# A consistent multivariate test of association based on ranks of distances

June 1, 2012

Ruth Heller

*Department of Statistics and Operations Research, Tel-Aviv university, Tel-Aviv, Israel. E-mail: ruheller@post.tau.ac.il*

Yair Heller

*E-mail: heller.yair@gmail.com*

Malka Gorfine

*Faculty of Industrial Engineering and Management, Technion – Israel Institute of Technology, Haifa, Israel. E-mail:gorfinm@ie.technion.ac.il*

### Abstract

We are concerned with the detection of associations between random vectors of any dimension. Few tests of independence exist that are consistent against all dependent alternatives. We propose a powerful test that is applicable in all dimensions and is consistent against all alternatives. The test has a simple form and is easy to implement. We demonstrate its good power properties in simulations and on examples.

## 1   Introduction

In modern applications, there is need to test for independence between random vectors. One example from genomics research is whether two groups of genes are associated. Another application is functional magnetic resonance imaging research, where voxels in the brain are measured over time under various experimental conditions, and it is of interest to discover whether sets of voxels that comprise different areas in the brain are functionally related.

Let $X \in \Re^p$ and $Y \in \Re^q$ be random vectors, where $p$ and $q$ are positive integers. We are interested in testing whether there is a relationship between the two vectors $X$ and $Y$. The null hypothesis states that the two vectors are independent,

$$H_0 : F_{XY} = F_X F_Y,$$

where the joint distribution of $(X, Y)$ is denoted by $F_{XY}$, and the distributions of $X$ and $Y$, respectively, by $F_X$ and $F_Y$. We are interested in the general alternative that

the vectors are dependent,

$$H_1 : F_{XY} \neq F_X F_Y.$$

There are $N$ independent copies $(x_i, y_i)$, $i = 1, \ldots, N$ from the joint distribution of $X$ and $Y$ for testing $H_0$. The dimensions of the vectors $p$ and $q$ may be much higher than $N$.

The purpose of this paper is to provide a powerful test of independence that is applicable in all dimensions, and is consistent against all alternatives. The test is based on the pairwise distances between the sample values of $X$ and of $Y$ respectively, $\{d_X(x_i, x_j) : i, j \in \{1, \ldots, N\}\}$, $\{d_Y(y_i, y_j) : i, j \in \{1, \ldots, N\}\}$. The only restriction on the distance metrics $d_X(\cdot, \cdot)$ and $d_Y(\cdot, \cdot)$ is that they are determined by norms. The test statistic is a function of ranks of these distances, and it can be expressed simply in closed form. It is proven to be consistent against all dependent alternatives.

Few multivariate tests of independence that are consistent against all alternatives are available to date. Fukumizu et al. (2008) suggest a test based on normalized cross-covariance operators on reproducing kernel Hilbert spaces. Bickel and Xu (2009) offer a test based on an approximation of Renyi correlation, since there is no explicit formula to compute the Renyi correlation. A very elegant test with a simple formula is provided in Szekely et al. (2007), and has been further investigated in Szekely and Rizzo (2009) and in the discussions that followed it. We revisit some of the examples of Szekely et al. (2007), and add new examples. In the examples considered our new test performs remarkably well in comparison to the test of Szekely et al. (2007).

## 2 The new test of independence

This section develops the new test of independence. To motivate the test, note that if $X$ and $Y$ are dependent and have a continuous joint density, then there exists a point $(x_0, y_0)$ in the sample space of $(X, Y)$, and radii $R_x$ and $R_y$ around $x_0$ and $y_0$, respectively, such that the joint distribution of $X$ and $Y$ is different than the product of the marginal distributions in the cartesian product of balls around $(x_0, y_0)$. Consider first an oracle that guesses such a point $(x_0, y_0)$ and radii $R_x$ and $R_y$.

Let $d(\cdot, \cdot)$ be the norm distance between two sample points, either in $X$ or in $Y$, so the distance between the vectors $x_i$ and $x_j$ from the distribution of $X$ is $d(x_i, x_j)$, and similarly the distance between the vectors $y_i$ and $y_j$ from the distribution of $Y$ is $d(y_i, y_j)$. Technically, this distance may be different for the samples of $X$ and for the samples of $Y$, but we omit this distinction for simplicity of notation. Consider the following two dichotomous random variables: $I\{d(x_0, X) \leq R_x\}$ and $I\{d(y_0, Y) \leq R_y\}$, where $I(\cdot)$ is the indicator function. We summarize the observed cross-classification of these two dichotomous random variables for the $N$ independent observations $k \in \{1, \ldots, N\}$ in Table 1, where $A_{11} = \sum_{k=1}^{N} I\{d(x_0, x_k) \leq R_x\} I\{d(y_0, y_k) \leq R_y\}$, $A_{12}, A_{21}, A_{22}$, defined similarly, and $A_{m\cdot}, A_{\cdot m} \quad m = 1, 2$, are the sum of the row or column, respectively.

Evidence against independence may be quantified by Pearson's chi-square test statistic, or the likelihood ratio test statistic, for $2 \times 2$ contingency tables. The test based on such a statistic is consistent, and its power for finite sample size depends on the choice of $(x_0, y_0)$, $R_x$ and $R_y$.

Table 1: The cross-classification of $I\{d(x_0, X) \leq R_x\}$ and $I\{d(y_0, Y) \leq R_y\}$

|  | $d(y_0, \cdot) \leq R_y$ | $d(y_0, \cdot) > R_y$ |  |
|---|---|---|---|
| $d(x_0, \cdot) \leq R_x$ | $A_{11}$ | $A_{12}$ | $A_{1\cdot}$ |
| $d(x_0, \cdot) > R_x$ | $A_{21}$ | $A_{22}$ | $A_{2\cdot}$ |
|  | $A_{\cdot 1}$ | $A_{\cdot 2}$ | $N$ |

Table 2: The cross-classification of $I\{d(x_i, X) \leq d(x_i, x_j)\}$ and $I\{d(y_i, Y) \leq d(y_i, y_j)\}$

|  | $d(y_i, \cdot) \leq d(y_i, y_j)$ | $d(y_i, \cdot) > d(y_i, y_j)$ |  |
|---|---|---|---|
| $d(x_i, \cdot) \leq d(x_i, x_j)$ | $A_{11}(i, j)$ | $A_{12}(i, j)$ | $A_{1\cdot}(i, j)$ |
| $d(x_i, \cdot) > d(x_i, x_j)$ | $A_{21}(i, j)$ | $A_{22}(i, j)$ | $A_{2\cdot}(i, j)$ |
|  | $A_{\cdot 1}(i, j)$ | $A_{\cdot 2}(i, j)$ | $N - 2$ |

Since we do not have an oracle that guesses well $(x_0, y_0)$, $R_x$ and $R_y$, in the sense that the test for independence by a $2 \times 2$ contingency tables will be powerful, we let the data guide us in these choices. For every sample point $i$, we choose it in its turn to be $(x_0, y_0)$. For every sample point $j \neq i$, we choose it in its turn to define $R_x = d(x_i, x_j)$ and $R_y = d(y_i, y_j)$. The $2 \times 2$ tables now comprise the remaining $N - 2$ points. The test aggregates the evidence against independence by summing over all $N(N - 1)$ test statistics from the $2 \times 2$ tables thus created.

