



# Independence test and canonical correlation analysis based on the alignment between kernel matrices for multivariate functional data

Tomasz Górecki<sup>1</sup> · Mirosław Krzyśko<sup>1,2</sup> · Waldemar Wołyński<sup>1</sup>

Published online: 10 November 2018  
© The Author(s) 2018

## Abstract

In the case of vector data, Gretton et al. (Algorithmic learning theory. Springer, Berlin, pp 63–77, 2005) defined Hilbert–Schmidt independence criterion, and next Cortes et al. (J Mach Learn Res 13:795–828, 2012) introduced concept of the centered kernel target alignment (KTA). In this paper we generalize these measures of dependence to the case of multivariate functional data. In addition, based on these measures between two kernel matrices (we use the Gaussian kernel), we constructed independence test and nonlinear canonical variables for multivariate functional data. We show that it is enough to work only on the coefficients of a series expansion of the underlying processes. In order to provide a comprehensive comparison, we conducted a set of experiments, testing effectiveness on two real examples and artificial data. Our experiments show that using functional variants of the proposed measures, we obtain much better results in recognizing nonlinear dependence.

**Keywords** Multivariate functional data · Functional data analysis · Correlation analysis · Canonical correlation analysis

## 1 Introduction

The theory and practice of statistical methods in situations where the available data are functions (instead of real numbers or vectors) is often referred to as Functional Data Analysis (FDA). The term Functional Data Analysis was already used by Ramsay and Dalzell (1991)

---

✉ Tomasz Górecki  
tomasz.gorecki@amu.edu.pl  
Mirosław Krzyśko  
mkrzysko@amu.edu.pl  
Waldemar Wołyński  
wozynski@amu.edu.pl

<sup>1</sup> Faculty of Mathematics and Computer Science, Adam Mickiewicz University, Umultowska 87, 61-614 Poznań, Poland

<sup>2</sup> Faculty of Management, President Stanisław Wojciechowski Higher Vocational State School, Nowy Świat 4, 62-800 Kalisz, Poland

two decades ago. This subject has become increasingly popular from the end of the 1990s and is now a major research field in statistics (Cuevas 2014). Good access to the large literature in this field comes from the books by Ramsay and Silverman (2002, 2005), Ferraty and Vieu (2006), and Horváth and Kokoszka (2012). Special issues devoted to FDA topics have been published by different journals, including *Statistica Sinica* 14(3) (2004), *Computational Statistics* 22(3) (2007), *Computational Statistics and Data Analysis* 51(10) (2007), *Journal of Multivariate Analysis* 101(2) (2010), *Advances in Data Analysis and Classification* 8(3) (2014).

The range of real world applications, where the objects can be thought of as functions, is as diverse as speech recognition, spectrometry, meteorology, medicine or clients segmentation, to cite just a few (Ferraty and Vieu 2003; James et al. 2009; Martin-Baragan et al. 2014; Devijver 2017).

The uncentered kernel alignment originally was introduced by Cristianini et al. (2001). Gretton et al. (2005) defined Hilbert–Schmidt Independence Criterion (HSIC) and the empirical HSIC. Centered kernel target alignment (KTA) was introduced by Cortes et al. (2012). This measure is a normalized version of HSIC. Zhang et al. (2011) gave an interesting kernel-based independence test. This independence testing method is closely related to the one based on the Hilbert–Schmidt independence criterion (HSIC) proposed by Gretton et al. (2008). Gretton et al. (2005) described a permutation-based kernel independence test. There is a lot of work in the literature for kernel alignment and its applications (good overview can be found in Wang et al. 2015).

This work is devoted to a generalization of these measures of dependence to the case of multivariate functional data. In addition, based on these measures we constructed independence test and nonlinear canonical correlation variables for multivariate functional data. These results are based on the assumption that the applied kernel function is Gaussian. Functional HSIC and KTA canonical correlation analysis can be viewed as natural nonlinear extension of functional canonical correlation analysis (FCCA). So, we propose two nonlinear functional CCA extensions that capture nonlinear relationship. Moreover, both algorithms are capable of extracting also linear dependency. Additionally, we show that functional KTA approach is only a normalized variant of HSIC coefficient also for functional data. Finally, we propose some interpretation of module weighting functions for functional canonical correlations.

Section 2 provides an overview of centered measures alignment for random vectors. They are defined by such concepts as: kernel function alignment, kernel matrix alignment, and Hilbert–Schmidt Independence Criterion (HSIC) and associations between them have been shown. Functional data can be seen as values of random processes. In our paper, the multivariate random function  $\mathbf{X}$  and  $\mathbf{Y}$  have special representation (8) in finite dimensional subspaces of the spaces of square integrable functions on the given intervals. In Sect. 3 we present kernel-based independence test. Section 4 discusses the concept of alignment for multivariate functional data. The kernel function, the alignment between two kernels functions, the centered kernel alignment (KTA) between two kernel matrices and the empirical Hilbert–Schmidt Independence Criterion (HSIC) are defined. The HSIC was used as the basis for an independence test. In Sect. 5 we present kernel-based independence test for multivariate functional data. In Sect. 5, based on the concept of alignment between kernel matrices, nonlinear canonical variables were constructed. It is a generalization of the results of Chang et al. (2013) for random vectors. In Sect. 5 we present an one artificial and two real examples which confirm the usefulness of proposed coefficients in detection of nonlinear dependency for group of variables.

## 2 An overview of kernel alignment and its applications

We introduce the following notational convention. Throughout this section,  $\mathbf{X}$  and  $\mathbf{Y}$  are random vectors, with domains  $\mathbb{R}^p$  and  $\mathbb{R}^q$ , respectively. Let  $P_{\mathbf{X}, \mathbf{Y}}$  be a joint probability measure on  $(\mathbb{R}^p \times \mathbb{R}^q, \Gamma \times \Lambda)$  (here  $\Gamma$  and  $\Lambda$  are the Borel  $\sigma$ -algebras on  $\mathbb{R}^p$  and  $\mathbb{R}^q$ , respectively), with associated marginal probability measures  $P_{\mathbf{X}}$  and  $P_{\mathbf{Y}}$ .

**Definition 1** (*Kernel functions*, Shawe-Taylor and Cristianini 2004) A kernel is a function  $k$  that for all  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^p$  satisfies

$$k(\mathbf{x}, \mathbf{x}') = \langle \boldsymbol{\varphi}(\mathbf{x}), \boldsymbol{\varphi}(\mathbf{x}') \rangle_{\mathcal{H}},$$

where  $\boldsymbol{\varphi}$  is a mapping from  $\mathbb{R}^p$  to an inner product feature space  $\mathcal{H}$

$$\boldsymbol{\varphi} : \mathbf{x} \rightarrow \boldsymbol{\varphi}(\mathbf{x}) \in \mathcal{H}.$$

We call  $\boldsymbol{\varphi}$  a feature map.

A kernel function can be interpreted as a kind of similarity measure between the vectors  $\mathbf{x}$  and  $\mathbf{x}'$ .

**Definition 2** (*Gram matrix*, Mercer 1909; Riesz 1909; Aronszajn 1950) Given a kernel  $k$  and inputs  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$ , the  $n \times n$  matrix  $\mathbf{K}$  with entries  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$  is called the Gram matrix (kernel matrix) of  $k$  with respect to  $\mathbf{x}_1, \dots, \mathbf{x}_n$ .

**Definition 3** (*Positive semi-definite matrix*, Hofmann et al. 2008) A real  $n \times n$  symmetric matrix  $\mathbf{K}$  with entries  $K_{ij}$  satisfying

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j K_{ij} \geq 0$$

for all  $c_i \in \mathbb{R}$  is called positive semi-definite.

**Definition 4** (*Positive semi-definite kernel*, Mercer 1909; Hofmann et al. 2008) A function  $k : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  which for all  $n \in \mathbb{N}$ ,  $\mathbf{x}_i \in \mathbb{R}^p$ ,  $i = 1, \dots, n$  gives rise to a positive semi-definite Gram matrix is called a positive semi-definite kernel.

This raises an interesting question: given a function of two variables  $k(\mathbf{x}, \mathbf{x}')$ , does there exist a function  $\boldsymbol{\varphi}(\mathbf{x})$  such that  $k(\mathbf{x}, \mathbf{x}') = \langle \boldsymbol{\varphi}(\mathbf{x}), \boldsymbol{\varphi}(\mathbf{x}') \rangle_{\mathcal{H}}$ ? The answer is provided by Mercer’s theorem (1909) which says, roughly, that if  $k$  is positive semi-definite then such a  $\boldsymbol{\varphi}$  exists.

Often, we will not know  $\boldsymbol{\varphi}$ , but a kernel function  $k : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  that encodes the inner product in  $\mathcal{H}$ , instead.

Popular positive semi-definite kernel functions on  $\mathbb{R}^p$  include the polynomial kernel of degree  $d > 0$ ,  $k(\mathbf{x}, \mathbf{x}') = (1 + \mathbf{x}^\top \mathbf{x}')^d$ , the Gaussian kernel  $k(\mathbf{x}, \mathbf{x}') = \exp(-\lambda \|\mathbf{x} - \mathbf{x}'\|^2)$ ,  $\lambda > 0$ , and the Laplace kernel  $k(\mathbf{x}, \mathbf{x}') = \exp(-\lambda \|\mathbf{x} - \mathbf{x}'\|)$ ,  $\lambda > 0$ . In this paper we use, the Gaussian kernel.

We start with the definition of centering and the analysis of its relevant properties.

### 2.1 Centered kernel functions

A feature mapping  $\phi: \mathbb{R}^p \rightarrow \mathcal{H}$  is centered by subtracting from it its expectation, that is transforming  $\phi(\mathbf{x})$  to  $\tilde{\phi}(\mathbf{x}) = \phi(\mathbf{x}) - E_{\mathbf{X}}[\phi(\mathbf{X})]$ , where  $E_{\mathbf{X}}$  denotes the expected value of  $\phi(\mathbf{X})$  when  $\mathbf{X}$  is distributed according to  $P_{\mathbf{X}}$ . Centering a positive semi-definite kernel function  $k: \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  consists centering in the feature mapping  $\phi$  associated to  $k$ . Thus, the centered kernel  $\tilde{k}$  associated to  $k$  is defined by

$$\begin{aligned} \tilde{k}(\mathbf{x}, \mathbf{x}') &= (\phi(\mathbf{x}) - E_{\mathbf{X}}[\phi(\mathbf{X})], \phi(\mathbf{x}') - E_{\mathbf{X}'}[\phi(\mathbf{X}')]) \\ &= k(\mathbf{x}, \mathbf{x}') - E_{\mathbf{X}}[k(\mathbf{X}, \mathbf{x}')] - E_{\mathbf{X}'}[k(\mathbf{x}, \mathbf{X}')] + E_{\mathbf{X}, \mathbf{X}'}[k(\mathbf{X}, \mathbf{X}')], \end{aligned}$$

assuming the expectations exist. Here, the expectation is taken over independent copies  $\mathbf{X}, \mathbf{X}'$  distributed according to  $P_{\mathbf{X}}$ . We see that,  $\tilde{k}$  is also a positive semi-definite kernel. Note also that for a centered kernel  $\tilde{k}$ ,  $E_{\mathbf{X}, \mathbf{X}'}[\tilde{k}(\mathbf{X}, \mathbf{X}')] = 0$ , that is, centering the feature mapping implies centering the kernel function.

### 2.2 Centered kernel matrices

Let  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  be a finite subset of  $\mathbb{R}^p$ . A feature mapping  $\phi(\mathbf{x}_i), i = 1, \dots, n$ , is centered by subtracting from it its empirical expectation, i.e., leading to  $\tilde{\phi}(\mathbf{x}_i) = \phi(\mathbf{x}_i) - \bar{\phi}$ , where  $\bar{\phi} = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i)$ . The kernel matrix  $\mathbf{K} = (K_{ij})$  associated to the kernel function  $k$  and the set  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  is centered by replacing it with  $\tilde{\mathbf{K}} = (\tilde{K}_{ij})$  defined for all  $i, j = 1, 2, \dots, n$  by

$$\tilde{K}_{ij} = K_{ij} - \frac{1}{n} \sum_{i=1}^n K_{ij} - \frac{1}{n} \sum_{j=1}^n K_{ij} + \frac{1}{n^2} \sum_{i,j=1}^n K_{ij}, \tag{1}$$

where  $K_{ij} = k(\mathbf{x}_i, \mathbf{x}_j), i, j = 1, \dots, n$ .

