

Miary zależności między dwoma wielowymiarowymi procesami stochastycznymi

Mirosław Krzyśko

mkrzysko@amu.edu.pl

Wydział Matematyki i Informatyki
Uniwersytet im. Adama Mickiewicza w Poznaniu



1 Wstęp

2 Wyglądanie danych

3 Współczynnik $I\rho V$

- Scałkowany współczynnik ρV
- Testowanie istotności współczynnika $I\rho V$

4 Korelacja odległościowa

- Korelacja odległościowa
- Testowanie istotności korelacji odległościowej

5 Jądrowy współczynnik zgodności

- Jądrowy współczynnik zgodności dla wektorów losowych
- Test istotności współczynnika zgodności dla wektorów losowych
- Jądrowy współczynnik zgodności dla procesów losowych
- Test istotności współczynnika zgodności dla procesów losowych

6 Przykład

7 Literatura

W ostatnich latach wiele uwagi poświęcono danym reprezentującym funkcje lub krzywe. Dane tego typu znane są w literaturze pod nazwą **danych funkcjonalnych** (Ramsay & Silverman, 2005). Przykłady danych funkcjonalnych można znaleźć w wielu dziedzinach takich, jak medycyna, ekonomia, meteorologia itd. W wielu zastosowaniach używamy metod statystycznych dla obiektów scharakteryzowanych wieloma cechami obserwowanymi w wielu momentach czasowych (dane podwójnie wielowymiarowe). Dane tego typu nazywamy **wielowymiarowymi danymi funkcjonalnymi**. Autorem pionierskiej pracy teoretycznej z tego zakresu był Besse (1979). W jego pracy zmienne losowe przyjmowały wartości w przestrzeni Hilberta. Saporta (1981) rozważał metody analizy czynnikowej (analiza składowych głównych i korelacji kanonicznych) w zastosowaniu do wielowymiarowych danych funkcjonalnych.

Założmy, że $\mathbf{X} \in L_2^p(I_1)$ i $\mathbf{Y} \in L_2^q(I_2)$ są procesami losowymi, gdzie $L_2(I)$ jest **przestrzenią Hilberta funkcji całkowlanych z kwadratem na przedziale I** .

Ponadto założmy, że

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

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

Założenie to nie powoduje utraty ogólności rozważań, ponieważ funkcjonalne współczynniki korelacji są wyznaczane jedynie na podstawie macierzy kowariancji procesów \mathbf{X} i \mathbf{Y} postaci

$$\Sigma(s, t) = \begin{bmatrix} \Sigma_{XX}(s, t) & \Sigma_{XY}(s, t) \\ \Sigma_{YX}(s, t) & \Sigma_{YY}(s, t) \end{bmatrix}, \quad s \in I_1, \quad t \in I_2.$$

W dalszym ciągu zakładamy, że każda składowa X_g procesu \mathbf{X} i każda składowa Y_h procesu \mathbf{Y} może być reprezentowana przez **skończoną liczbę funkcji bazowych** $\{\varphi_e\}$ i $\{\varphi_f\}$:

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

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

Stopień gładkości funkcji X_g i Y_h zależy od wartości E_g i F_h (małe wartości dają mniejsze wygładzenie funkcji).

Wprowadźmy następujące oznaczenia:

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

$$\beta = (\beta_{10}, \dots, \beta_{1F_1}, \dots, \beta_{q0}, \dots, \beta_{qF_q})',$$

$$\Phi_1(s) = \begin{bmatrix} \varphi'_{E_1}(s) & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \varphi'_{E_2}(s) & \dots & \mathbf{0} \\ \dots & \dots & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \dots & \varphi'_{E_p}(s) \end{bmatrix},$$

$$\Phi_2(t) = \begin{bmatrix} \varphi'_{F_1}(t) & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \varphi'_{F_2}(t) & \dots & \mathbf{0} \\ \dots & \dots & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \dots & \varphi'_{F_q}(t) \end{bmatrix},$$

gdzie $\varphi_{E_1}, \dots, \varphi_{E_p}$ i $\varphi_{F_1}, \dots, \varphi_{F_q}$ są **ortonormalnymi** funkcjami bazowymi przestrzeni $L_2(I_1)$ i $L_2(I_2)$ oraz $K_1 = E_1 + \dots + E_p$, $K_2 = F_1 + \dots + F_q$.

Używając notacji macierzowej procesy $\mathbf{X}(s)$ i $\mathbf{Y}(t)$ mają następującą reprezentację

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

Oznacza to, że realizacje procesów \mathbf{X} i \mathbf{Y} zawarte są w **skończenie wymiarowych podprzestrzeniach** przestrzeni $L_2^p(I_1)$ i $L_2^q(I_2)$ odpowiednio. Te podprzestrzenie będziemy oznaczali przez $\mathcal{L}_2^p(I_1)$ i $\mathcal{L}_2^q(I_2)$.