Specifically, for fixed observations $i$ and $j$, consider the dichotomous random variables: $I\{d(x_i, X) \leq d(x_i, x_j)\}$ and $I\{d(y_i, Y) \leq d(y_i, y_j)\}$. Table 2 summarizes the observed cross-classification of these two dichotomous random variables for the $N - 2$ independent observations $k \in \{1, \ldots, N\}, k \neq i, k \neq j$, where $A_{11}(i, j) = \sum_{k=1, k \neq i, k \neq j}^{N} I\{d(x_i, x_k) \leq d(x_i, x_j)\} I(d(y_i, y_k) \leq d\{y_i, y_j\})$, $A_{12}, A_{21}, A_{22}$ defined similarly, and $A_{m\cdot}, A_{\cdot m}, m = 1, 2$, are the sum of the row or column, respectively.

Let

$$S(i, j) = \frac{(N - 2)\{A_{12}(i, j)A_{21}(i, j) - A_{11}(i, j)A_{22}(i, j)\}^2}{A_{1\cdot}(i, j)A_{2\cdot}(i, j)A_{\cdot 1}(i, j)A_{\cdot 2}(i, j)}.$$

This is the classic test statistic for Pearson's chi square test for $2 \times 2$ contingency tables.

To test for independence between the two random vectors $X$ and $Y$, we suggest as a test statistic $T = \sum_{i=1}^{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} S(i, j)$. For $i$ and $j$ with 0 in at least one of the margins, we set $S(i, j) = 0$. The $p$-value from the permutation test based on the statistic $T$ is the fraction of replicates of $T$ under random permutations of the indices of the $Y$ sample, that are at least as large as the observed statistic.

We say a point $(x_0, y_0)$ is a point of dependence if the joint density of $X$ and $Y$ is different than the product of the marginal densities of $X$ and $Y$ at $(x_0, y_0)$, defined formally in equation (1) in the Appendix for the mixed case where the coordinates may be both discrete and continuous. Theorem 2.1 states that the test is consistent for discrete random vectors with countable support, as well as for continuous random vectors, and for random vectors where some of the coordinates are discrete and others

continuous, if the density of the continuous random vectors is continuous around a point of dependence.

**Theorem 2.1** *For dependent random vectors $(X, Y)$, $X \in \Re^p$ and $Y \in \Re^q$, denote the discrete and continuous coordinates of $X$ by $u \subseteq \{1, \ldots, p\}$ and $v = u^c$, respectively, and similarly the discrete and continuous coordinates of $Y$ by $s \subseteq \{1, \ldots, q\}$ and $t = s^c$, respectively. The permutation test based on the statistic $T$, with distances $d_X(\cdot, \cdot)$ and $d_Y(\cdot, \cdot)$ determined by norms, is consistent if either*

1. *$X$ and $Y$ are continuous, i.e. $u$ and $s$ are empty sets, and there exists a point of dependence $(x_0, y_0)$ for which the joint density is continuous.*

2. *At least one of $X$ or $Y$ has discrete coordinates in addition to the continuous coordinates, i.e. at least one of $u$ and $s$ is non-empty and both $v$ and $t$ are non-empty, and there exists a point of dependence $(x_0, y_0)$ for which (i) there exists a ball around the atom $\{x_0(u), y_0(s)\}$ that contains only this atom, and (ii) the joint density of the continuous coordinates conditional on the discrete coordinates is continuous.*

3. *Both $X$ and $Y$ are discrete, i.e. $v$ and $t$ are empty sets.*

4. *$X$ is discrete and $Y$ is continuous, i.e. $v$ and $s$ are empty sets, and there exists a point of dependence $(x_0, y_0)$ for which the conditional density of $Y$ given $X$ is continuous.*

See Appendix for a proof of case 2. The proofs of the other cases are very similar yet simpler, and they are given in the Supplementary Material.

## 2.1  Computational Complexity

For $N$ sample points, the naive implementation of the test will require an order of magnitude of $N^3$ operations. We provide an algorithm to efficiently calculate the score $T$ in order of magnitude $N^2 \log N$. This is done by providing an algorithm which for a given $i$ calculates $\{S(i,j) : j = 1, \ldots, N, j \neq i\}$ in order of magnitude $N \log N$. We shall show that we can calculate $\{A_{11}(i,j), A_{12}(i,j), A_{21}(i,j), A_{22}(i,j) : j = 1, \ldots, N, j \neq i\}$ in $O(N \log N)$.

For fixed $i$, let us look at all the distances from sample $i$ according to $X$ and let us sort the samples according to distance. Without loss of generality, renumber the indices of the $N-1$ sample points other than $i$ to be $1, \ldots, N-1$, so that the $j$th observation is the $j$th nearest to $i$ in $X$. Denote the order of the distance from $i$ in $Y$ by $\pi(1) \cdots \pi(N-1)$. So the $j$th observation is the $\pi(j)$th nearest to $i$ in $Y$. $\pi(\cdot)$ is a permutation of $1, \ldots, N-1$. The entries in the above Table 2 may be expressed as a function of $j$, $\pi(j)$ and $inv(j)$, where $inv(j)$ is defined as the number of inversions of $j$ in the permutation $\pi$, i.e. $inv(j)$ is the number indices $k \in \{1, \ldots, j-1\}$ such that $\pi(k) \in \{\pi(j) + 1, \ldots, N-1\}$. From the definition of $A_{12}(i,j)$ it follows that $A_{12}(i,j) = inv(j)$, and similarly $A_{22}(i,j) = N - \pi(j) - inv(j)$. Since $A_{1.}(i,j) = j - 1$, the remaining counts of the $2 \times 2$ contingency table for $S(i,j)$ are $A_{11} = j - 1 - inv(j)$, $A_{21} = \pi(j) + inv(j) - j - 1$. Therefore, it is enough to show that each of the following steps

4

Table 3: The power ($SE \times 100$) for a test at level 0.05 from a sample of size $N = 50$ from unusual bivariate relations. The results are based on 1000 simulations for rows $1-5$ and on 50000 simulations for the null setting in row 6.

| Distribution | Dcov | new test |
|---|---|---|
| W | 0.853 (1.1) | 1.000 (0.0) |
| Diamond | 0.037 (0.3) | 0.662 (1.5) |
| Parabola | 0.975 (0.5) | 0.998 (0.1) |
| 2 Parabolas | 0.303 (1.4) | 1.000 (0.0) |
| Circle | 0.000 (0.0) | 0.993 (0.3) |
| 4 independent clouds | 0.050 (0.1) | 0.050 (0.1) |

takes order of magnitude $N \log N$: (1) renumber the indices according to increasing distance in $X$ from $i$; (2) compute $\{\pi(j) : j = 1, \ldots, N, j \neq i\}$; (3) compute $\{inv(j) : j = 1, \ldots, N, j \neq i\}$. Since sorting takes order of magnitude $N \log N$, steps (1) and (2) are performed in the required computational time. It remains to show that (3) can be computed in order of magnitude $N \log N$. We show the algorithm in the Supplementary Material.

# 3 Simulations

In the simulations, we compare the performance of our test and the dCov test of Szekely and Rizzo (2009). We chose the latter test for two reasons. First, it is the only consistent test of simple form that is available. Second, the superiority of the dCov test over classical tests in Puri and Sen (1971) has been demonstrated in Szekely et al. (2007). Moreover, our aim is to investigate the performance of our test for non-monotone relationships, and these classical tests, or related tests for higher dimensions found in Taskinen et al. (2005), are ineffective for testing non-monotone types of dependence (Szekely et al., 2007).

In all simulations, the dCov test was applied by calling the function $dcov.test$ implemented in the R package *energy* (Szekely and Rizzo, 2009) with 10000 permutation samples. The Euclidean distance was used as a distance metric.

