

# Code Review funkcji grupy Gakubu

Jan Kruszewski

5 maja 2022

## 1 Wprowadzenie

W poniższym dokumencie przeprowadzam Code Review funkcji przygotowanej przez grupę Gakubu, której praca dotyczy pakietu AutoGluon. Celem funkcji było przygotowanie dowolnego tabularycznego zbioru danych opisującego binarny, problem predykcyjny do użycia frameworku AutoGluon.

## 2 Czy kod osiąga postawiony cel?

Postawiony cel sprowadzał się do 4 punktów: - wczytania zbiorów danych z benchmarku do przetestowania funkcji - przeprowadzenie preprocessingu - modelowanie

### 2.1 Wczytanie danych

W celu wczytania zbiorów zespół Gakub stworzył pomocniczą funkcję `create_data_frame_from_arff`, która przyjąwszy na wejściu adres url zwraca ramkę danych. Funkcja jest, krótka, czytelna i co najważniejsze spełnia swój cel.

### 2.2 Preprocessing

Preprocessing wykonywany jest wewnątrz głównej funkcji `make_kfold_cross_valiation`. Preprocessing obejmuje encoding zmiennych w formacie string za pomocą systemu Unicode UTF-8. Pozostały preprocessing i feature engineering został pozostawiony AutoGluon'owi, który wykonuje część pracy za użytkownika.

### 2.3 Modelowanie

Tworzenie modeli oraz ich walidacja wykonywana jest w funkcji `make_kfold_cross_valiation`. Na wejściu podaje poza argumentami funkcji `create_data_frame_from_arff` nazwę zbioru, nazwę zmiennej celu, ilość foldów do krosvalidacji oraz inne parametry dla krosvalidacji. Funkcja zwraca rezultaty krosvalidacji w postaci tablicy, zawierającej rezultaty dla kolejnych foldów w postaci nazwy modelu, wyniku dla metryki `roc_auc` oraz nazwy zbioru. Funkcja po przetestowaniu działa bezbłędnie, można by jednak moim zdaniem pokusić się o umieszczenie w wynikowej tablicy wyników dla również innych metryk takich jak np. `f1 Score` czy `accuracy`.

## 3 Czy w kodzie są jakieś oczywiste błędy logiczne?

Po dokładnej analizie kodu nie znalazłem żadnych błędów logicznych.

## 4 Czy wymagania z prezentacji zostały w pełni zaimplementowane?

W prezentacji grupa Gakubu przedstawiła wyniki swojej pracy w pełni wynikające z pracy wykonanej nad kodem. Wszystkie założenia zostały zaimplementowane w rzetelny sposób.

## 5 Czy kod jest zgodny z wytycznymi PEP8?

Zdecydowanie można odpowiedzieć twierdząco na zadane powyżej pytanie. Funkcje oraz zmienne zostały nazwane w odpowiedni sposób. Twórcy operują pustymi przestrzeniami by zwiększyć czytelność kodu. Nie przekraczane są maksymalne długości linii kodu (można to zauważyć np. przy deklaracji zmiennych przy funkcji `make_kfold_cross_valiation`). Co robi wrażenie autorzy kodu nie zapomnieli również o potrójnym cudzysłowie przy opisywaniu funkcji oraz nie dodawania spacji przy znaku '=' wewnątrz funkcji. Jedyną rzeczą do której można się przyklepić to drobne literówki.

## 6 Czy są jakieś obszary, w których kod mógłby zostać poprawiony?

Jak wspominałem wcześniej sądzę, że można dodać inne metryki do oceny jakości modeli. Nie widzę innych oczywistych sposobów poprawy kodu, ale z chęcią przeczytam co na ten temat do powiedzenia mają inne grupy.

## 7 Czy dokumentacja i komentarze są wystarczające?

Uważam, że objaśnienia funkcji oraz komentarze są zdecydowanie zaletą kodu zespołu Gakubu. Każda funkcja została opisana w jasny sposób. Komentarze wewnątrz kodu również pomagają zrozumieć dane fragmenty kodu.

## 8 Czy udało się odtworzyć zamieszczone przykłady w kodzie?

Przetestowałem funkcje grupy Gakubu na moim komputerze odpalając je dla 4 zbiorów danych z benchmarku. Uważałem uruchamianie jej dla wszystkich zbiorów za zbędne. Wszystko działało w odpowiedni sposób opisany powyżej.

## 9 Czy udało się użyć kodu na nowych danych?

Za pomocą kodu, co cieszy udało się stworzyć model dla innego binarnego zbioru danych 'liver disorders'.

## 10 Podsumowanie

Kod będący rezultatem kamienia milowego 2 jest przejrzysty oraz co najważniejsze działa. Kod spełnia wytyczne PEP8 co ułatwia jego czytanie. Pomijając drobne literówki w dokumentacji kod nie pozostawia dużego pola do krytyki.