

[< Return to Classroom](#)

Plagiarism Detector

审阅

代码审阅

HISTORY

符合要求

Dear student,

You did an excellent project.
Keep up the hard work.

All Required Files and Tests

The submission includes complete notebook files as `.ipynb`:
"2_Plagiarism_Feature_Engineering" and "3_Training_a_Model". And the test and helper files are included: "problem_unittests.py", "helpers.py". The submission also includes a training directory `source_sklearn` OR `source_pytorch`.

The two notebooks and the training directory are included

All the unit tests in project have passed.

All the unit tests passed.

Notebook 2: DataFrame Pre-Processing

The function `numerical_dataframe` should be complete, reading in the original `file_information.csv` file and returning a DataFrame of information with a numerical `Category` column and new, `Class` column.

The function `numerical_dataframe` works as expected.

There is no code requirement here, just make sure you run all required cells to create a `complete_df` that holds pre-processed file text data and `Datatype` information.

Perfect, the `preprocessed_df` data frame is created successfully.

Notebook 2: Features Created

The function `calculate_containment` should be complete, taking in the necessary information and returning a single, normalized containment value for a given answer file.

The function `calculate_containment` works as expected.

Provide an answer to the question about containment feature calculation.

Good answer! Indeed, there is no correlation between training and test data regarding the containment feature calculation.

The function `lcs_norm_word` should be complete, taking in two texts and returning a single, normalized LCS value.

Define an n-gram range to calculate multiple containment features. Run the code to calculate one LCS feature, and create a DataFrame that holds all of these feature calculations.

Notebook 2: Train and Test Files Created

Complete the function `train_test_data`. This should return only a *selection* of training and test features, and corresponding class labels.

The function `train_test_data` works as expected.

Select at least three features to use in your final training and test data.

Provide an answer that describes why you chose your final features.

Good justification of the choice.

Implement the `make_csv` function. The class labels for train/test data should be in the first column of the csv file; selected features in the rest of the columns. Run the rest of the cells to create `train.csv` and `test.csv` files.

Notebook 3: Data Upload

Upload the `train.csv` file to a specified directory in an S3 bucket.

Notebook 3: Training a Custom Model

Complete at least *one* of the `train.py` files by instantiating a model, and training it in the main if statement. If you are using a custom PyTorch model, you will have to complete the `model.py` file, as well (you do not have to do so if you choose to use an imported sklearn model).

Define a custom sklearn OR PyTorch estimator by passing in the required arguments.

Fit your estimator (from the previous rubric item) to the training data you stored in S3.

Notebook 3: Deploying and Evaluating a Model

Deploy the model and create a `predictor` by specifying a deployment instance.

The estimator is deployed. The instance ml.t2.medium is appropriate.

Pass test data to your deployed `predictor` and evaluate its performance by comparing its predictions to the true, class labels. Your model should get at least 90% test accuracy.

Provide an answer to the two model-related questions.

Notebook 3: Cleaning up Resources

Run the code to clean up your final model resources.

 下载项目

[返回 PATH](#)