

# Automatyczne Uczenie Maszynowe

## Projekt 1

Wiktoria Boguszevska i Mateusz Zacharecki

20 listopada 2024

Problemy do rozważenia:

- 1 ile iteracji każdej metody potrzebujemy żeby uzyskać stabilne wyniki optymalizacji,
- 2 określenie zakresów hiperparametrów dla poszczególnych modeli,
- 3 tunowalność poszczególnych algorytmów,
- 4 czy technika losowania punktów wpływa na różnice we wnioskach w punkcie 3. dotyczącym tunowalności algorytmów, czy występuje bias sampling.

# Wprowadzenie

Celem projektu jest przeanalizowanie tunowalności hiperparametrów 3 wybranych algorytmów uczenia maszynowego na co najmniej 4 zbiorach danych. Do tunowania modeli należy wykorzystać min. 2 różne techniki losowania punktów.

W projekcie wykorzystujemy 4 różne zbiory danych dla problemu klasyfikacji binarnej:

- Breast Cancer,
- Credit,
- Blood Transfusion,
- Bank marketing.

# Metody samplingu

Główna część projektu polegała na zastosowaniu różnych metod tunowalności hiperparametrów. Wykonaliśmy to dla 3 algorytmów: regresja logistyczna, KNN oraz XGBoost. Zastosowane przez nas metody samplingu to:

- Random Search,
- Bayes Search,
- Grid Search.

## Wyniki



# Wyniki dla poszczególnych metod

MODEL	LOGISTYCZNA				KNN				XGBOOST			
ZBIÓR	cancer	credit	blood	bank	cancer	credit	blood	bank	cancer	credit	blood	bank
TRAIN ACCURACY	0.9736	0.7175	0.7843	0.8827	0.967	0.7138	0.8027	0.8866	0.9363	0.7025	0.7792	0.8827
TEST ACCURACY	0.9825	0.68	0.72	0.8928	0.9912	0.705	0.7533	0.8961	0.9561	0.685	0.74	0.8928

**Tabela:** Wyniki otrzymane za pomocą metody Random Search.

MODEL	LOGISTYCZNA				KNN				XGBOOST			
ZBIÓR	cancer	credit	blood	bank	cancer	credit	blood	bank	cancer	credit	blood	bank
TRAIN ACCURACY	0.9736	0.7213	0.786	0.8836	0.967	0.7175	0.8094	0.8866	0.9758	0.7188	0.796	0.8858
TEST ACCURACY	0.9825	0.685	0.72	0.8961	0.9737	0.685	0.7533	0.8961	0.9649	0.685	0.76	0.9028

**Tabela:** Wyniki otrzymane za pomocą metody Bayes Search.

MODEL	LOGISTYCZNA				KNN				XGBOOST			
ZBIÓR	cancer	credit	blood	bank	cancer	credit	blood	bank	cancer	credit	blood	bank
TRAIN ACCURACY	0.9692	0.72	0.786	0.8836	0.967	0.7113	0.8044	0.8852	0.967	0.7075	0.7994	0.8852
TEST ACCURACY	0.9912	0.69	0.72	0.895	0.9737	0.705	0.7467	0.8983	0.9649	0.72	0.7533	0.8961

**Tabela:** Wyniki otrzymane za pomocą metody Grid Search.

## Analiza wyników

- 50 iteracji każdej z metod.
- 30 iteracji wewnątrz metody bayesian search.
- Najlepsza metoda: **bayesian search**.

# Zakres hiperparametrów

Algorytm	Hiperparametr	Typ	Dolna granica	Górna granica	Rozkład
LogisticRegression	<b>C</b>	numeric	$10^{-10}$	$10^{10}$	loguniform
	<b>l1_ratio</b>	numeric	0	1	uniform
KNeighborsClassifier	<b>n_neighbors</b>	integer	1	30	uniform
	<b>weights</b>	discrete	–	–	uniform
	<b>p</b>	integer	1	2	uniform
XGBClassifier	<b>booster</b>	discrete	–	–	uniform
	<b>eta</b>	numeric	$10^{-10}$	1	loguniform
	<b>n_estimators</b>	integer	1	500	uniform
	<b>max_depth</b>	integer	1	10	uniform
	<b>subsample</b>	numeric	0.1	1	uniform
	<b>colsample_bytree</b>	numeric	0	1	uniform
	<b>colsample_bylevel</b>	numeric	0	1	uniform
	<b>lambda</b>	numeric	$10^{-5}$	$10^2$	loguniform
	<b>alpha</b>	numeric	$10^{-5}$	$10^2$	loguniform
	<b>min_child_weight</b>	numeric	1	$10^7$	loguniform

**Tabela:** Zakresy hiperparametrów dla metody Random Search oraz Bayes Search.

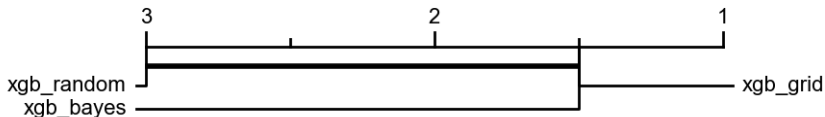
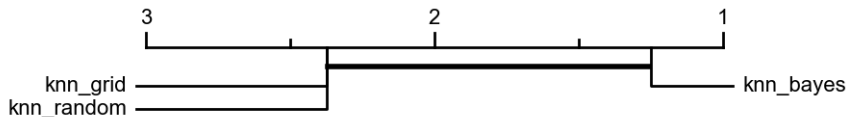
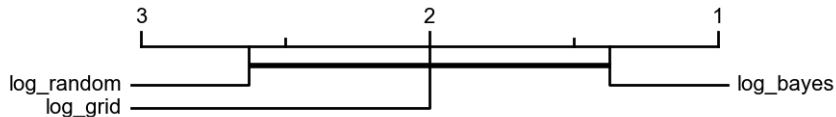
# Tunowalność poszczególnych algorytmów

Model	Defaultowe hiperparametry	Tunowalność
Regresja logistyczna	C: 6.563234996980466e-06, l1_ratio: 0.34345601404832493	$-3.5014 \cdot 10^{-6}$
kNN	weights: uniform, p: 2, n_neighbors: 18	-0.0004
XGBoost	alpha: 0.006066641915981054, booster: gbtree, colsample_bylevel: 0.141194832453028, colsample_bytree: 0.5588264346475861, eta: 2.4708290259874402e-08, lambda: 2.051302466178794, max_depth: 8, min_child_weight: 5.4627652310036785, n_estimators: 237, subsample: 0.2477228436396202	0

**Tabela:** Defaultowe hiperparametry i tunowalność dla metody Random Search.

- Test Manna-Whitneya do porównania różnic pomiędzy technikami losowania hiperparametrów.
- Porównanie metod randomized search oraz bayesian search.
- Różnica wyników z obu modeli jest istotna statystycznie, a wyniki pochodzą z różnych rozkładów.
- Występuje bias sampling.

# Critical Difference Diagrams



- [1] Philipp Probst, Anne-Laure Boulesteix, Bernd Bischl; Tunability: Importance of Hyperparameters of Machine Learning Algorithms



Dziękujemy za uwagę