

Metody redukcji wymiaru dla danych funkcyjnych

Mirosław Krzyśko i Łukasz Waszak

Wydział Matematyki i Informatyki
Uniwersytet im. Adama Mickiewicza, Poznań

STATYSTYKA – WIEDZA – ROZWÓJ. KONFERENCJA NAUKOWA
Z OKAZJI MIĘDZYNARODOWEGO ROKU STATYSTYKI
Łódź 17–18.10.2013

Dana jest cecha statystyczna Y . Jej wartości obserwujemy w J wybranych momentach czasowych, które nie muszą być równoodległe. Momenty czasowe oznaczamy będziemy przez t_j , a wartości cechy statystycznej Y zaobserwowane w tych momentach przez y_j , $j = 1, 2, \dots, J$. Wówczas nasze dane składają się z J par $\{t_i, y_i\}$. Dane dyskretne chcemy przekształcić w funkcję ciągłą $x(t)$, gdzie $t \in I$. Zbiór I jest zwartym przedziałem takim, że $t_i \in I$, dla $i = 1, 2, \dots, J$. Będziemy zakładać, że funkcja $x(t)$ jest postaci

$$x(t) = \sum_{k=0}^K c_k \varphi_k(t), t \in I,$$

gdzie $\{\varphi_k\}$ jest znanym układem bazowych funkcji ortonormalnych, natomiast c_0, \dots, c_K są nieznanymi współczynnikami, które należy wyestymować metodą najmniejszych kwadratów z danych $\{t_i, y_i\}$, $i=1,2,\dots,J$.

Przypomnijmy, że układ $\{\varphi_k\}$ w przestrzeni $L_2(I)$ jest układem ortonormalnym, gdy iloczyn skalarny każdych dwóch funkcji z tej przestrzeni jest równy:

$$\langle \varphi_i(t), \varphi_j(t) \rangle = \int_I \varphi_i(t) \varphi_j(t) dt = \delta_{ij}.$$

Przykładowymi układami bazowych funkcji ortonormalnych są: układ Fouriera oraz układ wielomianów Legendre'a. Ortonormalny układ Fouriera na przestrzeni $L_2([0, T])$ ma postać:

$$\varphi_0(t) = 1, \varphi_{2k-1}(t) = \sqrt{\frac{2}{T}} \sin \frac{2\pi kt}{T}, \varphi_{2k}(t) = \sqrt{\frac{2}{T}} \cos \frac{2\pi kt}{T},$$

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

Natomiast ortonormalny układ wielomianów Legendre'a na przestrzeni $L_2([-1, 1])$ ma postać:

$$\tilde{p}_k(x) = \sqrt{\frac{2k+1}{2}} p_k(x),$$

gdzie

$p_{k+1}(x) = \frac{1}{k+1}[(2k+1)xp_k(x) - kp_{k-1}(x)]$, $k \geq 1$, $p_1(x) = x$, $p_0(x) = 1$.
Każdy skończony przedział $[a, b]$ można przekształcić na przedział $[-1, 1]$ podstawiając:

$$x = \frac{2}{b-a}t - \frac{b+a}{b-a}, t \in [a, b], x \in [-1, 1].$$

Niech $\mathbf{y} = (y_1, y_2, \dots, y_J)'$, $\mathbf{c}_k = (c_0, \dots, c_K)'$ oraz niech Φ będzie macierzą $J \times (K + 1)$ o elementach $\varphi_k(t_j), j = 1, \dots, J, k = 0, \dots, K$. Poszukujemy zatem wektora $\hat{\mathbf{c}}$, dla którego funkcja celu

$$S(\mathbf{c}) = (\mathbf{y} - \Phi\mathbf{c})'(\mathbf{y} - \Phi\mathbf{c})$$

osiąga minimum. Poszukiwany wektor ma postać

$$\hat{\mathbf{c}} = (\Phi'\Phi)^{-1} \Phi'\mathbf{y}.$$

Stopień gładkości uzyskanej funkcji ciągłej $x(t)$ zależy od liczby wyrazów K kombinacji liniowej reprezentującej tę funkcję. Optymalna wartość K wybierana jest za pomocą bayesowskiego kryterium informacyjnego BIC postaci:

$$BIC(x(t)) = \ln \left(\frac{\mathbf{e}'\mathbf{e}}{2} \right) + K \left(\frac{\ln J}{J} \right),$$

gdzie $\mathbf{e} = (e_1, \dots, e_J)'$, $e_j = y_j - \sum_{k=0}^K \hat{c}_k \varphi_k(t_j)$, $j = 1, 2, \dots, J$. Za optymalną wartość K , przyjmujemy tę wartość, która minimalizuje $BIC(x(t))$. Kryterium BIC mierzy dokładność dopasowania przyjętego modelu funkcji $x(t)$ do danych $\{t_i, y_i\}$, $i = 1, 2, \dots, J$ (patrz Shmueli (2010)).

