. In each of these Figures, we plot DD plot of DG against DF where $F$ and $G$ are chosen appropriately, where DF and DG are the depth of the points with respect to $F$ and $G$ respectively. We use Simplicial depth as a depth function to plot the DD plot in figure 1, 2 and 3. The study reported here is based on Simplicial depth function. The DD plots have been plotted using 'depth' package available in R (R Core Team, 2016).



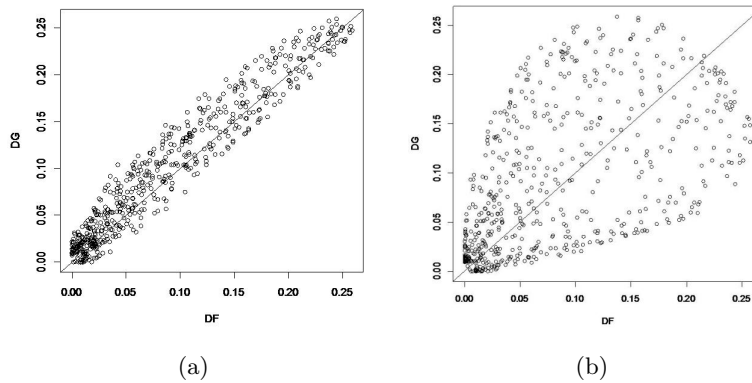(a)                                    (b)

Figure 1: DD plots of (a) $F = G = N_2(0, I_2)$ and (b) $F = N_2(0, I_2)$ and $G = N_2(0.5, I_2)$.



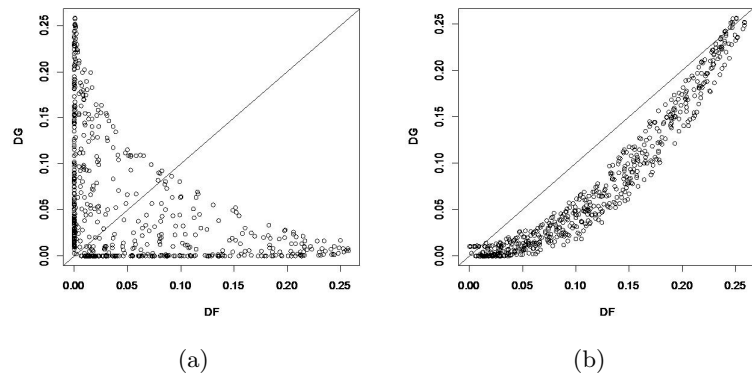(a)                                    (b)

Figure 2: DD plots of (a) $F = N_2(0, I_2)$ and $G = N_2(1.5, I_2)$ and (b) $F = N_2(0, I_2)$ and $G = N_2(0, 0.5I_2)$.
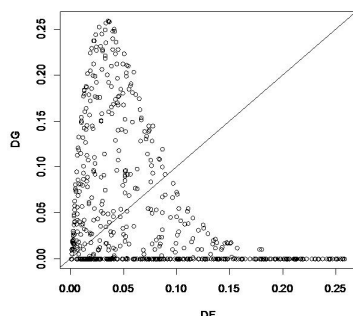
Figure 3: DD plot of $F = N_2(0, I_2)$ and $G = N_2(1, 0.1I_2)$.

In the following section, we describe $T$-based and $M$-based tests due to Li and Liu (2004).

# 3  $T$-based and $M$-based Tests

Li and Liu (2004) have proposed the $T$-based and the $M$-based tests for testing the equality of location parameters of two multivariate distributions by observing the DD plot introduced by Li and Liu (2004). These tests are completely nonparametric in nature.

Let $X = (X_1, X_2, ..., X_m)$ and $Y = (Y_1, Y_2, ..., Y_n)$, $X_i \in \mathbb{R}^p$, $Y_j \in \mathbb{R}^p$, $i = 1, 2, ..., m$, $j = 1, 2, ..., n$, be two data vectors observed from the distributions with CDF $F$ and $G$ respectively. Moreover, we assume that $F$ and $G$ are identical except for a possible location shift.

