

Warszawa, 16 stycznia 2019

# Metody ewolucyjne i uczenie się maszyn

## Projekt 2

### Dokumentacja końcowa

Na jednym z benchmarków dot. modelowania szeregów czasowych (np. NN3 lub M3) porównać metodę xgboost z modelem liniowym.

Prowadzący:  
dr inż. Paweł Zawistowski

Wykonali:  
Michał Kruszewski, Łukasz Błaszka

# Spis treści

1.	Cel . . . . .	3
2.	Założenia . . . . .	3
3.	Krótki wstęp teoretyczny . . . . .	3
3.1.	ARIMA . . . . .	3
3.2.	Metoda xgboost . . . . .	4
4.	Zbiór danych . . . . .	4
5.	Implementacja . . . . .	4
5.1.	ARIMA . . . . .	4
5.2.	XGBoost . . . . .	4
6.	Prezentacja wyników . . . . .	4

## 1. Cel

Celem projektu było praktyczne poznanie dowolnie wybranego modelu liniowego oraz metody xgboost, stosowanych do modelowania szeregów czasowych. Jako model liniowy wybrany został popularny model ARIMA (*ang. AutoRegressive Integrated Moving Average*). Realizacja celu oparta została na porównaniu obu metod na benchmark'u NN3.

## 2. Założenia

W specyfikacji wstępnej projektu znalazły się następujące założenia:

- Zadanie rozwiązane miało zostać w języku *R*.
- Jako benchmark testowy wybrany został nowszy benchmark NN3. Zawiera on mniej szeregów czasowych niż benchmark M3 (111 vs 3003), dlatego lepiej nadaje się do celów dydaktycznych (subiektywna ocena osób realizujących projekt).
- Jako model liniowy wybrany miał zostać popularny model ARIMA (*ang. AutoRegressive Integrated Moving Average*), uznawany za jedną z opcji do modelowania szeregów niestacjonarnych.
- W celu dopasowania najlepszego modelu ARIMA wykorzystana miała zostać metoda *auto.arima* pochodząca z pakietu *forecast*.
- W celu predykcji metodą xgboost wykorzystana miała zostać biblioteka *xgboost* (<https://github.com/dmlc/xgboost>).
- Ze względu na brak jakiegokolwiek doświadczenia z metodą xgboost, oraz liczbę jej parametrów, sposób optymalizacji określony miał zostać na etapie implementacji.

## 3. Krótki wstęp teoretyczny

### 3.1. ARIMA

Model ARIMA jest rozszerzeniem modelu ARMA, który poddawany jest różnicowaniu w celu usunięcia trendu. Model ARMA(p,q) dla stacjonarnego szeregu czasowego  $z_t$  ma postać:

$$z_t = \phi_1 z_{t-1} + \dots + \phi_p z_{t-p} + w_t + \psi_1 w_{t-1} + \dots + \psi_q w_{t-q} \quad (1)$$

gdzie  $z_t$  jest szeregiem stacjonarnym o zerowej średniej,  $\phi_i, \psi_i \in \mathbb{R} (\phi_p \neq 0, \psi_q \neq 0)$ , a  $w_i \sim N(0, \sigma^2)$ .

Różnicowanie potrzebne jest do osiągnięcia stacjonarności szeregu i polega na zastąpieniu szeregu czasowego  $z_t$  nowym szeregiem czasowym  $z'_t$  takim że:

$$z'_t = \nabla z_t = z_t - z_{t-1} \quad (2)$$

Trzeba dodać informację co jest  $x$  a co jest  $y$  podczas predykcji w ARIMie.

### 3.2. Metoda xgboost

Metoda xgboost (*Extreme Gradient Boosting*) wywodzi się z metody GBM (*ang. Gradient Boosting Machine*). Obie te metody opierają się o gradientowe wzmacnianie drzew (*ang. gradient boosted trees*). U ich podstaw znajduje się model zespołu drzew decyzyjnych (*ang. decision tree ensembles*). W zespole takim pojedyncze drzewo zazwyczaj słabo radzi sobie z klasyfikacją/regresją, dlatego do predykcji wykorzystywany jest cały zespół drzew. Idea wzmacniania gradientowego polega na tworzeniu ciągu prostych drzew, z których każde kolejne jest zbudowane do predykcji reszt generowanych przez poprzednie.

Różnica w idei działania xgboost a GDM jest niewielka. xgboost używa bardziej regularnej formalizacji modelu w celu kontroli przeuczenia (*over-fitting*). Dodatkowo xgboost jest mocno zoptymalizowany w obszarach szybkości obliczeń oraz zużycia pamięci.

Trzeba dodać informację co jest x a co jest y podczas predykcji w XGBoost.

## 4. Zbiór danych

Liczba szeregów czasowych znajdujących się w benchmark'u NN3 wynosi 111. Długość wektora testowego dla każdego z szeregów wynosi 18. Zbiór danych jest zaczerpnięty z jednorodnej populacji empirycznych biznesowych szeregów czasowych. Wśród danych dominują szeregi, które wydają się być mocno zaszumione (rys. 1a, 1b). Znaleźć tam również można szeregi czasowe z wyraźnym trendem (rys. 1c) lub wahaniami okresowymi (rys. 1d).

Szeregi czasowe w benchmark'u NN3 są stosunkowo krótkie (min 68, max 144), a sam benchmark uznawany jest przez ekspertów z dziedziny modelowania szeregów czasowych jako dość „prosty”.

