

DISTANCE COVARIANCE FOR STOCHASTIC PROCESSES

MUNEYA MATSUI, THOMAS MIKOSCH, AND GENNADY SAMORODNITSKY

ABSTRACT. The distance covariance of two random vectors is a measure of their dependence. The empirical distance covariance and correlation can be used as statistical tools for testing whether two random vectors are independent. We propose an analogs of the distance covariance for two stochastic processes defined on some interval. Their empirical analogs can be used to test the independence of two processes.

The authors of this paper would like to congratulate Tomasz Rolski on his 70th birthday. We would like to express our gratitude for his longstanding contributions to applied probability theory as an author, editor, and organizer. Tomasz kept applied probability going in Poland and beyond even in difficult historical times. The applied probability community, including ourselves, has benefitted a lot from his enthusiastic, energetic and reliable work.

Sto lat! Niech zyje nam! Zdrowia, szczescia, pomyslnosci!

1. DISTANCE COVARIANCE FOR PROCESSES ON $[0, 1]$

We consider a real-valued stochastic process $X = (X(t))_{t \in [0,1]}$ with sample paths in a measurable space S such that X is measurable as a map from its probability space into S . We assume that the probability measure P_X generated by X on S is uniquely determined by its finite-dimensional distributions. Examples include processes with continuous or càdlàg sample paths on $[0, 1]$. The probability measure P_X is then determined by the totality of the characteristic functions

$$\varphi_X(\mathbf{x}_k; \mathbf{s}_k) = \varphi_X^{(k)}(\mathbf{x}_k; \mathbf{s}_k) = \int_S e^{i(s_1 f(x_1) + \dots + s_k f(x_k))} P_X(df), \quad k \geq 1,$$

where $\mathbf{x}_k = (x_1, \dots, x_k)' \in [0, 1]^k$, $\mathbf{s}_k = (s_1, \dots, s_k)' \in \mathbb{R}^k$. In particular, for two such processes, X and Y , the measures P_X and P_Y coincide if and only if

$$\varphi_X(\mathbf{x}_k; \mathbf{s}_k) = \varphi_Y(\mathbf{x}_k; \mathbf{s}_k) \quad \text{for all } \mathbf{x}_k \in [0, 1]^k, \mathbf{s}_k \in \mathbb{R}^k, k \geq 1.$$

We now turn from the general question of identifying the distributions of X and Y to a more specific but related one: given two processes X, Y on $[0, 1]$ with values in S as above and defined on the same probability space, we intend to find some means to verify whether X and Y are independent. Motivated by the discussion above, we need to show that the joint law of (X, Y) on $S \times S$, denoted by $P_{X,Y}$, coincides with the product measure $P_X \otimes P_Y$. Assuming, once again, that a probability measure on $S \times S$ is determined by the finite-dimensional distributions (as is the case with the aforementioned examples), we need to show that the joint characteristic functions of (X, Y) factorize, i.e.,

$$\begin{aligned} \varphi_{X,Y}(\mathbf{x}_k; \mathbf{s}_k, \mathbf{t}_k) &= \int_{S^2} e^{i \sum_{j=1}^k (s_j f(x_j) + t_j h(x_j))} P_{X,Y}(df, dh) \\ (1.1) \quad &= \varphi_X(\mathbf{x}_k; \mathbf{s}_k) \varphi_Y(\mathbf{x}_k; \mathbf{t}_k), \quad \mathbf{x}_k \in [0, 1]^k, \mathbf{s}_k, \mathbf{t}_k \in \mathbb{R}^k, k \geq 1. \end{aligned}$$

1991 *Mathematics Subject Classification*. Primary 62E20; Secondary 62G20 62M99 60F05 60F25.

Key words and phrases. Empirical characteristic function, distance covariance, stochastic process, test of independence.

Muneya Matsui's research is partly supported by JSPS Grant-in-Aid for Young Scientists B (16K16023) and Nanzan University Pache Research Subsidy I-A-2 for the 2016 academic year. Thomas Mikosch's research is partly supported by the Danish Research Council Grant DFF-4002-00435. Gennady Samorodnitsky's research is partly supported by the ARO MURI grant W911NF-12-1-0385.

Clearly, this condition is hard to check and therefore we try to get a more compact equivalent condition which can also be used for some statistical test of independence between X and Y .

For this reason, we consider a unit rate Poisson process $N = (N(t))_{t \in [0,1]}$ with arrivals $0 < T_1 < T_2 < \dots < T_{N(1)} \leq 1$, write $\mathbf{T}_N = (T_1, \dots, T_{N(1)})'$ and, correspondingly $\mathbf{s}_N, \mathbf{t}_N$ for any vectors in $\mathbb{R}^{N(1)}$. Then, for any positive probability density function g on \mathbb{R} , we define

$$\begin{aligned}
& d(P_{X,Y}, P_X \otimes P_Y) \\
&= \mathbb{E}_N \left[\int_{\mathbb{R}^{2N(1)}} |\varphi_{X,Y}(\mathbf{T}_N; \mathbf{s}_N, \mathbf{t}_N) - \varphi_X(\mathbf{T}_N; \mathbf{s}_N) \varphi_Y(\mathbf{T}_N; \mathbf{t}_N)|^2 \prod_{j=1}^{N(1)} g(s_j) g(t_j) ds_N dt_N \right] \\
&= \sum_{k=1}^{\infty} \mathbb{P}(N(1) = k) \int_{[0,1]^k} \left[\int_{\mathbb{R}^{2k}} |\varphi_{X,Y}(\mathbf{x}_k; \mathbf{s}_k, \mathbf{t}_k) - \varphi_X(\mathbf{x}_k; \mathbf{s}_k) \varphi_Y(\mathbf{x}_k; \mathbf{t}_k)|^2 \right. \\
(1.2) \quad & \left. \prod_{j=1}^k g(s_j) g(t_j) ds_k dt_k \right] d\mathbf{x}_k,
\end{aligned}$$

where in the last step we used the order statistics property of the homogeneous Poisson process. Here we interpret the summand corresponding to $k = 0$ as zero, and we also suppress the dependence on g in the notation. Now, the right-hand integrals vanish if and only if (1.1) is satisfied for Lebesgue a.e. $\mathbf{x}_k, \mathbf{s}_k, \mathbf{t}_k$, hence if and only if (1.1) holds for any $\mathbf{x}_k, \mathbf{s}_k, \mathbf{t}_k$. We summarize:

Lemma 1.1. *If g is a positive probability density on \mathbb{R} then $d(P_{X,Y}, P_X \otimes P_Y) = 0$ if and only if $P_{X,Y} = P_X \otimes P_Y$.*

Remark 1.2. Lemma 1.1 can easily be extended in several directions.

1. The statement remains valid when the Poisson probabilities $(\mathbb{P}(N(1) = k))_{k \geq 1}$ are replaced by any summable sequence of nonnegative numbers with infinitely many positive terms.
2. Obvious modifications of Lemma 1.1 are valid for random fields X, Y on $[0, 1]^d$, say (in this case we can sample the values of the random fields at the points of a Poisson random measure on $[0, 1]^d$ whose mean measure is the d -dimensional Lebesgue measure). Moreover, the values of X, Y may be multivariate.
3. The positive probability density $\prod_{j=1}^k g(s_j) g(t_j)$ on \mathbb{R}^{2k} can be replaced by any positive measurable function provided the infinite series in (1.2) is finite. This idea will be exploited in Section 3 below.
4. Returning to our original problem about identifying the laws of X and Y , similar calculations show that the quantity

$$d(P_X, P_Y) = \sum_{k=1}^{\infty} \mathbb{P}(N(1) = k) \int_{[0,1]^k} \left[\int_{\mathbb{R}^k} |\varphi_X(\mathbf{x}_k; \mathbf{s}_k) - \varphi_Y(\mathbf{x}_k; \mathbf{s}_k)|^2 \prod_{j=1}^k g(s_j) ds_k \right] d\mathbf{x}_k$$

vanishes if and only if $X \stackrel{d}{=} Y$. The quantity $d(P_X, P_Y)$ can be taken as the basis for a goodness-of-fit test for the distributions of X and Y .

