

# Comparison of Multivariate Distributions Using Quantile-Quantile Plots and Related Tests<sup>1</sup>

Subhra Sankar Dhar, Biman Chakraborty and Probal Chaudhuri

## Abstract

The univariate quantile-quantile (Q-Q) plot is a well-known graphical tool for examining whether two data sets are generated from the same distribution or not. It is also used to determine how well a specified probability distribution fits a given sample. In this article, we develop and study a multivariate version of Q-Q plot based on spatial quantiles. The usefulness of the proposed graphical device is illustrated on different real and simulated data, some of which are fairly high dimensional. We also develop certain statistical tests that are related to the proposed multivariate Q-Q plots and study their asymptotic properties. The performance of those tests are compared with that of some other well-known tests for multivariate distributions available in the literature.

**Keywords and phrases:** characterization of distributions, consistency of tests, contiguous alternatives, Pitman efficacy, quantile difference plots, tests for distributions.

## 1 Introduction

Univariate quantile-quantile (Q-Q) plot is a diagnostic tool, which is widely used to assess the distributional similarities and differences between two independent samples (see, e.g., Gnanadesikan and Wilk (1968), Gnanadesikan (1977) and Chambers, Cleveland, Kleiner and Tukey (1983)). As discussed in Doksum (1974), Doksum and Sievers (1976) and Koenker (2005, pp. 31–32), there are some fundamental connections between Q-Q plot and statistical inference in two-sample problems under some semi-parametric treatment effect model. Q-Q plot is also a popular device for checking the appropriateness of a specified probability distribution for a given univariate data. While univariate

---

<sup>1</sup>Subhra Sankar Dhar (dsubhra@gmail.com) is a Ph.D. student and Probal Chaudhuri (probal@isical.ac.in) is a Professor at Theoretical Statistics & Mathematics Unit, Indian Statistical Institute, Calcutta, India. Biman Chakraborty is a Lecturer in School of Mathematics, University of Birmingham, U.K.

Q-Q plot has a long history as a graphical tool for data analysis, there are only limited attempts in the literature to generalize Q-Q plot for multivariate samples.

Easton and McCulloch (1990) proposed a multivariate Q-Q plot based on the idea of matching a multivariate data set with a multivariate reference sample using an appropriate permutation of the data. Their procedure is based on the permutation of the data that leads to minimum sum of Euclidean distances between the matching data points from the two given samples. They assumed equality of the sizes for the two samples, and in order to assess how well a specified probability distribution fits a given multivariate sample, they used samples simulated from the specified distribution. In order to solve the optimization problem involved in matching the two samples, they used an iterative algorithm.

For bivariate samples, Marden (1998) proposed a version of Q-Q plot, which is based on drawing arrows from the spatial quantiles in one sample to the corresponding spatial quantiles in another sample in a two-sample problem (or to the corresponding spatial quantiles of a specified probability distribution in a one-sample problem). Here  $d$ -dimensional spatial quantile (see, e.g., Breckling and Chambers (1988), Chaudhuri (1996) and Koltchinskii (1997))  $Q_F(\mathbf{u}) = (Q_{F,1}(\mathbf{u}), \dots, Q_{F,d}(\mathbf{u}))$  is defined as  $Q_F(\mathbf{u}) = \arg \min_{Q \in \mathbb{R}^d} E\{\Phi(\mathbf{u}, \mathbf{x} - Q) - \Phi(\mathbf{u}, \mathbf{x})\}$ , where  $\mathbf{x}$  is a random vector with a probability distribution  $F$  on  $\mathbb{R}^d$ ,  $\Phi(\mathbf{u}, \mathbf{t}) = \|\mathbf{t}\| + \langle \mathbf{u}, \mathbf{t} \rangle$ ,  $\mathbf{u}$  is indexing  $d$ -dimensional multivariate quantiles by elements of the open unit ball  $B^d = \{\mathbf{u} : \mathbf{u} \in \mathbb{R}^d, \|\mathbf{u}\| < 1\}$ ,  $\langle \cdot, \cdot \rangle$  is the usual inner product, and  $\|\cdot\|$  is the norm induced by the inner product. If we have a random sample  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , where  $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$ , then the empirical spatial quantile  $Q_{\mathcal{X}}(\mathbf{u}) = (Q_{\mathcal{X},1}(\mathbf{u}), \dots, Q_{\mathcal{X},d}(\mathbf{u}))$  will be obtained by replacing  $F$  by its empirical distribution function  $F_n$ . Marden (1998) considered the case  $d = 2$  only, and he also mentioned a test based on arrow lengths. However, such arrow plots can be conveniently drawn and visualize only for bivariate data.

Friedman and Rafsky (1979, 1981) proposed a different procedure for distributional comparison of two multivariate samples. Their methodology is based on the idea of minimal spanning tree. Liu, Parelius and Singh (1998) proposed an alternative visual-

ization methodology, namely, DD-plot, for comparing two multivariate data sets based on the concept of data depth. However, none of the graphical tools developed by Friedman and Rafsky (1981) and Liu et al. (1998) will coincide with the usual univariate Q-Q plot when they are applied to univariate data even though the univariate versions of those graphical tools are closely related to the usual Q-Q plot. Hence, strictly speaking, none of them can be taken as a natural multivariate extension of Q-Q plot.

In this article, we propose an extension of Q-Q plot using spatial quantiles for multivariate data. As we will see in the subsequent sections, these Q-Q plots are in many ways natural generalizations of univariate Q-Q plots. In particular, in the case of  $d$ -dimensional multivariate data, there will be  $d$  two-dimensional plots of points, where the points cluster around lines if and only if the two multivariate distributions under comparison are related to each other by location and homogeneous scale changes – and the slopes and the intercepts of those lines relate to the location shifts and the scale changes in the same way as in the case of univariate Q-Q plots. Here it is appropriate to point out that we have considered spatial quantiles instead of marginal quantiles (see Babu and Rao (1988)) or any other  $l_p$ -quantile (see Chakraborty (2001)) for constructing multivariate Q-Q plots because it is known that the spatial quantiles characterize the multivariate distributions (see Koltchinskii (1997)), and such a characterization property is not available for marginal quantile or any other  $l_p$ -quantile.

Unlike Friedman and Rafsky (1979, 1981) and Marden (1998), neither Liu et al. (1998) nor Easton and McCulloch (1990) mentioned any statistical test for comparing distributions of two multivariate samples. On the other hand, several distributional tests for univariate data, namely, Kolmogorov-Smirnov test, Cramer-von-Mises test (see, e.g., Serfling (1980)), Shapiro-Wilk test (see Shapiro and Wilk (1965)) and Anderson-Darling test (Anderson and Darling (1954)), etc. have been discussed in the literature. Multivariate extensions of a few of these tests have also been studied in the literature (see, e.g., Bickel (1969), Shorack and Wellner (1986), and Justel, Pena and Zamar (1997)). We also propose and study some statistical tests for distributions, which are motivated by our multivariate Q-Q plots. In our numerical and asymptotic study, such

tests turn out to have either comparable or superior performance when compared with other tests for distributions mentioned above.

## 2 Construction of multivariate quantile-quantile plots

Recall that if one has two univariate data sets with sample sizes  $n$  and  $m$ , in order to construct the Q-Q plot, one has to estimate the  $i/n$ -th quantile and the  $j/m$ -th quantile for  $i = 1, \dots, n$  and  $j = 1, \dots, m$  from each of the two data sets, and then plot the quantiles of one data set against the corresponding quantiles of another data set. In other words, here one has to match the quantiles of one data set with the corresponding quantiles of another data set, and for univariate data, this is done in a straightforward way. Note that when  $n \neq m$ , one has to compute  $(n + m)$  quantiles from each of the two data sets whereas when  $n = m$ , it is enough to compute only  $n$  quantiles from each data set. However, the problem of matching the quantile vectors for two multivariate data sets does not have any obvious solution. We have mentioned in the Introduction that Easton and McCulloch (1990) made some attempts to solve this matching problem for multivariate data. As we have already pointed out, there are several limitations of the procedure proposed by them, which is computationally quite complex. In this paper, we propose a solution to this matching problem, which is computationally much simpler, and it is partly motivated by the choices for the bases and the tips of the arrows used by Marden (1998) in his arrow plots.

Let  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  and  $\mathcal{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  be two independent  $d$ -dimensional data sets, where  $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$  and  $\mathbf{y}_j = (y_{j,1}, \dots, y_{j,d})$ . Suppose that  $\mathbf{u}_1, \dots, \mathbf{u}_n$  and  $\mathbf{u}_{n+1}, \dots, \mathbf{u}_{n+m}$  are the spatial ranks of the data sets  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively, where spatial rank of  $\mathbf{z} \in \mathbb{R}^d$  with respect to the data cloud formed by the observations  $\mathbf{x}_1, \dots, \mathbf{x}_n$  (or  $\mathbf{y}_1, \dots, \mathbf{y}_m$ ) is defined as  $n^{-1} \sum_{i:\mathbf{x}_i \neq \mathbf{z}} \|\mathbf{z} - \mathbf{x}_i\|^{-1} (\mathbf{z} - \mathbf{x}_i)$  (or  $m^{-1} \sum_{i:\mathbf{y}_i \neq \mathbf{z}} \|\mathbf{z} - \mathbf{y}_i\|^{-1} (\mathbf{z} - \mathbf{y}_i)$ ) (see, e.g., Mottonen and Oja (1995), Chaudhuri (1996) and Serfling (2004)). Note that here  $Q_{\mathcal{X}}(\mathbf{u}_k) = \mathbf{x}_k$  for  $k = 1, \dots, n$ , and  $Q_{\mathcal{Y}}(\mathbf{u}_k) = \mathbf{y}_k$  for  $k = n + 1, \dots, n + m$ , where  $Q_{\mathcal{X}}$  and  $Q_{\mathcal{Y}}$  are empirical spatial quantiles with respect to data cloud  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively. Now, we can match the two sets of quantiles by setting correspondence

between  $Q_{\mathcal{X}}(\mathbf{u}_k)$  and  $Q_{\mathcal{Y}}(\mathbf{u}_k)$  for all  $k = 1, \dots, n + m$  and consider the set of  $(n + m)$  points in  $\mathbb{R}^{2d}$  defined as  $S_{n,m}(\mathcal{X}, \mathcal{Y}) = \{(Q_{\mathcal{X}}(\mathbf{u}_k), Q_{\mathcal{Y}}(\mathbf{u}_k)) : k = 1, \dots, (n + m)\}$ . Note that when  $d = 1$ , our proposed multivariate matching coincides with the usual way of matching univariate quantiles, and the set of points  $S_{n,m}(\mathcal{X}, \mathcal{Y})$  is same as the set of points used in a univariate Q-Q plot.

We next describe the construction of Q-Q plots for a one-sample multivariate problem involving a  $d$ -dimensional data set  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , where  $\mathbf{x}_i = (x_{i,1}, \dots, x_{i,d})$  and a specified  $d$ -dimensional probability distribution  $F_0$ . Let  $\mathbf{u}_1, \dots, \mathbf{u}_n$  be the spatial ranks of the data points  $\mathbf{x}_i, i = 1, \dots, n$ . Suppose that  $Q_{F_0}(\mathbf{u}_k) = (Q_{F_0,1}(\mathbf{u}_k), \dots, Q_{F_0,d}(\mathbf{u}_k))$  is the  $\mathbf{u}_k$ -th spatial quantile of the specified distribution  $F_0$ , where  $k = 1, \dots, n$ . As in the preceding case, since  $Q_{\mathcal{X}}(\mathbf{u}_k) = \mathbf{x}_k$ , a natural way of matching quantiles of the data set with those of the specified probability distribution will be by setting the correspondence between  $\mathbf{x}_k$  and  $Q_{F_0}(\mathbf{u}_k)$  and forming the set  $S_n(\mathcal{X}, F_0) = \{(\mathbf{x}_k, Q_{F_0}(\mathbf{u}_k)) : k = 1, \dots, n\}$ , which is a subset of  $\mathbb{R}^{2d}$ . In particular, when  $d = 1$ ,  $S_n(\mathcal{X}, F_0)$  coincides with the set of points used in the construction of the univariate Q-Q plot for one-sample problem.

Theorem 2.1 establishes a necessary and sufficient condition for the points in  $S_{n,m}(\mathcal{X}, \mathcal{Y})$  and  $S_n(\mathcal{X}, F_0)$ , both of which consists of points from  $2d$ -dimensional Euclidean space  $\mathbb{R}^{2d}$ , to asymptotically cluster around some appropriate lower dimensional hyperplanes defined by  $d$  linear equations as the sample sizes grow to infinity.

**Theorem 2.1:** *Suppose that  $F$  and  $G$  have continuous and positive density functions, and  $\mathcal{X}$  and  $\mathcal{Y}$  are two independent data sets. Let  $n, m \rightarrow \infty$  in such a way that*

$$\lim_{n,m \rightarrow \infty} \frac{n}{(n + m)} \in (0, 1). \text{ Then, for every } \epsilon > 0 \text{ and } \delta > 0, \text{ we have}$$

$$\lim_{n,m \rightarrow \infty} P [S_{n,m}(\mathcal{X}, \mathcal{Y}, \delta) \in \{(\mathbf{v} \in \mathbb{R}^d, \mathbf{w} \in \mathbb{R}^d) : \|\mathbf{v} - \sigma \mathbf{w} - \boldsymbol{\mu}\| < \epsilon\}] = 1$$

*if and only if  $F(\mathbf{x}) = G((\mathbf{x} - \boldsymbol{\mu})/\sigma)$ , where  $\boldsymbol{\mu} \in \mathbb{R}^d$ ,  $\sigma > 0$ , and  $S_{n,m}(\mathcal{X}, \mathcal{Y}, \delta) = \{(Q_{\mathcal{X}}(\mathbf{u}_k), Q_{\mathcal{Y}}(\mathbf{u}_k)) : \|\mathbf{u}_k\| \leq 1 - \delta, k = 1, \dots, (n + m)\}$ . Further, for a one-sample problem, where  $F_0$  is a specified distribution function having a positive and continuous density, for every  $\epsilon > 0$  and  $\delta > 0$ , we have*

$$\lim_{n \rightarrow \infty} P [S_n(\mathcal{X}, F_0, \delta) \in \{(\mathbf{v} \in \mathbb{R}^d, \mathbf{w} \in \mathbb{R}^d) : \|\mathbf{v} - \sigma \mathbf{w} - \boldsymbol{\mu}\| < \epsilon\}] = 1$$

if and only if  $F(\mathbf{x}) = F_0((\mathbf{x} - \boldsymbol{\mu})/\sigma)$ , where  $\boldsymbol{\mu} \in \mathbb{R}^d$ ,  $\sigma > 0$  and  $S_n(\mathcal{X}, F_0, \delta) = \{(\mathbf{x}_k, Q_{F_0}(\mathbf{u}_k)) : \|\mathbf{u}_k\| \leq 1 - \delta, k = 1, \dots, n\}$ .

