

Measuring and testing mutual dependence for functional data

Tomasz Górecki, **Mirosław Krzyśko**, Waldemar Wołyński

Faculty of Mathematics and Computer Science
Adam Mickiewicz University, Poznań, Poland

The XLV annual conference "Statystyka Matematyczna"
Będlewo, 2 - 6 December 2019



- 1 *Introduction*
- 2 *Functional data model*
 - Stochastic processes
 - Smoothing
- 3 *$K = 2$ case*
 - Functional coefficient ρV
 - Functional distance correlation
- 4 *$K > 2$ case*
 - Coefficient of mutual correlation ρMV
 - Measures of multiple independence
- 5 *Example*

Let's consider the **problem of testing mutual independence** for random vectors $\mathbf{X}_1, \dots, \mathbf{X}_K$, $K \geq 2$, $\mathbf{X}_i \in L_2^{p_i}(I)$, $i = 1, \dots, K$, i.e. we are testing the following hypothesis:

$$H_0: \mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2 \perp\!\!\!\perp \dots \perp\!\!\!\perp \mathbf{X}_K$$

vs

$$H_1: \neg H_0.$$

This problem is connected with examining the significance of the **coefficient of mutual dependency**.

In our presentation we will present methods of measuring mutual independence for **functional data**. We will consider two or more vector random processes. Based on the obtained measures we will construct permutational tests to test mutual independence for functional data.

Let us assume that $\mathbf{X} \in L_2^p(I)$ is a random process, where $L_2(I)$ is a Hilbert space of square integrable functions on the interval I . Additionally, we also assume that

$$E(\mathbf{X}(t)) = \mathbf{0}, \quad t \in I.$$

We will further assume that each component X_g of the process \mathbf{X} can be represented by a **finite number of basis functions** $\{\varphi_e\}$:

$$X_g(t) = \sum_{e=0}^{B_g} \alpha_{ge} \varphi_e(t), s \in I, g = 1, 2, \dots, p.$$

The **degree of smoothness** of function X_g depends on the value B_g (a small value causes more smoothing of the function).

We introduce the following notation:

$$\boldsymbol{\alpha} = (\alpha_{10}, \dots, \alpha_{1B_1}, \dots, \alpha_{p0}, \dots, \alpha_{pB_p})^\top,$$

$$\boldsymbol{\Phi}(t) = \begin{bmatrix} \boldsymbol{\varphi}_1^\top(t) & \mathbf{0}^\top & \dots & \mathbf{0}^\top \\ \mathbf{0}^\top & \boldsymbol{\varphi}_2^\top(t) & \dots & \mathbf{0}^\top \\ \dots & \dots & \dots & \dots \\ \mathbf{0}^\top & \mathbf{0}^\top & \dots & \boldsymbol{\varphi}_p^\top(t) \end{bmatrix},$$

where $\boldsymbol{\varphi}_k(t) = (\varphi_0(t), \varphi_1(t), \dots, \varphi_{B_k}(t))^\top$, $k = 1, 2, \dots, p$ are **orthonormal basis functions** of space $L_2(I)$.

Using this matrix notation the random process \mathbf{X} can be represented as

$$\mathbf{X}(t) = \boldsymbol{\Phi}(t)\boldsymbol{\alpha},$$

where $\boldsymbol{\alpha} \in \mathbb{R}^{K+p}$, $K = B_1 + B_2 + \dots + B_p$.

Notation

Let $\mathbf{X} \in L_2^p(I)$ i $\mathbf{Y} \in L_2^q(I)$ be random processes with the following representations:

$$\mathbf{X}(t) = \Phi_1(t)\alpha, \quad \mathbf{Y}(s) = \Phi_2(s)\beta, \quad t, s \in I.$$

Additionally, let the covariance matrix of processes \mathbf{X} and \mathbf{Y} have the form

$$\Sigma(t, s) = \begin{bmatrix} \Sigma_{XX}(t, s) & \Sigma_{XY}(t, s) \\ \Sigma_{YX}(t, s) & \Sigma_{YY}(t, s) \end{bmatrix}, \quad t, s \in I.$$