Ponadto, dla wektorów α i β mamy:

$$E(\alpha) = \mathbf{0}, \quad E(\beta) = \mathbf{0}$$

oraz

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

Zazwyczaj proces jest obserwowany w skończonej liczbie momentów czasowych. Proces transformacji danych dyskretnych do danych funkcjonalnych jest wykonywany oddzielnie dla każdej realizacji każdej składowej procesu. Niech x_{gj} oznacza obserwowaną wartość składowej X_g , $g = 1, 2, \dots, p$ w j -tym momencie czasowym s_j , gdzie $j = 1, 2, \dots, J$. Podobnie, niech y_{hj} oznacza obserwowaną wartość składowej Y_h , $h = 1, 2, \dots, q$ w j -tym momencie czasowym t_j , gdzie $j = 1, 2, \dots, J$. Wtedy nasze dane składają się z pJ par (s_j, x_{gj}) oraz z qJ par (t_j, y_{hj}) . Współczynniki α_i i β_i są estymowane **metodą najmniejszych kwadratów**. Oznaczmy $\hat{\alpha}_i$ i $\hat{\beta}_i$ przez \mathbf{a}_i i \mathbf{b}_i , $i = 1, 2, \dots, n$.

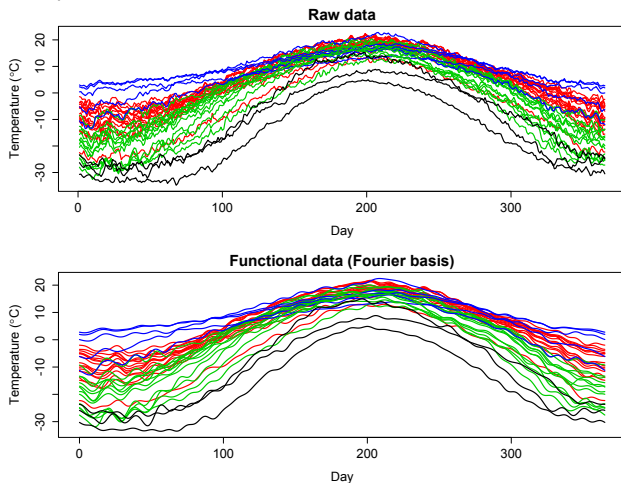
W rezultacie procesu transformacji otrzymujemy dane funkcjonalne postaci:

$$\mathbf{x}_i(s) = \Phi_1(s)\mathbf{a}_i, \quad \mathbf{y}_i(t) = \Phi_2(t)\mathbf{b}_i,$$

gdzie $s \in I_1$, $t \in I_2$, $i = 1, 2, \dots, n$. Niech $\mathbf{A} = (\mathbf{a}'_1, \mathbf{a}'_2, \dots, \mathbf{a}'_n)$ oraz $\mathbf{B} = (\mathbf{b}'_1, \mathbf{b}'_2, \dots, \mathbf{b}'_n)$. Wtedy

$$\hat{\Sigma}_{\alpha,\alpha} = \frac{1}{n}\mathbf{A}'\mathbf{A}, \quad \hat{\Sigma}_{\beta,\beta} = \frac{1}{n}\mathbf{B}'\mathbf{B}, \quad \hat{\Sigma}_{\alpha,\beta} = \frac{1}{n}\mathbf{A}'\mathbf{B}.$$

Średnie dzienne temperatury dla 35 stacji meteorologicznych w Kanadzie w latach 1960 – 1994. Dane pochodzą z książki Ramsay'a i Silvermana (2005), *Functional Data Analysis*, Springer, New York.



Kwadrat scałkowanego współczynnika ρV (Robert i Escoufier (1976)) definiujemy jako

$$(I\rho V)^2 = \frac{\|\Sigma_{XY}\|^2}{\|\Sigma_{XX}\| \|\Sigma_{YY}\|},$$

gdzie

$$\|\Sigma_{XY}\| = \sqrt{\int_{I_1} \int_{I_2} \text{tr}(\Sigma'_{XY}(s, t) \Sigma_{XY}(s, t)) ds dt}.$$

Własności:

- $0 \leq I\rho V \leq 1$.
- $I\rho V = 0$ wtedy i tylko wtedy, gdy dla każdych $g = 1, 2, \dots, p$ i $h = 1, 2, \dots, q$ takich, że $g \neq h$, składowa X_g procesu \mathbf{X} i składowa Y_h procesu \mathbf{Y} są nieskorelowane.