The centered kernel matrix  $\tilde{\mathbf{K}}$  is a positive semi-definite matrix. Also, as with the kernel function  $\frac{1}{n^2} \sum_{i,j} \tilde{K}_{ij} = 0$ .

Let  $\langle \cdot, \cdot \rangle_F$  denote the Frobenius product and  $\| \cdot \|_F$  the Frobenius norm defined for all  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$  by

$$\begin{aligned} \langle \mathbf{A}, \mathbf{B} \rangle_F &= \text{tr}(\mathbf{A}^\top \mathbf{B}), \\ \|\mathbf{A}\|_F &= (\langle \mathbf{A}, \mathbf{A} \rangle_F)^{1/2}. \end{aligned}$$

Then, for any kernel matrix  $\mathbf{K} \in \mathbb{R}^{n \times n}$ , the centered kernel matrix  $\tilde{\mathbf{K}}$  can be expressed as follows (Schölkopf et al. 1998):

$$\tilde{\mathbf{K}} = \mathbf{H} \mathbf{K} \mathbf{H}, \tag{2}$$

where  $\mathbf{H} = \mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^\top, \mathbf{1}_n \in \mathbb{R}^{n \times 1}$  denote the vector with all entries equal to one, and  $\mathbf{I}_n$  the identity matrix of order  $n$ . The matrix  $\mathbf{H}$  is called ‘‘centering matrix’’.

Since  $\mathbf{H}$  is the idempotent matrix ( $\mathbf{H}^2 = \mathbf{H}$ ), then we get for any two kernel matrices  $\mathbf{K}$  and  $\mathbf{L}$  based on the subset  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  of  $\mathbb{R}^p$  and the subset  $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$  of  $\mathbb{R}^q$ , respectively,

$$\langle \tilde{\mathbf{K}}, \tilde{\mathbf{L}} \rangle_F = \langle \mathbf{K}, \mathbf{L} \rangle_F = \langle \tilde{\mathbf{K}}, \mathbf{L} \rangle_F. \tag{3}$$

### 2.3 Centered kernel alignment

**Definition 5** (*Kernel function alignment*, Cristianini et al. 2001; Cortes et al. 2012) Let  $k$  and  $l$  be two kernel functions defined over  $\mathbb{R}^p \times \mathbb{R}^p$  and  $\mathbb{R}^q \times \mathbb{R}^q$ , respectively, such that  $0 < E_{\mathbf{X}, \mathbf{X}'}[\tilde{k}^2(\mathbf{X}, \mathbf{X}')] < \infty$  and  $0 < E_{\mathbf{Y}, \mathbf{Y}'}[\tilde{l}^2(\mathbf{Y}, \mathbf{Y}')] < \infty$ , where  $\mathbf{X}, \mathbf{X}'$  and  $\mathbf{Y}, \mathbf{Y}'$  are independent copies distributed according to  $P_{\mathbf{X}}$  and  $P_{\mathbf{Y}}$ , respectively. Then the alignment between  $k$  and  $l$  is defined by

$$\rho(k, l) = \frac{E_{\mathbf{X}, \mathbf{X}', \mathbf{Y}, \mathbf{Y}'}[\tilde{k}(\mathbf{X}, \mathbf{X}')\tilde{l}(\mathbf{Y}, \mathbf{Y}')] }{\sqrt{E_{\mathbf{X}, \mathbf{X}'}[\tilde{k}^2(\mathbf{X}, \mathbf{X}')] E_{\mathbf{Y}, \mathbf{Y}'}[\tilde{l}^2(\mathbf{Y}, \mathbf{Y}')]}}$$

We can define similarly the alignment between two kernel matrices  $\mathbf{K}$  and  $\mathbf{L}$  based on the finite subset  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  and  $\{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ , respectively.

**Definition 6** (*Kernel matrix alignment*, Cortes et al. 2012) Let  $\mathbf{K} \in \mathbb{R}^{n \times n}$  and  $\mathbf{L} \in \mathbb{R}^{n \times n}$  be two kernel matrices such that  $\|\tilde{\mathbf{K}}\|_F \neq 0$  and  $\|\tilde{\mathbf{L}}\|_F \neq 0$ . Then, the centered kernel target alignment (KTA) between  $\mathbf{K}$  and  $\mathbf{L}$  is defined by

$$\hat{\rho}(\mathbf{K}, \mathbf{L}) = \frac{(\tilde{\mathbf{K}}, \tilde{\mathbf{L}})_F}{\|\tilde{\mathbf{K}}\|_F \|\tilde{\mathbf{L}}\|_F}. \tag{4}$$

Here, by the Cauchy–Schwarz inequality,  $\hat{\rho}(\mathbf{K}, \mathbf{L}) \in [-1, 1]$  and in fact  $\hat{\rho}(\mathbf{K}, \mathbf{L}) \in [0, 1]$  when  $\mathbf{K}$  and  $\mathbf{L}$  are the kernel matrices of the positive semi-definite kernel  $\tilde{k}$  and  $\tilde{l}$ .

Gretton et al. (2005) defined Hilbert–Schmidt Independence Criterion (HSIC) as a test statistic to distinguish between null hypothesis  $H_0: P_{\mathbf{X}, \mathbf{Y}} = P_{\mathbf{X}} P_{\mathbf{Y}}$  (equivalently we may write  $\mathbf{X} \perp\!\!\!\perp \mathbf{Y}$ ) and alternative hypothesis  $H_1: P_{\mathbf{X}, \mathbf{Y}} \neq P_{\mathbf{X}} P_{\mathbf{Y}}$ .

**Definition 7** (*Reproducing kernel Hilbert space*, Riesz 1909; Mercer 1909; Aronszajn 1950) Consider a Hilbert space  $\mathcal{H}$  of functions from  $\mathbb{R}^p$  to  $\mathbb{R}$ . Then  $\mathcal{H}$  is a reproducing kernel Hilbert space (RKHS) if for each  $\mathbf{x} \in \mathbb{R}^p$ , the Dirac evaluation operator  $\delta_{\mathbf{x}}: \mathcal{H} \rightarrow \mathbb{R}$ , which maps  $f \in \mathcal{H}$  to  $f(\mathbf{x}) \in \mathbb{R}$ , is a bounded linear functional.

Let  $\varphi: \mathbb{R}^p \rightarrow \mathcal{H}$  be a map such that for all  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^p$  we have  $\langle \varphi(\mathbf{x}), \varphi(\mathbf{x}') \rangle_{\mathcal{H}} = k(\mathbf{x}, \mathbf{x}')$ , where  $k: \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$  is a unique positive semi-definite kernel. We will require in particular that  $\mathcal{H}$  be separable (it must have a complete, countable orthonormal system). We likewise define a second separable RKHS  $\mathcal{G}$ , with kernel  $l(\cdot, \cdot)$  and feature map  $\psi$ , on the separable space  $\mathbb{R}^q$ .

We may now define the mean elements  $\mu_{\mathbf{X}}$  and  $\mu_{\mathbf{Y}}$  with respect to the measures  $P_{\mathbf{X}}$  and  $P_{\mathbf{Y}}$  as those members of  $\mathcal{H}$  and  $\mathcal{G}$ , respectively, for which

$$\begin{aligned} \langle \mu_{\mathbf{X}}, f \rangle_{\mathcal{H}} &= E_{\mathbf{X}}[\langle \varphi(\mathbf{X}), f \rangle_{\mathcal{H}}] = E_{\mathbf{X}}[f(\mathbf{X})], \\ \langle \mu_{\mathbf{Y}}, g \rangle_{\mathcal{G}} &= E_{\mathbf{Y}}[\langle \psi(\mathbf{Y}), g \rangle_{\mathcal{G}}] = E_{\mathbf{Y}}[g(\mathbf{Y})], \end{aligned}$$

for all functions  $f \in \mathcal{H}$ ,  $g \in \mathcal{G}$ , where  $\varphi$  is the feature map from  $\mathbb{R}^p$  to the RKHS  $\mathcal{H}$ , and  $\psi$  maps from  $\mathbb{R}^q$  to  $\mathcal{G}$  and assuming the expectations exist.

Finally,  $\|\mu_{\mathbf{X}}\|_{\mathcal{H}}^2$  can be computed by applying the expectation twice via

$$\|\mu_{\mathbf{X}}\|_{\mathcal{H}}^2 = E_{\mathbf{X}, \mathbf{X}'}[\langle \varphi(\mathbf{X}), \varphi(\mathbf{X}') \rangle_{\mathcal{H}}] = E_{\mathbf{X}, \mathbf{X}'}[k(\mathbf{X}, \mathbf{X}')],$$

assuming the expectations exist. The expectation is taken over independent copies  $\mathbf{X}, \mathbf{X}'$  distributed according to  $P_{\mathbf{X}}$ . The means  $\mu_{\mathbf{X}}, \mu_{\mathbf{Y}}$  exist when positive semi-definite kernels  $k$  and  $l$  are bounded. We are now in a position to define the cross-covariance operator.

**Definition 8** (*Cross-covariance operator*, Gretton et al. 2005) The cross-covariance operator  $\mathbf{C}_{\mathbf{X},\mathbf{Y}}: \mathcal{G} \rightarrow \mathcal{H}$  associated with the joint probability measure  $P_{\mathbf{X},\mathbf{Y}}$  on  $(\mathbb{R}^p \times \mathbb{R}^q, \Gamma \times \Lambda)$  is a linear operator  $\mathbf{C}_{\mathbf{X},\mathbf{Y}}: \mathcal{G} \rightarrow \mathcal{H}$  defined as

$$\mathbf{C}_{\mathbf{X},\mathbf{Y}} = E_{\mathbf{X},\mathbf{Y}}[\phi(\mathbf{X}) \otimes \psi(\mathbf{Y})] - \mu_{\mathbf{X}} \otimes \mu_{\mathbf{Y}},$$

for all  $f \in \mathcal{H}$  and  $g \in \mathcal{G}$ , where the tensor product operator  $f \otimes g: \mathcal{G} \rightarrow \mathcal{H}, f \in \mathcal{H}, g \in \mathcal{G}$ , is defined as

$$(f \otimes g)h = f\langle g, h \rangle_{\mathcal{G}}, \text{ for all } h \in \mathcal{G}.$$

This is a generalization of the cross-covariance matrix between random vectors. Moreover, by the definition of the Hilbert–Schmidt (HS) norm, we can compute the HS norm of  $f \otimes g$  via

$$\|f \otimes g\|_{HS}^2 = \|f\|_{\mathcal{H}}^2 \|g\|_{\mathcal{G}}^2.$$

**Definition 9** (*Hilbert–Schmidt Independence Criterion*, Gretton et al. 2005) Hilbert–Schmidt Independence Criterion (HSIC) is the squared Hilbert–Schmidt norm (or Frobenius norm) of the cross-covariance operator associated with the probability measure  $P_{\mathbf{X},\mathbf{Y}}$  on  $(\mathbb{R}^p \times \mathbb{R}^q, \Gamma \times \Lambda)$ :

$$\text{HSIC}(P_{\mathbf{X},\mathbf{Y}}) = \|\mathbf{C}_{\mathbf{X},\mathbf{Y}}\|_F^2.$$

To compute it we need to express HSIC in terms of kernel functions (Gretton et al. 2005):

$$\begin{aligned} \text{HSIC}(P_{\mathbf{X},\mathbf{Y}}) &= E_{\mathbf{X},\mathbf{X}',\mathbf{Y},\mathbf{Y}'}[k(\mathbf{X}, \mathbf{X}')l(\mathbf{Y}, \mathbf{Y}')] \\ &\quad + E_{\mathbf{X},\mathbf{X}'}[k(\mathbf{X}, \mathbf{X}')]E_{\mathbf{Y},\mathbf{Y}'}[l(\mathbf{Y}, \mathbf{Y}')] \\ &\quad - 2E_{\mathbf{X},\mathbf{Y}}[E_{\mathbf{X}'}[k(\mathbf{X}, \mathbf{X}')]E_{\mathbf{Y}'}[l(\mathbf{Y}, \mathbf{Y}'])]. \end{aligned} \tag{5}$$

Here  $E_{\mathbf{X},\mathbf{X}',\mathbf{Y},\mathbf{Y}'}$  denotes the expectation over independent pairs  $(\mathbf{X}, \mathbf{Y})$  and  $(\mathbf{X}', \mathbf{Y}')$  distributed according to  $P_{\mathbf{X},\mathbf{Y}}$ .

It follows from (5) that the Frobenius norm of  $\mathbf{C}_{\mathbf{X},\mathbf{Y}}$  exists when the various expectations over the kernels are bounded, which is true as long as the kernels  $k$  and  $l$  are bounded.

