

Projekt 1

Analiza danych

Gabriela Matuszewska, Mateusz Wojtulewicz

1 Opis datasetu

Do analizy wybrano dataset "Adult income", który zawiera klasyfikacje przychodu rocznego ($< 50k$, $\geq 50k$) w zależności, od między innymi pochodzenia, wykształcenia, miejsca zamieszkania.

Liczba cech i ich typy oraz rozmiar:

```
df.shape
```

(48842, 15)

```
df.head()
```

	age	workclass	fnlwgt	education	educational-num	marital-status	occupation	relationship	race	gender	capital-gain	capital-loss	hours-per-week	native-country	income
0	25	Private	226802	11th	7	Never-married	Machine-op-inspct	Own-child	Black	Male	0	0	40	United-States	0
1	38	Private	89814	HS-grad	9	Married-civ-spouse	Farming-fishing	Husband	White	Male	0	0	50	United-States	0
2	28	Local-gov	336951	Assoc-acdm	12	Married-civ-spouse	Protective-serv	Husband	White	Male	0	0	40	United-States	1
3	44	Private	160323	Some-college	10	Married-civ-spouse	Machine-op-inspct	Husband	Black	Male	7688	0	40	United-States	1
4	18	Private	103497	Some-college	10	Never-married	Prof-specialty	Own-child	White	Female	0	0	30	United-States	0

2 Czyszczenie danych

Ze względu na specyfikacje zestawu danych konieczna była pewna modyfikacja tych danych. Brakujące wartości zostały zastąpione przez najczęściej występujące wartości. Ze względu na przeważająca liczbę `native-country="United-States"`, która mogłaby zakłócić wyniki postanowiono usunąć kolumnę `native-country`. Kolejnym spostrzeżeniem jest podobieństwo kolumn `education` i `educational-num` oraz `marital-status` i `relationship` dlatego zostały one zmapowane do jednej kolumny. Atrybut `workclass` zastępuje nam atrybut `occupation` dlatego możemy pozbyć się tej kolumny. Po usunięciu zbędnych kolumn oraz wypełnieniu brakujących wartości, a także zmapowaniu kategorii wypełnionych typem `"object"` na `int`, nasze dane wyglądają następująco:

3 Random Forest Classification

Las losowy to zbiór drzew klasyfikacyjnych o podziałach binarnych. Dla konkretnej obserwacji (wyrażonej jako wektor wejściowy), każde z drzew zwraca decyzję lub krótkie prawdopodobieństwa klasyfikacji, Prawdopodobieństwa z drzew wchodzących w skład lasu są traktowane jako głosy – > jako wynik zwracana jest decyzja która otrzymała najwięcej głosów (której średnie prawdopodobieństwo jest najwyższe).

3.1 Parametry

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100,
*, criterion='gini', max_depth=None, min_samples_split=2,
min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto',
max_leaf_nodes=None, min_impurity_decrease=0.0,
min_impurity_split=None, bootstrap=True, oob_score=False,
n_jobs=None, random_state=None, verbose=0, warm_start=False,
class_weight=None, ccp_alpha=0.0, max_samples=None)
```

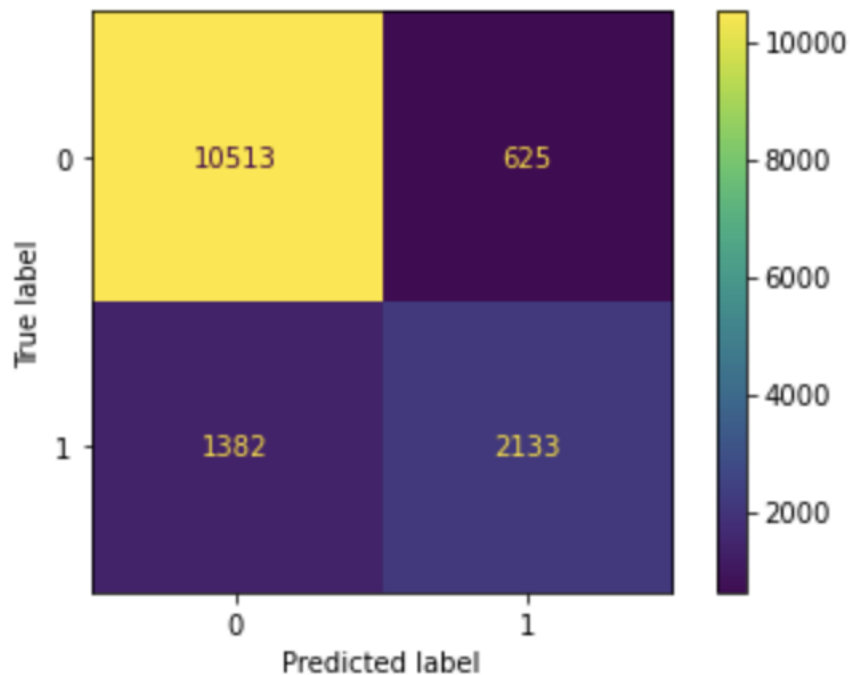
Opis najważniejszych parametrów:

1. `n_estimators`
2. `criterion` “gini”, “entropy”
3. `max_depth`
4. `min_samples_split`
5. `min_samples_leaf`
6. `min_weight_fraction_leaf`
7. `max_features` “auto”, “sqrt”, “log2”
8. `max_leaf_nodes`
9. `min_impurity_decrease`

4 Sposób wyboru zbioru testowego

5 Macierz pomyłek

Macierz pomyłek (confusion matrix) to narzędzie do oceny jakości klasyfikacji. Składa się z następujących klas: true-positive, true-negative, false-positive, false-negative. W naszym przypadku wygląda tak:



6 Walidacja Krzyżowa

Walidacja krzyżowa służy do testowania klasyfikatora i polega na wybraniu parametru n (zazwyczaj 10) i podzieleniu danych na n równych podzbiorów, z których każdy następnie przyjmujemy jako testowy, a pozostałe tworzą próbkę treningową. Wyniki klasyfikacji są sumowane.

6.1 Grid Search

TO DO: NA CZYM POLEGA

7 Normalizacja

Normalizacja polega na przeskalowaniu wartości do zakresu $[0,1]$:

$$X^* = \frac{X - \min(X)}{\max(x) - \min(X)}$$

8 Standaryzacja

Wartości mniejsze od średniej wartości po standaryzacji będą miały ujemne wartości, analogicznie wartości większe od średniej po standaryzacji będą miały

dodatnie wartości:

$$X^* = \frac{X - \mu(X)}{\sigma(x)}$$

9 PCA

PCA, czyli analiza głównych składowych, służy do redukcji liczby zmiennych (wymiarowości), tak aby najlepiej zachować strukturę danych lub do odkrycia prawidłowości między cechami. Opiera się o wyznaczenie osi zachowującej największą wartość wariancji zbioru uczącego. Polega na wyznaczeniu składowych będących kombinacją liniową badanych zmiennych.