Ponieważ

$$\Sigma_{XY}(s, t) = \Phi_1(s) \Sigma_{\alpha\beta} \Phi_2'(t)$$

to

$$\|\Sigma_{XY}\|^2 = \text{tr}(\Sigma'_{\alpha\beta} \Sigma_{\alpha\beta}).$$

Analogicznie

$$\|\Sigma_{XX}\|^2 = \text{tr}(\Sigma'_{\alpha\alpha} \Sigma_{\alpha\alpha}),$$

$$\|\Sigma_{YY}\|^2 = \text{tr}(\Sigma'_{\beta\beta} \Sigma_{\beta\beta}).$$

Stąd

$$(I\rho V)^2 = \frac{\text{tr}(\Sigma'_{\alpha\beta} \Sigma_{\alpha\beta})}{\sqrt{\text{tr}(\Sigma'_{\alpha\alpha} \Sigma_{\alpha\alpha}) \text{tr}(\Sigma'_{\beta\beta} \Sigma_{\beta\beta})}}.$$

Twierdzenie

Scałkowany współczynnik korelacji $I\rho V$ pary procesów losowych $\mathbf{X} \in \mathcal{L}_2^p(I_1)$ i $\mathbf{Y} \in \mathcal{L}_2^q(I_2)$ jest równoważny współczynnikowi korelacji ρV pary wektorów losowych α i β .

Współczynnik ρV możemy estymować na podstawie próby. Niech

$$\mathbf{W}_{\alpha\alpha} = \mathbf{A}\mathbf{A}', \quad \mathbf{W}_{\beta\beta} = \mathbf{B}\mathbf{B}'.$$

Wtedy

$$(\widehat{\rho V})^2 = (rV)^2 = \frac{\text{tr}(\mathbf{W}_{\alpha\alpha}\mathbf{W}_{\beta\beta})}{\sqrt{\text{tr}(\mathbf{W}_{\alpha\alpha}^2)\text{tr}(\mathbf{W}_{\beta\beta}^2)}}.$$

Zauważmy, że problem testowania hipotezy zerowej

$$H_0: I\rho V = 0$$

przeciwko

$$H_1: I\rho V \neq 0$$

dla pary procesów losowych $\mathbf{X} \in \mathcal{L}_2^p(I_1)$ i $\mathbf{Y} \in \mathcal{L}_2^q(I_2)$ jest równoważny z problemem testowania hipotezy zerowej

$$H_0: \rho V = 0$$

przeciwko

$$H_1: \rho V \neq 0$$

dla pary wektorów losowych α and β .

Przy prawdziwości hipotezy zerowej oraz przy założeniu, że łączny rozkład wektorów losowych należy do klasy rozkładów eliptycznych, graniczny rozkład współczynnika nrV jest równy rozkładowi zmiennej losowej

$$\frac{1+k}{\text{tr}(\Sigma_{\alpha\alpha}^2)\text{tr}(\Sigma_{\beta\beta}^2)} \sum_{l=1}^{K_1+p} \sum_{m=1}^{K_2+q} \lambda_l \gamma_m Z_{lm}^2,$$

gdzie k jest kurtozą rozkładu eliptycznego, $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{K_1+p}$ są wartościami własnymi macierzy kowariancji $\Sigma_{\alpha\alpha}$, $\gamma_1 \geq \gamma_2 \geq \dots \geq \gamma_{K_2+q}$ są wartościami własnymi macierzy kowariancji $\Sigma_{\beta\beta}$, Z_{lm} są niezależnymi zmiennymi losowymi o rozkładach $N(0, 1)$.

Ponieważ testy asymptotyczne w wielu sytuacjach zawodzą, sensowne wydaje się użycie **testu permutacyjnego** (Josse i inni (2008)).

Funkcja charakterystyczna łącznego rozkładu procesów losowych \mathbf{X} i \mathbf{Y} ma postać:

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

gdzie

$$\langle \mathbf{l}, \mathbf{X} \rangle_p = \int_{I_1} \mathbf{l}'(s) \mathbf{X}(s) ds, \quad \langle \mathbf{m}, \mathbf{Y} \rangle_q = \int_{I_2} \mathbf{m}'(t) \mathbf{Y}(t) dt.$$

Funkcje charakterystyczne rozkładów brzegowych procesów \mathbf{X} i \mathbf{Y} są postaci:

$$\varphi_{\mathbf{X}}(\mathbf{l}) = \varphi_{\mathbf{X},\mathbf{Y}}(\mathbf{l}, \mathbf{0}), \quad \varphi_{\mathbf{Y}}(\mathbf{m}) = \varphi_{\mathbf{X},\mathbf{Y}}(\mathbf{0}, \mathbf{m}).$$

Korelacja odległościowa

Założmy, że $\mathbf{X} \in \mathcal{L}_2^p(I_1)$ oraz $\mathbf{Y} \in \mathcal{L}_2^q(I_2)$. Wówczas procesy te mogą mieć reprezentację:

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

gdzie $\alpha \in \mathbb{R}^{K_1+p}$, $\beta \in \mathbb{R}^{K_2+q}$, natomiast funkcje wektorowe l i m mogą być reprezentowane jako

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

gdzie $\lambda \in \mathbb{R}^{K_1+p}$, $\mu \in \mathbb{R}^{K_2+q}$. Stąd

$$\langle l, \mathbf{X} \rangle_p = \lambda' \alpha, \quad \langle m, \mathbf{Y} \rangle_q = \mu' \beta,$$

oraz

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

gdzie $\varphi_{\alpha, \beta}(\lambda, \mu)$ jest łączną funkcją charakterystyczną wektorów losowych α i β .

