

# Automatyczne uczenie maszynowe – projekt 1

Magdalena Bartczak, Aleksandra Mach

## Spis treści

<b>1</b>	<b>Analizowane zbiory danych</b>	<b>2</b>
<b>2</b>	<b>Metody samplingu</b>	<b>2</b>
<b>3</b>	<b>Wyniki optymalizacji</b>	<b>4</b>
3.1	Las losowy . . . . .	4
3.2	XGBoost . . . . .	4
3.3	SVM . . . . .	5
<b>4</b>	<b>Wnioski</b>	<b>5</b>
<b>A</b>	<b>Proces uczenia</b>	<b>6</b>
<b>B</b>	<b>Tuning algorytmów</b>	<b>7</b>

# 1 Analizowane zbiory danych

Do naszego projektu wykorzystaliśmy cztery zbiory danych pobranych z *openml*, które opisują problem klasyfikacji binarnej.

1. **Zbiór 1: *cylinder-bands*** – zbiór rozmiaru  $540 \times 38$ , zawiera 18 zmiennych numerycznych i 19 zmiennych kategoriowych; zawiera braki danych w 263 obserwacjach; zmienna odpowiedzi *band\_type* przyjmuje dwie wartości: 'band' oraz 'noband'.
2. **Zbiór 2: *adult*** – zbiór rozmiaru  $48842 \times 15$ , zawiera 6 zmiennych numerycznych i 8 zmiennych kategoriowych; zmienna odpowiedzi *class* przyjmuje dwie wartości: ' $\leq 50K$ ' oraz ' $> 50K$ '.
3. **Zbiór 3: *labor*** – zbiór rozmiaru  $57 \times 17$ ; zawiera 8 zmiennych numerycznych i 8 zmiennych kategoriowych; zawiera braki danych w 56 obserwacjach; zmienna odpowiedzi *class* przyjmuje dwie wartości: 'good' oraz 'bad'.
4. **Zbiór 4: *sick*** – zbiór rozmiaru  $3772 \times 30$ ; zawiera 6 zmiennych numerycznych i 23 zmienne kategoriowe; zawiera braki danych w 3772 obserwacjach; zmienna odpowiedzi *Class* przyjmuje dwie wartości: 'sick' oraz 'negative'.

Każdy z wymienionych zbiorów został poddany preprocessingowi składającemu się z następujących kroków zawartych w odpowiednim pipeline:

- transformacje kolumn kategoriowych do kolumn numerycznych przy użyciu funkcji *OrdinalEncoder()* oraz *OneHotEncoder()*,
- wypełnianie braków danych przy użyciu funkcji *SimpleImputer()* z argumentem *strategy* = 'median' dla zmiennych numerycznych oraz *strategy* = 'most\_frequent' dla zmiennych kategoriowych,
- skalowanie zmiennych numerycznych przy użyciu funkcji *MinMaxScaler()*.

Pipeline został dopasowany (metoda *fit*) na odpowiednich zbiorach, a następnie zbiory te zostały odpowiednio przekształcone (metoda *transform*). Nie było podziału na zbiory treningowy i testowy ze względu na to, że w obu metodach samplingu krosvalidacja w nich zawarta sama dokonuje tych podziałów.

## 2 Metody samplingu

Zgodnie z treścią polecenia próbowaliśmy zdefiniować siatki hiperparametrów dla każdego z trzech algorytmów na podstawie artykułu. Ostatecznie, ze względu na ograniczenia obliczeniowe Google colab, w którym pracowaliśmy, postanowiliśmy ograniczyć te siatki do następujących:

- Las losowy:
  - *n\_estimators*: np.arange(50, 501)
  - *max\_depth*: np.arange(10,21)
  - *min\_samples\_leaf*: np.arange(1,11)
- XGBoost:

- *learning\_rate*: `uniform( $2^{-10}$ , 1)`
  - *max\_depth*: `np.arange(5,13)`
  - *num\_boost\_round*: `np.arange(50,501)`
  - *min\_child\_weight*: `uniform(1,  $2^7$ )`
- SVM:
    - *C*: `np.arange(1,12)`
    - *gamma*: `np.arange(1,12)`

Do tuningu hiperparametrów wybranych algorytmów zastosowaliśmy dwie metody losowania punktów:

