

Plan referatu

- 1 Zastosowania funkcjonalnej analizy danych w ekonomii
 - Odporna analiza ekonomiczna
 - Odporność procedury statystycznej
- 2 Odporna funkcjonalna analiza danych
 - Statystyczne funkcje głębi
 - Metody wykrywania funkcjonalnych outlierów
- 3 Wyzwania dla zastosowań odpornej FDA w ekonomii
 - Odporna analiza funkcjonalnych szeregów czasowych (FTS)
 - Hierarchiczne funkcjonalne szeregi czasowe (HFTS)
 - Propozycje odpornej predykcji HFTS
 - Popularne algorytmy analizy skupisk (AS) w FDA
 - Algorytm k-średnich w FDA
 - Odporna AS w FDA
 - Odporna analiza dyskryminacyjna
- 4 Przykłady empiryczne
 - Funkcjonalne outliery w dobowym zanieczyszczenie powietrza w Katowicach
 - Wykrywanie nietypowych zachowań internautów
 - Odporna predykcja HFTS opisującego zanieczyszczenie powietrza na Śląsku
 - Odporna AS dobowego zanieczyszczenia powietrza w Krakowie

Próba z modelu funkcjonalnej autoregresji FAR(1)

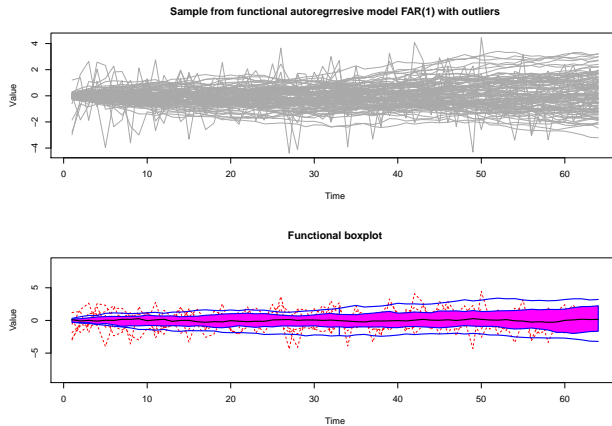
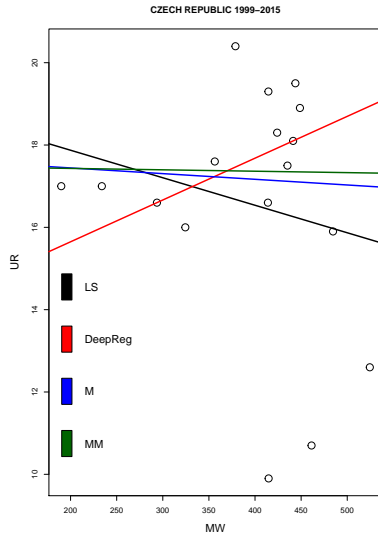
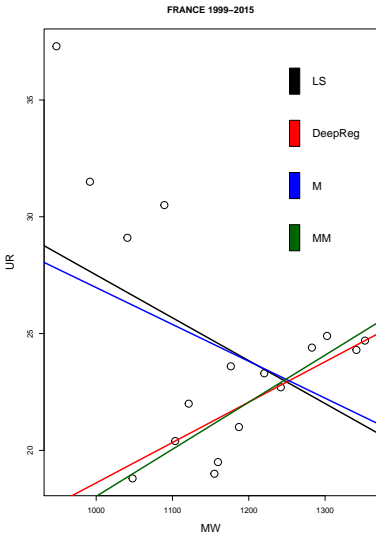


Figure 1: Próba 100 obs. z procesu FAR(1) z jądrem wykładniczym

Zależność pomiędzy stopą bezrobocia (UR) a płacą minimalną (MW) we Francji i Republice Czeskiej w okresie 1999-2016



Formalne koncepcje odporności

Niech pewna zmienna losowa X przyjmuje wartości z pewnej przestrzeni próbkowej (\mathbf{X}, \mathbb{B}) . Oznaczmy przez P jej rozkład prawdopodobieństwa (zob. Bosq 2000). Rozważmy odporność funkcjonału $T(P) = T(X)$. Funkcjonał ten jest szacowany na podstawie próby $X^n = \{X_1, \dots, X_n\}$ złożonej z n niezależnych kopii X . Ściślej, szacujemy T za pomocą funkcjonału empirycznego $T_n(P_n) = T_n(X_1, \dots, X_n)$, opierającego się o rozkład empiryczny P_n policzony z próby $X^n = \{X_1, \dots, X_n\}$.

Można zaproponować następującą modyfikację definicji F. Hampela jakościowej odporności funkcjonału T , która wykorzystuje pewną ustaloną metrykę d_P na rodzinie miar probabilistycznych \mathbb{P} :

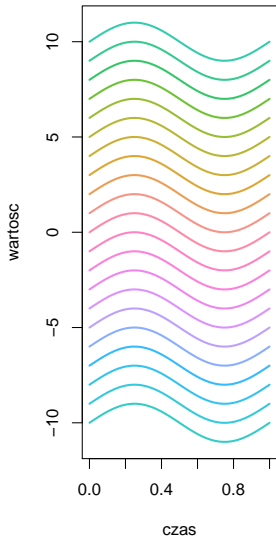
Definicja 1: Powiemy, że ciąg funkcjonałów $\{T_n\}$ jest **jakościowo odporny** w $P \in \mathbb{P}$, wtedy i tylko wtedy gdy dla dowolnego $\varepsilon > 0$, istnieje $\delta > 0$ oraz dodatnia liczba całkowita n_0 taka że, dla wszystkich $Q \in \mathbb{P}$ oraz $n > n_0$

$$d_P(P, Q) < \delta \Rightarrow d_T(L_P(T_n), L_Q(T_n)) < \varepsilon,$$

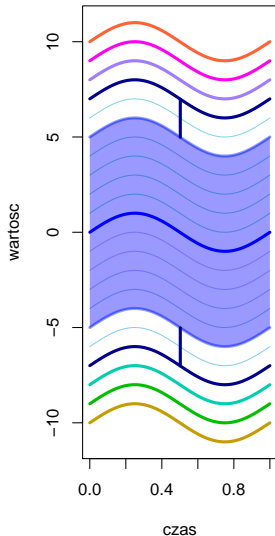
gdzie $P, Q \in \mathbb{P}$ oznaczają dwa rozkłady w przestrzeni Hilberta L^2 oraz L_P, L_Q oznaczają charakterystyki rozkładów P, Q . L_{P_n}, L_{Q_n} mogą oznaczać wartości miar jakości wyników analizy skupisk, wartości błędów klasyfikacji itp.

Funkcjonalne outliersy – wprowadzenie

Proba z procesu Gaussa



Funkc. wykres pudełkowy



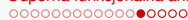
Funkcjonalne outliersy cd.

Zasadniczy problem każdej procedury FDA – właściwa analiza zmienności danych. Próby obserwacji – realizacji losowych funkcji – losowych elementów pewnych przestrzeni funkcyjnych na ogół nieskończone wymiarowych, rzeczywistych, ośrodkowych przestrzeni Banacha lub Hilberta. Zupełność i ośrodkowość przestrzeni są kluczowym założeniem – gwarantują, że sumowalna liniowa kombinacja elementów losowych dalej jest elementem losowym rozważanej przestrzeni a jednocześnie każdy element losowy jest kombinacją liniową co najwyżej przeliczalnej liczby elementów z pewnej bazy (ośrodka). Rozważamy losowe funkcje postaci

$$X : (\Omega, \mathcal{B}, \mathcal{P}) \rightarrow \mathcal{V},$$

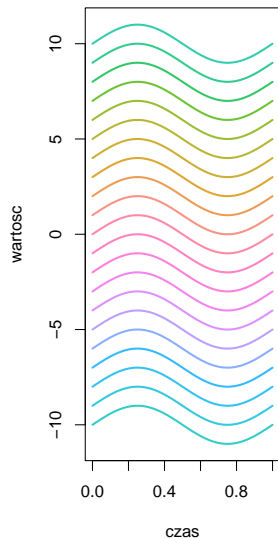
gdzie $(\Omega, \mathcal{B}, \mathcal{P})$ jest przestrzenią probabilistyczną oraz \mathcal{V} oznacza ośrodkową przestrzeń Banacha lub Hilberta wyposażoną w normę $\|\cdot\|$, w przypadku przestrzeni Hilberta, norma indukowana jest przez iloczyn skalarny. Dla wszystkich $\omega \in \Omega$ mamy

$$X_{\omega} : t \rightarrow X(\omega, t) \in \mathcal{V}$$

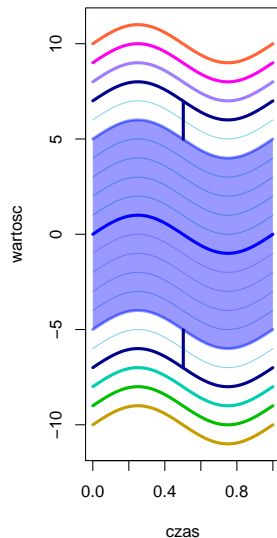


Wykrywanie funkcjonalnych outlierów typu "magnitude"

Proba z procesu Gaussa



Skor. funkcj. wyk. pud.