A major implication of Theorem 2.1 is that  $S_{n,m}(\mathcal{X}, \mathcal{Y})$  can be used, just like the univariate Q-Q plots, to determine whether the two multivariate samples have distributions that differ only in the location and the scale of the variables. Similarly, in the case of one-sample problem,  $S_n(\mathcal{X}, F_0)$  can be used to decide whether a given multivariate sample is generated from a specified probability distribution after some location and scale transformations.

We now propose our Q-Q plot for a  $d$ -dimensional multivariate data set as a collection of  $d$  two dimensional plots, where each plot corresponds to the graph for each component of the multivariate quantile. Then, as in the case of usual univariate Q-Q plots, Theorem 2.1 ensures that for all  $i = 1, \dots, d$ , the points in the  $i$ -th 2-dimensional plot will lie close to a straight line with slope  $\sigma$  and intercept  $\mu_i$  if and only if the data are obtained from the distributions  $F$  and  $G$  such that  $G(\mathbf{x}) = F((\mathbf{x} - \boldsymbol{\mu})/\sigma)$  (or,  $F(\mathbf{x}) = F_0((\mathbf{x} - \boldsymbol{\mu})/\sigma)$  for one-sample problem), where  $\mu_i$  is the  $i$ -th component of  $\boldsymbol{\mu}$ . Here it would be appropriate to point out that the “if part” of the above assertion is valid even if one uses marginal quantiles instead of spatial quantiles for the construction of  $S_{n,m}(\mathcal{X}, \mathcal{Y})$  or  $S_n(\mathcal{X}, F_0)$ . However, the “only if part” of the assertion will not hold in that case. This is the main reason to consider spatial quantiles instead of marginal quantiles for constructing multivariate Q-Q plots.

When  $d$  is large, it is not convenient to display and visually examine  $d$  different scatter plots, and we can plot  $(l, Q_{\mathcal{X},l}(\mathbf{u}_k) - Q_{\mathcal{Y},l}(\mathbf{u}_k))$  for all  $k = 1, \dots, (n + m)$ ,  $l = 1, \dots, d$ . Note that, in this way, we get a single two dimensional plot, where there are  $d$  vertical lines parallel to one another, and points are plotted along those lines. In Sections 4.1 and 4.2, we will demonstrate our Q-Q plots using some examples of real and simulated data, some of which are fairly high dimensional.

### 3 Tests for comparing multivariate distributions

The plots based on differences of quantiles for high dimensional data along with the characterization of distributions by spatial quantiles (Koltchinskii (1997)) motivated us to consider some tests for distributions based on Euclidean norm of the differences of spatial quantiles in one-and two-sample problems. Let  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be i.i.d. observations from an unknown distribution  $F$ , and suppose that we want to test  $H_0 : F = F_0 (\Leftrightarrow Q_F(\mathbf{u}) = Q_{F_0}(\mathbf{u}) \text{ for all } \mathbf{u})$  against the alternatives  $H_1 : F \neq F_0 (\Leftrightarrow Q_F(\mathbf{u}) \neq Q_{F_0}(\mathbf{u}) \text{ for some } \mathbf{u})$ , where  $F_0$  is a specified distribution with a continuous probability density function. In order to test  $H_0$  against  $H_1$ , we can use the test statistic  $V_n = n \int_{\mathbf{u} \in S(\delta)} \|\hat{Q}_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\|^2 d\mathbf{u}$ , where  $S(\delta) = \{\mathbf{u} : \|\mathbf{u}\| \leq 1 - \delta\}$ . The asymptotic behavior of the test based on  $V_n$  is stated in Theorem 3.1.

**Theorem 3.1:** *Let  $Z_1(\mathbf{u})$  be a Gaussian process with zero mean and covariance kernel*

$$\begin{aligned} k_1(\mathbf{u}_1, \mathbf{u}_2) &= [D_1^{F_0}\{Q(\mathbf{u}_1)\}]^{-1} [D_2^{F_0}\{Q(\mathbf{u}_1), Q(\mathbf{u}_2), \mathbf{u}_1, \mathbf{u}_2\}] [D_1^{F_0}\{Q(\mathbf{u}_2)\}]^{-1} \text{ if } d \geq 2 \\ &= \frac{(\frac{u_1+1}{2})(1 - \frac{u_2+1}{2})}{f_0(F_0^{-1}(\frac{u_1+1}{2}))f_0(F_0^{-1}(\frac{u_2+1}{2}))} \text{ if } d = 1 \text{ and for any } u_1 < u_2. \end{aligned}$$

Here  $D_1^{F_0}\{Q(\mathbf{u})\} = E_{F_0}[\|\mathbf{X} - Q(\mathbf{u})\|_2^{-1}\{I_d - \|\mathbf{X} - Q(\mathbf{u})\|_2^{-2}(\mathbf{X} - Q(\mathbf{u}))(\mathbf{X} - Q(\mathbf{u}))^T\}]$ ,  $D_2^{F_0}\{Q_1(\mathbf{u}), Q_2(\mathbf{v}), \mathbf{u}, \mathbf{v}\} = E_{F_0}[\{\|\mathbf{X} - Q_1(\mathbf{u})\|_2^{-1}(\mathbf{X} - Q_1(\mathbf{u})) + \mathbf{u}\}\{\|\mathbf{X} - Q_2(\mathbf{v})\|_2^{-1}(\mathbf{X} - Q_2(\mathbf{v})) + \mathbf{v}\}^T]$ , and  $f_0$  is the density function of  $F_0$ . A test based on  $V_n$ , which rejects  $H_0$  when  $V_n > c_1(\alpha)$ , where  $c_1(\alpha)$  is the  $(1 - \alpha)$ -th quantile ( $0 < \alpha < 1$ ) of the distribution of  $\int_{\mathbf{u} \in S(\delta)} \|Z_1(\mathbf{u})\|^2 d\mathbf{u}$  will have asymptotic size  $\alpha$ , and for  $0 < \delta < 1$  appropriately large, it will be a consistent test in the sense that when  $F \neq F_0$ , the asymptotic power of the test will be 1.

In order to implement our test, we need to compute  $V_n$ , and one simple and convenient way to evaluate the integral that appears in the test statistic is by averaging of the integrands over i.i.d. Monte-Carlo replications obtained from random generations of  $\mathbf{u}$  from the uniform distribution on  $S(\delta)$ . Now, for computing  $c_1(\alpha)$ , one can simulate observations from the Gaussian process  $Z_1(\mathbf{u})$  and use the empirical approximation

for the distribution of  $\int_{\mathbf{u} \in S(\delta)} \|Z_1(\mathbf{u})\|^2 d\mathbf{u}$ . Alternatively, in view of the asymptotically Gaussian distribution of the process  $\{\hat{Q}_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\}$  for  $\mathbf{u} \in S(\delta)$  and the well-known orthogonal decomposition of a finite dimensional multivariate normal distribution, the distribution of the test statistics  $V_n$  under the null hypothesis can be approximated by a weighted sum of chi-square random variables each with one degree of freedom. In our numerical work, we have computed  $c_1(\alpha)$  by generating several Monte-Carlo replications from a weighted sum of chi-square variables with weights estimated from the empirical covariance kernel of the corresponding Gaussian process.

Let us next consider a two-sample problem with two independent sets of i.i.d. observations  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  and  $\mathcal{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_m\}$  from unknown distributions  $F$  and  $G$  having continuous probability density functions defined on  $\mathbb{R}^d$ , respectively, and we want to test  $H_0^* : F = G (\Leftrightarrow Q_F(\mathbf{u}) = Q_G(\mathbf{u}) \text{ for all } \mathbf{u})$  against  $H_1^* : F \neq G (\Leftrightarrow Q_F(\mathbf{u}) \neq Q_G(\mathbf{u}) \text{ for some } \mathbf{u})$ . In order to test  $H_0^*$  against  $H_1^*$ , one can use the test statistic  $T_{n,m} = (n + m) \int_{\mathbf{u} \in S(\delta)} \|\hat{Q}_F(\mathbf{u}) - \hat{Q}_G(\mathbf{u})\|^2 d\mathbf{u}$ . The asymptotic behavior of the test based on  $T_{n,m}$  is stated in Theorem 3.2.

**Theorem 3.2:** *Suppose that  $n, m \rightarrow \infty$  in such a way that  $\lim_{n,m \rightarrow \infty} \frac{n}{n+m} = \lambda \in (0, 1)$ . Let  $Z_2(\mathbf{u})$  be a Gaussian process with mean function 0 and covariance kernel*

$$\begin{aligned} k_2(\mathbf{u}_1, \mathbf{u}_2) &= \frac{[D_1^F\{Q(\mathbf{u}_1)\}]^{-1} [D_2^F\{Q(\mathbf{u}_1), Q(\mathbf{u}_2), \mathbf{u}_1, \mathbf{u}_2\}] [D_1^F\{Q(\mathbf{u}_2)\}]^{-1}}{\lambda(1-\lambda)} && \text{if } d \geq 2 \\ &= \frac{(\frac{u_1+1}{2})(1 - \frac{u_2+1}{2})}{\lambda(1-\lambda)f(F^{-1}(\frac{u_1+1}{2}))f(F^{-1}(\frac{u_2+1}{2}))} && \text{if } d = 1 \text{ and for any } u_1 < u_2. \end{aligned}$$

Here  $D_1^F\{Q(\mathbf{u})\} = E_F[\|\mathbf{X} - Q(\mathbf{u})\|_2^{-1} \{I_d - \|\mathbf{X} - Q(\mathbf{u})\|_2^{-2} (\mathbf{X} - Q(\mathbf{u}))(\mathbf{X} - Q(\mathbf{u}))^T\}]$ ,  $D_2^F\{Q_1(\mathbf{u}), Q_2(\mathbf{v}), \mathbf{u}, \mathbf{v}\} = E_F[\{\|\mathbf{X} - Q_1(\mathbf{u})\|_2^{-1} (\mathbf{X} - Q_1(\mathbf{u})) + \mathbf{u}\} \{\|\mathbf{X} - Q_2(\mathbf{v})\|_2^{-1} (\mathbf{X} - Q_2(\mathbf{v})) + \mathbf{v}\}^T]$ , and  $f$  is the density function of  $F$ . A test based on  $T_{n,m}$ , which rejects  $H_0^*$  when  $T_{n,m} > c_2(\alpha)$ , where  $c_2(\alpha)$  is the  $(1-\alpha)$ -th quantile ( $0 < \alpha < 1$ ) of the distribution of  $\int_{\mathbf{u} \in S(\delta)} \|Z_2(\mathbf{u})\|^2 d\mathbf{u}$ , will have asymptotic size  $\alpha$ , and for  $0 < \delta < 1$  appropriately large, it will be a consistent test in the sense that when  $F \neq G$ , its asymptotic power will be 1.

For numerical work, one can compute  $T_{n,m}$  and the critical value  $c_2(\alpha)$  for the two-

sample problem in a similar way as in the case of  $V_n$  and  $c_1(\alpha)$  in the one-sample problem. In Section 4.3, we have compared the performance of our tests with that of some other well-known tests studied in the literature.

## 4 Examples and Numerical results

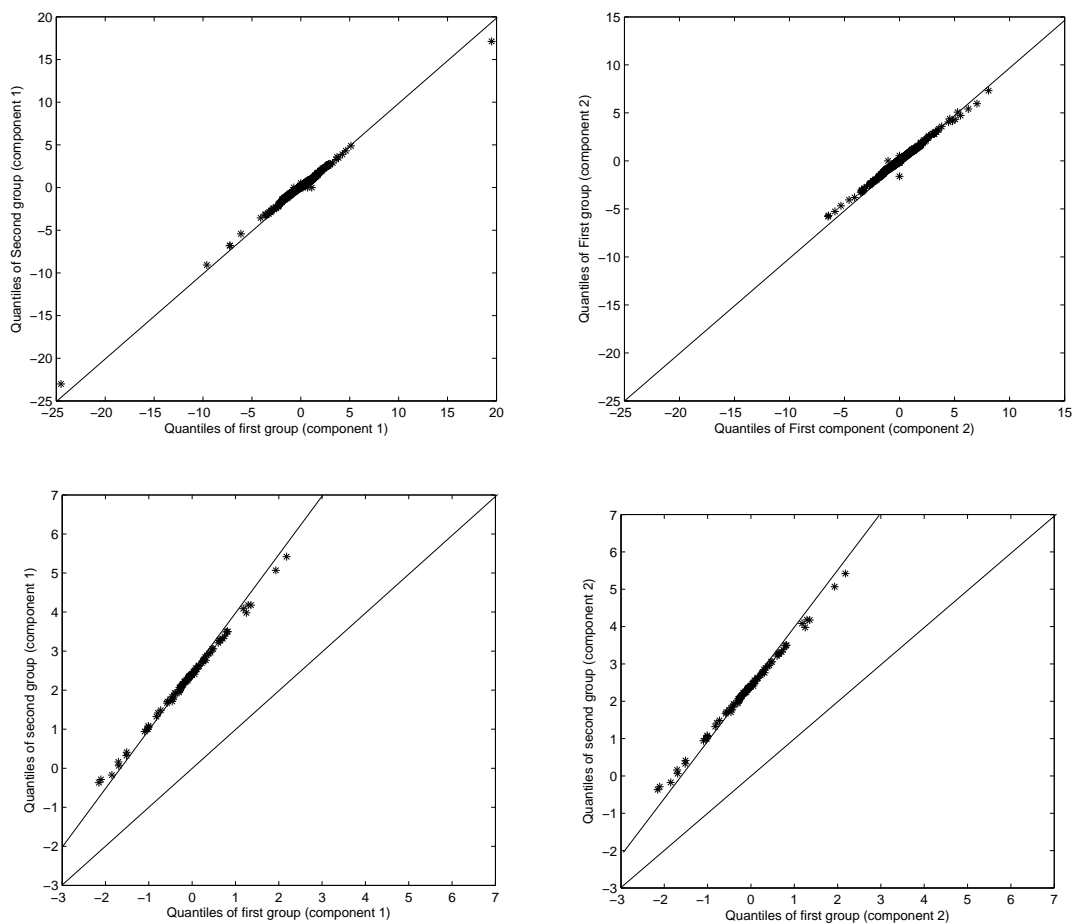
### 4.1 Comparing multivariate distributions using Q-Q plots

Here we demonstrate our methodology of Q-Q plots, which we have discussed in Section 2, using some examples. We generated 100 i.i.d. observations from  $F = N_2(\boldsymbol{\mu}_1, \Sigma_1)$  and  $G = N_2(\boldsymbol{\mu}_2, \Sigma_2)$  distributions, where  $N_2(\boldsymbol{\mu}, \Sigma)$  denotes the bivariate normal distribution with location parameter  $\boldsymbol{\mu}$  and scatter matrix  $\Sigma$ . First, we considered  $\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2 = (0, 0)$ , and  $\Sigma_1 = \Sigma_2 = I_2 = 2 \times 2$  identity matrix. Next, we considered  $\boldsymbol{\mu}_1 = (0, 0)$ ,  $\boldsymbol{\mu}_2 = (2, 2)$ ,  $\Sigma_1 = I_2$ , and  $\Sigma_2 = 2^2 I_2$ , i.e.,  $F$  and  $G$  are related with each other by a location and a scale transformation. The Q-Q plots for the above cases are displayed in Figure 1. As expected, in each plot in the first row of Figure 1, the points are clustered tightly around the  $45^\circ$  line passing through the origin, while in each of the plots in the second row, the points are tightly clustered around a line through the point  $(0, 2)$  having slope 2.

