

## Two Notes on Energy Distance \*

CHEN Minqiong

(*Xinhua College of Sun Yat-Sen University, Guangzhou, 510520, China*)

**Abstract:** The definition of Brownian distance is presented and it's proved that Brownian distance coincides with the energy distance with respect to Brownian motion. Energy distance for dependent random vectors is also given and the asymptotic distribution is derived under null hypothesis. A simple numerical simulation result shows that the method for paired-sample test based on energy distance can distinguish the distributions of the paired variables more effectively than the classical t-test and Wilcoxon signed rank test.

**Keywords:** energy distance; Brownian distance; paired-sample test; asymptotic distribution; numerical simulation

**2010 Mathematics Subject Classification:** 62G99

**Citation:** CHEN M Q. Two notes on energy distance [J]. Chinese J Appl Probab Statist, 2018, 34(5): 463-474.

### §1. Introduction: Energy Distance

Székely<sup>[1]</sup> introduced a new concept named energy distance to measure the difference between two independent probability distributions. If  $X$  and  $Y$  are independent random vectors in  $\mathbb{R}^p$  with cumulative distribution functions (cdf)  $F$  and  $G$  respectively, then the energy distance between the distributions  $F$  and  $G$  is defined as

$$\varepsilon(F, G) = 2\mathbf{E}\|X - Y\| - \mathbf{E}\|X - X'\| - \mathbf{E}\|Y - Y'\|, \quad (1)$$

where  $X'$  is an i.i.d. copy of  $X$ , and  $Y'$  is an iid copy of  $Y$ .  $\mathbf{E}$  is the expected value, and  $\|\cdot\|$  denotes the Euclidean norm. One can also write  $\varepsilon(F, G)$  as  $\varepsilon(X, Y)$ , and call it be the energy distance of  $X$  and  $Y$ . Székely<sup>[1]</sup> proved that for real-valued random variables this distance is exactly twice Harald Cramér's distance, that is

$$2 \int_{-\infty}^{\infty} [F(t) - G(t)]^2 dt = 2\mathbf{E}\|X - Y\| - \mathbf{E}\|X - X'\| - \mathbf{E}\|Y - Y'\|.$$

\*The project was supported by the Natural Science Foundation of Guangdong Province for Young Innovative Talents (Grant No. 2014KQNCX253).

E-mail: C\_skirt@163.com.

Received December 13, 2016. Revised August 30, 2017.

In higher dimensions, however, the two distances are different because the energy distance is rotation invariant while Cramér's distance is not. The equality becomes

$$2\mathbb{E}\|X - Y\| - \mathbb{E}\|X - X'\| - \mathbb{E}\|Y - Y'\| = \frac{1}{c_p} \int_{\mathbb{R}^p} \frac{|\widehat{f}(t) - \widehat{g}(t)|^2}{\|t\|^{p+1}} dt, \quad (2)$$

where  $\widehat{f}$  is  $X$ 's characteristic function and  $\widehat{g}$  be  $Y$ 's characteristic function,  $c_p = \pi^{(p+1)/2} / \Gamma[(p+1)/2]$ . Thus  $\varepsilon(F, G) \geq 0$  with equality to zero if and only if  $F = G$ . This property makes it possible to use  $\varepsilon(F, G)$  for testing goodness-of-fit, homogeneity, etc. in a consistent way. Rizzo<sup>[2]</sup> discussed the theoretical background of homogeneity test for two multivariate populations based on energy distance. He considered the test

$$H_0 : F_1 = F_2$$

versus alternative  $F_1 \neq F_2$ . Suppose  $\mathbf{X} = \{X_1, X_2, \dots, X_m\}$  be an i.i.d. sample from the distribution of  $F_1$  in  $\mathbb{R}^p$ , and  $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_n\}$  be an i.i.d. sample from the distribution of  $F_2$  in  $\mathbb{R}^p$ . The sample version of  $\varepsilon(F_1, F_2)$  was defined as:

$$V_{m,n} := \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n \|X_i - Y_j\| - \frac{1}{m^2} \sum_{k,l=1}^m \|X_k - X_l\| - \frac{1}{n^2} \sum_{u,v=1}^n \|Y_u - Y_v\|. \quad (3)$$

If one denotes

$$h(x_1, x_2; y_1, y_2) := \|x_1 - y_1\| + \|x_2 - y_2\| - \|x_1 - x_2\| - \|y_1 - y_2\|,$$

then, obviously  $V_{m,n}$  is a two-sample  $V$ -statistic with kernel  $h$

$$V_{m,n} = \frac{1}{m^2 n^2} \sum_{i,k=1}^m \sum_{j,l=1}^n h(X_i, X_k; Y_j, Y_l).$$

Under the hypothesis of  $F_1 = F_2$ ,  $V_{m,n}$  is first-order degenerated, which leads  $[mn/(m+n)]V_{m,n}$  to converges in distribution to a quadratic form of centered Gaussian random variables<sup>[2]</sup>. So one can choose  $\varepsilon_{m,n} := [mn/(m+n)]V_{m,n}$  as the test statistic for  $H_0$ , and reject  $H_0$  for large values of  $\varepsilon_{m,n}$ . Rizzo<sup>[2]</sup> also presented the procedure for practical implementation via bootstrap method.