Wykrywanie outlierów w FTS–aktualny stan wiedzy

- Febrero i in. (2007, 2008): Pierwszy bootstrapowy test wykrywania do outlierów dla przypadku niezależnych obserwacji o tym samym rozkładzie.
- Funkcjonalny boxplot i skorygowany boxplot Sun i Genton (2011, 2012), Tarabelloni (2017).
- Outliergram zaproponowany przez Arribas-Gil i Romo (2014).
- Rana i in. (2015): W przypadku FTS setup, lokalne trendy wynikające ze struktury zależności mogą maskować obecność outlierów–zaproponowali modyfikację procedury Febrero i in.(2008).
- Nagy i in. (2017): Oryginalne podejście do klasyfikacji i wykrywania outlierów wykorzystujące pochodne funkcji i tzw. głębie całkowite.

Outlieroqram–wykrywanie outlierów typu "shape"

Dysponujemy próbą $X^n = \{X_1, \dots, X_n\}$ funkcji określonych na przedziale $I \subset \mathbb{R}$.
Modified Epigraphic Index (MEI) funkcji Z względem próby X^n :

$$MEI(Z|X^n) = \frac{1}{n} \sum_{i=1}^n \frac{\lambda(\{t \in I : Z(t) \leq X_i(t)\})}{\lambda(I)},$$

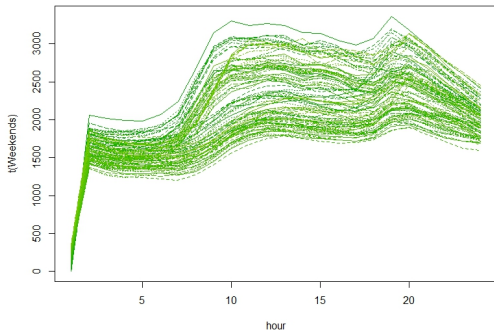
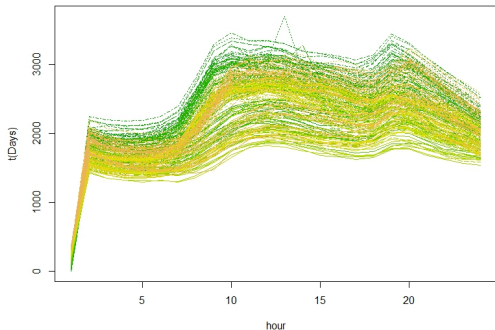
gdzie λ oznacza miarę Lebesgue'a. Arribas- Gil i Romo (2014) pokazali, że MEI i MBD (zmodyfikowana głębina pasma) pozostają w relacji, którą można wykorzystać do wykrywania outlierów typu "shape":

$$MBD(Z, X^n) \leq a_0 + a_1 MEI(Z|X^n) + a_2 n^2 MEI^2(Z|X^n),$$

gdzie $a_0 = a_2 = -2/n(n-1)$, $a_1 = 2(n+1)/(n-1)$.

Outliery typu "shape" mają niewielką wartość MBD w porównaniu do tej opisanej za pomocą parabolicznego trendu, rozbieżność pomiędzy prawą i lewą stroną nierówności pozwala je zidentyfikować.

Obserwacje funkcjonalne pokazujące zużycie energii elektrycznej w dni robocze (po lewej) i w weekendy (po prawej) w 2016 w Danii



Hierarchiczne funkcjonalne szeregi czasowe (HFTS)

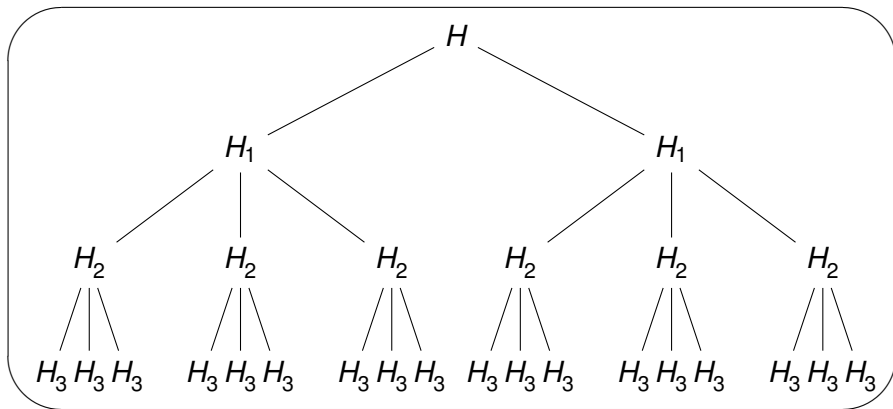
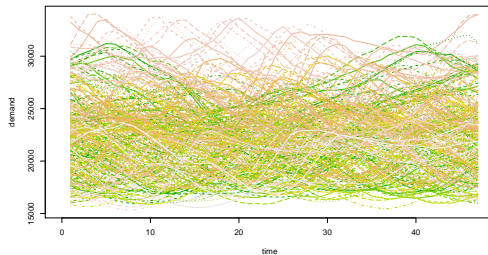


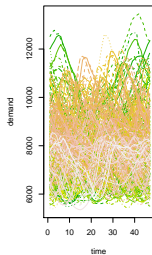
Figure 11: Przykładowa struktura hierarchicznego szeregu czasowego

Przykład hierarchicznego szeregu czasowego (HFTS)

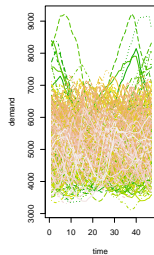
Electricity Demand Australia Total



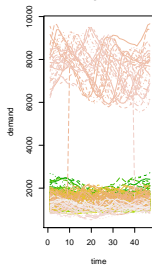
NSW



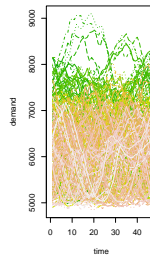
VIC



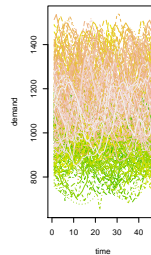
SA



QLD



TAS



Znane metody prognozowania HFTS

METODY:

Propozycja Hyndmana i Shanga (2017), rodzina naszych propozycji na bazie ruchomej funkcjonalnej mediany (2017a,b,2018), ruchoma funkcjonalna średnia, metoda naiwna

KRYTERIA OCENY PREDYKTORA HFTS:

- * trafność predykcji
- * zgodność (statystyczna)
- * zgodność agregacyjna (ang. aggregate consistency termin merytoryczny, H&S)
- * nieobciążoność
- * odporność
- * złożoność obliczeniowa
- * sposób oceny niepewności prognozy

Predyktor referencyjny–propozycja Hyndmana i Shanga (2017)

Podejście wywodzi się z wcześniejszych prac autorów dotyczących szeregów hierarchicznych 1D. Uzgodniona prognoza ma postać

$$\hat{X}_{n+1}(t) = F(\hat{x}^{level_1}, \hat{x}^{level_2}, \dots, \hat{x}^{level_L}),$$

gdzie \hat{x}^{level_i} oznacza prognozy uzyskane dla FTS na poziomie i oraz F jest pewną postacią estymatora uogólnionej metody najmniejszych kwadratów (GMNK).

Propozycja Hyndma i Shanga część 2

Struktura HFTS w chwili t :

$$X_t = S_t b_t,$$

gdzie wektor X_t zawiera wszystkie szeregi na wszystkich poziomach hierarchii;
 b_t jest wektorem najbardziej zdezagregowanych szeregów (tzn., na najniższym poziomie hierarchii);

S_t jest macierzą zależności pomiędzy wektorami.

Prognoza przybiera postać

$$\hat{X}_{n+1} = S_{n+1} \beta_{n+1} + \epsilon_{n+1},$$

gdzie \hat{X}_{n+1} jest macierzą prognoz dla wszystkich szeregów na wszystkich poziomach hierarchii. $\beta_{n+1} = E[b_{n+1} | X_1, \dots, X_n]$ jest nieznaną wartością oczekiwaną rozkładu prognozy dla szeregów na najniższym poziomie oraz ϵ_{n+1} oznacza błędy związane z uzgadnianiem prognoz.

Propozycja Hyndmana i Shanga część 3

Prognozy dla poszczególnych poziomów uzyskuje się stosując wysoce efektywną jednak nieodporną metodę, w której FTS są przekształcane do rodzin jednowymiarowych szeregów czasowych współrzędnych funkcjonalnych składowych głównych (ang. functional principal component scores) (por. Kosiorowski, 2014), β_{n+1} w podejściu Shanga i Hyndmana (2017) szacowane są za pomocą uogólnionej metody najmniejszych kwadratów GMNK:

$$\hat{\beta}_{n+1} = \left(S_{n+1}^T W^{-1} S_{n+1} \right)^{-1} S_{n+1}^T W^{-1} \hat{X}_{n+1},$$

gdzie W jest macierzą diagonalną, z wariancjami prognoz dla szeregów. Finalnie, prognozę uzyskujemy z równania

$$\bar{X}_{n+1} = S_{n+1} \hat{\beta}_{n+1}.$$

Zalety i wady metody Hyndmana i Shanga

Prognozy metodą H&S są agregacyjnie zgodne (spełniają ograniczenia narzucone przez sposób agregacji szeregów) oraz są nieobciążone w sensie wartości oczekiwanej.