**Definition 10** (*Empirical HSIC*, Gretton et al. 2005) Let  $S = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_n, \mathbf{y}_n)\} \subseteq \mathbb{R}^p \times \mathbb{R}^q$  be a series of  $n$  independent observations drawn from  $P_{\mathbf{X},\mathbf{Y}}$ . An estimator of HSIC, written  $\text{HSIC}(S)$ , is given by

$$\text{HSIC}(S) = \frac{1}{n^2} \langle \mathbf{K}, \tilde{\mathbf{L}} \rangle_F, \tag{6}$$

where  $\mathbf{K} = (k(\mathbf{x}_i, \mathbf{x}_j)), \mathbf{L} = (l(\mathbf{y}_i, \mathbf{y}_j)) \in \mathbb{R}^{n \times n}$ .

Comparing (4) and (6) and using (3), we see that the centered kernel target alignment (KTA) is simply a normalized version of  $\text{HSIC}(S)$ .

In two seminar papers on Székely et al. (2007) and Székely and Rizzo (2009) introduced the distance covariance (dCov) and distance correlation (dCor) as powerful measures of dependence.

For column vectors  $\mathbf{s} \in \mathbb{R}^p$  and  $\mathbf{t} \in \mathbb{R}^q$ , denote by  $\|\mathbf{s}\|_p$  and  $\|\mathbf{t}\|_q$  the standard Euclidean norms on the corresponding spaces. For jointly distributed random vectors  $\mathbf{X} \in \mathbb{R}^p$  and  $\mathbf{Y} \in \mathbb{R}^q$ , let

$$f_{\mathbf{X},\mathbf{Y}}(\mathbf{s}, \mathbf{t}) = E_{\mathbf{X},\mathbf{Y}}[\exp[i\langle \mathbf{s}, \mathbf{X} \rangle_p + i\langle \mathbf{t}, \mathbf{Y} \rangle_q]],$$

be the joint characteristic function of  $(\mathbf{X}, \mathbf{Y})$ , and let  $f_{\mathbf{X}}(\mathbf{s}) = f_{\mathbf{X}, \mathbf{Y}}(\mathbf{s}, \mathbf{0})$  and  $f_{\mathbf{Y}}(\mathbf{t}) = \varphi_{\mathbf{X}, \mathbf{Y}}(\mathbf{0}, \mathbf{t})$  be the marginal characteristic functions of  $\mathbf{X}$  and  $\mathbf{Y}$ , where  $\mathbf{s} \in \mathbb{R}^p$  and  $\mathbf{t} \in \mathbb{R}^q$ . The distance covariance between  $\mathbf{X}$  and  $\mathbf{Y}$  is the nonnegative number  $\nu(\mathbf{X}, \mathbf{Y})$  defined by

$$\nu^2(\mathbf{X}, \mathbf{Y}) = \frac{1}{C_p C_q} \int_{\mathbb{R}^{p+q}} \frac{|f_{\mathbf{X}, \mathbf{Y}}(\mathbf{s}, \mathbf{t}) - f_{\mathbf{X}}(\mathbf{s})f_{\mathbf{Y}}(\mathbf{t})|^2}{\|\mathbf{s}\|_p^{p+1} \|\mathbf{t}\|_q^{q+1}} d\mathbf{s} d\mathbf{t},$$

and  $|z|$  denotes the modulus of  $z \in \mathbb{C}$  and

$$C_p = \frac{\pi^{\frac{1}{2}(p+1)}}{\Gamma(\frac{1}{2}(p+1))}.$$

The distance correlation between  $\mathbf{X}$  and  $\mathbf{Y}$  is the nonnegative number defined by

$$\mathcal{R}(\mathbf{X}, \mathbf{Y}) = \frac{\nu(\mathbf{X}, \mathbf{Y})}{\sqrt{\nu(\mathbf{X}, \mathbf{X})\nu(\mathbf{Y}, \mathbf{Y})}}$$

if both  $\nu(\mathbf{X}, \mathbf{X})$  and  $\nu(\mathbf{Y}, \mathbf{Y})$  are strictly positive, and defined to be zero otherwise. For distributions with finite first moments, the distance correlation characterizes independence in that  $0 \leq \mathcal{R}(\mathbf{X}, \mathbf{Y}) \leq 1$  with  $\mathcal{R}(\mathbf{X}, \mathbf{Y}) = 0$  if and only if  $\mathbf{X}$  and  $\mathbf{Y}$  are independent.

Sejdinovic et al. (2013) demonstrated that distance covariance is an instance of the Hilbert–Schmidt Independence Criterion. Górecki et al. (2016, 2017) showed an extension of the distance covariance and distance correlation coefficients to the functional case.

### 2.4 Kernel-based independence test

Statistical tests of independence have been associated with a broad variety of dependence measures. Classical tests such as Spearman’s  $\rho$  and Kendall’s  $\tau$  are widely applied, however they are not guaranteed to detect all modes of dependence between the random variables. Contingency table-based methods, and in particular the power-divergence family of test statistics (Read and Cressie 1988) are the best known general purpose tests of independence, but are limited to relatively low dimensions, since they require a partitioning of the space in which random variable resides. Characteristic function-based tests (Feuerverger 1993; Kankainen 1995) have also been proposed. They are more general than kernel-based tests, although to our knowledge they have been used only to compare univariate random variables.

Now, we describe how HSIC can be used as an independence measure, and as the basis for an independence test. We begin by demonstrating that the Hilbert–Schmidt norm can be used as a measure of independence, as long as the associated RKHSs are universal.

A continuous kernel  $k$  on a compact metric space is called universal if the corresponding RKHS  $\mathcal{H}$  is dense in the class of continuous functions of the space.

Denote by  $\mathcal{H}, \mathcal{G}$  RKHSs with universal kernels  $k, l$  on the compact domains  $\mathcal{X}$  and  $\mathcal{Y}$  respectively. We assume without loss of generality that  $\|f\|_\infty \leq 1$  and  $\|g\|_\infty \leq 1$  for all  $f \in \mathcal{H}$  and  $g \in \mathcal{G}$ . Then Gretton et al. (2005) proved that  $\|\mathbf{C}_{\mathbf{X}, \mathbf{Y}}\|_{HS} = 0$  if and only if  $\mathbf{X}$  and  $\mathbf{Y}$  are independent. Examples of universal kernels are Gaussian kernel and Laplacian kernel, while the linear kernel  $k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^\top \mathbf{x}'$  is not universal—the corresponding HSIC tests only linear relationships, and a zero cross-covariance matrix characterizes independence only for multivariate Gaussian distributions. Working with the infinite dimensional operator with universal kernels, allows us to identify any general nonlinear dependence (in the limit) between any pair of vectors, not just Gaussians.

We recall that in this paper we use the Gaussian kernel. We now consider the asymptotic distribution of statistics (6).

We introduce the null hypothesis  $H_0: \mathbf{X} \perp\!\!\!\perp \mathbf{Y}$  ( $\mathbf{X}$  is independent of  $\mathbf{Y}$ , i.e.,  $P_{\mathbf{X},\mathbf{Y}} = P_{\mathbf{X}}P_{\mathbf{Y}}$ ). Suppose that we are given the i.i.d. samples  $S_{\mathbf{X}} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  and  $S_{\mathbf{Y}} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$  for  $\mathbf{X}$  and  $\mathbf{Y}$ , respectively. Let  $\tilde{\mathbf{K}}$  and  $\tilde{\mathbf{L}}$  be the centered kernel matrices associated to the kernel function  $k$  and the sets  $S_{\tilde{\mathbf{X}}}$  and  $S_{\tilde{\mathbf{Y}}}$ , respectively. Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$  be the eigenvalues of the matrix  $\tilde{\mathbf{K}}$  and let  $\mathbf{v}_1, \dots, \mathbf{v}_n$  be a set of orthonormal eigenvectors corresponding to these eigenvalues. Let  $\lambda'_1 \geq \lambda'_2 \geq \dots \geq \lambda'_n \geq 0$  be the eigenvalues of the matrix  $\tilde{\mathbf{L}}$  and let  $\mathbf{v}'_1, \dots, \mathbf{v}'_n$  be a set of orthonormal eigenvectors corresponding to these eigenvalues. Let  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ ,  $\Lambda' = \text{diag}(\lambda'_1, \dots, \lambda'_n)$ ,  $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_n)$  and  $\mathbf{V}' = (\mathbf{v}'_1, \dots, \mathbf{v}'_n)$ . Suppose further that we have the eigenvalue decomposition (EVD) of the centered kernel matrices  $\tilde{\mathbf{K}}$  and  $\tilde{\mathbf{L}}$ , i.e.,  $\tilde{\mathbf{K}} = \mathbf{V}\Lambda\mathbf{V}^\top$  and  $\tilde{\mathbf{L}} = \mathbf{V}'\Lambda'(\mathbf{V}')^\top$ .

Let  $\Psi = (\Psi_1, \dots, \Psi_n) = \mathbf{V}\Lambda^{1/2}$  and  $\Psi' = (\Psi'_1, \dots, \Psi'_n) = \mathbf{V}'(\Lambda')^{1/2}$ , i.e.,  $\Psi_i = \sqrt{\lambda_i}\mathbf{v}_i$ ,  $\Psi'_i = \sqrt{\lambda'_i}\mathbf{v}'_i$ ,  $i = 1, \dots, n$ .

The following result is true (Zhang et al. 2011): under the null hypothesis that  $\mathbf{X}$  and  $\mathbf{Y}$  are independent, the statistic (6) has the same asymptotic distribution as

$$Z_n = \frac{1}{n^2} \sum_{i,j=1}^n \lambda_{i,n} \lambda'_{j,n} Z_{ij}^2, \tag{7}$$

where  $Z_{ij}^2$  are i.i.d.  $\chi_1^2$ -distributed variables,  $n \rightarrow \infty$ .

Note that the data-based test statistic HSIC (or its probabilistic counterpart) is sensible to dependence/independence and therefore can be used as a test statistic. Also important is the knowledge of its asymptotic distribution. These facts inspire the following dependence/independence testing procedure. Given the sample  $S_{\mathbf{X}}$  and  $S_{\mathbf{Y}}$ , one first calculates the centered kernel matrices  $\tilde{\mathbf{K}}$  and  $\tilde{\mathbf{L}}$  and their eigenvalues  $\lambda_i$  and  $\lambda'_i$ , and then evaluates the statistic  $\text{HSIC}(S)$  according to (6). Next, the empirical null distribution of  $Z$  under the null hypothesis can be simulated in the following way: one draws i.i.d. random samples from the  $\chi_1^2$ -distributed variables  $Z_{ij}^2$ , and then generates samples for  $Z$  according to (7). Finally the  $p$  value can be found by locating  $\text{HSIC}(S)$  in the simulated null distribution.

A permutation-based test is described in Gretton et al. (2005). In the first step they propose to calculate the test statistic  $T$  (HSIC or KTA) for the given data. Next, keeping the order of the first sample we randomly permute the second sample a large number of times, and recompute the selected statistic each time. This destroys any dependence between samples simulating a draw from the product of marginals, making the empirical distribution of the permuted statistics behave like the null distribution of the test statistic. For a specified significance level  $\alpha$ , we calculate threshold  $t_\alpha$  in the right tail of the null distribution. We reject  $H_0$  if  $T > t_\alpha$ . This test was proved to be consistent against any fixed alternative. It means that for any fixed significance level  $\alpha$ , the power goes to 1 as the sample size tends to infinity.

### 2.5 Functional data

In recent years methods for representing data by functions or curves have received much attention. Such data are known in the literature as the functional data (Ramsay and Silverman 2005; Horváth and Kokoszka 2012; Hsing and Eubank 2015). Examples of functional data can be found in various application domains, such as medicine, economics, meteorology and many others. Functional data can be seen as the values of random process  $X(t)$ . In practice, the



values of the observed random process  $X(t)$  are always recorded at discrete times  $t_1, \dots, t_J$ , less frequently or more densely spaced in the range of variability of the argument  $t$ . So we have a time series  $\{x(t_1), \dots, y(t_J)\}$ . However, there are many reasons to model these series as elements of functional space., because the functional data has many advantages over other ways of representing the time series.