Then

$$\Sigma_{XX}(t, s) = \Phi_1(t)\Sigma_{\alpha\alpha}\Phi_1^\top(s),$$

$$\Sigma_{XY}(t, s) = \Phi_1(t)\Sigma_{\alpha\beta}\Phi_2^\top(s),$$

$$\Sigma_{YY}(t, s) = \Phi_2(t)\Sigma_{\beta\beta}\Phi_2^\top(s),$$

where

$$\Sigma = \begin{bmatrix} \Sigma_{\alpha\alpha} & \Sigma_{\alpha\beta} \\ \Sigma_{\beta\alpha} & \Sigma_{\beta\beta} \end{bmatrix}$$

is a covariance matrix of vectors α i β .

Functional coefficient ρ_V is a nonnegative number given by

$$\rho_{V_{\mathbf{X}, \mathbf{Y}}} = \frac{\|\boldsymbol{\Sigma}_{XY}\|_F}{\sqrt{\|\boldsymbol{\Sigma}_{XX}\|_F \|\boldsymbol{\Sigma}_{YY}\|_F}},$$

where

$$\|\boldsymbol{\Sigma}_{XY}\|_F = \sqrt{\int_I \int_I \text{tr}(\boldsymbol{\Sigma}_{XY}^\top(t, s) \boldsymbol{\Sigma}_{XY}(t, s)) dt ds.}$$

If processes \mathbf{X} and \mathbf{Y} have the form

$$\mathbf{X}(t) = \Phi_1(t)\boldsymbol{\alpha}, \quad \mathbf{Y}(s) = \Phi_2(s)\boldsymbol{\beta}, \quad t, s \in I,$$

then

$$\rho V_{\mathbf{X}, \mathbf{Y}} = \rho V_{\boldsymbol{\alpha}, \boldsymbol{\beta}}.$$

Coefficient ρV for random vectors was proposed by **Escoufier (1973)**.

Escoufier, G.Y. (1973): Le traitement des variables vectorielles, Biometrics 29, 751-760.

Let's assume that the joint distribution of random vectors α and β is $p + q$ dimensional normal distribution.

Problem of testing the null hypothesis

$$H_0: X \perp\!\!\!\perp Y$$

is equivalent to the problem of testing the null hypothesis

$$H_0: \rho V_{X,Y} = 0.$$

which is equivalent to the problem of testing the null hypothesis

$$H_0: \rho V_{\alpha,\beta} = 0$$

for a pair of random vectors α and β . To verify this hypothesis we use a **permutational test**.

For random vectors $\alpha \in \mathbb{R}^p$ and $\beta \in \mathbb{R}^q$ Székely et al. (2007) introduced the **distance correlation** (dCorr) as a nonnegative number given by

$$\text{dCorr}(\alpha, \beta) = \frac{\text{dCov}(\alpha, \beta)}{\sqrt{\text{dCov}(\alpha, \alpha) \text{dCov}(\beta, \beta)}},$$

where

$$\text{dCov}(\alpha, \beta) = \|\phi_{\alpha, \beta}(\lambda, \mu) - \phi_{\alpha}(\lambda)\phi_{\beta}(\mu)\|_w,$$

and

$$\|f\|_w = \sqrt{\iint |f(\lambda, \mu)|^2 w(\lambda, \mu) d\lambda d\mu}.$$

Székely, G.J., Rizzo, M.L., Bakirov, N.K. (2007): Measuring and testing dependence by correlation of distances, The Annals of Statistics 35, 2769-2794.

Functional distance correlation

For jointly distributed random processes $\mathbf{X} \in L_2^p(I)$ and $\mathbf{Y} \in L_2^q(I)$, let

$$\phi_{\mathbf{X}, \mathbf{Y}}(l, \mathbf{m}) = E\{\exp[i \langle l, \mathbf{X} \rangle + i \langle \mathbf{m}, \mathbf{Y} \rangle]\}$$

be the **joint characteristic function** of (\mathbf{X}, \mathbf{Y}) . If processes \mathbf{X} and \mathbf{Y} have the form

$$\mathbf{X}(t) = \Phi_1(t)\alpha, \quad \mathbf{Y}(s) = \Phi_2(s)\beta, \quad t, s \in I,$$

and if

$$l(s) = \Phi_1(s)\lambda, \quad \mathbf{m}(t) = \Phi_2(t)\mu,$$

then

$$\phi_{\mathbf{X}, \mathbf{Y}}(l, \mathbf{m}) = E\{\exp[i\lambda'\alpha + i\mu'\beta]\} = \phi_{\alpha, \beta}(\lambda, \mu).$$