For more applications of energy distance on the classical statistic problems such as multivariate normality test, hierarchical clustering, dependence test and extension of analysis of variance, etc., see [3–6]. Székely and Rizzo summarized all those applications in [7]. Contrary to the classical methods, these methods based on energy distance are simple to calculate, applicable for more widely distributed types of data and can deal with multi-variables.

Rizzo<sup>[2]</sup> tells us that one can distinguish any two probability distributions on Euclidean space in a very simple way, just by computing the difference of the “average between sample distances” and “average within sample distances” of the two samples  $\mathbf{X} = \{X_1, X_2, \dots, X_m\}$  and  $\mathbf{Y} = \{Y_1, Y_2, \dots, Y_n\}$ . Here the distance only refers to the Euclidean distance. Lyons<sup>[8]</sup> generalized the notion of energy distance to probability distributions on metric spaces. Let  $(M, d)$  be a metric space with its Borel sigma algebra  $\mathcal{B}(M)$ ,  $\mathcal{P}(M)$  denotes the collection of all probability measures on the measurable space  $(M, \mathcal{B}(M))$ . If  $\mu$  and  $\nu$  are probability measures in  $\mathcal{P}(M)$ , then the energy distance of  $\mu$  and  $\nu$  can be defined as

$$\varepsilon(\mu, \nu) = 2\mathbb{E}[d(X, Y)] - \mathbb{E}[d(X, X')] - \mathbb{E}[d(Y, Y')],$$

provided that these expectations exist, where  $X, X' \stackrel{\text{i.i.d.}}{\sim} \mu$  and  $Y, Y' \stackrel{\text{i.i.d.}}{\sim} \nu$ . However the result “ $\varepsilon(\mu, \nu) \geq 0$  with equality to zero if and only if  $\mu = \nu$ ” does not always hold. Lyons<sup>[8]</sup> indicated that it holds unless the metric space  $(M, d)$  is strong negative. We describe this point more clearly as following proposition.

**Proposition 1** If  $(M, d)$  is a metric space of negative type, then, one has  $\varepsilon(\mu, \nu) \geq 0$ . Moreover, if  $(M, d)$  is strong negative, then  $\varepsilon(\mu, \nu) = 0$  if and only if  $\mu = \nu$ .

We give a explicit proof of this proposition.

**Proof** If  $(M, d)$  is a metric space of negative type, then there exists a Hilbert space  $H$  and a map  $\phi : M \rightarrow H$  such that  $\forall x, x' \in M, d(x, x') = \|\phi(x) - \phi(x')\|^2$  holds<sup>[10]</sup>, so

$$\begin{aligned} \mathbb{E}d(X, Y) &= \mathbb{E}\|\phi(X) - \phi(Y)\|^2 \\ &= \mathbb{E}\|\phi(X) - \mathbb{E}[\phi(X)] + \mathbb{E}[\phi(X)] - \mathbb{E}[\phi(Y)] + \mathbb{E}[\phi(Y)] - \phi(Y)\|^2 \\ &= \mathbb{E}\{\|\phi(X) - \mathbb{E}[\phi(X)]\|^2 + \|\mathbb{E}[\phi(X)] - \mathbb{E}[\phi(Y)]\|^2 + \|\phi(Y) - \mathbb{E}[\phi(Y)]\|^2 \\ &\quad + 2\langle \phi(X) - \mathbb{E}[\phi(X)], \mathbb{E}[\phi(X)] - \mathbb{E}[\phi(Y)] \rangle \\ &\quad + 2\langle \phi(X) - \mathbb{E}[\phi(X)], \mathbb{E}[\phi(Y)] - \phi(Y) \rangle \\ &\quad + 2\langle \mathbb{E}[\phi(X)] - \mathbb{E}[\phi(Y)], \mathbb{E}[\phi(Y)] - \phi(Y) \rangle\} \\ &= \mathbb{E}\|\phi(X) - \mathbb{E}[\phi(X)]\|^2 + \mathbb{E}\|\mathbb{E}[\phi(X)] - \mathbb{E}[\phi(Y)]\|^2 + \mathbb{E}\|\phi(Y) - \mathbb{E}[\phi(Y)]\|^2, \end{aligned}$$

and obviously

$$\begin{aligned} \mathbb{E}\|\phi(X) - \phi(X')\|^2 &= 2\mathbb{E}\|\phi(X) - \mathbb{E}[\phi(X)]\|^2, \\ \mathbb{E}\|\phi(Y) - \phi(Y')\|^2 &= 2\mathbb{E}\|\phi(Y) - \mathbb{E}[\phi(Y)]\|^2. \end{aligned}$$