We now consider the reference distribution  $F_0 = N_2(\mathbf{0}, I_2)$ . Marden (1998) provided a result for computing spatial quantiles of the multivariate normal distribution, and we have used that for our calculation. We generated 100 i.i.d. observations from each of  $N_2(\boldsymbol{\mu}, \Sigma)$  and  $C_2(\boldsymbol{\mu}, \Sigma)$ , where  $N_2(\boldsymbol{\mu}, \Sigma)$  and  $C_2(\boldsymbol{\mu}, \Sigma)$  denote the bivariate normal and Cauchy (with p.d.f.  $f(\mathbf{x}) = c\{1 + (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\}^{-1}$ ) distributions having location parameter  $\boldsymbol{\mu}$  and scatter matrix  $\Sigma$ , respectively. In the following Q-Q plots, we consider  $\boldsymbol{\mu} = \mathbf{0}$  and  $\Sigma = I_2$ .

It is clearly evident from the diagrams in the first row of Figure 2 that the reference distribution (i.e., standard bivariate normal) fits the data well as the points in those diagrams are tightly clustered around the  $45^\circ$  line passing through the origin. On the other hand, in each scatter plot in the second row, the points are significantly deviating

Figure 1: Q-Q plots for two-sample problems when samples are generated from the same distribution (first row), and from distributions having different locations and scales (second row).

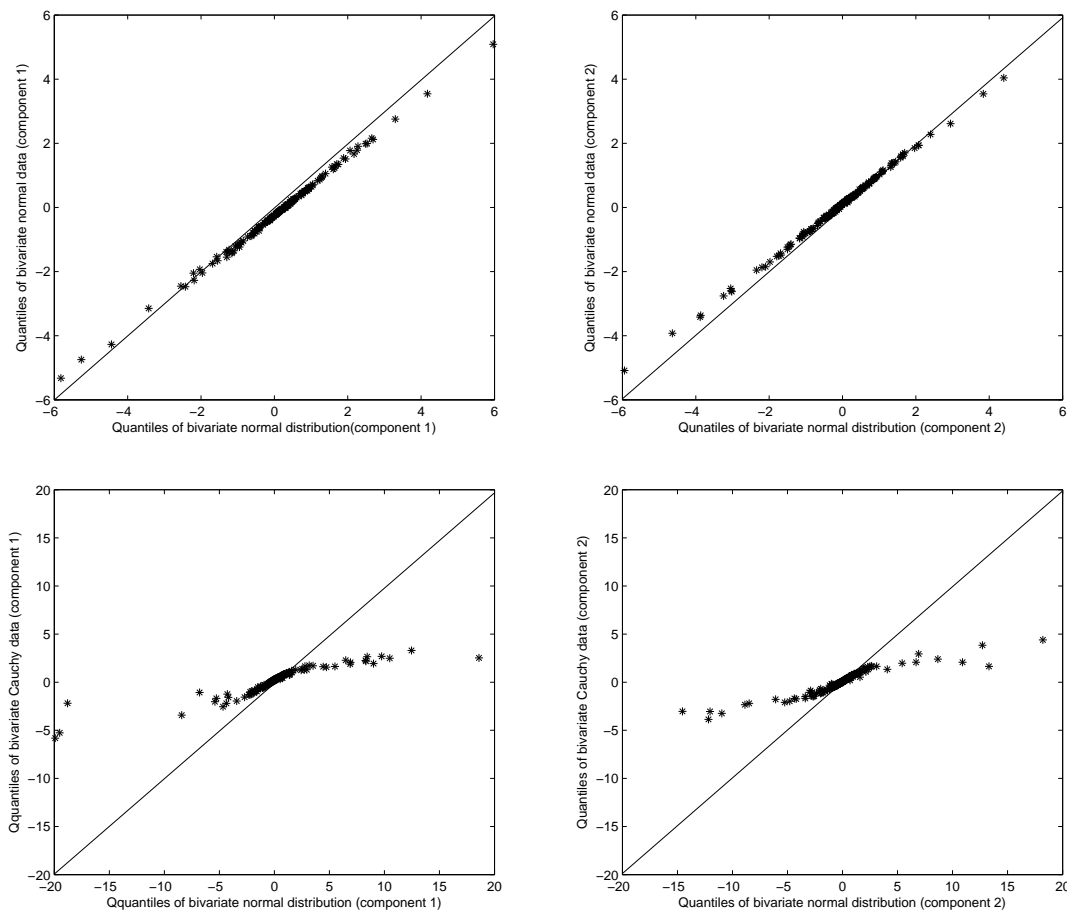


from the  $45^\circ$  line passing through the origin, and the points are actually clustered around a nonlinear curve. This indicates that the reference distribution does not fit the data well, and the distribution generating the data cannot be obtained from the reference distribution by location and scale transformations.

#### 4.1.1 Analysis of real data

Here we use three real data sets for illustrating our methodology. The hemophilia data can be obtained from the “rrcov” package in the software *R*, sunspot number data is available in <http://www.ngdc.noaa.gov/stp/SOLAR/ftpsunspotnumber.html>, and sea level pressures data can be obtained from <http://www.cpc.noaa.gov/data/indices/darwin>

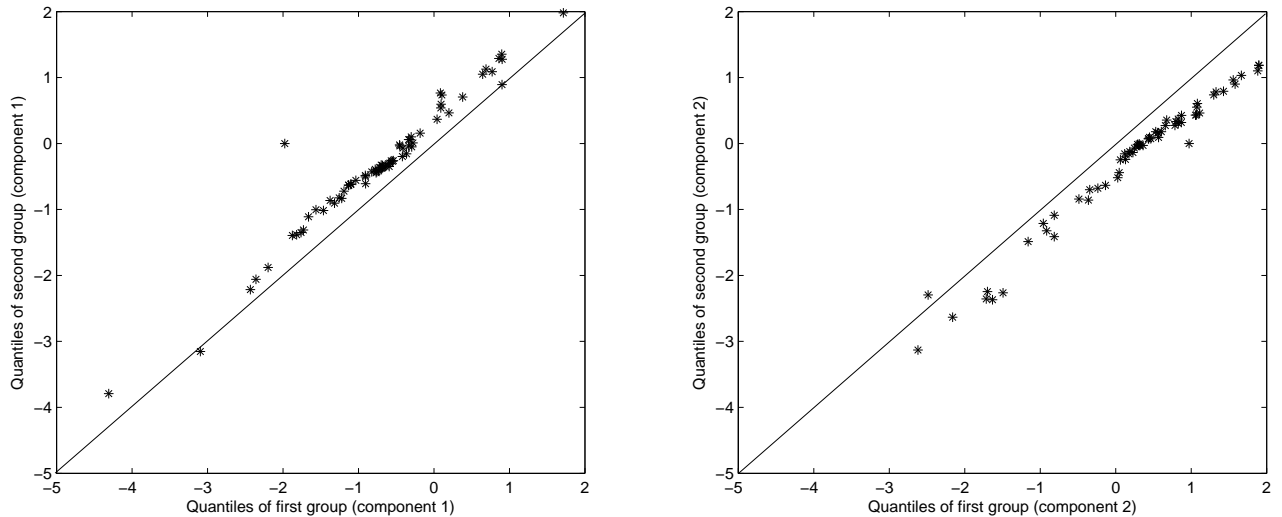
Figure 2: Q-Q plots for one-sample problems when the reference distribution is bivariate normal and the data are generated from bivariate normal (first row) and bivariate Cauchy (second row) distributions, respectively.



and <http://www.cpc.noaa.gov/data/indices/tahti>.

**Hemophilia data:** The hemophilia data set contains two groups. Among the 75 women considered, 30 are non-carriers (first group), and the remaining 45 of them are known as Hemophilia A carriers (second group). In Figure 3, we display Q-Q plots for these two groups of data along with the straight lines with slope = 1 and passing through the origin. There are visible deviations of the points in the diagrams from the straight lines drawn. In fact, the points are not clustered around any straight line, and they exhibit nonlinear patterns. This indicates that the distributions associated with two groups are not related by any location and scale transformation.

Figure 3: Q-Q plots for hemophilia data.



**Monthly sunspot number data:** The sunspot number data contains monthly average number of sunspots during the period 1749 to 2009. As data for 1749 and 2009 are incomplete, we have carried out our analysis on the observations for the remaining 259 (1750 to 2008) years. We divided the data into two samples. One sample contains six dimensional data corresponding to the six months January, February, March, October, November and December, and the other one consists of six dimensional data corresponding to the months April, May, June, July, August and September. The motivation behind splitting the data into two parts corresponding to the periods October-March and April-September comes from the fact that one equinox in a year occurs on March 20/21 and another on September 22/23. In Figure 4, we have six scatter plots. In each scatter plot, the points appear to lie on a straight line with a slope around 2 and an intercept close to zero. This is an indication that two multivariate data sets corresponding to the two periods October-March and April-September have distributions that differ only in the scales of the variables. The two distributions appear to have the same location, and one distribution can be obtained from the other by a scale transformation using the factor = 2. This fact was further confirmed when we carried

Figure 4: Q-Q plots for monthly sunspot number data.

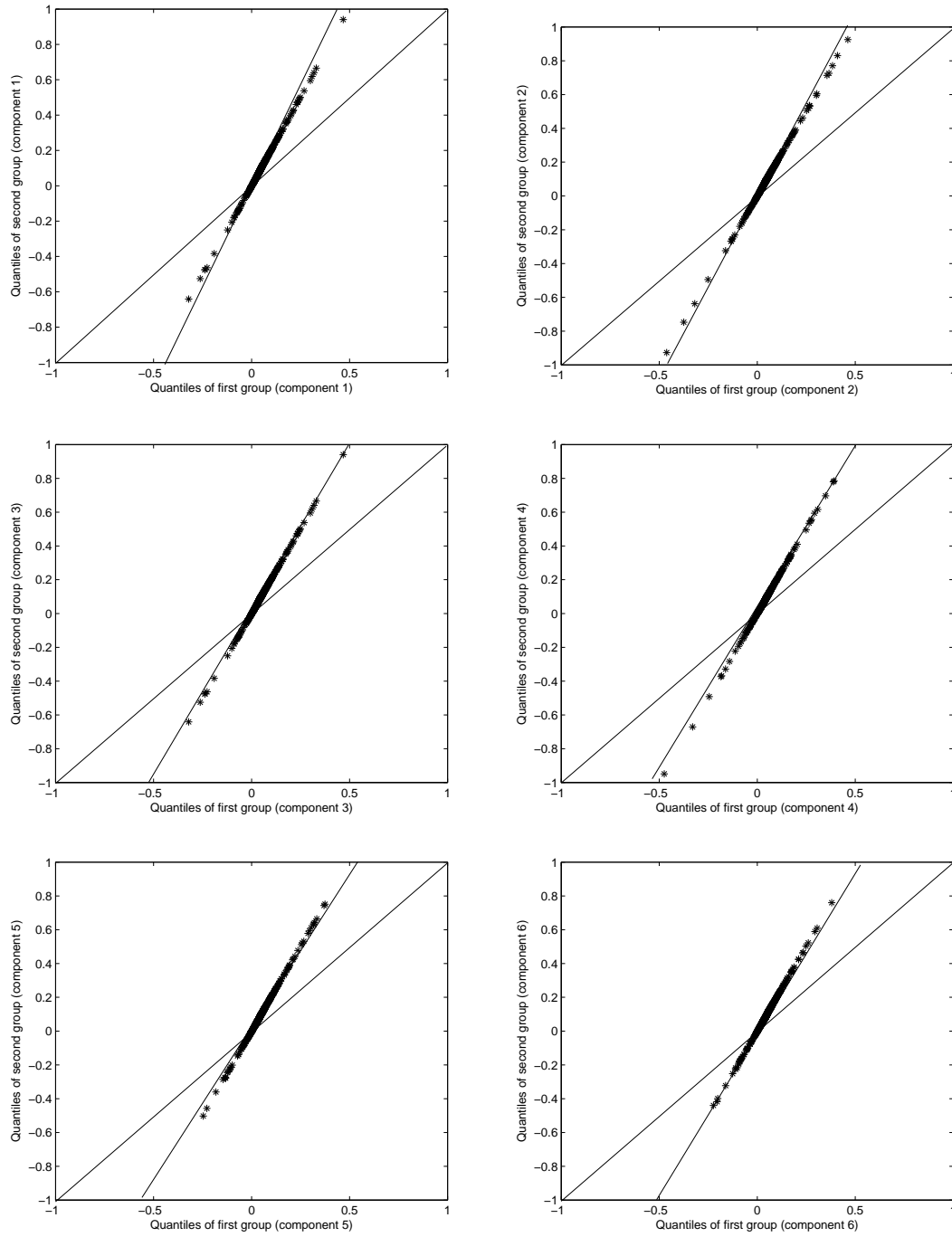
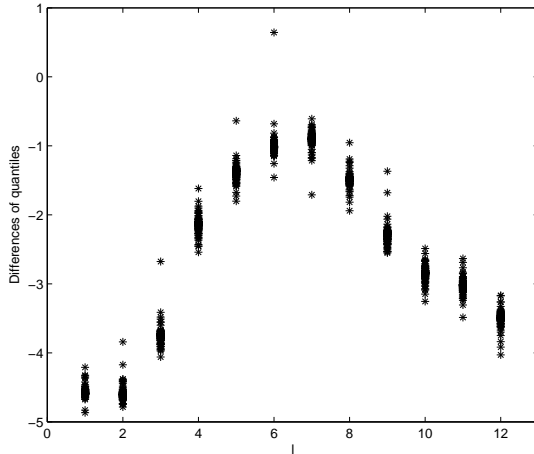


Figure 5: Quantile difference plot for data on sea level pressures.