In what follows, we refer to the quantities $d(P_{X,Y}, P_X \otimes P_Y)$ as *distance covariance* between the stochastic processes X and Y . This name is motivated by work on *distance covariance* for random vectors $\mathbf{X} \in \mathbb{R}^p, \mathbf{Y} \in \mathbb{R}^q$ (possibly of different dimensions) defined by

$$T(\mathbf{X}, \mathbf{Y}) = \int_{\mathbb{R}^{p+q}} |\varphi_{\mathbf{X},\mathbf{Y}}(\mathbf{s}, \mathbf{t}) - \varphi_{\mathbf{X}}(\mathbf{s}) \varphi_{\mathbf{Y}}(\mathbf{t})|^2 \mu(ds, dt),$$

where μ is a (possibly infinite) measure on \mathbb{R}^{p+q} ; see for example [1, 2, 6, 7, 9]. The last mentioned authors coined the names *distance covariance* and *distance correlation* for the standardized version $R(\mathbf{X}, \mathbf{Y}) = T(\mathbf{X}, \mathbf{Y}) / \sqrt{T(\mathbf{X}, \mathbf{X})T(\mathbf{Y}, \mathbf{Y})}$; they chose some special infinite measures μ which lead

to an elegant form of $T(\mathbf{X}, \mathbf{Y})$ and $R(\mathbf{X}, \mathbf{Y})$; see Section 3 for more information on this approach. The goal in the aforementioned literature was to find a statistical tool for testing independence between the vectors \mathbf{X} and \mathbf{Y} using the fact that $R(\mathbf{X}, \mathbf{Y}) = 0$ if and only if \mathbf{X}, \mathbf{Y} are independent provided μ has a positive Lebesgue density on \mathbb{R}^{p+q} . The sample versions $T_n(\mathbf{X}, \mathbf{Y})$ and $R_n(\mathbf{X}, \mathbf{Y}) = T_n(\mathbf{X}, \mathbf{Y}) / \sqrt{T_n(\mathbf{X}, \mathbf{X})T_n(\mathbf{Y}, \mathbf{Y})}$, constructed from an iid sample $(\mathbf{X}_i, \mathbf{Y}_i), i = 1, \dots, n$, of copies of (\mathbf{X}, \mathbf{Y}) , are then used as test statistics for checking independence of \mathbf{X} and \mathbf{Y} .

For stochastic processes X, Y on $[0, 1]$ one might be tempted to test their independence based on independent observations $\mathbf{X}_i = (X_i(x_1), \dots, X_i(x_k))'$, $\mathbf{Y}_i = (Y_i(x_1), \dots, Y_i(x_k))'$, $i = 1, \dots, n$ of the processes X, Y at the locations \mathbf{x}_k in $[0, 1]^k$. However, [8] observed that the empirical distance correlation $R_n(\mathbf{X}, \mathbf{Y})$ has the tendency to be very close to 1 even for relatively small values k . Our approach avoids the high dimensionality of the vectors \mathbf{X}_i and \mathbf{Y}_i by randomizing their dimension k .

Our paper is organized as follows. In Section 2 we study some of the theoretical properties of the distance covariance between two stochastic processes X, Y on $[0, 1]$ where we assume that g is a positive probability density. We find a tractable representation of this distance covariance from which we derive the corresponding sample version. In Section 3 we choose the non-integrable weight function g from [6]. Again, we find a suitable representation of this distance covariance, derive the corresponding sample version and show that it is a consistent estimator of its deterministic counterpart. In Section 4 we conduct a small simulation study based on the sample distance correlation introduced in Section 2. We compare the small sample behavior of the sample distance correlation with the corresponding sample distance correlation of [6] for independent and dependent Brownian and fractional Brownian sample paths.

2. PROPERTIES OF DISTANCE COVARIANCE

2.1. Distance correlation. In the context of stochastic processes X, Y one may be interested in standardizing the distance covariance $T(X, Y) = d(P_{X,Y}, P_X \otimes P_Y)$, i.e., in the *distance correlation*

$$R(X, Y) = \frac{T(X, Y)}{\sqrt{T(X, X)T(Y, Y)}}.$$

However, it is not obvious that $R(X, Y)$ assumes only values between 0 and 1. This property is guaranteed by a Cauchy-Schwarz argument.

Lemma 2.1. *Assume that $g(s) = g(-s)$. Then $0 \leq R(X, Y) \leq 1$.*

Proof. Let (X', Y') be an independent copy of (X, Y) . Applying the Cauchy-Schwarz inequality first to the k -dimensional integral with respect to the produce of k copies of g , then to the expectation with respect to the law of (X, Y) , next with respect to the Lebesgue measure on $[0, 1]^k$ and, finally, with respect to the law of N , and using the symmetry of the density g , we obtain

$$\begin{aligned} & d(P_{X,Y}, P_X \otimes P_Y) \\ &= \sum_{k=1}^{\infty} \mathbb{P}(N(1) = k) \int_{[0,1]^k} d\mathbf{x}_k \\ & \quad \mathbb{E} \left[\int_{\mathbb{R}^{2k}} \left[\left(e^{i \sum_{j=1}^k s_j X_j} - \varphi_X(\mathbf{x}_k; \mathbf{s}_k) \right) \left(e^{i \sum_{j=1}^k t_j Y_j} - \varphi_Y(\mathbf{x}_k; \mathbf{t}_k) \right) \right. \right. \\ & \quad \left. \left. \left(e^{-i \sum_{j=1}^k s_j X'_j} - \varphi_X(\mathbf{x}_k; -\mathbf{s}_k) \right) \left(e^{-i \sum_{j=1}^k t_j Y'_j} - \varphi_Y(\mathbf{x}_k; -\mathbf{t}_k) \right) \right] \prod_{j=1}^k g(s_j)g(t_j) ds_k dt_k \right] \\ & \leq \sum_{k=1}^{\infty} \mathbb{P}(N(1) = k) \int_{[0,1]^k} d\mathbf{x}_k \end{aligned}$$