Thus, we have

$$\begin{aligned} & 2\mathbf{E}d(X, Y) - \mathbf{E}d(X, X') - \mathbf{E}d(Y, Y') \\ &= 2\mathbf{E}\|\phi(X) - \phi(Y)\|^2 - \mathbf{E}\|\phi(X) - \phi(X')\|^2 - \mathbf{E}\|\phi(Y) - \phi(Y')\|^2 \\ &= 2\mathbf{E}\|\mathbf{E}[\phi(X)] - \mathbf{E}[\phi(Y)]\|^2, \end{aligned}$$

which means  $\varepsilon(X, Y) \geq 0$ .

Moreover, if  $(M, d)$  is strong negative, then, one can derive that the map  $\beta_\phi(\mu) : \mu \mapsto \int \phi(x)d\mu(x)$  is injective<sup>[8]</sup>, which implies  $\mathbf{E}[\phi(X)] = \mathbf{E}[\phi(Y)]$  if and only if  $\mu = \nu$ .  $\square$

Lyons<sup>[8]</sup> illustrated that all Euclidean spaces have strong negative type and proved that every separable Hilbert space is of strong negative type. Thus even if the observations are complex objects, like functions, texts, and graphs etc., one can use their real-valued nonnegative distances for inference.

In this paper, we give the next two notes on energy distance: Brownian distance and energy distance for dependent random vectors.

## §2. Brownian Distance

The notion of covariance with respect to a stochastic process, named Brownian distance covariance, was introduced by Székely and Rizzo<sup>[9]</sup>, and it was shown that the population distance covariance coincides with Brownian covariance. Similarly, we will present the notion of Brownian distance, and show the consistent of Brownian distance and energy distance.

**Definition 2** Let  $X$  and  $Y$  be two  $\mathbb{R}^p$ -valued random variables, with distributions  $F$  and  $G$  respectively. Suppose  $W$  be a Brownian motion with expectation zero and covariance function

$$\text{Cov}(W(t), W(s)) = \|s\| + \|t\| - \|s - t\| \quad (4)$$

on  $\mathbb{R}^p$ , independent of  $X$  and  $Y$ . The Brownian distance between  $X$  and  $Y$  is the non-negative number whose square is:

$$\text{BD}_W^2(X, Y) = \mathbf{E}\{[W(X) - W(Y)][W(X') - W(Y')]\}, \quad (5)$$

where  $X'$  is an iid copy of  $X$ , and  $Y'$  is an i.i.d. copy of  $Y$ .

For two independent real-valued random variables  $X, Y$ , the square of their mean's difference is

$$(\mathbf{E}X - \mathbf{E}Y)^2 = (\mathbf{E}X - \mathbf{E}Y)(\mathbf{E}X' - \mathbf{E}Y') = \mathbf{E}(X - Y)\mathbf{E}(X' - Y') = \mathbf{E}(X - Y)(X' - Y').$$

So, Brownian distance can be understood as the generalized definition of square of mean difference by replacing identity process  $I(t)$  with Brownian motion  $W(t)$ .

**Theorem 3** If  $X$  and  $Y$  are two independent  $\mathbb{R}^p$ -valued random variables, with  $E\|X\| + E\|Y\| < \infty$ , then  $E\{[W(X) - W(Y)][W(X') - W(Y')]\}$  is nonnegative and finite, and

$$BD_W^2(X, Y) = 2E\|X - Y\| - E\|X - X'\| - E\|Y - Y'\|.$$

**Proof** Observe that

$$\begin{aligned} BD_W^2(X, Y) &= E\{[W(X) - W(Y)][W(X') - W(Y')]\} \\ &= E\{E\{[W(X) - W(Y)][W(X') - W(Y')] \mid W\}\} \\ &= E\{E\{[W(X) - W(Y)] \mid W\}E\{[W(X') - W(Y')] \mid W\}\} \\ &= E\{E[W(X) \mid W] - E[W(Y) \mid W]\}^2, \end{aligned}$$