Let $\mu_1$ and $\mu_2$ be the location parameters of $F$ and $G$ respectively. The problem under consideration is to test

$$H_0 : \mu_1 = \mu_2 \quad \text{Vs} \quad H_1 : \mu_1 \neq \mu_2.$$

It is equivalent to test

$$H_0 : \theta = 0 \quad \text{Vs} \quad H_1 : \theta \neq 0,$$

where $\theta = \mu_1 - \mu_2$. That is $\theta$ is the shift in location parameters of two multivariate distributions.

## 3.1  The $T$-based Test

In the presence of location shift in two distribution, the DD plot has a leaf shaped figure (Figure 1(b), Figure 2(a)) with the leaf stem anchoring at the lower left corner point $(0, 0)$ and the cusp lying on the diagonal line pointing towards the upper right corner. On the basis of this observation, Li and Liu (2004) constructed the test statistic which

is the distance between the origin (0,0) and the cusp point. Li and Liu (2004) suggested the following procedure to calculate the distance between the cusp point and the origin (0,0).

For $(a_1, b_1)$ and $(a_2, b_2)$ in $\in \mathbb{R}^2$, define

$(a_1, b_1) \geq (a_2, b_2)$     if $a_1 \geq a_2$ and $b_1 \geq b_2$,
$(a_1, b_1) < (a_2, b_2)$     otherwise.

Define the set Q as

$Q = \{z \in X \cup Y : there\, does\, not\, exist\, w \in X \cup Y\, s.t.$

$$(D(w, F_m), D(w, G_n)) \geq (D(z, F_m), D(z, G_n))\}.$$

Then the cusp point is the point $(D(z_c, F_m), D(z_c, G_n))$ that satisfies $z_c \in Q$ and $|D(z_c, F_m) - D(z_c, G_n)| \leq |D(z, F_m) - D(z, G_n)|$ for all $z \in Q$. Let $T = (D(z_c, F_m) + D(z_c, G_n))/2$. The distance between the origin (0,0) and the cusp point is approximately $\sqrt{2}T$. Li and Liu (2004) used $T$ as a test statistic instead of using $\sqrt{2}T$ and smaller the value of $T$ indicates the larger shift in location. The p-value of the test is obtained by using the Fisher's permutation test. Let

$$P_B^T = \frac{\sum_{i=1}^{B} I_{(T_i^* \leq T_{obs})}}{B},$$

where $I(.)$ is the indicator function, $T_{obs}$ is the observed value of test statistic $T$ calculated from the original combined sample, $B$ is the number of times the combined sample $X \cup Y$ is permuted and $T_i^*$ is the value of test statistic $T$ corresponding to $i^{th}$ permuted combined sample, $i = 1, 2, ..., B$.

## 3.2 The $M$-based Test

Li and Liu (2004) developed another test for testing the equality of location parameters of two multivariate distributions based on the deepest point. In the theory of data depth, the location parameter is the point having maximum depth. Therefore if the two distributions $F$ and $G$ are identical then they should have the same deepest point. If there is a shift in location then the deepest point corresponding to the distribution $F$ would not be the deepest point corresponding to the distribution $G$. In fact, the deepest point of $F$ will have a smaller depth value with respect to $G$. $M$-based test statistic due to Li and Liu (2004) is given by,

$$M = min\{D(v_1, F_m), D(u_1, G_n)\},$$

where $v_1$ is the deepest point of $X \cup Y$ corresponding to $G_n$, and $u_1$ is the deepest point of $X \cup Y$ corresponding to $F_m$. Here larger the location difference, smaller the value of $M$. The p-value of the test is obtained by using the Fisher's permutation test. Let

$$P_B^M = \frac{\sum_{i=1}^{B} I_{(M_i^* \leq M_{obs})}}{B},$$

where $I(.)$ is the indicator function, $M_{obs}$ is the observed value of test statistic $M$ calculated from the original combined sample, $B$ is the number of times the combined sample $X \cup Y$ is permuted and $M_i^*$ is the value of test statistic $M$ corresponding to $i^{th}$ permuted combined sample, $i = 1, 2, ..., B$.