We define the **functional distance correlation** in the following way

$$\text{dCorr}(\mathbf{X}, \mathbf{Y}) = \text{dCorr}(\boldsymbol{\alpha}, \boldsymbol{\beta}).$$

Górecki T., Krzyśko M., Wołyński W. (2017): Correlation analysis for multivariate functional data. In: Data Science, Studies in Classification, Data Analysis, and Knowledge Organization, Springer International Publishing, 243-258.

Górecki T., Krzyśko M., Wołyński W. Independence test and canonical correlation analysis based on the alignment between kernel matrices for multivariate functional data. Artificial Intelligence Review (in press).

Problem of testing the null hypothesis

$$H_0: \mathbf{X} \perp\!\!\!\perp \mathbf{Y}$$

is equivalent to the problem of testing the null hypothesis

$$H_0: \text{dCorr}(\mathbf{X}, \mathbf{Y}) = 0,$$

which is equivalent to the problem of testing the null hypothesis

$$H_0: \text{dCorr}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = 0$$

for a pair of random vectors $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$. To verify this hypothesis we use a **permutational test**.

Let $\mathbf{X}_1 \in L_2^{p_1}(I)$, $\mathbf{X}_2 \in L_2^{p_2}(I)$, \dots , $\mathbf{X}_K \in L_2^{p_K}(I)$ be random processes with the following representation:

$$\mathbf{X}_1(t) = \Phi_1(t)\alpha_1, \mathbf{X}_2(t) = \Phi_2(t)\alpha_2, \dots, \mathbf{X}_K(t) = \Phi_K(t)\alpha_K, t \in I.$$

Additionally, let the matrix of covariance for vectors $\alpha_1, \alpha_2, \dots, \alpha_K$ has the form:

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} & \cdots & \Sigma_{1K} \\ \Sigma_{21} & \Sigma_{22} & \cdots & \Sigma_{2K} \\ \vdots & \vdots & & \vdots \\ \Sigma_{K1} & \Sigma_{K2} & \cdots & \Sigma_{KK} \end{bmatrix}.$$

Assuming joint $p_1 + p_2 + \dots + p_K$ dimensional normal distribution of vectors $\alpha_1, \alpha_2, \dots, \alpha_K$, the problem of testing the null hypothesis

$$H_0: \mathbf{X}_1 \perp \mathbf{X}_2 \perp \dots \perp \mathbf{X}_K$$

is equivalent to the problem of testing the null hypothesis

$$H_0: \sum_{i < j} \|\Sigma_{ij}\|_F = 0.$$

We define **coefficient of mutual correlation ρ_{MV}** as a positive number given by

$$\rho^2_{MV} = \frac{2}{K(K-1)} \sum_{i < j} \rho^2 V(\mathbf{X}_i, \mathbf{X}_j).$$

Assuming joint $p_1 + p_2 + \dots + p_k$ dimensional normal distribution of vectors $\alpha_1, \alpha_2, \dots, \alpha_K$, the problem of testing the null hypothesis

$$H_0: \mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2 \perp\!\!\!\perp \dots \perp\!\!\!\perp \mathbf{X}_K$$

is equivalent to the problem of testing the null hypothesis

$$H_0: \rho_{MV} = 0.$$

Measures of multiple independence

Let $\text{Corr}(\mathbf{X}, \mathbf{Y})$ be **some measure of dependence** for random vectors \mathbf{X} and \mathbf{Y} with property: $\text{Corr}(\mathbf{X}, \mathbf{Y}) = 0$ if and only if random vectors \mathbf{X} and \mathbf{Y} are independent.