Bazując na idei **kowariancji odległościowej** pomiędzy dwoma wektorami losowymi (Székely i inni (2007)), możemy zdefiniować **kowariancję odległościową** pomiędzy procesami losowymi \mathbf{X} i \mathbf{Y} jako

$$\nu_{\mathbf{X}, \mathbf{Y}} = \nu_{\alpha, \beta},$$

gdzie

$$\nu_{\alpha, \beta}^2 = \frac{1}{C_{K_1+p} C_{K_2+q}} \int_{\mathbb{R}^{K_1+K_2+p+q}} \frac{|\varphi_{\alpha, \beta}(\boldsymbol{\lambda}, \boldsymbol{\mu}) - \varphi_{\alpha}(\boldsymbol{\lambda}) \varphi_{\beta}(\boldsymbol{\mu})|^2}{\|\boldsymbol{\lambda}\|_{K_1+p}^{K_1+p+1} \|\boldsymbol{\mu}\|_{K_2+q}^{K_2+q+1}} d\boldsymbol{\lambda} d\boldsymbol{\mu},$$

gdzie

$$C_r = \frac{\pi^{\frac{1}{2}(r+1)}}{\Gamma(\frac{1}{2}(r+1))}.$$

Funkcjonalna korelacja odległościowa pomiędzy procesami losowymi \mathbf{X} i \mathbf{Y} jest liczbą nieujemną zdefiniowaną jako

$$\mathcal{R}_{\mathbf{X},\mathbf{Y}} = \mathcal{R}_{\mathbf{Y},\mathbf{X}} = \frac{\nu_{\mathbf{X},\mathbf{Y}}}{\sqrt{\nu_{\mathbf{X},\mathbf{X}}\nu_{\mathbf{Y},\mathbf{Y}}}}$$

jeżeli $\nu_{\mathbf{X},\mathbf{X}}$ i $\nu_{\mathbf{Y},\mathbf{Y}}$ są dodatnie i zero w przeciwnym przypadku.
Dla rozkładów ze skończonymi pierwszymi momentami mamy

$$0 \leq \mathcal{R}_{\mathbf{X},\mathbf{Y}} \leq 1$$

i $\mathcal{R}_{\mathbf{X},\mathbf{Y}} = 0$ wtedy i tylko wtedy gdy \mathbf{X} i \mathbf{Y} są niezależne. $\mathcal{R}_{\mathbf{X},\mathbf{Y}}$ jest estymatorem zgodnym.

Musimy teraz oszacować $\mathcal{R}_{\mathbf{X}, \mathbf{Y}}$ na podstawie $(\mathbf{a}_k, \mathbf{b}_k)$, $k = 1, 2, \dots, n$.

Niech

$$\bar{\mathbf{a}} = \frac{1}{n} \sum_{i=1}^n \mathbf{a}_i, \quad \bar{\mathbf{b}} = \frac{1}{n} \sum_{i=1}^n \mathbf{b}_i,$$

$$\tilde{\mathbf{a}}_k = \mathbf{a}_k - \bar{\mathbf{a}}, \quad \tilde{\mathbf{b}}_k = \mathbf{b}_k - \bar{\mathbf{b}}, \quad k = 1, 2, \dots, n$$

oraz

$$\mathbf{A} = (a_{kl}), \quad \mathbf{B} = (b_{kl}), \\ \tilde{\mathbf{A}} = (A_{kl}), \quad \tilde{\mathbf{B}} = (B_{kl}),$$

gdzie

$$a_{kl} = \|\mathbf{a}_k - \mathbf{a}_l\|_{K_1+p}, \quad b_{kl} = \|\mathbf{b}_k - \mathbf{b}_l\|_{K_2+q}, \\ A_{kl} = \|\tilde{\mathbf{a}}_k - \tilde{\mathbf{a}}_l\|_{K_1+p}, \quad B_{kl} = \|\tilde{\mathbf{b}}_k - \tilde{\mathbf{b}}_l\|_{K_2+q}, \quad k, l = 1, 2, \dots, n.$$

Zachodzą związki

$$\tilde{\mathbf{A}} = \mathbf{H}\mathbf{A}\mathbf{H} \text{ oraz } \tilde{\mathbf{B}} = \mathbf{H}\mathbf{B}\mathbf{H},$$

gdzie

$$\mathbf{H} = \mathbf{I}_n - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^T$$

jest macierzą centrującą.