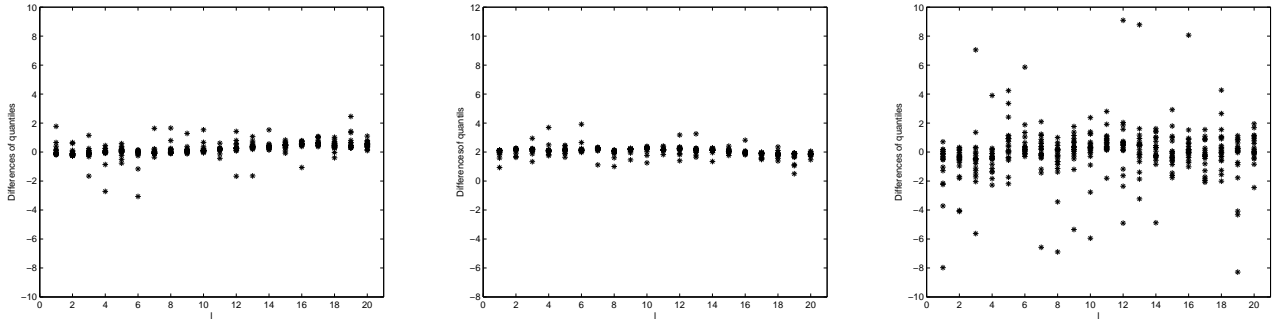
out alternative statistical analysis of the data such as comparison of marginal quantiles and direct comparison of the means and the variances of the variables.

## 4.2 Plots for high dimensional real and simulated data

When dimension of the data is large, in such cases, as we have already mentioned in Section 2, we can plot  $(l, Q_{X,l}(\mathbf{u}_k) - Q_{Y,l}(\mathbf{u}_k))$  for all  $k = 1, \dots, (n + m)$ ,  $l = 1, \dots, d$  in a single 2-dimensional plot with  $d$  vertical lines parallel to one another.

**Data on sea level pressures:** This data set consists of monthly sea level pressures from two different islands in southern Pacific ocean, namely, Darwin ( $13^\circ S$ ,  $131^\circ E$ ) and Tahiti ( $17^\circ S$ ,  $149^\circ W$ ) during the period of 1850 to 2008. Here we have a two-sample problem with each sample corresponding to an island and containing 159 twelve dimensional observations. Here we plot (see Figure 5)  $(l, Q_{X,l}(\mathbf{u}_k) - Q_{Y,l}(\mathbf{u}_k))$  for all  $k = 1, \dots, 159 + 159 = 318$ ,  $l = 1, \dots, 12$ . It is amply indicated in Figure 5 that the distributions corresponding to two samples are significantly different, and there are differences in their locations as well as scales. On each vertical line, the points are distributed in such a way that the center of each distribution is significantly different from zero. Further, the spreads of the distribution of the points on different vertical lines appear to be quite different.

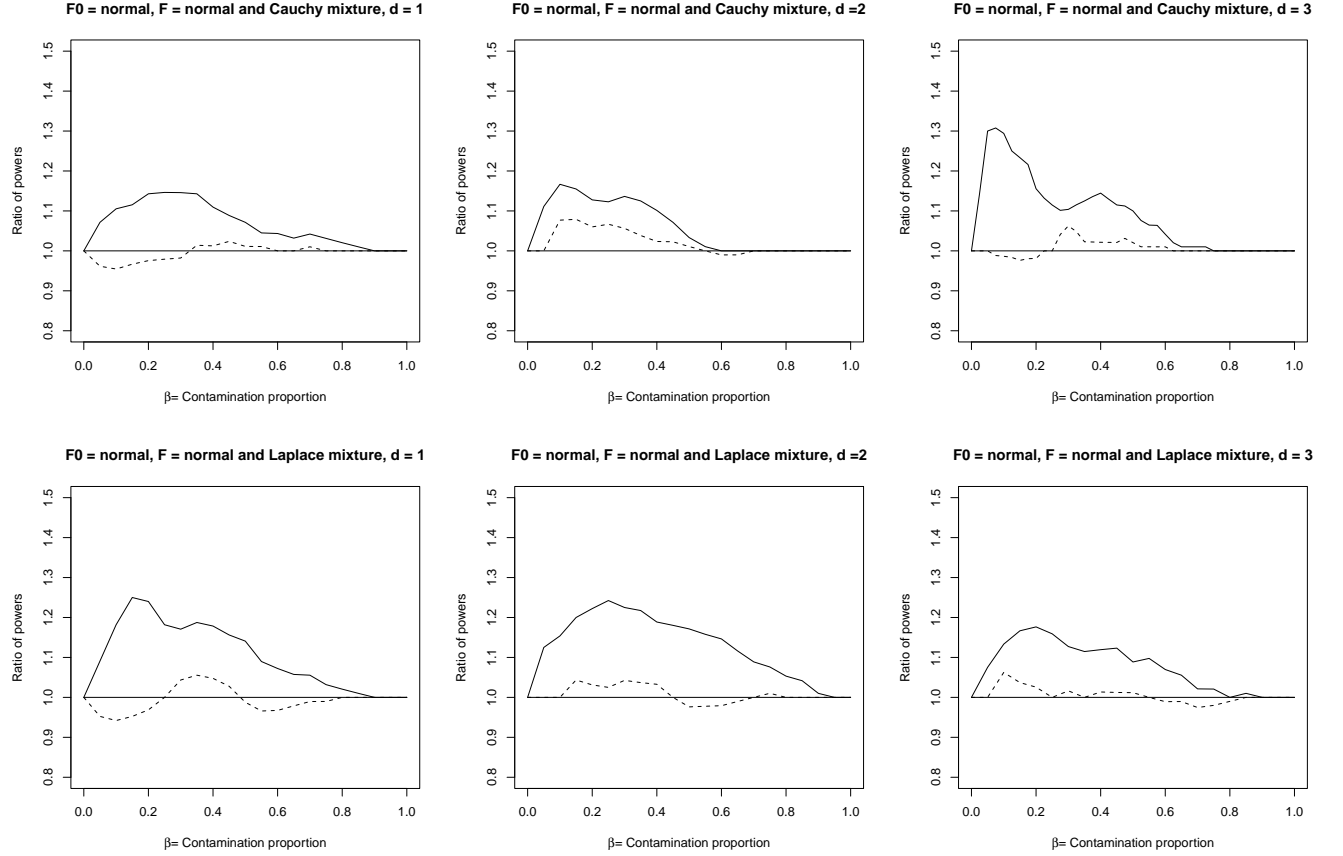
Figure 6: Quantile difference plots when two Brownian motions are identical (first diagram from the left), two processes having different mean functions but same covariance kernels (second diagram) and two processes having the same mean function but different covariance kernels (third diagram).



**Observations simulated from Gaussian process:** We generated 10 i.i.d. observations from each of three pairs of Gaussian processes with different choices of mean functions:  $m_1(t)$ ,  $m_2(t)$ , and covariance kernels:  $k_1(s, t)$ ,  $k_2(s, t)$ , where  $s, t \in [0, 1]$ . In our study, we consider an equally spaced partition  $\{t_1, \dots, t_{20}\}$  of  $[0, 1]$  and sample the observations at those time points. First, we consider  $m_1(t) = m_2(t) = 0$  and  $k_1(s, t) = k_2(s, t) = \min(s, t)$ . Next, we consider  $m_1(t) = 0$ ,  $m_2(t) = 2$ , and  $k_1(s, t) = k_2(s, t) = \min(s, t)$ , and in the third case, our choices of parameters are  $m_1(t) = m_2(t) = 0$ ,  $k_1(s, t) = \min(s, t)$ , and  $k_2(s, t) = 2 \min(s, t)$ . The plots of quantile differences for the above cases are displayed in Figure 6. In the first diagram (from the left) in Figure 6, the points in each vertical line are tightly clustered around the horizontal axis passing through the origin, which indicates that the samples are obtained from similar distributions. On the other hand, the difference between the distributions in their locations and scales, are reflected in the other two diagrams.

Note that this alternative approach based on quantile differences is related to the arrow plots considered by Marden (1998) for bivariate data. However, Marden's plots are limited to bivariate data while our quantile difference plots can be conveniently used for multivariate data with dimensions two or larger.

Figure 7: Graphs of ratios of empirical powers based on 1000 Monte-Carlo replications and  $n = 10$  for one-sample problem. *Solid* curves denote the ratios between the powers of our test (numerator) and those of KS test (denominator) while *dotted* curves denote those for our test (numerator) and CVM test (denominator).



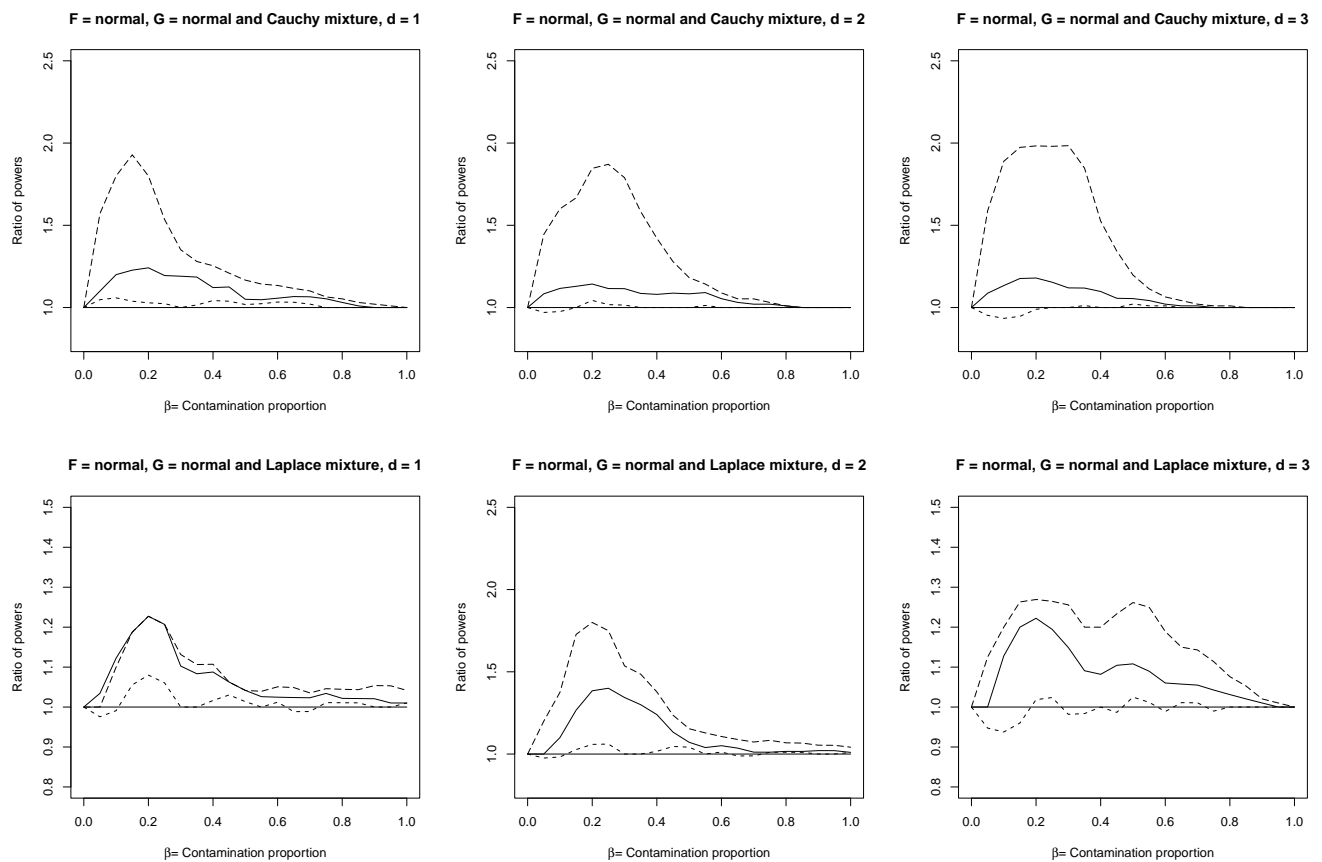
### 4.3 Level and power study for the proposed tests

A few multivariate tests for distributions have been proposed for one- and two-sample problems, and among those tests, Kolmogorov-Smirnov test (we denote it as KS test) and Cramer-von-Mises test (we denote it as CVM test) are possibly most well-known (see Bickel (1969), Justel, Pena and Zamar (1997) and Shorack and Wellner (1986)).

