

CMPT459 Fall 2020
Data Mining
Martin Ester
TAs: Madana Krishnan Vadakandara Krishnan
and Rhea Rodriguez

Milestone 2 of the Course Project

Deadline: November 19

Total marks: 100

Introduction

In the first milestone, you have completed the data pre-processing steps. You should now have a *dataset* that is cleaned to an extent, obtained by merging the *individual* cases and *locations* datasets.

In the second milestone, you will be building various classification models and use metrics to evaluate the performance of the models. Each group member has to build one classification model. Thus each group will have 2 or 3 models built, depending on the group size. Remind yourselves that the problem statement for this project is to predict the outcome of a case.

Tasks

2.1 Splitting dataset (10 marks)

Split the dataset into training dataset and testing dataset. Train to test ratio should be 80:20.

2.2 Build models (40 marks)

As mentioned above, each individual in a group has to build a classification model. The models can be of any type, and you can use any existing Python (or R) libraries (ex: Scikit-Learn) to build them. One out of the two/three models **MUST** be a variant of the boosting tree (ex: XGBoost, AdaBoost, LightGBM). The other one/two models could include SVMs, KNN, Decision Trees, Random Forests, MLPs, Naive Bayes or any other classifiers that you think could be appropriate for the problem statement.

Save each trained model to your disk (if you choose .pkl file, the models would look like `xgb_classifier.pkl`, `rf_classifier.pkl`), and include the models in submission.

2.3 Evaluation (30 marks)

Load the saved models from task 2.2. Once loaded, evaluate the models on both the training set and the test set. Choose appropriate metrics to do the evaluation. Report the scores of your metrics for the training and the test set and interpret them. Which of the metric(s) are most important for this classification problem?

2.4 Overfitting (20 marks)

It is not uncommon for the classification models to overfit. Do you observe overfitting in the models that you trained? How do you check for overfitting? Explain steps taken by using plots and/or metrics

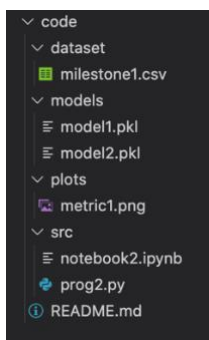
evaluated. You can vary the values of at least one hyperparameter, train models for different values of that hyperparameter and then compare the performance metric on training and test dataset.

Submission (Code + Report)

A). Code

Submit a 'code.zip' file with the following contents. It should contain the code that you have written for the above tasks. The code could contain how train-test split was done, steps performed to fit the classification models, saving the models as .pkl files, loading the .pkl files to perform evaluation, different evaluation metrics used. Make sure to submit the saved models in the submission. You can include any plots that you might generate.

Please note that in order for the TAs to run the code for milestone 2, you would need to provide the *dataset* that you generated from milestone 1. Include the *dataset* within a folder named 'dataset' as shown in the figure. The structure can look like this.



B). Report

Briefly explain your work for task 2.1. Explain in more detail your work for tasks 2.2, 2.3 and 2.4. Reports should **NOT** be more than 2 pages. Submit a 'report.pdf' file.

NOTE:

- For milestone 2, you can vary the values for at least one of the hyperparameters to check for overfitting. In milestone 3, you will be tuning all the combinations of hyperparameters to find the best model.
- Include the code to save and load the models in the submission. Ensure to submit the two/three saved models (.pkl files).
- Please include steps for execution in the README.md file.