

AutoML - FLAML

Group KTR

1 Outline of introduction

Sources that we will be quoting:

- <https://microsoft.github.io/FLAML/docs/reference/model> – everything about the model
- https://new.openml.org/search?type=benchmark&sort=tasks_included&study_type=taskid=218 – data sets for evaluation

2 Description of FLAML framework

FLAML (Fast and Lightweight AutoML) is a Python library, developed by Microsoft under Open Source license, used for Auto Machine Learning (AutoML). According to its creators, FLAML is a lightweight package that finds the right machine learning model for given task in an automatic, effective and economic way. Hence it should free a user from the necessity of choosing an optimal machine learning model and its hyperparameters. This framework supports both conventional machine learning methods - quickly finds qualitative models while using low amount of computational resources for common tasks, such as classification or regression - and deep learning.

The main class of the package is `flaml.AutoML`. Its constructor takes kwarg `settings`, which includes time budget, metrics, task (classification, regression, etc.), list of models to be used and others. The most important method, just like in scikit-learn package, is `fit()`, which takes data frame of features `X` and label `y` (it can also overwrite settings). While calling this method suitable models are fitted with data and optimization of hyperparameters is performed on them. With help of `best_estimator` and `best_config` attributes it is possible to obtain the best model and its configuration of hyperparameters acquired by FLAML.

Apart from `fit()` we can call `predict()` and `predict_proba()` methods, which allows compatibility with scikit-learn package.

Another advantage of FLAML is availability of personalising software - the package can handle an user defined function (most conveniently by using functions compatible with scikit-learn), so ones including `fit()` and `predict()` methods) and use an exterior metrics.

Module `tune` of this library enables acquiring the right parameters for given model. It is used internally by `flaml.AutoML`, but it can also be used directly by the user, meeting one of the conditions:

- user's task is not one of the built-in `flaml.AutoML` tasks
- user's input data cannot be presented as features data frame and label column
- the user wants to tune a function that does not qualify as a machine learning procedure

A considerable disadvantage of the package is its incomplete documentation. For example, preprocessing done by the package is not described in any place (there was only a mention of `DataTransformer` class in SDK), which led us to do our own preprocessing.

3 Actions necessary to engage to make the package work properly

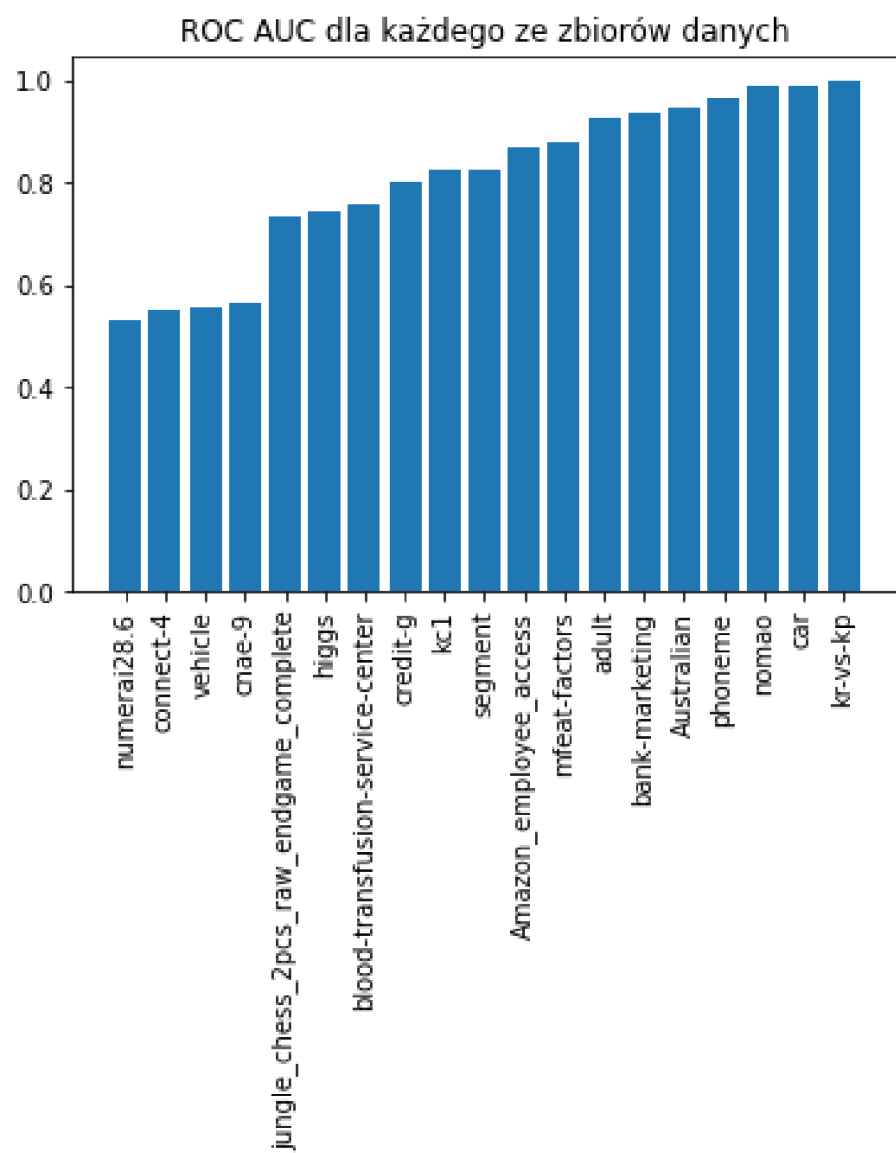
Due to poor documentation, we reckoned that the package does not perform preprocessing of data. Hence to eliminate undefined behaviour we had to perform preprocessing, so that the user does not have to handle this matter. Actions that we decided to include are:

- removing columns that are IDs
- imputation by mode
- replacing outliers with ceiling value
- normalization with `QuantileTransformer`
- scaling with `StandardScaler`
- encoding of categorical data with `WoEEncoder`

We chose that kind of preprocessing not only to guarantee the lack of undefined actions in the case of binary classification, but also to avoid overfitting and to acquire better models. It turns out however that FLAML includes necessary preprocessing for the package to work. Similarly to ours, it removes columns having zero variance and imputes missing data either with a new category in case of categorical features, or with a median in case of numerical ones. It also fragments dates to smaller components.

4 Evaluation results

As an initial evaluation we checked the performance of a function created by us on twenty data sets from OpenML, included in a task available under the link (we have not used each of the data sets included). 13 out of these data sets represent binary classification. On each of them we performed the `fit()` method from a `OurFlaml` class created by us (for that to happen we had to download each of the data sets split it into features data frame and label column). We appended the best ROC AUC score obtained by the function, the time to train the best configuration and the name of the data set to appropriate lists. Then we visualized the results:



Czas trenowania najlepszej konfiguracji dla każdego ze zbiorów danych