1. They easily cope with the problem of missing observations, an inevitable problem in many areas of research. Unfortunately, most data analysis methods require complete time series. One solution is to delete a time series that has missing values from the data, but this can lead to , and generally leads to, loss of information. Another option is to use one of many statistical methods to predict the missing values, but then the results will depend on the interpolation method. In contrast to this type of solutions, in the case of functional data, the problem of missing observations is solved by expressing time series in the form of a set of continuous functions.
2. The functional data naturally preserve the structure of observations, i.e. they maintain the time dependence of the observations and take into account the information about each measurement.
3. The moments of observations do not have to be evenly spaced in individual time series.
4. Functional data avoids the curse of dimensionality. When the number of time points is greater than the number of time series considered, most statistical methods will not give satisfactory results due to overparametrization. In the case of functional data, this problem can be avoided because the time series are replaced with a set of continuous functions independent of the number of time points in which observations are measured.

In most of the papers on functional data analysis, objects are characterized by only one feature observed at many time points. In several applications there is a need to use statistical methods for objects characterized by many features observed at many time points (double multivariate data). In this case, such data are transformed into multivariate functional data.

Let us assume that  $\mathbf{X} = (X_1, \dots, X_p)^\top = \{\mathbf{X}(s), s \in I_1\} \in L_2^p(I_1)$  and  $\mathbf{Y} = (Y_1, \dots, Y_q)^\top = \{\mathbf{Y}(t), t \in I_2\} \in L_2^q(I_2)$  are random processes, where  $L_2(I)$  is a space of square integrable functions on the interval  $I$ . We also assume that

$$E(\mathbf{X}(s)) = \mathbf{0}, s \in I_1, \quad E(\mathbf{Y}(t)) = \mathbf{0}, t \in I_2.$$

We will further assume that each component  $X_g$  of the random process  $\mathbf{X}$  and  $Y_h$  of the random process  $\mathbf{Y}$  can be represented by a finite number of orthonormal basis functions  $\{\varphi_e\}$  and  $\{\varphi_f\}$  of space  $L_2(I_1)$  and  $L_2(I_2)$ , respectively:

$$X_g(s) = \sum_{e=0}^{E_g} \alpha_{ge} \varphi_e(s), s \in I_1, g = 1, 2, \dots, p,$$

$$Y_h(t) = \sum_{f=0}^{F_h} \beta_{hf} \varphi_f(t), t \in I_2, h = 1, 2, \dots, q,$$

where  $\alpha_{ge}$  and  $\beta_{hf}$  are the random coefficients. The degree of smoothness of processes  $X_g$  and  $Y_h$  depends on the values  $E_g$  and  $F_h$  respectively (small values imply more smoothing). The optimum values for  $E_g$  and  $F_h$  are selected using Bayesian Information Criterion (BIC) (see Górecki et al. 2018). As basis functions we can use e.g. the Fourier basis system or spline functions.

We introduce the following notation:

$$\boldsymbol{\alpha} = (\alpha_{10}, \dots, \alpha_{1E_1}, \dots, \alpha_{p0}, \dots, \alpha_{pE_p})^\top,$$

$$\begin{aligned} \boldsymbol{\beta} &= (\beta_{10}, \dots, \beta_{1F_1}, \dots, \beta_{q0}, \dots, \beta_{qF_q})^\top, \\ \boldsymbol{\varphi}_{E_g}(s) &= (\varphi_0(s), \dots, \varphi_{E_g}(s))^\top, s \in I_1, g = 1, 2, \dots, p, \\ \boldsymbol{\varphi}_{F_h}(t) &= (\varphi_0(t), \dots, \varphi_{F_h}(t))^\top, \\ t \in I_2, h = 1, 2, \dots, q, \boldsymbol{\Phi}_1(s) &= \begin{bmatrix} \boldsymbol{\varphi}_{E_1}^\top(s) & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\varphi}_{E_2}^\top(s) & \dots & \mathbf{0} \\ \dots & \dots & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \dots & \boldsymbol{\varphi}_{E_p}^\top(s) \end{bmatrix}, \\ \boldsymbol{\Phi}_2(t) &= \begin{bmatrix} \boldsymbol{\varphi}_{F_1}^\top(t) & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\varphi}_{F_2}^\top(t) & \dots & \mathbf{0} \\ \dots & \dots & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \dots & \boldsymbol{\varphi}_{F_q}^\top(t) \end{bmatrix}, \end{aligned}$$

where  $\boldsymbol{\alpha} \in \mathbb{R}^{K_1+p}$ ,  $\boldsymbol{\beta} \in \mathbb{R}^{K_2+q}$ ,  $\boldsymbol{\Phi}_1 \in \mathbb{R}^{p+(K_1+p)}$ ,  $\boldsymbol{\Phi}_2 \in \mathbb{R}^{q+(K_2+q)}$ ,  $K_1 = E_1 + \dots + E_p$ ,  $K_2 = F_1 + \dots + F_p$ .

Using the above matrix notation the random processes  $\mathbf{X}$  and  $\mathbf{Y}$  can be represented as:

$$\mathbf{X}(s) = \boldsymbol{\Phi}_1(s)\boldsymbol{\alpha}, s \in I_1, \mathbf{Y}(t) = \boldsymbol{\Phi}_2(t)\boldsymbol{\beta}, t \in I_2, \tag{8}$$

where  $E(\boldsymbol{\alpha}) = \mathbf{0}$ ,  $E(\boldsymbol{\beta}) = \mathbf{0}$ .

This means that the values of random processes  $\mathbf{X}$  and  $\mathbf{Y}$  are in finite dimensional subspaces of  $L_2^p(I_1)$  and  $L_2^q(I_2)$ , respectively. We will denote these subspaces by  $\mathcal{L}_2^p(I_1)$  and  $\mathcal{L}_2^q(I_2)$ .

Typically data are recorded at discrete moments in time. The process of transformation of discrete data to functional data is performed for each realization and each variable separately. Let  $x_{gj}$  denote an observed value of the feature  $X_g$ ,  $g = 1, 2, \dots, p$  at the  $j$ th time point  $s_j$ , where  $j = 1, 2, \dots, J$ . Similarly, let  $y_{hj}$  denote an observed value of feature  $Y_h$ ,  $h = 1, 2, \dots, q$  at the  $j$ th time point  $t_j$ , where  $j = 1, 2, \dots, J$ . Then our data consist of  $pJ$  pairs of  $(s_j, x_{gj})$  and of  $qJ$  pairs of  $(t_j, y_{hj})$ . Let  $\mathbf{X}_1, \dots, \mathbf{X}_n$  and  $\mathbf{Y}_1, \dots, \mathbf{Y}_n$  be independent trajectories of random processes  $\mathbf{X}$  and  $\mathbf{Y}$  having the representation (8).

The coefficients  $\boldsymbol{\alpha}_i$  and  $\boldsymbol{\beta}_i$  are estimated by the least squares method. Let us denote these estimates by  $\mathbf{a}_i$  and  $\mathbf{b}_i$ ,  $i = 1, 2, \dots, n$ .

As a result, we obtain functional data of the form:

$$\mathbf{X}_i(s) = \boldsymbol{\Phi}_1(s)\mathbf{a}_i, \mathbf{Y}_i(t) = \boldsymbol{\Phi}_2(t)\mathbf{b}_i, \tag{9}$$

where  $s \in I_1, t \in I_2, \mathbf{a}_i \in \mathbb{R}^{K_1+p}, \mathbf{b}_i \in \mathbb{R}^{K_2+q}$ ,  $K_1 = E_1 + \dots + E_p, K_2 = F_1 + \dots + F_q$ , and  $i = 1, 2, \dots, n$ .

Górecki and Smaga (2017) described a multivariate analysis of variance (MANOVA) for functional data. In the paper by Górecki et al. (2018), three basic methods of dimension reduction for multidimensional functional data are given: principal component analysis, canonical correlation analysis, and discriminant coordinates.

### 3 Alignment for multivariate functional data

#### 3.1 The alignment between two kernel functions and two kernel matrices for multivariate functional data

Let  $\mathbf{x}(s) \in \mathcal{L}_2^p(I_1)$ ,  $s \in I_1$ , where  $\mathcal{L}_2^p(I_1)$  is a finite-dimensional space of continuous square-integrable vector functions over interval  $I_1$ .

Let

$$k^* : \mathcal{L}_2^p(I_1) \times \mathcal{L}_2^p(I_1) \rightarrow \mathbb{R}$$

be a kernel function on  $\mathcal{L}_2^p(I_1)$ . As already mentioned, in this paper we use the Gaussian kernel. For the multivariate functional data this kernel has the form:

$$k^*(\mathbf{x}(s), \mathbf{x}'(s)) = \exp(-\lambda_1 \|\mathbf{x}(s) - \mathbf{x}'(s)\|^2), \lambda_1 > 0.$$

But from (9), and by the orthonormality of the basis functions, we have:

$$\begin{aligned} \|\mathbf{x}(s) - \mathbf{x}'(s)\|^2 &= \int_{I_1} (\mathbf{x}(s) - \mathbf{x}'(s))^\top (\mathbf{x}(s) - \mathbf{x}'(s)) ds \\ &= \|\mathbf{a} - \mathbf{a}'\|^2. \end{aligned}$$

Hence

$$k^*(\mathbf{x}(s), \mathbf{x}'(s)) = k(\mathbf{a}, \mathbf{a}')$$

and

$$k^*(\mathbf{y}(t), \mathbf{y}'(t)) = k(\mathbf{b}, \mathbf{b}'),$$

where  $\mathbf{a}$  and  $\mathbf{b}$  are vectors occurring in the representation (9) of vector functions  $\mathbf{x}(s), s \in I_1, \mathbf{y}(t), t \in I_2$ .

For a given subset  $\{\mathbf{x}_1(s), \dots, \mathbf{x}_n(s)\}$  of  $\mathcal{L}_2^p(I_1)$  and the given kernel function  $k^*$  on  $\mathcal{L}_2^p(I_1) \times \mathcal{L}_2^p(I_1)$ , the matrix  $\mathbf{K}^*$  of size  $n \times n$ , which has its  $(i, j)$ th element  $K_{ij}^*(s)$ , given by  $K_{ij}^*(s) = k^*(\mathbf{x}_i(s), \mathbf{x}_j(s)), s \in I_1$ , is called the kernel matrix of the kernel function  $k^*$  with respect to the set  $\{\mathbf{x}_1(s), \dots, \mathbf{x}_n(s)\}, s \in I_1$ .

**Definition 11** (*Kernel function alignment for functional data*) Let  $\tilde{k}^*$  and  $\tilde{l}^*$  be two kernel functions defined over  $\mathcal{L}_2^p(I_1) \times \mathcal{L}_2^p(I_1)$  and  $\mathcal{L}_2^q(I_2) \times \mathcal{L}_2^q(I_2)$ , respectively, such that  $0 < E_{\mathbf{X}, \mathbf{X}'}[\tilde{k}^{*2}(\mathbf{X}, \mathbf{X}')] < \infty$  and  $0 < E_{\mathbf{Y}, \mathbf{Y}'}[\tilde{l}^{*2}(\mathbf{Y}, \mathbf{Y}')] < \infty$ , where  $\mathbf{X}, \mathbf{X}'$  and  $\mathbf{Y}, \mathbf{Y}'$  are independent copies distributed according to  $P_{\mathbf{X}}$  and  $P_{\mathbf{Y}}$ , respectively. Then the alignment between  $\tilde{k}^*$  and  $\tilde{l}^*$  is defined by

$$\rho(\tilde{k}^*, \tilde{l}^*) = \frac{E_{\mathbf{X}, \mathbf{Y}}[\tilde{k}^*(\mathbf{X}, \mathbf{X}')\tilde{l}^*(\mathbf{Y}, \mathbf{Y}')] }{\sqrt{E_{\mathbf{X}}[\tilde{k}^{*2}(\mathbf{X}, \mathbf{X}')] E_{\mathbf{Y}}[\tilde{l}^{*2}(\mathbf{Y}, \mathbf{Y}')]}}. \tag{10}$$

We can define similarly the alignment between two kernel matrices  $\tilde{\mathbf{K}}^*$  and  $\tilde{\mathbf{L}}^*$  based on the subset  $\{\mathbf{x}_1(s), \dots, \mathbf{x}_n(s)\}, s \in I_1$ , and  $\{\mathbf{y}_1(t), \dots, \mathbf{y}_n(t)\}, t \in I_2$ , of  $\mathcal{L}_2^p(I_1)$  and  $\mathcal{L}_2^q(I_2)$ , respectively.