where  $E[W(X) \mid W] = \int_{\mathbb{R}^d} W(t)dF(t)$  and  $F(t)$  is the distribution function of  $X$ , and similar means for  $E[W(Y) \mid W]$ . So,  $E\{[W(X) - W(Y)][W(X') - W(Y')]\}$  is always nonnegative. For finiteness, note that

$$E[W^2(X)] = E\{E[W^2(X) \mid X]\} = E(2\|X\|) = 2E\|X\| < \infty.$$

Therefore, we have

$$\begin{aligned} BD_W^2(X, Y) &= E\{E[W(X) \mid W] - E[W(Y) \mid W]\}^2 \\ &\leq 2E\{E[W(X) \mid W]^2 + E[W(Y) \mid W]^2\} \\ &\leq 2E\{E[W^2(X) \mid W] + E[W^2(Y) \mid W]\} \\ &= 2\{E[W^2(X)] + E[W^2(Y)]\} = 4(E\|X\| + E\|Y\|) < \infty. \end{aligned}$$

On the other hand,

$$\begin{aligned} BD_W^2(X, Y) &= E\{[W(X) - W(Y)][W(X') - W(Y')]\} \\ &= E[W(X)W(X') + W(Y)W(Y') - W(X)W(Y') - W(Y)W(X')] \\ &= E\{E\{[W(X)W(X') + W(Y)W(Y') - W(X)W(Y') - W(Y)W(X')] \mid X, X', Y, Y'\}\} \\ &= E[\|X\| + \|X'\| - \|X - X'\| + \|Y\| + \|Y'\| - \|Y - Y'\| - 2(\|X\| + \|Y'\| - \|X - Y'\|)] \\ &= 2E\|X - Y\| - E\|X - X'\| - E\|Y - Y'\|. \end{aligned}$$

This completes the proof.  $\square$

Theorem 3 indicates that the properties of Brownian distance for random vectors  $X$  and  $Y$  are therefore the same properties for energy distance. We have the result as follows.

**Remark 4** Let  $X$  and  $Y$  be two  $\mathbb{R}^p$ -valued random variables, with distributions  $F$  and  $G$  respectively,  $W$  be a Brownian motion with expectation zero and covariance function (4) on  $\mathbb{R}^p$ , independent of  $X$  and  $Y$ , then, one has:

$$\text{BD}_W(X, Y) = 0 \text{ if and only if } X \stackrel{\text{D}}{=} Y,$$

where  $X \stackrel{\text{D}}{=} Y$  refers to  $X$  and  $Y$  are identically distributed.

### §3. Energy Distance for Paired Variables

From the introduction in section one, we can see that although the concept of energy distance is defined for independent variables, it's also applicable for dependent variables. In this section, we discuss the energy distance for paired variables which is a special case of dependent variables.

**Definition 5** Let  $X$  and  $Y$  be two  $\mathbb{R}^p$ -valued random variables, suppose  $(X, Y)$  has joint distribution  $H$ , with marginal distribution  $F$  on  $X$  and  $G$  on  $Y$  respectively. Assume that  $\text{E}\|X\| + \text{E}\|Y\| < \infty$ , the energy distance between  $X$  and  $Y$  is defined as

$$\epsilon(X, Y) := \text{E}\|X - Y'\| + \text{E}\|Y - X'\| - \text{E}\|X - X'\| - \text{E}\|Y - Y'\|, \quad (6)$$

where  $(X', Y')$  is an i.i.d. copy of  $(X, Y)$ .

Similar to equality (2) we have following result.

**Theorem 6** Denote  $\hat{f}$  as  $X$ 's characteristic function and  $\hat{g}$  as  $Y$ 's characteristic function, then, one has the equality

$$\text{E}\|X - Y'\| + \text{E}\|Y - X'\| - \text{E}\|X - X'\| - \text{E}\|Y - Y'\| = \frac{1}{c_p} \int_{\mathbb{R}^p} \frac{|\hat{f}(t) - \hat{g}(t)|^2}{\|t\|^{p+1}} dt, \quad (7)$$

where  $c_p = \pi^{(p+1)/2} / \Gamma[(p+1)/2]$ . Thus  $\epsilon(X, Y) \geq 0$  with equality to zero if and only if  $X$  and  $Y$  are identically distributed.

