

# Advanced Machine Learning Coursework: Kaggle Competition NYC Taxis

The Unicorn Hunters  
(Christian Frantzen, Guilherme Barreiro Vieira,  
Mortimer Sotom and Théo Verhelst )

May 8, 2018

## Abstract

This is the abstract

## 1 Introduction

## 2 Data Analysis

The first step of a machine learning project is to explore and analyse the data, in order to better understand the problem. Our dataset is composed of 11 features: a unique identifier for each row, the identifier of the taxi company, the pickup time, the pickup and drop-off locations, the number of passengers, and a boolean flag indicated if the trip data has been stored on-board or directly sent to the data server.

Figures XX to XX show the distribution of some of these features. We did not show the pickup month, minute and seconds, as their graphs are uniformly distributed and therefore not visually informative. Outliers have been removed in order to properly display the pickup and drop-off location distribution, as well as the trip duration. These outliers will be discussed in the next section. We can see in figure 1 and 2 that the location seems to be roughly normally distributed, and the logarithm of the trip duration also appears to be normally distributed in figure 7. The smaller bumps outside of the main bell in the location curves correspond to trips to the city airport.

It is also important to make sure that the training and test sets are independently and identically distributed. For this, we compared the above distributions with those from the test set. If the distribution from the training and testing set significantly overlap, then we can consider that the I.I.D. assumption is verified. As shown for example in figure and , it is indeed the case.

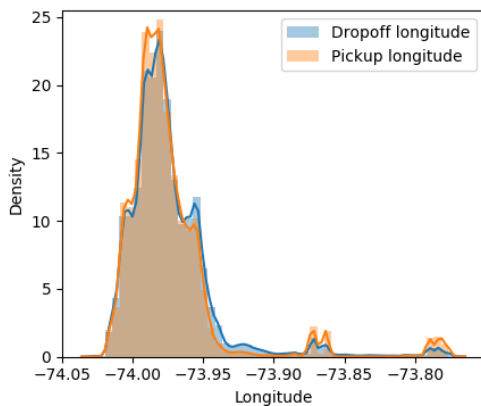


Figure 1: Distribution of the longitude.

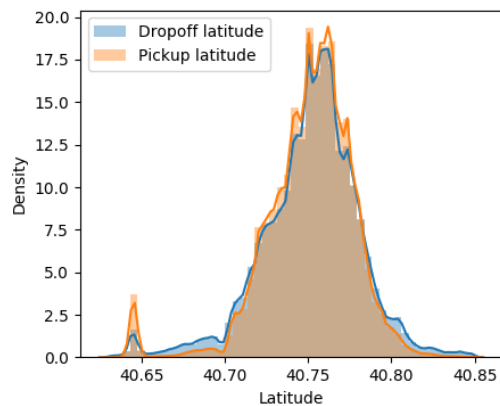


Figure 2: Distribution of the latitude.

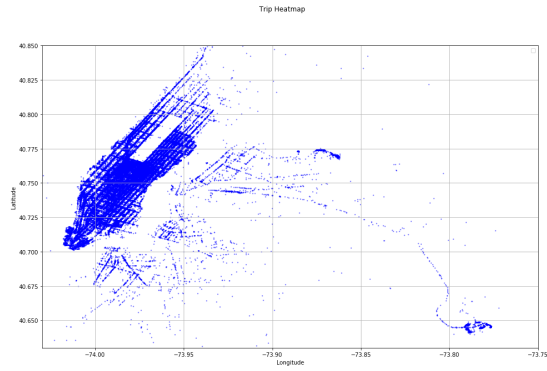


Figure 3: Trip locations as distributed on a map.