We consider first the six simulated examples of unusual bivariate distributions in Newton (2009). These examples mimic those at the wikipedia.org page on Pearson correlation, see Supplementary Material for details. The example of 4 independent clouds is an example of a null distribution. Table 3 shows the power comparison between dCov and the new test for $N = 50$ sample points and a significance level $\alpha = 0.05$. Large differences are observed. The most pronounced difference is observed for the circle relation, where the power of the new test is 0.993 yet dCov has no power to detect the relation. For the diamond relation, the new test has a power of 0.662 yet the power of dCov is 0.037. The tests based on Pearson and Spearman correlations had a power of at most 0.16 in all examples.

Szekely et al. (2007) considered multivariate examples and compared them to like-

Table 4: The power ($SE \times 100$) of a test at level 0.05 per sample size from a 5 dimensional joint distribution, where $X \sim N(0, I_{5 \times 5})$ and $Y = \log(X^2)$ or $Y = (Y_1, \ldots, Y_5)$ has coordinates $Y_j = X_j \cdot \epsilon_j$, where $\epsilon_j \sim N(0, 1)$ independent of $X_j$. The results are based on 1000 simulations.

| Sample size | $Y = \log(X^2)$ | | $Y_j = X_j \cdot \epsilon_j$ | |
|---|---|---|---|---|
| | dCov | new test | dCov | new test |
| N=20 | 0.172 (1.2) | 0.299 (1.4) | 0.335 (1.5) | 0.554 (1.6) |
| N=30 | 0.290 (1.4) | 0.595 (1.6) | 0.384 (1.5) | 0.792 (1.3) |
| N=40 | 0.436 (1.6) | 0.819 (1.2) | 0.417 (1.6) | 0.920 (0.9) |
| N=50 | 0.629 (1.5) | 0.945 (0.7) | 0.443 (1.6) | 0.968 (0.6) |

Table 5: The power ($SE \times 100$) of a test at level 0.05 per sample size from a 5 dimensional joint distribution, where $Y_j = \beta_1 X_j + \beta_2 X_j^2 + \epsilon_j, j = 1, \ldots, m_1$ and $Y_j = \epsilon_j, j = m_1 + 1, \ldots, 5$, with $\epsilon_j \sim N(0, \sigma^2)$ independent of $X_j \sim N(0, 1)$. The results are based on 1000 simulations.

| $m_1$ | $\beta_1$ | $\beta_2$ | $\sigma^2$ | dCov | | new test | |
|---|---|---|---|---|---|---|---|
| | | | | N=20 | N=30 | N=20 | N=30 |
| 0 | 0 | 0 | 1 | 0.040 (0.6) | 0.047 (0.7) | 0.051 (0.7) | 0.047 (0.7) |
| 2 | 1 | 4 | 9 | 0.501 (1.6) | 0.637 (1.5) | 0.669 (1.5) | 0.984 (0.4) |
| 2 | 3 | 2.5 | 9 | 0.841 (1.2) | 0.963 (0.6) | 0.706 (0.5) | 0.998 (0.1) |

lihood ratio type of tests. In the following two examples from Szekely et al. (2007), none of the likelihood ratio type of tests considered performed well. Using our notation, the distribution of $X = (X_1, \ldots, X_5)$ is standard multivariate normal with 5 dimensions. First, let $Y$ be equal to $\log(X^2)$. Columns 2 and 3 of Table 4 shows the power of a test at level 0.05 for dCov as well as for the new test. The new test has a power of 0.82 for $N = 40$ sample points, whereas the power of dCov is 0.436. Second, let $Y = (Y_1, \ldots, Y_5)$ have coordinates $Y_j = X_j \cdot \epsilon_j$, where $\epsilon_j$ are independent standard normal variables and independent of $X_j$. Columns 4 and 5 of Table 4 show the power of a test at level 0.05 for dCov as well as for the new test. The new test has a power of 0.968 for $N = 50$ sample points, whereas the power of dCov is 0.443.

A more sophisticated scenario, which includes both a monotone and non-monotone component, is the following: $Y_j = \beta_1 X_j + \beta_2 X_j^2 + \epsilon_j, j = 1, \ldots, m_1$ and $Y_j = \epsilon_j, j = m_1 + 1, \ldots, 5$, with $\epsilon_j \sim N(0, \sigma^2)$ and $X_j \sim N(0, 1)$ for all $j$. Table 7 shows the power of a test at level 0.05 for dCov as well as for the new test for various values of $\beta_1, \beta_2, \sigma^2, m_1 \in \{0, 2\}$. Further results in 100 dimensions are included in the Supplementary Material. When $\beta_2$ is large relative to $\beta_1$, the power of the new test is better than that of dCov.

Finally, we consider an example where $X$ and $Y$ are both of dimension 1000, from a mixture distribution with 10 equally likely components. In the $i$th component, $i \in$

Table 6: The power ($SE \times 100$) of a test at level 0.05 per sample size from the joint distribution of 10 mixture components for random vectors of dimension 1000, each component is centered around a different mean and is either multivariate Cauchy or multivariate t with 3 degrees of freedom. The results are based on 200 simulations.

| Sample size | t (3df) | | Cauchy | |
|---|---|---|---|---|
| | dCov | new test | dCov | newtest |
| N=50 | 0.100 (2.1) | 0.570 (3.5) | 0.040 (1.4) | 0.130 (2.4) |
| N=100 | 0.190 (2.8) | 0.980 (1.0) | 0.050 (1.5) | 0.185 (2.7) |
| N=200 | 0.345 (3.4) | 1.000 (0.0) | 0.075 (1.9) | 0.390 (3.5) |
| N=300 | 0.620 (3.2) | 1.000 (0.0) | 0.020 (1.0) | 0.580 (3.5) |

$\{1, \ldots, 10\}$, $(X, Y)$ are the random variables $\{\mu_x(i) + \epsilon, \mu_y(i) + \eta\}$, where $\mu_x(i)$ and $\mu_y(i)$ are sampled (once) from the 1000 dimensional multivariate standard normal distribution, and $(\epsilon, \eta)$ are sampled independently from the multivariate Cauchy or multivariate $t$ with 3 degrees of freedom, with the identity correlation matrix. The dependency of $X$ and $Y$ is through the fixed pairs $\{\mu_x(i), \mu_y(i)\}, i = 1, \ldots, 10$ such that the data consists of 10 clouds around these pairs. See Supplementary Material for details. Table 6 shows the power of a test at level 0.05 for dCov as well as for the new test. The new test has a power of one for $N = 200$ sample points in the multivariate $t$ distribution, whereas the power of dCov is 0.23. For the multivariate cauchy distribution, dCov has no power even at $N = 300$, as expected since dCov is consistent only for distributions with finite first moments (Szekely et al., 2007). The power of the new test is 0.58 for $N = 300$ sample points. Moreover, for the multivariate normal distribution, the power for both tests is one for $N = 50$ sample points.

## 4    An example

In a homogeneous population, the dependence between single nucleotide polymorphysms (SNPs) on the same chromosome is weaker the farther the SNPs are from each other due to recombination (Lander and Schork, 1994). A question of interest is whether SNPs across chromosomes are independent. To answer this question we examined the DNA of a sample of 97 unrelated individuals of Han Chinese in Beijing, China, available from the HapMap project (The International HapMap Consortium, 2003). This sample is regarded to be of relatively homogeneous ancestry, since donors were required to have at least three Han Chinese grandparents. For the purpose of this example, we limit ourselves to chromosomes 21 and 22 and ask whether the SNPs on chromosome 21 are independent of the SNPs on chromosome 22. We first preprocessed the data by removing subjects with more than 30% missing SNPs on a chromosome, SNPs with missing subjects, and SNPs with minor allele frequency below 0.05. After preprocessing, 43 subjects remained. For each subject we had a vector of dimension 31,858 of SNPs from chromosome 21, and a vector of dimension 36,264 of SNPs from chromosome 22. The Euclidean distance was used as a distance metric. Our proposed

test was highly significant, with a $p$-value below $1 \times 10^{-4}$. The dCov test was also significant, with a $p$-value of $6 \times 10^{-4}$.