Let $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_K)$, and

$$\mathbf{X}_{c+} = (\mathbf{X}_{c+1}^\top, \dots, \mathbf{X}_K^\top)^\top, \quad c = 1, \dots, K-1,$$

$$\mathbf{X}_{c-} = (\mathbf{X}_1^\top, \dots, \mathbf{X}_{c-1}^\top, \mathbf{X}_{c+1}^\top, \dots, \mathbf{X}_K^\top)^\top, \quad c = 1, \dots, K.$$

We define **the coefficients of multiple independence** as

$$\mathcal{R}(\mathbf{X}) = \frac{1}{K-1} \sum_{c=1}^{K-1} \text{Corr}^2(\mathbf{X}_c, \mathbf{X}_{c+}),$$

and

$$\mathcal{S}(\mathbf{X}) = \frac{1}{K} \sum_{c=1}^K \text{Corr}^2(\mathbf{X}_c, \mathbf{X}_{c-}).$$

Theorem

$$\mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2 \perp\!\!\!\perp \cdots \perp\!\!\!\perp \mathbf{X}_K \iff \mathcal{R}(\mathbf{X}) = 0,$$

$$\mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2 \perp\!\!\!\perp \cdots \perp\!\!\!\perp \mathbf{X}_K \iff \mathcal{S}(\mathbf{X}) = 0.$$

Remark

In the place of Corr we can put eg. ρV or dCorr.

Hence, the problem of testing the null hypothesis

$$H_0: \mathbf{X}_1 \perp\!\!\!\perp \mathbf{X}_2 \perp\!\!\!\perp \cdots \perp\!\!\!\perp \mathbf{X}_K$$

is equivalent to the problem of testing the null hypothesis

$$H_0: \mathcal{R}(\mathbf{X}) = 0 \ (\mathcal{S}(\mathbf{X}) = 0).$$

Example

Let

$$X_t = \varepsilon_{1t},$$

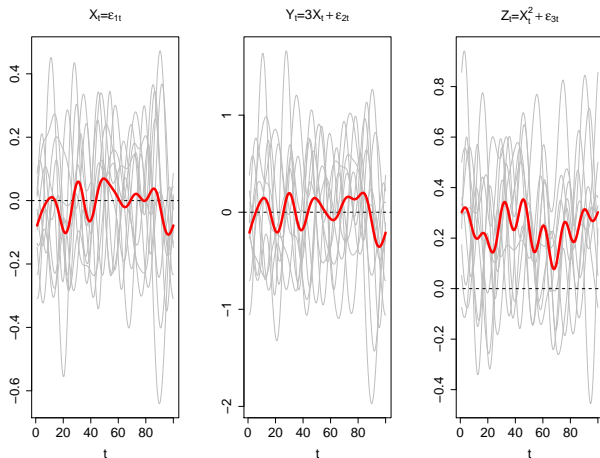
$$Y_t = 3X_t + \varepsilon_{2t},$$

$$Z_t = X_t^2 + \varepsilon_{3t},$$

where ε_{1t} , ε_{2t} and ε_{3t} are independent random variables with $N(0, 0.25)$ distribution. We generated 1000 random realizations for each process with length 100. To smooth the data we used Fourier series with 15 elements.

Processes X_t and Y_t are **linearly dependent**. Processes X_t , Z_t and Y_t , Z_t are **non-linearly dependent**.

Example



10 randomly selected realizations of processes X_t , Y_t and Z_t (functional means in red).

Example

Coefficient	Processes	p -value
ρMV	X_t, Y_t, Z_t	0.014
$\mathcal{R} - \rho V$	X_t, Y_t, Z_t	0.006
$\mathcal{S} - \rho V$	X_t, Y_t, Z_t	0.027
$\mathcal{R} - d\text{Corr}$	X_t, Y_t, Z_t	0.007
$\mathcal{S} - d\text{Corr}$	X_t, Y_t, Z_t	0.016
ρV	X_t vs Y_t	0.001
	X_t vs Z_t	0.367
	Y_t vs Z_t	0.481
dCorr	X_t vs Y_t	0.001
	X_t vs Z_t	0.003
	Y_t vs Z_t	0.025