Niech $\tilde{\mathbf{A}} \circ \tilde{\mathbf{B}} = (A_{kl}B_{kl})$ będzie **iloczynem Hadamarda** macierzy $\tilde{\mathbf{A}}$ oraz $\tilde{\mathbf{B}}$.

Na podstawie wyników Székely'ego i innych (2007), mamy

$$\hat{\nu}_{\mathbf{X},\mathbf{Y}}^2 = \frac{1}{n^2} \sum_{k,l=1}^n A_{kl} B_{kl}.$$

Funkcjonalna korelacja odległościowa z próby jest zatem zdefiniowana jako

$$\hat{\mathcal{R}}_{\mathbf{X},\mathbf{Y}} = \hat{\mathcal{R}}_{\mathbf{a},\mathbf{b}} = \frac{\hat{\nu}_{\mathbf{a},\mathbf{b}}}{\sqrt{\hat{\nu}_{\mathbf{a},\mathbf{a}}\hat{\nu}_{\mathbf{b},\mathbf{b}}}}$$

jeżeli $\hat{\nu}_{\mathbf{a},\mathbf{a}}$ oraz $\hat{\nu}_{\mathbf{b},\mathbf{b}}$ są dodatnie i zero w przeciwnym przypadku.

Problem testowania niezależności pomiędzy procesami losowymi $\mathbf{X} \in \mathcal{L}_2^p(I_1)$ oraz $\mathbf{Y} \in \mathcal{L}_2^q(I_2)$ jest równoważny problemowi testowania $H_0: \mathcal{R}_{\mathbf{X}, \mathbf{Y}} = 0$. Székely i inni (2007) pokazali, że przy prawdziwości H_0 , $n\hat{\mathcal{R}}_{\mathbf{X}, \mathbf{Y}}$ zbiega do

$$\sum_{j=1}^{\infty} \nu_j Z_j^2,$$

gdzie Z_j są niezależnymi zmiennymi losowymi o rozkładzie $N(0, 1)$, oraz ν_j zależy od rozkładu (\mathbf{X}, \mathbf{Y}) . W praktyce używamy jednak **testów permutacyjnych** w celu zbadania istotności korelacji odległościowej (Josse & Holmes (2016)).

Niech $\mathbf{X} \in \mathbb{R}^p$ i $\mathbf{Y} \in \mathbb{R}^q$ będą **wektorami losowymi**. Oznaczmy przez $P_{\mathbf{X}}$, $P_{\mathbf{Y}}$ i $P_{\mathbf{X},\mathbf{Y}}$ miary probabilistyczne wektorów \mathbf{X} , \mathbf{Y} i (\mathbf{X}, \mathbf{Y}) odpowiednio określonych na przestrzeniach \mathbb{R}^p , \mathbb{R}^q i $\mathbb{R}^p \times \mathbb{R}^q$. Niech $S_{\mathbf{X}} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ będzie **próbą** z rozkładu $P_{\mathbf{X}}$. Niech

$$k : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$$

będzie rzeczywistą funkcją ciągłą, zwaną **jądrem** i niech $\mathbf{K} = (k_{ij})$, gdzie $k_{ij} = k(\mathbf{x}_i, \mathbf{x}_j)$ będzie **macierzą jądrową**, $i, j = 1, \dots, n$. Dalej będziemy przyjmować, że k jest **jądrem gaussowskim**, tj.

$$k(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\lambda \|\mathbf{x}_i - \mathbf{x}_j\|^2),$$

$$\lambda > 0.$$

Dla każdej macierzy jądrowej $\mathbf{K} \in \mathbb{R}^{n \times n}$, przez **jądrową macierz scentrowaną** rozumiemy macierz postaci

$$\tilde{\mathbf{K}} = \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right) \mathbf{K} \left(\mathbf{I}_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T \right),$$

gdzie $\mathbf{1}_n \in \mathbb{R}^{n \times 1}$ oznacza wektor złożony z samych jedynek. Macierz $\tilde{\mathbf{K}}$ jest macierzą nieujemnie określoną.

