

Machine Learning Engineer Nanodegree

New York City Taxi Trip Duration

Cosa Santos

June 25, 2019

1 Domain Background

The chosen project for the Capstone Project is the [New York City Taxi Duration Trip](#) competition from [Kaggle](#). The challenge is to build a model that predicts the total ride duration of taxi trips in New York City.

Solving problems alike is not a new challenge [1,2], however, it was not until recently that the amount of data has increased enough, and become more precise with the rise of high-tech smartphones and GPS use and real-time tracking, that solutions were made possible.

1.0.1 2. Problem Statement

By considering historical information such as pick up and drop off date, hour and geo-localization, and trip duration regarding taxi trips, the objective is to predict the duration of each taxi trip given in a specific test set, which does not present either drop off time, nor trip duration. The prediction will be performed through a supervised learning regressor to be defined given the data structure.

1.0.2 3. Datasets and Inputs

The challenge provides two [data sets](#). * train.csv - the training set, which contains 1.458.644 observations * test.csv - the testing set, which contains 625.134 observations

The training set contains eleven features:

- `id` - a unique identifier for each trip
- `vendor_id` - a code indicating the provider associated with the trip record, *categorical*
- `pickup_datetime` - date and time when the meter was engaged, *categorical*
- `dropoff_datetime` - date and time when the meter was disengaged
- `passenger_count` - the number of passengers in the vehicle (driver entered value)
- `pickup_longitude` - the longitude where the meter was engaged
- `pickup_latitude` - the latitude where the meter was engaged
- `dropoff_longitude` - the longitude where the meter was disengaged
- `dropoff_latitude` - the latitude where the meter was disengaged
- `store_and_fwd_flag` - This flag indicates whether the trip record was held in vehicle memory before sending to the vendor because the vehicle did not have a connection to the server
 - Y=store and forward; N=not a store and forward trip
- `trip_duration` - duration of the trip in seconds

The testing set doesn't have neither the `trip_duration`, target variable, nor the `dropoff_datetime`. So, the model should be trained on the training set, and predict the `trip_duration` for the testing set observations. The predictions should be [submitted to the Kaggle's challenge](#) for calculating the final score.

The variables `pickup_latitude`, `pickup_longitude`, `dropoff_latitude` and `dropoff_longitude` will be used to calculate the trip distance between the pickup and dropoff locations. Also, they are important and will be analysed because where in NY the each trip starts and ends influences greatly the velocity in which the taxi can travel.

The variables `pickup_datetime` are important, and also will be analysed profoundly, because the moment on time in which the trips occur also influences greatly the the taxi velocity.

The remaining variables will be analysed in order to conclude about their influence on the target variable.

The data sets provided by the Kaggle challenge were made available in [Big Query on Google Cloud Platform](#), originally published by [NYC Taxi and Limousine Commission \(TLC\)](#).

2 Solution Statement

The solution will be provided by a regression method. The optimization will be performed after a preprocessing of the data and through minimization of the root mean square logarithm error (RMSLE), cross validation for detecting over and/or underfitting, and grid-search for fine tuning hyperparameters.

3 Benchmark Model and Result

The benchmark model will be the raw linear regression as it is well known and simple.

The benchmark result will be the [winning](#) score of the Kaggle's challenge when the winning team submitted the provided testing set. The challenge's metric is the RMSLE, and the winning score is 0.28976.

4 Evaluation Metrics

The evaluation metric is the one from the Kaggle's challenge, [RMSLE](#), in order to make the score comparable to the benchmark result. The metric is quite similar to the [root mean squared error](#), however, considers the squared difference of the logarithm values of prediction (\hat{y}) and true values(y).

$$\text{RMSLE}(y, \hat{y}) = \sqrt{\frac{1}{N} \sum_k \log(1 + y_k) - \log(1 + \hat{y}_k)}$$

5 Project Design

The project will be segregated in four major parts: data preprocessing, model selection, benchmark comparison and conclusion.

5.1 Data Preprocessing

Data preprocessing will be segregated into two sections: Generals and Specifics.

5.1.1 Generals

In this section, data quality will be considered and performed.

- Are there any missing values?
- The features are formatted on a convenient way?
- Is there any unnecessary information?

Along with that, from the features `pickup_latitude`, `pickup_longitude`, `dropoff_latitude` and `dropoff_longitude`, the feature `distance` will be calculated.

5.1.2 Specifics

In this section, the data will be deeply explored with the objective of understanding the features influences, finding transformations and detecting and dropping outliers in order to make the training and prediction better.

Continuous data First step will look for skewness of the continuous data and propose transformations that will be Logarithmic or Box-Cox and, finally. A probability plot against a normal distribution will help the decision in which transformation to perform.

Second step, take a good look at the data scatter plotting continuous features against each other and looking at their distributions to search for outliers, that will be studied further to decide whether they should be dropped or not.

Finally, the data will be normalized to the interval $[0, 1]$.

Categorical data The influence of categorical data will be studied separately. First, `datetime` related will be segregated into `month`, `day`, `hour`, `day of the week` and binary variables that flag if it's a holiday or not, in order to search for periodical/special behaviour. Second, the other categorical features such as `passenger_count` and `vendor_id`. They will be transform to dummy variables through *one-hot-encoding*. For these analysis, the velocity, distance over `trip_duration`, will be taken into account as the trip duration may depend on traffic, that depends greatly on date and time.

5.2 Model Selection

This parte will train the models and consider the RMSLE scoring for choosing the best one. As the testing set provided doesn't have the target feature, `trip_duration`, the training set will be splitted into training and testing set for the model selection.

5.2.1 Initial Model Evaluation

For the initial model evaluation, a set of different regressors will be trained in different data set sizes, using cross validation, in order to verify underfitting and overfitting between regressors. No tuning of hyperparameters will be performed.

As there will be a lot of dummy binary variables after the *one-hot-encoding* transformation, ensemble methods will be a good option for the final model. In this step, will be trained:

- Simple Linear Regression (benchmark model)

- Decision Trees
- Rain Forest
- Bagging
- Gradient Boosting
- Extreme Boosting

5.2.2 Model Tunning

The best two models from the previous step will have their hyperparameters tuned through gridsearch. The results of the best set of parameters for each model will be compared and the final model decided.

5.2.3 Final Validation

Once chosen the best model, it must pass through a set of final validations in order to check its robustness.

5.3 Benchmark Comparison

Once the model has been selected and trained on the whole provided training set, the predictions over the provided testint set will be submitted to the private [leaderboard](#) on Kaggle. The result will be compared with the bechmark model, Linear Regression, and the benchmark result, which is 0.28976.

5.4 Conclusion

The last section it is the conclusion with a reflexion about the data, the model and a proposal possible improvements.