$$\begin{aligned}
& \left(\mathbb{E} \left[\left| \int_{\mathbb{R}^k} \left(e^{i \sum_{j=1}^k s_j X_j} - \varphi_X(\mathbf{x}_k; \mathbf{s}_k) \right) \left(e^{-i \sum_{j=1}^k s_j X'_j} - \varphi_X(\mathbf{x}_k; -\mathbf{s}_k) \right) \prod_{j=1}^k g(s_j) ds_k \right|^2 \right] \right)^{1/2} \\
& \times \left(\mathbb{E} \left[\left| \int_{\mathbb{R}^k} \left(e^{i \sum_{j=1}^k t_j Y_j} - \varphi_Y(\mathbf{x}_k; \mathbf{t}_k) \right) \left(e^{-i \sum_{j=1}^k t_j Y'_j} - \varphi_Y(\mathbf{x}_k; -\mathbf{t}_k) \right) \prod_{j=1}^k g(t_j) dt_k \right|^2 \right] \right)^{1/2} \\
& = \sum_{k=1}^{\infty} \mathbb{P}(N(1) = k) \int_{[0,1]^k} d\mathbf{x}_k \\
& \quad \left[\int_{\mathbb{R}^{2k}} \left| \varphi_{X,X}(\mathbf{x}_k; \mathbf{s}_k, \mathbf{t}_k) - \varphi_X(\mathbf{x}_k; \mathbf{s}_k) \varphi_X(\mathbf{x}_k; \mathbf{t}_k) \right|^2 \prod_{j=1}^k g(s_j) g(t_j) ds_k dt_k \right]^{1/2} \\
& \quad \times \left[\int_{\mathbb{R}^{2k}} \left| \varphi_{Y,Y}(\mathbf{x}_k; \mathbf{s}_k, \mathbf{t}_k) - \varphi_Y(\mathbf{x}_k; \mathbf{s}_k) \varphi_Y(\mathbf{x}_k; \mathbf{t}_k) \right|^2 \prod_{j=1}^k g(s_j) g(t_j) ds_k dt_k \right]^{1/2} \\
& \leq \sum_{k=1}^{\infty} \mathbb{P}(N(1) = k) \\
& \quad \left[\int_{[0,1]^k} d\mathbf{x}_k \int_{\mathbb{R}^{2k}} \left| \varphi_{X,X}(\mathbf{x}_k; \mathbf{s}_k, \mathbf{t}_k) - \varphi_X(\mathbf{x}_k; \mathbf{s}_k) \varphi_X(\mathbf{x}_k; \mathbf{t}_k) \right|^2 \prod_{j=1}^k g(s_j) g(t_j) ds_k dt_k \right]^{1/2} \\
& \quad \times \left[\int_{[0,1]^k} d\mathbf{x}_k \int_{\mathbb{R}^{2k}} \left| \varphi_{Y,Y}(\mathbf{x}_k; \mathbf{s}_k, \mathbf{t}_k) - \varphi_Y(\mathbf{x}_k; \mathbf{s}_k) \varphi_Y(\mathbf{x}_k; \mathbf{t}_k) \right|^2 \prod_{j=1}^k g(s_j) g(t_j) ds_k dt_k \right]^{1/2} \\
& \leq (d(P_{X,X}, P_X \otimes P_X) d(P_{Y,Y}, P_Y \otimes P_Y))^{1/2}.
\end{aligned}$$

This proves that $0 \leq R(X, Y) \leq 1$. \square

2.2. Representations. Our next goal is to find explicit expressions for $d(P_{X,Y}, P_X \otimes P_Y)$. We observe that

$$\begin{aligned}
& \left| \varphi_{X,Y}(\mathbf{x}_k; \mathbf{s}_k, \mathbf{t}_k) - \varphi_X(\mathbf{x}_k; \mathbf{s}_k) \varphi_Y(\mathbf{x}_k; \mathbf{t}_k) \right|^2 \\
& = \left| \varphi_{X,Y}(\mathbf{x}_k; \mathbf{s}_k) \right|^2 + \left| \varphi_X(\mathbf{x}_k; \mathbf{s}_k) \right|^2 \left| \varphi_Y(\mathbf{x}_k; \mathbf{t}_k) \right|^2 - 2 \operatorname{Re} \left\{ \varphi_{X,Y}(\mathbf{x}_k; \mathbf{s}_k, \mathbf{t}_k) \varphi_X(\mathbf{x}_k; -\mathbf{s}_k) \varphi_Y(\mathbf{x}_k; -\mathbf{t}_k) \right\}.
\end{aligned}$$

This expression suggests to decompose (1.2) into 3 distinct parts, the first one being

$$\begin{aligned}
& \sum_{k=1}^{\infty} \frac{e^{-1}}{k!} \int_{[0,1]^k} \left[\int_{\mathbb{R}^{2k}} \left| \varphi_{X,Y}(\mathbf{x}_k; \mathbf{s}_k, \mathbf{t}_k) \right|^2 \prod_{j=1}^k g(s_j) g(t_j) ds_k dt_k \right] d\mathbf{x}_k \\
& = \int_{S^2} \sum_{k=1}^{\infty} \frac{e^{-1}}{k!} \left(\int_{[0,1]^k} \left[\int_{\mathbb{R}^{2k}} e^{i \sum_{r=1}^k (s_r (f(x_r) - f'(x_r)) + t_r (h(x_r) - h'(x_r)))} \right. \right. \\
& \quad \left. \left. \prod_{j=1}^k g(s_j) g(t_j) ds_k dt_k \right] d\mathbf{x}_k \right) P_{X,Y}(d(f, h)) P_{X,Y}(d(f', h')) \\
& = \int_{S^2} \sum_{k=1}^{\infty} \frac{e^{-1}}{k!} \left(\int_{[0,1]} \left[\int_{\mathbb{R}} e^{is(f(x) - f'(x))} g(s) ds \int_{\mathbb{R}} e^{it(h(x) - h'(x))} g(t) dt \right] dx \right)^k \\
& \quad P_{X,Y}(d(f, h)) P_{X,Y}(d(f', h')) \\
& = e^{-1} \int_{S^2} \left[\exp \left(\int_{[0,1]} \left[\int_{\mathbb{R}^2} e^{is(f(x) - f'(x)) + it(h(x) - h'(x))} g(s) g(t) ds dt \right] dx \right) - 1 \right]
\end{aligned}$$

$$P_{X,Y}(d(f, h)) P_{X,Y}(d(f', h')).$$

Similar calculations yield

$$\begin{aligned} d(P_{X,Y}, P_X \otimes P_Y) &= e^{-1} \int_{S^2} \left[\exp \left(\int_{[0,1]} \int_{\mathbb{R}} e^{is(f(x)-f'(x))} g(s) ds \int_{\mathbb{R}} e^{is(h(x)-h'(x))} g(s) ds dx \right) \right. \\ &\quad \times \left[P_{X,Y}(d(f, h)) P_{X,Y}(d(f', h')) + P_X \otimes P_Y(d(f, h)) P_X \otimes P_Y(d(f', h')) \right. \\ &\quad \left. \left. - P_{X,Y}(d(f, h)) P_X \otimes P_Y(d(f', h')) - P_{X,Y}(d(f', h')) P_X \otimes P_Y(d(f, h)) \right] \right]. \end{aligned}$$

We summarize our results:

Lemma 2.2. *The distance covariance between the processes X, Y on $[0, 1]$ with values in S can be written in the following form:*

$$\begin{aligned} e^1 d(P_{X,Y}, P_X \otimes P_Y) &= \mathbb{E} \left[\exp \left(\int_{[0,1]} \int_{\mathbb{R}} e^{is(X(x)-X'(x))} g(s) ds \int_{\mathbb{R}} e^{is(Y(x)-Y'(x))} g(ds) ds dx \right) \right] \\ &\quad + \mathbb{E} \left[\exp \left(\int_{[0,1]} \int_{\mathbb{R}} e^{is(X(x)-X'(x))} g(s) ds \int_{\mathbb{R}} e^{is(Y''(x)-Y'''(x))} g(s) ds dx \right) \right] \\ &\quad - 2 \operatorname{Re} \mathbb{E} \left[\exp \left(\int_{[0,1]} \int_{\mathbb{R}} e^{is(X(x)-X'(x))} g(s) ds \int_{\mathbb{R}} e^{is(Y(x)-Y''(x))} g(s) ds dx \right) \right], \end{aligned}$$

where (X', Y') is an independent copy of (X, Y) and Y'', Y''' are independent copies of Y which are also independent of X, X', Y, Y' .