Niech

$$\langle \mathbf{A}, \mathbf{B} \rangle_F = \text{tr}(\mathbf{A}^T \mathbf{B})$$
$$\|\mathbf{A}\|_F = (\langle \mathbf{A}, \mathbf{A} \rangle_F)^{1/2}$$

będą **iloczynem i normą Frobeniusa** macierzy $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$. Dla każdych macierzy jądrowych opartych na próbach $S_{\mathbf{x}} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ oraz $S_{\mathbf{y}} = \{\mathbf{y}_1, \dots, \mathbf{y}_n\}$ zachodzą równości:

$$\langle \tilde{\mathbf{K}}_{\mathbf{x}}, \tilde{\mathbf{K}}_{\mathbf{y}} \rangle_F = \langle \mathbf{K}_{\mathbf{x}}, \mathbf{K}_{\mathbf{y}} \rangle_F = \langle \tilde{\mathbf{K}}_{\mathbf{x}}, \mathbf{K}_{\mathbf{y}} \rangle_F .$$

Gretton i inni (2005) zdefiniowali współczynnik **HSIC** (ang. **Hilbert-Schmidt Independence Criterion**) postaci

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

gdzie

$$\mathbf{C}_{\mathbf{X}, \mathbf{Y}} = E_{\mathbf{X}, \mathbf{Y}}[\langle \varphi(\mathbf{X}) - E_{\mathbf{X}}[\varphi(\mathbf{X})], \psi(\mathbf{Y}) - E_{\mathbf{Y}}[\psi(\mathbf{Y})] \rangle_{\mathcal{H}}]$$

jest operatorem **kowariancji** między wektorami \mathbf{X} i \mathbf{Y} , będącym uogólnieniem macierzy kowariancji dla wektorów losowych oraz $\|\mathbf{A}\|_{HS} = \text{tr}(\mathbf{A}^T \mathbf{A})$.

W przestrzeniach rzeczywistych norma Hilberta-Schmidta nazywana jest normą Frobeniusa.

Dla jądra gaussowskiego $\|\mathbf{C}_{\mathbf{X},\mathbf{Y}}\|_{HS} = 0$ wtedy i tylko wtedy, gdy \mathbf{X} i \mathbf{Y} są niezależne.

Estymator z próby $S_{\mathbf{X},\mathbf{Y}}$ pobranej z rozkładu $P_{\mathbf{X},\mathbf{Y}}$ współczynnika HSIC ma postać

$$\text{HSIC}(S_{\mathbf{X},\mathbf{Y}}) = \frac{1}{n^2} \langle \tilde{\mathbf{K}}_{\mathbf{X}}, \tilde{\mathbf{K}}_{\mathbf{Y}} \rangle_F .$$

Cortes i inni (2012) zdefiniowali **współczynnik zgodności między macierzami jądrowymi (KTA – ang. Kernel Target Alignment)** \mathbf{K}_x i \mathbf{K}_y postaci

$$\rho(\mathbf{K}_x, \mathbf{K}_y) = \frac{\langle \tilde{\mathbf{K}}_x, \tilde{\mathbf{K}}_y \rangle_F}{\|\tilde{\mathbf{K}}_x\|_F \cdot \|\tilde{\mathbf{K}}_y\|_F},$$

gdzie $\rho(\mathbf{K}_x, \mathbf{K}_y) \in [0, 1]$. Porównując $\text{HSIC}(S_{x,y})$ z $\rho(\mathbf{K}_x, \mathbf{K}_y)$ widzimy, że współczynnik zgodności $\rho(\mathbf{K}_x, \mathbf{K}_y)$ jest po prostu znormalizowaną wersją współczynnika $\text{HSIC}(S_{x,y})$.

Współczynnik zgodności ρ jest podstawą do weryfikacji hipotezy o niezależności wektorów \mathbf{X} i \mathbf{Y} . Wprowadzimy **dwie hipotezy**:

$H_0: P_{\mathbf{X},\mathbf{Y}} = P_{\mathbf{X}}P_{\mathbf{Y}}$, że \mathbf{X} oraz \mathbf{Y} są niezależne i $H_1: \neg H_0$ są zależne.

Zhang i inni (2011) pokazali, że przy prawdziwości H_0 statystyka $\text{HSIC}(S_{\mathbf{x},\mathbf{y}})$ ma ten sam rozkład asymptotyczny co statystyka

$$Z = \frac{1}{n^2} \sum_{i,j=1}^n \lambda_{\mathbf{x},i} \lambda_{\mathbf{y},j} Z_{ij}^2,$$

gdzie Z_{ij}^2 są niezależnymi zmiennymi losowymi o rozkładzie χ_1^2 , $\lambda_{\mathbf{x},i}$ jest i -tą nieujemną wartością własną macierzy $\tilde{\mathbf{K}}_{\mathbf{x}}$ oraz $\lambda_{\mathbf{y},j}$ jest j -tą nieujemną wartością własną macierzy $\tilde{\mathbf{K}}_{\mathbf{y}}$. Na tej podstawie został opracowany test hipotezy H_0 przeciwko H_1 .

Z drugiej strony **Doran i inni (2014)** opracowali test permutacyjny badania niezależności macierzy jądrowych. Dokonałiśmy jego adaptacji do naszej sytuacji.