In order to test  $H_0$  against  $H_1$ , the KS test statistic and the CVM test statistic are  $T_n^{(1)} = \sup_{\mathbf{x}} \sqrt{n} |F_n(\mathbf{x}) - F_0(\mathbf{x})|$  and  $T_n^{(2)} = n \int_{\mathbf{x}} [F_n(\mathbf{x}) - F_0(\mathbf{x})]^2 dF_0(\mathbf{x})$ , respectively, where  $F_n(\mathbf{x})$  is the empirical distribution function of  $F(\mathbf{x})$ . For testing  $H_0^*$  against  $H_1^*$ , the KS test statistic and the CVM test statistic are  $T_{n,m}^{(1)} = \sup_{\mathbf{x}} \sqrt{n+m} |F_n(\mathbf{x}) - G_m(\mathbf{x})|$

and  $T_{n,m}^{(2)} = (n+m) \int_{\mathbf{x}} [F_n(\mathbf{x}) - G_m(\mathbf{x})]^2 dM(\mathbf{x})$ , respectively, where  $M(\mathbf{x}) = (1-\lambda)F(\mathbf{x}) + \lambda G(\mathbf{x})$ , and  $F_n$  and  $G_m$  are empirical distribution functions of  $F$  and  $G$ , respectively.

Figure 8: Graphs of ratios of empirical powers based on 1000 Monte-Carlo replications and  $n = m = 10$  for two-sample problem at 5% nominal level. The numerator in each of the ratio is power of our test while the denominators of ratios represented by *solid*, *dotted* and *dashed* curves are powers of KS, CVM and BF tests, respectively.

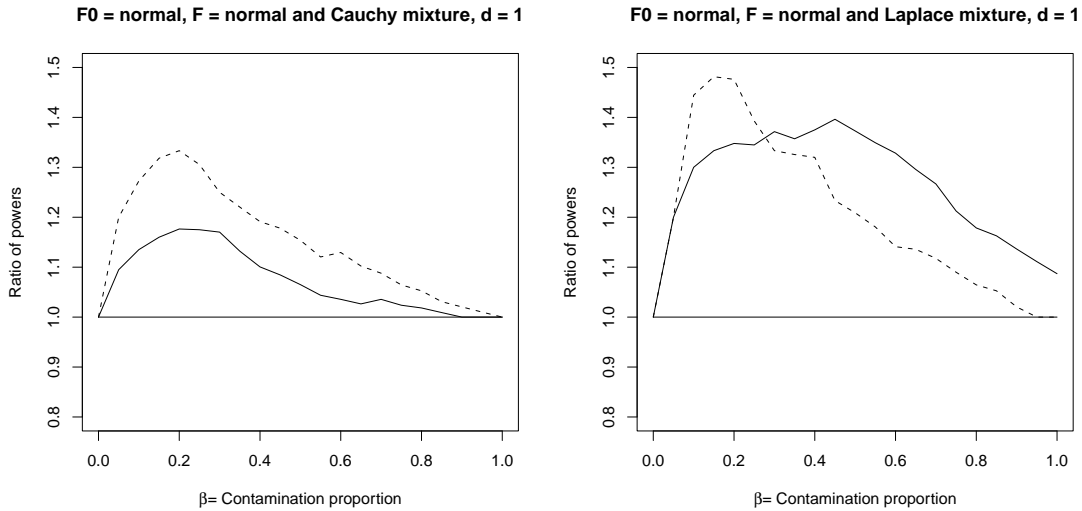


Baringhaus and Franz (2010) proposed some tests for distributions in multivariate two-sample problem based on transformation of interpoint distances. In our numerical study, following their recommendation, we have considered the test (we denote it as BF) based on the logarithmic transformation of interpoint distances. The computations for BF test has been done using the “Cramer” package in the Statistical software *R*, and the KS and the CVM tests for multivariate data have been implemented using the asymptotic distributions (see, e.g., Bickel (1969) and Shorack and Wellner (1986)) of the corresponding test statistics. For univariate data, the computations of the KS and

the CVM tests are done using the “stats” and the “CvM2SL2Test” packages in the Statistical software *R*, respectively.

Here we carry out a simulation study of the levels and the powers of our tests and some other well-known tests for multivariate data. In the case of two-sample problem, we have considered  $F = N_d(\mathbf{0}, I_d)$  and  $G = (1 - \beta)N_d(\mathbf{0}, I_d) + \beta C_d(\mathbf{0}, I_d)$  or  $(1 - \beta)N_d(\mathbf{0}, I_d) + \beta L_d(\mathbf{0}, I_d)$ , where  $\beta \in [0, 1]$ ,  $L_d(\mathbf{0}, I_d)$ ,  $N_d(\mathbf{0}, I_d)$  and  $C_d(\mathbf{0}, I_d)$  are the  $d$ -dimensional Laplace, normal and Cauchy distributions with location  $\mathbf{0}$  and scatter  $I_d$ , respectively. In the case of one-sample problem, on the other hand, we have considered  $F_0 = N_d(\mathbf{0}, I_d)$  and  $F = (1 - \beta)N_d(\mathbf{0}, I_d) + \beta C_d(\mathbf{0}, I_d)$  or  $(1 - \beta)N_d(\mathbf{0}, I_d) + \beta L_d(\mathbf{0}, I_d)$ .

Figure 9: Graphs of ratios of empirical powers based on 1000 Monte-Carlo replications and  $n = 10$  for one-sample problem at 5% nominal level. *Solid* curves denote the ratios between the powers of our test (numerator) and those of AD test (denominator) while *dotted* curves denote those for our test (numerator) and SW test (denominator).



In the diagrams in Figures 7 and 8, we have plotted the ratio between the power of a test and that of our test for different values of contamination proportion  $\beta$ . It is evident from Figures 7 and 8 that our test is comparable with CVM test in terms of their power and levels in all the cases considered in our simulation study. Our test appears to be significantly more powerful than KS and BF test when data are obtained from mixture distributions.

Table 1: Comparison of empirical powers of our two-sample test and MST-run test when  $F = N_d(\mathbf{0}, I_d)$  and  $G = N_d(\frac{\Delta}{\sqrt{d}}\mathbf{1}_d, \sigma I_d)$  at 5% nominal level based on 100 Monte-Carlo replications and  $n = m = 100$ . Here  $d$ ,  $\Delta$  and  $\sigma$  are as in Friedman and Rafsky (1981, p.706), and  $\mathbf{1}_d = (1, \dots, 1)_{1 \times d}$ .

	$d = 1$	$d = 2$	$d = 5$	$d = 10$	$d = 20$
	$\Delta = 0.3, \sigma = 1$	$\Delta = 0.5, \sigma = 1$	$\Delta = 0.75, \sigma = 1$	$\Delta = 1.0, \sigma = 1$	$\Delta = 1.2, \sigma = 1$
Our test	0.332	0.554	0.701	0.832	0.993
MST-run test	0.180	0.350	0.640	0.780	0.860
	$\Delta = 0, \sigma = 1.3$	$\Delta = 0, \sigma = 1.2$	$\Delta = 0, \sigma = 1.2$	$\Delta = 0, \sigma = 1.1$	$\Delta = 0, \sigma = 1.075$
Our test	0.251	0.171	0.263	0.075	0.144
MST-run test	0.240	0.140	0.210	0.090	0.130

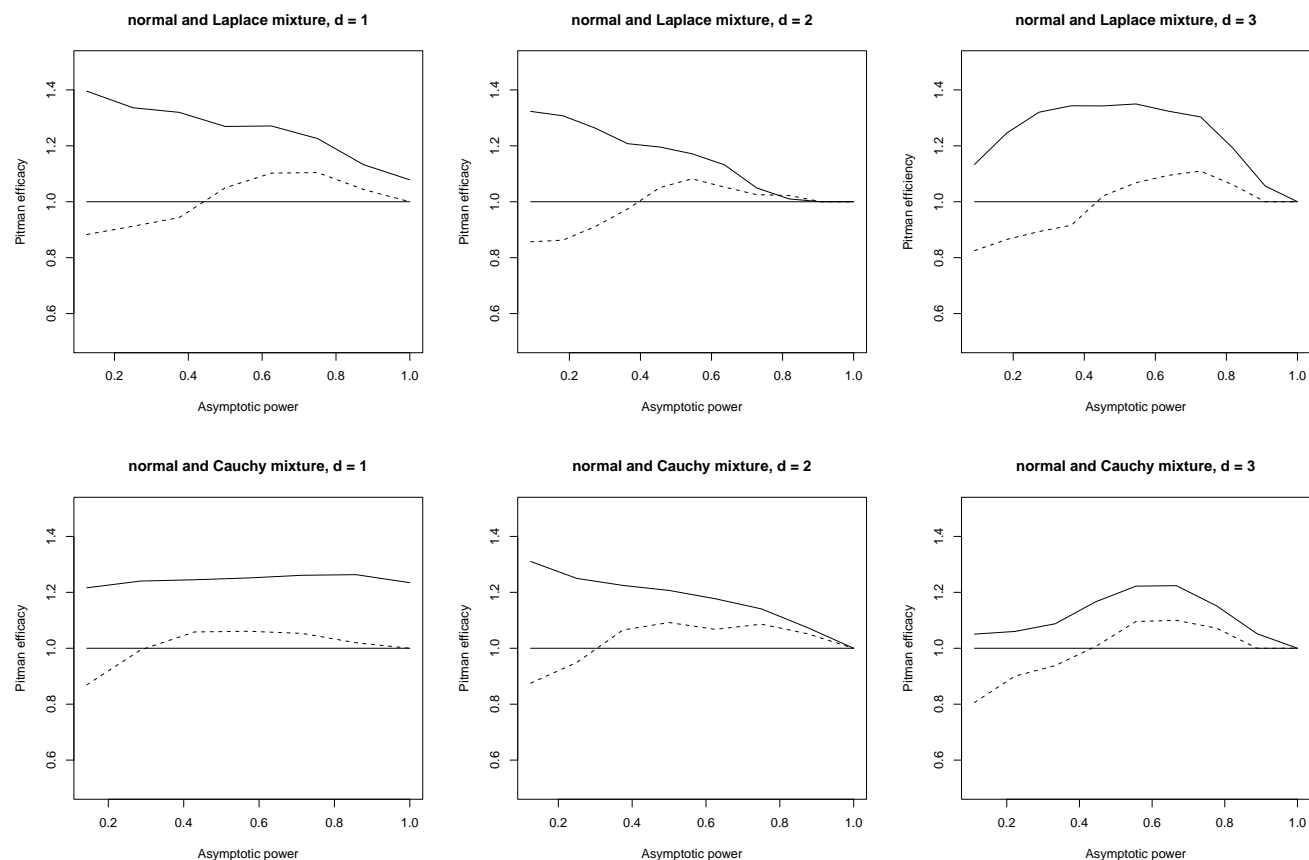
Shapiro-Wilks test (we denote it as SW, see Shapiro and Wilk (1965)) and Anderson-Darling test (we denote it as AD test, see Anderson and Darling (1954)) are two other well-known tests for checking the Gaussianity of univariate data, and it will be appropriate to compare the performance of our test with that of the SW and the AD tests in the univariate case. The SW and the AD tests for have been implemented using the asymptotic distributions (see, e.g., Leslie, Stephens and Fotopoulos (1986) and Shorack and Wellner (1986)) of the corresponding test statistics. The computations of the SW and the AD tests are done using the “stats” and “nortest” packages in the Statistical software *R*. It is evident from the Figure 9 that both of the SW and the AD tests are outperformed by our test.

Friedman and Rafsky (1979) proposed a multivariate generalizations of the Wald-Wolfowitz run test using the idea of minimum spanning tree (we denote it as MST-run). We have compared the powers of our two-sample test with those of the MST-run test in different distributions (namely, multivariate normal with different locations and scales) considered by Friedman and Rafsky (1981, p.706). The results are reported in Table 1, and it is clear that the MST-run test does not perform well compared to our test (see figures in Table 1).

## 5 Asymptotic power study under contiguous alternatives

The performance of our tests is comparable with that of CVM test, and it is slightly better than that of KS test and substantially better than that of BF, SW and AD tests in finite sample situations, as we have seen in Section 4.3. Since our tests, the KS and the CVM tests are all consistent, a natural question is how the asymptotic powers of our tests and the KS and the CVM tests compare with one another under contiguous alternatives.

Figure 10: Pitman efficacy of our test relative to KS test (*solid curve*) and CVM test (*dotted curve*) for one-sample problem at 5% nominal level.



In the case of one-sample problem, the null hypothesis is given by  $H_0 : F(\mathbf{x}) = F_0(\mathbf{x})$ , and we consider a sequence of contiguous alternatives  $H_n : F_n(\mathbf{x}) = (1 - \delta/\sqrt{n})F_0(\mathbf{x}) +$

$(\delta/\sqrt{n})H(\mathbf{x})$  for a fixed  $\delta > 0$ .

**Theorem 5.1:** *Assume that  $F_0$  and  $H$  have continuous densities  $f_0$  and  $h$ , respectively. Also, if  $f_0$  is positive and  $E_{f_0} \left\{ \frac{h(\mathbf{x})}{f_0(\mathbf{x})} - 1 \right\}^2 < \infty$ , then the sequence of alternatives  $H_n$  form a contiguous sequence. Under the sequence of such alternatives, the asymptotic power of the test based on  $V_n$  is given by  $P_\delta[\int_{\mathbf{u} \in S(\delta)} \|Z'_1(\mathbf{u})\|^2 d\mathbf{u} > c_1(\alpha)]$ , where  $Z'_1(\mathbf{u})$  is a Gaussian process with mean function*

$$\begin{aligned} m_1(\mathbf{u}, \delta) &= \delta[D_1\{Q_{F_0}(\mathbf{u})\}]^{-1} E_H \left\{ \frac{\mathbf{x} - Q_{F_0}(\mathbf{u})}{\|\mathbf{x} - Q_{F_0}(\mathbf{u})\|} + \mathbf{u} \right\} \quad \text{if } d \geq 2 \\ &= \delta E_H \left\{ \frac{\left(\frac{u+1}{2}\right) - 1_{\{x \leq F^{-1}(\frac{u+1}{2})\}}}{f(F^{-1}(\frac{u+1}{2}))} \right\} \quad \text{if } d = 1 \end{aligned}$$

and covariance kernel  $k_1(\mathbf{u}_1, \mathbf{u}_2)$ . Here  $k_1(\mathbf{u}_1, \mathbf{u}_2)$  and  $c_1(\alpha)$  are as defined in Theorem 3.1 such that  $P_{\delta=0}[\int_{\mathbf{u} \in S(\delta)} \|Z'_1(\mathbf{u})\|^2 d\mathbf{u} > c_1(\alpha)] = \alpha$ . Further, under the sequence of those alternatives, the asymptotic powers of the tests based on  $T_n^{(1)}$  and  $T_n^{(2)}$  are given by  $P_\delta[\sup_{\mathbf{x}} |Z''_1(\mathbf{x})| > c_1^*(\alpha)]$  and  $P_\delta[\int_{\mathbf{x}} \{Z''_1(\mathbf{x})\}^2 dF_0(\mathbf{x}) > c_1^{**}(\alpha)]$ , respectively, where  $Z''_1(\mathbf{x})$  is a Gaussian process with mean function  $m'_1(\mathbf{x}, \delta) = \delta\{H(\mathbf{x}) - F_0(\mathbf{x})\}$  and covariance kernel  $k_3(\mathbf{x}_1, \mathbf{x}_2) = F_0(\mathbf{x}_1)[1 - F_0(\mathbf{x}_2)]$ . Here  $c_1^*(\alpha)$  and  $c_1^{**}(\alpha)$  satisfy  $P_{\delta=0}[\sup_{\mathbf{x}} |Z''_1(\mathbf{x})| > c_1^*(\alpha)] = \alpha$  and  $P_{\delta=0}[\int_{\mathbf{x}} \{Z''_1(\mathbf{x})\}^2 dF_0(\mathbf{x}) > c_1^{**}(\alpha)] = \alpha$ .