## 5   Final remarks

Pearson's chi-squared test statistic was originally proposed as an approximation to the log-likelihood ratio statistic, in our context

$$S_{LR}(i,j) = 2\sum_{k=1}^{2}\sum_{l=1}^{2} A_{kl}(i,j)\log[A_{kl}(i,j)/\{\frac{A_{\cdot l}(i,j)A_{k\cdot}(i,j)}{N-2}\}].$$

An alternative test statistic for independence may therefore be $T_{LR} = \sum_{i=1}^{N}\sum_{\substack{j=1 \\ j\neq i}}^{N} S_{LR}(i,j)$. In the simulation results considered, the permutation test with this test statistic had very similar power to the power of the suggested test.

After discovering that the random vectors are dependent, a natural question to ask is which sub-vectors are dependent. This can be done using multiple comparisons procedures, similar to post-hoc testing in the analysis of variance (Scheffe, 1959). Moreover, the larger the value of $S(i,j)$, the stronger the dependence between the variables $I\{d(x_i,X) \leq d(x_i,x_j)\}$ and $I\{d(y_i,Y) \leq d(y_i,y_j)\}$. Informally, if $S(i,j)$ is large and $d(x_i,x_j)$ and $d(y_i,y_j)$ are small, this suggests that the random vectors $X$ and $Y$ are dependent in balls of size $d(x_i,x_j)$ and $d(y_i,y_j)$ around $x_i$ and $y_i$. We plan to explore methods of localizing the dependency in future work.

## Acknowledgement

## Supplementary material

Supplementary material includes the proofs of cases 3 and 4 of the theorem, the algorithm for implementing the test in order of magnitude $N^2 \log(N)$, further simulations, and an additional one-dimensional real data example.

## Appendix

We shall prove Theorem 2.1 for the case where the index sets $u, v, s, t$ are all non-empty, since it is straightforward to adapt the proof to the cases where $u$ or $s$ are empty sets.

From henceforth, for notational convenience we shall repress the conditioning event and denote the joint and marginal densities conditional on the discrete coordinate values as $h\{x(v), y(t)\}, f\{x(v)\}$, and $g\{y(t)\}$ in place of $h\{x(v), y(t) \mid X(u) = x(u), Y(s) = y(s)\}, f\{x(v) \mid X(u) = x(u)\}$, and $g\{y(t) \mid Y(s) = y(s)\}$. Moreover, we denote $p\{x(u), y(s)\} = Pr\{X(u) = x(u), Y(s) = y(s)\}, p\{x(u)\} = Pr\{X(u) = x(u)\}$, and $p\{y(s)\} = Pr\{Y(s) = y(s)\}$.

If $H_0$ is false, and the point of dependence $(x_0, y_0)$ satisfies properties (i) and (ii) of Theorem 2.1. Without loss of generality, suppose

$$p\{x_0(u), y_0(s)\}h\{x_0(v), y_0(t)\} > p\{x_0(u)\}f\{x_0(v)\}p\{y_0(s)\}g\{y_0(t)\}. \quad (1)$$

Let $R_d$ be a positive constant smaller than both the radius of the ball around $x_0(u)$ that contains only $x_0(u)$, and the radius of the ball around $y_0(s)$ that contains only the point $y_0(s)$. Then the set $\{(x, y) : d(x, x_0) < R_d, d(y, y_0) < R_d\}$ contains only points with discrete coordinates $x(u) = x_0(u), y(s) = y_0(s)$. Moreover, since the joint density conditional on $\{x_0(u), y_0(s)\}$ is continuous, there exists a radius $R_c$ such that $p\{x_0(u), y_0(s)\}h\{x(v), y(t)\} > p\{x_0(u)\}f\{x(v)\}p\{y_0(s)\}g\{y(t)\}$ for all points $(x, y)$ in the set $\{(x, y) : d(x, x_0) < R_c, d(y, y_0) < R_c, x(u) = x_0(u), y(s) = y_0(s)\}$. Let $R = \min\{R_d, R_c\}$ and $\mathcal{A} = \{(x, y) : d(x, x_0) < R, d(y, y_0) < R\}$. Then the set $\mathcal{A}$ has positive probability, for all points $(x, y) \in \mathcal{A}$ the discrete coordinates are $x(u) = x_0(u)$ and $y(s) = y_0(s)$, and moreover

$$\min_{\mathcal{A}}[p\{x(u), y(s)\}h\{x(v), y(t)\} - p\{x(u)\}f\{x(v)\}p\{y(s)\}g\{y(t)\}] > 0.$$

Denote this minimum by the positive constant $c$.

Clearly the following two subsets of $\mathcal{A}$ have positive probability as well:

$$\mathcal{A}_1 = \{(x, y) : d(x, x_0) < R/8, d(y, y_0) < R/8\}$$

and

$$\mathcal{A}_2 = \{(x, y) : 3R/8 < d(x, x_0) < R/2, 3R/8 < d(y, y_0) < R/2\}.$$

Denote the probabilities of $\mathcal{A}_1$ and $\mathcal{A}_2$ by $f_1$ and $f_2$ respectively. Therefore, we expect $(Nf_1)(Nf_2)$ pairs of sample points $i$ and $j$ such that $(x_i, y_i) \in \mathcal{A}_1$ and $(x_j, y_j) \in \mathcal{A}_2$. For these sample points $i$ and $j$,

$$3R/8 \le d(x_j, x_0) \le d(x_j, x_i) + d(x_i, x_0) \le d(x_j, x_i) + R/8 \quad (2)$$

where the second inequality is the triangle inequality, and the first and third inequalities follow since $(x_j, y_j) \in \mathcal{A}_2$ and $(x_i, y_i) \in \mathcal{A}_1$. It follows from (2) that

$$d(x_i, x_j) \ge R/4, \quad d(y_i, y_j) \ge R/4. \quad (3)$$

Moreover, if a sample point $k$ is closer to $i$ than to $j$ both in the $X$ vector and in the $Y$ vector, then it is within the $x$ and $y$ spheres of radius $R$:

**Lemma .1** *If $d(x_k, x_i) < d(x_i, x_j)$, then $d(x_k, x_0) \le R$. Similarly, if $d(y_k, y_i) < d(y_i, y_j)$, then $d(y_k, y_0) \le R$.*

Proof: Since the proof follows the same steps for $x_k$ and $y_k$, we only show it for the $x$ coordinates. The result follows by applying the triangle inequality several times,

$$
\begin{aligned}
d(x_k, x_0) \quad & \leq d(x_k, x_i) + d(x_i, x_0) \leq d(x_j, x_i) + d(x_i, x_0) \\
& \leq d(x_j, x_0) + 2d(x_i, x_0) \leq R/2 + 2R/8 = 6R/8 \leq R.
\end{aligned}
$$

The consequence of Lemma .1 is that for all such samples $k$, $(x_k, y_k) \in \mathcal{A}$.