**Definition 12** (*Kernel matrix alignment for functional data*) Let  $\tilde{\mathbf{K}}^* \in \mathbb{R}^{n \times n}$  and  $\tilde{\mathbf{L}}^* \in \mathbb{R}^{n \times n}$  be two kernel matrices such that  $\|\tilde{\mathbf{K}}^*\|_F \neq 0$  and  $\|\tilde{\mathbf{L}}^*\|_F \neq 0$ . Then, the centered kernel target alignment (KTA) between  $\tilde{\mathbf{K}}^*$  and  $\tilde{\mathbf{L}}^*$  is defined:

$$\hat{\rho}(\mathbf{K}^*, \mathbf{L}^*) = \frac{\langle \tilde{\mathbf{K}}^*, \tilde{\mathbf{L}}^* \rangle_F}{\|\tilde{\mathbf{K}}^*\|_F \|\tilde{\mathbf{L}}^*\|_F}. \tag{11}$$

If  $\tilde{\mathbf{K}}^*$  and  $\tilde{\mathbf{L}}^*$  are positive semi-definite matrices, then  $\hat{\rho}(\mathbf{K}^*, \mathbf{L}^*) \in [0, 1]$ . We have

$$\hat{\rho}(\mathbf{K}^*, \mathbf{L}^*) = \hat{\rho}(\mathbf{K}, \mathbf{L}),$$

where  $\mathbf{K}$  is the matrix of size  $n \times n$ , which has its  $(i, j)$ th element  $K_{ij}$ , given by  $K_{ij} = k(\mathbf{a}_i, \mathbf{a}_j)$ .

### 3.2 Kernel-based independence test for multivariate functional data

**Definition 13** (*Empirical HSIC for functional data*) The empirical HSIC for functional data is defined as

$$\text{HSIC}(S^*) = \frac{1}{n^2} \langle \mathbf{K}^*, \mathbf{L}^* \rangle_F,$$

where  $S^* = \{(\mathbf{x}_1(s), \mathbf{y}_1(t)), \dots, (\mathbf{x}_n(s), \mathbf{y}_n(t))\}$ ,  $s \in I_1$ ,  $t \in I_2$ ,  $\mathbf{K}^*$  and  $\mathbf{L}^*$  are kernel matrices based on the subsets  $\{\mathbf{x}_1(s), \dots, \mathbf{x}_n(s)\}$ ,  $s \in I_1$ , and  $\{\mathbf{y}_1(t), \dots, \mathbf{y}_n(t)\}$ ,  $t \in I_2$  of  $\mathcal{L}_2^p(I_1)$  and  $\mathcal{L}_2^q(I_2)$ , respectively.

But  $\mathbf{K}^* = \mathbf{K}$ , where  $\mathbf{K}$  is the kernel matrix of size  $n \times n$ , which has its  $(i, j)$ th element  $K_{ij}$  given by  $K_{ij} = k(\mathbf{a}_i, \mathbf{a}_j)$ , where  $\mathbf{a}_1, \dots, \mathbf{a}_n$  are vectors occurring in the representation (9) vector functions  $\mathbf{X}(s)$ ,  $s \in I_1$ . Analogously,  $\mathbf{L}^* = \mathbf{L}$ , where  $\mathbf{L}$  is the kernel matrix of size  $n \times n$ , which has its  $(i, j)$ th element  $L_{ij}$  given by  $L_{ij} = l(\mathbf{b}_i, \mathbf{b}_j)$ , where  $\mathbf{b}_1, \dots, \mathbf{b}_n$  are vectors occurring in the representation (9) vector functions  $\mathbf{Y}(t)$ ,  $t \in I_2$ . Hence

$$\text{HSIC}(S^*) = \text{HSIC}(S_v),$$

where  $S_v = \{(\mathbf{a}_1, \mathbf{b}_1), \dots, (\mathbf{a}_n, \mathbf{b}_n)\}$ .

Note also that the null hypothesis  $H_0: \mathbf{X} \perp \mathbf{Y}$  of independence of the random processes  $\mathbf{X}$  and  $\mathbf{Y}$  is equivalent to the null hypothesis  $H_0: \boldsymbol{\alpha} \perp \boldsymbol{\beta}$  of independence of random vectors  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta}$  occurring in the representation (8) random processes  $\mathbf{X}$  and  $\mathbf{Y}$ . We can therefore use the tests described in Section 2.4, replacing  $\mathbf{x}$  and  $\mathbf{y}$  by  $\mathbf{a}$  and  $\mathbf{b}$ .

### 3.3 Canonical correlation analysis based on the alignment between kernel matrices for multivariate functional data

In classical canonical correlation analysis (Hotelling 1936), we are interested in the relationship between two random vectors  $\mathbf{X}$  and  $\mathbf{Y}$ . In the functional case we are interested in the relationship between two random functions  $\mathbf{X}$  and  $\mathbf{Y}$ . Functional canonical variables  $U$  and  $V$  for random processes  $\mathbf{X}$  and  $\mathbf{Y}$  are defined as follows

$$U = \langle \mathbf{u}, \mathbf{X} \rangle = \int_{I_1} \mathbf{u}^\top(s) \mathbf{X}(s) ds, \tag{12}$$

$$V = \langle \mathbf{v}, \mathbf{Y} \rangle = \int_{I_2} \mathbf{v}^\top(t) \mathbf{Y}(t) dt, \tag{13}$$

where the vector functions  $\mathbf{u}$  and  $\mathbf{v}$  are called the vector weight functions and are of the form

$$\mathbf{u}(s) = \Phi_1(s) \mathbf{u}, \quad \mathbf{v}(t) = \Phi_2(t) \mathbf{v} \quad \text{where } \mathbf{u} \in \mathbb{R}^{K_1+p}, \quad \mathbf{v} \in \mathbb{R}^{K_2+q}. \tag{14}$$

Classically the weight functions  $\mathbf{u}$  and  $\mathbf{v}$  are chosen to maximize the sample correlation coefficient (Górecki et al. 2018):

$$\rho = \frac{\text{Cov}(U, V)}{\sqrt{\text{Var}(U) \text{Var}(V)}}. \tag{15}$$

The sample correlation coefficient between the variables  $U$  and  $V$  is now replaced by a centered kernel target alignment (KTA) between kernel matrices  $\mathbf{K}$  and  $\mathbf{L}$  based on the projected data  $\langle \mathbf{u}(s), \mathbf{x}_i(s) \rangle_{\mathcal{H}}$  and  $\langle \mathbf{v}(t), \mathbf{y}_i(t) \rangle_{\mathcal{H}}$ , i.e. their  $(i, j)$ th entry are

$$K_{i,j} = k(\langle \mathbf{u}(s), \mathbf{x}_i(s) \rangle_{\mathcal{H}}, \langle \mathbf{u}(s), \mathbf{x}_j(s) \rangle_{\mathcal{H}}), s \in I_1,$$

and

$$L_{i,j} = l(\langle \mathbf{v}(t), \mathbf{y}_i(t) \rangle_{\mathcal{H}}, \langle \mathbf{v}(t), \mathbf{y}_j(t) \rangle_{\mathcal{H}}), t \in I_2,$$

respectively,  $i, j = 1, \dots, n$ :

$$\hat{\rho}(\mathbf{u}(s), \mathbf{v}(t)) = \frac{\text{tr}(\tilde{\mathbf{K}}^\top \tilde{\mathbf{L}})}{\sqrt{\text{tr}(\tilde{\mathbf{K}}^\top \tilde{\mathbf{K}}) \text{tr}(\tilde{\mathbf{L}}^\top \tilde{\mathbf{L}})}} \tag{16}$$

subject to

$$\|\mathbf{u}(s)\| = \|\mathbf{v}(t)\| = 1. \tag{17}$$

But

$$\begin{aligned} \mathbf{K}_{i,j} &= k\left(\int_{I_1} \mathbf{u}^\top(s) \mathbf{x}_i(s) ds, \int_{I_1} \mathbf{u}^\top(s) \mathbf{x}_j(s) ds\right) \\ &= k(\mathbf{u}^\top \mathbf{a}_i, \mathbf{u}^\top \mathbf{a}_j) = \mathbf{K}_{i,j}^{(\mathbf{u})} \end{aligned}$$

and

$$\begin{aligned} \mathbf{L}_{i,j} &= l\left(\int_{I_2} \mathbf{v}^\top(t) \mathbf{y}_i(t) dt, \int_{I_2} \mathbf{v}^\top(t) \mathbf{y}_j(t) dt\right) \\ &= l(\mathbf{v}^\top \mathbf{b}_i, \mathbf{v}^\top \mathbf{b}_j) = \mathbf{L}_{i,j}^{(\mathbf{v})}, \end{aligned}$$

where  $\mathbf{a}_i$  and  $\mathbf{b}_i$  are vectors occurring in the representation (9) vectors functions  $\mathbf{x}(s), s \in I_1, \mathbf{y}(t), t \in I_2, i = 1, \dots, n, \mathbf{u} \in \mathbb{R}^{K_1+p}, \mathbf{v} \in \mathbb{R}^{K_2+q}$ .

Thus, the choice of weighting functions  $\mathbf{u}(s)$  and  $\mathbf{v}(t)$  so that the coefficient (16) has a maximum value subject to (17) is equivalent to the choice of vectors  $\mathbf{u} \in \mathbb{R}^{K_1+p}$  and  $\mathbf{v} \in \mathbb{R}^{K_2+q}$  such that the coefficient

$$\hat{\rho}(\mathbf{u}, \mathbf{v}) = \frac{\text{tr}(\tilde{\mathbf{K}}_{\mathbf{u}}^\top \tilde{\mathbf{L}}_{\mathbf{v}})}{\sqrt{\text{tr}(\tilde{\mathbf{K}}_{\mathbf{u}}^\top \tilde{\mathbf{K}}_{\mathbf{u}}) \text{tr}(\tilde{\mathbf{L}}_{\mathbf{v}}^\top \tilde{\mathbf{L}}_{\mathbf{v}})}} \tag{18}$$

has a maximum value subject to

$$\|\mathbf{u}\| = \|\mathbf{v}\| = 1, \tag{19}$$

where  $\mathbf{K}_{\mathbf{u}} = (K_{i,j}^{(\mathbf{u})}), \mathbf{L}_{\mathbf{v}} = (L_{i,j}^{(\mathbf{v})}), i, j = 1, \dots, n$ .

In order to maximize the coefficient of (18) we can use the result of Chang et al. (2013). Authors used a gradient descent algorithm, with modified gradient to ensure the unit length constraint is satisfied at each step (Edelman et al. 1998). Optimal step-sizes were found numerically using the Nelder-Mead method. This article employs the Gaussian kernel exclusively while other kernels are available. The bandwidth parameter  $\lambda$  of the Gaussian kernel was chosen using the ‘‘median trick’’ (Song et al. 2010), i.e. the median Euclidean distance between all pairs of points.

The coefficients of the projection of the  $i$ th value  $\mathbf{x}_i(t)$  of random process  $\mathbf{X}$  on the  $k$ th functional canonical variable are equal to

$$\hat{U}_{ik} = \langle \mathbf{u}_k, \mathbf{x}_i \rangle = \int_{I_1} \mathbf{u}_k^\top(s) \mathbf{x}_i(s) ds = \mathbf{a}_i^\top \mathbf{u}_k,$$

analogously the coefficients of the projection of the  $i$ th value  $y_i(t)$  of random process  $\mathbf{Y}_t$  on the  $k$ th functional canonical variable are equal to

$$\hat{V}_{ik} = \mathbf{b}_i^\top \mathbf{v}_k,$$

where  $i = 1, \dots, n, k = 1, \dots, \min(\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B}))$ , where  $\mathbf{A} \in \mathbb{R}^{n \times p}$  and  $\mathbf{B} \in \mathbb{R}^{n \times q}$ , where the  $i$ th rows are  $\mathbf{a}_i$  and  $\mathbf{b}_i$ , respectively, which have column means of zero.

As we mentioned earlier KTA is a normalized variant of HSIC. Hence, we can repeat the above reasoning for HSIC criterion. However, we should remember that both approaches are not equivalent and we can obtain different results.

## 4 Experiments

Let us recall some and introduce another symbols:

- KTA—centered kernel target alignment,
- HSIC—Hilbert–Schmidt Independence Criterion,
- FCCA—classical functional canonical correlation analysis (Ramsay and Silverman 2005; Horváth and Kokoszka 2012),
- HSIC.FCCA—functional canonical correlation analysis based on HSIC,
- HSIC.KTA—functional canonical correlation analysis based on KTA.