Założmy teraz, że mamy n niezależnych par wartości $\{t_{ij}, y_{ij}\}$, $i = 1, \dots, J, j = 1, \dots, n$. Te dane dyskretne przekształcamy w funkcje ciągłe postaci:

$$x_j(t) = \sum_{k=0}^{K_j} \hat{c}_{jk} \varphi_k(t), j = 1, \dots, n, t \in I.$$

Spośród wszystkich K_1, K_2, \dots, K_n wybieramy wartość modalną (powiedzmy K) i przyjmujemy, że każda funkcja $x_j(t)$ ma postać

$$x_j(t) = \sum_{k=0}^K \hat{c}_{jk} \varphi_k(t), j = 1, \dots, n, t \in I.$$

Uzyskany w ten sposób zbiór funkcji $\{x_1(t), \dots, x_n(t)\}$ nazywa się zbiorem danych funkcjonalnych (patrz Ramsay i Silverman (2005)).

Do tej pory zajmowaliśmy się danymi dotyczącymi pojedynczej cechy statystycznej. Rozważania nasze możemy uogólnić na przypadek $p \geq 2$ cech statystycznych. Wówczas nasze dane będą się składały z n niezależnych funkcji wektorowych postaci

$\mathbf{x}_j(t) = [x_{1j}(t), x_{2j}(t), \dots, x_{pj}(t)]'$, $j = 1, \dots, n$. Dane $\{\mathbf{x}_1(t), \dots, \mathbf{x}_n(t)\}$ nazywać będziemy wielozmiennymi danymi funkcjonalnymi .

PCA dla wielozmiennych danych funkcjonalnych

Wielozmienne dane funkcjonalne możemy traktować jako realizacje wielozmiennych procesów stochastycznych $\mathbf{X}(t)$ z ciągłym parametrem $t \in I$, gdzie $\mathbf{X}(t) = [X_1(t), X_2(t), \dots, X_p(t)]'$, $p \geq 2$. Będziemy zakładać, że $E(\mathbf{X}(t)) = \mathbf{0}$ oraz, że $\mathbf{X}(t) \in L_2(I)^p$, gdzie $L_2(I)$ jest przestrzenią Hilberta funkcji całkownych z kwadratem na przedziale I wyposażoną w iloczyn skalarny:

$$\langle u, v \rangle = \int_I u(t)v(t)dt.$$

Weźmy pod uwagę przypadek gdy i -ta składowa procesu stochastycznego $\mathbf{X}(t)$ może być reprezentowana za pomocą skończonej liczby ortonormalnych funkcji bazowych $\{\varphi_k\}$

$$X_i(t) = \sum_{k=0}^{K_i} c_{ik} \varphi_k(t), t \in I, i = 1, 2, \dots, p,$$

gdzie c_{ik} są zmiennymi losowymi takimi, że $E(c_{ik}) = 0, \text{Var}(c_{ik}) < \infty$, $i = 1, 2, \dots, p, k = 0, \dots, K_i$.

Niech

$$\mathbf{c} = (\mathbf{c}'_1, \mathbf{c}'_2, \dots, \mathbf{c}'_p),$$

$$\Phi(t) = \begin{bmatrix} \varphi'_{K_1}(t) & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \varphi'_{K_2}(t) & \dots & \mathbf{0} \\ \dots & \dots & \dots & \dots \\ \mathbf{0} & \mathbf{0} & \dots & \varphi'_{K_p}(t) \end{bmatrix},$$

gdzie $\mathbf{c}_i = (c_{i0}, \dots, c_{iK_i})'$, $\varphi_{K_i}(t) = (\varphi_0(t), \dots, \varphi_{K_i}(t))'$, $i = 1, \dots, p$.

Wówczas

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

Naszym celem jest znalezienie iloczynu skalarnego

$$U = \langle \mathbf{u}(t), \mathbf{X}(t) \rangle$$

mającego maksymalną wariancję dla wszystkich $\mathbf{u}(t) \in L_2(I)^p$ takich, że $\langle \mathbf{u}(t), \mathbf{u}(t) \rangle = 1$.

