

Analiza tunowalności hiperparametrów

Martyna Sadowska, Zuzanna Zalewska, Aleksandra Zawadka

1 Wstęp

Celem projektu jest przeanalizowanie i porównanie różnych metod strojenia hiperparametrów oraz wyznaczenie nowej domyślnej konfiguracji dla każdego z rozważanych algorytmów. Wykorzystane zostały trzy modele uczenia maszynowego: Decision Tree, Random Forest oraz XGBoost. Każdy z nich został zoptymalizowany przy użyciu dwóch technik doboru punktów: losowego przeszukiwania (Random Search) oraz optymalizacji bayesowskiej (Bayesian Optimization), w celu zbadania możliwości dostrajania poszczególnych algorytmów. Do pomiaru jakości działania modeli wykorzystano metrykę ROC-AUC, umożliwiającą porównanie wydajności klasyfikatorów niezależnie od progu decyzyjnego. Dzięki temu zmiany w ROC-AUC są porównywalne między różnymi zbiorami danych i nie wymagają dodatkowego skalowania, co czyni ją odpowiednią miarą do oceny jakości klasyfikacji w tym projekcie. Podstawą naszych badań jest artykuł [1].

2 Zbiory danych

Do przeprowadzenia badań wykorzystano następujące zbiory danych klasyfikacyjnych:

1. Fitness Classification Dataset
2. Heart Failure Prediction Dataset
3. Placement Prediction Dataset
4. Travel Insurance Prediction Data
5. Customer Churn Dataset

Każdy ze zbiorów został przetworzony w jednolity sposób, czyli usunięto brakujące wartości, przeprowadzono kodowanie zmiennych kategorycznych oraz normalizację cech numerycznych.

3 Definicja i estymacja tunowalności hiperparametrów

W analizie tunowalności hiperparametrów rozważamy każdy algorytm uczenia maszynowego jako funkcję $\hat{f}(X, \theta)$, gdzie θ jest wektorem hiperparametrów należących do przestrzeni poszukiwań Θ . Aby ocenić jakość wybranego zestawu hiperparametrów, definiujemy funkcję straty $L(Y, \hat{f}(X, \theta))$ oraz odpowiadające jej oczekiwane ryzyko:

$$R(\theta) = E[L(Y, \hat{f}(X, \theta)) | P],$$

gdzie P jest nieznanym rozkładem danych.

Dla m różnych zbiorów danych lub rozkładów P_1, \dots, P_m otrzymujemy m funkcji ryzyka

$$R^{(j)}(\theta) = E[L(Y, \hat{f}(X, \theta)) | P_j], \quad j = 1, \dots, m.$$

Najlepsza konfiguracja hiperparametrów dla zbioru danych j definiowana jest jako

$$\theta^{(j)*} := \arg \min_{\theta \in \Theta} R^{(j)}(\theta).$$

Wartości domyślne, stosowane w oprogramowaniu, powinny dobrze sprawdzać się na różnych zbiorach danych. Możemy zdefiniować optymalne wartości domyślne na podstawie eksperymentów na wielu zbiorach danych jako

$$\theta^* := \arg \min_{\theta \in \Theta} g(R^{(1)}(\theta), \dots, R^{(m)}(\theta)),$$

gdzie g jest funkcją agregującą (w tym przypadku średnią).

Tunowalność algorytmu na zbiorze danych j definiujemy jako różnicę między ryzykiem przy wartości domyślnej a ryzykiem przy najlepszej możliwej konfiguracji:

$$d^{(j)} := R^{(j)}(\theta^*) - R^{(j)}(\theta^{(j)*}).$$

4 Siatka hiperparametrów

Dla każdego algorytmu wybrano zestaw kluczowych hiperparametrów, które mają istotny wpływ na wydajność modelu. Wartości hiperparametrów były losowane z określonych przedziałów (zostało to przedstawione w tabeli 1), co umożliwiała badanie tunowalności modeli. Zakresy hiperparametrów używane do strojenia powinny obejmować wartości bliskie optymalnym konfiguracjom z wysokim prawdopodobieństwem.