Example 2.3. Let g be the density of a suitably scaled symmetric α -stable law on \mathbb{R} , $\alpha \in (0, 2]$. Then

$$\int_{\mathbb{R}} e^{is(f(x)-f'(x))} g(s) ds = e^{-|f(x)-f'(x)|^\alpha},$$

and so for a uniform random variable U on $(0, 1)$ which is independent of X, Y, X', Y', Y'', Y''' ,

$$\begin{aligned} d(P_{X,Y}, P_X \otimes P_Y) &= e^{-1} \mathbb{E} \left[\exp \left(\mathbb{E}_U e^{-|X(U)-X'(U)|^\alpha - |Y(U)-Y'(U)|^\alpha} \right) \right. \\ &\quad \left. + \exp \left(\mathbb{E}_U e^{-|X(U)-X'(U)|^\alpha - |Y''(U)-Y'''(U)|^\alpha} \right) \right. \\ (2.1) \quad &\quad \left. - 2 \exp \left(\mathbb{E}_U e^{-|X(U)-X'(U)|^\alpha - |Y(U)-Y''(U)|^\alpha} \right) \right]. \end{aligned}$$

2.3. Sample distance covariance. Let $(X_1, Y_1), \dots, (X_n, Y_n)$ be an iid sample with distribution $P_{X,Y}$ and let $P_{n,X,Y}$ be the corresponding empirical distribution with marginals $P_{n,X}$ and $P_{n,Y}$. Then we can define the *sample distance covariance* given by

$$\begin{aligned} e^1 d(P_{n,X,Y}, P_{n,X} \otimes P_{n,Y}) &= \frac{1}{n^2} \sum_{j_1=1}^n \sum_{j_2=1}^n \exp \left(\int_{[0,1]} \int_{\mathbb{R}} e^{is(X_{j_1}(x)-X_{j_2}(x))} g(s) ds \int_{\mathbb{R}} e^{is(Y_{j_1}(x)-Y_{j_2}(x))} g(s) ds dx \right) \\ &\quad + \frac{1}{n^4} \sum_{j_1=1}^n \sum_{j_2=1}^n \sum_{j_3=1}^n \sum_{j_4=1}^n \exp \left(\int_{[0,1]} \int_{\mathbb{R}} e^{is(X_{j_1}(x)-X_{j_2}(x))} g(s) ds \int_{\mathbb{R}} e^{is(Y_{j_3}(x)-Y_{j_4}(x))} g(s) ds dx \right) \\ &\quad - 2 \operatorname{Re} \frac{1}{n^3} \sum_{j_1=1}^n \sum_{j_2=1}^n \sum_{j_3=1}^n \exp \left(\int_{[0,1]} \int_{\mathbb{R}} e^{is(X_{j_1}(x)-X_{j_2}(x))} g(s) ds \int_{\mathbb{R}} e^{is(Y_{j_1}(x)-Y_{j_3}(x))} g(s) ds dx \right). \end{aligned}$$

Example 2.4. Assume that g is the density of a suitably scaled symmetric α -stable random variable. Then

$$e^1 d(P_{n,X,Y}, P_{n,X} \otimes P_{n,Y})$$

$$\begin{aligned}
&= \frac{1}{n^2} \sum_{j_1=1}^n \sum_{j_2=1}^n \exp \left(\int_{[0,1]} e^{-|X_{j_1}(x)-X_{j_2}(x)|^\alpha - |Y_{j_1}(x)-Y_{j_2}(x)|^\alpha} dx \right) \\
&\quad + \frac{1}{n^4} \sum_{j_1=1}^n \sum_{j_2=1}^n \sum_{j_3=1}^n \sum_{j_4=1}^n \exp \left(\int_{[0,1]} e^{-|X_{j_1}(x)-X_{j_2}(x)|^\alpha - |Y_{j_3}(x)-Y_{j_4}(x)|^\alpha} dx \right) \\
&\quad - \frac{1}{n^3} \sum_{j_1=1}^n \sum_{j_2=1}^n \sum_{j_3=1}^n \exp \left(\int_{[0,1]} e^{-|X_{j_1}(x)-X_{j_2}(x)|^\alpha - |Y_{j_1}(x)-Y_{j_3}(x)|^\alpha} dx \right).
\end{aligned}$$

Remark 2.5. The form of the sample distance covariance indicates that one needs to involve numerical methods for its calculation. In addition, in general we cannot assume that the sample paths of (X_i, Y_i) are completely observed. Then we need to approximate the path-dependent integrals appearing in the exponents of the expressions above by appropriate sums on a grid. These problems are not studied further in this paper.

The following result is an immediate consequence of the strong law of large numbers for U -statistics (see [3]) and the observation that $d(P_{n,X,Y}, P_{n,X} \otimes P_{n,Y})$ is a linear combination of U -statistics.

Proposition 2.6. *Assume that $((X_i, Y_i))_{i=1, \dots, n}$ is an iid sequence of S^2 -valued random elements. Then*

$$d(P_{n,X,Y}, P_{n,X} \otimes P_{n,Y}) \xrightarrow{\text{a.s.}} d(P_{X,Y}, P_X \otimes P_Y), \quad n \rightarrow \infty.$$

3. DISTANCE COVARIANCE WITH INFINITE WEIGHT MEASURES

So far we assumed that g is a positive Lebesgue density. In the aforementioned literature (see for example [6]) positive weight functions g were used which are not integrable over \mathbb{R} . In what follows, we consider an approach with suitable non-integrable weight functions which lead to a distance covariance for stochastic processes. To begin, note that if the function g is not necessarily integrable but is symmetric, then appealing to (1.2) and using the symmetry of both the cosine function and the function g we have

$$\begin{aligned}
&d(P_{X,Y}, P_X \otimes P_Y) \\
&= \sum_{k=1}^{\infty} \mathbb{P}(N(1) = k) \int_{[0,1]^k} \mathbb{E} \left[\int_{\mathbb{R}^{2k}} \left(\cos(\mathbf{s}'_k(\mathbf{X}_k - \mathbf{X}'_k)) \cos(\mathbf{t}'_k(\mathbf{Y}_k - \mathbf{Y}'_k)) \right. \right. \\
&\quad \left. \left. + \cos(\mathbf{s}'_k(\mathbf{X}_k - \mathbf{X}'_k)) \cos(\mathbf{t}'_k(\mathbf{Y}''_k - \mathbf{Y}'''_k)) - 2 \cos(\mathbf{s}'_k(\mathbf{X}_k - \mathbf{X}'_k)) \cos(\mathbf{t}'_k(\mathbf{Y}_k - \mathbf{Y}''_k)) \right) \right. \\
&\quad \left. \prod_{j=1}^k g(s_j) g(t_j) ds_k dt_k \right] d\mathbf{x}_k,
\end{aligned} \tag{3.1}$$

where $\mathbf{X}_k = (X(x_1), \dots, X(x_k))'$, $\mathbf{Y}_k = (Y(x_1), \dots, Y(x_k))'$ and $(\mathbf{X}'_k, \mathbf{Y}'_k)$ is an independent copy of $(\mathbf{X}_k, \mathbf{Y}_k)$ while $\mathbf{Y}''_k, \mathbf{Y}'''_k$ are iid copies of \mathbf{Y}_k independent of everything else. Since

$$\cos u \cos v = 1 - (1 - \cos u) - (1 - \cos v) + (1 - \cos u)(1 - \cos v), \tag{3.2}$$

we have

$$\begin{aligned}
&d(P_{X,Y}, P_X \otimes P_Y) \\
&= \sum_{k=1}^{\infty} \mathbb{P}(N(1) = k) \int_{[0,1]^k} \mathbb{E} \left[\int_{\mathbb{R}^{2k}} \left((1 - \cos(\mathbf{s}'_k(\mathbf{X}_k - \mathbf{X}'_k))) (1 - \cos(\mathbf{t}'_k(\mathbf{Y}_k - \mathbf{Y}'_k))) \right. \right. \\
&\quad \left. \left. + (1 - \cos(\mathbf{s}'_k(\mathbf{X}_k - \mathbf{X}'_k))) (1 - \cos(\mathbf{t}'_k(\mathbf{Y}''_k - \mathbf{Y}'''_k))) \right) \right] d\mathbf{x}_k
\end{aligned}$$

$$-2(1 - \cos(\mathbf{s}'_k(\mathbf{X}_k - \mathbf{X}'_k)))(1 - \cos(\mathbf{t}'_k(\mathbf{Y}_k - \mathbf{Y}''_k))) \prod_{j=1}^k g(s_j)g(t_j) ds_k dt_k \Big] d\mathbf{x}_k.$$