Niech

$$\lambda_1 = \sup_{\mathbf{u}(t) \in L_2(I)^P} \text{Var}(\langle \mathbf{u}(t), \mathbf{X}(t) \rangle) = \text{Var}(\langle \mathbf{u}_1(t), \mathbf{X}(t) \rangle),$$

gdzie $\langle \mathbf{u}_1(t), \mathbf{u}_1(t) \rangle = 1$. Iloczyn skalarny $U_1 = \langle \mathbf{u}_1(t), \mathbf{X}(t) \rangle$ nosi nazwę pierwszej funkcjonalnej składowej głównej, a funkcja wektorowa $\mathbf{u}_1(t)$ jest nazywana pierwszą wektorową funkcją wagową. Następnie poszukujemy drugiej funkcjonalnej składowej głównej $U_2 = \langle \mathbf{u}_2(t), \mathbf{X}(t) \rangle$, która maksymalizuje $\text{Var}(\langle \mathbf{u}(t), \mathbf{X}(t) \rangle)$, spełnia warunek $\langle \mathbf{u}_2(t), \mathbf{u}_2(t) \rangle = 1$ i jest nieskorelowana z pierwszą funkcjonalną składową główną U_1 , tj. spełnia warunek $\langle \mathbf{u}_1(t), \mathbf{u}_2(t) \rangle = 0$.

Ogólnie k -ta funkcjonalna składowa główna $U_k = \langle \mathbf{u}_k(t), \mathbf{X}(t) \rangle$ spełnia warunki

$$\lambda_k = \sup_{\mathbf{u}(t) \in L_2(I)^p} \text{Var}(\langle \mathbf{u}(t), \mathbf{X}(t) \rangle) = \text{Var}(\langle \mathbf{u}_k(t), \mathbf{X}(t) \rangle),$$

$$\langle \mathbf{u}_i(t), \mathbf{u}_j(t) \rangle = \delta_{ij}, \quad i, j = 1, \dots, k.$$

Parę $(\lambda_k, \mathbf{u}_k(t))$ nazywać będziemy k -tym układem głównym procesu losowego $\mathbf{X}(t)$.

Wcześniej pokazaliśmy, że proces $\mathbf{X}(t)$ ma reprezentację $\mathbf{X}(t) = \Phi(t)\mathbf{c}$, $t \in I$. Weźmy teraz pod uwagę składowe główne wektora losowego \mathbf{c} . k -ta składowa główna $U_k^* = \langle \mathbf{u}_k, \mathbf{c} \rangle$ tego wektora spełnia warunki:

$$\begin{aligned}\gamma_k &= \sup_{\mathbf{u} \in \mathbb{R}^{K+p}} \text{Var}(\langle \mathbf{u}, \mathbf{c} \rangle) = \sup_{\mathbf{u} \in \mathbb{R}^{K+p}} \mathbf{u}' \text{Var}(\mathbf{c}) \mathbf{u} \\ &= \sup_{\mathbf{u} \in \mathbb{R}^{K+p}} \mathbf{u}' \Sigma \mathbf{u}, \\ \mathbf{u}'_i \mathbf{u}_j &= \delta_{ij},\end{aligned}$$

gdzie $i, j = 1, \dots, k$, $K = K_1 + K_2 + \dots + K_p$. Parę (γ_k, \mathbf{u}_k) nazywać będziemy k -tym układem głównym wektora losowego \mathbf{c} .

Wyznaczenie k -tego układu głównego wektora losowego \mathbf{c} sprowadza się zatem do znalezienia wartości własnych i odpowiadających im wektorów własnych macierzy kowariancji $\mathbf{\Sigma}$ tego wektora normowanych tak, by $\mathbf{u}'_i \mathbf{u}_j = \delta_{ij}$. Zauważmy, że $\text{tr}(\mathbf{\Sigma}) = \sum_{i=1}^K \text{Var}(c_i) = \sum_{i=1}^K \gamma_i$, czyli suma wariancji zmiennych pierwotnych c_1, \dots, c_K jest równa sumie wariancji składowych głównych. Zatem wskaźnik

$$\frac{\gamma_1 + \dots + \gamma_r}{\gamma_1 + \dots + \gamma_K} \cdot 100\%$$

jest miarą wyjaśnienia całkowitej zmienności składowych wektora \mathbf{c} przez r pierwszych składowych głównych.

