

Wybór zmiennych w klasyfikacji dla wielowymiarowych danych funkcjonalnych

Tomasz Górecki, Mirosław Krzyśko, Waldemar Wołyński

Wydział Matematyki i Informatyki
Uniwersytet im. Adama Mickiewicza w Poznaniu

II Kongres Statystyki Polskiej
Warszawa 10-12.07.2018



Rozważamy **zagadnienie klasyfikacji obiektu** opisanego p -wymiarowym procesem stochastycznym $\mathbf{X} = (X_1, X_2, \dots, X_p)^\top$ do jednej z q klas. Ponadto zakładamy, że dysponujemy n -elementową próbą uczącą

$$\{(\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2), \dots, (\mathbf{x}_n, \mathbf{y}_n)\},$$

gdzie $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ są niezależnymi realizacjami procesu losowego \mathbf{X} , a $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n$ etykietami takimi, że $\mathbf{y}_i \in \mathbb{R}^q$, $i = 1, 2, \dots, n$.

W prezentacji przedstawimy algorytm pozwalający na **redukcję wymiaru** wektora losowego \mathbf{X} z jednoczesnym zachowaniem efektywności procedur klasyfikacyjnych.

- 1 Model danych funkcjonalnych.
- 2 Funkcjonalna kowariancja odległościowa.
- 3 Funkcjonalny współczynnik niezależności HSIC.
- 4 Algorytm wyboru zmiennych.
- 5 Przykład.
- 6 Literatura.

Założmy, że $\mathbf{X} \in L_2^p(I)$, gdzie $L_2(I)$ jest przestrzenią Hilberta funkcji całkowlanych z kwadratem na przedziale I . Ponadto założmy, że $E(\mathbf{X}) = \mathbf{0}$.

Z powyższego wynika, że każda składowa procesu \mathbf{X} może być przedstawiona w następującej postaci:

$$X_k(t) = \sum_{b=0}^{\infty} \alpha_{kb} \varphi_b(t), \quad t \in I,$$

przy czym funkcje $\varphi_1, \varphi_2, \dots$ tworzą bazę w przestrzeni $L_2(I)$.

W praktyce posługujemy się przybliżoną reprezentacją wykorzystującą jedynie skończoną liczbę pierwszych funkcji bazowych.

Założmy zatem, że k -ta składowa procesu \mathbf{X} ma następującą reprezentację:

$$X_k(t) = \sum_{b=0}^{B_k} \alpha_{kb} \varphi_b(t), \quad t \in I,$$

gdzie liczba B_k decyduje o stopniu gładkości funkcji X_k (im mniejsza wartość B_k tym większy stopień wygładzania).

Przyjmijmy następujące oznaczenia:

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

oraz

$$\Phi(t) = \begin{bmatrix} \varphi_1^\top(t) & \mathbf{0}^\top & \dots & \mathbf{0}^\top \\ \mathbf{0}^\top & \varphi_2^\top(t) & \dots & \mathbf{0}^\top \\ \dots & \dots & \dots & \dots \\ \mathbf{0}^\top & \mathbf{0}^\top & \dots & \varphi_p^\top(t) \end{bmatrix},$$

gdzie $\varphi_k(t) = (\varphi_0(t), \varphi_1(t), \dots, \varphi_{B_k}(t))^\top$, $k = 1, 2, \dots, p$.

Używając powyższej notacji macierzowej, proces \mathbf{X} może być zapisany w następujący sposób:

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

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

Dane funkcjonalne

Współczynniki $\alpha = (\alpha_1^\top, \alpha_2^\top, \dots, \alpha_p^\top)$, $\alpha_k \in \mathbb{R}^{B_k+1}$ szacujemy na bazie n niezależnych realizacji $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$ procesu losowego \mathbf{X} za pomocą metody najmniejszych kwadratów, w taki sposób aby zminimalizować funkcję

$$f(\mathbf{a}_k) = \sum_{j=1}^{J_k} (x_{kj} - \varphi_k^\top(t_{kj})\mathbf{a}_k)^2,$$

gdzie $\mathbf{a}_k = (a_{k0}, a_{k1}, \dots, a_{kB_k})^\top$, $k = 1, 2, \dots, p$.