## 4 Proposed Tests

In the $M$-based test, Li and Liu (2004) consider only single deepest point for constructing the $M$-based test statistic. The test based on single deepest point considers a single data point. There is scope for improving the performance of this test by incorporating few more data points while constructing the test. This can be achieved by considering more than one deepest point. We propose the following two test statistic which are based on $k$ $(k \geq 2)$ deepest points for above hypothesis testing problem which can be considered as extensions of the previously discussed $M$-based test.

Suppose the set $U$ consists of the $k$ deepest points in $X \cup Y$ with respect to $F_m$ and the set $V$ consists of the $k$ deepest points in $X \cup Y$ with respect to $G_n$. Then we define two test statistic as follows,

- $M_1$-**based test statistic**

$$M_1 = min\{\tfrac{1}{k} \sum_{i=1}^{k} D(u_i, G_n), \tfrac{1}{k} \sum_{i=1}^{k} D(v_i, F_m)\},$$

- $M_2$-**based test statistic**

$$M_2 = \tfrac{1}{k} \sum_{i=1}^{k} (\min_i(D(u_i, G_n), D(v_i, F_m))),$$

where $u_i$ is the $i^{th}$ point of the set $U$ and $v_i$ is the $i^{th}$ point of the set $V$. Here for both of these two test statistic, larger the location difference, smaller the value of $M_1$ as well as of $M_2$. Therefore we propose two tests based on the above defined two statistic. Each test rejects $H_0$ for smaller value of the corresponding statistic.

The p-value of the proposed tests are obtained by using the Fisher's permutation test. Let

$$P_B^{M_1} = \frac{\sum_{i=1}^{B} I_{(M_{1i}^* \leq M_{1obs})}}{B},$$

where $I(.)$ and $B$ are defined as earlier, $M_{1obs}$ is the observed value of test statistic $M_1$ calculated from the original combined sample and $M_{1i}^*$ is the value of test statistic $M_1$ corresponding to $i^{th}$ permuted combined sample, $i = 1, 2, ..., B$. Similarly, we can calculate the p-value for test statistic $M_2$.

## 5 Performance of Tests

We have carried out extensive simulation study to assess the performance of two proposed tests, T-based, M-based and Hotelling $T^2$ tests for a bivariate data. The performance of proposed tests has been evaluated in terms of power for two Bivariate symmetric distributions (Bivariate normal, Bivariate Cauchy) as well as two Bivariate skewed distributions with pattern 1 and pattern 2 (Bivariate skew normal; Azzalini, 2005), bivariate skew-t

distribution (Azzalini and Capitanio, 2003). In the simulation study, the number of observations generated from each distribution F and G are taken to be m=n=100 and the original sample is permuted B=500 times. The power of $T$-based, $M$-based, Hotelling $T^2$, $M_1$-based and $M_2$-based tests are obtained by the proportion of the simulated p-values less than equal to the level of significance $\alpha = 0.05$. Here 1000 simulations are used for reporting the power and also results are reported for various values of k=2,3,4,5. Distributions used in the simulation study are listed in Table-1.

Table 1: Distributions used in the simulation study

| Distribution | Parameters |
|---|---|
| Symmetric normal | $N_2(\xi, \Omega = I)$ |
| Symmetric cauchy | $Cauchy(\xi, \Omega = I)$ |
| Skew-normal Pattern 1 | $SN_2(\xi, \Omega = I, a = (10, 4)^T)$ |
| Skew-normal Pattern 2 | $SN_2(\xi, \Omega = I, a = (4, 10)^T)$ |
| Skew-t Pattern 1 | $ST_2(\xi, \Omega = I, a = (10, 4)^T, v = 1)$ |
| Skew-t Pattern 2 | $ST_2(\xi, \Omega = I, a = (10, 4)^T, v = 3)$ |