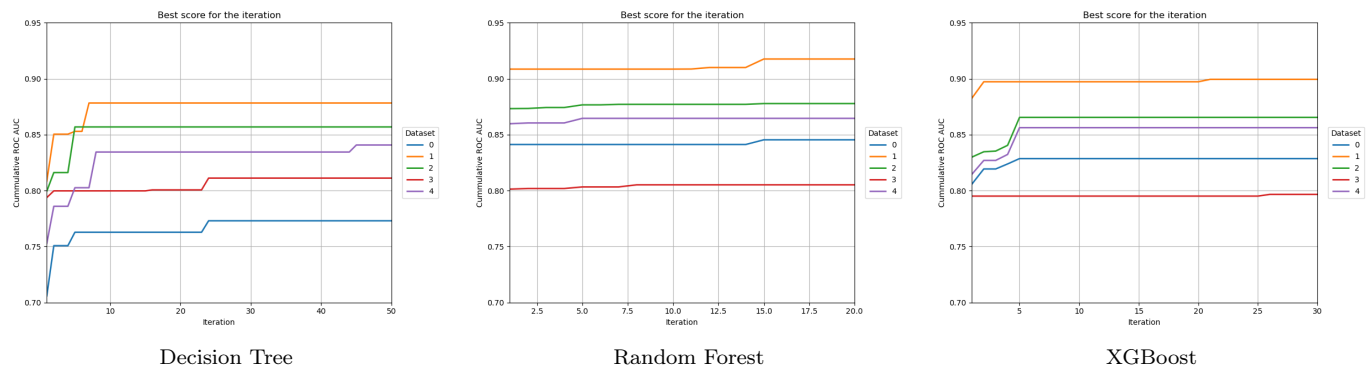
Algorytm	Hiperparametr	Typ	Dolna granica	Górna granica
Decision Tree	<i>max_depth</i>	całkowity	1	30
	<i>min_samples_split</i>	całkowity	2	20
	<i>min_samples_leaf</i>	całkowity	1	10
	<i>max_features</i>	liczbowy	0.1	0.9
Random Forest	<i>n_estimators</i>	całkowity	100	2000
	<i>max_depth</i>	całkowity	5	30
	<i>min_samples_leaf</i>	całkowity	1	10
	<i>max_features</i>	liczbowy	0.1	0.9
XGBoost	<i>n_estimators</i>	całkowity	100	5000
	<i>learning_rate</i>	liczbowy	0.01	0.3
	<i>max_depth</i>	całkowity	1	15
	<i>subsample</i>	liczbowy	0.1	0.9

Tabela 1: Hiperparametry algorytmów.

5 Random Search

Pierwszą metodą losowania punktów, która została użyta był Random Search. Dla każdego modelu zastosowaliśmy pięciokrotną walidację krzyżową i w zależności od wydajności modelu od 20 do 50 iteracji.

Na wykresie 1 widać, że dla Decision Tree istotne było wykonać przynajmniej 10 iteracji, natomiast później poprawa wyniku jest już nieznaczna. Dla Random Forest nowe iteracje nie wносиły dużej poprawy wartości ROC-AUC w żadnej z kolejnych iteracji. Dla XGBoost pierwsze 5 iteracji znacząco poprawiło wynik, natomiast później nie ma już widocznej zmiany.



Rysunek 1: Najlepsze wartości ROC-AUC otrzymane do danej iteracji dla trzech modeli po wyszukiwaniu losowym.

