

AutoML - Praca Domowa 1

Adrianna Wieczorek i Joanna Witczak

November 21, 2023

Contents

1	Wstęp	2
2	Opis Eksperymentu	2
2.1	Losowanie z siatki za pomocą losowania z rozkładu jednostajnego	2
2.2	Optymalizacja Bayesowska	3
2.3	Testy	4
3	Wyniki	4
3.1	Stabilizacja RandomForest	4
3.2	Stabilizacja GradientBoosting	5
3.3	Stabilizacja SVM	5
3.4	Tunowalność	6
4	Wnioski	6

1 Wstęp

W ramach projektu zbadaliśmy tunowalność algorytmów wykorzystywanych w problemie klasyfikacji binarnej. Eksperyment przeprowadziliśmy dla poniższych metod:

- RandomForest
- SVM
- GradientBoosting.

W ramach eksperymentu porównaliśmy dwa sposoby wyboru najlepszych hiperparametrów z siatki. Pierwsza metoda opiera się o losowanie punktów z rozkładu jednostajnego – RandomizedSearchCV, druga metoda bazuje na optymalizacji bayesowskiej. Do oceny jakości modeli wykorzystaliśmy dokładność.

2 Opis Eksperymentu

2.1 Losowanie z siatki za pomocą losowania z rozkładu jednostajnego

W przeprowadzonym przez nas eksperymencie, badaliśmy jakość dopasowania algorytmu na 4 zbiorach danych: “credit-g”, “monks-problems-1”, “monks-problems-2”, “pc1”. Wykorzystane przez nas zbiory dostępne są w openML.

W celu poprawienia jakości danych, zbudowaliśmy pipeline. W ramach niego przeprowadzone są następujące czynności:

- Imputacja braków danych, w przypadku zmiennych ciągłych uzupełniane są one przez wartość średnią, a dla zmiennych kategorycznych jest to najczęściej występująca wartość.
- Skalowanie przeprowadzone dla zmiennych ciągłych. Jest ono szczególnie istotne w przypadku metody SVM.
- One-Hot-Encoding.

W kolejnym kroku przeprowadziliśmy optymalizację hiperparametrów dla każdego algorytmu za pomocą RandomizedSearchCV. Jest to metoda opierająca się na losowaniu punktów z rozkładu jednostajnego. Dzięki niej nie musimy badać wszystkich możliwych kombinacji hiperparametrów. W ramach naszego eksperymentu rozważyliśmy przeszukiwanie siatki w 50 krokach. Okazało się to być wystarczające, gdyż stabilizacja wyników następowała stosunkowo szybko - około 20 iteracji.

Każda z wybranych siatek (dla danego algorytmu) wykorzystaliśmy do dopasowania modelu dla każdego zbioru i wyznaczenia średnio najlepszej siatki hiperparametrów. Porównując wyniki dokładności uzyskane za jej pomocą dla każdego zbioru danych z “najlepszą” siatką wybraną dla tego zbioru, otrzymaliśmy różnice rzędu drugiego miejsca po przecinku. Zdarzyły się przypadki gdy średnio najlepsza siatka zwróciła wynik lepszy, niż ta dobrana w sposób spersonalizowany.

Dla lasu losowego rozważyliśmy następujące hiperparametry dla metody Random Search:

- n_estimators, czyli liczba drzew decyzyjnych w lesie losowym,
- min_samples_split, czyli parametr informujący o minimalnej wymaganej liczbie obserwacji w danym węźle drzewa decyzyjnego, aby nastąpił podział węzła,
- warm_start, czyli parametr odpowiadający za to, czy w kolejnej iteracji używamy rozwiązania z poprzedniego dopasowania i do tego dokładamy kolejne drzewa (wartość True), czy dopasowujemy zupełnie nowy las (wartość False),
- max_depth, czyli maksymalna głębokość drzewa (najdłuższa ścieżka między węzłem korzenia a węzłem liścia),
- min_samples_leaf, czyli minimalna liczba próbek, które powinny być obecne w węźle liścia.