Metoda jest bardzo wymagająca pod względem obliczeniowym i programistycznym (optymalizacja MNNK przy rzadkich macierzach eksperymentu, uogólnione odwrotności wielkich macierzy).

Metodę da się jednak "uodpornić" (por Kosiorowski i in., 2017b, 2018b).

Metoda Shanga i Hyndmana jest wyposażona w bezpośredni wewnętrzny mechanizm uzgadniania prognoz, dla porównania w naszej metodzie podwójnej mediany, uzgodnienie jest produktem ubocznym własności konkretnej mediany funkcjonalnej, nieprzechodności głębi, która ją indukuje.

Metoda S&H redukuje zagadnienie predykcji FTS do zagadnienia pewnej formy funkcjonalnej regresji (ang. functional principal component regression – funkcje przedstawia się w bazie empirycznych funkcjonalnych składowych a następnie stosuje się metody prognozowania stacjonarnych szeregów czasowych 1D do współrzędnych tej reprezentacji (por. Kosiorowski, 2014)).

Rodzina propozycji na bazie ruchomej funkcjonalnej mediany

Mamy próbę N funkcji, tzn., $X^N = \{x_i(t), i = 1, 2, \dots, N\}$ oraz $t \in [0, T]$. Niech $FD(y|X^N)$ oznacza funkcjonalną głębieść $y(t)$ z próby X^N . Medianę z próby definiujemy jako

$$MED_{FD}(X^N) = \arg \max_{i=1, \dots, N} FD(x_i|X^N).$$

Uwaga 1: jeżeli maksimum przyjmuje więcej niż jedna funkcja, to jako medianę przyjmujemy ich punktową średnią.

Uwaga 2: Różne funkcjonalne głębieść mogą prowadzić do różnych median, uwypuklających różne aspekty kształtu/położenia większości danych.

Propozycja metody podwójnej mediany

W propozycji korzystamy z ruchomej funkcjonalnej mediany:

$$\hat{x}_{n+1}(t) = MED_{FD}(W_{n,k})$$

gdzie $W_{n,k}$ jest ruchomym oknem o długości k z końcem w chwili n , tzn.,

$$W_{n,k} = \{x_{n-k+1}, \dots, x_n\}.$$

Uwaga: można zastosować układ wag, w którym najbliższe obserwacje "ważą więcej" niż obserwacje bardziej odległe w czasie.

Propozycja metody podwójnej mediany– szczegóły

Pierwszy krok: obliczamy ruchomą funkcjonalną medianę indukowaną przez głąbię MBD dla każdego węzła na najniższym poziomie hierarchii tzn.,

$$MED_{MBD}(W_{n,k}).$$

W przykładzie empirycznym dla każdego miasta w chwili t , liczymy funkcjonalną medianę indukowaną przez MBD na podstawie ruchomego okna o długości 10:

$$\hat{x}_{t+1}^{miasto} = MED_{MBD}\{x_t^{miasto}, x_{t-1}^{miasto}, \dots, x_{t-9}^{miasto}\}.$$

Drugi krok: dla poziomu hierarchii wyższego o jeden liczymy funkcjonalną medianę z median policzonych w kroku pierwszym.

Powtarzamy drugi krok do chwili, gdy dotrzemy do najwyższego poziomu hierarchii.

W przykładzie empirycznym drugi krok jest krokiem ostatnim, otrzymujemy prognozę dla $t = 10, \dots, 181$:

$$\hat{x}_{t+1} = MED_{MBD}\{\hat{x}_{t+1}^{miasto_1}, \dots, \hat{x}_{t+1}^{miasto_5}\}.$$

Typowe podejścia do AS w FDA

- Obserwacje funkcjonalne wyrażamy w ustalonym układzie funkcji bazowych a następnie stosujemy klasyczne metody AS do współrzędnych tych reprezentacji – Abraham i in. (2003), Serban i Wasserman (2005), James i Sugar (2003) and Luan and Li (2003) – wyniki AS są wrażliwe na wybór bazy, liczbę jej elementów, wybór węzłów).
- Dane rzutowane są na skończenie wymiarową przestrzeń rozpiętą przez kilka szczególnie ważnych empirycznych składowych głównych a następnie stosuje się klasyczne metody AS do współrzędnych obserwacji w układzie współrzędnych tych składowych Song i in (2007).
- Praca bezpośrednio z surowymi danymi empirycznymi przy nałożeniu pewnych ograniczeń co do gładkości funkcji Ma i in. (2006), Rocci i Gattone - zredukowane funkcjonalne k-średnich.

Algorytm k-średnich dla danych funkcyjnych

- Niech $\{x_1(t), x_2(t), \dots, x_l(t)\}$ będzie zbiorem funkcyjnych obserwacji, gdzie $t \in \Gamma$, a Γ jest to przedział na prostej \mathbb{R} . Funkcje należą do ośrodkowej przestrzeni Hilberta.
- Załóżmy, że funkcyjne obserwacje $x_i(t)$ dane są wzorem

$$x_i(t) = \sum_{g=1}^G u_{ig} m_g(t) + \varepsilon_i(t), \quad i = 1, \dots, l, \quad (3)$$

gdzie

- m_g – gładkie nieznanne funkcje (centroidy),
- $\varepsilon_i(t)$ – nieobserwowalne błędy losowe o średniej zero,
- $u_{ig} \in \{0, 1\}$, $\sum_g u_{ig} = 1$ dla każdego i oraz $u_{ig} = 1$, jeżeli x_i należy do g -tego skupiska.



Parametry u_{ig} i m_g są szacowane minimalizując wyrażenie

$$\sum_{i,g} u_{ig} \int_{\Gamma} [x_i(t) - m_g(t)]^2 dt = \sum_{i,g} u_{ig} \|x_i(t) - m_g(t)\|^2 \rightarrow \min . \quad (4)$$

Odporna analiza skupisk w FDA

- Zasadnicza kwestia – jak rozumieć odporność procedury AS?
- Odporna wersja algorytmu Warda.
- Metoda k-lokalnych median.

Algorytm Warda dla danych funkcjonalnych

Metoda Warda należy do metod grupowania aglomeracyjnego. Algorytm ten łączy grupy w taki sposób, aby miara niejednorodności wewnątrz skupienia nie wzrosła „zbyt wiele”.

Niech A i B będą niepustymi podzbiórmi przestrzeni L_2 . Wówczas odległość pomiędzy zbiorami A i B zdefiniowana jest następująco

$$d(A, B) = \sum_{i \in A \cup B} \|f_i - \bar{f}_{A \cup B}\|^2 - \sum_{i \in A} \|f_i - \bar{f}_A\|^2 - \sum_{i \in B} \|f_i - \bar{f}_B\|^2,$$

gdzie $\|\cdot\|$ oznacza normę w przestrzeni L_2 , a \bar{f}_A jest centroidem skupienia A .

- Zauważmy, że odległość Warda określa o ile wzrośnie rozproszenie wewnątrz nowego skupienia po połączeniu skupień A i B .
- Odległość $d(A, B)$ nazywamy kosztem połączenia skupień A i B (ang. *merging cost*).
- jeżeli $A = \{f_1, \dots, f_{n_A}\}$ oraz $B = \{g_1, \dots, g_{n_B}\}$, to odległość pomiędzy dwoma skupieniami można zapisać jako

$$d(A, B) = \frac{n_A n_B}{n_A + n_B} \|\bar{f} - \bar{g}\|^2$$

- Odległość w metodzie Warda jest zawsze dodatnia oraz $d(A, B) = 0$ wtedy i tylko wtedy, gdy $\bar{f}_A = \bar{f}_B$.

Kryterium grupowania – metoda Warda

Szukamy podziału zbioru danych $\mathbf{C} \in \Pi_k$ minimalizującego wyrażenie

$$\sum_{C \in \mathbf{C}} \sum_{f \in C} \|f - \bar{f}_C\|^2,$$

gdzie

$$\Pi_k = \{\mathbf{C} : \mathbf{C} \text{ – podział zbioru danych i } \#\mathbf{C} = k\}$$

Rozszerzenie metody Warda – Szlachtowska (2018)

Metoda Warda jest oparta na odległości pomiędzy skupieniami lub równoważnie na odległości pomiędzy centroidami skupień. Niestety, metoda oparta na odległości pomiędzy centroidami nie jest efektywna dla skupień o równych średnich. W przestrzeni euklidesowej próbuje się modyfikować algorytm Warda poprzez modyfikacje odległości pomiędzy skupieniami. Nasze rozważania dotyczą przestrzeni L_2 oraz danych funkcjonalnych określonych na tej przestrzeni. Wyniki zawarte w pracy Székely i Rizzo przenieśliśmy dla danych funkcjonalnych.