Moreover, all points that are within the $x$ and $y$ spheres of radius $R/8$ are closer to $i$ than the point $j$:

**Lemma .2** *If $d(x_k, x_0) < R/8$, then $d(x_k, x_i) < d(x_i, x_j)$. Similarly, if $d(y_k, y_0) < R/8$, then $d(y_k, y_i) < d(y_i, y_j)$.*

Proof: Since the proof follows the same steps for $x_k$ and $y_k$, we only show it for the $x$ coordinates. Applying the triangle inequality, $d(x_k, x_i) \leq d(x_k, x_0) + d(x_i, x_0) \leq R/8 + R/8 = R/4$. The result follows from (3). Therefore, if $(x_k, y_k) \in \mathcal{A}_1$, then $k$ is closer to $i$ than to $j$ in both $X$ and $Y$.

By the law of large numbers, almost surely

$$
\lim_{N \to \infty} \frac{A_{11}(i, j)}{N - 2} = p\{x_0(u), y_0(s)\} \int_{\mathcal{A}_3} h\{x(v), y(t)\} dx(v) dy(t) \tag{4}
$$

$$
\lim_{N \to \infty} \frac{A_{1\cdot}(i, j)}{N - 2} = p\{x_0(u)\} \int_{\mathcal{A}_4} f\{x(v)\} dx(v) \tag{5}
$$

$$
\lim_{N \to \infty} \frac{A_{\cdot 1}(i, j)}{N - 2} = p\{y_0(s)\} \int_{\mathcal{A}_5} g\{y(t)\} dy(t) \tag{6}
$$

where $\mathcal{A}_3 = \{(x, y) : d(x, x_i) < d(x_i, x_j), d(y, y_i) < d(y_i, y_j)\}$, $\mathcal{A}_4 = \{x : d(x, x_i) < d(x_i, x_j)\}$, and $\mathcal{A}_5 = \{y : d(y, y_i) < d(y_i, y_j)\}$.

Recall that $S(i, j) = \sum_{k=1}^{2} \sum_{l=1}^{2} \{A_{k,l}(i, j) - A_{k\cdot}(i, j)A_{\cdot l}(i, j)/(N - 2)\}^2 / \{A_{k\cdot}(i, j)A_{\cdot l}(i, j)/(N - 2)\}$. It is enough to look at the term with $l = 1$ and $k = 1$ in $S(i, j)$, i.e. the term

$$
S_1(i, j) = \frac{\{A_{11}(i, j) - A_{1\cdot}(i, j)A_{\cdot 1}(i, j)/(N - 2)\}^2}{A_{1\cdot}(i, j)A_{\cdot 1}(i, j)/(N - 2)}.
$$

It follows that $S(i, j) \geq S_1(i, j)$, and therefore that our test statistic $T \geq \sum_{i=1}^{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} S_1(i, j)$.

By Slutzky's theorem and the continuous mapping theorem, almost surely

$$
\begin{aligned}
\lim_{N \to \infty} \frac{S_1(i, j)}{N - 2} &= \lim_{N \to \infty} \frac{1}{N - 2} \frac{\{A_{11}(i, j) - A_{1\cdot}(i, j)A_{\cdot 1}(i, j)/(N - 2)\}^2}{A_{1\cdot}(i, j)A_{\cdot 1}(i, j)/(N - 2)} \\
&= \frac{(\int_{\mathcal{A}_3} [p\{x_0(u), y_0(s)\} h\{x(v), y(t)\} - p\{x_0(u)\} f\{x(v)\} p\{y_0(s)\} g\{y(t)\}] dx(v) dy(t))^2}{\int_{\mathcal{A}_3} [p\{x_0(u)\} f\{x(v)\} p\{y_0(s)\} g\{y(t)\}] dx(v) dy(t)} \tag{7}
\end{aligned}
$$

We shall show that this limit can be bound from below by a positive constant that depends on $(x_0, y_0)$ but not on $i$ and $j$. From Lemma .1 it follows that $\mathcal{A}_3 \subseteq \mathcal{A}$, and

10

from Lemma .2 it follows that $\mathcal{A}_1 \subseteq \mathcal{A}_3$, and therefore a positive lower bound on the numerator of (7) can be obtained:

$$\int_{\mathcal{A}_3} [p\{x_0(u), y_0(s)\}h\{x(v), y(t)\} - p\{x_0(u)\}f\{x(v)\}p\{y_0(s)\}g\{y(t)\}]dx(v)dy(t)$$

$$\geq c \int_{\mathcal{A}_3} dx(v)dy(t) \geq c \int_{\mathcal{A}_1} dx(v)dy(t).$$

Moreover, $\int_{\mathcal{A}_3} \{p\{x_0(u)\}f\{x(v)\}p\{y_0(s)\}g\{y(t)\}\}dx(v)dy(t) \leq 1$. Therefore, denoting the lower bound by $c' = \{c \int_{\mathcal{A}_1} dx(v)dy(t)\}^2$, it follows that $S_1(i,j)/(N-2)$ converges almost surely to a constant larger than $c' > 0$. Therefore, $S_1(i,j) > (N-2)c'/2$ with probability going to 1 as $N \to \infty$. Since, moreover, the number of pairs of points $i$ and $j$ such that $(x_i, y_i) \in \mathcal{A}_1$ and $(x_j, y_j) \in \mathcal{A}_2$, divided by $f_1 f_2 N^2$, converges almost surely to 1, it follows that there exists a constant $\delta$ such that $\lim_{N \to \infty} Pr(T > \delta N^3) = 1$.

Under the null hypothesis, for large enough sample size $N$, $S(i,j)$ is distributed $\chi^2$ with 1 degree of freedom. Therefore, the null expectation of $T$ is approximately $N(N-1)$, and the null variance is bounded above by a term of order $N^4$ (more precisely, by $\{N(N-1)\}^2 2$). Since $\sum_{i=1}^{N} \sum_{\substack{j=1 \\ j \neq i}}^{N} S(i,j)$ is of order of magnitude of $N^3$, it follows that $T$ will be rejected with probability 1.

# References

Bickel, P. and Xu, Y. (2009). Discussion of: Brownian distance covariance. *The Annals of Applied Statistics*, 3 (4):1266–1269.

Fukumizu, K., Gretton, A., Sun, X., and Scholkopf, B. (2008). Kernel measures of conditional dependence. In *Advances in Neural Information Processing Systems*, pages 489–496.

Lander, E. and Schork, N. (1994). Genetic dissection of complex traits. *Science*, 265:2037–2048.

Newton, M. (2009). Introducing the discussion paper by Szekely and Rizzo. *The Annals of Applied Statistics*, 3 (4):1233–1235.

Puri, M. and Sen, P. (1971). *Nonparametric Methods in Multivariate Analysis*. John Wiley & Sons, Inc, New York.

Scheffe, H. (1959). *The Analysis of Variance*. John Wiley & Sons, Inc, New York.

Szekely, G. and Rizzo, M. (2009). Brownian distance covariance. *The Annals of Applied Statistics*, 3 (4):1236–1265.

Szekely, G., Rizzo, M., and Bakirov, N. (2007). Measuring and testing independence by correlation of distances. *The Annals of Statistics*, 35:2769–2794.

Taskinen, S., Oja, H., and Randles, R. (2005). Multivariate nonparametric tests of independence. *American Statistical Association*, 100 (471):916–925.

The International HapMap Consortium (2003). The International Hapmap Project. *Nature*, 426:789–796.

# A    Supplementary Material

## A.1    Proofs

The proof of case 1 is omitted, since it is very similar to the more complex case 2. The proofs of the countable case 3, and the mixed case where one random vector is discrete and the other continuous, are given, respectively, in Sections A.1.1 and A.1.2 below.

