

Analiza kanoniczna dla danych funkcjonalnych

Waldemar Wołyński

Wydział Matematyki i Informatyki UAM Poznań

Poznań, 18 października 2017

WSTĘP

Dane pochodzą z **EUROSTATU** i dotyczą spożycia indywidualnego według celu w 27 wybranych krajach europejskich w latach 2000-2010.

Klasyfikacja Spożycia Indywidualnego według Celu (COICOP) obejmuje 12 grup wydatków:

- 1 Artykuły żywnościowe i napoje bezalkoholowe,
- 2 Napoje alkoholowe i tytoń,
- 3 Odzież i obuwie,
- 4 Mieszkanie, woda, elektryczność, gaz i inne paliwa,
- 5 Wyposażenie wnętrza, sprzęty domowe i bieżące utrzymanie budynku,
- 6 Opieka zdrowotna,
- 7 Transport,
- 8 Łączność,
- 9 Wypoczynek i kultura,
- 10 Szkolnictwo,
- 11 Hotele, kawiarnie i restauracje,
- 12 Różne towary i usługi.

Wydatki uwzględnione w grupie 2 zostały podzielone na wydatki na napoje alkoholowe (cecha Y_1) oraz wydatki na wyroby tytoniowe (cecha Y_2). Jedenaście pozostałych grup wydatków przyjęto jako cechy X_1, X_2, \dots, X_{11} .

Powstały w ten sposób dwie nowe grupy cech:

- Grupa I: cechy Y_1, Y_2 ,
- Grupa II: cechy X_1, \dots, X_{11} .

Dysponujemy wartościami każdej z cech dla 27 jednostek (kraje) i 11 momentów czasowych (lata 2000-2010).

Jak zbadać korelację pomiędzy badanymi grupami cech?

CZĘŚĆ 1

Niech $X \in \mathbf{R}$ i $Y \in \mathbf{R}$ będą zmiennymi losowymi. Ponadto założmy, że $E(X) = E(Y) = 0$.

Współczynnikiem korelacji (liniowej Pearsona) nazywamy liczbę:

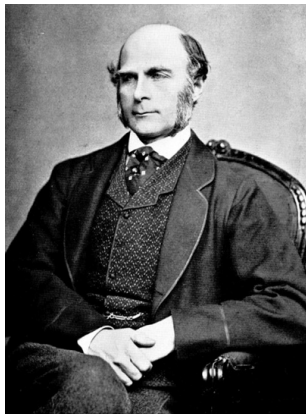
$$\rho_{X,Y} = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}.$$

Własności:

- 1 $-1 \leq \rho_{X,Y} \leq 1$,
- 2 Jeżeli $\rho_{X,Y} = 0$, to brak korelacji liniowej pomiędzy zmiennymi X i Y ,
- 3 Jeżeli $|\rho_{X,Y}| = 1$, to istnieje liniowa zależność pomiędzy zmiennymi X i Y .



Karl Pearson
1857 - 1936



Francis Galton
1822 - 1911

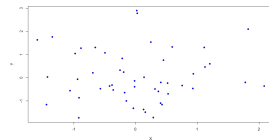
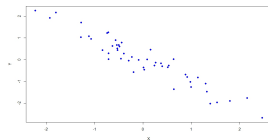
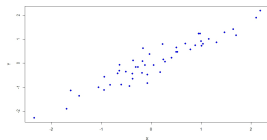
Niech

$\mathbf{x} = (x_1, \dots, x_n)'$ - obserwacje (wycentrowane) zmiennej X ,

$\mathbf{y} = (y_1, \dots, y_n)'$ - obserwacje (wycentrowane) zmiennej Y .

Współczynnik korelacji (z próby):

$$r_{X,Y} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i^2}}.$$



Wykres rozrzutu (diagram korelacyjny)

Niech $\mathbf{X} \in \mathbf{R}^p$ i $\mathbf{Y} \in \mathbf{R}^q$ będą wektorami losowymi losowymi. Ponadto załóżmy, że $E(\mathbf{X}) = E(\mathbf{Y}) = \mathbf{0}$.