Székely G.J., Rizzo M.L. (2005), "Hierarchical Clustering Via Joint Between-Within Distances: Extending Ward's Minimum Variance Method", *Journal of Classification*, 22(2), 151–183.

Odległość między skupieniami

Niech $A = \{f_1, \dots, f_{n_A}\}$ i $B = \{g_1, \dots, g_{n_B}\}$ będą niepustymi podzbiórmi przestrzeni L_2 . Wówczas odległość pomiędzy zbiorami A i B definiujemy następująco

$$d(A, B) = \frac{n_A n_B}{n_A + n_B} \left(\frac{2}{n_A n_B} \sum_{i=1}^{n_A} \sum_{j=1}^{n_B} \|f_i - g_j\|^\alpha \right. \\ \left. - \frac{1}{n_A^2} \sum_{i=1}^{n_A} \sum_{j=1}^{n_A} \|f_i - f_j\|^\alpha - \frac{1}{n_B^2} \sum_{i=1}^{n_B} \sum_{j=1}^{n_B} \|g_i - g_j\|^\alpha \right), \quad (5)$$

gdzie $\|\cdot\|$ oznacza odległość w przestrzeni L_2 , a $\alpha \in (0, 2]$.

Można wykazać, że dla każdego $\alpha \in (0, 2]$ w zmodyfikowanym algorytmie Warda odległość między skupieniami może być obliczana za pomocą rekurencyjnej formuły Lance'a-Williamsa.

Twierdzenie

Niech A, B, C będą rozłącznymi, niepustymi skończonymi podzbiorami przestrzeni L_2 takimi, że

$$d(A, B) \leq \min\{d(A, C), d(B, C)\}.$$

Wówczas rekurencyjna formuła dla odległości $d(A \cup B, C)$ wyraża się wzorem

$$d(A \cup B, C) = \frac{n_A + n_C}{n_A + n_B + n_C} d(A, C) + \frac{n_B + n_C}{n_A + n_B + n_C} d(B, C) - \frac{n_C}{n_A + n_B + n_C} d(A, B) \quad (6)$$

- W przypadku $\alpha = 2$ jest proporcjonalna do kwadratu odległości w przestrzeni L_2 . W rezultacie szukamy takiego grupowania, które minimalizuje wariancję w skupieniach, czyli minimalizuje funkcję celu w metodzie Warda. Należy się zatem spodziewać, że odległość $d(A, B)$ jest wyrażona jako ważony kwadrat odległości między centroidami (średnimi) skupień A i B .

Twierdzenie

Niech $A = \{f_1, \dots, f_{n_A}\}$ i $B = \{g_1, \dots, g_{n_B}\}$ będą podzbiorami przestrzeni L^2 .
Wówczas

$$d(A, B) = \frac{2n_A n_B}{n_A + n_B} \|\bar{f} - \bar{g}\|^2,$$

gdzie \bar{f} jest centroidem skupienia A , a \bar{g} jest centroidem skupienia B .

Metoda k-lokalnych median Szlachtowska i Kosiorowski (2017)

- Globalna vs. lokalna funkcja głębi – Paindaveine i Van Bever (2013)
- Lokalność głębi funkcjonalnej – Squera et al. (2016); nasze propozycje na bazie podejścia Paindaveine i Van Bever (2013)
- Odporna analiza skupisk z wykorzystaniem głębi lokalnych

Algorytm k -lokalnych median dla danych funkcjonalnych

Wybór parametrów:

- 1 Wybieramy parametr k , gdzie oznacza wymaganą liczbę skupień.
- 2 Wybieramy parametr β (domyślna wartość β to 0,2). Za pomocą parametru β możemy rozpatrzeć problem na różnych „poziomach rozdzielczości”, tj. zmieniając wartość parametru β możemy kontrolować dokładność podziału. Z drugiej strony obliczone wartości lokalnej głębi dla wszystkich punktów są pomocne przy wyborze odpowiedniej liczby skupień.

Algorytm k -lokalnych median dla danych funkcjonalnych cd.

W rozważanym algorytmie grupowania rozważamy dwie miary podobieństwa/bliskości. Jako odporną miarę podobieństwa użyjemy lokalną zmodyfikowaną głębnię pasma. Kryterium oparte na lokalnej głębi stanowi pierwszą część funkcji celu, w której mierzymy jakość grupowania ze względu na odporność

$$F = \sum_{i=1}^k \sum_{j \in C_i} LD^\beta(f_j, P_{n,i}),$$

gdzie $P_{n,i}$ oznacza rozkład empiryczny w i -tym skupieniu. Natomiast jako miarę przydziału do skupień będziemy używać odległości L_2 .

W pierwszym kroku wybieramy wartości parametrów k i β , a następnie obliczamy wartości lokalnej głębi MBD dla wszystkich obserwacji w zbiorze danych. Następnie wybieramy k centroidów c_1, \dots, c_k spełniających warunki

$$\sum_{i=1}^k LD^{\beta}(c_i, P_n) \rightarrow \max$$

$$\sum_{i,j=1}^k \|c_i - c_j\|_2 \rightarrow \max$$

W drugim kroku tworzymy nowe skupienia w taki sposób, że dla każdej krzywej f obliczamy odległości w normie L_2 od wszystkich centroidów. Następnie na podstawie wartości $clus(f)$ przypisujemy obserwację do najbliższego jej centroidu, tj.

$$clus(f) = \{i : d(f, c_i) = \min\{d(f, c_1), d(f, c_2), \dots, d(f, c_k)\}\},$$

gdzie d oznacza odległość w przestrzeni L_2 . Jeśli do zbiorze $clus(f)$ należą co najmniej dwie indeksy, tj. istnieją co najmniej dwa centroidy, do których obserwacja f jest równo odległa, wówczas przypisujemy krzywą funkcjonalną do skupienia o niższej wartości indeksu.

W trzecim kroku dla nowo powstałych skupień obliczamy lokalne mediany MBD względem empirycznych rozkładów poszczególnych skupień. Mediana ta traktowana jest jako nowy centroid. Powtarzamy krok drugi i trzeci, aż centroidy nie zmieniają się lub dopóki tylko 1% punktów zmieni skupienie.

Przycięty algorytm k -lokalnych median dla danych funkcjonalnych

Inspiracja algorytm tclust (zob. algorytm tclust, Fritz, Garcia-Escudero, Mayo-Isacar, *J Stat Soft* 2006, Szlachtowska et. al. 2016)

Niech $d(f, g) = \|f - g\|_2$ oznacza odległość L_2 między funkcjami f i g . W każdej iteracji algorytmu, dla każdej funkcji f obliczamy jej odległość do wszystkich centroidów. W rezultacie otrzymujemy ciąg odległości, który możemy posortować rosnąco

$$d(f, c_{(1)}) \leq d(f, c_{(2)}) \leq \dots \leq d(f, c_{(k)}).$$

Niech $clus(f)$ oznacza indeks skupienia, do którego przypisano krzywą funkcjonalną f . Wówczas

$$clus(f) = \{i : d(f, c_i) = d(f, c_{(1)})\}.$$

Użytkownik wybiera parametr γ , tj. wielkość przycinania. Dla każdej obserwacji g obliczamy stopień przynależności do skupienia, do którego została dana obserwacja przypisana. W tym celu obliczamy odległość obserwacji od centroidu danego skupienia. Obliczone wartości ustawiamy malejąco

$$d(f, c_f) \geq d(g, c_g) \geq \dots \geq d(h, c_h).$$

Następnie odrzucamy γ część obserwacji o najwyższych wartościach stopnia afiliacji, tj. najbardziej odległe od centroidów skupień, do których zostały przypisane.

Ocena jakości grupowania

Dla tej metody podobnie jak w algorytmie tclust, można obliczyć współczynniki dyskryminacyjne dla każdej obserwacji (przyciętej i nieprzyciętej).

Jakość decyzji przypisania do skupienia dla nieprzyciętej obserwacji f do skupienia można ocenić porównując stopień przynależności $d(f, c_f) = d(f, c_{(1)})$ to drugiego najlepszego przypisania $d(f, c_{(2)})$. Stąd

$$DF(f) = \log \left(\frac{d(f, c_{(2)})}{d(f, c_f)} \right) = \log \left(\frac{d(f, c_{(2)})}{d(f, c_{(1)})} \right)$$

dla nieprzyciętej obserwacji f .

Aby obliczyć wartości współczynników dyskryminacyjnych dla krzywych przyciętych ustawiamy obserwacje $f_{(1)}, \dots, f_{(n)}$ w kolejności malejącej ze względu na ich wartości $d(f_{(i)}, c_{f_{(i)}})$.

Nie jest trudno zauważyć, że $f_{(1)}, \dots, f_{(\lceil \gamma n \rceil)}$ to przycięte obserwacje, które nie są przypisane do żadnego skupienia.

Niemniej jednak dla przyciętej obserwacji f można obliczyć stopień przynależności $d(f, c_f)$ do najbliższego jej skupienia.

