

Wygładzanie i klasyfikacja krzywych funkcjonalnych

Karolina Sopata

6 czerwca 2022

Wygładzanie splajnu przy użyciu "kary na szorstkość"

- ▶ rozszerzenia bazy mogą stanowić dobre przybliżenie danych funkcjonalnych
- ▶ popularną i efektywną metodą jest metoda najmniejszych kwadratów
- ▶ podejście oparte na karze na „szorstkość” jest jedną z metod aproksymacji danych funkcjonalnych
- ▶ metoda ta daje lepsze wyniki zwłaszcza przy estymacji pochodnych
- ▶ metoda oparta na optymalizacji kryterium dopasowania

Określanie szorstkości

$[D^2x(t)]^2$ - krzywizna funkcji

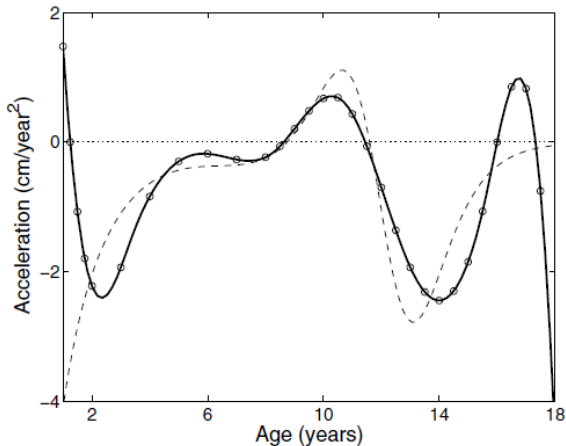
Naturalna miara szorstkości:

$$\text{PEN}_2(x) = \int [D^2x(s)]^2 ds$$

Generalnie:

$$\text{PEN}_m(x) = \int [D^m x(s)]^2 ds$$

Przykład



Rysunek: Krzywa przyspieszenia wzrostu, źródło: Ramsay, Functional Data Analysis, Fig. 4.3.

Kryterium wyboru parametru wygładzającego

$x(t)$ - wektor wartości otrzymanych przy użyciu funkcji x na wektorze argumentów t

Kompromis pomiędzy gładkością a dopasowaniem do danych:

$$\text{PENSSE}_\lambda(x|y) = [y - x(t)]'W[y - x(t)]^2 + \lambda \times \text{PEN}_2(x)$$

de Boor (2002)

Krzywa x , która daje najmniejszą wartość powyższego kryterium jest splajnem sześciennym z węzłami w punktach t_j , dla $j=1, \dots, n$.

Wektor współczynników \mathbf{c}

1. Przypadek bez parametru wygładzającego.
Postać funkcji:

$$\mathbf{x}(t) = \sum_k^K c_k \phi_k(t) = \mathbf{c}' \boldsymbol{\phi}(t) = \boldsymbol{\phi}'(t) \mathbf{c},$$

Estymacja współczynników:

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

Predykcja zmiennej objaśnianej:

$$\hat{\mathbf{y}} = \boldsymbol{\Phi} (\boldsymbol{\Phi}' \mathbf{W} \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}' \mathbf{W} \mathbf{y} = \mathbf{S}_{\phi} \mathbf{y},$$

Wektor współczynników c

2. Przypadek z parametrem wygładzającym.

$$\text{PEN}_m(x) = \int [D^m x(s)]^2 ds$$

$$= \int [D^m c' \phi(s)]^2 ds$$

$$= \int c' D^m \phi(s) D^m \phi'(s) c ds$$

$$= c' \left[\int D^m \phi(s) D^m \phi'(s) ds \right] c$$

$$= c' R c ,$$

gdzie:

$$R = \int D^m \phi(s) D^m \phi'(s) ds .$$

2. Przypadek z parametrem wygładzającym, c.d.

$$\text{PENSSSE}_m(\mathbf{y}|\mathbf{c}) = (\mathbf{y} - \Phi\mathbf{c})'\mathbf{W}(\mathbf{y} - \Phi\mathbf{c}) + \lambda\mathbf{c}'\mathbf{R}\mathbf{c}$$

$$-2\Phi'\mathbf{W}\mathbf{y} + \Phi'\mathbf{W}\Phi\mathbf{c} + \lambda\mathbf{R}\mathbf{c} = 0,$$