Założenie to nie powoduje utraty ogólności rozważań, ponieważ korelacje kanoniczne wyznaczone są jedynie na podstawie macierzy kowariancji wektorów \mathbf{X} i \mathbf{Y} postaci

$$\Sigma = \begin{bmatrix} \Sigma_{XX} & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_{YY} \end{bmatrix}.$$

Zmienne kanoniczne wektorów losowych \mathbf{X} i \mathbf{Y} definiujemy następująco:

$$U = \langle \mathbf{u}, \mathbf{X} \rangle = \sum_{i=1}^p u_i X_i,$$

$$V = \langle \mathbf{v}, \mathbf{Y} \rangle = \sum_{j=1}^q v_j Y_j,$$

gdzie \mathbf{u} i \mathbf{v} są wektorami wagowymi dobranymi tak, aby zmaksymalizować współczynnik korelacji $\rho_{U,V}$ przy dodatkowym warunku ograniczającym $\text{Var}(U) = \text{Var}(V) = 1$.

Maksymalny współczynnik korelacji $\rho_{U,V}$ nazywamy **współczynnikiem korelacji kanonicznej**.

Zagadnienie maksymalizacji współczynnika korelacji $\rho_{U,V}$ sprowadza się do rozwiązania równań:

$$\begin{aligned} & |\Sigma_{XX}^{-1}\Sigma_{XY}\Sigma_{YY}^{-1}\Sigma_{YX} - \rho^2 I_p|, \\ & |\Sigma_{YY}^{-1}\Sigma_{YX}\Sigma_{XX}^{-1}\Sigma_{XY} - \rho^2 I_q|. \end{aligned}$$

Liczba niezerowych pierwiastków tych równań jest równa $s = \text{rzęd}(\Sigma_{XY})$.

Otrzymujemy s par nieskorelowanych **zmiennych kanonicznych** (U_i, V_i) , $i = 1, \dots, s$.

Pierwiastki kwadratowe z $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_s^2$ są **współczynnikami korelacji kanonicznej**.

Wektory własne odpowiadające odpowiadające pierwiastkom równań są **wektorami wagowymi**.

O wkładzie poszczególnych składowych wektorów losowych \mathbf{X} oraz \mathbf{Y} w budowę zmiennych kanonicznych można wnioskować na podstawie wektorów wagowych \mathbf{u} oraz \mathbf{v} .

Wkład j -tej składowej wektora losowego \mathbf{X} w budowę zmiennej kanonicznej U określamy wzorem:

$$\frac{|u_j|}{\sum_{i=1}^p |u_i|} \times 100\%, \quad j = 1, \dots, p.$$

Analogicznie wnioskujemy dla zmiennej kanonicznej V .



Harold Hotelling
1895 - 1973

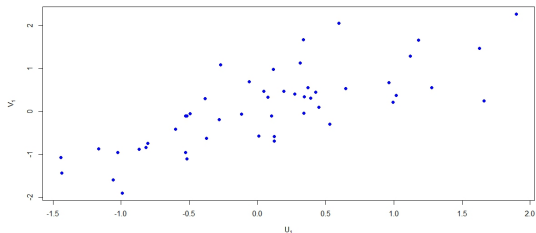
Niech

$\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_n)'$ - obserwacje (wycentrowane) wektora X ,

$\mathbf{y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)'$ - obserwacje (wycentrowane) wektora Y .

Wtedy

$$\hat{\Sigma}_{XX} = \frac{1}{n} \mathbf{x}'\mathbf{x}, \quad \hat{\Sigma}_{YY} = \frac{1}{n} \mathbf{y}'\mathbf{y}, \quad \hat{\Sigma}_{XY} = \frac{1}{n} \mathbf{x}'\mathbf{y} = \hat{\Sigma}'_{YX}.$$



Rzut na pierwszą parę zmiennych kanonicznych

- 1 Więcej niż dwie grupy zmiennych.
- 2 Nieliniowe zmienne kanoniczne.
- 3 Dane geograficznie ważone.
- 4 Dane czasowo-przestrzenne.

Niech $\mathbf{X} \in L_2^p(I)$ i $\mathbf{Y} \in L_2^q(I)$ są wektorowymi procesami losowymi, gdzie $L_2(I)$ jest **przestrzenią Hilberta funkcji całkowlanych z kwadratem na przedziale I** . Ponadto założmy, że $E(\mathbf{X}(t)) = E(\mathbf{Y}(t)) = \mathbf{0}$, $t \in I$.