## 5. Implementacja

### 5.1. ARIMA

Do uzupełnienia.

### 5.2. XGBoost

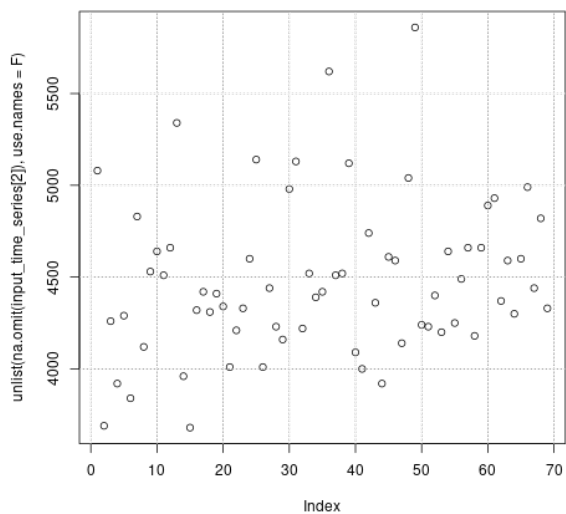
Do uzupełnienia.

## 6. Prezentacja wyników

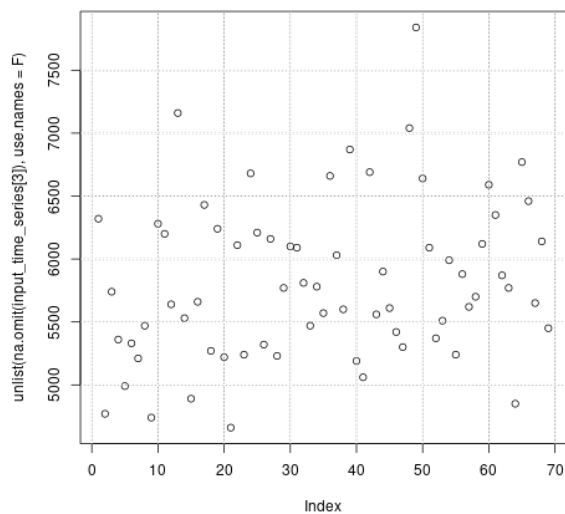
Porównanie modelu ARIMA z metodą xgboost dokonane zostanie na podstawie przebiegu funkcji gęstości prawdopodobieństwa względnych błędów. Wartości błędów nie będą uśredniane dla poszczególnych szeregów, co oznacza, że długość wektora błędów dla którego zostanie przedstawiony przebieg funkcji gęstości prawdopodobieństwa wynosić będzie  $111 \times 18 = 1998$ .

Zagregowana funkcja gęstości prawdopodobieństwa błędów może ukrywać niektóre właściwości stosowanych modeli, dlatego zostanie ukazany również, na wspólnym wykresie, przebieg błędów dla wszystkich szeregów czasowych.

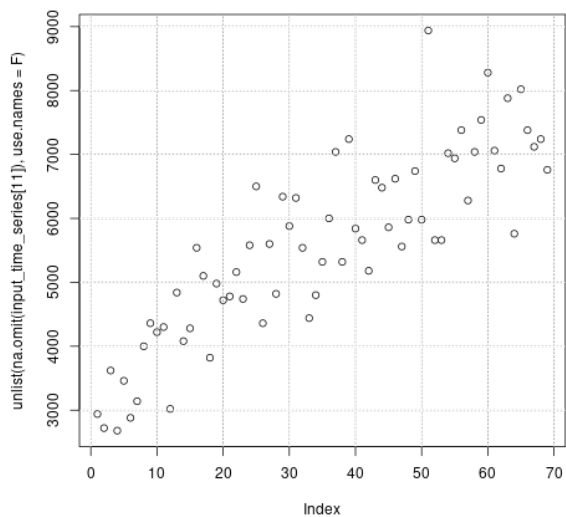
Wyniki strojenia parametrów metody xgboost zaprezentowane zostaną w analogiczny sposób. Dodatkowo zamieszczone mogą zostać fragmenty wydruków stosowanych funkcji, niosące wartościowe informacje.



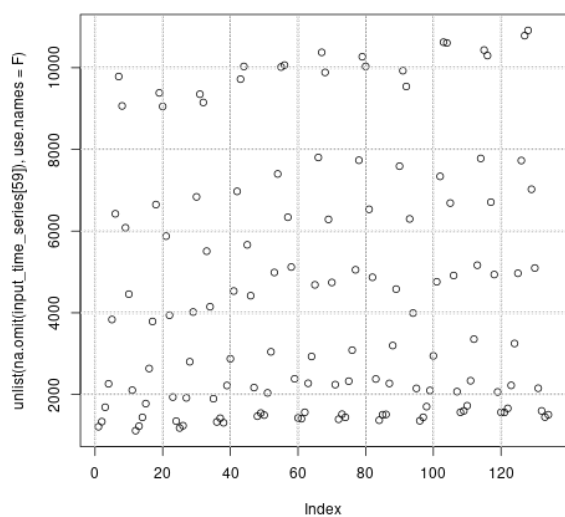
(a) Dane zaszumione



(b) Dane zaszumione



(c) Trend



(d) Wahania okresowe

Rysunek 1: Przykładowe szeregi czasowe z benchmark'u NN3.