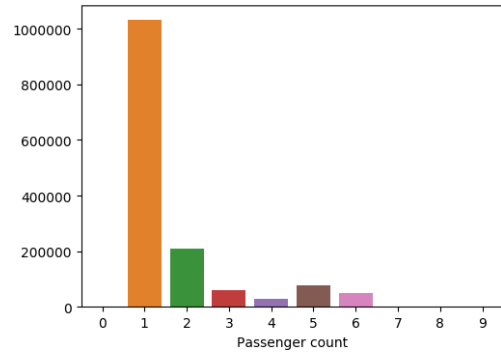


Figure 4: Distribution of the number of passengers.

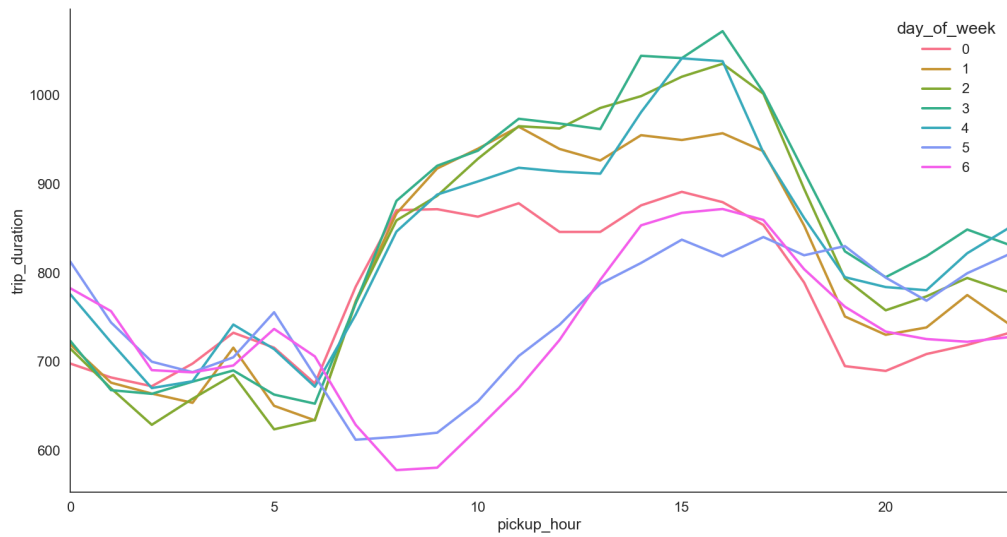


Figure 5: Distribution of the pickup time, for different days of the week.

### 3 Preprocessing

### 4 Feature Selection

### 5 Methods for Model Selection

### 6 Results

### 7 Conclusion

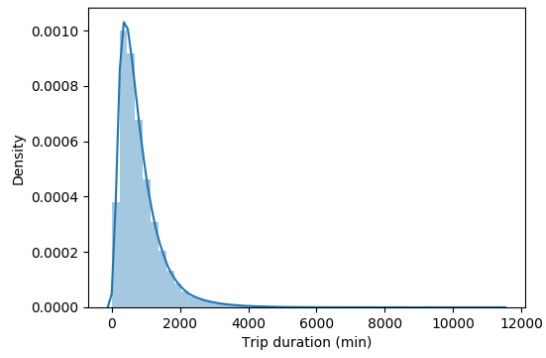


Figure 6: Distribution of the trip duration, in seconds.

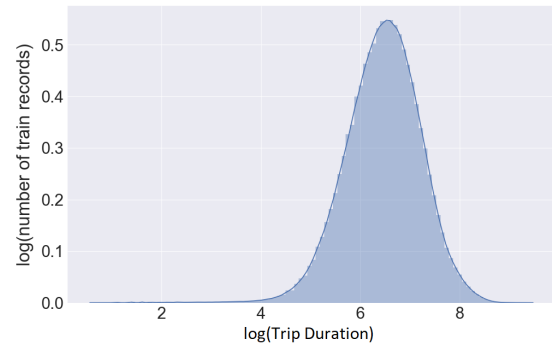


Figure 7: Distribution of the logarithm of the trip duration, in seconds.

## References

## A Additional Figures