The parameter $\xi$ denotes the location parameter, $\Omega$ denotes the dispersion parameter, $a$ denotes the shape parameter (or skewness parameter) and $v$ denotes the degrees of freedom. From all these distributions, the first random sample of size 100 is generated with parameter $\xi = (0,0)^T$ and dispersion parameter $\Omega$ is an identity matrix of order 2 and second random sample of size 100 is generated with parameter $\xi = (\mu, \mu)^T$ and dispersion parameter $\Omega$ is an identity matrix of order 2. Details regarding shape parameter $a$ and degrees of freedom $v$ are provided in Table-1. We provide powers of all these discussed tests for different values of $\mu = 0.0, 0.1, 0.2, 0.3, 0.4, 0.5$. R-software is used for simulation studies.

Table-8 provides powers for $T$-based, $M$-based, Hotelling $T^2$ and proposed tests when $F$ is bivariate Cauchy distribution with parameters $((0,0), I_2)$ and $G$ is bivariate normal distribution with parameters $((\mu_1, \mu_2), I_2)$ with sample sizes m=n=100 and Table-9 provides powers for $T$-based, $M$-based, Hotelling $T^2$ and proposed tests when $F$ is trivariate Cauchy distribution with parameters $((0,0,0), I_3)$ and $G$ is trivariate normal distribution with parameters $((\mu_1, \mu_2, \mu_3), I_3)$ with sample sizes m=n=50.

Table 2: Power comparison of $T$-based, $M$-based, Hotelling $T^2$ and proposed tests when underlying distribution is bivariate normal with sample sizes $m = n = 100$ for simplicial depth function.

|  |  | $\mu$ | 0.0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 |
|---|---|---|---|---|---|---|---|---|
|  | $T$-based | | 0.046 | 0.094 | 0.267 | 0.505 | 0.773 | 0.950 |
|  | $M$-based | | 0.046 | 0.106 | 0.282 | 0.547 | 0.810 | 0.954 |
|  | Hotelling $T^2$ | | 0.059 | 0.142 | 0.413 | 0.769 | 0.957 | 0.995 |
| k=2 | $M_1$-based | | 0.051 | 0.099 | 0.310 | 0.567 | 0.846 | 0.972 |
| | $M_2$-based | | 0.052 | 0.091 | 0.317 | 0.584 | 0.847 | 0.966 |
| k=3 | $M_1$-based | | 0.048 | 0.108 | 0.324 | 0.598 | 0.857 | 0.973 |
| | $M_2$-based | | 0.055 | 0.108 | 0.334 | 0.613 | 0.864 | 0.974 |
| k=4 | $M_1$-based | | 0.046 | 0.110 | 0.317 | 0.609 | 0.851 | 0.976 |
| | $M_2$-based | | 0.048 | 0.107 | 0.324 | 0.633 | 0.864 | 0.979 |
| k=5 | $M_1$-based | | 0.045 | 0.113 | 0.337 | 0.614 | 0.859 | 0.984 |
| | $M_2$-based | | 0.045 | 0.110 | 0.337 | 0.628 | 0.867 | 0.983 |

Table 3: Power comparison of $T$-based, $M$-based, Hotelling $T^2$ and proposed tests when underlying distribution is bivariate cauchy with sample sizes $m = n = 100$ for simplicial depth function.