Next, in the case of two-sample problem, the null hypothesis is given by  $H_0 : F(\mathbf{x}) = G(\mathbf{x})$ , and we consider a sequence of alternatives  $H_{n,m} : G_{n,m}(\mathbf{x}) = (1 - \delta/\sqrt{n+m})F(\mathbf{x}) + (\delta/\sqrt{n+m})H(\mathbf{x})$  for a fixed  $\delta > 0$ .

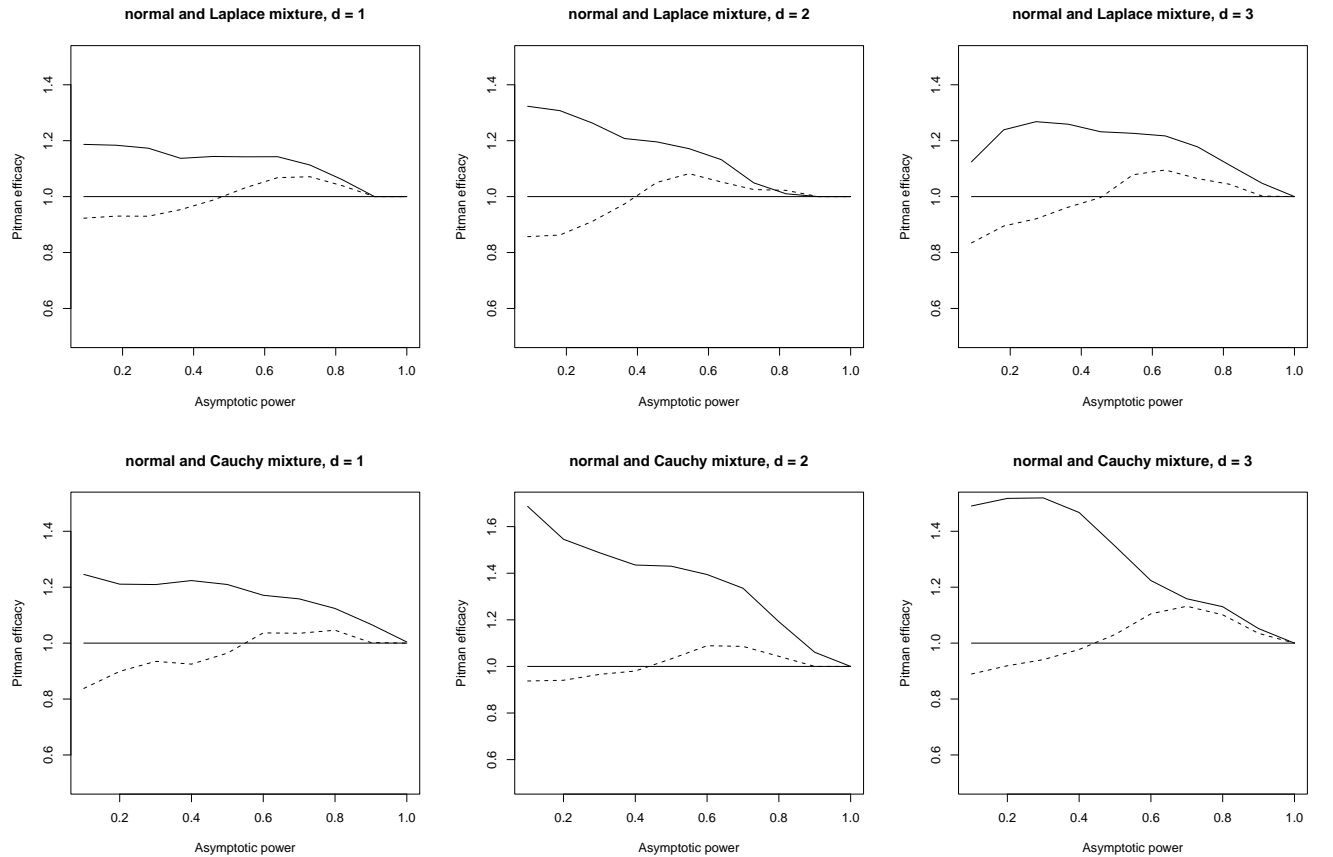
**Theorem 5.2:** *Assume that  $F$  and  $H$  have continuous densities  $f$  and  $h$ , respectively. Also, let  $f$  is positive,  $E_f \left\{ \frac{h(\mathbf{y})}{f(\mathbf{y})} - 1 \right\}^2 < \infty$  and  $n, m \rightarrow \infty$  in such a way that  $\lim_{n,m \rightarrow \infty} \frac{n}{n+m} = \lambda \in (0, 1)$ . Then, the sequence of densities associated with alternatives  $H_{n,m}$  form a contiguous sequence. Under the sequence of such alternatives, the asymptotic power of the test based on  $T_{n,m}$  is given by  $P_\delta[\int_{\mathbf{u} \in S(\delta)} \|Z'_2(\mathbf{u})\|^2 d\mathbf{u} > c_2(\alpha)]$ , where  $Z'_2(\mathbf{u})$  is a Gaussian process with mean function*

$$m_2(\mathbf{u}, \delta) = -\delta[D_1^F(Q(\mathbf{u}))]^{-1} E_h \left\{ \frac{\mathbf{y} - Q_F(\mathbf{u})}{\|\mathbf{y} - Q_F(\mathbf{u})\|} + \mathbf{u} \right\} \quad \text{if } d \geq 2$$

$$= -\delta E_h \left\{ \frac{\left(\frac{u+1}{2}\right) - 1_{\{y \leq F^{-1}(\frac{u+1}{2})\}}}{f(F^{-1}(\frac{u+1}{2}))} \right\} \quad \text{if } d = 1$$

and covariance kernel  $k_2(\mathbf{u}_1, \mathbf{u}_2)$ . Here  $k_2(\mathbf{u}_1, \mathbf{u}_2)$  and  $c_2(\alpha)$  are as defined in Theorem 3.2 such that  $P_{\delta=0}[\int_{\mathbf{u} \in S(\delta)} \|Z'_2(\mathbf{u})\|^2 d\mathbf{u} > c_2(\alpha)] = \alpha$ . Further, under the sequence of those alternatives, the asymptotic powers of the tests based on  $T_{n,m}^{(1)}$  and  $T_{n,m}^{(2)}$  are given by  $P_\delta[\sup_{\mathbf{x}} |Z''_2(\mathbf{x})| > c_2^*(\alpha)]$  and  $P_\delta[\int_{\mathbf{x}} \{Z''_2(\mathbf{x})\}^2 dM(\mathbf{x}) > c_2^{**}(\alpha)]$ , respectively, where  $Z''_2(\mathbf{x})$  is a Gaussian process with mean function  $m'_2(\mathbf{x}, \delta) = -\delta\{H(\mathbf{x}) - F(\mathbf{x})\}$  and covariance kernel  $k_4(\mathbf{x}_1, \mathbf{x}_2) = \frac{F(\mathbf{x}_1)(1-F(\mathbf{x}_2))}{\lambda(1-\lambda)}$ . Here  $c_2^*(\alpha)$  and  $c_2^{**}(\alpha)$  satisfy  $P_{\delta=0}[\sup_{\mathbf{x}} |Z''_2(\mathbf{x})| > c_2^*(\alpha)] = \alpha$  and  $P_{\delta=0}[\int_{\mathbf{x}} \{Z''_2(\mathbf{x})\}^2 dM(\mathbf{x}) > c_2^{**}(\alpha)] = \alpha$ .

Figure 11: Pitman efficacy of our test relative to KS test (*solid curve*) and CVM test (*dotted curve*) for two-sample problem at 5% nominal level.



Theorems 5.1 and 5.2 are basic results for computing the asymptotic powers of the

tests under contiguous alternatives, and they enable us to derive Pitman efficacies of our tests relative to KS and CVM tests. The Pitman efficacy (see, e.g., Serfling (1980) and Lehmann and Romano (2005)) of our test relative to another test for varying choices of asymptotic power determined by  $\delta$  is given by  $(\delta'/\delta)^2$ , where  $\delta$  and  $\delta'$  are such that the asymptotic power of our test under contiguous alternatives like  $(1 - \delta/\sqrt{n})F_0(\mathbf{x}) + (\delta/\sqrt{n})H(\mathbf{x})$  or  $(1 - \delta/\sqrt{n+m})F(\mathbf{x}) + (\delta/\sqrt{n+m})H(\mathbf{x})$  is same as the asymptotic power of the other test under contiguous alternatives like  $(1 - \delta'/\sqrt{n})F_0(\mathbf{x}) + (\delta'/\sqrt{n})H(\mathbf{x})$  or  $(1 - \delta'/\sqrt{n+m})F(\mathbf{x}) + (\delta'/\sqrt{n+m})H(\mathbf{x})$ , respectively.

In the diagrams in Figures 10 and 11, we have plotted the Pitman efficacy of our test for different values of asymptotic power. It is clearly indicated by the diagrams in Figures 10 and 11 that our test and the CVM test asymptotically outperform KS test in terms of Pitman efficacy in all the cases considered here. However, between our test and the CVM test, one has superior performances in some cases while the other has superior performance in other cases, and there are only small differences in their performance.

## 6 Appendix: Proofs

*Proof of Theorem 2.1:* Note that for  $d = 1$ , the centered and normalized two-sample quantile process  $\sqrt{n+m}\{(\hat{Q}_X(u), \hat{Q}_Y(u)) - (Q_F(u), Q_G(u))\}$ , has its sample paths lying in the space of real valued right continuous functions defined on the interval  $S(\delta) = \{u : |u| \leq 1 - \delta\}$ . However, for  $d \geq 2$ , in view of the results in Chaudhuri (1996) and Koltchinskii (1997), the sample paths of the centered and normalized stochastic process  $\sqrt{n+m}\{(\hat{Q}_X(\mathbf{u}), \hat{Q}_Y(\mathbf{u})) - (Q_F(\mathbf{u}), Q_G(\mathbf{u}))\}$  lie in the space of  $\mathbb{R}^{2d}$  valued continuous functions defined on the  $S(\delta) = \{\mathbf{u} : \|\mathbf{u}\| \leq 1 - \delta\}$ . Suppose that we have the usual supremum norm on the spaces of continuous and right continuous functions, which makes those spaces appropriate metric spaces for deriving weak convergence results for those quantile processes. Then, for  $d \geq 2$  and under the conditions stated in the theorem, it follows from Chaudhuri (1996) and Koltchinskii (1997) that for  $\mathbf{u} \in$

$S(\delta)$ , the vector valued stochastic process  $\sqrt{n+m}\{(\hat{Q}_{\mathcal{X}}(\mathbf{u}), \hat{Q}_{\mathcal{Y}}(\mathbf{u})) - (Q_F(\mathbf{u}), Q_G(\mathbf{u}))\}$  converges weakly to a Gaussian process with zero mean. Further, for  $d = 1$ , a similar weak convergence result for the centered and normalized two-sample univariate quantile process follows from Bickel (1967). These weak convergence results imply that for all  $d \geq 1$  and any  $0 < \delta < 1$ , we have

$$\sup_{\mathbf{u} \in S(\delta)} \|\{(\hat{Q}_{\mathcal{X}}(\mathbf{u}), \hat{Q}_{\mathcal{Y}}(\mathbf{u})) - (Q_F(\mathbf{u}), Q_G(\mathbf{u}))\}\| = o_p(1). \quad (1)$$

It follows from the location and the scale equivariance of spatial quantiles (see Chaudhuri (1996, Section 2.2)) that  $Q_G(\mathbf{u}) = \sigma Q_F(\mathbf{u}) + \boldsymbol{\mu}$ , which implies that the set of points  $\{(Q_F(\mathbf{u}), Q_G(\mathbf{u})) : \|\mathbf{u}\| < 1\}$  will be a hyperplane in  $\mathbb{R}^d$  defined by  $d$  linear equations of the form  $v_i = \sigma w_i + \mu_i$  for  $i = 1, \dots, d$ , where  $v = (v_1, \dots, v_d)$ ,  $w = (w_1, \dots, w_d) \in \mathbb{R}^d$ . This fact along with the uniform convergence result in (1) leads to the proof of the “if part” of the theorem in the two-sample case.

Next, since  $\lim_{n,m \rightarrow \infty} P[S_{n,m}(\mathcal{X}, \mathcal{Y}, \delta) \in \{(\mathbf{v} \in \mathbb{R}^d, \mathbf{w} \in \mathbb{R}^d) : \|\mathbf{v} - \sigma \mathbf{w} - \boldsymbol{\mu}\| < \epsilon\}] = 1$  for every  $\epsilon > 0$ , we have  $Q_G(\mathbf{u}) = \sigma Q_F(\mathbf{u}) + \boldsymbol{\mu}$  in view of (1). It now follows from the characterization of distributions based on spatial quantiles (see Koltchinskii (1997), Corollary 2.9) and the location and scale equivariance of spatial quantiles (see Chaudhuri (1996, Section 2.2)) that  $F(\mathbf{x}) = G((\mathbf{x} - \boldsymbol{\mu})/\sigma)$ . This completes the proof of “only if part” of the theorem in the two-sample case.