Zatem jakość decyzji o przycięciu tej obserwacji można ocenić porównując $d(f, c_f)$ oraz

$$d\left(f_{(\lceil \gamma n \rceil + 1)}, c_{f_{(\lceil \gamma n \rceil + 1)}}\right),$$

gdzie $f_{(\lceil \gamma n \rceil + 1)}$ oznacza nieprzyciętą obserwację o największej wartości $d(\cdot, c_{(\cdot)})$.
Stąd

$$DF(f) = \log \left(\frac{d(f, c_f)}{d\left(f_{(\lceil \gamma n \rceil + 1)}, c_{f_{(\lceil \gamma n \rceil + 1)}}\right)} \right)$$

dla przyciętej obserwacji f .

W rezultacie czynniki dyskryminacyjne $DF(f)$ są uzyskiwane dla każdej obserwacji w zbiorze danych, czy jest to obserwacja przycięta, czy nie. Obserwacje o małych wartościach $DF(f)$ (czyli o wartościach zbliżonych do zera) wskazują na wątpliwe przypisanie do skupienia lub na wątpliwą decyzję dotyczące przycinania.

Odporne klasyfikatory w przypadku FDA

Rozumienie odporności klasyfikatora w dalszym ciągu stanowi otwarty problem. Spośród ciekawych propozycji warto zauważyć Cuevas i Romo (1993), którzy badali jakościową odporność bootstrapowej aproksymacji estymatora o postaci funkcjonálu statystycznego T . Pokazali, że jednostajna ciągłość funkcjonálu statystycznego T jest wystarczającym warunkiem na jakościową odporność jego bootstrapowej aproksymacji. Niech $L_n(F) = L_n(T; F)$ będzie rozkładem próbkowym statystyki $T_n(X_1, \dots, X_n)$ gdzie próba pochodzi z F , oraz $L[L_n(F_n)]$ jest rozkładem próbkowym uogólnionej statystyki $L_n(F_n)$ w przestrzeni odpowiednich miar probabilistycznych. W definicji Cuevasa i Romo (1993) "dla danego ciągu T_n statystyk generowanych przez funkcjonal statystyczny T , ciąg jego bootstrapowych oszacowań $\{L_n(F_n)\}$ jest jakościowo odporny w F wtedy i tylko wtedy gdy, ciąg transformacji $\{G \rightarrow L[L_n(G_n)]\}$ jest asymptotycznie równociągły." Definicja ta została wykorzystana m. in. przez Christmanna i in. oraz przez nas (Kosiorowski i in. 2018)

Definicja 3: Przez **jakościową odporność** rozumiemy równociągłość rozkładu statystyki, gdy zwiększa się wielkość próby.

Zatem jakościowa odporność klasyfikatora wiąże się z ciągłością względem słabej topologii stosownej przestrzeni, pewnej charakterystyki jego rozkładu (por. Rudin, 1991).

Klasyfikatory dla danych funkcjonalnych

- W obrębie FDA stosuje się m. in. metody bazujące na zasadzie k - najbliższych sąsiadów. Obserwacja klasyfikowana jest do klasy obiektów których jest najwięcej w pewnym sąsiedztwie tej. Kluczową i wciąż nierozwiązaną kwestią jest wybór metryki za pomocą której definiuje się sąsiedztwa oraz liczba sąsiadów branych pod uwagę k (Ferraty & Vieu, 2006); metody przestrzeni Hilberta z jądrem reprodukującym (RKHS) (por. Schoelkopf i Smola, 2002); metody bazujące na koncepcji najbliższego centroidu. W ramach tej rodziny metod przyporządkowujemy obserwacji etykietę klasy, której centroid jest jej najbliższy. Jako centroid można przyjąć funkcjonalną średnią bądź pewną funkcjonalną medianę. Kolejna grupa metod wykorzystuje statystyczne funkcje głębokości (por. Cuevas i Fraiman, 2009).
- Kosiorowski, Bocian (2015), Kosiorowski, Mielczarek, Rydlewski (2017) klasyfikator SVM via DD plot.
- Kosiorowski, Mielczarek, Rydlewski (2018c) nowy klasyfikator nie nawiązujący do metody reprodukującego jądra przestrzeni Hilberta.

Zalety naszej propozycji

- Przestrzeń L_2 nie jest przestrzenią z jądrem reprodukującym. Najczęściej w przypadku klasyfikatorów dla danych funkcjonalnych mamy do czynienia z tzw. "trickiem jądrowym" tzn. z przekształceniem danych z przestrzeni L_2 na dane z pewnej przestrzeni, najczęściej skończonej wymiarowej przestrzeni Hilberta z jądrem reprodukującym.
- Dobór tej przestrzeni z jądrem reprodukującym ma charakter arbitralny i nie jest znana autorom odporna metoda "dopasowania" jądra do danych funkcjonalnych. Klasyfikacja danych zależy od "właściwego doboru" jądra.
- Nasza metoda jest pozbawiona tych wad, ponieważ zaproponowany klasyfikator nie dokonuje żadnych transformacji danych funkcjonalnych.
- Klasyfikator jest bardzo łatwy w implementacji i stabilny numerycznie.

Zaproponowany przez autorów klasyfikator dla danych funkcjonalnych jest liniową kombinacją funkcyjonałów $g_j: L_2(\Omega) \rightarrow \mathbb{R}$ takich, że

$$g_j(X_k) = \delta_{kj},$$

czyli tzw. funkcyjonałów normujących.

Klasyfikator g jest określony wzorem

$$g = \sum_{j=1}^m Y_j (1 - Y_j b) g_j,$$

spełnia on żądane własności.

Okazują się więc, że wystarczy wskazać funkcjonały g_i . Przyjmijmy następujące oznaczenia dla dowolnego układu wektorów $\{W_1, \dots, W_m\}$ z przestrzeni Hilberta z iloczynem skalarzym $\langle \cdot, \cdot \rangle$

$$M(W_1, \dots, W_m) = \det [\langle W_i, W_j \rangle]_{i=1\dots m, j=1\dots m}.$$

Liczba $M(W_1, \dots, W_m)$ jest więc wyznacznikiem macierzy Gramma dla układu wektorów $\{W_1, \dots, W_m\}$.

Przypomnijmy, że układ wektorów jest liniowo niezależny wtedy i tylko wtedy liczba $M(W_1, \dots, W_m)$ jest większa od zera, jest natomiast równa zero jeżeli układ wektorów $M(W_1, \dots, W_m)$ jest liniowo zależny.

Przy tak przyjętych oznaczeniach szukane funkcjonały g_i można określić następującym wzorem dla dowolnego $Y \in L_2(\Omega)$

$$g_i(Y) = \frac{M(X_1, \dots, X_{i-1}, Y, X_{i+1}, \dots, X_m)}{M(X_1, \dots, X_m)}.$$

Tak więc wyznaczenie klasyfikatora dla danych funkcjonalnych sprowadza się do wyliczenia odpowiednich wyznaczników. Implementacja klasyfikatora jest więc stosunkowo łatwa.

Przeprowadzone przez autorów badania empiryczne potwierdziły większą skuteczność tego klasyfikatora od klasyfikatorów jądrowych dla danych internetowych oraz dla danych obrazujących zapotrzebowanie na energię elektryczną oraz dla danych pochodzących z pewnego serwisu internetowego podzielonego na podserwisy.

Np. dla serwisu internetowego ryzyko empiryczne naszej metody wyniosło 20%, podczas gdy dla metod RHKS wyniosło 16% (jądro gaussowskie), 18% (jądro wielomianowe) i 24% (jądro sigmoidalne), dla metod bazujących na głębi (głębia Fraimana-Muniza) aż 46%.

Zaletą naszej metody była tu szybkość - była około 10 razy szybsza.

Zanieczyszczenie powietrza w Katowicach cd.

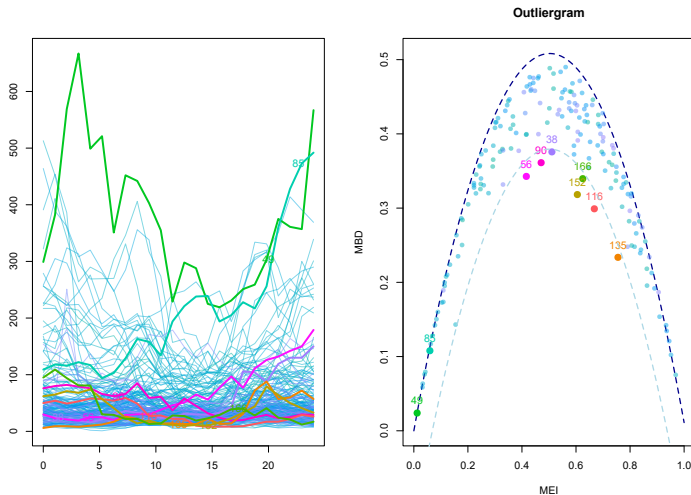


Figure 14: Outlierogram

Dobowe zanieczyszczenie powietrza w Katowicach cd.