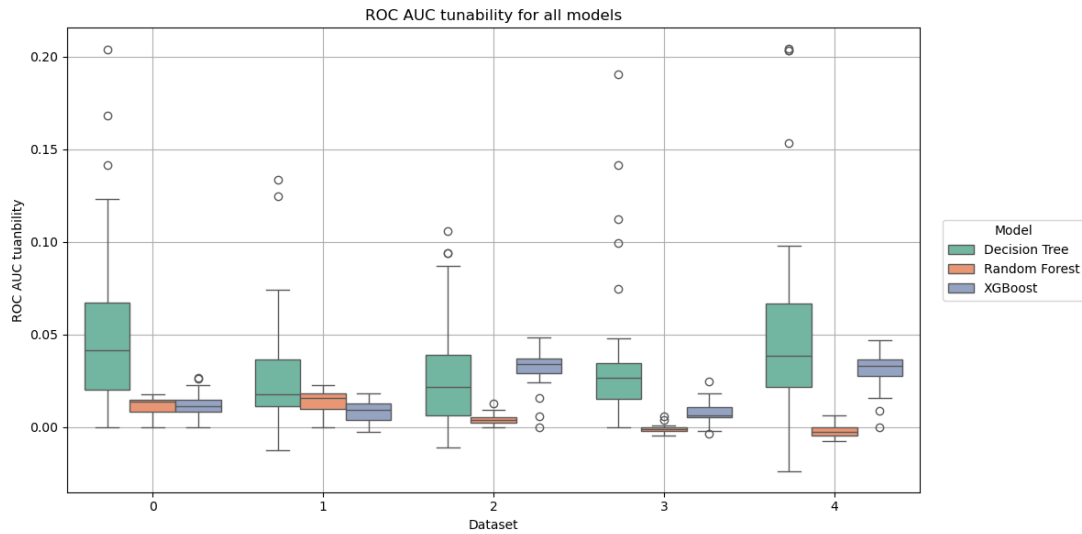
Po wykonaniu wyszukiwania losowego dla każdego z modeli uzyskaliśmy hiperparametry, które dawały średnio najlepsze wyniki na wszystkich zbiorach danych. Modele, które zostały naszymi nowymi defaultami są umieszczone w tabeli 2.

Decision Tree		Random Forest		XGBoost	
<i>max_depth</i> :	9.0	<i>max_depth</i> :	12.0	<i>learning_rate</i> :	0.1421
<i>max_features</i> :	0.5565	<i>max_features</i> :	0.1980	<i>max_depth</i> :	5.0
<i>min_samples_leaf</i> :	9.0	<i>min_samples_leaf</i> :	6.0	<i>n_estimators</i> :	139.0
<i>min_samples_split</i> :	15.0	<i>n_estimators</i> :	795.0	<i>subsample</i> :	0.7411
Mean Test Score:	0.8226	Mean Test Score:	0.8598	Mean Test Score:	0.8481

Rysunek 2: Najlepsze modele otrzymane po wyszukiwaniu losowym

Tutaj tunowalność jest liczona jako różnica między wynikiem optymalnym a uzyskanym w danej iteracji. Przedstawiliśmy ją na wykresie 3 dla każdego zbioru danych i dla każdego modelu. Widać na nim, że wyniki dla Decision Tree

znacząco różniły się na przestrzeni iteracji i wybrany przez nas model ma lepsze wyniki niż modele o innych wylosowanych hiperparametrach. Tunowalność Random Foresta jest najmniejsza z naszych modeli. Wyniki w poszczególnych iteracjach tylko nieznacznie się od siebie różnią. XGBoost także nie daje dużych wyników tunowalności, ale jednak wyniki najlepszego modelu są trochę lepsze od reszty.

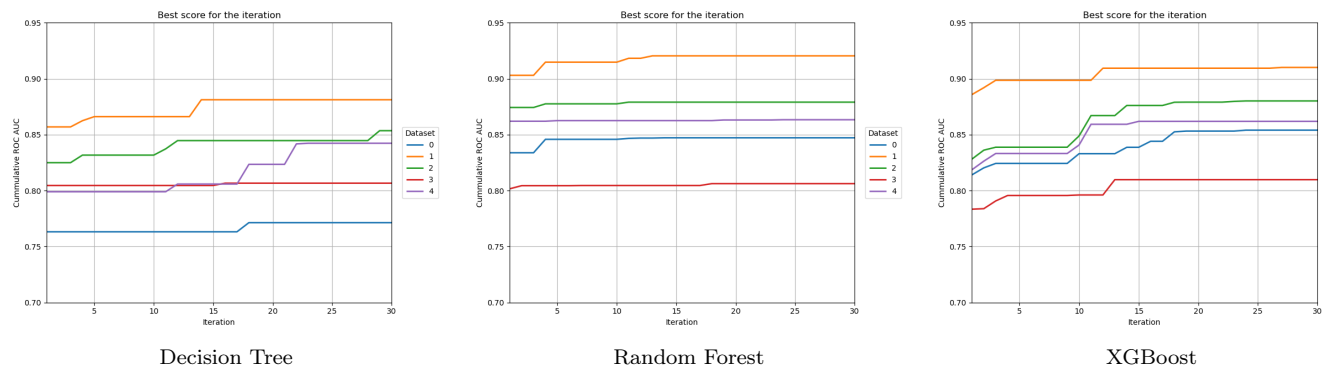


Rysunek 3: Tunowalność wszystkich testowanych modeli.

