

Tunowalność Algorytmów ML

Automatyczne Uczenie Maszynowe – Projekt 1

Igor Kołodziej Hubert Kowalski Julia Przybytniowska

Spis treści

01	Cel projektu	05	Wyniki eksperymentów
02	Zbiory danych	06	Tunowalność
03	Modele i siatki hiperparametrów Techniki	07	Bias Sampling
04	losowania punktów	08	Podsumowanie

Cel projektu



Cel



- Optymalne hiperparametry **domyślne** lepsze niż wartości *default* w *scikit-learn*
- 2 techniki losowania punktów:
 - Rozkład jednostajny
 - Metoda bayesowska
- Tunowalność algorytmów ML na 4 zbiorach danych [1]





Zbiory Danych



Wybrane **zbiory danych**





1000 obserwacji, 20 zmiennych objaśniających, 0 wartości brakujących

Diabetes (ID 37)

768 obserwacji, 8 zmiennych objaśniających, 0 wartości brakujących

Wine (<u>ID 43980</u>)

2554 obserwacji, 11 zmiennych objaśniających, 0 wartości brakujących

Shrutime (<u>ID 45062</u>)

10000 obserwacji, 10 zmiennych objaśniających, 0 wartości brakujących



Preprocessing danych

Zmienna celu



Label Encoder

Zmienne numeryczne



Standard Scaler





One Hot Encoding

Column Transformer



Ewaluacja modeli

Metryka

Testowane: F1 Score, AUC, Accuracy



Wybrane: AUC

- Problem przeuczenia modelu
 - CV
 - 5 fold
- Powtarzalność wyników:
 - RANDOM_SEEDS

Modele i ich siatki hiperparametrów



Logistic Regression (ElasticNet)



Algorytm	Hyperparametr	Dolny zakres	Górny zakres
ElasticNet	C	1e-4	1e4
	penalty	elasticnet	
	l1_ratio	1e-4	1.0
	class_weight	balanced	
	max_iter	1500	
	solver	saga	

Tabela 1: Siatka hiperparametrów dla ElasticNet.





Algorytm	Hiperparametr	Dolny zakres Górny zakres	
Extra Trees	$n_{estimators}$	10	1000
	criterion	gini, entropy, log_loss	
	bootstrap	True	
	$\max_samples$	0.5	1.0
	\max_{features}	0.1	0.9
	$\min_{\text{samples_leaf}}$	0.05	0.25

Tabela 2: Siatka hiperparametrów dla Extra Trees Classifier.



Algorytm	Hiperparametr	Dolny zakres	Górny zakres
	$n_{estimators}$	10	2000
	learning_rate	1e-4	0.4
	subsample	0.25	1.0
	booster	gbtree	
XGBoost	\max_{depth}	1	15
AGDOOSt	$\min_{child_{weight}}$	1	128
	$colsample_bytree$	0.2	1.0
	$colsample_bylevel$	0.2	1.0
	reg_alpha	1e-4	512.0
	reg_lambda	1e-3	1e3

Tabela 3: Siatka hiperparametrów dla XGBoost Classifier.

Techniki losowania punktów



Metody losowania punktów

Technika losowania	Liczba iteracji	
Random Search	300	
Tree-Structured Parzen Estimator (Bayes Search)	(punkty startowe) 3 * 100 = 300	



Biblioteka implementująca



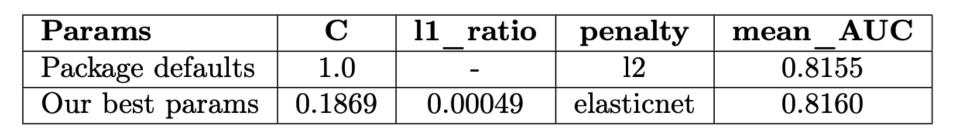


(O) OPTUNA



Wyniki eksperymentów

Optymalne wartości domyślne hiperparamerów - porównanie dla Logistic Regression







Wyniki optymalizacji modeli – ile możemy zyskać?



Model	Random Search mean best AUC	Bayesian Search mean best AUC	Default
LogisticRegression	0.8165	0.8166	0.8159
ExtraTreesClassifier	0.7979	0.8047	0.7946
XGBClassifier	0.8433	0.8527	0.8403

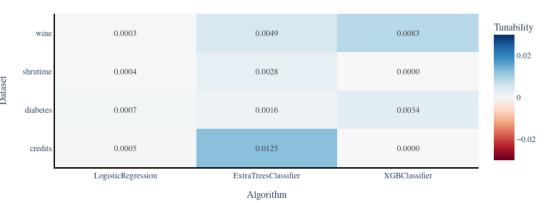




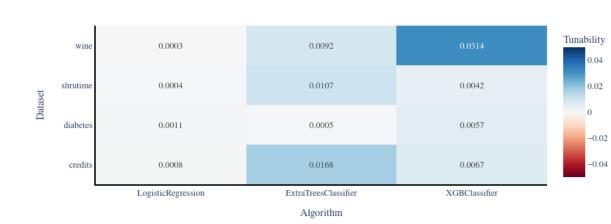


Tunowalność

Random Search Tunability Heatmap



Bayesian Search Tunability Heatmap



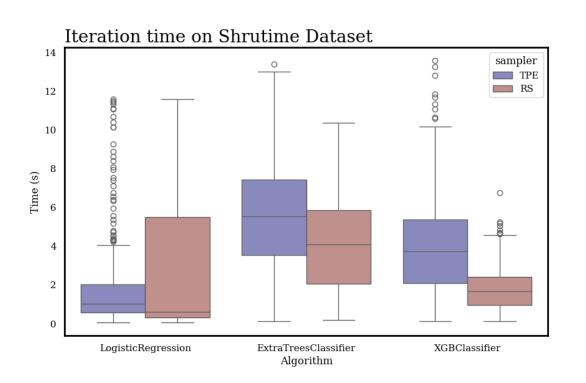


Czas strojenia - przykład dla zbioru Shrutime

Sampler	Algorytm	Czas strojenia (min)	AUC
	ExtraTreesClassifier	20.422524	0.8099
RS	LogisticRegression	15.317961	0.8328
	XGBClassifier	9.063305	0.8645
	ExtraTreesClassifier	28.442875	0.8194
TPE	LogisticRegression	9.626361	0.8328
	XGBClassifier	20.624290	0.8684