Estymator wektora współczynników \mathbf{c} :

$$\hat{\mathbf{c}} = (\Phi'\mathbf{W}\Phi + \lambda\mathbf{R})^{-1}\Phi'\mathbf{W}\mathbf{y}$$

Metoda uogólnionej walidacji krzyżowej

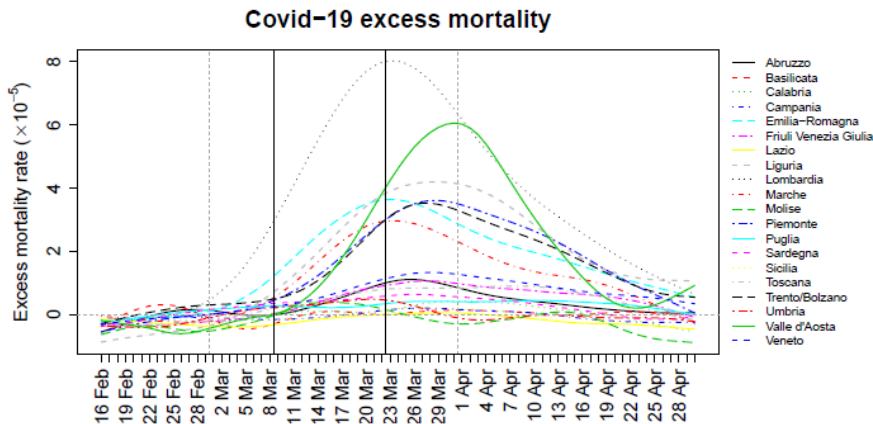
Kryterium:

$$\text{GCV}(\lambda) = \frac{n^{-1} \text{SSE}}{[n^{-1} \text{trace}(\mathbf{I} - \mathbf{S}_{\phi, \lambda})]^2}$$

Równoważna forma:

$$\text{GCV}(\lambda) = \left(\frac{n}{n - df(\lambda)} \right) \left(\frac{\text{SSE}}{n - df(\lambda)} \right)$$

Wygładzanie krzywych - przykład



Rysunek: Wygładzone krzywe śmiertelności na Covid-19 we Włoszech, źródło: M. Cremona, Probabilistic K-mean with local alignment for clustering and discovery in functional data

Klasyfikacja metodą K-means

- ▶ Ustalenie liczby klas.
- ▶ Inicjalizacja centroidów: tasujemy zbiór danych, a następnie wybieramy losowo K punktów.
- ▶ Iteracje: do momentu, aż nie będzie żadnych zmian w centroidach, tzn. przydział punktów do klastrów nie będzie się zmieniał
 - ▶ obliczenie sumy kwadratów odległości między punktami a wszystkimi centroidami
 - ▶ przypisanie każdego punktu danych do najbliższego klastra (centroidu)
 - ▶ obliczenie centroidów klastrów, jest to średnia wartość wszystkich punktów należących do każdego klastra

K-means algorytm

Funkcja celu:

$$J = \sum_{i=1}^m \sum_{k=1}^K w_{ik} \|x^i - \mu_k\|^2$$

K-means algorithm

Expectation step:

$$\frac{\partial J}{\partial w_{ik}} = \sum_{i=1}^m \sum_{k=1}^K \|x^i - \mu_k\|^2$$
$$\Rightarrow w_{ik} = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_j \|x^i - \mu_j\|^2 \\ 0 & \text{otherwise.} \end{cases}$$

Maximization step:

$$\frac{\partial J}{\partial \mu_k} = 2 \sum_{i=1}^m w_{ik} (x^i - \mu_k) = 0$$
$$\Rightarrow \mu_k = \frac{\sum_{i=1}^m w_{ik} x^i}{\sum_{i=1}^m w_{ik}}$$

Klasyfikacja metodą fuzzy C-means

- ▶ ustalenie c - liczby klastrów, wybór parametru rozmycia m (zwykle m $1.25 < m < 2$) i inicjalizacja macierzy partycji $U^{(0)}$
- ▶ wyliczenie centroidów $C^{(k)}$ przy ustalonym $U^{(k)}$

$$c_j = \frac{\sum_{i=1}^N u_{ij}^m \cdot x_i}{\sum_{i=1}^N u_{ij}^m}$$

- ▶ aktualizacja macierzy partycji $U^{(k)}, U^{(k+1)}$

$$u_{ij} = \frac{1}{\sum_{k=1}^c \left(\frac{\|x_i - c_j\|}{\|x_i - c_k\|} \right)^{\frac{2}{m-1}}}$$