### A.1.1    Proof of the countable case 3

Suppose $X \in \Re^p$ and $Y \in \Re^q$ are both discrete with countable support. $H_0$ is false implies that there exists at least one pair of atoms $(x_0, y_0)$ such that $Pr(X = x_0, Y = y_0) > Pr(X = x_0)Pr(Y = y_0)$. We expect $NPr(X = x_0, Y = y_0)$ points to have values $(x_0, y_0)$. Let $i$ and $j$ be two such points. By the law of large numbers, almost surely

$$\lim_{N \to \infty} \frac{A_{11}(i,j)}{N - 2} = Pr(X = x_0, Y = y_0), \quad \lim_{N \to \infty} \frac{A_{1\cdot}(i,j)}{N - 2} = Pr(X = x_0), \quad \lim_{N \to \infty} \frac{A_{\cdot 1}(i,j)}{N - 2} = Pr(Y = y_0).$$

Recall that

$$S(i,j) = \sum_{k=1}^{2} \sum_{l=1}^{2} \{A_{k,l}(i,j) - A_{k\cdot}(i,j)A_{\cdot l}(i,j)/(N - 2)\}^2 / \{A_{k\cdot}(i,j)A_{\cdot l}(i,j)/(N-2)\}.$$

It is enough to look at the term with $l = 1$ and $k = 1$ in $S(i,j)$, i.e. the term

$$S_1(i,j) = \frac{\{A_{11}(i,j) - A_{1\cdot}(i,j)A_{\cdot 1}(i,j)/(N - 2)\}^2}{A_{1\cdot}(i,j)A_{\cdot 1}(i,j)/(N - 2)}.$$

It follows that $S(i,j) \geq S_1(i,j)$, and therefore that our test statistic $T \geq \sum_{i=1}^{N} \sum_{\substack{j \neq i \\ j=1}}^{N} S_1(i,j)$.

By Slutzky's theorem, almost surely

$$
\begin{aligned}
\lim_{N \to \infty} \frac{S_1(i,j)}{N - 2} &= \lim_{N \to \infty} \frac{1}{N - 2} \frac{\{A_{11}(i,j) - A_{1\cdot}(i,j)A_{\cdot 1}(i,j)/(N - 2)\}^2}{A_{1\cdot}(i,j)A_{\cdot 1}(i,j)/(N - 2)} \\
&= \frac{\{Pr(X = x_0, Y = y_0) - Pr(X = x_0)Pr(Y = y_0)\}^2}{Pr(X = x_0)Pr(Y = y_0)}.
\end{aligned}
$$

It follows that $S_1(i,j)/(N - 2)$ converges almost surely to a positive constant $c' > 0$. Therefore, $S_1(i,j) > (N - 2)c'/2$ with probability going to 1 as $N \to \infty$. Since we have order of magnitude of $N^2$ pairs of points $i$ and $j$ that satisfy the

12

inequality $S_1(i,j) > (N-2)c'/2$, it follows that there exists a constant $\delta$ such that $\lim_{N\to\infty} Pr(T > \delta N^3) = 1$. By the same argument as in the last paragraph of the Appendix in the main text, it therefore follows that $T$ will be rejected with probability 1.

### A.1.2 Proof of mixed case 4

Suppose $X \in \Re^p$ is discrete with countable support, and $Y \in \Re^q$ has a continuous density given $X$, denoted by $h(y \mid X = x)$, and a marginal density $g(y)$. $H_0$ is false implies that there exists at least one pair of points $x_0, y_0$ such that $Pr(X = x_0)h(Y = y_0 \mid X = x_0) > Pr(X = x_0)g(Y = y_0)$. Since $h(\cdot \mid X = x_0)$ is continuous, there exists a radius $R$ such that $Pr(X = x_0)h(Y = y \mid X = x_0) > Pr(X = x_0)g(Y = y)$ for $(x,y) \in \mathcal{A} = \{(x,y) : x = x_0, d(y, y_0) < R\}$. The set $\mathcal{A}$ has positive probability, and moreover

$$\min_{\mathcal{A}}\{Pr(X = x_0)h(Y = y \mid X = x_0) - Pr(X = x_0)g(Y = y)\} > 0.$$

Denote this minimum by the positive constant $c$.

Clearly the following two subsets of $\mathcal{A}$ have positive probability as well:

$$\mathcal{A}_1 = \{(x,y) : x = x_0, d(y, y_0) < R/8\}$$

and

$$\mathcal{A}_2 = \{(x,y) : x = x_0, 3R/8 < d(y, y_0) < R/2\}.$$

Denote the probabilities of $\mathcal{A}_1$ and $\mathcal{A}_2$ by $f_1$ and $f_2$ respectively. Therefore, we expect $(Nf_1)(Nf_2)$ pairs of sample points $i$ and $j$ such that $(x_i, y_i) \in \mathcal{A}_1$ and $(x_j, y_j) \in \mathcal{A}_2$.

For these sample points $i$ and $j$, $d(y_i, y_j) \geq R/4$. From Lemma 1 in the Appendix, if $d(y_k, y_i) < d(y_i, y_j)$, then $d(y_k, y_0) \leq R$. From Lemma 2 in the Appendix, if $d(y_k, y_0) < R/8$, then $d(y_k, y_i) < d(y_i, y_j)$. Therefore, if $(x_k, y_k) \in \mathcal{A}_1$, then $k$ is closer to $i$ than to $j$ in $Y$.

By the law of large numbers, almost surely

$$\lim_{N\to\infty} \frac{A_{11}(i,j)}{N-2} = Pr(X = x_0) \int_{\mathcal{A}_3} h(y \mid X = x_0)dy \tag{8}$$

$$\lim_{N\to\infty} \frac{A_{1\cdot}(i,j)}{N-2} = Pr(X = x_0) \tag{9}$$

$$\lim_{N\to\infty} \frac{A_{\cdot 1}(i,j)}{N-2} = \int_{\mathcal{A}_4} g(y)dy \tag{10}$$

where $\mathcal{A}_3 = \{(x,y) : x = x_0, d(y, y_i) < d(y_i, y_j)\}$, , and $\mathcal{A}_4 = \{y : d(y, y_i) < d(y_i, y_j)\}$ .

Recall that

$$S(i,j) = \sum_{k=1}^{2}\sum_{l=1}^{2} \{A_{k,l}(i,j) - A_{k\cdot}(i,j)A_{\cdot l}(i,j)/(N-2)\}^2/\{A_{k\cdot}(i,j)A_{\cdot l}(i,j)/(N-2)\}.$$

13

It is enough to look at the term with $l = 1$ and $k = 1$ in $S(i, j)$, i.e. the term

$$S_1(i, j) = \frac{\{A_{11}(i, j) - A_{1\cdot}(i, j)A_{\cdot 1}(i, j)/(N-2)\}^2}{A_{1\cdot}(i, j)A_{\cdot 1}(i, j)/(N-2)}.$$

It follows that $S(i, j) \geq S_1(i, j)$, and therefore that our test statistic $T \geq \sum_{i=1}^{N} \sum_{\substack{j \neq i \\ j=1}}^{N} S_1(i, j)$.