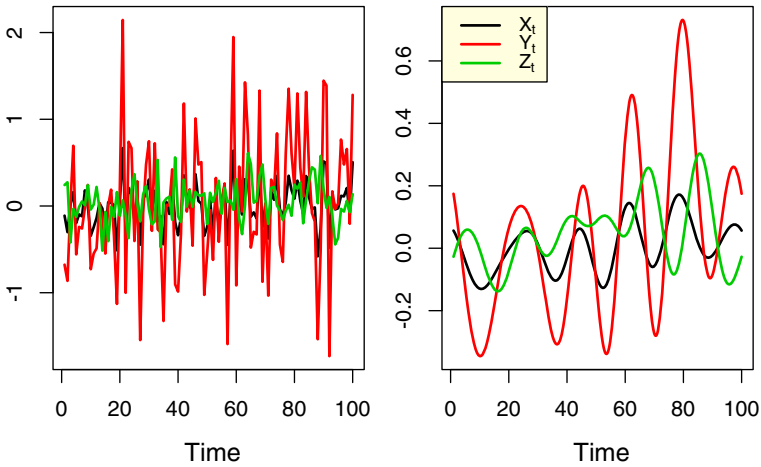
### 4.1 Simulation

We generated random processes along with some noises to test the performance of the introduced measures. Random processes are specified by

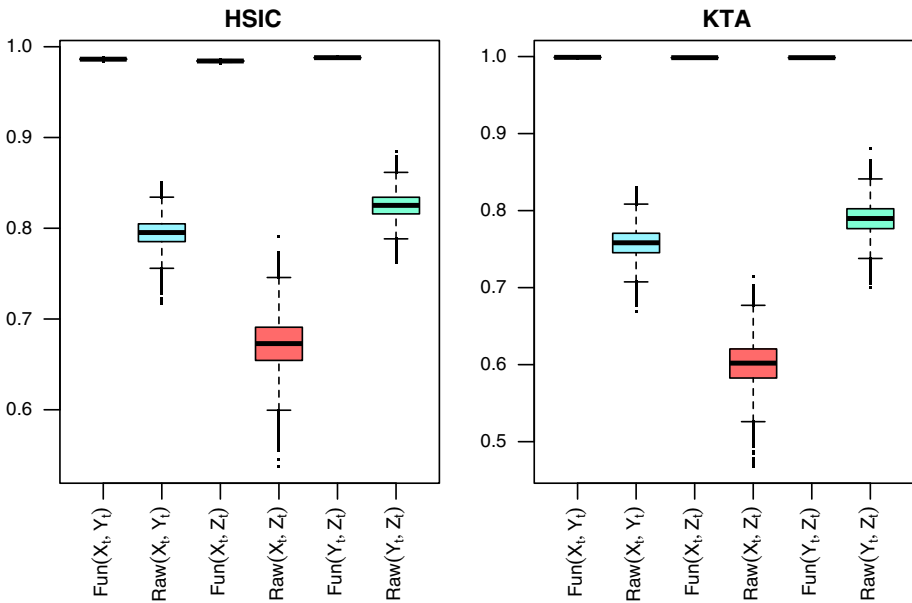
$$\begin{aligned} X_t &= \varepsilon_t, \\ Y_t &= 3X_t + \eta_t, \\ Z_t &= X_t^2 + \xi_t, \end{aligned}$$

where  $\varepsilon_t$ ,  $\eta_t$  and  $\xi_t$  are jointly independent random variables from Gaussian distribution with 0 mean and 0.25 variance. We generated processes of length 100.  $N = 10000$  samples are generated for all processes. The objective is to examine how well functional variants (Fourier basis with 15 basis functions) of KTA and HSIC measures perform compared to measures used on raw data (artificially generated at discrete time stamps). Here, raw data are represented as vector data of generated trajectories, so raw data are three 10000 by 100 dimensional matrices (one for  $X_t$ , second for  $Y_t$  and third for  $Z_t$ ). On the other hand, functional data are three 10000 by 15 dimensional matrices (coefficients of Fourier basis). Here,  $X_t$  and  $Y_t$  are linearly dependent, whereas  $X_t$ ,  $Z_t$  and  $Y_t$ ,  $Z_t$  are nonlinearly dependent (Fig. 1).

From Fig. 2 and Table 1, we see that the proposed extension of HSIC and KTA coefficients to functional data gives larger values of coefficients than the variants for raw time series. Unfortunately, it is not possible to perform inference based only on the values of coefficients. We have to apply tests. In Fig. 3 and in Table 2, we observe that when we use functional variants of the proposed measures, we obtain much better results in recognizing nonlinear dependence. Linear dependence between  $X_t$  and  $Y_t$  was easily recognized by each method (100% of correct decisions— $p$  values below 5%). Results of functional KTA and HSIC are very similar. Non-functional measures HSIC and KTA give only 7.2% and 6.7% correct decisions ( $p$  values below 5%) for relationship  $X_t$ ,  $Z_t$  and  $Y_t$ ,  $Z_t$ , respectively. On the



**Fig. 1** Sample trajectories of  $X_t$ ,  $Y_t$  and  $Z_t$  time series for raw (left plot) and functional (right plot) representation



**Fig. 2** Raw and functional HSIC and KTA coefficients for artificial time series

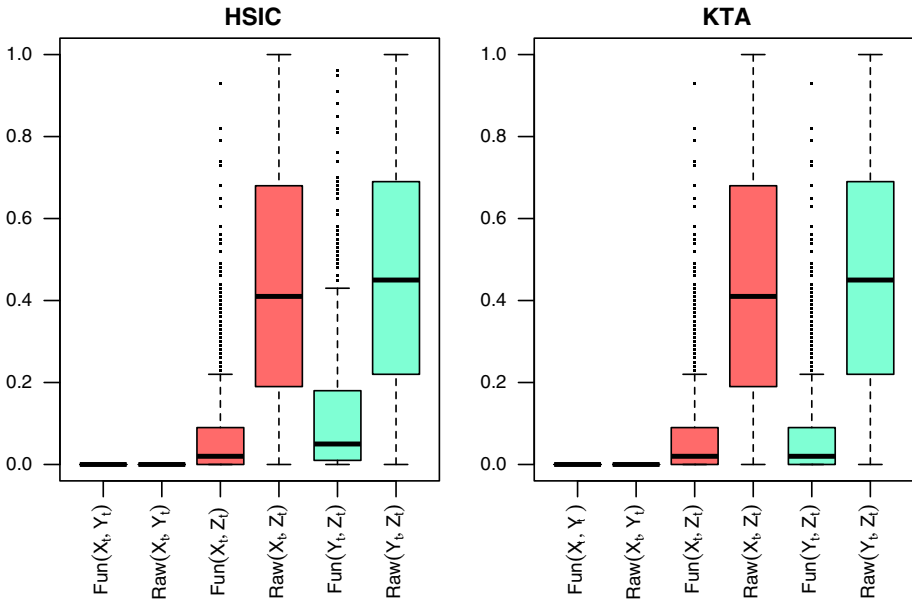
other hand, functional variants recognize dependency ( $p$  values below 5%) in 63.3% (both measures) for  $X_t$ ,  $Y_t$  and in 47.8% (HSIC), 63.3% (KTA) for  $Y_t$ ,  $Z_t$ .

### 4.2 Univariate example

As a first real example we used average daily temperature (in Celsius degrees) for each day of the year and average daily rainfall (in mm) for each day of the year rounded to 0.1 mm at

**Table 1** Average raw and functional HSIC and KTA coefficients for artificial time series (number in brackets means standard deviation)

	$(X_t, Y_t)$	$(X_t, Z_t)$	$(Y_t, Z_t)$
<i>Raw</i>			
HSIC	0.795 (0.015)	0.672 (0.027)	0.825 (0.014)
KTA	0.758 (0.019)	0.601 (0.028)	0.789 (0.019)
<i>Functional</i>			
HSIC	0.986 (0.000)	0.984 (0.001)	0.988 (0.000)
KTA	0.999 (0.000)	0.999 (0.000)	0.999 (0.000)



**Fig. 3**  $p$  Values from permutation-based tests for raw and functional variants of HSIC and KTA coefficients

**Table 2** Average  $p$  values from permutation-based tests for raw and functional variants of HSIC and KTA coefficients (number in brackets means standard deviation)

	$(X_t, Y_t)$	$(X_t, Z_t)$	$(Y_t, Z_t)$
<i>Raw</i>			
HSIC	0.000 (0.000)	0.445 (0.290)	0.458 (0.282)
KTA	0.000 (0.000)	0.445 (0.290)	0.458 (0.282)
<i>Functional</i>			
HSIC	0.000 (0.000)	0.077 (0.129)	0.125 (0.169)
KTA	0.000 (0.000)	0.077 (0.129)	0.077 (0.129)

35 different weather stations in Canada from 1960 to 1994. Each station belongs to one of four climate zone: Arctic (3 stations—blue color on plots), Atlantic (15—red color on plots), Continental (12—black color on plots) or Pacific (5—green color on plots) zone (Fig. 4). This data set comes from Ramsay and Silverman (2005).

In the first step, we smoothed data. We used the Fourier basis with various values of the smoothing parameter (number of basis functions) from 3 to 15. We can observe the effect



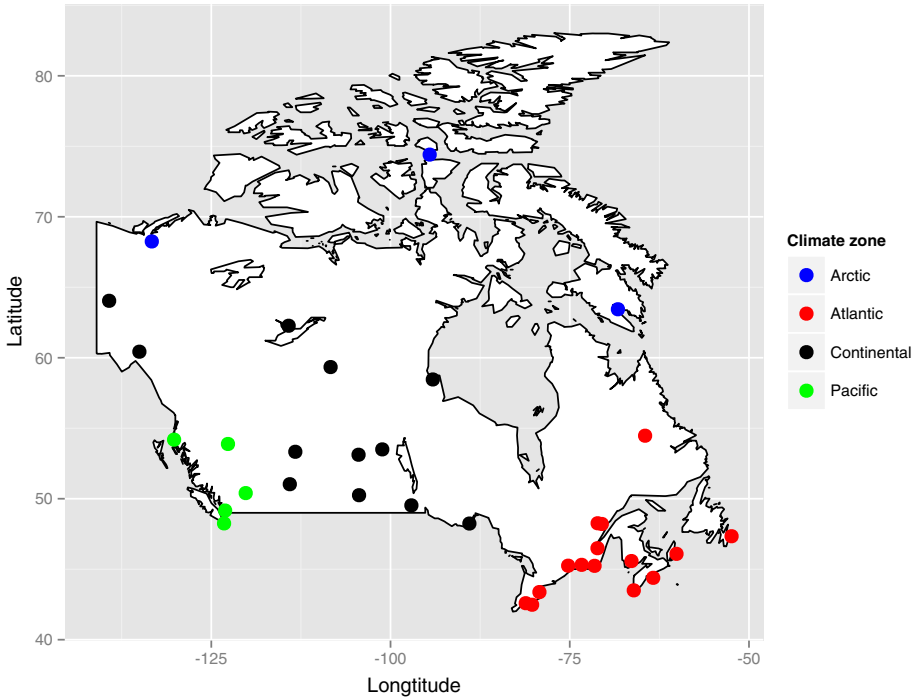


Fig. 4 Location of Canadian weather stations

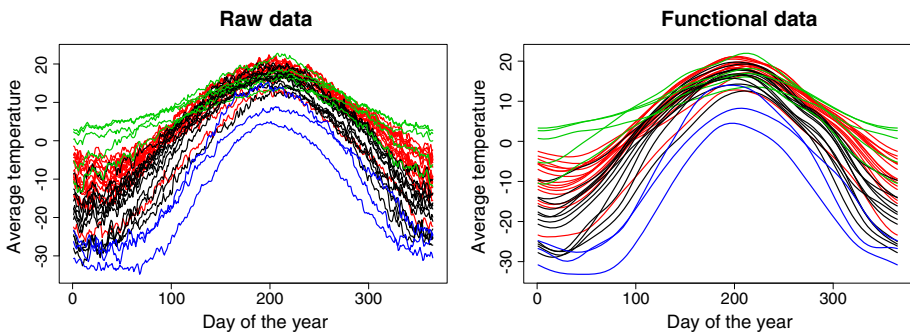


Fig. 5 Raw and functional temperature for Canadian weather stations

of smoothing in Figs. 5 and 6 (for Fourier basis with 15 basis functions). We decided to use the Fourier basis for two reasons: it has excellent computational properties, especially if the observations are equally spaced, and it is natural for describing periodic data, such as the annual weather cycles. Here, raw data are two 35 by 100 dimensional matrices (one for temperature and second for precipitation). On the other hand, functional data are two 35 by 15 dimensional matrices (coefficients of Fourier basis).