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

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

Funkcjonalne zmienne kanoniczne procesów losowych \mathbf{X} i \mathbf{Y} definiujemy następująco:

$$U = \langle \mathbf{u}, \mathbf{X} \rangle = \int_I \mathbf{u}'(t) \mathbf{X}(t) dt,$$

$$V = \langle \mathbf{v}, \mathbf{Y} \rangle = \int_I \mathbf{v}'(t) \mathbf{Y}(t) dt,$$

gdzie $\mathbf{u} \in L_2^p(I)$ i $\mathbf{v} \in L_2^q(I)$ są funkcjami wagowymi dobranymi tak, aby zmaksymalizować współczynnik korelacji $\rho_{U,V}$ przy dodatkowym warunku ograniczającym $\text{Var}(U) = \text{Var}(V) = 1$.

Maksymalny współczynnik korelacji $\rho_{U,V}$ nazywamy **funkcjonalnym współczynnikiem korelacji kanonicznej**.

Funkcjonalne korelacje kanoniczne - problem!!!

Maksymalizacja funkcjonalnego współczynnika korelacji kanonicznej nie daje zadawalających wyników. Dowolnie wybierając funkcję wagową \mathbf{u} możemy znaleźć funkcję wagową \mathbf{v} taką, aby funkcjonalny współczynnik korelacji kanonicznej $\rho_{\mathbf{X},\mathbf{Y}}$ był równy 1.

Rozwiązaniem jest modyfikacja warunków ograniczających poprzez wprowadzenie dodatkowej funkcji kary:

$$PEN_2(\mathbf{u}) = \int_I \left(\frac{\partial^2 \mathbf{u}(t)}{\partial t^2} \right)' \frac{\partial^2 \mathbf{u}(t)}{\partial t^2} dt$$

oraz

$$PEN_2(\mathbf{v}) = \int_I \left(\frac{\partial^2 \mathbf{v}(t)}{\partial t^2} \right)' \frac{\partial^2 \mathbf{v}(t)}{\partial t^2} dt.$$

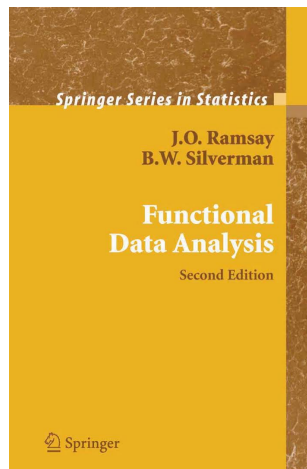
Współczynnik kary PEN_2 służy do oszacowania gładkości funkcji. Kwadrat drugiej pochodnej funkcji w momencie t jest nazywany jej **krzywizną** w punkcie t . Naturalną miarą krzywizny funkcji jest zatem scałkowany kwadrat drugiej pochodnej tej funkcji.

Po wprowadzeniu kary warunki ograniczające przyjmują postać:

$$\text{Var}(U) + \lambda PEN_2(\mathbf{u}) = 1,$$

$$\text{Var}(V) + \lambda PEN_2(\mathbf{v}) = 1.$$

Besse (1979);
Ramsay & Silverman (1997,
2005);
Horváth & Kokoszka (2012).

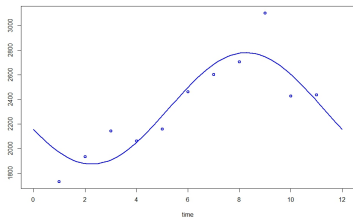


W dalszym ciągu zakładamy, że każda składowa X_g procesu \mathbf{X} i każda składowa Y_h procesu \mathbf{Y} może być reprezentowana przez **skończoną liczbę funkcji bazowych** $\{\varphi_e\}$ i $\{\varphi_f\}$:

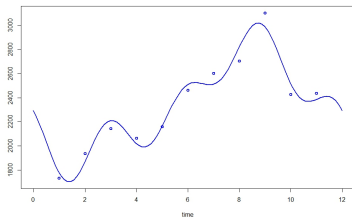
$$X_g(t) = \sum_{e=0}^{E_g} \alpha_{ge} \varphi_e(t), t \in I, g = 1, 2, \dots, p,$$

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

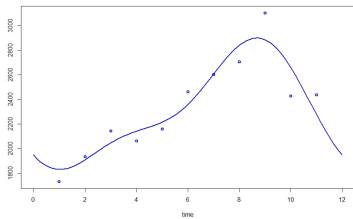
Stopień gładkości funkcji X_g i Y_h zależy od wartości E_g i F_h (małe wartości dają mniejsze wygładzenie funkcji).