W rezultacie procesu transformacji n niezależnych realizacji procesu losowego \mathbf{X} otrzymujemy **dane funkcjonalne** postaci:

$$\mathbf{x}_i(t) = \Phi(t)\mathbf{a}_i,$$

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

Kowariancja odległościowa

Funkcja charakterystyczna łącznego rozkładu procesu losowego $\mathbf{X} \in L_2^p(I)$ oraz wektora losowego $\mathbf{Y} \in \mathbb{R}^q$ ma postać

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

Przyjmijmy, że (Ramsay, Silverman (2005)) wektor funkcji wagowych l oraz proces losowy \mathbf{X} należą do tej samej podprzestrzeni przestrzeni $L_2^p(I)$, tzn. funkcja l może być przedstawiona w postaci

$$l(t) = \Phi(t)\lambda,$$

gdzie $\lambda \in \mathbb{R}^{K+p}$.

Zatem

$$\langle l, \mathbf{X} \rangle_p = \lambda' \alpha$$

oraz

$$f_{\mathbf{X}, \mathbf{Y}}(l, \mathbf{m}) = f_{\alpha, \mathbf{Y}}(\lambda, \mathbf{m}),$$

gdzie $f_{\alpha, \mathbf{Y}}(\lambda, \mathbf{m})$ jest funkcją charakterystyczną łącznego rozkładu pary wektorów losowych (α, \mathbf{Y}) .

Kowariancja odległościowa

Bazując na pomysle **kowariancji odległościowej** pomiędzy dwoma wektorami losowymi (Székely i inni (2007)), możemy zdefiniować **funkcjonalną kowariancję odległościową** pomiędzy procesem losowym \mathbf{X} a wektorem losowym \mathbf{Y} jako nieujemną liczbę $d\text{Cov}_{\mathbf{X},\mathbf{Y}}$ taką, że

$$d\text{Cov}_{\mathbf{X},\mathbf{Y}}^2 = \frac{1}{C_{K+p}C_q} \int_{\mathbb{R}^{K+p+q}} \frac{|f_{\alpha,\mathbf{Y}}(\boldsymbol{\lambda}, \mathbf{m}) - f_{\alpha}(\boldsymbol{\lambda})f_{\mathbf{Y}}(\mathbf{m})|^2}{\|\boldsymbol{\lambda}\|_{K+p}^{K+p+1} \|\mathbf{m}\|_q^{q+1}} d\boldsymbol{\lambda}d\mathbf{m},$$

gdzie

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

oraz $f_{\alpha}(\boldsymbol{\lambda})$, $f_{\mathbf{Y}}(\mathbf{m})$ są funkcjami charakterystycznymi rozkładów brzegowych.

Dla rozkładów o skończonych pierwszych momentach, funkcjonalna kowariancja odległościowa charakteryzuje niezależność w taki sposób, że $d\text{Cov}_{\mathbf{X}, \mathbf{Y}} = 0$ wtedy i tylko wtedy, gdy \mathbf{X} i \mathbf{Y} są **niezależne**.

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

$$d\text{Cov}_{\mathbf{X}, \mathbf{Y}}^2 = \frac{1}{n^2} \sum_{k, l=1}^n A_{kl} B_{kl},$$

gdzie

$$A_{kl} = \|\tilde{\mathbf{a}}_k - \tilde{\mathbf{a}}_l\|_{K+p}, \quad B_{kl} = \|\tilde{\mathbf{y}}_k - \tilde{\mathbf{y}}_l\|_q,$$
$$\tilde{\mathbf{a}}_k = \mathbf{a} - \bar{\mathbf{a}}, \quad \tilde{\mathbf{y}}_k = \mathbf{y} - \bar{\mathbf{y}}, \quad k, l = 1, 2, \dots, n.$$

Zakładamy, że jedynie pewna liczba składowych procesu losowego \mathbf{X} ma wpływ na wektor losowy \mathbf{Y} . Wybieramy istotne składowe w taki sposób, aby funkcjonalna kowariancja odległościowa $d\text{Cov}_{\mathbf{X},\mathbf{Y}} = d\text{Cov}_{\boldsymbol{\alpha},\mathbf{Y}}$ była duża.

Prawdziwe jest następujące twierdzenie

Twierdzenie (Kong i inni (2015))

Niech $\mathbf{X}, \mathbf{Z} \in \mathbb{R}^p$ i $\mathbf{Y} \in \mathbb{R}^q$ oraz załóżmy, że wektor losowy \mathbf{Z} jest niezależny od wektorów losowych (\mathbf{X}, \mathbf{Y}) . Wtedy

$$d\text{Cov}_{(\mathbf{X},\mathbf{Z}),\mathbf{Y}} \leq d\text{Cov}_{\mathbf{X},\mathbf{Y}} .$$

Zastosowaliśmy twierdzenie Konga jako **regułę stopu** w procedurze wyboru zmiennych. Procedura składa się z następujących kroków:

- 1 Wyznaczamy brzegową kowariancję odległościową dla X_k , $k = 1, \dots, p$ i Y .
- 2 Porządkujemy zmienne w kolejności malejących kowariancji odległościowych. Oznaczmy uporządkowane zmienne jako $X_{(1)}, X_{(2)}, \dots, X_{(p)}$. Zaczynamy od $\mathbf{X}_S = \{X_{(1)}\}$.
- 3 Dla k od 2 do p , dodajemy $X_{(k)}$ do \mathbf{X}_S jeżeli $d\text{Cov}_{\mathbf{X}_S, Y}$ rośnie o więcej niż zadana wartość progowa (ε). Zatrzymujemy procedurę w przeciwnym razie.

Niech $\mathbf{Z} \in \mathbb{R}^P$ będzie wektorem losowym. Ponadto, niech

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

będzie rzeczywistą funkcją ciągłą, zwaną **jądrem** oraz $\mathbf{K} = (k_{ij})$, gdzie $k_{ij} = k(\mathbf{z}_i, \mathbf{z}_j)$ będzie **macierzą jądrową**, $i, j = 1, 2, \dots, n$.

W dalszym ciągu wszystkie jądra będą **jądrami gaussowskimi**, tzn.

$$k(\mathbf{z}_i, \mathbf{z}_j) = \exp(-\lambda \|\mathbf{z}_i - \mathbf{z}_j\|^2), \quad \lambda > 0.$$

Jądrowy współczynnik zgodności

Niech $\mathbf{X} \in L_2^p(I)$ będzie procesem losowym, $\mathbf{x}_1, \dots, \mathbf{x}_n$ jego niezależnymi realizacjami postaci:

$$\mathbf{x}_i(t) = \Phi(t)\mathbf{a}_i, \quad t \in I, \quad i = 1, 2, \dots, n.$$

Wtedy

$$k_{\mathbf{X}}(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\lambda_1 \|\mathbf{x}_i - \mathbf{x}_j\|^2) = \exp(-\lambda_1 \|\mathbf{a}_i - \mathbf{a}_j\|^2) = k_{\mathbf{a}}(\mathbf{a}_i, \mathbf{a}_j),$$

dla $i, j = 1, 2, \dots, n$.

Analogicznie, niech $\mathbf{Y} \in \mathbb{R}^q$ będzie wektorem losowym, $\mathbf{y}_1, \dots, \mathbf{y}_n$ jego niezależnymi realizacjami.

Wtedy

$$k_{\mathbf{Y}}(\mathbf{y}_i, \mathbf{y}_j) = \exp(-\lambda_2 \|\mathbf{y}_i - \mathbf{y}_j\|^2),$$

dla $i, j = 1, 2, \dots, n$.

Jądrowy współczynnik zgodności

Niech $\phi: L_2^p(I) \rightarrow \mathcal{H}$, gdzie \mathcal{H} jest przestrzenią Hilberta, będzie przekształceniem takim, że dla dowolnych $\mathbf{x}, \mathbf{x}' \in L_2^p(I)$:

$$\langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle_{\mathcal{H}} = k_{\mathbf{X}}(\mathbf{x}, \mathbf{x}').$$

Analogicznie, niech $\psi: \mathbb{R}^q \rightarrow \mathcal{G}$, gdzie \mathcal{G} jest przestrzenią Hilberta, będzie przekształceniem takim, że dla dowolnych $\mathbf{y}, \mathbf{y}' \in \mathbb{R}^q$:

$$\langle \psi(\mathbf{y}), \psi(\mathbf{y}') \rangle_{\mathcal{G}} = k_{\mathbf{Y}}(\mathbf{y}, \mathbf{y}').$$

Operator kowariancji pomiędzy procesem losowym \mathbf{X} , a wektorem losowym \mathbf{Y} definiujemy następująco:

$$\text{Cov}_{\mathbf{X}, \mathbf{Y}} = \text{E}_{\mathbf{X}, \mathbf{Y}}[(\phi(\mathbf{X}) - \text{E}_{\mathbf{X}}(\phi(\mathbf{X}))) \otimes (\psi(\mathbf{Y}) - \text{E}_{\mathbf{Y}}(\psi(\mathbf{Y})))].$$

Jądrowy współczynnik zgodności

Bazując na definicji **współczynnika niezależności HSIC** pomiędzy dwoma wektorami losowymi (Gretton i inni (2005)), możemy zdefiniować **funkcjonalny współczynnik niezależności HSIC** pomiędzy procesem losowym \mathbf{X} a wektorem losowym \mathbf{Y} jako

$$\text{HSIC}_{\mathbf{X},\mathbf{Y}} = \|\text{Cov}_{\mathbf{X},\mathbf{Y}}\|_{HS}^2.$$