**Proof** Note that  $\hat{f}(t) = \text{E}e^{it'X}$ ,  $\hat{g}(t) = \text{E}e^{it'Y}$ , so

$$\begin{aligned} |\hat{f}(t) - \hat{g}(t)|^2 &= [\hat{f}(t) - \hat{g}(t)][\overline{\hat{f}(t) - \hat{g}(t)}] \\ &= \hat{f}(t)\overline{\hat{f}(t)} + \hat{g}(t)\overline{\hat{g}(t)} - \hat{f}(t)\overline{\hat{g}(t)} - \overline{\hat{f}(t)}\hat{g}(t) \\ &= \text{E}e^{it'X}\text{E}e^{-it'X'} + \text{E}e^{it'Y}\text{E}e^{-it'Y'} - \text{E}e^{it'X}\text{E}e^{-it'Y'} - \text{E}e^{-it'X'}\text{E}e^{it'Y} \\ &= \text{E}e^{it'(X-X')} + \text{E}e^{it'(Y-Y')} - \text{E}e^{it'(X-Y')} - \text{E}e^{it'(Y-X')} \\ &= 1 - \text{E}e^{it'(X-Y')} + 1 - \text{E}e^{it'(X'-Y)} - (1 - \text{E}e^{it'(X-X')}) - (1 - \text{E}e^{it'(Y-Y')}). \end{aligned}$$

Therefore, by integral on both sides, we get

$$\frac{1}{c_p} \int_{\mathbb{R}^p} \frac{|\widehat{f}(t) - \widehat{g}(t)|^2}{\|t\|^{p+1}} dt = \mathbb{E}\|X - Y'\| + \mathbb{E}\|Y - X'\| - \mathbb{E}\|X - X'\| - \mathbb{E}\|Y - Y'\|,$$

via the integral equation<sup>[7]</sup>

$$\int_{\mathbb{R}^p} \frac{1 - \cos\langle t, x \rangle}{\|t\|^{p+1}} dt = c_p \|x\|.$$

Thus, we have  $\epsilon(X, Y) \geq 0$  with equality to zero if and only if  $X$  and  $Y$  are identically distributed.  $\square$

**Remark 7** Let  $X$  be a  $\mathbb{R}^p$ -valued random variable, with distribution  $F$ , suppose  $\mathbb{E}\|X\| < \infty$ , then  $\mathbb{E}\|X + X'\| \geq \mathbb{E}\|X - X'\|$ , where  $X'$  is an i.i.d. copy of  $X$ ,  $\mathbb{E}\|X + X'\| = \mathbb{E}\|X - X'\|$  if and only if  $X$  is diagonally symmetric.

**Proof** Wang et al.<sup>[12]</sup> gives the proof of this result for  $X$  be univariate. Generally, we can put  $Y = -X$  to get

$$\begin{aligned} \epsilon(X, Y) &= \epsilon(X, -X) \\ &= \mathbb{E}\|X - (-X')\| + \mathbb{E}\|X' - (-X)\| - \mathbb{E}\|X - X'\| - \mathbb{E}\|-X - (-X')\| \\ &= 2(\mathbb{E}\|X + X'\| - \mathbb{E}\|X - X'\|). \end{aligned}$$

Thus, according to Theorem 6 we have  $\mathbb{E}\|X + X'\| - \mathbb{E}\|X - X'\| \geq 0$  and  $\mathbb{E}\|X + X'\| - \mathbb{E}\|X - X'\| = 0$  if and only if  $X$  and  $-X$  are identically distributed, which means that  $X$  is diagonally symmetric.  $\square$

We now give the empirical energy distance for paired variables.

**Definition 8** Let  $W_i = (X_i, Y_i)$ ,  $i = 1, 2, \dots, n$  be a sample from the distribution of  $(X, Y)$ , and denote  $\mathbf{W} = (\mathbf{X}, \mathbf{Y}) = \{W_1, W_2, \dots, W_n\}$ , the sample energy distance of  $(X, Y)$  is given as

$$\epsilon_n(\mathbf{X}, \mathbf{Y}) := \frac{1}{C_n^2} \sum_{i < j} (\|X_i - Y_j\| + \|Y_i - X_j\| - \|X_i - X_j\| - \|Y_i - Y_j\|).$$

If we denote

$$h(w_1, w_2) := h((x_1, y_1), (x_2, y_2)) = \|x_1 - y_2\| + \|y_1 - x_2\| - \|x_1 - x_2\| - \|y_1 - y_2\|,$$

then  $\epsilon_n(\mathbf{X}, \mathbf{Y})$  has the form of an U-statistic, with kernel  $h$

$$\epsilon_n(\mathbf{X}, \mathbf{Y}) = \frac{1}{C_n^2} \sum_{i < j} h(W_i, W_j).$$

Obviously,  $\epsilon_n(\mathbf{X}, \mathbf{Y})$  is an unbiased estimator of  $\epsilon(X, Y)$ .