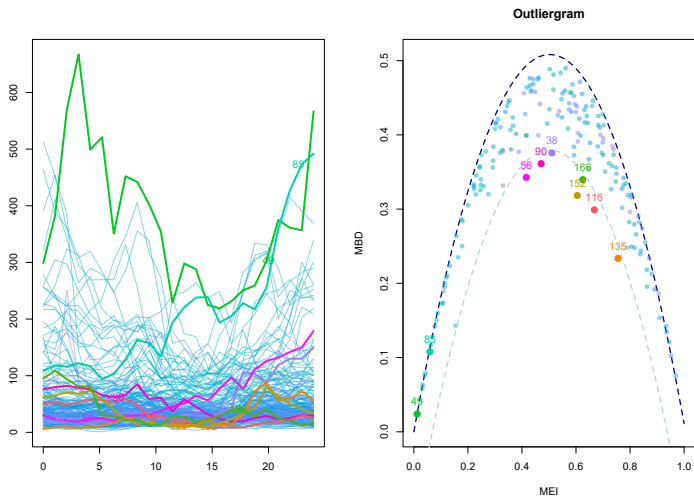


Figure 15: Outlierogram

Zanieczyszczenie powietrza w Katowicach cd.

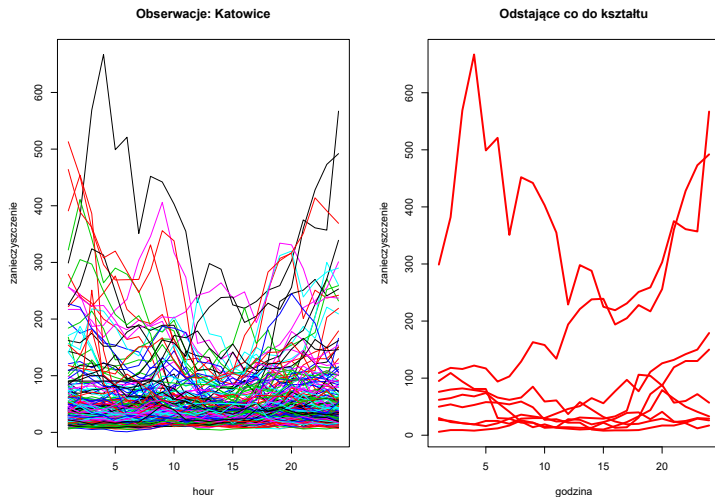


Figure 16: Outliers typu "shape"

Zanieczyszczenie powietrza w Katowicach cd.

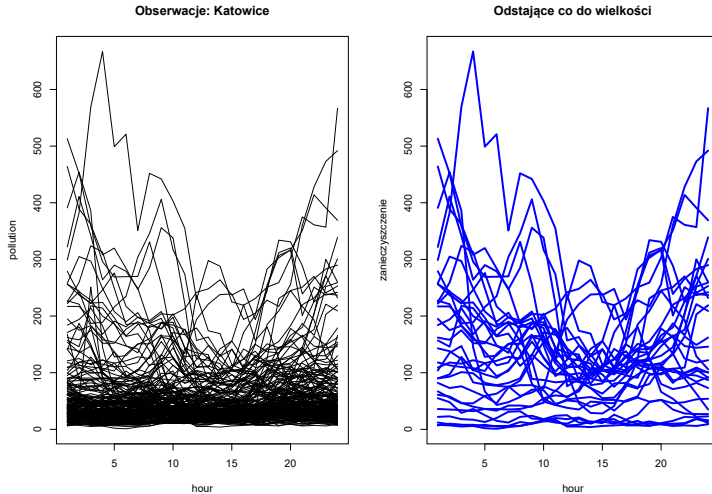
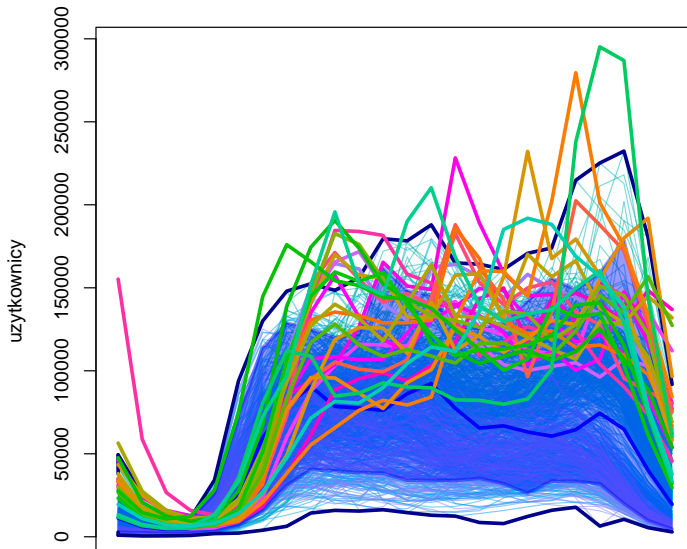


Figure 17: Outliers typu "magnitude"

Monitorowanie zachowań internautów

Uzytkownicy intrnetu, skorygow. boxplot



Monitorowanie zachowań internautów

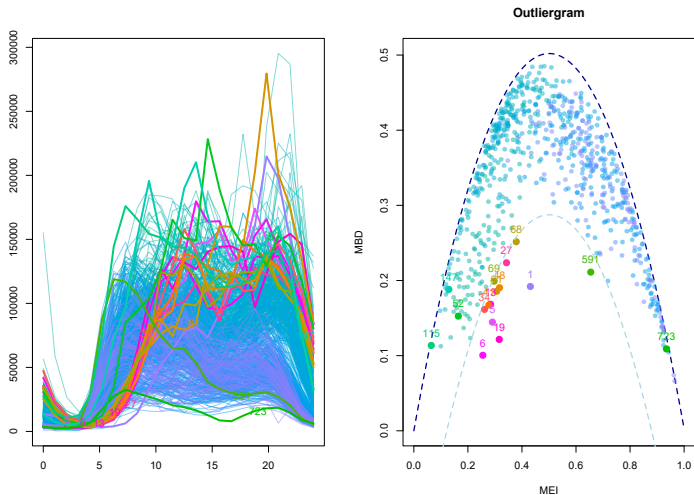


Figure 19: Outliergram, użytkownicy Internetu

Monitorowanie zachowań internautów

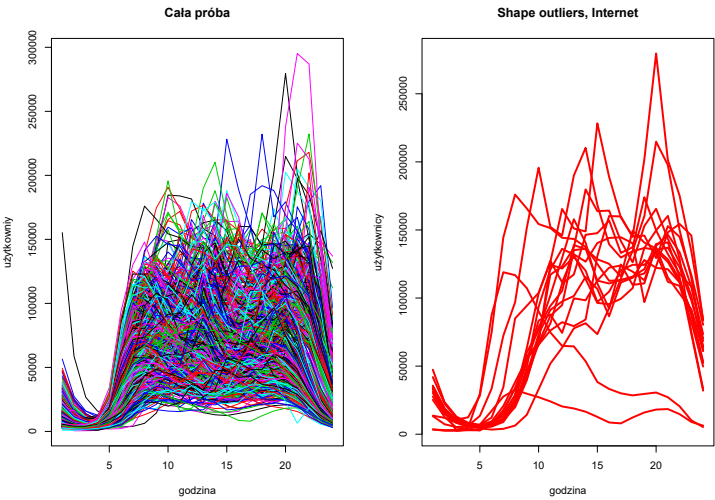


Figure 20: Outliery typu shape, użytkownicy Internetu

Monitorowanie zachowań internautów

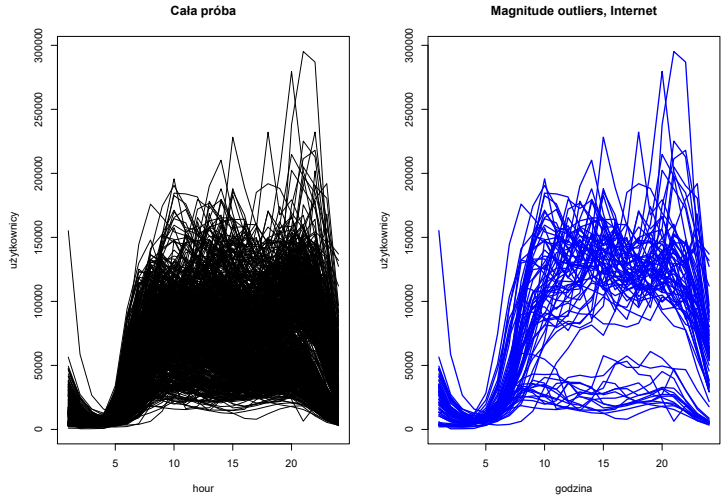


Figure 21: Outliery typu magnitude, użytkownicy Internetu

Predykcja HFTS w maksymalizacji dobrobytu społecznego

Zanieczyszczenie powietrza ma negatywny wpływ na zdrowie człowieka. Substancje niebezpieczne mogą wchodzić ze sobą w interakcje, wpływ może zależeć od grupy wiekowej i pory dnia.