From the plots we can observe that the level of smoothness seems big enough. Additionally, we can observe some relationship between average temperature and precipitation. Namely, for weather stations with large average temperature, we observe relatively bigger average precipitation while for Arctic stations with lowest average temperatures we observe

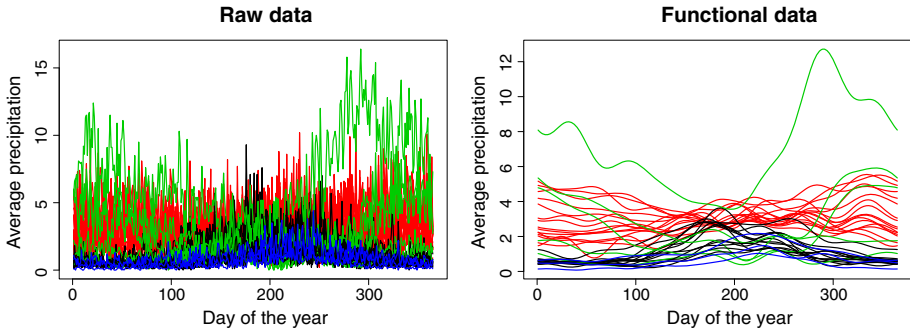
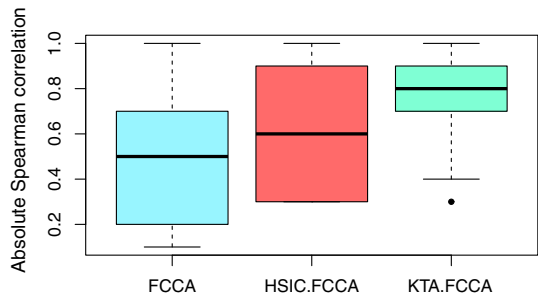


Fig. 6 Raw and functional precipitation for Canadian weather stations

Fig. 7 Absolute Spearman correlation coefficient for the first set of functional canonical variables



the smallest average precipitation. So we can expect some relationship between average temperature and average precipitation for Canadian weather stations.

In the next step, we calculated the values of described earlier coefficients, the values of which are presented in Fig. 8. We observe quite big values of HSIC and KTA, but it is impossible to infer dependency from these values. We see that the values of HSIC and KTA coefficients are stable (both do not depend on basis size).

To statistically confirm the association between temperature and precipitation we performed some simulation study. This study based on Chang et al. (2013) simulation. Finding a good nonlinear dependency measure is not trivial. KTA and HSIC are not on the same scale. As Chang et al. (2013) we used Spearman correlation coefficient. We performed 50 random splits with the inclusion of 25 samples to identify models. The Spearman correlation coefficient was then calculated using the remaining 10 samples for each of 50 splits. As we know the strongest signal between the temperature and precipitation for Canadian weather stations is nonlinear (Chang et al. 2013). From Fig. 7 we can observe that HSIC.FCCA and KTA.FCCA have produced larger absolute Spearman coefficients than FCCA. Such results suggest that HSIC.FCCA & KTA.FCCA can be viewed as natural nonlinear extensions of CCA also in the case of multivariate functional data.

Finally, we performed permutation-based tests for HSIC and KTA coefficients. The results are presented in Fig. 8. All tests rejected  $H_0$  ( $p$  values close to 0) for all basis sizes, so we can infer that we have some relationship between average temperature and average precipitation for Canadian weather stations. Unfortunately, we know nothing about the strength and direction of the dependency. Only a visual inspection of the plots suggests that there is a strong and positive relationship.

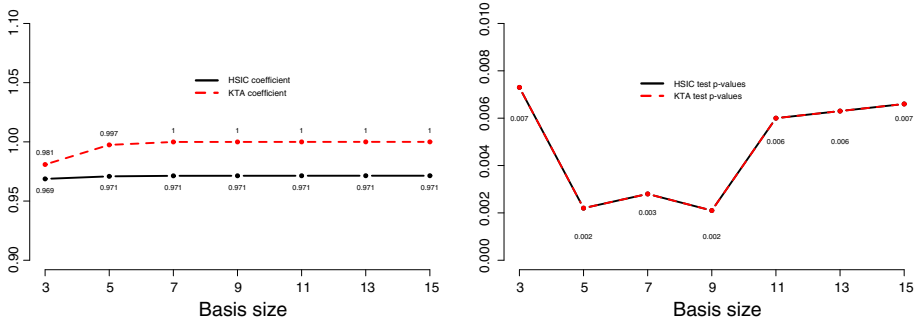


Fig. 8 HSIC and KTA coefficients and  $p$  values of permutation-based tests for Canadian weather data

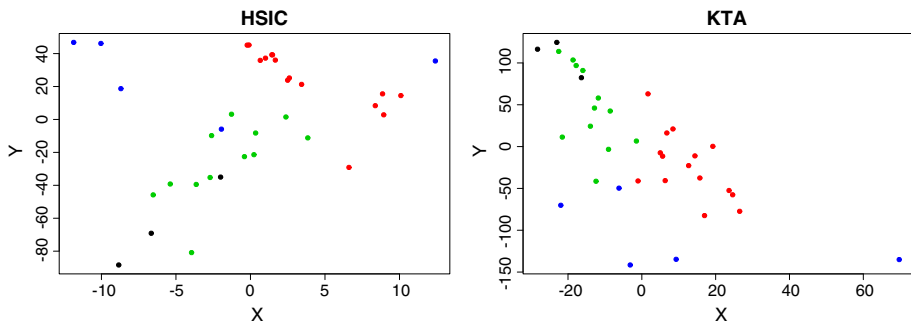


Fig. 9 Projection of the 35 Canadian weather stations on the plane  $(\hat{U}_1, \hat{V}_1)$

The relative positions of the 35 Canadian weather stations in the system  $(\hat{U}_1, \hat{V}_1)$  of functional canonical variables are shown in Fig. 9. It seems that for both coefficients the weather stations group reasonably.

### 4.3 Multivariate example

The described method was employed here to cluster the twelve groups (pillars) of variables of 38 European countries in the period 2008-2015. The list of countries used in the dependency analysis is contained in Table 3. Table 4 describes the pillars used in the analysis. For this purpose, use was made of data published by the World Economic Forum (WEF) in its annual reports (<http://www.weforum.org>). Those are comprehensive data, describing exhaustively various socio-economic conditions or spheres of individual states (Górecki et al. 2016). The data were transformed into functional data. Calculations were performed using the Fourier basis. In view of a small number of time periods ( $J = 7$ ), for each variable the maximum number of basis components was taken to be equal to five. Here, raw data are twelve matrices (one for each pillar). Dimensions of matrices are different and depend on the number of variables in the pillar. Eg. for the first pillar we have 16 (number of variables) \* 7 (number of time points) = 112 columns, hence dimensionality of the matrix for this pillar is 38 by 112. Similarly for the others. On the other hand, functional data are twelve matrices with 38 rows and appropriate number of columns (coefficients of Fourier basis). Number of columns for functional data eg. for the first pillar we calculate as 16 (number of variables) \* 5 (number of basis elements) = 80.

**Table 3** Countries used in analysis, 2008–2015

1	Albania (AL)	14	Greece (GR)	27	Poland (PL)
2	Austria (AT)	15	Hungary (HU)	28	Portugal (PT)
3	Belgium (BE)	16	Iceland (IS)	29	Romania (RO)
4	Bosnia and Herzegovina (BA)	17	Ireland (IE)	30	Russian Federation (RU)
5	Bulgaria (BG)	18	Italy (IT)	31	Serbia (XS)
6	Croatia (HR)	19	Latvia (LV)	32	Slovak Republic (SK)
7	Cyprus (CY)	20	Lithuania(LT)	33	Slovenia (SI)
8	Czech Republic (CZ)	21	Luxembourg (LU)	34	Spain (ES)
9	Denmark (DK)	22	Macedonia FYR (MK)	35	Sweden (SE)
10	Estonia (EE)	23	Malta (MT)	36	Switzerland (CH)
11	Finland (FI)	24	Montenegro (ME)	37	Ukraine (UA)
12	France (FR)	25	Netherlands (NL)	38	United Kingdom (GB)
13	Germany (DE)	26	Norway (NO)		

**Table 4** Pillars used in analysis, 2008–2015

	Pillar	Number of variables
G1	Institutions	16
G2	Infrastructure	6
G3	Macroeconomic environment	2
G4	Health and primary education	7
G5	Higher education and training	6
G6	Goods market efficiency	10
G7	Labor market efficiency	6
G8	Financial market development	5
G9	Technological readiness	4
G10	Market size	4
G11	Business sophistication	9
G12	Innovation	5

Tables 5 and 6 contain the values of functional HSIC and KTA coefficients. As expected, they are all close to one. But high values of these coefficients do not necessarily mean that there is a significant relationship between the two groups of variables. We can expect association between groups of pillars. However, it is really hard to guess what groups are associated.

Similarly to the Canadian weather example we performed small simulation study for pillars G5 and G6. From Fig. 10 we can observe that HSIC.FCCA and KTA.FCCA have produced larger absolute Spearman coefficients than FCCA. This result suggest that proposed measures have better characteristic in discovering nonlinear relationship for this example.

We performed permutation-based tests for the HSIC and KTA coefficients discussed above. For most of tests,  $p$  values were close to zero, on the basis of which it can be inferred that there is some significant relationship between the groups (pillars) of variables. Table 7 contains the  $p$  values obtained for each test. We have exactly the same  $p$  values for both methods. Now,

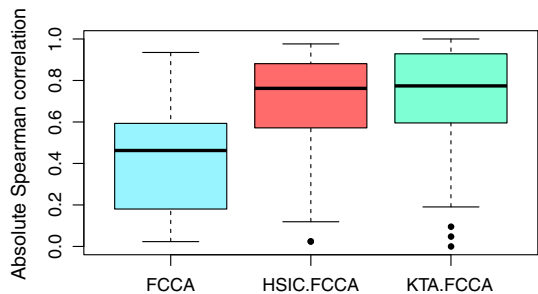
**Table 5** Functional HSIC coefficients

	1	2	3	4	5	6	7	8	9	10	11
2	0.9736										
3	0.9736	0.9737									
4	0.9736	0.9737	0.9737								
5	0.9708	0.9706	0.9706	0.9706							
6	0.9728	0.9727	0.9727	0.9727	0.9753						
7	0.9687	0.9683	0.9683	0.9683	0.9799	0.9780					
8	0.9730	0.9730	0.9730	0.9730	0.9725	0.9740	0.9721				
9	0.9736	0.9737	0.9737	0.9737	0.9706	0.9727	0.9683	0.9730			
10	0.9736	0.9737	0.9737	0.9737	0.9706	0.9727	0.9683	0.9730	0.9737		
11	0.9714	0.9711	0.9711	0.9711	0.9785	0.9755	0.9828	0.9726	0.9711	0.9711	
12	0.9688	0.9683	0.9683	0.9683	0.9778	0.9741	0.9897	0.9715	0.9783	0.9683	0.9830

**Table 6** Functional KTA coefficients

	1	2	3	4	5	6	7	8	9	10	11
2	1.0000										
3	1.0000	1.0000									
4	1.0000	1.0000	1.0000								
5	0.9918	0.9916	0.9916	0.9916							
6	0.9980	0.9978	0.9978	0.9978	0.9951						
7	0.9741	0.9736	0.9736	0.9936	0.9801	0.9821					
8	0.9991	0.9990	0.9990	0.9990	0.9933	0.9989	0.9772				
9	1.0000	1.0000	1.0000	1.0000	0.9916	0.9978	0.9736	0.9990			
10	1.0000	1.0000	1.0000	1.0000	0.9916	0.9978	0.9736	0.9990	1.0000		
11	0.9927	0.9924	0.9924	0.9924	0.9947	0.9957	0.9833	0.9936	0.9924	0.9924	
12	0.9793	0.9788	0.9788	0.9788	0.9831	0.9834	0.9794	0.9917	0.9788	0.9788	0.9887

**Fig. 10** Absolute Spearman correlation coefficient for the first set of functional canonical variables for pillars G5 & G6



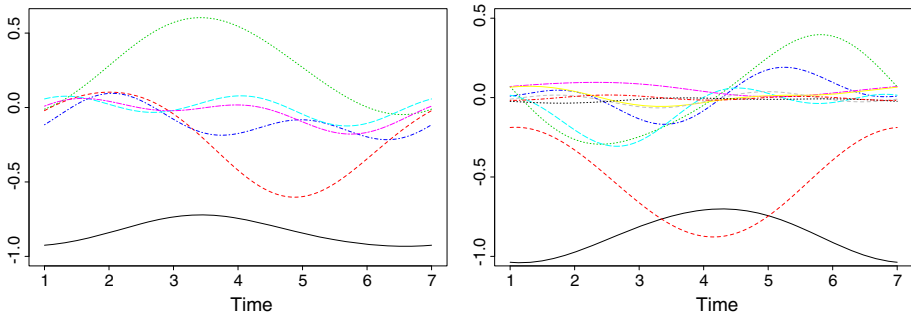
we can observe that some groups are independent ( $\alpha = 0.05$ ): G1 & G3, G3 & G6, G3 & G8, G3 & G11, G3 & G12, G4 & G9.