By Slutzky's theorem and the continuous mapping theorem, almost surely

$$\lim_{N \to \infty} \frac{S_1(i, j)}{N-2} = \lim_{N \to \infty} \frac{1}{N-2} \frac{\{A_{11}(i, j) - A_{1\cdot}(i, j)A_{\cdot 1}(i, j)/(N-2)\}^2}{A_{1\cdot}(i, j)A_{\cdot 1}(i, j)/(N-2)}$$

$$= \frac{Pr(X = x_0)[\int_{\mathcal{A}_3}\{h(y \mid X = x_0)dy - g(y)\}dy]^2}{\int_{\mathcal{A}_4} g(y)dy}$$

It follows that $S_1(i, j)/(N-2)$ converges almost surely to a positive constant $c' > 0$. Therefore, $S_1(i, j) > (N-2)c'/2$ with probability going to 1 as $N \to \infty$. Since we expect $(Nf_1)(Nf_2)$ pairs of sample points $i$ and $j$ that satisfy the inequality $S_1(i, j) > (N-2)c'/2$, it follows that there exists a constant $\delta$ such that $\lim_{N \to \infty} Pr(T > \delta N^3) = 1$. By the same argument as in the last paragraph in the Appendix of the main text, it therefore follows that $T$ will be rejected with probability 1.

## A.2  Computational Complexity

In this Section we give a $C$ implementation of the computation of $\{inv(j) : j = 1, \ldots, N, j \neq i\}$ in order of magnitude $N \log N$. The algorithm uses an adaptation of the classic merge sort algorithm. The basic idea is to split the array in half and sort each half while counting the number of inversions for each element in each half. In the merging stage of both halves, if an element in the right side is smaller than an element in the left side, it means that the number of inversions for the smaller element should be updated by adding to it the number of elements on the left side which are larger than it. The complexity of this algorithm $T(N)$ respects the recursion $T(N) = 2T(N/2) + O(N)$ and therefore it is $T(N) = O(N \log N)$. The C code is given below.

```
int Inversions(int *permutation, int *source, int
*inversion_count,int dim) {
    if (dim==1)
        return 0;
    else{
        Inversions(permutation, source, inversion_count, dim/2);
        Inversions(&permutation[dim/2], &source[dim/2], inversion_count,dim/2);
        Merge(permutation, source, inversion_count, dim);
    }
    return 0;
}

int Merge(int *permutation, int *source, int *inversion_count, int
dim) {
    int i;
```

```
    int left[MAX_DIM], right[MAX_DIM], left_source[MAX_DIM], right_source[MAX_DIM];
    int left_index=0, right_index=0;
    for (i=0;i<dim/2;i++){
        left[i]=permutation[i];
        left_source[i]=source[i];
    }
    for(i=0;i<dim/2;i++){
        right[i]=permutation[i+dim/2];
        right_source[i]=source[i+dim/2];
    }
    for(i=0;i<dim;i++){
        if ( (left_index<dim/2) && (right_index<dim/2)){
            if (left[left_index]<right[right_index]){
                permutation[i]=left[left_index];
                source[i]=left_source[left_index];
                left_index++;
            }
            else{
                permutation[i]=right[right_index];
                source[i]=right_source[right_index];
                printf("adding %d invs to %d\n", dim/2-left_index, source[i]);
                inversion_count[source[i]]+=(dim/2-left_index);
                right_index++;
            }
        }
        else{
            if (left_index<dim/2){
                permutation[i]=left[left_index];
                source[i]=left_source[left_index];
                left_index++;
            }
            if (right_index<dim/2){
                permutation[i]=right[right_index];
                source[i]=right_source[right_index];
                right_index++;
            }

        }
    }
    return 0;
}
```

## A.3   Simulations

In the simulations presented in the main text, we first considered the six simulated examples of unusual bivariate distributions. Figure 1 shows the scatter plots for a sample of size $N = 50$ from each of these distributions.

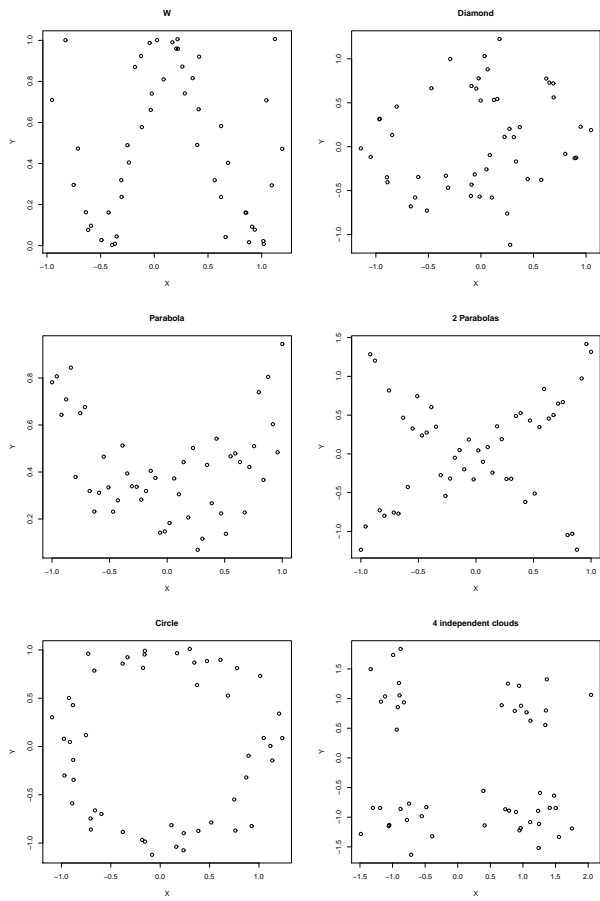In the simulations presented in the main text, the last example was of a mixture

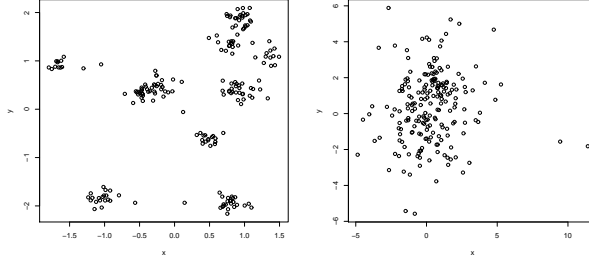Figure 1: Six simulated examples of unusual bivariate distributions; a sample of size N=50 from each distribution.

Figure 2: A scatter plot of the first coordinate in the mixture distribution of 10 components, where each coordinate has a t distribution with 3df around a different center. Left panel, noise 10 times smaller than generated; Right panel, noise used in the simulation.

distribution in 1000 dimensions. Figure 2 shows the first coordinate of $X$ and $Y$ in a setting where the standard deviation of the noise is 10 times smaller than actually generated (Left), as well as with the actual noise used in the simulation (Right panel), for the multivariate $t$ distribution with 3df.

A more sophisticated scenario in 100 dimensions, which includes both a monotone and non-monotone component, is the following: $Y_j = \beta_1 X_j + \beta_2 X_j^2 + \epsilon_j, j \in I_1$ and $Y_j = \epsilon_j, j \in \{1, \ldots, 100\} \backslash I_1$, with $\epsilon_j \sim N(0, \Sigma_X)$ and $X \sim N(0, \Sigma_X)$. The covariance matrix $\Sigma_X$ is block diagonal, with symmetric correlation of 0.9 in the first block, 0.8 in the second block, etc. The last block has 0 correlation, and the diagonal entries of $\Sigma_X$ are 1. In the null setting where $I_1 = \emptyset$, the empirical power for the new test, based on 1000 simulations, was 0.046, 0.043, and 0.051 for $N = 30, 40$, and 50, respectively. Table 7 shows the power of a test at level 0.05 for dCov as well as for the new test for $\beta_1 = 1, \beta_2 = 4, \sigma^2 = 9$, and two configurations of $I_1$. The power of the new test is better than that of dCov in the settings considered, in which the non-monotone part of the relationship has a stronger effect than the monotone part of the relationship. Moreover, the power of both tests is larger in the first setting, of strong dependence between the coordinates of $X$, than in the second setting, where the dependence across coordinates is weaker, since in the first setting the highly associated components of $X$ cause dependence between each coordinate of $Y$ with several coordinates of $X$.