NASZ CEL: maksymalizacja "zagregowanej" użyteczności pewnej lokalnej społeczności w pewnym okresie czasu, co symbolicznie zapiszemy jako (por. Fleurbaey i Maniquet, 2011) :

$$U_{Total} = \sum_{i=1}^{365} \int_{[0^{00}, 24^{00}]} U_i(W_{PM10}(t), C_{PM10reduc}(t)) dt, \quad (7)$$

gdzie i wskazuje indeks doby, W_{PM10} oznacza społeczny dobrobyt związany z redukcją emisji pyłu PM10 (dodatnie i ujemne efekty zewnętrzne wyrażone w ustalonej walucie) and C_{PM10} oznacza koszt redukcji emisji pyłu PM10 wyrażony w ustalonej walucie.

Zakładamy, że

$$W_{PM10} = F(Air_{qual}, ENV_{polit}, INF_{qual}, POP_{param}),$$

oraz

$$C_{PM10} = G(C_{fixed}, C_{var}, C_{political}, Pred_{qual}).$$

Oznacza to, że dobrobyt związany z pyłem PM10 jest funkcją jakości powietrza (wycenioną na podstawie kosztów hospitalizacji z powodu chorób płuc, wydatków związanych z alergiami (Air_{qual}), przyjazności otoczenia lokalnego (ENV_{polit}), jakości lokalnego systemu informowania o jakości powietrza i zagrożeniach (INF_{qual}) oraz socio-demograficznych parametrów społeczności (POP_{param}).

Predykcja HFTS w maksymalizacji dobrobytu społecznego

Koszty związane z redukcją pyłu PM10 wiążą się z **kosztami stałymi** obejmującymi inwestycje w nowe technologie (C_{fixed}), i **kosztami zmiennymi** związane ze zmianami warunków pogodowych, które skutkują większym bądź mniejszym popytem na energię elektryczną i ciepłą (C_{var}), *koszty polityczne* związane z przechodzeniem z popularnych systemów bazujących na węglu na "czyste" systemy opierające się np. na energii jądrowej i **koszty związane z jakością predykcji zanieczyszczenia powietrza** ($Pred_{qual}$).

Przykład 4 HFTS w predykcji jakości powietrza na Śląsku

W badaniach empirycznych skupiliśmy się na obszarze Śląska, gdzie zlokalizowanych jest 28 stacji automatycznego pomiaru jakości powietrza nadzorowanych przez Wojewódzki Inspektorat Ochrony Środowiska (WIOŚ, w Katowicach).

Analizowaliśmy dane pochodzące z 5 stacji pomiarowych. Na surowe dane składają się po 181 trajektorii dla każdej z pięciu stacji pokazujące koncentrację pyłu PM10 w atmosferze w $\mu g/m^3$.

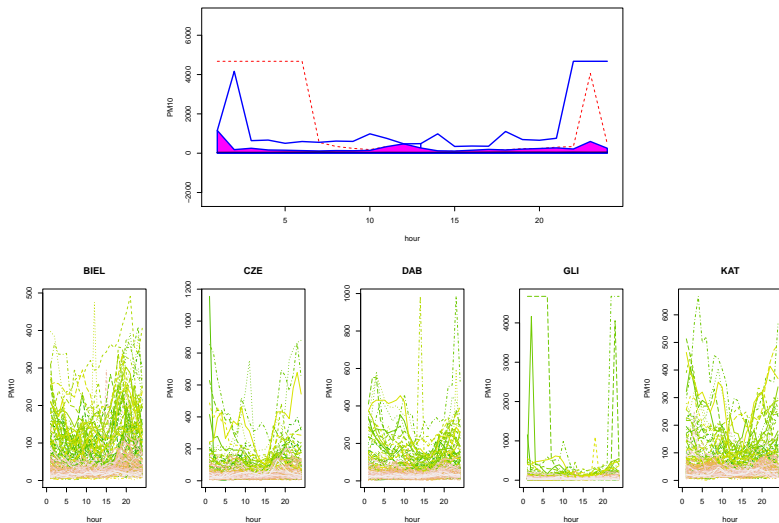


Figure 23: Surowe dane 181 trajektorii koncentracji PM10 dla pięciu stacji w atmosferze w $\mu\text{g}/\text{m}^3$. Funkcjonalny boxplot przedstawia wszystkie krzywe.

Przykład "naiwny" – ruchoma średnia punktowa

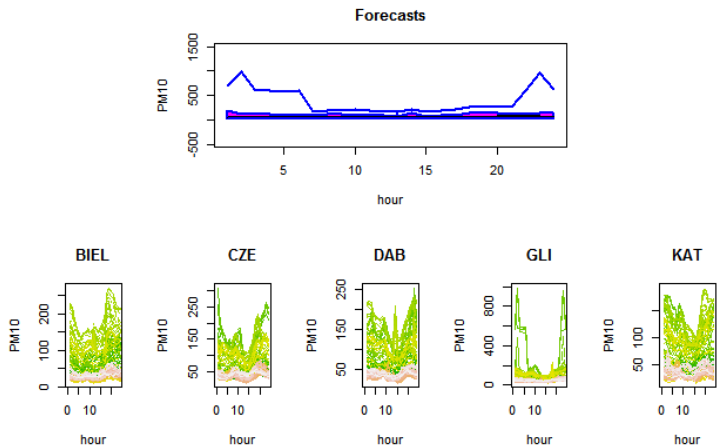


Figure 24: Przewidywania koncentracji PM10 uzyskane za pomocą metody ruchomej średniej funkcjonalnej w $\mu g/m^3$ (długość okna wynosi 10 obs.) dla pięciu stacji.

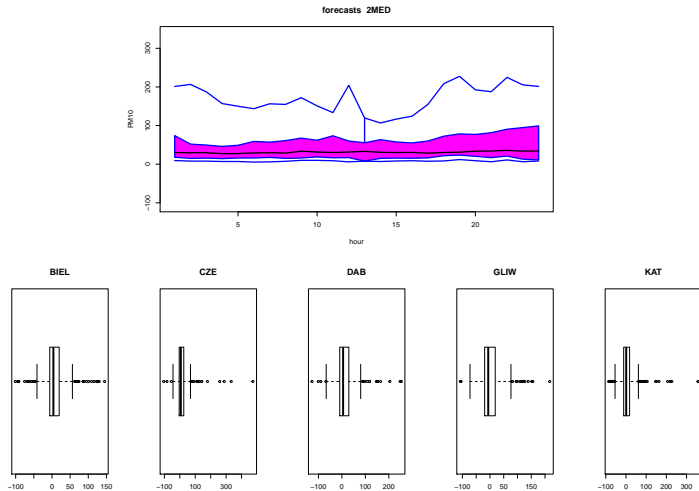
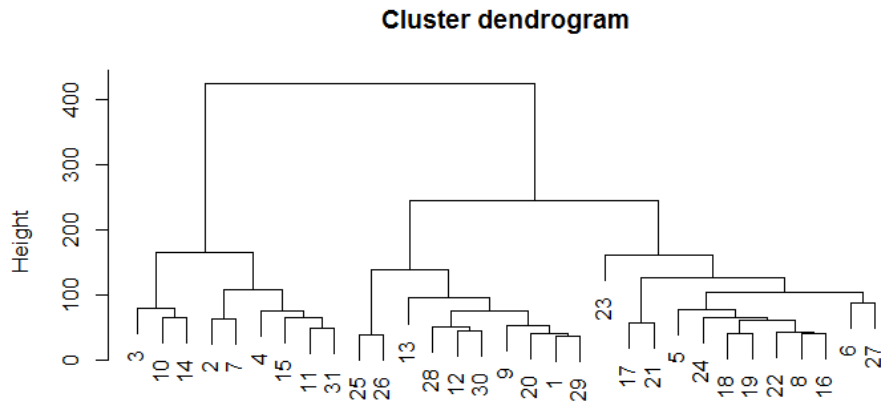


Figure 25: Prognozy koncentracji PM10 uzyskane za pomocą metody podwójnej mediany MBD w $\mu g/m^3$ (długość okna wynosi 10 obs.) dla pięciu stacji. Boxploty przedstawiają uśrednione miary różnic pomiędzy obserwowanymi i przewidywanymi trajektoriami.

Przykład 5 AS dobowego zanieczyszczenia powietrza w Krakowie

Do analizy wybrano również **dane dotyczące jakości powietrza ze stacji Aleja Krasińskiego w Krakowie** za okres od 1 stycznia do 31 grudnia 2015 r. Wyniki pomiarów siedmiu substancji są prezentowane co godzinę na stronie <http://monitoring.krakow.pios.gov.pl/>.

Należy podkreślić, że w/w okresie nie ma wszystkich pomiarów. Zatem pojawia się **problem brakujących danych**. Dane uzupełniono na dwa sposoby: uzupełniając średnią z danej godziny oraz uzupełniając medianą z danej godziny. W obu wariantach zastosowano algorytm k-średnich dla danych funkcjonalnych w celu zbadania jego odporności na problem brakujących danych.



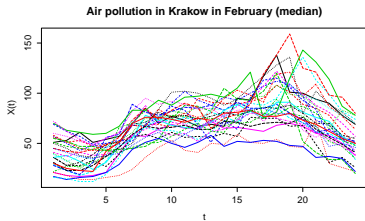
Pollination of nitrogen dioxide in December 2015, Kraków
hclust (*, "ward.D2")

Figure 26: Zanieczyszczenie nitrogen dioxide w grudniu 2015, Kraków.