Twierdzenie

k -ty układ główny $(\lambda_k, \mathbf{u}_k(t))$ procesu losowego $\mathbf{X}(t)$ oraz k -ty układ główny (γ_k, \mathbf{u}) wektora losowego \mathbf{c} związane są zależnościami:

$$\lambda_k = \gamma_k, \quad \mathbf{u}_k(t) = \Phi(t)\mathbf{u}_k, \quad t \in I, k = 1, \dots, K + p,$$

gdzie $K = K_1 + K_2 + \dots + K_p$.

Macierz kowariancji Σ nie jest znana i musi być zastąpiona estymatorem z próby losowej. Niech $\hat{\mathbf{x}}_1(t), \hat{\mathbf{x}}_2(t), \dots, \hat{\mathbf{x}}_n(t)$ będą n niezależnymi realizacjami postaci $\hat{\mathbf{x}}_i(t) = \Phi(t)\hat{\mathbf{c}}_i$ procesu losowego $\mathbf{X}(t)$, gdzie wektory $\hat{\mathbf{c}}_i$ są scentrowane, $i = 1, 2, \dots, n$. Niech $\hat{\mathbf{C}} = (\hat{\mathbf{c}}_1, \hat{\mathbf{c}}_2, \dots, \hat{\mathbf{c}}_n)'$. Wówczas

$$\hat{\Sigma} = \frac{1}{n} \hat{\mathbf{C}}' \hat{\mathbf{C}}.$$

Jeżeli $n > K$, to macierz $\hat{\Sigma}$ z prawdopodobieństwem 1 jest dodatnio określona. Niech $\hat{\gamma}_1 \geq \hat{\gamma}_2 \geq \dots \geq \hat{\gamma}_s$ będą niezerowymi wartościami własnymi macierzy $\hat{\Sigma}$, a $\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_s$ odpowiadającymi im wektorami własnymi, gdzie $s = \text{rzęd}(\hat{\Sigma})$.

Ponadto k -ty układ główny procesu losowego $\mathbf{X}(t)$ wyznaczony z próby ma postać:

$$(\hat{\lambda}_k = \hat{\gamma}_k, \hat{\mathbf{u}}_k(t) = \Phi(t)\hat{\mathbf{u}}_k), \quad k = 1, \dots, s.$$

Współrzędna rzutu i -tej realizacji $\hat{\mathbf{x}}_i(t)$ procesu $\mathbf{X}(t)$ na k -tą funkcjonalną składową główną z próby jest równa:

$$\begin{aligned} \hat{U}_{ik} &= \langle \hat{\mathbf{u}}_k(t), \hat{\mathbf{x}}_i(t) \rangle = \langle \Phi(t)\hat{\mathbf{u}}_k, \Phi(t)\hat{\mathbf{c}}_i \rangle \\ &= \hat{\mathbf{u}}_k' \langle \Phi(t), \Phi(t) \rangle \hat{\mathbf{c}}_i = \hat{\mathbf{u}}_k' \hat{\mathbf{c}}_i, \end{aligned}$$

dla $i = 1, 2, \dots, n, k = 1, 2, \dots, s$. Zatem współrzędne rzutu i -tej realizacji $\hat{\mathbf{x}}_i(t)$ procesu $\mathbf{X}(t)$ na płaszczyznę dwóch pierwszych funkcjonalnych składowych głównych z próby są równe $(\hat{\mathbf{u}}_1' \hat{\mathbf{c}}_i, \hat{\mathbf{u}}_2' \hat{\mathbf{c}}_i), i = 1, 2, \dots, n$.

Przykład. Dane udostępnione przez prof. Waldemara Ratajczaka z Wydziału Nauk Geograficznych i Geologicznych UAM.

Liczba cech: $p = 6$

Liczba województw: $n = 16$

Liczba lat obserwacji: $J = 10$

Przyjmujemy:

Przedział czasowy: $[0, T] = [0, 10]$

Momenty czasowe: $t_1 = 0,5(2002), t_2 = 1,5(2003), \dots, t_{10} = 9,5(2011)$

Tabela : Województwa.