## A.4 A univariate example

Szekely and Rizzo (2009) examined the Saviotti aircraft data of Saviotti (1996), that records six characteristics of aircraft designs during the twentieth century. They consider two variables, wing span (m) and speed (km/h) for the 230 designs of the third (of three) periods. This example and the data (aircraft) are from Bowman and Azzalini (1997). They showed that the dCov test of independence of log(Speed) and log(Span) in period 3 is significant (p-value $\leq 0.00001$), while the Pearson correlation test is not significant (p-value = 0.8001). Our proposed test is also highly significant (p-value

17

Table 7: The power of a test at level $0.05$ per sample size from a 100 dimensional joint distribution, where $Y_j = X_j + 4X_j^2 + \epsilon_j, j \in I_1$ and $Y_j = \epsilon_j, j \in \{1, \ldots, 100\} \backslash I_1$, with $\epsilon_j \sim N(0, 9)$. The results are based on 1000 simulations.

| $I_1$ | Sample size | dCov | new test |
|---|---|---|---|
| $\{1, \ldots, 10, 51, \ldots, 55\}$ | $N = 30$ | 0.382 | 0.629 |
| | $N = 40$ | 0.456 | 0.782 |
| | $N = 50$ | 0.541 | 0.879 |
| $\{41, \ldots, 50, 91, \ldots, 100\}$ | $N = 30$ | 0.246 | 0.243 |
| | $N = 40$ | 0.271 | 0.340 |
| | $N = 50$ | 0.293 | 0.474 |
| | $N = 60$ | 0.359 | 0.553 |
| | $N = 70$ | 0.369 | 0.626 |
| | $N = 80$ | 0.433 | 0.673 |

$\leq 0.00001$). Moreover, if we take a random sample of 30 observations and apply the dCov test and the proposed test to this small random sample, then we typically get smaller $p$-values using our proposed test than using the $dCov$ test. Specifically, repeating the testing of a random sample of 30 observations 100 times, the p-value of our proposed test was below 0.05 for 58/100 simulation runs, whereas for dCov only for 18/100 simulation runs. Figure 3 shows the scatter plot of wing span vs. speed on the log scale for a sample of 30 points. The relationship appears fan-like. For this particular sample, the $p$-value from the $dCov$ test and our proposed test were 0.21 and 0.03, respectively. Figure 4 shows the distribution of the 100 $p$-values for each of the tests.

# References

Bowman, A. and Azzalini, A. (1997). *Applied Smoothing Techniques for Data Analysis: The Kernel Approach with S-Plus Illustrations*. Oxford Univ. Press, Oxford.

Saviotti, P. (1996). *Technological Evolution, Variety and Economy*. Edward Elgar, Cheltenham.

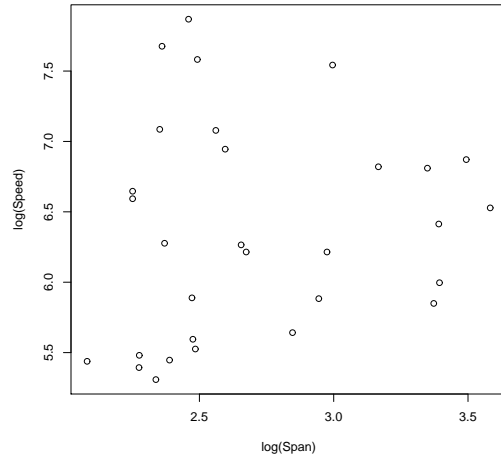Szekely, G. and Rizzo, M. (2009). Brownian distance covariance. *The Annals of Applied Statistics*, 3 (4):1236–1265.

Figure 3: The scatter of wing span vs. speed on the log scale for a sample of 30 points. The $p$-value from the $dCov$ test and our proposed test were 0.21 and 0.03, respectively.
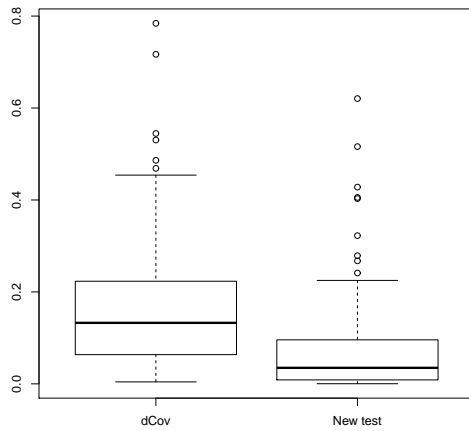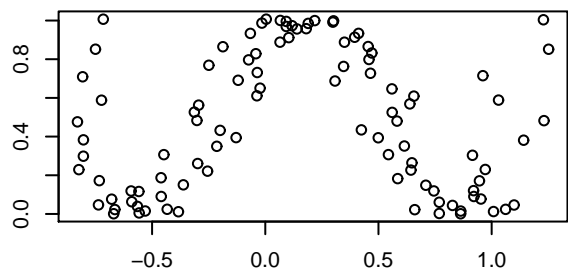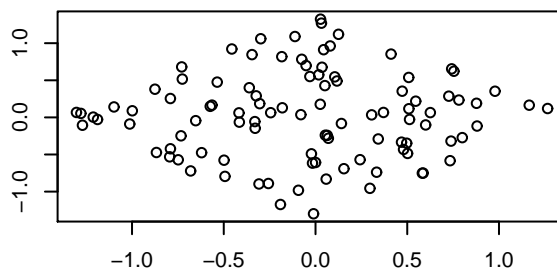


Figure 4: The boxplots of the 100 $p$-values for dCov and the proposed test based on a random sample of 30 points from the Aircraft data.
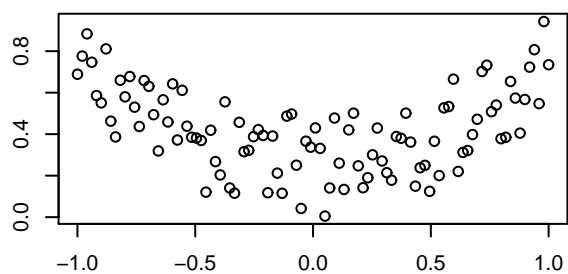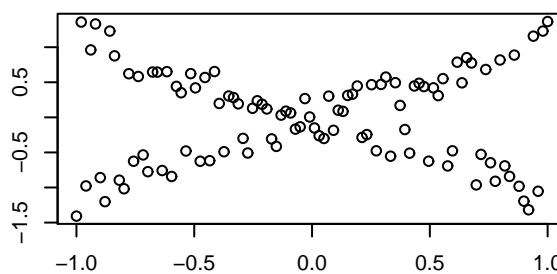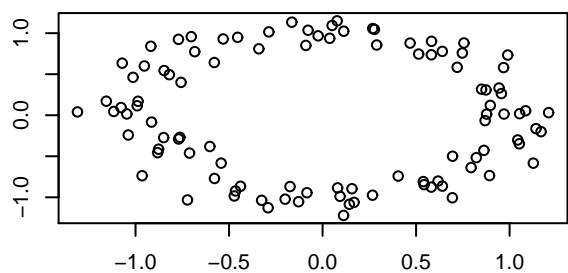
**W**

**Diamond**

**Parabola**

**Hyperbola**

**Circle**

**4 Clouds**