**Proposition 9** The empirical energy distance  $\epsilon_n(\mathbf{X}, \mathbf{Y})$  almost surely converges to the energy distance  $\epsilon(X, Y)$ . That is,  $\epsilon_n(\mathbf{X}, \mathbf{Y}) \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \epsilon(X, Y)$ .

**Proof** As to  $h(W_1, W_2)$ , we have

$$\mathbf{E}|h(W_1, W_2)| \leq 4 * (\mathbf{E}\|X_1\| + \mathbf{E}\|Y_1\|) < \infty,$$

due to the suppose  $\mathbf{E}\|X\| + \mathbf{E}\|Y\| < \infty$  and  $\mathbf{E}[h(W_1, W_2)] = \epsilon(X, Y)$ . According to [11], we can obtain that

$$\epsilon_n(\mathbf{X}, \mathbf{Y}) \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \epsilon(X, Y).$$

Hence, the result follows.  $\square$

**Theorem 10** Let  $X$  and  $Y$  be two  $\mathbb{R}^p$ -valued random variables, with  $\mathbf{E}\|X\| + \mathbf{E}\|Y\| < \infty$ , then

(i) If  $X$  and  $Y$  are identically distributed, then

$$n\epsilon_n(\mathbf{X}, \mathbf{Y}) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \sum_{v=1}^{\infty} \lambda_v (Z_v^2 - 1),$$

where  $Z_v$  are independent standard normal random variables,  $\lambda_v$  are nonnegative constants that depend on the distribution of  $(X, Y)$ .

(ii) If  $X$  and  $Y$  are not identically distributed, then

$$n\epsilon_n(\mathbf{X}, \mathbf{Y}) \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \infty.$$

**Proof** (i) When  $X$  and  $Y$  are identically distributed, we have

$$\mathbf{E}h(W_1, W_2) = \mathbf{E}\|X_1 - Y_2\| + \mathbf{E}\|X_2 - Y_1\| - \mathbf{E}\|X_1 - X_2\| - \mathbf{E}\|Y_1 - Y_2\| = 0,$$

and

$$\begin{aligned} \bar{h}_1(w_1) &:= \mathbf{E}h(w_1, W_2) = \mathbf{E}h((x_1, y_1), (X_2, Y_2)) \\ &= \mathbf{E}\|x_1 - Y_2\| + \mathbf{E}\|X_2 - y_1\| - \mathbf{E}\|x_1 - X_2\| - \mathbf{E}\|y_1 - Y_2\| = 0, \end{aligned}$$

$$\bar{h}_1(w_2) := \mathbf{E}h(W_1, w_2) = 0,$$

$$\bar{h}_2(w_1, w_2) = \|x_1 - y_2\| + \|x_2 - y_1\| - \|x_1 - x_2\| - \|y_1 - y_2\|,$$

which means that  $\epsilon_n(\mathbf{X}, \mathbf{Y})$  is a degenerate U-statistic of order 1, therefore, it converges in distribution to a quadratic form:

$$n\epsilon_n(\mathbf{X}, \mathbf{Y}) \xrightarrow[n \rightarrow \infty]{\mathcal{D}} \sum_{v=1}^{\infty} \lambda_v (Z_v^2 - 1),$$

where  $Z_v$  are independent standard normal random variables,  $\lambda_v$  are the eigenvalues of the equation  $\int \bar{h}_2(w_1, w_2) f(w_2) dH(w_2) = \lambda f(w_1)$ .

(ii) The result is obviously because  $\epsilon_n(\mathbf{X}, \mathbf{Y}) \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \epsilon(X, Y) > 0$  when  $X$  and  $Y$  are not identically distributed.  $\square$

Theorem 10 implies that we can reject the null hypothesis of  $X$  and  $Y$  are identically distributed when  $n\epsilon_n(\mathbf{X}, \mathbf{Y})$  is too large.

#### §4. Paired-Sample Test Based on Energy Distance

In this section, we consider the the problem of testing the equality of distributions for paired variables based on the results we discussed in Section 3. For example, if we give two different treatments  $A, B$  to the same research object, denote  $X = (X_1, X_2, \dots, X_p)$  as the records of  $p$  features under  $A$  treatment, and  $Y = (Y_1, Y_2, \dots, Y_p)$  as the records under  $B$  treatment. If we want to compare the effect of this two different treatment, then, we have to test the equality of distributions of  $X$  and  $Y$ . In another case, if we give a treatment to a research object, also denote  $X$  the records before treatment, and  $Y$  the records after treatment. Again, to test whether the treatment is effective or not, we need to test the equality of distributions of  $X$  and  $Y$ .