Niech \mathbf{X}_s i \mathbf{Y}_t będą procesami losowymi o wcześniej opisanej reprezentacji.
Dla procesu \mathbf{X}_s jądro gaussowskie jest równe

$$k_{\mathbf{X}_s}(\mathbf{x}(s), \mathbf{x}'(s)) = \exp(-\lambda_1 \|\mathbf{x}(s) - \mathbf{x}'(s)\|^2).$$

Ale

$$\begin{aligned} \|\mathbf{x}(s) - \mathbf{x}'(s)\|^2 &= \int_{I_1} (\mathbf{x}(s) - \mathbf{x}'(s))^T (\mathbf{x}(s) - \mathbf{x}'(s)) ds \\ &= (\boldsymbol{\alpha} - \boldsymbol{\alpha}')^T \left(\int_{I_1} \boldsymbol{\Phi}_1^T(s) \boldsymbol{\Phi}_1(s) ds \right) (\boldsymbol{\alpha} - \boldsymbol{\alpha}') \\ &= (\boldsymbol{\alpha} - \boldsymbol{\alpha}')^T (\boldsymbol{\alpha} - \boldsymbol{\alpha}') = \|\boldsymbol{\alpha} - \boldsymbol{\alpha}'\|^2. \end{aligned}$$

Stąd

$$k_{\mathbf{X}_s}(\mathbf{x}(s), \mathbf{x}'(s)) = k_{\boldsymbol{\alpha}}(\boldsymbol{\alpha}, \boldsymbol{\alpha}').$$

Podobnie

$$k_{\mathbf{Y}_t}(\mathbf{y}(t), \mathbf{y}'(t)) = k_{\beta}(\beta, \beta').$$

A to oznacza, że

$$\rho(K_{\mathbf{X}_s}, K_{\mathbf{Y}_t}) = \rho(K_{\alpha}, K_{\beta}) = \frac{\langle \tilde{K}_{\alpha}, \tilde{K}_{\beta} \rangle_F}{\|\tilde{K}_{\alpha}\|_F \|\tilde{K}_{\beta}\|_F}.$$

Zatem współczynnik zgodności procesów losowych \mathbf{Y}_s i \mathbf{Y}_t można wyrazić przez współczynnik zgodności wektorów losowych α i β występujących w ich reprezentacji.

Zauważmy również, że hipoteza zerowa $H_0 : \mathbf{X}_s \perp \mathbf{Y}_t$ niezależności procesów losowych \mathbf{Y}_s i \mathbf{Y}_t jest równoważna hipotezie $H_0 : \boldsymbol{\alpha} \perp \boldsymbol{\beta}$ niezależności wektorów losowych $\boldsymbol{\alpha}$ i $\boldsymbol{\beta}$ występujących w reprezentacji tych procesów. Możemy zatem posłużyć się testami przedstawionymi wcześniej.

Jako przykład posłużyły nam dane Światowego Forum Ekonomicznego (<http://www.weforum.org>) dotyczące wskaźników socjo-ekonomicznych 38 europejskich krajów w latach 2008-2015.

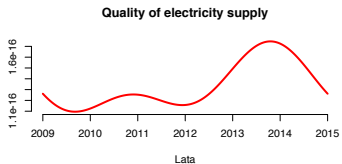
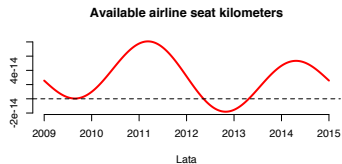
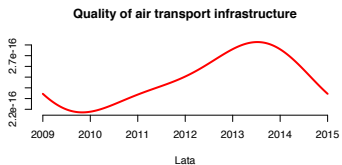
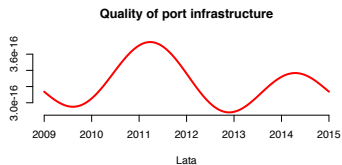
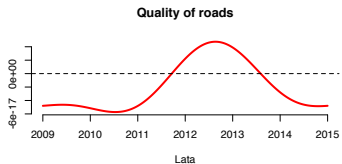
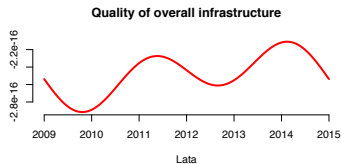
1	Albania (AL)	14	Greece (GR)	27	Poland (PL)
2	Austria (AT)	15	Hungary (HU)	28	Portugal (PT)
3	Belgium (BE)	16	Iceland (IS)	29	Romania (RO)
4	Bosnia and Herzegovina (BA)	17	Ireland (IE)	30	Russian Federation (RU)
5	Bulgaria (BG)	18	Italy (IT)	31	Serbia (XS)
6	Croatia (HR)	19	Latvia (LV)	32	Slovak Republic (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)		

Kraje użyte w analizie.