1. ***Random Search*** – z wykorzystaniem funkcji *RandomizedSearchCV* z pakietu *scikit-learn*, przeszukuje losowo zdefiniowaną przestrzeń hiperparametrów modelu, oceniając skuteczność każdej kombinacji za pomocą krosvalidacji, aby znaleźć optymalne parametry, przyspieszając proces przeszukiwania w porównaniu do pełnego przeglądu siatki. Algorytm losuje określoną liczbę kombinacji parametrów, ocenia ich wydajność i zwraca najlepsze parametry dla ostatecznego dopasowania modelu.
2. ***Bayes Optimization*** – z wykorzystaniem pakietu *SMAC*, wykorzystuje model statystyczny, zwany procesem Gaussa, do adaptacyjnego przeszukiwania przestrzeni hiperparametrów, selekcji nowych punktów na podstawie aktualnych ocen, minimalizując liczbę prób potrzebnych do znalezienia optymalnych parametrów modelu. To iteracyjny proces, gdzie każda iteracja dostosowuje model na podstawie wyników poprzednich prób, skoncentrowanych na obszarach przestrzeni parametrów, które mają największe szanse zawierania optymalnych ustawień.

Napisałyśmy dwie funkcje *random\_search* oraz *bayesian\_opt* odpowiedzialne za przeszukiwanie siatki hiperparametrów oraz wybranie najlepszego zestawu przy użyciu odpowiedniej metody samplingu.

1. ***random\_search*** – przyjmuje argumenty *model*, gdzie *model*  $\in$  [*'rf\_classifier'*, *'xgb\_classifier'*, *'SVM\_classifier'*], *X\_train*, *y\_train* oraz *liczba\_iteracji*; kolejne kroki w funkcji:
  - (i) w zależności od wybranego modelu jest definiowana odpowiednia siatka hiperparametrów, z której będą losowane kombinacje,
  - (ii) użycie funkcji *RandomizedSearchCV* z pakietu *scikit-learn* do przeprowadzenia 3-krotnej krosvalidacji ze scoringiem *roc\_auc*,
  - (iii) dopasowanie danych *X\_train* i *y\_train*,
  - (iv) funkcja zwraca wyniki przeprowadzonej krosvalidacji dla każdej iteracji algorytmu oraz zestaw wybranych optymalnych hiperparametrów.
2. ***bayesian\_opt*** – przyjmuje argumenty *model*, gdzie *model*  $\in$  [*'rf\_classifier'*, *'xgb\_classifier'*, *'SVM\_classifier'*], *X*, *y*, *liczba\_prob*, *seed*; kolejne kroki w funkcji:
  - (i) w zależności od wybranego modelu jest definiowana odpowiednia siatka hiperparametrów, z której będą losowane kombinacje,

- (ii) definiowana jest funkcja celu *train*, w której przeprowadzana jest 3-krotna krosvalidacja przy użyciu wbudowanej funkcji *cross\_val\_score* przy użyciu scoringu *roc\_auc\_score*, a zwracany jest średni wynik dla każdej wylosowanej kombinacji parametrów,
- (iii) funkcja celu jest optymalizowana i zapisywane są wyniki z każdej wylosowanej kombinacji parametrów; pakiet *smac* pozwala nam na dodanie maksymalnej liczby prób (tj. losowanych kombinacji parametrów), ale może zakończyć działanie wcześniej, gdy wyniki zaczną zbiegać,
- (iv) funkcja zwraca wszystkie zapisywane wyniki oraz zestaw parametrów, który maksymalizuje AUC.

### 3 Wyniki optymalizacji

Każdy z wymienionych algorytmów został zastosowany do każdego z czterech wspomnianych zbiorów danych. Do każdego z nich metody samplingu dobierały optymalne parametry przeszukując tą samą siatkę zdefiniowaną wcześniej. W tym rozdziale prezentujemy najlepsze parametry wybrane dla każdego z algorytmów zarówno dla każdego z 4 zbiorów indywidualnie jak i dla wszystkich zbiorów razem (defaultowe) znalezione przy pomocy dwóch wspomnianych metod samplingu.