The classical methods for testing equality of distributions for paired-sample include univariate t-test, Hotelling  $T^2$  test, sign test, and Wilcoxon signed rank test, etc. Univariate t-test and Hotelling  $T^2$  test are applicable to  $(X, Y)$  with joint distribution of two-dimensional or multidimensional normal distributions, and in fact are testing the equality of means for  $X, Y$ . When  $(X, Y)$  does not meet the joint normal distribution assumption, one can use nonparametric methods such as sign test and Wilcoxon signed rank test. As we all know, the sign test is actually to test whether the distribution of  $X - Y$  is with zero as median or not, while Wilcoxon signed rank test, is to test the symmetry of distribution of  $X - Y$ . So they are not consistent tests. Moreover, sign test and Wilcoxon signed rank test are only applicable for  $X, Y$  be both univariate.

The statistic

$$\epsilon_n(\mathbf{X}, \mathbf{Y}) = \frac{1}{C_n^2} \sum_{i < j} h(W_i, W_j)$$

presented in Section 3 gives a new idea to test the equality of distributions for paired-sample of arbitrary dimensions. Theorem 10 shows that  $n\epsilon_n(\mathbf{X}, \mathbf{Y})$  converges in distribution to a quadratic form  $\sum_{v=1}^{\infty} \lambda_v (Z_v^2 - 1)$  under null hypothesis of  $X$  and  $Y$  are identically distributed. But the  $\lambda$ s depend on the distribution of  $(X, Y)$ , so it's

difficult to compute the p-values by this result. Moreover, we need to be cautious of its serious limitations in practice, because the sample size in practice may be too small.

In practice, we consider the bootstrap method as an alternative approximation to the null distribution of  $n\epsilon_n(\mathbf{X}, \mathbf{Y})$ . Noticing that under  $H_0$ ,  $(X, Y)$  has the same distribution as  $(Y, X)$ , so we resample from the sample  $D_n := \{(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n), (Y_1, X_1), (Y_2, X_2), \dots, (Y_n, X_n)\}$ . Denote  $\mathbf{W}_n^* = \{(X_i^*, Y_i^*)\}_{i=1}^n$  the bootstrap sample obtained by sampling with replacement from  $D_n$ , and the bootstrap statistic is  $n\epsilon_n^*$ . Repeat this procedure for  $B$  times to obtain  $n\epsilon_{nk}^*$ ,  $k = 1, 2, \dots, B$ , then the p-value of the test is given by

$$p = \left[ 1 + \sum_{k=1}^B I(|n\epsilon_{nk}^*| > |n\epsilon_n|) \right] / (B + 1).$$

We conduct a simple Monte Carlo simulation to demonstrate the performance of the method we proposed here in comparison to t-test and Wilcoxon signed rank test for paired sample. Consider the following 6 examples:

Example 1:  $(X, Y) \sim N_2(\mu, \Sigma)$ , with  $\mu = (0, 0)'$ ,  $\Sigma = \begin{pmatrix} 1 & 0.2 \\ 0.2 & 1 \end{pmatrix}$ ;

Example 2:  $X \sim U(0, 6)$ ,  $Y = 6 - X$ ;

Example 3:  $X \sim B(10, 0.5)$ ,  $Y = 10 - X$ ;

Example 4:  $(X, Y) \sim N_2(\mu, \Sigma)$ , with  $\mu = (1, 1)'$ ,  $\Sigma = \begin{pmatrix} 1 & 0.8 \\ 0.8 & 4 \end{pmatrix}$ ;

Example 5:  $X \sim \text{Exp}(4)$ ,  $Y = 0.5e^{-4X}$ ;

Example 6:  $X \sim U(0, 1)$ ,  $Z \sim N(0, 1)$ ,  $Y = X - Z$ .

Examples 1–3 consider the cases that  $X$  and  $Y$  are identically distributed. Examples 4–6 are the cases that  $X$  and  $Y$  are not identically distributed, with Example 4 be the case that  $X$  and  $Y$  have the same expectation, and Example 6 be the case  $X - Y$  are symmetric distributed. We conduct the simulation with sample sizes  $n = 30, 50, 100, 200, 400$ , and the number of bootstrap resamples is 199. Each experiment is based on 500 replications, and we use the significance level of 0.05. Simulation results are summarized in Table 1, where **pctest** refers to the method we proposed in this paper.

The numerical simulation results of Examples 1–3 show that **pctest** can control the Type I error very well when  $X$  and  $Y$  are identically distributed, while Examples 4–6 illustrate that **pctest** can identify the distribution differences of the paired variables more effectively than the classical t-test and Wilcoxon signed rank test.