6 Optymalizacja Bayesowska

Druga metodą losowania punktów opierała się na technice bayesowskiej. Dla każdego modelu zastosowaliśmy pięciokrotną walidację krzyżową i 30 iteracji.

Na wykresie 4 widać, że dla Decision Tree najwięcej wносиły iteracje 10-25, ale później też były iteracje poprawiające wynik. Dla Random Forest pierwsze kilka iteracji podnosiły wartość AUC-ROC, ale kolejne już nic nie zmieniły. XGBoost potrzebował ok. 20 iteracji, aby osiągnąć optymalny wynik.



Rysunek 4: Najlepsze wartości ROC-AUC otrzymane do danej iteracji dla trzech modeli po optymalizacji bayesowskiej.

Po wykonaniu optymalizacji bayesowskiej dla każdego z modeli i dla każdego zbioru danych uzyskaliśmy hiperparametry, które dawały najlepsze wyniki. Modele, które zostały naszymi nowymi domyślnymi konfiguracjami są umieszczone w tabeli 6 w załączniku.

Analizując tunowalność modeli po optymalizacji bayesowskiej porównaliśmy wyniki ROC-AUC dla najlepszych znalezionych przez nas modeli z wynikami otrzymanymi w domyślnych modelach systemowych. W tabeli 2 umieszczone są otrzymane przez nas wyniki. Widać, że dla wszystkich zbiorów danych wyniki Decision Tree zostały poprawione. Dla XGBoosta niektóre zbiory uzyskały lepsze wyniki niż modele domyślne, ale niektóre poprawiły się nieznacznie. Namniejszą poprawę widać w modelu Random Forest. Poprawa wyniku dla każdego zbioru jest niewielka.