Next we replace the product kernels $\prod_{j=1}^k g(s_j)$ above by other positive measurable functions on \mathbb{R}^k . Inspired by [6] we choose the functions

$$g_k(\mathbf{s}) = c_k |\mathbf{s}|^{-k-\alpha}, \quad \mathbf{s} \in \mathbb{R}^k, \quad \alpha \in (0, 2),$$

where the constant $c_k = c_k(\alpha) > 0$ is such that

$$\int_{\mathbb{R}^k} (1 - \cos(\mathbf{s}'\mathbf{x})) g_k(\mathbf{s}) d\mathbf{s} = |\mathbf{x}|^\alpha, \quad x \in \mathbb{R}^k.$$

The corresponding distance covariance between X and Y becomes:

$$\begin{aligned} & d(P_{X,Y}, P_X \otimes P_Y) \\ &= \sum_{k=1}^{\infty} \mathbb{P}(N(1) = k) \int_{[0,1]^k} \mathbb{E} \left[\int_{\mathbb{R}^{2k}} \left((1 - \cos(\mathbf{s}'_k(\mathbf{X}_k - \mathbf{X}'_k)))(1 - \cos(\mathbf{t}'_k(\mathbf{Y}_k - \mathbf{Y}'_k))) \right. \right. \\ & \quad \left. \left. + (1 - \cos(\mathbf{s}'_k(\mathbf{X}_k - \mathbf{X}'_k)))(1 - \cos(\mathbf{t}'_k(\mathbf{Y}''_k - \mathbf{Y}'''_k))) \right. \right. \\ & \quad \left. \left. - 2(1 - \cos(\mathbf{s}'_k(\mathbf{X}_k - \mathbf{X}'_k)))(1 - \cos(\mathbf{t}'_k(\mathbf{Y}_k - \mathbf{Y}''_k))) \right) \right. \\ & \quad \left. g_k(\mathbf{s}_k)g_k(\mathbf{t}_k) d\mathbf{s}_k d\mathbf{t}_k \right] d\mathbf{x}_k. \end{aligned}$$

By Fubini's theorem and the order statistics property of the Poisson process,

$$\begin{aligned} & d(P_{X,Y}, P_X \otimes P_Y) \\ &= \sum_{k=1}^{\infty} \mathbb{P}(N(1) = k) \int_{[0,1]^k} \left(\mathbb{E}[|\mathbf{X}_k - \mathbf{X}'_k|^\alpha | \mathbf{Y}_k - \mathbf{Y}'_k|^\alpha] + \mathbb{E}[|\mathbf{X}_k - \mathbf{X}'_k|^\alpha] \mathbb{E}[|\mathbf{Y}''_k - \mathbf{Y}'''_k|^\alpha] \right. \\ & \quad \left. - 2\mathbb{E}[|\mathbf{X}_k - \mathbf{X}'_k|^\alpha | \mathbf{Y}_k - \mathbf{Y}''_k|^\alpha] \right) d\mathbf{x}_k \\ &= \mathbb{E}[|\mathbf{X}_N - \mathbf{X}'_N|^\alpha | \mathbf{Y}_N - \mathbf{Y}'_N|^\alpha] + \mathbb{E}[|\mathbf{X}_N - \mathbf{X}'_N|^\alpha | \mathbf{Y}''_N - \mathbf{Y}'''_N|^\alpha] \\ & \quad - 2\mathbb{E}[|\mathbf{X}_N - \mathbf{X}'_N|^\alpha | \mathbf{Y}_N - \mathbf{Y}''_N|^\alpha] \\ &= I_1 + I_2 - 2I_3, \end{aligned}$$

where $\mathbf{X}_N = (X(T_1), \dots, X(T_{N(1)}))'$, $\mathbf{Y}_N = (Y(T_1), \dots, Y(T_{N(1)}))'$, etc. In particular, all the expectations are finite if

$$(3.3) \quad \sup_{0 \leq x \leq 1} \mathbb{E}[|X(x)|^\alpha + |Y(x)|^\alpha + |X(x)Y(x)|^\alpha] < \infty.$$

An empirical version of I_1 is then given by

$$\hat{I}_1 = \frac{1}{l_n} \frac{1}{n^2} \sum_{1 \leq i, j \leq n} \sum_{k=1}^{l_n} |\mathbf{X}_{i, N_k} - \mathbf{X}_{j, N_k}|^\alpha |\mathbf{Y}_{i, N_k} - \mathbf{Y}_{j, N_k}|^\alpha,$$

where $((X_k, Y_k))$ are iid copies of (X, Y) independent of the iid copies (N_i) of the homogeneous Poisson process N . The empirical versions \hat{I}_2, \hat{I}_3 of I_2, I_3 are defined in an analogous way. The integer sequence (l_n) is such that $l_n \rightarrow \infty$.

In view of the strong law of large numbers for U -statistics, for fixed l , as $n \rightarrow \infty$,

$$\frac{1}{l} \sum_{k=1}^l A_{nk} = \frac{1}{l} \frac{1}{n^2} \sum_{1 \leq i, j \leq n} \sum_{k=1}^l |\mathbf{X}_{i, N_k} - \mathbf{X}_{j, N_k}|^\alpha |\mathbf{Y}_{i, N_k} - \mathbf{Y}_{j, N_k}|^\alpha$$

$$\xrightarrow{\text{a.s.}} \frac{1}{l} \sum_{k=1}^l \mathbb{E}[|\mathbf{X}_{N_k} - \mathbf{X}'_{N_k}|^\alpha |\mathbf{Y}_{N_k} - \mathbf{Y}'_{N_k}|^\alpha | N_k] := \frac{1}{l} \sum_{k=1}^l A_k.$$

Therefore, we can choose a sequence $\epsilon_n \downarrow 0$ such that

$$\mathbb{P}\left(\frac{1}{l} \left| \sum_{k=1}^l (A_{nk} - A_k) \right| > \epsilon_n\right) \rightarrow 0$$

and then also choose an integer sequence (r_n) such that $r_n \rightarrow \infty$ and

$$r_n \mathbb{P}\left(\frac{1}{l} \left| \sum_{k=1}^l (A_{nk} - A_k) \right| > \epsilon_n\right) \rightarrow 0.$$

Note that the sequence (r_n) can be chosen to be monotone and such that $r_n - r_{n-1} \in \{0, 1\}$ for each n . Then

$$\mathbb{P}\left(\frac{1}{r_n l} \left| \sum_{s=1}^{r_n} \sum_{k=(s-1)l+1}^{sl} (A_{nk} - A_k) \right| > \epsilon_n\right) \leq \mathbb{P}\left(\frac{1}{l} \sup_{s=1, \dots, r_n} \left| \sum_{k=(s-1)l+1}^{sl} (A_{nk} - A_k) \right| > \epsilon_n\right) \rightarrow 0.$$

This means that

$$\frac{1}{r_n l} \sum_{k=1}^{r_n l} (A_{nk} - A_k) \xrightarrow{\mathbb{P}} 0, \quad n \rightarrow \infty.$$

However, by the strong law of large numbers, as $n \rightarrow \infty$,

$$\frac{1}{r_n l} \sum_{k=1}^{r_n l} A_k \xrightarrow{\text{a.s.}} \mathbb{E}[A_1] = \mathbb{E}[|\mathbf{X}_N - \mathbf{X}'_N|^\alpha |\mathbf{Y}_N - \mathbf{Y}'_N|^\alpha].$$

Hence, for every l there is an (r_n) such that

$$\frac{1}{r_n l} \sum_{k=1}^{r_n l} A_{nk} \xrightarrow{\mathbb{P}} \mathbb{E}[A_1], \quad n \rightarrow \infty.$$