Table 1: Siatka hiperparametrów dla Random Forest

Parametr	Wartości
model__n_estimators	[20, 50, 100, 200, 500, 700, 1000, 1500, 2000]
model__min_samples_split	[1e-5, 1e-4, 0.0001, 0.001, 0.01, 0.1]
model__warm_start	[True, False]
model__max_depth	[10, 30, 50, 70, 100, 150, 200, 500, 1000]
model__min_samples_leaf	[2, 3, 5, 10, 15, 20]

Dla Gradient Boosting rozważaliśmy następujące hiperparametry dla metody Random Search:

- n_estimators, czyli liczba drzew decyzyjnych w modelu,
- subsample, czyli frakcja obserwacji, która ma zostać wykorzystana do dopasowania poszczególnych uczących się modeli,
- loss, czyli funkcja straty do optymalizacji,
- max_depth, czyli maksymalna głębokość drzewa (najdłuższa ścieżka między węzłem korzenia a węzłem liścia),
- min_samples_leaf, czyli minimalna liczba próbek, które powinny być obecne w węźle liścia.

Table 2: Siatka hiperparametrów dla Gradient Boosting

Parametr	Wartości
model__n_estimators	[10, 50, 100, 150, 500, 1000, 1500, 2000, 2500]
model__subsample	[1e-5, 0.005, 0.01, 0.05, 0.1, 0.5, 0.7, 0.9]
model__max_depth	[15, 50, 100, 150, 500, 1000, 1500, 2000, 2500]
model__min_samples_leaf	[1, 3, 5, 10, 15, 20]
model__loss	['log_loss', 'exponential']

Dla SVM rozważaliśmy następujące hiperparametry dla metody Random Search:

- C, czyli hiperparametr regularyzacji, gdzie siła regularyzacji jest odwrotnie proporcjonalna do C,
- kernel, czyli hiperparametr określający typ jądra, jakie ma zostać użyte w algorytmie,
- gamma, czyli współczynnik jądra dla jader typu rbf, poly i sigmoid,
- tol, czyli tolerancja dla kryterium zatrzymania.

Table 3: Siatka hiperparametrów dla SVM

Parametr	Wartości
model__C	[0.001, 0.005, 0.01, 1, 10, 20, 50, 100, 200, 500, 1000]
model__kernel	['rbf', 'sigmoid', 'poly']
model__gamma	[0.01, 0.1, 0.5, 1, 5, 20, 50, 100]
model__tol	[1e-6, 1e-5, 1e-4, 1e-3, 1e-2, 0.1]

2.2 Optymalizacja Bayesowska

Tak jak w przypadku zrandomizowanego wyboru hiperparametrów pracowaliśmy na danych oczyszczonych i wystandaryzowanych. Dzięki sposobowi działania tego algorytmu mogliśmy przeszukiwać przestrzeń ciągłą, a nie dyskretną. Optymalizacja bayesowska to technika, która wykorzystuje modele probabilistyczne do tunowania hiperparametrów modeli. Na początku tworzony jest model zastępczy, który przybliża funkcję celu. Następnie na podstawie tego przybliżenia wybierane są punkty w przestrzeni hiperparametrów do ewaluacji. Po każdej ewaluacji model jest aktualizowany, a proces iteracyjnie powtarza się, dążąc do znalezienia optymalnej konfiguracji hiperparametrów minimalizującej funkcję celu.

Rozważaliśmy te same hiperparametry oraz te same modele co w poprzednim kroku.

Jako punkt odniesienia przyjęliśmy średnio najlepszą siatkę parametrów otrzymana metodą Randomized-SearchCV.

2.3 Testy

Przeprowadziliśmy test Kołmogorowa Smirnowa do porównania rozkładów dokładności otrzymanych w kolejnych iteracjach (maksimum do danej iteracji) dla optymalizacji bayesowskiej i RandomizedSearchCV. W ich wyniku otrzymaliśmy, że na podstawie dostępnych danych możemy stwierdzić, że rozkłady dokładności na poziomie istotności 0.05 są istotnie różne.

W związku z tym zdecydowaliśmy się sprawdzić, czy wyniki otrzymywane dla metody optymalizacji bayesowskiej są istotnie lepsze od wyników otrzymywanych w sposób losowy.

Przeprowadziliśmy test Wilcoxona, w wyniku którego otrzymaliśmy, że w 9 na 12 przypadków przyjmujemy hipotezę alternatywną (na poziomie istotności 0.05), mówiącą o tym, że dokładność otrzymywana metodą optymalizacji bayesowskiej jest lepsza niż, ta za pomocą losowania z rozkładu jednostajnego.

Warto podkreślić, że testy wykonywaliśmy oddzielnie dla każdej metody i każdego zbioru danych.