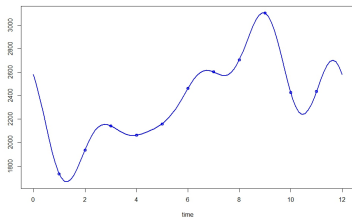
$E = 2$



$E = 8$



$E = 4$



$E = 10$

Wprowadźmy następujące oznaczenia:

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

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

$$\Phi_1(t) = \begin{bmatrix} \varphi'_{E_1}(t) & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \varphi'_{E_2}(t) & \dots & \mathbf{0} \\ \dots & \dots & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \dots & \varphi'_{E_p}(t) \end{bmatrix},$$

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

gdzie $\varphi_{E_1}, \dots, \varphi_{E_p}$ i $\varphi_{F_1}, \dots, \varphi_{F_q}$ są wektorami ortonormalnych funkcji bazowych przestrzeni $L_2(I)$.

Używając notacji macierzowej procesy \mathbf{X} i \mathbf{Y} mają następującą reprezentację

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

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

Ponadto, dla wektorów $\boldsymbol{\alpha}$ i $\boldsymbol{\beta}$ mamy:

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

oraz

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

- Wielomianowa;
- Fouriera;
- Funkcji sklejanych (spline);
- Funkcji falkowych (wavelet).

Baza Fouriera zdefiniowana jest w następujący sposób:

$$\varphi_0(t) = \frac{1}{\sqrt{T}},$$

$$\varphi_{2k-1}(t) = \sqrt{\frac{2}{T}} \sin \frac{2\pi kt}{T},$$

$$\varphi_{2k}(t) = \sqrt{\frac{2}{T}} \cos \frac{2\pi kt}{T},$$

gdzie $t \in I = [0, T]$, $k = 1, 2, \dots$

Podstawowa zależność

$$\mathbf{X} \longleftrightarrow \boldsymbol{\alpha}$$

$$\text{proces} \longleftrightarrow \text{wektor}$$

Jeżeli $\mathbf{X} \in \mathcal{L}_2^p(I)$ i $\mathbf{Y} \in \mathcal{L}_2^q(I)$, to

$$\Sigma_{XX}(s, t) = \Phi_1(s) \Sigma_{\alpha\alpha} \Phi_1'(t),$$

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

$$\Sigma_{YY}(s, t) = \Phi_2(s) \Sigma_{\beta\beta} \Phi_2'(t).$$

Niech $\mathbf{u} \in \mathcal{L}_2^p(I)$ i $\mathbf{v} \in \mathcal{L}_2^q(I)$.

Wtedy

$$\mathbf{u}(t) = \Phi_1(t)\mathbf{u}, \quad \mathbf{v}(t) = \Phi_2(t)\mathbf{v},$$

gdzie $t \in I$.

Zatem

$$\text{Cov}(U, V) = \mathbf{u}'\Sigma_{\alpha\beta}\mathbf{v},$$

$$\text{Var}(U) = \mathbf{u}'\Sigma_{\alpha\alpha}\mathbf{u}, \quad \text{Var}(V) = \mathbf{v}'\Sigma_{\beta\beta}\mathbf{v},$$

$$\text{PEN}_2(u) = \mathbf{u}'\mathbf{R}_1\mathbf{u}, \quad \text{PEN}_2(v) = \mathbf{v}'\mathbf{R}_2\mathbf{v},$$

gdzie

$$\mathbf{R}_i = \int_I \left(\frac{\partial^2 \Phi_i(t)}{\partial t^2} \right)' \frac{\partial^2 \Phi_i(t)}{\partial t^2} dt, \quad i = 1, 2.$$