|  |  | $\mu$ | 0.0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 |
|---|---|---|---|---|---|---|---|---|
|  | $T$-based | | 0.050 | 0.094 | 0.171 | 0.351 | 0.561 | 0.793 |
|  | $M$-based | | 0.058 | 0.090 | 0.169 | 0.366 | 0.568 | 0.801 |
|  | Hotelling $T^2$ | | 0.019 | 0.023 | 0.026 | 0.033 | 0.038 | 0.045 |
| k=2 | $M_1$-based | | 0.057 | 0.093 | 0.189 | 0.396 | 0.587 | 0.818 |
| | $M_2$-based | | 0.064 | 0.101 | 0.183 | 0.386 | 0.601 | 0.822 |
| k=3 | $M_1$-based | | 0.059 | 0.092 | 0.175 | 0.382 | 0.595 | 0.821 |
| | $M_2$-based | | 0.061 | 0.093 | 0.185 | 0.378 | 0.609 | 0.819 |
| k=4 | $M_1$-based | | 0.057 | 0.097 | 0.185 | 0.393 | 0.601 | 0.816 |
| | $M_2$-based | | 0.059 | 0.098 | 0.191 | 0.389 | 0.611 | 0.817 |
| k=5 | $M_1$-based | | 0.062 | 0.086 | 0.170 | 0.391 | 0.601 | 0.828 |
| | $M_2$-based | | 0.059 | 0.091 | 0.181 | 0.385 | 0.608 | 0.823 |

Table 4: Power comparison of $T$-based, $M$-based, Hotelling $T^2$ and proposed tests when underlying distribution is bivariate skew-normal distribution, pattern 1 with sample sizes $m = n = 100$ for simplicial depth function.

|  |  | $\mu$ | 0.0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 |
|---|---|---|---|---|---|---|---|---|
|  |  | $T$-based | 0.051 | 0.156 | 0.517 | 0.905 | 0.995 | 1.000 |
|  |  | $M$-based | 0.051 | 0.176 | 0.587 | 0.914 | 0.992 | 0.999 |
|  |  | Hotelling $T^2$ | 0.047 | 0.262 | 0.783 | 0.995 | 1.000 | 1.000 |
| k=2 |  | $M_1$-based | 0.052 | 0.181 | 0.641 | 0.946 | 0.998 | 1.000 |
|  |  | $M_2$-based | 0.049 | 0.189 | 0.659 | 0.952 | 0.998 | 1.000 |
| k=3 |  | $M_1$-based | 0.049 | 0.189 | 0.644 | 0.959 | 1.000 | 1.000 |
|  |  | $M_2$-based | 0.055 | 0.200 | 0.671 | 0.964 | 1.000 | 1.000 |
| k=4 |  | $M_1$-based | 0.046 | 0.190 | 0.656 | 0.967 | 1.000 | 1.000 |
|  |  | $M_2$-based | 0.048 | 0.201 | 0.684 | 0.972 | 1.000 | 1.000 |
| k=5 |  | $M_1$-based | 0.044 | 0.202 | 0.667 | 0.962 | 0.999 | 1.000 |
|  |  | $M_2$-based | 0.049 | 0.216 | 0.686 | 0.973 | 1.000 | 1.000 |

Table 5: Power comparison of $T$-based, $M$-based, Hotelling $T^2$ and proposed tests when underlying distribution is bivariate skew-normal distribution, pattern 2 with sample sizes $m = n = 100$ for simplicial depth function.

|  |  | $\mu$ | 0.0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 |
|---|---|---|---|---|---|---|---|---|
|  |  | $T$-based | 0.047 | 0.162 | 0.535 | 0.898 | 0.997 | 1.000 |
|  |  | $M$-based | 0.053 | 0.202 | 0.603 | 0.935 | 0.995 | 1.000 |
|  |  | Hotelling $T^2$ | 0.054 | 0.262 | 0.791 | 0.989 | 1.000 | 1.000 |
| k=2 |  | $M_1$-based | 0.044 | 0.220 | 0.654 | 0.951 | 0.998 | 1.000 |
|  |  | $M_2$-based | 0.044 | 0.231 | 0.670 | 0.954 | 0.998 | 1.000 |
| k=3 |  | $M_1$-based | 0.042 | 0.220 | 0.671 | 0.955 | 1.000 | 1.000 |
|  |  | $M_2$-based | 0.047 | 0.224 | 0.688 | 0.962 | 1.000 | 1.000 |
| k=4 |  | $M_1$-based | 0.042 | 0.218 | 0.668 | 0.957 | 1.000 | 1.000 |
|  |  | $M_2$-based | 0.044 | 0.237 | 0.685 | .968 | 1.000 | 1.000 |
| k=5 |  | $M_1$-based | 0.049 | 0.214 | 0.673 | 0.957 | 1.000 | 1.000 |
|  |  | $M_2$-based | 0.054 | 0.219 | 0.699 | 0.963 | 1.000 | 1.000 |