Model	Dane	p-wartość (K-S)	p-wartość (Wilcoxon)
RandomForest	dane 1	0	0
RandomForest	dane 2	0.00004	0.0005
RandomForest	dane 3	0	0
RandomForest	dane 4	0	1.0
SVM	dane 1	0	0
SVM	dane 2	0	0.00005
SVM	dane 3	0	0.00382
SVM	dane 4	0	1
GradientBoosting	dane 1	0	0
GradientBoosting	dane 2	0.00192	0.0008
GradientBoosting	dane 3	0	0
GradientBoosting	dane 4	0	1

Table 4: Wyniki testów Wilcoxona i K-S

3 Wyniki

3.1 Stabilizacja RandomForest

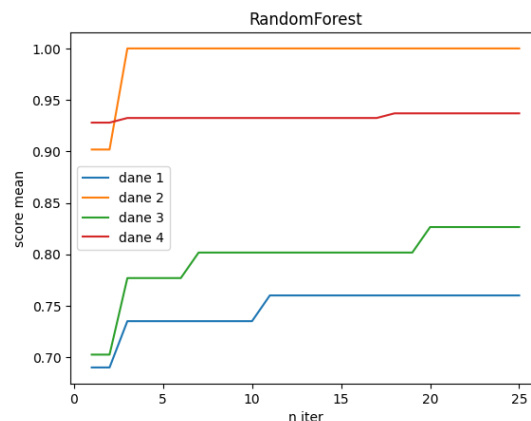
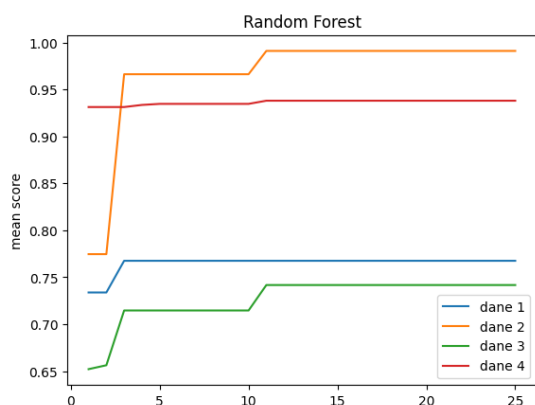


Figure 1: Stabilizacja RandomForest dla Random search.

Figure 2: Stabilizacja RandomForest dla optymalizacji bayesowskiej.

3.2 Stabilizacja GradientBoosting

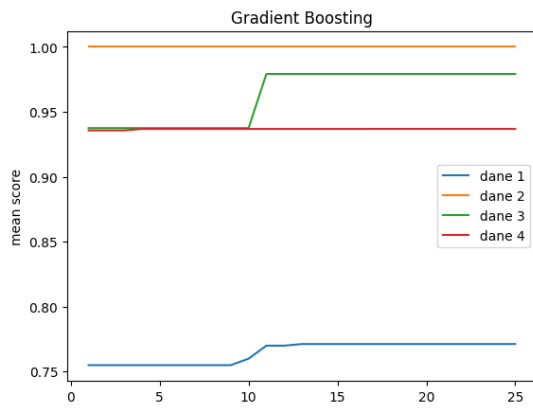


Figure 3: Stabilizacja GradientBoosting dla Random search.

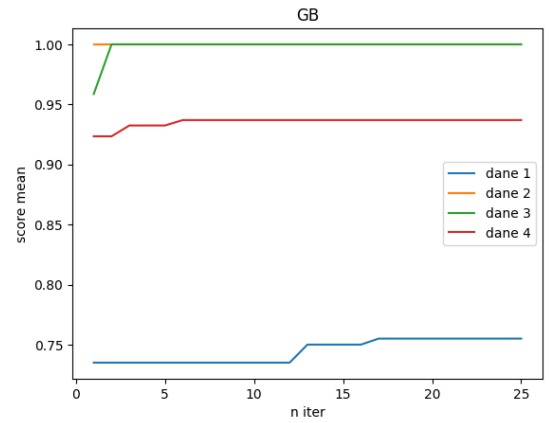


Figure 4: Stabilizacja GradientBoosting dla optymalizacji bayesowskiej.