Klasyfikacja metodą probKMA

- ▶ Rozważamy zbiór N krzywych (d -wymiarowych) $x_i : R \rightarrow R^d$, $i = 1, \dots, N$.
- ▶ Celem jest zidentyfikowanie K (d -wymiarowych) centrów klas $v_k : (0, c_k) \rightarrow R^d$.
- ▶ Dla każdej krzywej znajdujemy przesunięcie, tak aby minimalizować odległość każdej z nich od centrum klasy (v_k), gdzie funkcja przesunięcia jest klasy $W = \{h : t \rightarrow t + s, s \in \mathbb{R}\}$.
- ▶ Każdej rozważanej krzywej przypisujemy prawdopodobieństwo należenia do danej klasy. Funkcja prawdopodobieństwa zdefiniowana jest w następujący sposób:

$$p_k : \{x_1, \dots, x_N\} \rightarrow [0, 1]$$

gdzie $p_k(x_i) = p_{k,i}$.

Probabilistic K-means Algorithm

Spotykamy się z następującym problemem optymalizacji: chcemy znaleźć K centrów klas (v_1, \dots, v_K) , macierz prawdopodobieństw P oraz macierz przesunięć S , tak aby zminimalizować uogólnione kryterium najmniejszych kwadratów zdefiniowane w kontekście analizy danych funkcjonalnych.

Kryterium zdefiniowane jest następująco:

$$G_m(P, S, v_1, \dots, v_K) = \sum_{i=1}^N \sum_{k=1}^K (p_{k,i})^m d^2(\tilde{x}_{i,S_{k,i}}, v_k), \quad (1)$$

gdzie $m > 1$ jest pewnym ustalonym parametrem kontrolującym stopień rozmycia, natomiast $\tilde{x}_{i,S_{k,i}} = x \circ h_{k,i}$ przesuniętą krzywą; $d(\cdot, \cdot)$ jest odległością krzywej od centrum klasy.

Probabilistic K-means Algorithm

Inicjalizacja:

Zadanie liczby klas K oraz długości centroidów, wylosowanie początkowych prawdopodobieństw należenia do danej klasy oraz przesunięć.

Iteracje:

- ▶ wyliczenie centroidów $\hat{v}_k^{(it)}$
- ▶ wybór przesunięcia $\hat{s}_{k,i}^{(it)}$ oraz przypisanie krzywych do klas
- ▶ obliczenie prawdopodobieństwa $\hat{p}_{k,i}^{(it)}$ należenia krzywych do danej klasy

Kryterium stopu:

$$BC_k = -\log \left(\sum_{i=1}^N \sqrt{p_{k,i}^{(it)} p_{k,i}^{(it-1)}} \right)$$

Probabilistic K-means Algorithm

Rozważmy zbiór $B = \{i \in \{1, \dots, N\} \mid d(\tilde{x}_{i, \hat{s}_{k,i}}, \hat{v}_k) > 0, \forall k\}$ - zbiór indeksów krzywych, których odległość od każdego centrum klasy jest większa od 0. Przy ustalonych wartościach $\hat{v}_1, \dots, \hat{v}_K$ oraz \hat{S} estymator $\hat{P} = [\hat{p}_{k,i}]$, który globalnie minimalizuje kryterium (1) przyjmuje następującą postać:

$$\hat{p}_{k,i} = \left(\sum_{l=1}^K \left(\frac{d^2(\tilde{x}_{i, \hat{s}_{k,i}}, \hat{v}_k)}{d^2(\tilde{x}_{i, \hat{s}_{l,i}}, \hat{v}_l)} \right)^{\frac{1}{m-1}} \right)^{-1}, \quad (2)$$

dla $i \in B$ oraz:

$$\hat{p}_{k,i} = \begin{cases} 0, & \text{dla } k : d^2(\tilde{x}_{i, \hat{s}_{k,i}}, \hat{v}_k) > 0 \\ \in [0, 1], & \text{dla } k : d^2(\tilde{x}_{i, \hat{s}_{k,i}}, \hat{v}_k) = 0 \end{cases}, \quad (3)$$