Nr	Województwo
1	ŁÓDZKIE
2	MAZOWIECKIE
3	MAŁOPOLSKIE
4	ŚLĄSKIE
5	LUBELSKIE
6	PODKARPACKIE
7	PODLASKIE
8	ŚWIĘTOKRZYSKIE
9	LUBUSKIE
10	WIELKOPOLSKIE
11	ZACHODNIOPOMORSKIE
12	DOLNOŚLĄSKIE
13	OPOLSKIE
14	KUJAWSKO-POMORSKIE
15	POMORSKIE
16	WARMIŃSKO-MAZURSKIE

Tabela : Cechy.

Nr	Cechy dotyczące ochrony środowiska
1	Emisja zanieczyszczeń gazowych (t/km^2)
2	Emisja zanieczyszczeń pyłowych (kg/km^2)
3	Odpady wytworzone w ciągu roku (t/km^2)
4	Ścieki odprowadzane ogółem (* $dam^3/1000mieszk.$)
5	Ścieki przemysłowe odprowadzane (* $dam^3/1000mieszk.$)
6	Zużycie wody na potrzeby gospodarki narodowej i ludności (* $dam^3/1000mieszk.$)

* dam^3 (dekametr sześcienny) = $1000m^3$

Tabela : Wartości współczynników wagowej funkcji wektorowej pierwszej funkcjonalnej składowej głównej.

	u_0	u_1	u_2	u_3	u_4	u_5	u_6	u_7	u_8	Udział %
1	0,6947	-0,0007	-0,0179	-0,0056	0,0189	-0,0220	-0,0104	-0,0064	-0,0046	40,48
2	0,2927	0,0794	0,0094	0,04010	0,0138	0,0114	-0,0101	0,0299	0,0011	17,94
3	0,6443	0,0551	-0,0088	0,0086	0,0129	-0,0010	-0,0074	0,0147	-0,0002	37,65
4	0,0008	0,0001	0,0000	0,0001	0,0000	0,0001	0,0000	0,0001	0,0000	0,00
5	-0,0317	0,0001	0,0009	-0,0005	-0,0009	0,0003	0,0003	0,0009	-0,0004	1,86
6	-0,0355	0,0013	0,0011	0,0006	-0,0008	0,0003	0,0004	0,0013	-0,0004	2,07

$$u_i(t) = \mathbf{u}_i' \boldsymbol{\varphi}(t), \quad i = 1, 2, \dots, 6,$$

$$\mathbf{u}_i = (u_{i0}, u_{i1}, \dots, u_{iK})', \quad \boldsymbol{\varphi}(t) = (\varphi_0(t), \varphi_1(t), \dots, \varphi_K(t))'.$$

Nierówność Cauchy'ego-Schwarza $|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|, \quad \|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}.$

$$|u_i(t)| = |\mathbf{u}_i' \boldsymbol{\varphi}(t)| \leq \|\mathbf{u}_i\| \|\boldsymbol{\varphi}(t)\|, \quad \|\mathbf{u}_i\| = \sqrt{\sum_{k=0}^K u_{ik}^2}.$$

Rysunek : Wykres 6 funkcji wagowych pierwszej funkcjonalnej składowej głównej.

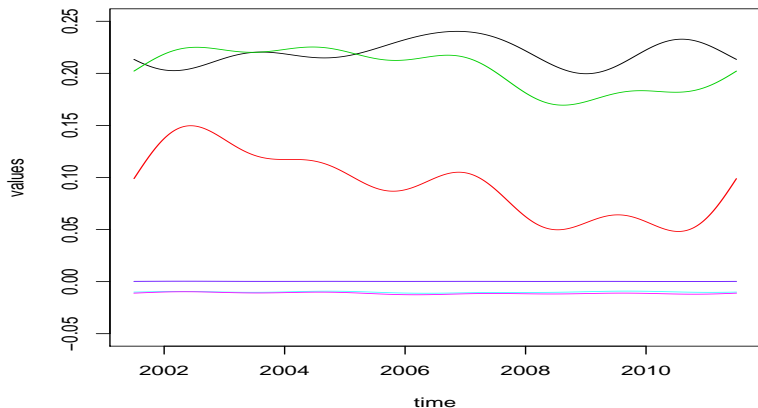
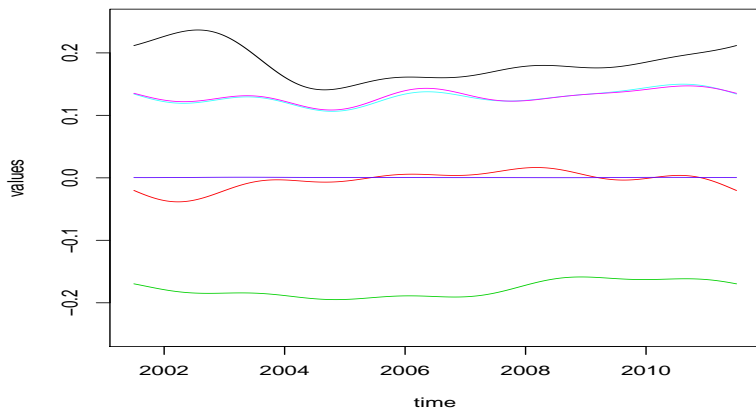