#### 3.1 Las losowy

Tabela 1: Optymalne parametry wybrane przez las losowy

	RandomSearch					Bayes				
	db 1	db 2	db 3	db 4	all	db 1	db 2	db 3	db 4	all
min_samples_leaf	6	1	2	1	10	4	2	3	1	2
max_depth	17	18	15	18	20	16	19	23	19	19
n_estimators	84	394	241	394	237	328	155	299	260	480

#### 3.2 XGBoost

Tabela 2: Optymalne parametry wybrane przez XGBoost

	RandomSearch					Bayes				
	db 1	db 2	db 3	db 4	all	db 1	db 2	db 3	db 4	all
learning_rate	0.002	0.001	0.01	0.001	0.002	0.11	0.12	0.70	0.27	0.7
max_depth	6	10	9	5	6	12	6	7	11	10
num_boost_round	447	257	121	489	447	369	65	126	94	106
min_child_weight	9.89	21.2	37.87	36.74	9.89	19.55	48.65	3.58	34.25	24.81

### 3.3 SVM

Tabela 3: Optymalne parametry wybrane przez SVM

	RandomSearch					Bayes				
	db 1	db 2	db 3	db 4	all	db 1	db 2	db 3	db 4	all
C	9	6	5	5	7	11	14	18	17	15
gamma	2	4	1	1	3	1	2	2	3	5

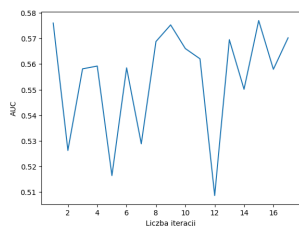
## 4 Wnioski

Podsumowując przedstawione wyniki można stwierdzić, że proces doboru optymalnych parametrów jest bardzo istotnym krokiem w procedurach uczenia maszynowego. Chcąc przeszukać odpowiednią dużą siatkę hiperparametrów, aby udało się nam znaleźć jak najlepszą kombinację, musimy liczyć się ze sporym czasem obliczeń. Podobnie, ograniczeniem jest też liczba iteracji, czyli liczba losowanych kombinacji parametrów. Zwiększając ją, zwiększamy szansę na wylosowanie jak najlepszej kombinacji. Mimo to zauważyliśmy, że dla optymalizacji bayesowskiej stabilne wyniki otrzymujemy już po kilkunastu iteracjach dla każdego zbioru. Generalnie optymalne parametry wybrane dla wszystkich zbiorów wspólnie różnią się od tych wybieranych dla każdego zbioru indywidualnie, a wyniki dla parametrów defaultowych są zazwyczaj nieco mniejsze niż najlepsze wyniki dla każdego ze zbiorów z osobna. Jednak różnice te są na tyle małe, że uważamy, że dobór parametrów defaultowych może być dobrym i uniwersalnym rozwiązaniem.

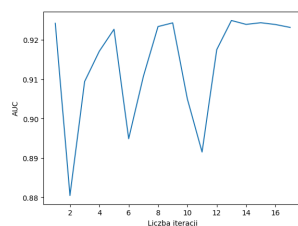
Dodatkowo przedstawiamy wybrane wizualizacje procesów uczenia oraz tunowalności algorytmów.

## A Proces uczenia

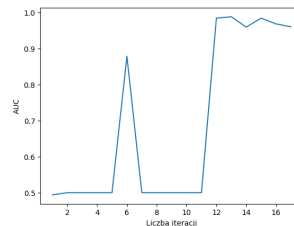
### 1. XGBoost – optymalizacja bayesowska



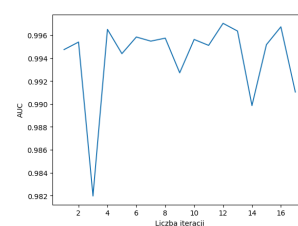
(a) Zbiór 1



(b) Zbiór 2

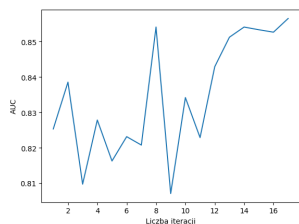


(c) Zbiór 3

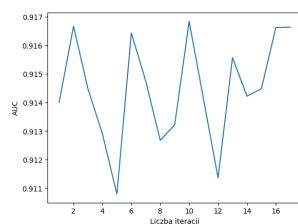


(d) Zbiór 4

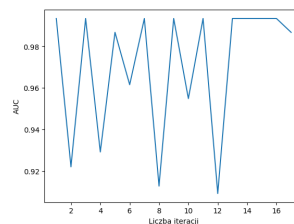
### 2. Random Forest – optymalizacja bayesowska