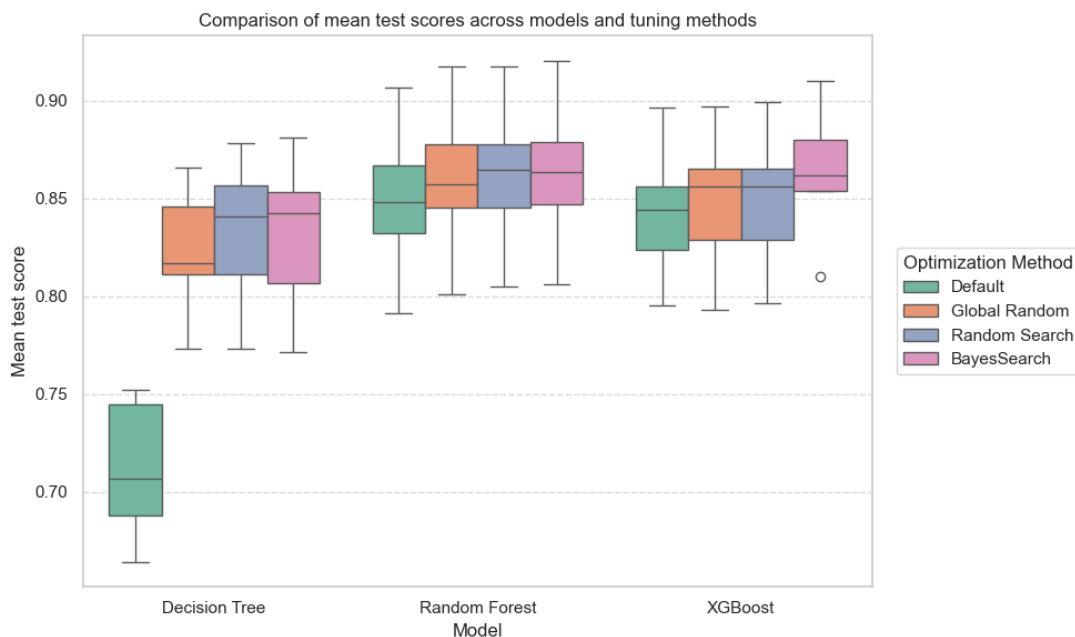
7 Porównanie metod samplingu

Na wykresie 5 widać, że największa poprawa została osiągnięta dla Decision Tree. BayesSearch okazał się trochę skuteczniejszy od RandomSearch dla modeli Random Forest i XGBoost. W ogólności oba sposoby samplingu popra-

Zbiór danych	Model	Wynik ROC-AUC po optymalizacji bayesowskiej	Poprawa wyniku domyślnego
0	Decision Tree	0.7714	0.1073
1	Decision Tree	0.8812	0.1290
2	Decision Tree	0.8536	0.1468
3	Decision Tree	0.8067	0.0618
4	Decision Tree	0.8423	0.1543
0	Random Forest	0.8472	0.0152
1	Random Forest	0.9204	0.0140
2	Random Forest	0.8791	0.0121
3	Random Forest	0.8063	0.0148
4	Random Forest	0.8634	0.0150
0	XGBoost	0.8540	0.0304
1	XGBoost	0.9101	0.0137
2	XGBoost	0.8801	0.0243
3	XGBoost	0.8098	0.0148
4	XGBoost	0.8619	0.0180

Tabela 2: Porównanie wyników modeli po optymalizacji bayesowskiej i poprawy względem wartości domyślnych.

wiają wyniki w porównaniu do Default. Dokładne wyniki różnic w tych metodach zostały przedstawione w tabeli 3 w załączniku.



Rysunek 5: Porównanie metod samplingu: Default - model z domyślnymi hiperparametrami; Global Random - najlepszy model globalnie według RandomSearch; Random Search - najlepszy model dla danego zbioru; BayesSearch - najlepszy model dla danego zbioru danych.

8 Podsumowanie

Celem projektu było zbadanie tunowalności poszczególnych algorytmów. Porównano skuteczność dwóch metod strojenia hiperparametrów - Random Search i optymalizacji bayesowskiej - dla modeli Decision Tree, Random Forest oraz XGBoost. Obie techniki poprawiły wyniki względem konfiguracji domyślnych, przy czym największą tunowalność zaobserwowano dla Decision Tree. Random Forest i XGBoost okazały się mniej wrażliwe na zmianę hiperparametrów, co potwierdza dobre dopasowanie ich ustawień domyślnych. Optymalizacja bayesowska osiągała nieco lepsze rezultaty niż wyszukiwanie losowe, zwłaszcza dla bardziej złożonych modeli.

Literatura

- [1] Philipp Probst, Anne-Laure Boulesteix & Bernd Bischl. *Tunability: Importance of Hyperparameters of Machine Learning Algorithms*. *Journal of Machine Learning Research*, 20 (2019), 1–32.