Funkcjonalne zmienne kanoniczne U i V dobieramy tak, zmaksymalizować

$$\text{Cov}(U, V) = \mathbf{u}'\boldsymbol{\Sigma}_{\alpha\beta}\mathbf{v},$$

przy warunkach

$$\mathbf{u}'(\boldsymbol{\Sigma}_{\alpha\alpha} + \lambda\mathbf{R}_1)\mathbf{u} = \mathbf{v}'(\boldsymbol{\Sigma}_{\beta\beta} + \lambda\mathbf{R}_2)\mathbf{v} = 1.$$

Procesy \mathbf{X} i \mathbf{Y} są obserwowane 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 t_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 (t_j, x_{gj}) oraz z qJ par (t_j, y_{hj}) . Współczynniki α_i i β_i są estymowane **metodą najmniejszych kwadratów**. Oznaczmy te estymatory przez \mathbf{a}_i i \mathbf{b}_i , $i = 1, 2, \dots, n$. W rezultacie procesu transformacji otrzymujemy **dane funkcjonalne** postaci:

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

gdzie $t \in I$, $i = 1, 2, \dots, n$

Niech

$$\mathbf{A} = (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n)' \quad \mathbf{B} = (\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n)'.$$

Wtedy

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

Zatem dla funkcjonalnego współczynnika korelacji kanonicznej mamy:

$$r_{X,Y} = \max_{\mathbf{u}, \mathbf{v}} \mathbf{u}' \hat{\Sigma}_{\alpha\beta} \mathbf{v},$$

przy warunku

$$\mathbf{u}'(\hat{\Sigma}_{\alpha\alpha} + \lambda \mathbf{R}_1) \mathbf{u} = \mathbf{v}'(\hat{\Sigma}_{\beta\beta} + \lambda \mathbf{R}_2) \mathbf{v} = 1.$$

O wkładzie poszczególnych składowych wektorowych procesów losowych \mathbf{X} oraz \mathbf{Y} w budowę zmiennych kanonicznych można wnioskować na podstawie wektorowych funkcji wagowych \mathbf{u} oraz \mathbf{v} .

Niech P_j będzie polem pod modułem funkcji u_j na przedziale I . Wkład j -tej składowej procesu losowego \mathbf{X} w budowę zmiennej kanonicznej U określamy wzorem:

$$P_j^* = \frac{P_j}{\sum_{i=1}^p P_i} \times 100\%, \quad j = 1, \dots, p.$$

Analogicznie wnioskujemy dla zmiennej kanonicznej V .

CZĘŚĆ 2

Dwie grupy cech:

- Grupa I: cechy Y_1, Y_2 - wydatki na napoje alkoholowe i wyroby tytoniowe,
- Grupa II: cechy X_1, \dots, X_{11} - wydatki na pozostałe artykuły konsumpcyjne.

Dysponujemy wartościami każdej z cech dla 27 jednostek (kraje) i 11 momentów czasowych (lata 2000-2010).

Przykład

W celu ujednoczenia wartości rozpatrywanych cech, które mają różne przedziały zmienności, przeprowadzono ich **unitaryzację**.

Niech x_{kij} będzie wartością cechy X_k zaobserwowaną w i -tym kraju oraz j -tym roku.

Wówczas zunitaryzowana wartość z_{kij} wartości x_{kij} ma postać:

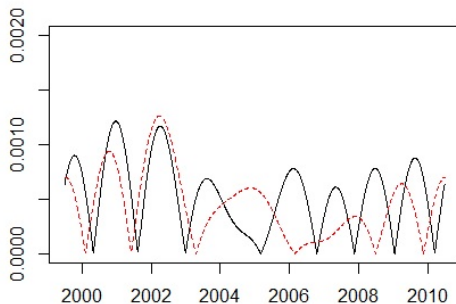
$$z_{kij} = \frac{x_{kij} - \min_i x_{kij}}{r_{kj}}$$

gdzie

$$r_{kj} = \max_i x_{kij} - \min_i x_{kij}$$

jest rozstępem k -tej cechy w j -tym roku.

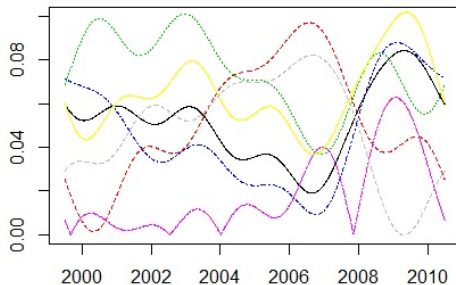
Po unitaryzacji rozstęp wszystkich cech we wszystkich latach jest stały i równy 1, natomiast wariancje i kowariancje w danym roku są proporcjonalne do wariancji i kowariancji cech bez unitaryzacji w tym roku.