We conclude that

$$\begin{aligned} & \sup_{lr_{n-1} \leq v \leq lr_n} \left| \frac{1}{v} \sum_{k=1}^v A_{nk} - \frac{1}{lr_n} \sum_{k=1}^{lr_n} A_{nk} \right| \\ & \leq \frac{r_n - r_{n-1}}{lr_{n-1}r_n} \sum_{k=1}^{lr_n} A_{nk} + \frac{1}{lr_n} \sum_{k=lr_{n-1}+1}^{lr_n} A_{nk}. \end{aligned}$$

The right-hand side converges in probability to zero, hence we have the law of large numbers for \hat{I}_1 . Similar arguments apply to \hat{I}_2, \hat{I}_3 . We summarize:

Proposition 3.1. *Let $\alpha \in (0, 2)$ and assume that (3.3) holds. Then for any integer sequence (l_n) with $l_n \rightarrow \infty$,*

$$\begin{aligned} d(P_{n,X,Y}, P_{n,X} \otimes P_{n,Y}) &= \frac{1}{l_n} \frac{1}{n^2} \sum_{1 \leq i, j \leq n} \sum_{k=1}^{l_n} |\mathbf{X}_{i, N_k} - \mathbf{X}_{j, N_k}|^\alpha |\mathbf{Y}_{i, N_k} - \mathbf{Y}_{j, N_k}|^\alpha \\ &+ \frac{1}{l_n} \frac{1}{n^2} \sum_{1 \leq i, j \leq n} \sum_{k=1}^{l_n} |\mathbf{X}_{i, N_k} - \mathbf{X}_{j, N_k}|^\alpha \frac{1}{l_n} \frac{1}{n^2} \sum_{1 \leq i, j \leq n} \sum_{k=1}^{l_n} |\mathbf{Y}_{i, N_k} - \mathbf{Y}_{j, N_k}|^\alpha \end{aligned}$$

$$\begin{aligned}
 & -2 \frac{1}{l_n} \frac{1}{n^3} \sum_{1 \leq i, j, l \leq n} \sum_{k=1}^{l_n} |\mathbf{X}_{i, N_k} - \mathbf{X}_{j, N_k}|^\alpha |\mathbf{Y}_{i, N_k} - \mathbf{Y}_{l, N_k}|^\alpha \\
 & \xrightarrow{\mathbb{P}} d(P_{X, Y}, P_X \otimes P_Y).
 \end{aligned}$$

4. A SIMULATION STUDY

In what follows, we conduct a small simulation study for the sample distance correlation $R_n(X, Y)$ from Section 2 for the standard normal density g . This choice implies that

$$\begin{aligned}
 & e^1 d(P_{n, X, Y}, P_{n, X} \otimes P_{n, Y}) \\
 &= \frac{1}{n^2} \sum_{j_1=1}^n \sum_{j_2=1}^n \exp \left(\int_{[0,1]} e^{-|X_{j_1}(x) - X_{j_2}(x)|^2/2 - |Y_{j_1}(x) - Y_{j_2}(x)|^2/2} dx \right) \\
 &+ \frac{1}{n^4} \sum_{j_1=1}^n \sum_{j_2=1}^n \sum_{j_3=1}^n \sum_{j_4=1}^n \exp \left(\int_{[0,1]} e^{-|X_{j_1}(x) - X_{j_2}(x)|^2/2 - |Y_{j_3}(x) - Y_{j_4}(x)|^2/2} dx \right) \\
 &- \frac{1}{n^3} \sum_{j_1=1}^n \sum_{j_2=1}^n \sum_{j_3=1}^n \exp \left(\int_{[0,1]} e^{-|X_{j_1}(x) - X_{j_2}(x)|^2/2 - |Y_{j_1}(x) - Y_{j_3}(x)|^2/2} dx \right).
 \end{aligned}$$

As a matter of fact, simulations of this quantity are highly complex. We choose a moderate sample size $n = 100$ and approximate the integrals on $[0, 1]$ by their Riemann sums at an equidistant grid with mesh $1/50$. For (X, Y) , we take a bivariate Brownian motion (B_1, B_2) with correlation $\rho \in [0, 1]$, i.e.,

$$\text{cov}(B_1(s), B_2(t)) = \rho \min_{s, t \in [0, 1]} \{s, t\}, \quad s, t \in [0, 1],$$

and a bivariate fractional Brownian motion (W_1, W_2) with correlation $\rho \in [0, 1]$, i.e.,

$$\text{cov}(W_1(s), W_2(t)) = \frac{\rho}{2} \{|s|^{2H} + |t|^{2H} - |t - s|^{2H}\}, \quad s, t \in [0, 1],$$

where we assume that the Hurst parameters of W_1 and W_2 are the same; see [4] for more general cross-correlation structures of vector-fractional Brownian motions.

We compare the behavior of the sample distance correlation $R_n(X, Y)$ of the aforementioned stochastic processes with the corresponding sample distance correlation from [6]; we denote it by $R_n^{\text{Sz}}(\mathbf{X}, \mathbf{Y})$. We calculate the sample distance correlation $R_n^{\text{Sz}}(\mathbf{X}, \mathbf{Y})$ based on $n = 100$ iid simulations of the vector $(\mathbf{X}, \mathbf{Y}) = (X(i/50), Y(i/50))_{i=1, \dots, 50}$. The calculation of $R_n(X, Y)$ and $R_n^{\text{Sz}}(\mathbf{X}, \mathbf{Y})$ is based on the same simulated sample paths $((X_i, Y_i))_{i=1, \dots, n}$.

Figures 1–3 are based on 40 independent simulations of $R_n(X, Y)$ and $R_n^{\text{Sz}}(\mathbf{X}, \mathbf{Y})$. The 3 left (right) histograms show $R_n(X, Y)$ ($R_n^{\text{Sz}}(\mathbf{X}, \mathbf{Y})$) for 3 different choices of processes (X, Y) . Although it is difficult to judge from such a small simulation study with rather special stochastic processes, these graphs give one the impression that both sample distance correlations capture the independence or dependence of the processes X and Y quite well. The quantities $R_n^{\text{Sz}}(\mathbf{X}, \mathbf{Y})$ have the tendency to be larger than $R_n(X, Y)$.

Finally, we consider two independent piecewise constant processes X and Y on $[0, 1]$ which assume iid standard normal values on the intervals $((i-1)/50, i/50]$, $i = 1, 2, \dots, 50$. This is essentially the setting of [8] who chose independent vectors of iid normal random variables for the construction of $R_n^{\text{Sz}}(\mathbf{X}, \mathbf{Y})$. In the right histogram of Figure 4 one can see that $R_n^{\text{Sz}}(\mathbf{X}, \mathbf{Y})$ is typically far from zero. This was observed in [8] who studied the case when the dimension of the vectors is large compared to the sample size. On the other hand, our measure $R_n(X, Y)$ is quite in agreement with the independence hypothesis.

Of course, more investigations are needed in order to find out about the strengths and weaknesses of the distance covariances and correlation for processes introduced in this paper. One of the main

problems will be to find reliable confidence bands for the estimator $R_n(X, Y)$. This is work in progress.