Table 6: Power comparison of $T$-based, $M$-based, Hotelling $T^2$ and proposed tests when underlying distribution is bivariate skew-t distribution, pattern 1 with sample sizes $m = n = 100$ for simplicial depth function.

|  |  | $\mu$ | 0.0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 |
|---|---|---|---|---|---|---|---|---|
|  | $T$-based |  | 0.040 | 0.119 | 0.353 | 0.740 | 0.938 | 0.991 |
|  | $M$-based |  | 0.049 | 0.147 | 0.451 | 0.811 | 0.961 | 0.999 |
|  | Hotelling $T^2$ |  | 0.040 | 0.128 | 0.281 | 0.561 | 0.801 | 0.928 |
| k=2 | $M_1$-based |  | 0.052 | 0.137 | 0.483 | 0.847 | 0.976 | 1.000 |
|  | $M_2$-based |  | 0.050 | 0.158 | 0.499 | 0.854 | 0.976 | 1.000 |
| k=3 | $M_1$-based |  | 0.060 | 0.170 | 0.513 | 0.867 | 0.982 | 0.999 |
|  | $M_2$-based |  | 0.052 | 0.177 | 0.527 | 0.882 | 0.985 | 0.999 |
| k=4 | $M_1$-based |  | 0.055 | 0.162 | 0.528 | 0.882 | 0.981 | 1.000 |
|  | $M_2$-based |  | 0.053 | 0.168 | 0.536 | 0.889 | 0.986 | 1.000 |
| k=5 | $M_1$-based |  | 0.055 | 0.155 | 0.520 | 0.903 | 0.988 | 1.000 |
|  | $M_2$-based |  | 0.051 | 0.170 | 0.548 | 0.905 | 0.989 | 1.000 |

Table 7: Power comparison of $T$-based, $M$-based, Hotelling $T^2$ and proposed tests when underlying distribution is bivariate skew-t distribution, pattern 2 with sample sizes $m = n = 100$ for simplicial depth function.

|  |  | $\mu$ | 0.0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 |
|---|---|---|---|---|---|---|---|---|
|  | $T$-based |  | 0.053 | 0.080 | 0.194 | 0.371 | 0.603 | 0.799 |
|  | $M$-based |  | 0.046 | 0.082 | 0.251 | 0.482 | 0.748 | 0.911 |
|  | Hotelling $T^2$ |  | 0.016 | 0.017 | 0.023 | 0.031 | 0.036 | 0.048 |
| k=2 | $M_1$-based |  | 0.054 | 0.097 | 0.272 | 0.541 | 0.804 | 0.940 |
|  | $M_2$-based |  | 0.052 | 0.092 | 0.274 | 0.537 | 0.805 | 0.944 |
| k=3 | $M_1$-based |  | 0.055 | 0.096 | 0.284 | 0.575 | 0.819 | 0.948 |
|  | $M_2$-based |  | 0.060 | 0.093 | 0.295 | 0.589 | 0.822 | 0.950 |
| k=4 | $M_1$-based |  | 0.051 | 0.095 | 0.308 | 0.584 | 0.841 | 0.952 |
|  | $M_2$-based |  | 0.055 | 0.085 | 0.300 | 0.593 | 0.840 | 0.955 |
| k=5 | $M_1$-based |  | 0.047 | 0.111 | 0.314 | 0.604 | 0.846 | 0.960 |
|  | $M_2$-based |  | 0.051 | 0.103 | 0.306 | 0.611 | 0.852 | 0.960 |

Table 8: Power comparison of $T$-based, $M$-based, Hotelling $T^2$ and proposed tests when $F : Cauchy((0,0), I_2)$ and $G : N_2((\mu_1, \mu_2), I_2)$ with sample sizes $m = n = 100$ for simplicial depth function.