3.3 Stabilizacja SVM

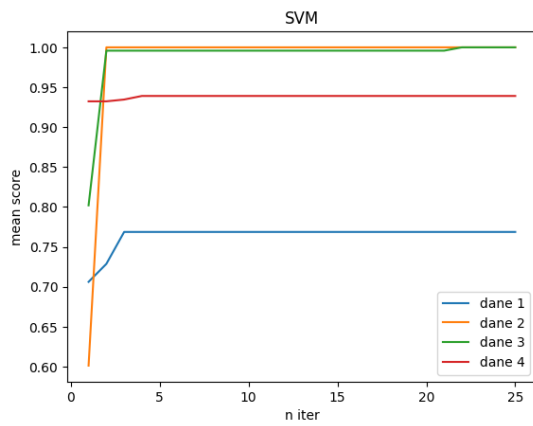


Figure 5: Stabilizacja SVM dla Random search.

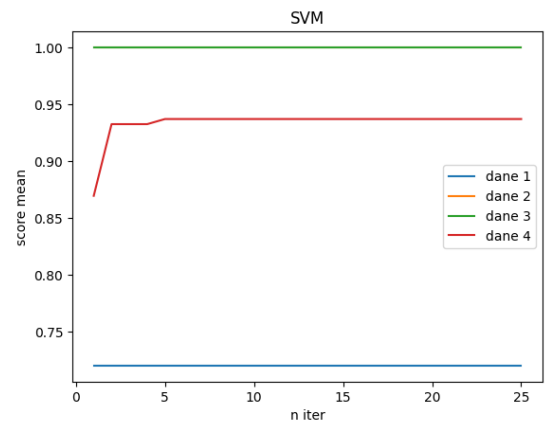


Figure 6: Stabilizacja SVM dla optymalizacji bayesowskiej.

3.4 Tunowalność

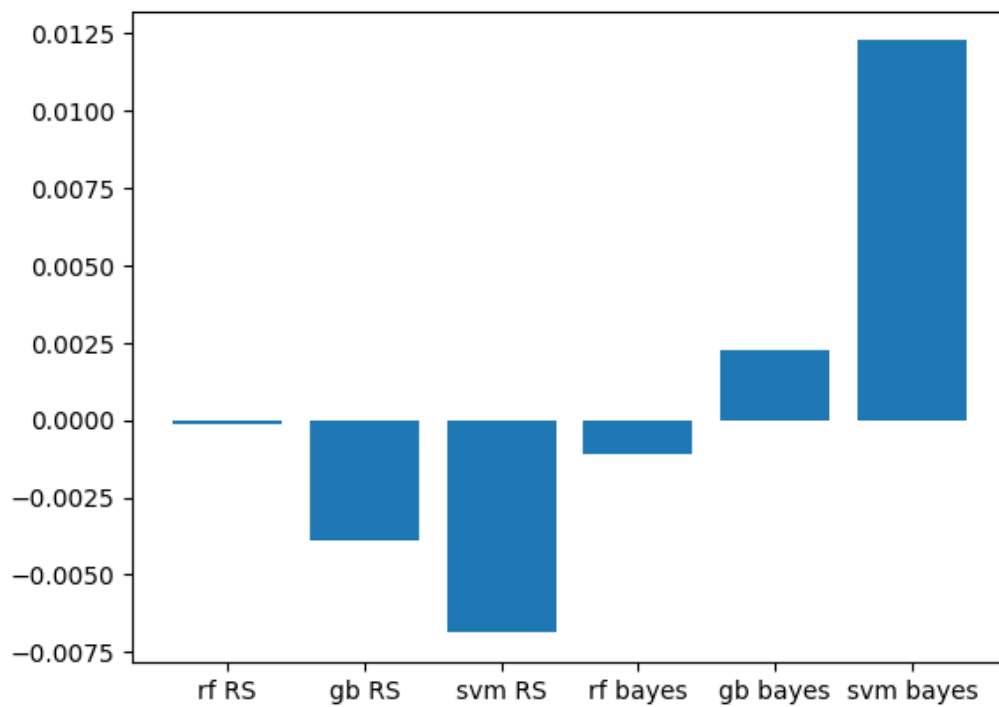


Figure 7: Tunowalność algorytmów.

4 Wnioski

Za pomocą optymalizacji bayesowskiej jesteśmy w stanie otrzymać nieznacznie szybszą zbieżność do optymalnego rozwiązania oraz jesteśmy w stanie otrzymać średnio lepsze wyniki w tej samej liczbie kroków niż za pomocą losowania siatek z rozkładu jednostajnego.

Badając tunowalność algorytmów zaobserwowaliśmy, że w 3 na 6 przypadków tunowalność przyjęła wartość ujemną. Sugeruje to, że dokładność otrzymana za pomocą średnio najlepszej siatki była znacząco gorsza od wyników otrzymanych za pomocą najlepszej siatki.