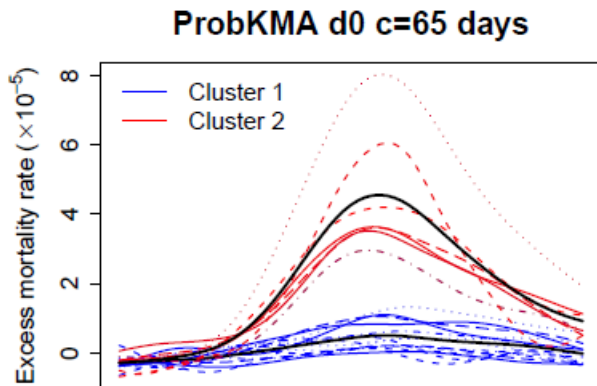
dla $i \notin B$.

Probabilistic K-means Algorithm

Przy ustalonych wartościach prawdopodobieństw należenia do danej klasy w macierzy \hat{P} oraz przesunięć w macierzy \hat{S} , unikalny estymator centrum klasy, który minimalizuje funkcję $G_m(\cdot, \hat{P}, \hat{S})$ ma postać:

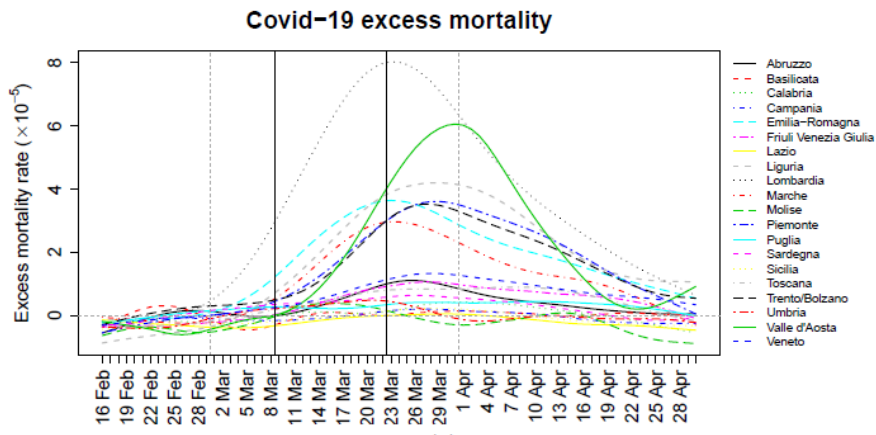
$$\hat{v}_k = \frac{\sum_{i=1}^N (\hat{p}_{k,i})^m \tilde{x}_{i, \hat{s}_{k,i}}}{\sum_{i=1}^N (\hat{p}_{k,i})^m}, \quad (4)$$

ProbKMA na danych Covid-19



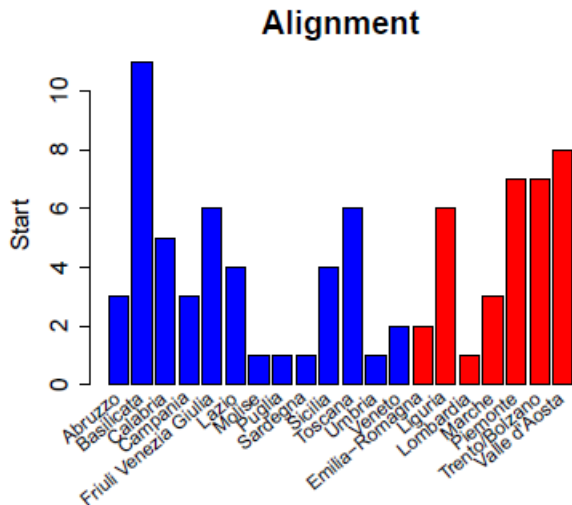
Rysunek: Klasyfikacja krzywych algorytmem probKMA, źródło: M. Cremona, Probabilistic K-mean with local alignment for clustering and discovery in functional data.

Wyglądanie krzywych



Rysunek: źródło: M. Cremona, Probabilistic K-mean with local alignment for clustering and discovery in functional data.

ProbKMA na danych Covid-19



Rysunek: Przesunięcia (liczba dni) krzywych, źródło: M. Cremona, Probabilistic K-mean with local alignment for clustering and discovery in functional data.

Dziękuję za uwagę!