| | $(\mu_1, \mu_2)$ | (0.1,0) | (0,0.2) | (0.1,0.2) | (0.3,0.3) |
|---|---|---|---|---|---|
| | $T$-based | 0.057 | 0.078 | 0.074 | 0.160 |
| | $M$-based | 0.087 | 0.107 | 0.135 | 0.277 |
| | Hotelling $T^2$ | 0.019 | 0.028 | 0.032 | 0.049 |
| k=2 | $M_1$-based | 0.101 | 0.145 | 0.166 | 0.346 |
| | $M_2$-based | 0.092 | 0.151 | 0.171 | 0.360 |
| k=3 | $M_1$-based | 0.115 | 0.151 | 0.185 | 0.402 |
| | $M_2$-based | 0.102 | 0.153 | 0.185 | 0.410 |
| k=4 | $M_1$-based | 0.126 | 0.177 | 0.206 | 0.438 |
| | $M_2$-based | 0.114 | 0.167 | 0.195 | 0.441 |
| k=5 | $M_1$-based | 0.129 | 0.190 | 0.221 | 0.469 |
| | $M_2$-based | 0.113 | 0.177 | 0.196 | 0.457 |

Table 9: Power comparison of $T$-based, $M$-based, Hotelling $T^2$ and proposed tests when $F : Cauchy((0,0,0), I_3)$ and $G : N_3((\mu_1, \mu_2, \mu_3), I_3)$ with sample sizes $m = n = 50$ for simplicial depth function.

| | $(\mu_1, \mu_2, \mu_3)$ | (0.0, 0.0, 0.1) | (0.0, 0.2, 0.0) | (0.0, 0.1, 0.2) | (0.1, 0.2, 0.3) |
|---|---|---|---|---|---|
| | $T$-based | 0.062 | 0.064 | 0.078 | 0.103 |
| | $M$-based | 0.096 | 0.096 | 0.138 | 0.173 |
| | Hotelling $T^2$ | 0.026 | 0.029 | 0.031 | 0.051 |
| k=2 | $M_1$-based | 0.107 | 0.143 | 0.151 | 0.221 |
| | $M_2$-based | 0.084 | 0.124 | 0.134 | 0.218 |
| k=3 | $M_1$-based | 0.097 | 0.134 | 0.151 | 0.215 |
| | $M_2$-based | 0.088 | 0.116 | 0.135 | 0.214 |
| k=4 | $M_1$-based | 0.107 | 0.120 | 0.148 | 0.197 |
| | $M_2$-based | 0.095 | 0.098 | 0.143 | 0.198 |
| k=5 | $M_1$-based | 0.096 | 0.110 | 0.134 | 0.192 |
| | $M_2$-based | 0.087 | 0.106 | 0.123 | 0.187 |

It is clear from the power comparison Table-2 to Table-9 that the proposed $M_1$-based and $M_2$-based tests give better performance in terms of power as compared to the $T$-based, $M$-based and Hotelling $T^2$ tests for skewed multivariate distributions as well as

multivariate cauchy distribution with a simplicial depth function. Proposed tests also give comparable results to Hotelling $T^2$, when the underlying distribution is bivariate normal. As such there is no criterion defined to choose an optimal value of $k$. However $k = 5$ appears to be reasonably good choice for majority of distributions. Between $M_1$-based and $M_2$-based tests, we recommend $M_2$-based test, as it has more power than $M_1$-based test for most of the distributions.

# 6 Application to Real Life Data

We consider Iris dataset (Fisher, 1936), which contains 150 observations each 50 for setosa, versicolor and virginica with four variables sepal length, sepal width, petal length and petal width. These are three populations corresponding to setosa, versicolor and virginica respectively. We select only two populations namely setosa and versicolor for illustration. The location parameters consists of values of sepal length, sepal width, petal length and petal width in the respective populations.