Załącznik

Decision Tree		Random Forest		XGBoost	
Dataset 0:	0.7714	Dataset 0:	0.8472	Dataset 0:	0.8540
<i>max_depth</i>	4.0	<i>max_depth</i>	14.0	<i>learning_rate</i>	0.01
<i>max_features</i>	0.8463	<i>max_features</i>	0.1045	<i>max_depth</i>	1.0
<i>min_samples_leaf</i>	10.0	<i>min_samples_leaf</i>	9.0	<i>n_estimators</i>	4142.0
<i>min_samples_split</i>	19.0	<i>n_estimators</i>	1987.0	<i>subsample</i>	0.4153
Dataset 1:	0.8812	Dataset 1:	0.9204	Dataset 1:	0.9101
<i>max_depth</i>	5.0	<i>max_depth</i>	6.0	<i>learning_rate</i>	0.01
<i>max_features</i>	0.1258	<i>max_features</i>	0.1162	<i>max_depth</i>	15.0
<i>min_samples_leaf</i>	4.0	<i>min_samples_leaf</i>	1.0	<i>n_estimators</i>	564.0
<i>min_samples_split</i>	9.0	<i>n_estimators</i>	1228.0	<i>subsample</i>	0.1000
Dataset 2:	0.8536	Dataset 2:	0.8791	Dataset 2:	0.8801
<i>max_depth</i>	8.0	<i>max_depth</i>	8.0	<i>learning_rate</i>	0.01
<i>max_features</i>	0.2860	<i>max_features</i>	0.1000	<i>max_depth</i>	1.0
<i>min_samples_leaf</i>	7.0	<i>min_samples_leaf</i>	10.0	<i>n_estimators</i>	5000.0
<i>min_samples_split</i>	5.0	<i>n_estimators</i>	2000.0	<i>subsample</i>	0.1000
Dataset 3:	0.8067	Dataset 3:	0.8063	Dataset 3:	0.8098
<i>max_depth</i>	12.0	<i>max_depth</i>	30.0	<i>learning_rate</i>	0.01
<i>max_features</i>	0.9000	<i>max_features</i>	0.8982	<i>max_depth</i>	10.0
<i>min_samples_leaf</i>	9.0	<i>min_samples_leaf</i>	5.0	<i>n_estimators</i>	100.0
<i>min_samples_split</i>	20.0	<i>n_estimators</i>	143.0	<i>subsample</i>	0.8337
Dataset 4:	0.8423	Dataset 4:	0.8634	Dataset 4:	0.8619
<i>max_depth</i>	7.0	<i>max_depth</i>	21.0	<i>learning_rate</i>	0.01
<i>max_features</i>	0.6300	<i>max_features</i>	0.4371	<i>max_depth</i>	15.0
<i>min_samples_leaf</i>	10.0	<i>min_samples_leaf</i>	8.0	<i>n_estimators</i>	100.0
<i>min_samples_split</i>	2.0	<i>n_estimators</i>	2000.0	<i>subsample</i>	0.2992

Rysunek 6: Najlepsze modele otrzymane po optymalizacji bayesowskiej dla poszczególnych zbiorów danych.

Dataset	Mean Test Score Default	Mean Test Score Global Random	Mean Test Score Random	Mean Test Score Bayes	Diff Default vs Global	Diff Default vs Random	Diff Default vs Bayes	Model
0	0.6641	0.7731	0.7731	0.7714	0.1090	0.1090	0.1073	Decision Tree
1	0.7522	0.8659	0.8783	0.8812	0.1137	0.1261	0.1290	Decision Tree
2	0.7067	0.8459	0.8569	0.8536	0.1391	0.1502	0.1468	Decision Tree
3	0.7448	0.8112	0.8112	0.8067	0.0664	0.0664	0.0618	Decision Tree
4	0.6880	0.8171	0.8408	0.8423	0.1291	0.1527	0.1543	Decision Tree
0	0.8320	0.8454	0.8454	0.8472	0.0134	0.0134	0.0152	Random Forest
1	0.9065	0.9176	0.9176	0.9204	0.0111	0.0111	0.0140	Random Forest
2	0.8670	0.8778	0.8778	0.8791	0.0109	0.0109	0.0121	Random Forest
3	0.7915	0.8008	0.8052	0.8063	0.0094	0.0137	0.0148	Random Forest
4	0.8483	0.8573	0.8646	0.8634	0.0090	0.0163	0.0150	Random Forest
0	0.8236	0.8286	0.8286	0.8540	0.0050	0.0050	0.0304	XGBoost
1	0.8964	0.8968	0.8994	0.9101	0.0004	0.0031	0.0137	XGBoost
2	0.8559	0.8654	0.8654	0.8801	0.0096	0.0096	0.0243	XGBoost
3	0.7950	0.7932	0.7966	0.8098	-0.0018	0.0016	0.0148	XGBoost
4	0.8439	0.8563	0.8563	0.8619	0.0123	0.0123	0.0180	XGBoost

Tabela 3: Porównanie wyników modeli w zależności od metody optymalizacji hiperparametrów. Default - model z domyślnymi hiperparametrami; Global Random - najlepszy model globalnie według RandomSearch; Random Search - najlepszy model dla danego zbioru; BayesSearch - najlepszy model dla danego zbioru danych według BayesSearch