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

Ponadto dla jądra gaussowskiego prawdziwe jest następujące twierdzenie:

Twierdzenie

Niech $\mathbf{X}, \mathbf{Z} \in \mathbb{R}^p$ i $\mathbf{Y} \in \mathbb{R}^q$ oraz załóżmy, że wektor losowy \mathbf{Z} jest niezależny od wektorów losowych (\mathbf{X}, \mathbf{Y}) . Wtedy

$$\text{HSIC}_{(\mathbf{X},\mathbf{Z}),\mathbf{Y}} \leq \text{HSIC}_{\mathbf{X},\mathbf{Y}}.$$

Estymator współczynnika zgodności HSIC ma postać:

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

gdzie $\tilde{\mathbf{K}}_{\mathbf{X}} = \mathbf{H}\mathbf{K}_{\mathbf{X}}\mathbf{H}$, $\tilde{\mathbf{K}}_{\mathbf{Y}} = \mathbf{H}\mathbf{K}_{\mathbf{Y}}\mathbf{H}$ oraz $\mathbf{H} = (\mathbf{I}_n - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top)$.

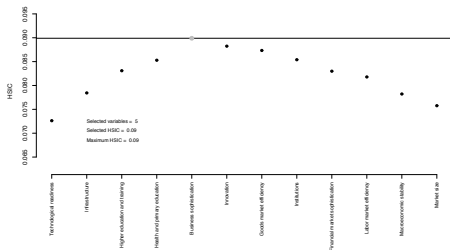
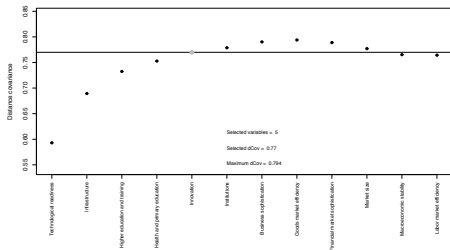
Zatem, w algorytmie wyboru zmiennych, zamiast kowariancją odległościową możemy posłużyć się współczynnikiem HSIC.

Rozważmy dane pochodzące ze stron **Światowego Forum Ekonomicznego (WEF)**. Dane dotyczące wskaźników socjo-ekonomicznych 115 krajów z lat 2008-2017. Ekspert WEF podzielili analizowane kraje na pięć grup.

W naszych rozważaniach wzięliśmy pod uwagę 12 zmiennych, tzw. filarów:

1. Institutions
2. Infrastructure
3. Macroeconomic stability
4. Health and primary education
5. Higher education and training
6. Goods market efficiency
7. Labor market efficiency
8. Financial market sophistication
9. Technological readiness
10. Market size
11. Business sophistication
12. Innovation

Przykład



Wybrane zmienne:

2. Infrastructure
4. Health and primary education
5. Higher education and training
9. Technological readiness
12. Innovation

2. Infrastructure
4. Health and primary education
5. Higher education and training
9. Technological readiness
11. Business sophistication

Ocena poprawności klasyfikacji






Baza Fouriera: $B_k = 5$.

Metoda oceny: LOO CV.

Reguła stopu: $\varepsilon = 0.05$.

Klasyfikator	Wybrane zmienne (5)	Wszystkie zmienne (12)
LDA	71,30%	66,09%
kNN ($k = 1, \dots, 8$)	77,39%	71,30%
Naiwny Bayes (normal)	69,57%	65,22%
Naiwny Bayes (kernel)	67,83%	62,61%
Multinom	60,87%	56,52%

Bibliography

-  GÓRECKI, T., KRZYŚKO, M., WASZAK, Ł., WOŁYŃSKI, W. (2014): Methods of reducing dimension for functional data. *Statistics in Transition new series* 15, 231–242.
-  GRETTON, A., BOUSQUET, O., SMOLA, A., SCHÖLKOPF, B., (2005): Measuring statistical dependence with Hilbert-Schmidt norms. In: *Algorithmic Learning Theory* (S. Jain, H.U. Simon and E. Tomita, eds.). *Lecture Notes in Computer Science* 3734, 63–77. Springer, Berlin.
-  KONG, J., WANG, S., WAHBA G. (2015): Using distance covariance for improved variable selection with application to learning genetic risk models. *Statistics in Medicine* 34, 1708–1720.
-  RAMSAY, J.O., SILVERMAN, B.W. (2005): *Functional Data Analysis*. Springer, New York.
-  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.