We are interested in testing equality of location parameters of these two populations. Multivariate normality test for setosa and versicolor data based on Shapiro test gives p-value 0.07906 and 0.00574 respectively. Therefore, sepal length, sepal width, petal length and petal width corresponding to versicolor population do not follow four variate normal distribution and Hotelling $T^2$ test is not appropriate in this case. Therefore, we use proposed tests to evaluate whether there is shift in location parameters of distribution of setosa and versicolor. The p-values for the proposed tests based on $B = 500$ permutations are reported in the following Table.

Table 10: $T$-based, $M$-based, $M_1$-based and $M_2$-based p-values for the Iris dataset based on $B = 500$ permutations using simplicial depth function

|       | Test          | p-value |
|-------|---------------|---------|
|       | $T$-based     | 0.000   |
|       | $M$-based     | 0.148   |
| k=2   | $M_1$-based   | 0.034   |
|       | $M_2$-based   | 0.036   |
| k=3   | $M_1$-based   | 0.014   |
|       | $M_2$-based   | 0.018   |
| k=4   | $M_1$-based   | 0.006   |
|       | $M_2$-based   | 0.008   |
| k=5   | $M_1$-based   | 0.002   |
|       | $M_2$-based   | 0.004   |

It is clear from the Table-10 that all the p-values of the proposed and $T$-based tests indicates that setosa and versicolor populations do not have same location but $M$-based test fails to conclude that setosa and versicolor populations do not have same location.

# 7 Conclusion

In this paper, we use data depth approach for comparing location parameters of two multivariate distributions. The proposed tests are purely nonparametric tests. They have a better performance in terms of power as compared to the existing $M$-Based and $T$-based test for symmetric as well as skewed multivariate distributions. Notion of data depth is useful for testing location and/or scale of two multivariate distributions.

# Acknowledgement

# References

Azzalini, A. (2005). The skew-normal distribution and related multivariate families. *Scandinavian Journal of Statistics*, 32(2):159–188.

Azzalini, A. and Capitanio, A. (2003). Distributions generated by perturbation of symmetry with emphasis on a multivariate skew t-distribution. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(2):367–389.

Chenouri, S. and Small, C. G. (2012). A nonparametric multivariate multisample test based on data depth. *Electronic Journal of Statistics*, 6:760–782.

Donoho, D. L. and Gasko, M. (1992). Breakdown properties of location estimates based on halfspace depth and projected outlyingness. *The Annals of Statistics*, pages 1803–1827.

Dovoedo, Y. and Chakraborti, S. (2015). Power of depth-based nonparametric tests for multivariate locations. *Journal of Statistical Computation and Simulation*, 85(10):1987–2006.

Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2):179–188.

Li, J., Ban, J., and Santiago, L. S. (2011). Nonparametric tests for homogeneity of species assemblages: a data depth approach. *Biometrics*, 67(4):1481–1488.

Li, J. and Liu, R. Y. (2004). New nonparametric tests of multivariate locations and scales using data depth. *Statistical Science*, pages 686–696.

Liu, R. Y. (1990). On a notion of data depth based on random simplices. *The Annals of Statistics*, 18(1):405–414.

Liu, R. Y., Parelius, J. M., and Singh, K. (1999). Multivariate analysis by data depth: descriptive statistics, graphics and inference,(with discussion and a rejoinder by liu and singh). *The annals of statistics*, 27(3):783–858.

Mahalanobis, P. (1936). Mahalanobis distance. In *Proceedings National Institute of Science of India*, volume 49, pages 234–256.

R Core Team (2016). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.

Rousson, V. (2002). On distribution-free tests for the multivariate two-sample location-scale model. *Journal of multivariate analysis*, 80(1):43–57.

Singh, K. (1991). A notion of majority depth. *Unpublished document.*

Tukey, J. W. (1975). Mathematics and the picturing of data. In *Proceedings of the international congress of mathematicians*, volume 2, pages 523–531.

Zuo, Y. and Serfling, R. (2000). General notions of statistical depth function. *Annals of statistics*, pages 461–482.