**Table 1** Proportion of rejections for Examples 1–6 at 5% level

model	method	$n = 30$	$n = 50$	$n = 100$	$n = 200$	$n = 400$
Example 1	t.test	0.054	0.030	0.052	0.030	0.040
	Wilcox.test	0.058	0.036	0.050	0.038	0.044
	<b>ptest</b>	<b>0.048</b>	<b>0.030</b>	<b>0.046</b>	<b>0.022</b>	<b>0.050</b>
Example 2	t.test	0.072	0.074	0.060	0.048	0.058
	Wilcox.test	0.072	0.068	0.060	0.048	0.062
	<b>ptest</b>	<b>0.062</b>	<b>0.066</b>	<b>0.052</b>	<b>0.044</b>	<b>0.058</b>
Example 3	t.test	0.080	0.068	0.070	0.048	0.050
	Wilcox.test	0.072	0.068	0.062	0.044	0.044
	<b>ptest</b>	<b>0.062</b>	<b>0.078</b>	<b>0.054</b>	<b>0.042</b>	<b>0.052</b>
Example 4	t.test	0.060	0.048	0.040	0.040	0.050
	Wilcox.test	0.048	0.038	0.050	0.034	0.042
	<b>ptest</b>	<b>0.574</b>	<b>0.884</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>
Example 5	t.test	0.046	0.050	0.038	0.042	0.042
	Wilcox.test	0.072	0.092	0.132	0.234	0.406
	<b>ptest</b>	<b>0.098</b>	<b>0.176</b>	<b>0.390</b>	<b>0.916</b>	<b>1.000</b>
Example 6	t.test	0.034	0.044	0.030	0.038	0.050
	Wilcox.test	0.032	0.040	0.022	0.038	0.056
	<b>ptest</b>	<b>0.994</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>	<b>1.000</b>

## References

- [1] SZÉKELY G J. E-statistics: the energy of statistical samples [R]. Technical Report No. 02-16, Bowling Green, Ohio: Bowling Green State University, Department of Mathematics and Statistics, 2002.
- [2] RIZZO M L. A test of homogeneity for two multivariate populations [C] // *2002 Proceedings of the American Statistical Association, Physical and Engineering Sciences Section* [CD-ROM], Alexandria, VA: American Statistical Association, 2002.
- [3] SZÉKELY G J, RIZZO M L. A new test for multivariate normality [J]. *J Multivariate Anal*, 2005, **93(1)**: 58–80.
- [4] SZÉKELY G J, RIZZO M L. Hierarchical clustering via joint between-within distances: extending Ward's minimum variance method [J]. *J Classification*, 2005, **22(2)**: 151–183.
- [5] SZÉKELY G J, RIZZO M L, BAKIROV N K. Measuring and testing dependence by correlation of distances [J]. *Ann Statist*, 2007, **35(6)**: 2769–2794.
- [6] RIZZO M L, SZÉKELY G J. DISCO analysis: a nonparametric extension of analysis of variance [J]. *Ann Appl Stat*, 2010, **4(2)**: 1034–1055.
- [7] SZÉKELY G J, RIZZO M L. Energy statistics: a class of statistics based on distances [J]. *J Statist Plann Inference*, 2013, **143(8)**: 1249–1272.
- [8] LYONS R. Distance covariance in metric spaces [J]. *Ann Probab*, 2013, **41(5)**: 3284–3305.

- [9] SZÉKELY G J, RIZZO M L. Brownian distance covariance [J]. *Ann Appl Stat*, 2009, **3**(4): 1236–1265.
- [10] SCHOENBERG I J. On certain metric spaces arising from Euclidean spaces by a change of metric and their imbedding in Hilbert space [J]. *Ann Math*, 1937, **38**(4): 787–793.
- [11] LEE A J. U-Statistics: Theory and Practice [M]. New York: Marcel Dekker, Inc., 1990.
- [12] WANG F, CUI H J, JIN J. Symmetric center test by using U-statistics [J]. *J Beijing Normal Univ (Nat Sci)*, 2009, **45**(1): 17–21. (in Chinese)

## 关于能量距离的两点注记

陈 敏 琼

(中山大学新华学院, 广州, 510520)

**摘 要:** 本文讨论了能量距离的两个问题. 类似 Brownian 协方差的讨论提出了 Brownian 距离的定义, 并证明了 Brownian 距离与能量距离的一致性. 给出了配对变量的能量距离的表示, 并探讨了将能量距离用于配对样本同分布的检验问题时原假设下的渐近分布理论. 最后通过一个简单的数值模拟说明基于能量距离的配对样本的分布差异的检验方法比传统的  $t$  检验及 Wilcoxon 符号秩检验更有效.

**关键词:** 能量距离; Brownian 距离; 配对变量; 渐近分布; 数值模拟

**中图分类号:** O212.7