Arguing in a very similar way as in the proof of two-sample problem, one can establish the result for the one-sample problem.  $\square$

*Proof of Theorem 3.1:* The asymptotic power of the test is given by

$$\lim_{n \rightarrow \infty} P_{H_1}[V_n > c_1(\alpha)] = \lim_{n \rightarrow \infty} P_{H_1} \left[ n \int_{\mathbf{u} \in S(\delta)} \|\hat{Q}_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\|^2 d\mathbf{u} > c_1(\alpha) \right],$$

where  $c_1(\alpha)$  is the critical value determined from the level  $(0 < \alpha < 1)$  of the test based on the distribution of  $\int_{\mathbf{u} \in S(\delta)} \|Z_1(\mathbf{u})\|^2 d\mathbf{u}$ , and  $Z_1(\mathbf{u})$  is as defined in the statement of the theorem.

Now, we have

$$\begin{aligned}
& \lim_{n \rightarrow \infty} P_{H_1} \left[ n \int_{\mathbf{u} \in S(\delta)} \|\hat{Q}_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\|^2 d\mathbf{u} > c_1(\alpha) \right] \\
&= \lim_{n \rightarrow \infty} P_{H_1} \left[ n \int_{\mathbf{u} \in S(\delta)} \|\{\hat{Q}_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\} - \{Q_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\}\|^2 d\mathbf{u} \right. \\
&> c_1(\alpha) + n \left[ \int_{\mathbf{u} \in S(\delta)} \{Q_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\}' \{Q_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\} d\mathbf{u} \right. \\
&\left. \left. - 2 \int_{\mathbf{u} \in S(\delta)} \{\hat{Q}_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\}' \{Q_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\} d\mathbf{u} \right] \right].
\end{aligned}$$

It follows from Chaudhuri (1996) and Koltchinskii (1997) in the case  $d \geq 2$  that  $\sqrt{n}\{\hat{Q}_F(\mathbf{u}) - Q_F(\mathbf{u})\}$  converges weakly to the Gaussian process  $Z_1(\mathbf{u})$  for  $\mathbf{u} \in S(\delta)$ . For  $d = 1$ , a similar weak convergence result follows from Bickel (1967). So, it now follows using the continuity of the integral functional that  $V_n$  converges in distribution to  $\int_{\mathbf{u} \in S(\delta)} \|Z_1(\mathbf{u})\|^2 d\mathbf{u}$ . Also, note that

$$\begin{aligned}
& c_1(\alpha) + n \left[ \int_{\mathbf{u} \in S(\delta)} \{Q_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\}' \{Q_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\} d\mathbf{u} \right. \\
&\left. - 2 \int_{\mathbf{u} \in S(\delta)} \{\hat{Q}_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\}' \{Q_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\} d\mathbf{u} \right] \rightarrow -\infty
\end{aligned}$$

as  $n \rightarrow \infty$  since  $\hat{Q}_F(\mathbf{u}) \xrightarrow{p} Q_F(\mathbf{u})$  uniformly in  $\mathbf{u} \in S(\delta)$  and the continuous functions  $Q_F(\mathbf{u})$  and  $Q_{F_0}(\mathbf{u})$  satisfy  $Q_F(\mathbf{u}) \neq Q_{F_0}(\mathbf{u})$  for some  $\mathbf{u} \in S(\delta)$  if  $H_1$  is true, and  $0 < \delta < 1$  is appropriately large. Hence,  $P_{H_1}[V_n > c_1(\alpha)] \rightarrow 1$  as  $n \rightarrow \infty$ . This completes the proof.  $\square$

*Proof of Theorem 3.2:* Arguing in a similar way as in the proof of Theorem 3.1 and using the independence of the two samples, if  $n, m \rightarrow \infty$  in such a way that  $\lambda = \lim_{n, m \rightarrow \infty} \frac{n}{n+m} \in (0, 1)$ , under  $H_0^*$ , one can show that  $T_{n,m}$  converges in distribution to  $\int_{\mathbf{u} \in S(\delta)} \|Z_2(\mathbf{u})\|^2 d\mathbf{u}$ . Next, the asymptotic power of the test is given by  $P_{H_1^*}[T_{n,m} > c_2(\alpha)]$ , where  $c_2(\alpha)$  is the critical value determined from the level ( $0 < \alpha < 1$ ) of the test based on the distribution of  $\int_{\mathbf{u} \in S(\delta)} \|Z_2(\mathbf{u})\|^2 d\mathbf{u}$ . Using a very similar argument as in the proof of Theorem 3.1, one can establish that  $P_{H_1^*}[T_{n,m} > c_2(\alpha)] \rightarrow 1$  as  $n, m \rightarrow \infty$ .  $\square$

*Proof of Theorem 5.1:* The likelihood ratio for testing  $H_0$  against  $H_n$  is

$$\begin{aligned}
L_n &= \sum_{i=1}^n \log \frac{f_n(\mathbf{x}_i)}{f_0(\mathbf{x}_i)} = \sum_{i=1}^n \log \frac{(1-\beta)f_0(\mathbf{x}_i) + \beta h(\mathbf{x}_i)}{f_0(\mathbf{x}_i)} = \sum_{i=1}^n \log \left[ 1 + \beta \left\{ \frac{h(\mathbf{x}_i)}{f_0(\mathbf{x}_i)} - 1 \right\} \right] \\
&= \frac{\delta}{\sqrt{n}} \sum_{i=1}^n \left\{ \frac{h(\mathbf{x}_i)}{f_0(\mathbf{x}_i)} - 1 \right\} - \frac{\delta^2}{2n} \sum_{i=1}^n \left\{ \frac{h(\mathbf{x}_i)}{f_0(\mathbf{x}_i)} - 1 \right\}^2 + R_n \quad (\text{using } \beta = \frac{\delta}{\sqrt{n}}) \\
&= \frac{\delta}{\sqrt{n}} \sum_{i=1}^n k_i - \frac{\delta^2}{2} \times \frac{1}{n} \sum_{i=1}^n k_i^2 + R_n \quad (\text{here } k_i = \frac{h(\mathbf{x}_i)}{f_0(\mathbf{x}_i)} - 1). \tag{2}
\end{aligned}$$

Note that  $R_n \xrightarrow{p} 0$  as  $n \rightarrow \infty$ . Further, using straightforward applications of C.L.T. and W.L.L.N., we have that the first term of (2) is asymptotically normal with mean zero and variance  $\delta^2 \text{var}(k_i) = \delta^2 \sigma^2$ , and the second term of (2) will converge in probability to  $\frac{\delta^2}{2} \sigma^2$ . So, using Slutsky's theorem, we have  $L_n$  asymptotically normally distributed with mean  $-\frac{\delta^2}{2} \sigma^2$  and variance  $\delta^2 \sigma^2$ . This ensures the contiguity of the sequence  $H_n$  using the corollary to Lecam's First Lemma in Hajek and Sidak (1967, p. 204).

Now, we consider  $\mathbf{u}_1, \dots, \mathbf{u}_k \in S(\delta)$  and  $\mathbf{x}_1, \dots, \mathbf{x}_l \in \mathbb{R}^d$ . Then, under  $H_0$ , one can establish that the joint distribution of  $\sqrt{n}\{\hat{Q}_F(\mathbf{u}_1) - Q_{F_0}(\mathbf{u}_1), \dots, \hat{Q}_F(\mathbf{u}_k) - Q_{F_0}(\mathbf{u}_k), F_n(\mathbf{x}_1) - F_0(\mathbf{x}_1), \dots, F_n(\mathbf{x}_l) - F_0(\mathbf{x}_l), L_n\}$  is asymptotically multivariate Gaussian using the Bahadur type linear expansion of spatial quantile  $\hat{Q}_F(\mathbf{u})$  (see Chaudhuri (1996) for  $d \geq 2$  and Serfling (1980) for  $d = 1$ ), the expansion of log likelihood ratio  $L_n$  (see (2) above) and the fact that  $\{F_n(\mathbf{x}) - F_0(\mathbf{x})\}$  is a simple average of i.i.d. random variables. Note that for any  $p = 1, \dots, k$ , when  $d \geq 2$ , the asymptotic covariance between  $\sqrt{n}\{\hat{Q}_F(\mathbf{u}_p) - Q_{F_0}(\mathbf{u}_p)\}$  and  $\sqrt{n}L_n$  is

$$\begin{aligned}
m_1(\mathbf{u}_p, \delta) &= \frac{\delta}{n} E_{f_0} \left[ \sum_{i=1}^n \left\{ D_1[Q_{F_0}(\mathbf{u}_p)]^{-1} \left\{ \frac{\mathbf{x}_i - Q_{F_0}(\mathbf{u}_p)}{\|\mathbf{x}_i - Q_{F_0}(\mathbf{u}_p)\|} + \mathbf{u}_p \right\} \right\} \times \left\{ \frac{h(\mathbf{x}_i)}{f_0(\mathbf{x}_i)} - 1 \right\} \right] \\
&= \delta [D_1\{Q_{F_0}(\mathbf{u}_p)\}]^{-1} E_h \left\{ \frac{\mathbf{x} - Q_{F_0}(\mathbf{u}_p)}{\|\mathbf{x} - Q_{F_0}(\mathbf{u}_p)\|} + \mathbf{u}_p \right\} \quad (\text{as } E_{f_0} \left\{ \frac{\mathbf{x} - Q_{F_0}(\mathbf{u}_p)}{\|\mathbf{x} - Q_{F_0}(\mathbf{u}_p)\|} + \mathbf{u}_p \right\} = \mathbf{0}).
\end{aligned}$$

For  $d = 1$ , the asymptotic covariance between  $\sqrt{n}\{\hat{Q}_F(u_p) - Q_{F_0}(u_p)\}$  and  $\sqrt{n}L_n$  is

$$m_1(u_p, \delta) = \frac{\delta}{n} E_{f_0} \left[ \sum_{i=1}^n \left\{ \frac{\binom{u_p+1}{2} - 1_{\{x_i \leq F^{-1}(\frac{u_p+1}{2})\}}}{f(F^{-1}(\frac{u_p+1}{2}))} \right\} \times \left\{ \frac{h(x_p)}{f_0(x_p)} - 1 \right\} \right]$$

$$= \delta E_h \left\{ \frac{\left(\frac{u_{p+1}}{2}\right) - 1_{\{x \leq F^{-1}(\frac{u_{p+1}}{2})\}}}{f(F^{-1}(\frac{u_{p+1}}{2}))} \right\} \quad (\text{as } E_{f_0} \left\{ \frac{\left(\frac{u_{p+1}}{2}\right) - 1_{\{x \leq F^{-1}(\frac{u_{p+1}}{2})\}}}{f(F^{-1}(\frac{u_{p+1}}{2}))} \right\} = 0).$$

Also, one can show that for any  $j = 1, \dots, l$ , the asymptotic covariance between  $\sqrt{n}\{F_n(\mathbf{x}_l) - F_0(\mathbf{x}_l)\}$  and  $\sqrt{n}L_n$  is  $m'_1(\mathbf{x}_j, \delta) = \delta\{H(\mathbf{x}_j) - F_0(\mathbf{x}_j)\}$ .

Now, by a straightforward application of Lecam's third lemma (see Hajek and Sidak (1967), p.208), one can establish that under contiguous alternatives,  $\sqrt{n}\{\hat{Q}_F(\mathbf{u}_1) - Q_{F_0}(\mathbf{u}_1), \dots, \hat{Q}_F(\mathbf{u}_k) - Q_{F_0}(\mathbf{u}_k)\}$  is asymptotically  $kd$ -dimensional multivariate normal with mean having  $d$ -dimensional  $p$ -th block  $m_1(\mathbf{u}_p, \delta)$  and  $kd \times kd$ -dimensional covariance matrix obtained from the covariance kernel  $k_1$  as given in the statement of Theorem 3.1. Further, the spatial quantile process satisfies tightness condition under contiguous alternatives in view of the fact that it is tight under  $H_0$ . The tightness under  $H_0$  follows from the weak convergence of the spatial quantile process (see Chaudhuri (1996) and Koltchinskii (1997) for  $d \geq 2$  and Bickel (1967) for  $d = 1$ ). So, the spatial quantile process  $\sqrt{n}\{\hat{Q}_F(\mathbf{u}) - Q_{F_0}(\mathbf{u})\}$  converges to  $Z'_1(\mathbf{u})$ , where  $Z'_1(\mathbf{u})$  is a Gaussian process with mean function  $m_1(\mathbf{u}, \delta)$  and covariance kernel  $k_1(\mathbf{u}_1, \mathbf{u}_2)$ . Hence, under  $H_n$ , the asymptotic power of the test based on  $V_n$  is  $P_\delta[\int_{\mathbf{u} \in S(\delta)} \|Z'_1(\mathbf{u})\|^2 d\mathbf{u} > c_1(\alpha)]$ .