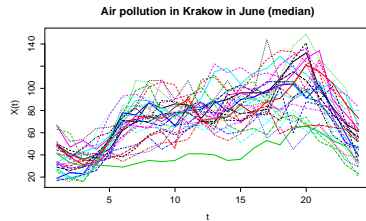
Tabela przedstawia liczbę obserwacji w skupisku w zależności od zadanej liczby skupisk k od 2 do 5 (pomiaru w grudniu 2015).

k	1	2	3	4	5
2	22	9			
3	10	9	12		
4	10	6	3	12	
5	10	6	3	11	1

Table 1: Liczba obserwacji w skupisku .Metoda – metoda Warda „Ward.D2 ”



(a) Luty 2015 roku



(b) Czerwiec 2015 roku

Figure 28: Krzywe funkcjonalne obrazujące zapylenie dla poszczególnych dni, w których brakujące dane zostały uzupełnione medianą dla danej godziny.

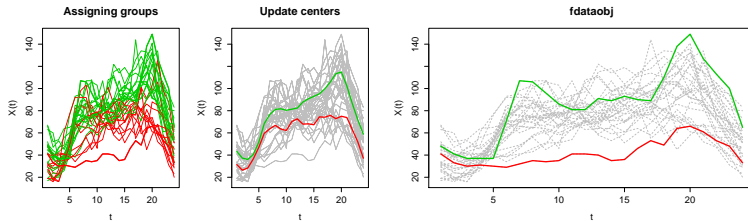


Figure 29: Środki skupień dla krzywych funkcjonalnych obrazujące zapylenie dla poszczególnych dni w czerwcu, w których brakujące dane zostały uzupełnione **medianą wartości zapylenia** w danej godzinie, parametr **$k=2$** .

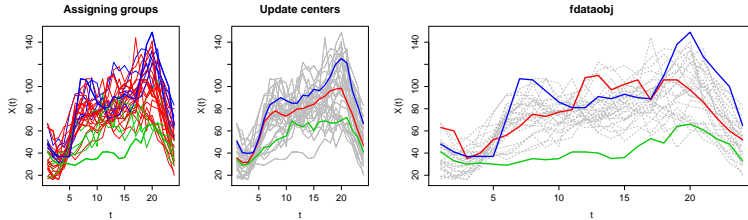


Figure 30: Środki skupień dla krzywych funkcjonalnych obrazujące zapylenie dla poszczególnych dni w czerwcu, w których brakujące dane zostały uzupełnione **medianą wartości zapylenia** w danej godzinie, parametr **k=3**.

Ustawiając $k = 5$ zauważamy, że obserwacja 23 tworzy skupisko. To Wigilia (Świąt Bożego Narodzenia). Można zauważyć największą koncentrację w południe zmniejszającą się po południu. Około godziny 18 następuje dalszy wzrost koncentracji NO₂ w powietrzu, jednak istotnie mniejszy niż w godzinach rannych.

- Pierwsze skupisko to weekendy, święta i okres tuż po świętach.
- Trzecie skupisko, to dni przed świętami odznaczające się szczególnie wzmożonym ruchem związanym z wyjazdami do rodzin spoza Krakowa.

Konkluzje:

- Najmniejsze stężenie tlenków azotu obserwujemy w weekendy i święta (pierwsze skupisko).
- Największą koncentrację obserwujemy w dni robocze (drugie skupisko).
- Tuż przed Świętami, w trakcie wyjazdów, zanieczyszczenie pozostaje na średnim poziomie.

Badanie potwierdziło hipotezę, że koncentracja tlenków azotu w atmosferze wiąże się z natężeniem ruchu na krakowskich drogach.

Pyły zawieszane – algorytm k lokalnych median

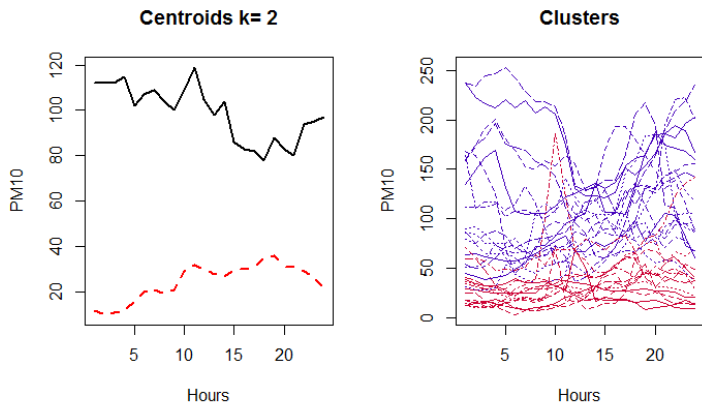











Figure 31: Mediany funkcjonalne poszczególnych skupień. Metoda k lokalnych median dla danych funkcjonalnych, grudzień 2016 r., pył zawieszony PM10.







Podsumowanie—wyzwania dla zastosowań FDA w ekonomii

- Wyzwanie 1: Sensowne merytoryczne modelowanie zjawisk ekonomicznych za pomocą modeli FDA oraz sensowna merytorycznie interpretacja wyników procedur FDA.
- Wyzwanie 2: Postępowanie w przypadku występowania różnego rodzaju braków w danych funkcjonalnych?
- Wyzwanie 3: Znalezienie teoretycznych podstaw dla metod oceny niepewności prognozy dla FTS - np. teoria dla ruchomej mediany funkcjonalnej.
- Wyzwanie 4: Znalezienie odpowiednika outlierogramu dla innych niż MBD głębi funkcjonalnych.
- Wyzwanie 5: Prognozowanie niestacjonarnych FTS.
- Wyzwanie 6: Wyodrębnienie w klasie "shape outliers" podtypów szczególnie użytecznych w analizie cykli gospodarczych.
- Wyzwanie 7: Koncepcja odporności w przypadku AS i klasyfikacji obiektów.

Wybrane pozycje literatury

-  Arribas-Gil A, Romo J (2014) Shape outlier detection and visualization for functional data: the outliergram. *Biostatistics*, 15(4):603–619.
-  Bosq D (2000). *Linear processes in function spaces*. Springer.
-  Górecki, T., Hörmann, S., Horváth, L. & Kokoszka, P. (2017). Testing Normality of Functional Time Series. *Journal of Time Series Analysis*. doi:10.1111/jtsa.12281.
-  Górecki T (2015) „Wprowadzenie do analizy danych funkcjonalnych”, referat plenarny na 7th International Scientific Conference Faculty of Management Cracow University of Economics.
-  Górecki T., Krzyśko M.K., Waszak Ł., Wołyński W., (2016), Selected statistical methods of data analysis for multivariate functional data, *Statistical Papers*, DOI10.1007/s00362-016-0757-8
-  Horváth L, Kokoszka P (2012) *Inference for functional data with applications*. Springer, New York
-  Kosiorowski D (2012). „Wstęp do statystyki odpornej: kurs z wykorzystaniem środowiska R.” Wydawnictwo Uniwersytetu Ekonomicznego w Krakowie.
-  Kosiorowski D, Zawadzki Z (2017) *DepthProc: An R Package for Robust Exploration of Multidimensional Economic Phenomena*. arXiv:1408.4542
-  Kosiorowski D, Rydlewski JP , Snarska M (2017a) Detecting a structural change in functional time series using local Wilcoxon statistic. *Stat Papers*. DOI 10.1007/s00362-017-0891-y

Wybrane pozycje literatury

-  Kosiorowski D, Mielczarek D, Rydlewski JP, (2018a) Forecasting of a Hierarchical Functional Time Series on Example of Macromodel for the Day and Night Air Pollution in Silesia Region — A Critical Overview. Central European Journal of Economic Modelling and Econometrics 10: 53-73.
-  Kosiorowski D., Rydlewski J. P., & Zawadzki Z., (2018b). Functional outliers detection by the example of air quality monitoring. Statistical Review (in Polish, to appear).
-  Kosiorowski D, Mielczarek D, Rydlewski JP, Snarska M, (2018c). Generalized exponential smoothing in prediction of hierarchical time series. Statistics in Transition, June 2018 Vol. 19, No. 2, pp. xxx-xxx. (to appear).
-  Kosiorowski D, Mielczarek D, Rydlewski JP, (2018d). Outliers in Functional Time Series – Challenges for Theory and Applications of Robust Statistics, In M. Papież & S. Śmiech (eds.), The 12 th Professor Aleksander Zeliaś International Conference on Modelling and Forecasting of Socio-Economic Phenomena, Conference Proceedings, Cracow: Foundation of the Cracow University of Economics, pp. 209–218.
-  López-Pintado S, Romo J (2009) On the concept of depth for functional data. J. Amer. Statist. Assoc. 104: 718-734.
-  Nagy S, Gijbels I, Hlubinka D (2017) Depth-Based Recognition of Shape Outlying Functions. Journal of Computational and Graphical Statistics, DOI: 10.1080/10618600.2017.1336445