The graphs of the components of the vector weight function for the first functional canonical variables of the processes are shown in Fig. 11. From Fig. 11 (left) it can be seen that the greatest contribution in the structure of the first functional canonical correlation ( $U_1$ ) comes

**Table 7** Functional HSIC & KTA  $p$  values permutation-based tests (only non-zero)

	1	2	3	4	5	6	7	8	9	10	11
2	0.0142										
3	<b>0.0714</b>	0.0332									
4		0.0042	0.0343								
5		0.0001	0.0268								
6		0.0157	<b>0.0772</b>								
7		0.0009	0.0061								
8		0.0294	<b>0.0636</b>								
9	0.0030	0.0055	0.0198	<b>0.0640</b>	0.0002	0.0003	0.0009	0.0040			
10		0.0059	0.0294	0.0021					0.0055		
11		0.0039	<b>0.1034</b>						0.0008		
12		0.0008	<b>0.0563</b>						0.0044		

$p$ -values greater than usual level of significance 5% are given in bold



**Fig. 11** Weight functions for first functional canonical variable  $U_1$  (left) and  $V_1$  (right)

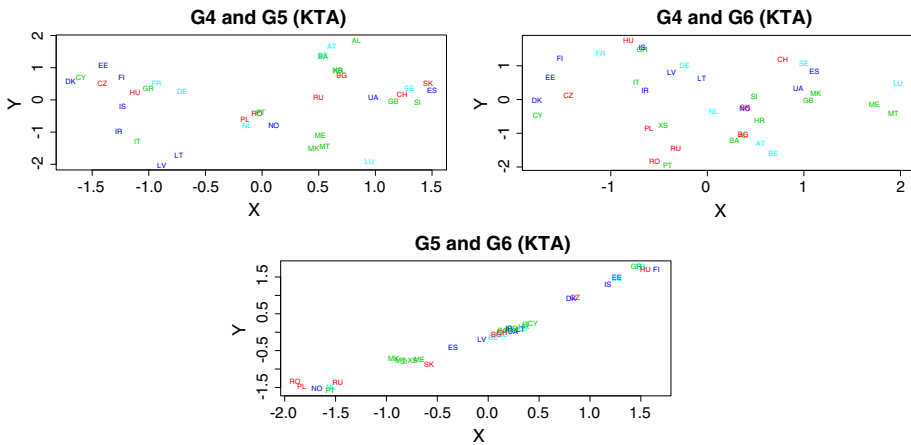
from “black” process, and this holds for all of the observation years considered. Figure 11 (right) shows that, on specified time intervals, the greatest contribution in the structure of the first second functional canonical correlation ( $V_1$  comes alternately from the processes “black” and “red dotted”). The total contribution of a particular original process in the structure of a given functional canonical correlation is equal to the area under the module weighting function corresponding to this process. These contributions for the components are given in Table 8.

Figure 12 contains the relative positions of the 38 European countries in the system ( $\hat{U}_1, \hat{V}_1$ ) of functional canonical variables for selected groups of variables. The high correlation of the first two functional canonical variables can be seen in Fig. 12 for two pillars G5 and G6. For KTA criterion, the countries with the highest value of functional canonical variables  $U_1$  and  $V_1$  are: Finland (FI), France (FR), Hungary (HU), Greece (GR), Estonia (EE), Germany (DE), Iceland (IS), Czech Republic (CZ) and Denmark (DK). The countries with the lowest value of functional canonical variables  $U_1$  and  $V_1$  are: Romania (RO), Poland (PL), Norway (NO), Portugal (PT), Netherlands (NL) and Russian Federation (RU). Other countries belong to the intermediate group.

During the numerical calculation process we used R software (R Core Team 2018) and packages *fda* (Ramsay et al. 2018) and *hsicCCA* (Chang 2013).

**Table 8** Sorted areas under module weighting functions

No.	Area	Proportion (in %)
<i>First functional canonical variable (G5)</i>		
1	5.008	51.74
2	1.724	17.81
3	1.567	16.19
4	0.713	7.36
5	0.351	3.63
6	0.317	3.27
<i>First functional canonical variable (G6)</i>		
1	5.187	44.77
2	3.194	27.56
3	1.287	11.11
4	0.580	5.00
5	0.511	4.41
6	0.323	2.79
7	0.206	1.77
8	0.152	1.31
9	0.091	0.78
10	0.057	0.49



**Fig. 12** Selected projections of the 38 European countries on the plane  $(\hat{U}_1, \hat{V}_1)$ . Regions used for statistical processing purposes by the United Nations Statistics Division: blue square—Northern Europe, cyan square—Western Europe, red square—Eastern Europe, green square—Southern Europe. (Color figure online)

### 5 Conclusions

We proposed an extension of two dependency measures for two sets of variables for multivariate functional data. We proposed to use tests to examine the significance of results because the values of proposed coefficients are rather hard to interpret. Additionally, we presented the methods of constructing nonlinear canonical variables for multivariate functional data using HSIC and KTA coefficients. Tested on two real examples, the proposed method has proven

useful in investigating the dependency between two sets of variables. Examples confirm usefulness of our approach in revealing the hidden structure of co-dependence between groups of variables.

During the study of proposed coefficients we discovered that the size of basis (smoothing parameter) is rather unimportant, the values (and  $p$  values for tests) do not depend on the basis size.

Of course, the performance of the methods needs to be further evaluated on additional real and artificial data sets. Moreover, we can examine the behavior of coefficients (and tests) for different bases like B-splines or wavelets (when data are not periodic, the Fourier basis could fail). This constitutes the direction of our future research.

**Acknowledgements** The authors are grateful to editor and two anonymous reviewers for giving many insightful and constructive comments and suggestions which led to the improvement of the earlier manuscript.

**Open Access** This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

## References

- Aronszajn N (1950) Theory of reproducing kernels. *Trans Am Math Soc* 68:337–404
- Chang B (2013) hsicCCA: Canonical Correlation Analysis based on Kernel Independence Measures. R package version 1.0. <https://CRAN.R-project.org/package=hsicCCA>
- Chang B, Kruger U, Kustra R, Zhang J (2013) Canonical correlation analysis based on hilbert-schmidt independence criterion and centered kernel target alignment. In: *Proceedings of the 30th international conference on machine learning*, Atlanta, Georgia. *JMLR: W and CP* 28(2), 316–324
- Cortes C, Mohri M, Rostamizadeh A (2012) Algorithms for learning kernels based on centered alignment. *J Mach Learn Res* 13:795–828
- Cristianini N, Shawe-Taylor J, Elisseeff A, Kandola JS (2001) On kernel-target alignment. In: *NIPS-2001*, 367–373
- Cuevas A (2014) A partial overview of the theory of statistics with functional data. *J Stat Plan Inference* 147:1–23
- Devijver E (2017) Model-based regression clustering for high-dimensional data: application to functional data. *Adv Data Anal Classif* 11(2):243–279
- Edelman A, Arias TA, Smith S (1998) The geometry of algorithms with orthogonality constraints. *SIAM J Matrix Anal Appl* 20(2):303–353
- Ferraty F, Vieu P (2003) Curves discrimination: a nonparametric functional approach. *Comput Stat Data Anal* 44(1–2):161–173
- Ferraty F, Vieu P (2006) *Nonparametric functional data analysis: theory and practice*. Springer, Berlin
- Feuerverger A (1993) A consistent test for bivariate dependence. *Int Stat Rev* 61(3):419–433
- Górecki T, Krzyśko M, Ratajczak W, Wołyński W (2016) An extension of the classical distance correlation coefficient for multivariate functional data with applications. *Stat Transit* 17(3):449–9466
- Górecki T, Krzyśko M, Wołyński W (2017) Correlation analysis for multivariate functional data. In: Palumbo F, Montanari A, Montanari M (eds) *Data science. Studies in classification, data analysis, and knowledge organization*. Springer, Berlin, pp 243–258
- Górecki T, Krzyśko M, Waszak Ł, Wołyński W (2018) Selected statistical methods of data analysis fir multivariate functional data. *Stat Papers* 59:153–182
- Górecki T, Smaga Ł (2017) Multivariate analysis of variance for functional data. *J Appl Stat* 44:2172–2189
- Gretton A., Bousquet O., Smola A., and Schölkopf B., (2005): Measuring statistical dependence with Hilbert–Schmidt norms. In: Jain S, Simon HU, Tomita E (eds) *Algorithmic learning theory. Lecture notes in computer science* 3734, 63–77. Springer
- Gretton A, Fukumizu K, Teo CH, Song L, Schölkopf B, Smola AJ (2008) A kernel statistical test of independence. In: Platt JC, Koller D, Singer Y, Roweis S (eds) *Advances in neural information processing systems*. Curran, Red Hook, pp 585–592
- Hofmann T, Schölkopf B, Smola AJ (2008) Kernel methods in machine learning. *Ann Stat* 36(3):1171–1220



- Horváth L, Kokoszka P (2012) Inference for functional data with applications. Springer, Berlin
- Hotelling H (1936) Relations between two sets of variates. *Biometrika* 28:321–377
- Hsing T, Eubank R (2015) Theoretical foundations of functional data analysis, with an introduction to linear operators. Wiley, Hoboken
- James GM, Wang JW, Zhu J (2009) Functional linear regression that's interpretable. *Ann Stat* 37(5):2083–2108
- Kankainen A (1995) Consistent testing of total independence based on the empirical characteristic function. Ph.D. thesis, University of Jyväskylä
- Martin-Baragan B, Lillo R, Romo J (2014) Interpretable support vector machines for functional data. *Eur J Oper Res* 232:146–155
- Mercer J (1909) Functions of positive and negative type and their connection with the theory of integral equations. *Philos Trans R Soc Lond Ser A* 209:415–446
- R Core Team (2018) R: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>
- Ramsay JO, Dalzell CJ (1991) Some tools for functional data analysis (with discussion). *J R Stat Soc Ser B* 53(3):539–572
- Ramsay JO, Silverman BW (2002) Applied functional data analysis. Springer, New York
- Ramsay JO, Silverman BW (2005) Functional data analysis, 2nd edn. Springer, Berlin
- Ramsay JO, Wickham H, Graves S, Hooker G (2018) fda: Functional data analysis. R package version 2.4.8. <https://CRAN.R-project.org/package=fda>
- Read T, Cressie N (1988) Goodness-of-fit statistics for discrete multivariate analysis. Springer, Berlin
- Riesz F (1909) Sur les opérations fonctionnelles linéaires. *Comptes rendus hebdomadaires des séances de l'Académie des sciences* 149:974–977
- Sejdicinovic D, Sriperumbudur B, Gretton A, Fukumizu K (2013) Equivalence of distance-based and RKHS-based statistics in hypothesis testing. *Ann Stat* 41(5):2263–2291
- Schölkopf B, Smola AJ, Müller KR (1998) Nonlinear component analysis as a kernel eigenvalue problem. *Neural Comput* 10:1299–1319
- Shawe-Taylor J, Cristianini N (2004) Kernel methods for pattern analysis. Cambridge University Press, Cambridge
- Song L, Boots B, Siddiqi S, Gordon G, Somla A (2010) Hilbert space embeddings of hidden Markov models. In: Proceedings of the 26th international conference on machine learning (ICML2010)
- Székely GJ, Rizzo ML, Bakirov NK (2007) Measuring and testing dependence by correlation of distances. *Ann Stat* 35(6):2769–2794
- Székely GJ, Rizzo ML (2009) Brownian distance covariance. *Ann Appl Stat* 3(4):1236–1265
- Wang T, Zhao D, Tian S (2015) An overview of kernel alignment and its applications. *Artif Intell Rev* 43(2):179–192
- Zhang K, Peters J, Janzing D, Schölkopf B (2011) Kernel-based conditional independence test and application in causal discovery. In: Cozman FG, Pfeffer A (eds) Proceedings of the 27th conference on uncertainty in artificial intelligence, AUAI Press, Corvallis, OR, USA, 804–813