Similarly, using Lecam's third lemma, one can show that under contiguous alternatives  $\sqrt{n}\{F_n(\mathbf{x}_1) - F_0(\mathbf{x}_1), \dots, F_n(\mathbf{x}_l) - F_0(\mathbf{x}_l)\}$  is asymptotically  $l$ -dimensional multivariate normal with mean having  $j$ -th component  $m'_1(\mathbf{x}_j, \delta)$  and  $l \times l$ -dimensional covariance matrix obtained from the covariance kernel  $k_3$  as given in the statement of the theorem. Now, it follows from the the finite dimensional asymptotic distribution and the tightness of the process  $\sqrt{n}\{F_n(\mathbf{x}) - F_0(\mathbf{x})\}$  under contiguous alternatives that the stochastic process  $\sqrt{n}\{F_n(\mathbf{x}) - F_0(\mathbf{x})\}$  converges to  $Z''_1(\mathbf{x})$ , where  $Z''_1(\mathbf{x})$  is a Gaussian process with mean function  $m'_1(\mathbf{x}, \delta)$  and covariance kernel  $k_3(\mathbf{x}_1, \mathbf{x}_2)$ . Consequently, under  $H_n$ , the asymptotic powers of the tests based on  $T_n^{(1)}$  and  $T_n^{(2)}$  are  $P_\delta[\sup_{\mathbf{x}} |Z''_1(\mathbf{x})| > c_1^*(\alpha)]$  and  $P_\delta[\int_{\mathbf{x}} \{Z''_1(\mathbf{x})\}^2 dF_0(\mathbf{x}) > c_1^{**}(\alpha)]$ , respectively. This completes the proof.  $\square$

*Proof of Theorem 5.2:* The likelihood ratio for testing  $H_0$  against  $H_{n,m}$  is

$$\begin{aligned}
L_{n,m} &= \log \frac{\prod_{i=1}^n f(\mathbf{x}_i) \prod_{j=1}^m \{(1-\beta)f(\mathbf{y}_j) + \beta h(\mathbf{y}_j)\}}{\prod_{i=1}^n f(\mathbf{x}_i) \prod_{j=1}^m f(\mathbf{y}_j)} \\
&= \sum_{j=1}^m \log \left\{ 1 + \frac{\delta}{\sqrt{n+m}} \left( \frac{h(\mathbf{y}_j)}{f(\mathbf{y}_j)} - 1 \right) \right\} \quad (\text{using } \beta = \frac{\delta}{\sqrt{n+m}}) \\
&= \frac{\delta}{\sqrt{n+m}} \sum_{j=1}^m k'_j - \frac{\delta^2}{2(n+m)} \times \sum_{j=1}^m k_j'^2 + R_{n,m} \quad (\text{here } k'_j = \frac{h(\mathbf{y}_j)}{f(\mathbf{y}_j)} - 1). \quad (3)
\end{aligned}$$

Note that  $R_{n,m} \xrightarrow{p} 0$  as  $n, m \rightarrow \infty$ . As a consequence of C.L.T., W.L.L.N. and Slutsky's theorem, it follows from (3) that  $L_{n,m}$  is asymptotically normal with mean  $-\frac{\delta^2}{2}(1-\lambda)E_F\left(\frac{h(\mathbf{y})}{f(\mathbf{y})} - 1\right)^2$  and variance  $\delta^2(1-\lambda)E_F\left(\frac{h(\mathbf{y})}{f(\mathbf{y})} - 1\right)^2$ . This fact ensures the contiguity of the alternative sequence of densities using the corollary to Lecam's First Lemma in Hajek and Sidak (1967, p. 204).

Now, here also, we consider  $\mathbf{u}_1, \dots, \mathbf{u}_k \in S(\delta)$  and  $\mathbf{x}_1, \dots, \mathbf{x}_l \in \mathbb{R}^d$ . Then, under  $H_0$ , one can establish that the joint distribution of  $\sqrt{n+m}\{\hat{Q}_F(\mathbf{u}_1) - \hat{Q}_G(\mathbf{u}_1), \dots, \hat{Q}_F(\mathbf{u}_k) - \hat{Q}_G(\mathbf{u}_k), F_n(\mathbf{x}_1) - G_m(\mathbf{x}_1), \dots, F_n(\mathbf{x}_l) - G_m(\mathbf{x}_l), L_{n,m}\}$  is asymptotically multivariate Gaussian. This asymptotic normality is a consequence of the difference based on two independent samples, Bahadur type linear expansion of the difference of spatial quantiles  $\hat{Q}_F(\mathbf{u}) - \hat{Q}_G(\mathbf{u})$  (see Chaudhuri (1996) for  $d \geq 2$  and Serfling (1980) for  $d = 1$ ), the expansion of log likelihood ratio  $L_{n,m}$  (given in (3)) and the fact that  $F_n(\mathbf{x})$  and  $G_m(\mathbf{x})$  are simple averages of i.i.d. random variables. Note that for any  $p = 1, \dots, k$ , when  $d \geq 2$ , the asymptotic covariance between  $\sqrt{n+m}\{\hat{Q}_F(\mathbf{u}_p) - \hat{Q}_G(\mathbf{u}_p)\}$  and  $\sqrt{n+m}L_{n,m}$  is

$$\begin{aligned}
m_2(\mathbf{u}_p, \delta) &= E_f \left[ \sqrt{n+m} \left[ \sum_{i=1}^n \left\{ D_1[Q_F(\mathbf{u}_p)]^{-1} \left\{ \frac{\mathbf{x}_i - Q_F(\mathbf{u}_p)}{\|\mathbf{x}_i - Q_F(\mathbf{u}_p)\|} + \mathbf{u}_p \right\} \right\} \right. \right. \\
&\quad \left. \left. - \sum_{j=1}^m \left\{ D_1[Q_F(\mathbf{u}_p)]^{-1} \left\{ \frac{\mathbf{y}_j - Q_F(\mathbf{u}_p)}{\|\mathbf{y}_j - Q_F(\mathbf{u}_p)\|} + \mathbf{u}_p \right\} \right\} \right] \times \frac{\delta}{\sqrt{n+m}} \sum_{j=1}^m \left\{ \frac{h(\mathbf{y}_j)}{f(\mathbf{y}_j)} - 1 \right\} \right] \\
&= -\sqrt{n+m}E_f \left[ \sum_{j=1}^m \left\{ D_1[Q_F(\mathbf{u}_p)]^{-1} \left\{ \frac{\mathbf{y}_j - Q_F(\mathbf{u}_p)}{\|\mathbf{y}_j - Q_F(\mathbf{u}_p)\|} + \mathbf{u}_p \right\} \right\} \right]
\end{aligned}$$

$$\begin{aligned}
& \times \frac{\delta}{\sqrt{n+m}} \sum_{j=1}^m \left\{ \frac{h(\mathbf{y}_j)}{f(\mathbf{y}_j)} - 1 \right\} \Big] \quad (\text{since } \mathbf{x} \text{ and } \mathbf{y} \text{ are independent}) \\
& = -\delta [D_1^F(Q(\mathbf{u}))]^{-1} E_h \left\{ \frac{\mathbf{y} - Q_F(\mathbf{u}_p)}{\|\mathbf{y} - Q_F(\mathbf{u}_p)\|} + \mathbf{u}_p \right\} \quad (\text{as } E_f \left\{ \frac{\mathbf{y} - Q_F(\mathbf{u}_p)}{\|\mathbf{y} - Q_F(\mathbf{u}_p)\|} + \mathbf{u}_p \right\} = \mathbf{0}).
\end{aligned}$$

For any  $p = 1, \dots, k$ , when  $d = 1$ , as in the case of  $d \geq 2$ , one can show that the asymptotic covariance between  $\sqrt{n+m}\{\hat{Q}_F(u_p) - \hat{Q}_G(u_p)\}$  and  $\sqrt{n+m}L_{n,m}$  is  $m_2(u_p, \delta) = -\delta E_h \left\{ \frac{(\frac{u_p+1}{2})^{-1} \mathbf{1}_{\{y \leq F^{-1}(\frac{u_p+1}{2})\}}}{f(F^{-1}(\frac{u_p+1}{2}))} \right\}$ . Arguing in a similar way as in the proof of Theorem 5.1, one can establish that under  $H_{n,m}$ , the process  $\sqrt{n}\{\hat{Q}_F(\mathbf{u}) - \hat{Q}_G(\mathbf{u})\}$  converges to  $Z'_2(\mathbf{u})$ , where  $Z'_2(\mathbf{u})$  is a Gaussian process with mean function  $m_2(\mathbf{u}, \delta)$  and  $k_2(\mathbf{u}_1, \mathbf{u}_2)$ , which is defined in the statement of Theorem 3.2. Hence, the asymptotic power of the test based on  $T_{n,m}$  is  $P_\delta[\int_{\mathbf{u} \in S(\delta)} \|Z'_2(\mathbf{u})\|^2 d\mathbf{u} > c_2(\alpha)]$ .

Also, one can show that for any  $j = 1, \dots, l$ , the asymptotic covariance between  $\sqrt{n+m}\{F_n(\mathbf{x}_l) - G_m(\mathbf{x}_l)\}$  and  $\sqrt{n+m}L_{n,m}$  is  $m'_2(\mathbf{x}_j, \delta) = -\delta\{H(\mathbf{x}_j) - F(\mathbf{x}_j)\}$ . Now, it follows from the the finite dimensional asymptotic distribution and the tightness of the process  $\sqrt{n+m}\{F_n(\mathbf{x}) - G_m(\mathbf{x})\}$  under contiguous alternatives that the stochastic process  $\sqrt{n+m}\{F_n(\mathbf{x}) - G_m(\mathbf{x})\}$  converges to  $Z''_2(\mathbf{x})$ , where  $Z''_2(\mathbf{x})$  is a Gaussian process with mean function  $m'_2(\mathbf{x}, \delta)$  and covariance kernel  $k_4(\mathbf{x}_1, \mathbf{x}_2)$ . Consequently, under  $H_{n,m}$ , the asymptotic powers of the tests based on  $T_{n,m}^{(1)}$  and  $T_{n,m}^{(2)}$  are  $P_\delta[\sup_{\mathbf{x}} |Z''_2(\mathbf{x})| > c_2^*(\alpha)]$  and  $P_\delta[\int_{\mathbf{x}} \{Z''_2(\mathbf{x})\}^2 dM(\mathbf{x}) > c_2^{**}(\alpha)]$ , respectively. This completes the proof.  $\square$

**Acknowledgment:** The research of the first author is partially supported by a grant from Council of Scientific and Industrial Research (CSIR), Government of India.

## References

- [1] Anderson, T. W. and Darling, D. A. (1954) A test of goodness of fit. *Journal of the American Statistical Association*, **49**, 765–769.
- [2] Babu, G. J. and Rao, C. R. (1988) Joint asymptotic distribution of marginal quantile functions in samples from a multivariate population. *Journal of Multivariate Analysis*, **27**, 15–23.
- [3] Baringhaus, L. and Franz, C. (2010) Rigid motion invariant two-sample tests. *Statistica Sinica*, **20**, 1333–1361.
- [4] Bickel, P. J. (1967) Some contributions to the theory of order statistics. *Proceedings of the Fifth Berkeley Symposium Mathematical Statistics and Probability* **1**, 575–591.
- [5] Bickel, P. J. (1969) A distribution free version of Smirnov two sample test in the p-variate case. *The Annals of Statistics*, **40**, 1–23.

- [6] Breckling, J. and Chambers, R. (1988) M-Quantiles. *Biometrika*, **75**, 761–777.
- [7] Chakraborty, B. (2001) On Affine Equivariant Multivariate Quantiles. *The Annals of the Institute of Statistical Mathematics*, **53**, 380–403.
- [8] Chambers, J., Cleveland, W., Kleiner, B. and Tukey, P. (1983) *Graphical Methods for Data Analysis*, Wadsworth.
- [9] Chaudhuri, P. (1996) On a Geometric Notion of Quantiles for Multivariate Data. *Journal of the American Statistical Association*, **91**, 862–872.
- [10] Doksum, K. (1974) Empirical Probability Plots and Statistical Inference for Nonlinear Models in the Two-Sample Case. *The Annals of Statistics*, **2**, 267–277.
- [11] Doksum, K. and Sievers, G. (1976) Plotting with Confidence: Graphical Comparisons of Two Populations. *Biometrika*, **63**, 421–434.
- [12] Easton, G. S. and McCulloch, R. E. (1990) A Multivariate Generalization of Quantile-Quantile Plots. *Journal of the American Statistical Association*, **85**, 376–386.
- [13] Friedman, J. H. and Rafsky, L. C. (1979) Multivariate Generalizations of the Wald-Wolfowitz and Smirnov Two-Sample Tests. *The Annals of Statistics*, **7**, 697–717.
- [14] Friedman, J. H. and Rafsky, L. C. (1981) Graphics for the Multivariate Two-Sample problem. *Journal of the American Statistical Association*, **76**, 277–287.
- [15] Gnanadesikan, R. and Wilk, M.B. (1968) Probability plotting methods for the analysis of data. *Biometrika*, **55**, 1–17.
- [16] Gnanadesikan, R. (1977) *Methods for Statistical Analysis of Multivariate Observations*. Wiley & Sons.
- [17] Hajek, J and Sidak, Z. (1967) *Theory of Rank Tests*. Academic press, New York.
- [18] Justel, A., Pena, D. and Zamar, R. (1997) A multivariate Kolmogorov-Smirnov test of goodness of fit. *Statistics and Probability Letters*, **35**, 251–259.
- [19] Koenker, R. (2005) *Quantile Regression*. Econometric Society Monographs.
- [20] Koltchinskii, V. (1997) M-Estimation, Convexity and Quantiles. *The Annals of Statistics*, **25**, 435–477.
- [21] Leslie, J. R., Stephens, M. A. and Fotopoulos, S. (1986) Asymptotic Distribution of the Shapiro-Wilk W for Testing for Normality. *The Annals of Statistics*, **14**, 1497–1506.
- [22] Lehmann, E. L. and Romano, J. H. (2005) *Testing Statistical Hypotheses*. Springer.
- [23] Liu, R., Parelius, J. M. and Singh, K. (1999) Multivariate Analysis by Data Depth: Descriptive Statistics, Graphics and Inference. *The Annals of Statistics*, **27**, 783–840.
- [24] Marden, J. (1998) Bivariate Q-Q plots and spider web plots. *Statistica Sinica*, **8**, 813–826.
- [25] Mottonen, J. and Oja, H. (1995) Multivariate spatial sign and rank methods. *Journal of Non-parametric Statistics*, **5**, 201–213.
- [26] Serfling, R. (1980) *Approximation Theorems of Mathematical Statistics*. John Wiley & Sons, New York.
- [27] Serfling, R. (2004) Nonparametric multivariate descriptive measures based on spatial quantiles. *Journal of Statistical Planning and Inference*, **123**, 259–278.
- [28] Shapiro, S. S. and Wilk, M. B. (1965) Analysis of variance test for normality (complete samples). *Biometrika*, **52**, 591–611.
- [29] Shorack, G. R. and Wellner, J. A. (1986) *Empirical Processes with Applications to Statistics*. John Wiley & Sons, New York.