Czas pojedynczej iteracji

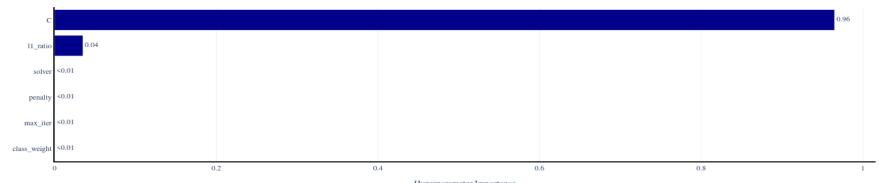




Istotność hiperparametrów



Logistic Regression Hyperparameter Importances based on Joint Tuning



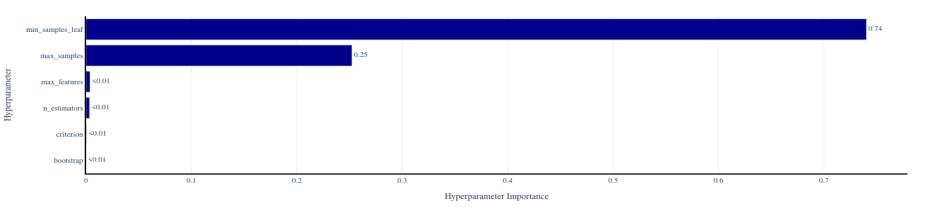
Hyperparameter Importance



Istotność hiperparametrów



Extra Trees Hyperparameter Importances based on Joint Tuning

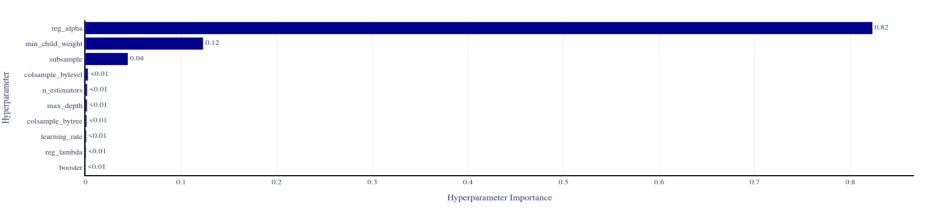




Istotność hiperparametrów

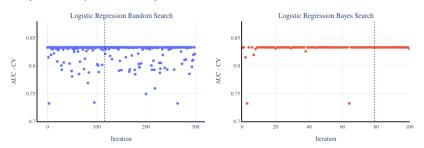


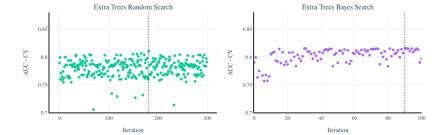
XGBoost Hyperparameter Importances based on Joint Tuning

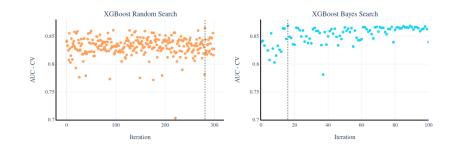


Bias Sampling

Optimization History for All Models and Optimizers on Shrutime Dataset





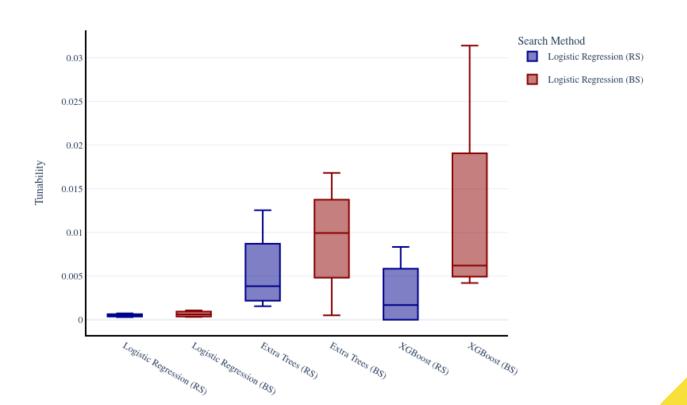




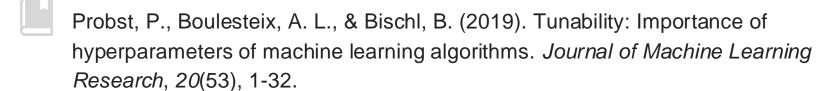


Wpływ techniki losowania na tunowalność

Tunability of Models with Random Search and Bayesian Search







- Watanabe, S. (2023). Tree-structured parzen estimator: Understanding its algorithm components and their roles for better empirical performance. arXiv preprint arXiv:2304.11127.
- Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019, July). Optuna: A next-generation hyperparameter optimization framework. In Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining (pp. 2623-2631).



Dziękujemy za uwagę!

XGBoost Contour Plot for learning_rate ad reg_alpha on Diabetes Dataset