Wykres funkcji wagowych pierwszej zmiennej kanonicznej dla procesu Y

j	P_j^*
1 - napoje alkoholowe	54,7
2 - wyroby tytoniowe	45,3

Wskaźniki P_j^* odpowiadające zmiennym procesowi Y .



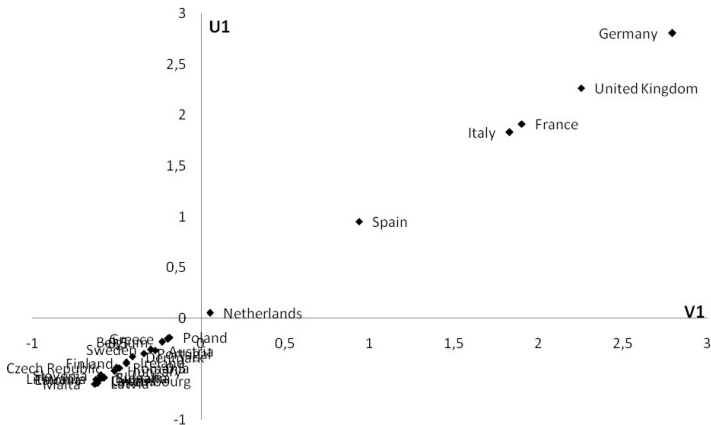
Wykres funkcji wagowych pierwszej zmiennej kanonicznej dla procesu X

Przykład

j	P_j^*
1 - artykuły żywnościowe i napoje bezalkoholowe	10,29
2 - odzież i obuwie	10,12
3 - mieszkanie, woda, elektryczność, gaz i inne paliwa	9,28
4 - wyposażenie wnętrza, sprzęty domowe i bieżące utrzymanie budynku	9,36
5 - opieka zdrowotna	7,23
6 - transport	9,96
7 - łączność	9,19
8 - wypoczynek i kultura	9,27
9 - szkolnictwo	8,16
10 - hotele, kawiarnie i restauracje	9,26
11 - różne towary i usługi	7,54

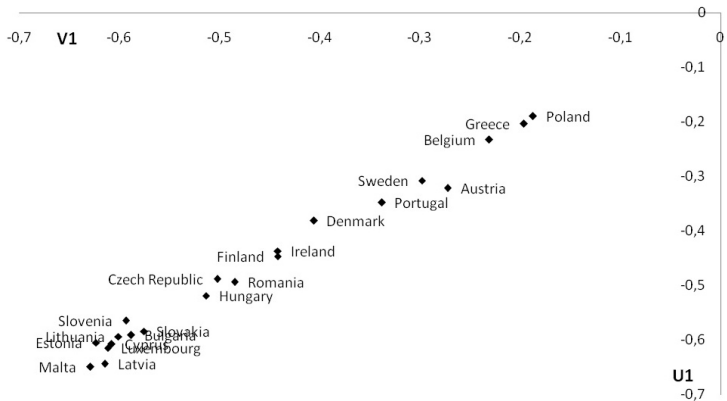
Wskaźniki P_j^* odpowiadające zmiennym procesu X .

Przykład



Położenie 27 krajów w układzie pierwszych zmiennych kanonicznych

Przykład



Powiększona lewa dolna ćwiartka wykresu

- 1 Górecki, T., Krzyśko, M., Waszak, Ł., Wołyński, W. (2016): Selected statistical methods of data analysis for multivariate functional data, *Statistical Papers*, Published online.
- 2 Horváth, L., Kokoszka, P. (2012): *Inference for Functional Data with Applications*, Springer.
- 3 Ramsay, J.O., Silverman, B.W. (2005): *Functional Data Analysis, Second Edition*, Springer.
- 4 Deręgowski, K., Krzyśko, M., Waszak, Ł., Wołyński, W. (2017): Zastosowanie funkcjonalnej analizy kanonicznej w badaniu zależności między wydatkami konsumpcyjnymi w europejskich gospodarstwach domowych, *Wiadomości statystyczne* **672**, 19-37.

KONIEC