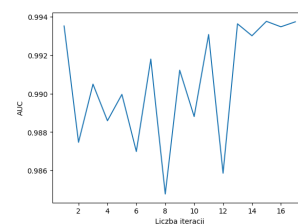
(a) Zbiór 1



(b) Zbiór 2

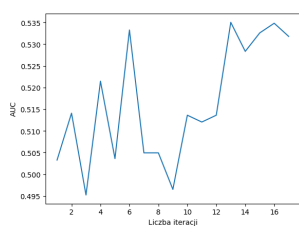


(c) Zbiór 3

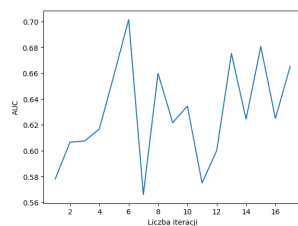


(d) Zbiór 4

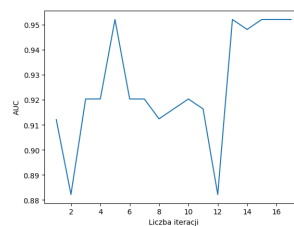
### 3. SVM – optymalizacja bayesowska



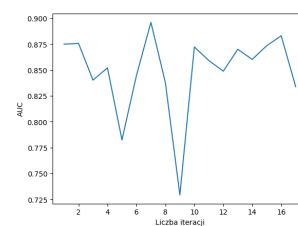
(a) Zbiór 1



(b) Zbiór 2

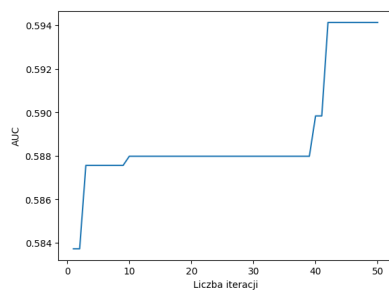


(c) Zbiór 3

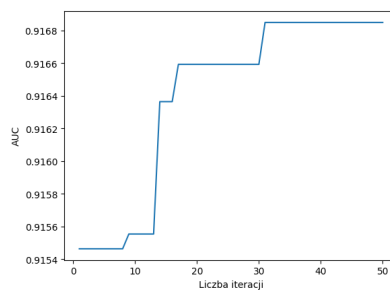


(d) Zbiór 4

## 4. Random Forest – RandomSearch



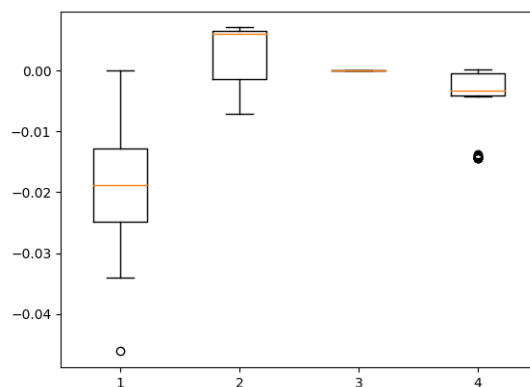
(a) Zbiór 1



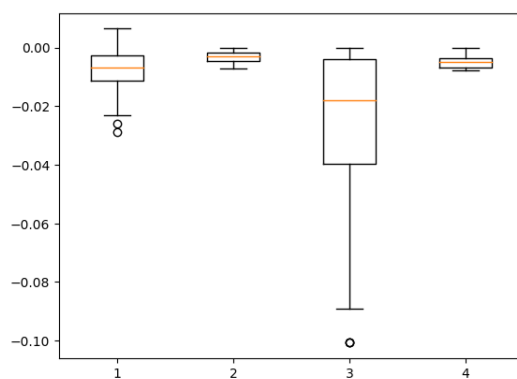
(b) Zbiór 2

## B Tuning algorytmów

## 1. Tuning dla 4 zbiorów – XGBoost, RandomSearch



## 2. Tuning dla 4 zbiorów – Random Forest, RandomSearch



## 3. Tuning dla 4 zbiorów – SVM, RandomSearch