Tabela : Wartości współczynników wagowej funkcji wektorowej drugiej funkcjonalnej składowej głównej.

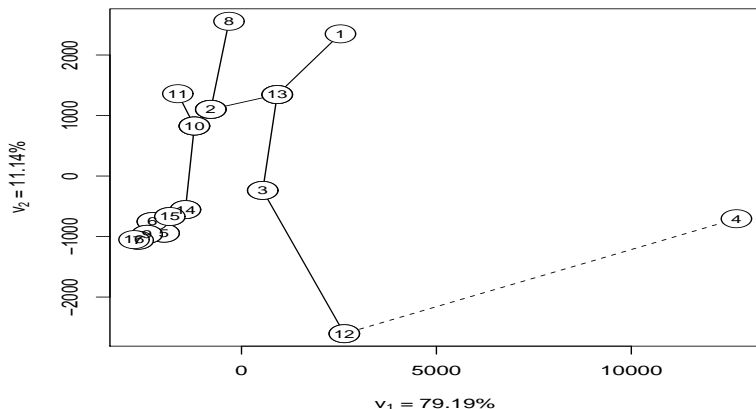
	u_0	u_1	u_2	u_3	u_4	u_5	u_6	u_7	u_8	Udział %
1	0,5771	-0,0005	0,0680	0,0392	0,0193	0,0165	-0,0105	-0,0091	-0,0114	28,89
2	-0,0122	-0,0227	-0,0289	-0,0087	-0,0092	-0,0140	0,0010	-0,0157	0,0005	2,27
3	-0,5636	-0,0314	0,0210	-0,0042	-0,0070	0,0018	0,0017	-0,0078	0,0034	27,99
4	0,0017	0,0004	0,0003	-0,0001	-0,0002	-0,0002	-0,0003	-0,0001	-0,0002	0,09
5	0,4080	-0,0210	0,0120	-0,0081	0,0096	-0,0065	-0,0155	-0,0117	0,0055	20,28
6	0,4120	-0,0164	0,0081	-0,0069	0,0122	-0,0041	-0,0163	-0,0115	0,0072	20,48

Rysunek : Wykres 6 funkcji wagowych drugiej funkcjonalnej składowej głównej.



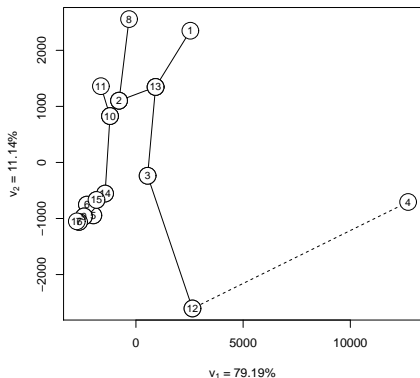
Przykład - Ochrona środowiska i rozwój zrównoważony

Rysunek : Położenie 16 województw w układzie dwóch pierwszych funkcjonalnych składowych głównych wraz z rozpiętym dendrytem.





Przykład - Ochrona środowiska i rozwój zrównoważony

4 - ŚLĄSKIE
12 - DOLNOŚLĄSKIE
3 - MAŁOPOLSKIE
13 - OPOLSKIE
1 - ŁÓDZKIE
2 - MAZOWIECKIE
10 - WIELKOPOLSKIE
11 - ZACHODNIOPOMORSKIE
8 - ŚWIĘTOKRZYSKIE



Rysunek : Położenie 16 województw w układzie dwóch pierwszych funkcjonalnych składowych głównych wraz z rozpiętym dendrytem.

-  Ramsay, J.O., Silverman, B.W. (2005). Functional Data Analysis, Second Edition, Springer.
-  Shmueli, G. (2010): To explain or to predict?, Statistical Science 25(3), 289-310.