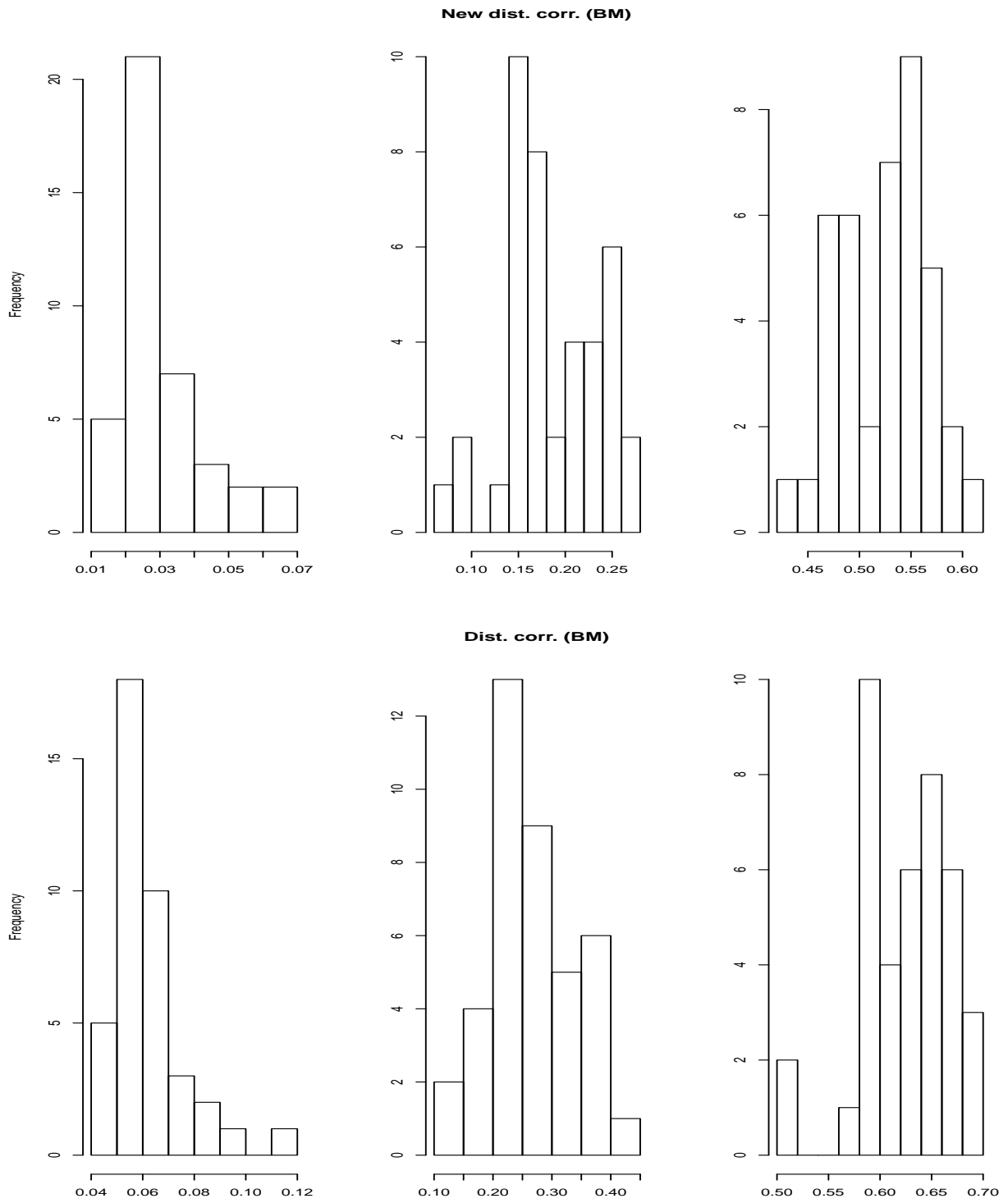


FIGURE 1. Histograms of $R_n(B_1, B_2)$ (top) and $R_n^{Sz}(B_1, B_2)$ (bottom) based on 40 samples. The correlations of B_1 and B_2 are respectively $\rho = 0, 0.5, 0.8$, from left to right.

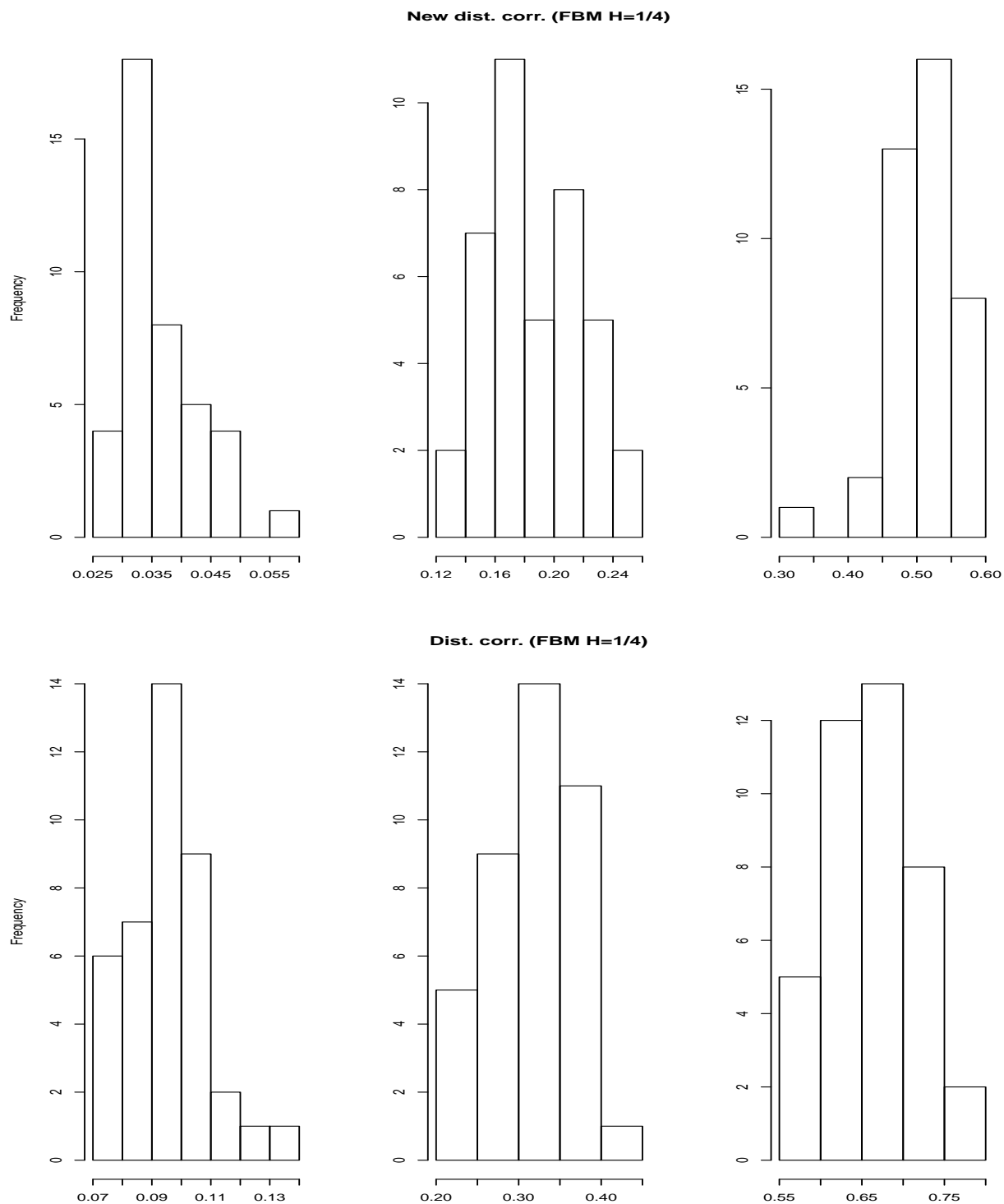


FIGURE 2. Histograms of $R_n(W_1, W_2)$ (top) and $R_n^{Sz}(W_1, W_2)$ (bottom) for $H = 0.25$ based on 40 samples. The correlations of W_1 and W_2 are respectively $\rho = 0, 0.5, 0.8$, from left to right.