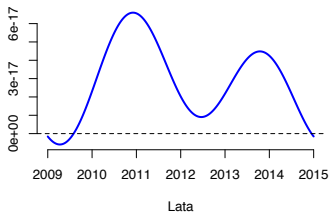
X (Infrastructure)	Y (Market size)
1. Quality of overall infrastructure	1. Domestic market size index
2. Quality of roads	2. Foreign market size index
3. Quality of port infrastructure	3. GDP valued at PPP
4. Quality of air transport infrastructure	4. Exports as a percentage of GDP
5. Available airline seat kilometers	
6. Quality of electricity supply	

Zmienne użyte w analizie

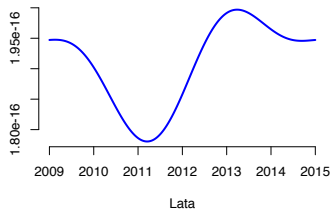
W pierwszym kroku wygładziliśmy dane. Została użyta baza **Fouriera**. Ponieważ liczba punktów czasowych była niewielka ($J = 7$), dla każdej zmiennej ustaliliśmy maksymalną **wielkość bazy na 5**. Na kolejnych dwóch wykresach przedstawione są średnie wartości analizowanych zmiennych dla danych funkcjonalnych.



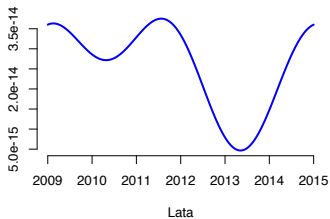
Domestic market size index



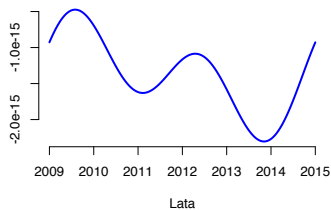
Foreign market size index



GDP valued at PPP









Exports as a percentage of GDP









W następnym kroku zostały policzone opisane wcześniej współczynniki i p -wartości testów permutacyjnych dla ich istotności.

	Wartość współczynnika	p wartość
ρV	0.941	0.0001
dcor	0.965	0.0001
HSIC	0.974	0.0046
KTA	1.000	0.0046

Zgodnie z oczekiwaniami wszystkie współczynniki są bliskie 1. Jednakże ich wysokie wartości nie oznaczają jeszcze statystycznej istotności relacji pomiędzy grupami zmiennych. Dopiero małe p -wartości wskazują na ich istotność.

-  CORTES C., MOHRI M., ROSTAMIZADEH A. (2012): Algorithms for learning kernels based on centered alignment. *Journal of Machine Learning Research* 13, 795–828.
-  DORAN G., MUANDET K., ZHANG K., SCHÖLKOPF B. (2014): A Permutation-based kernel conditional independence test. *Proceedings of the 30th Conference on Uncertainty in Artificial Intelligence* 132–141.
-  GÓRECKI, T., KRZYŚKO, M., WOŁYŃSKI, W. (2017): Correlation analysis for multivariate functional data. *Studies in Classification Data Analysis, and Knowledge Organization: Data Science* 243–258.
-  GÓRECKI, T., KRZYŚKO, M., WOŁYŃSKI, W. (2017): Independence test and canonical correlation analysis based on the alignment between kernel matrices for multivariate functional data – submitted.
-  GRETTON A., BOUSQUET O., SMOLA A., SCHÖLKOPF B. (2005): Measuring statistical dependence with Hilbert-Schmidt norms. *Lecture Notes in Computer Science* 3734, 63–77.
-  HORVÁTH, L., KOKOSZKA, P. (2012): *Inference for Functional Data with Applications*, Springer.

-  JOSSE, J., HOLMES, S. (2016): Measuring multivariate association and beyond. *Statistics Surveys* 10, 132–167.
-  JOSSE, J., PAGES, J., HUSSON, F. (2008): Testing the significance of the *RV* coefficient. *Computational Statistics and Data Analysis* 53, 82–91.
-  LEURGANS, S.E., MOYEED, R.A., SILVERMAN, B.W. (1993): Canonical correlation analysis when the data are curves. *Journal of the Royal Statistical Society: Series B* 55, 725–740.
-  ROBERT, P., ESCOUFIER, Y. (1976): A unifying tool for linear multivariate statistical methods: the *RV*-coefficient. *Applied Statistics* 25, 257–265.
-  SZÉKELY, G.J., RIZZO, M.L., BAKIROV, N.K. (2007): Measuring and testing dependence by correlation of distances. *The Annals of Statistics* 35(6), 2769–2794.
-  ZHANG K., PETERS J., JANZING D., SCHÖLKOPF B., (2011): Kernel-based conditional independence test and application in causal discovery. *Proceedings of the 27th Conference on Uncertainty in Artificial Intelligence* 804–813.