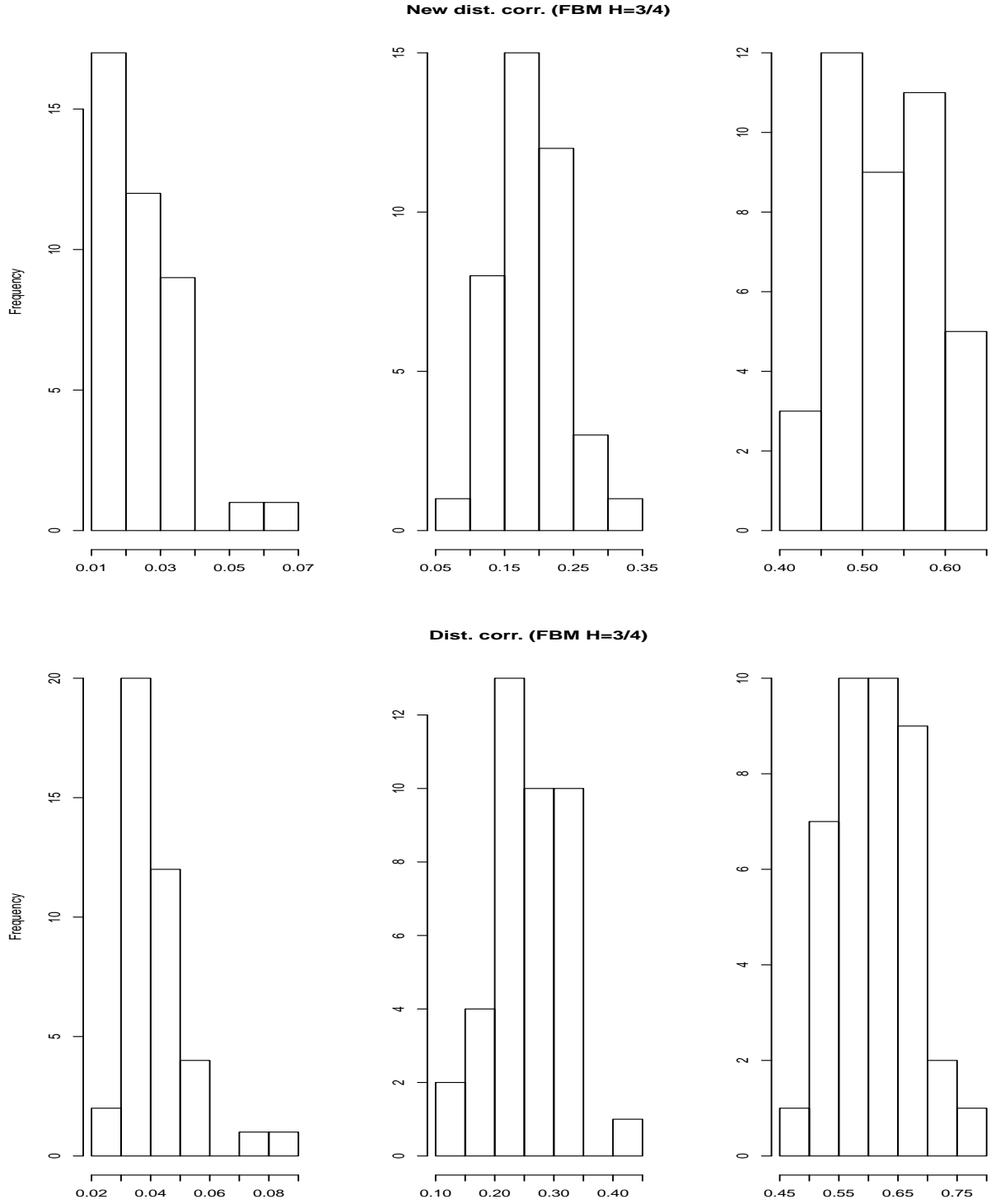


FIGURE 3. Histograms of $R_n(W_1, W_2)$ (top) and $R_n^{Sz}(W_1, W_2)$ (bottom) for $H = 0.75$ based on 40 samples. The correlations of W_1 and W_2 are respectively $\rho = 0, 0.5, 0.8$, from left to right.

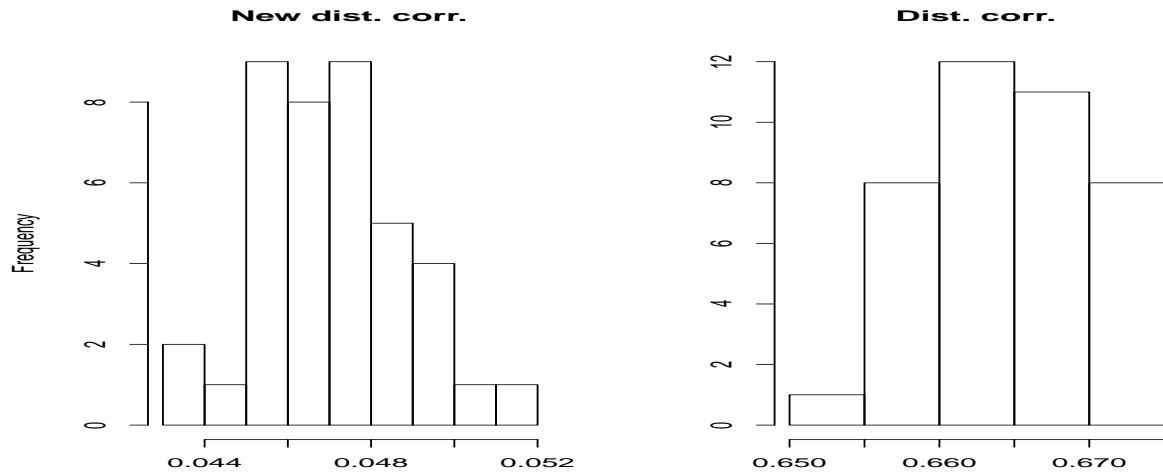


FIGURE 4. Histograms of $R_n(X, Y)$ (left) and $R_n^{Sz}(X, Y)$ (right) based on 40 samples, where X and Y are independent piecewise constant processes based on iid normal random variables.

REFERENCES

- [1] DAVIS, R.A., MATSUI, M., MIKOSCH, T. AND WAN, P. (2016) Applications of distance correlation to time series. Technical report.
- [2] FEUERVERGER, A. (1993) A consistent test for bivariate dependence. *Int. Stat. Rev.* **61**, 419–433.
- [3] HOFFMANN-JØRGENSEN, J. (1994) *Probability with a View Towards Statistics*. Chapman & Hall, New York.
- [4] LAVANCIER, F., PHILIPPE, A. AND SURGAILIS D. (2009) Covariance function of vector self-similar processes. *Statist. Probab. Lett.* **79**, 2415–2421.
- [5] LYONS, R. (2013) Distance covariance in metric spaces. *Ann. Probab.* **41**, 3284–3305.
- [6] SZÉKELY, G.J., RIZZO, M.L. AND BAKIROV, N.K. (2007) Measuring and testing dependence by correlation of distances. *Ann. Statist.* **35**, 2769–2794.
- [7] SZÉKELY, G.J. AND RIZZO, M.L. (2009) Brownian distance covariance. *Ann. Appl. Stat.* **3**, 1236–1265.
- [8] SZÉKELY, G.J. AND RIZZO, M.L. (2013) The distance correlation t -test of independence in high dimension. *J. Multivariate Anal.* **117**, 193–213.
- [9] SZÉKELY, G.J. AND RIZZO, M.L. (2014) Partial distance correlation with methods for dissimilarities. *Ann. Statist.* **42**, 2382–2412.

DEPARTMENT OF BUSINESS ADMINISTRATION, NANZAN UNIVERSITY, 18 YAMAZATO-CHO, SHOWA-KU, NAGOYA 466-8673, JAPAN.

E-mail address: mmuneya@gmail.com

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF COPENHAGEN, UNIVERSITETSPARKEN 5, DK-2100 COPENHAGEN, DENMARK

E-mail address: mikosch@math.ku.dk

SCHOOL OF OPERATIONS RESEARCH AND INFORMATION ENGINEERING, CORNELL UNIVERSITY, 220 RHODES HALL, ITHACA, NY 14853, U.S.A.

E-mail address: gennady@orie.cornell.edu