# MVP

# TLC TRIP RECORD DATA

Done by:

**Arwa Alolyani    Maitha Alqahtani**

Submit to: Mr.Ali El-kassas

In our study, we selected a sample of basic data containing more than 6 million rows and created a million-row sample.
Our plans are to divide the data into three sections: training, verification, and testing.
Testing and Verification contains 200,000 rows and 18 columns, while training contains 600,000 rows and 18 columns.
This data is lacking coordinates for each location, so the NYC Taxi Zones (GeoJSON) file was downloaded from NYC OpenData.
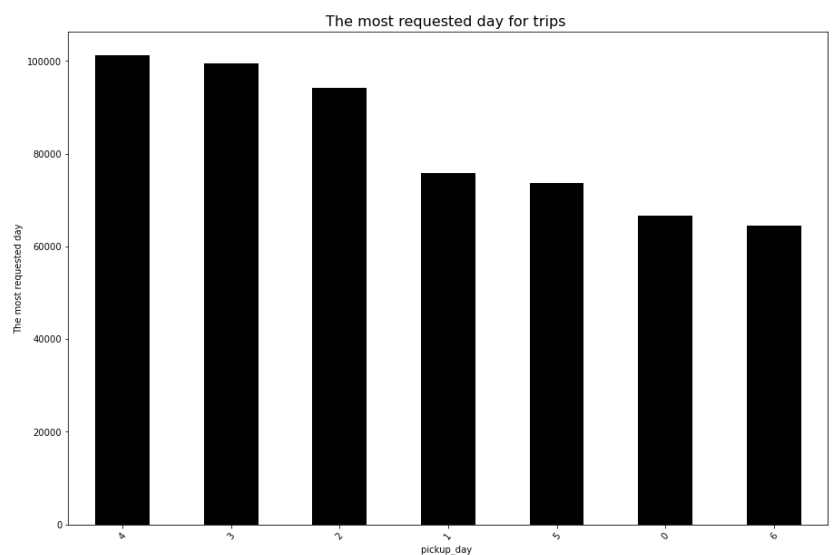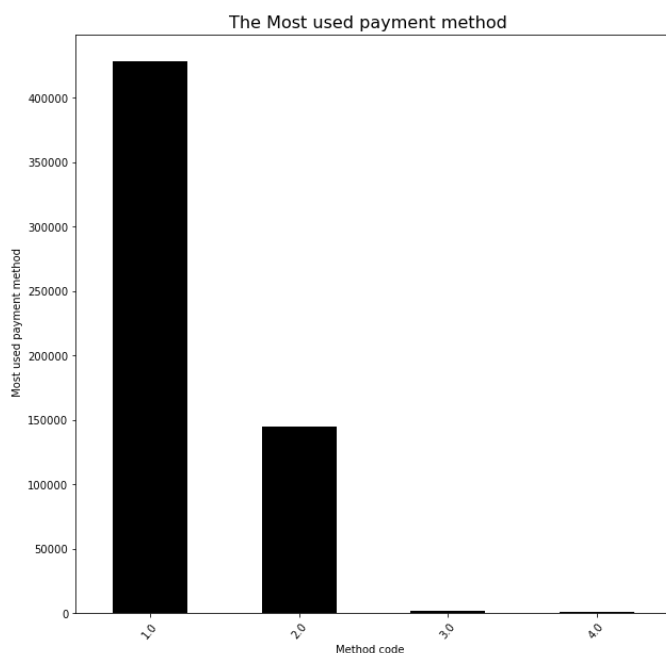
## process:

To begin with, we cleaned the data by removing outliers and nulls, converting datetime from object type to datetime type, adding new columns for day and hour, and then extracted latitude and longitude coordinates by geometric column transformation from NYC Taxi Zones database.
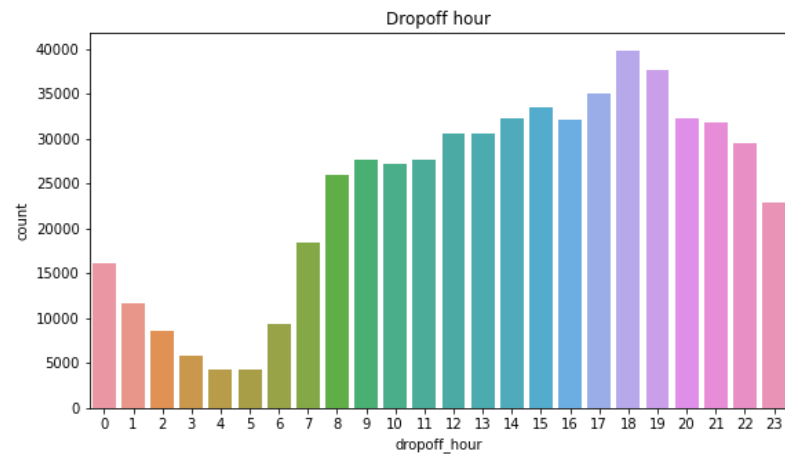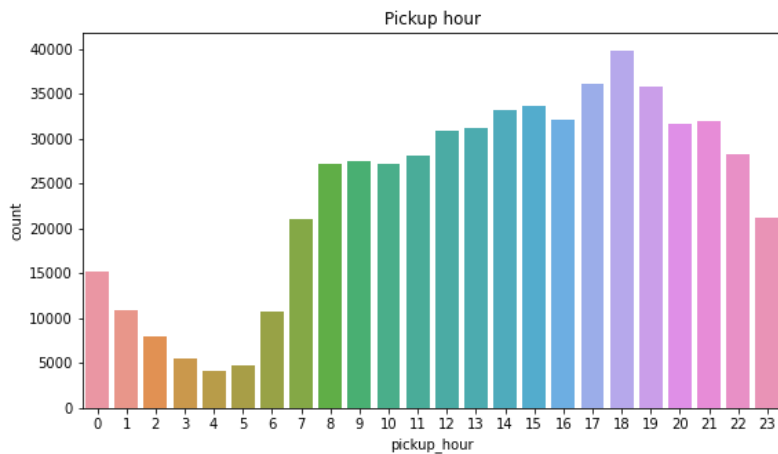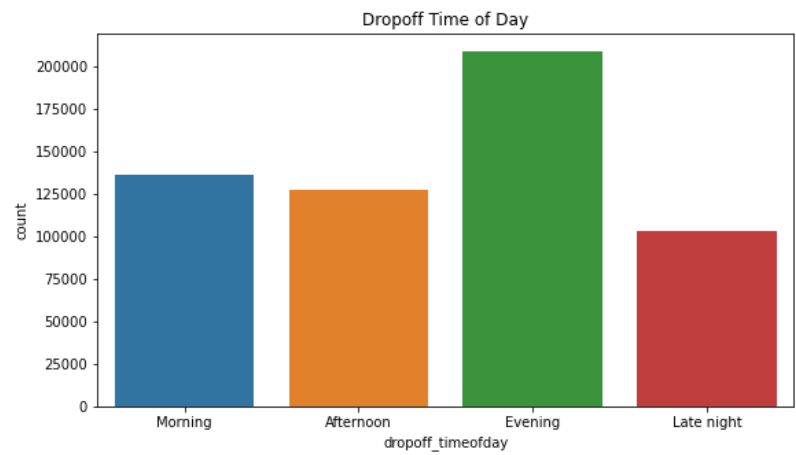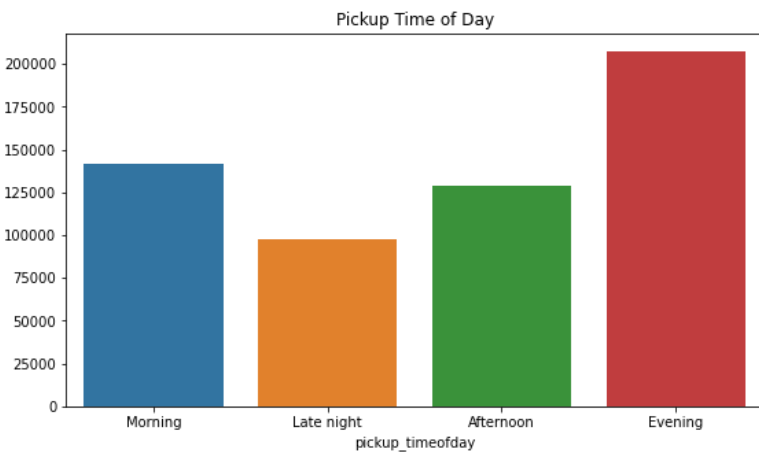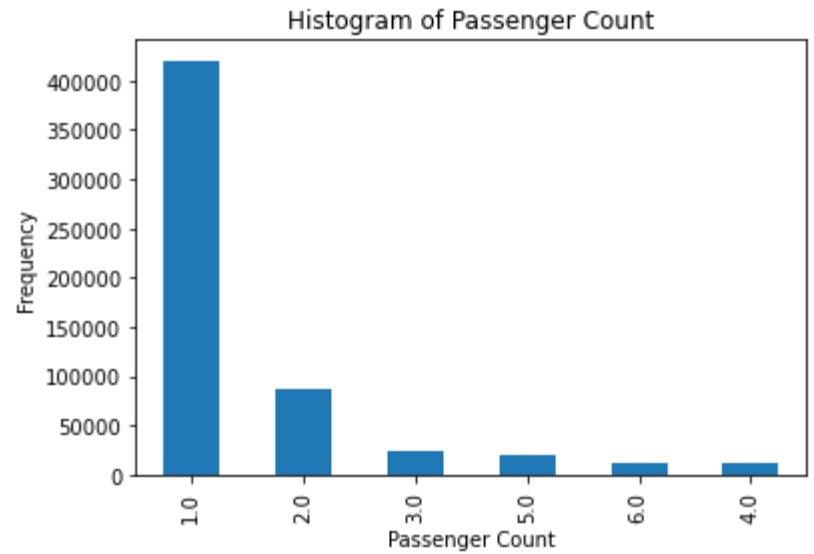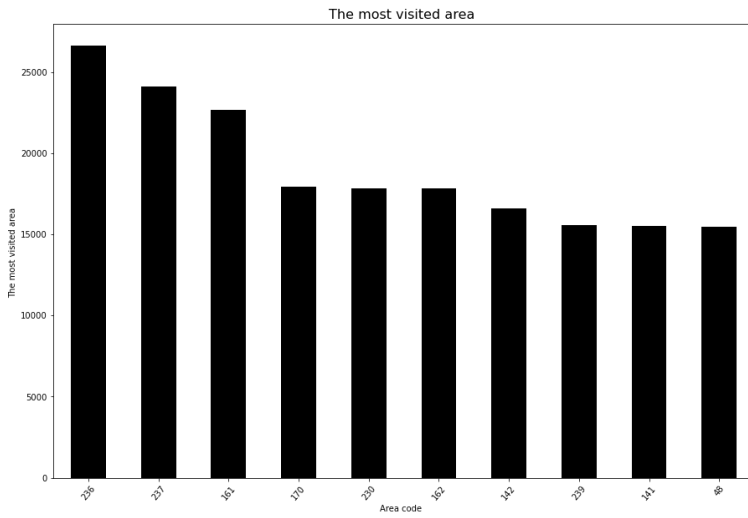The second stage of data extraction analysis:

The following data was extracted
1-    Preparing passengers for each flight
2-    Total price for trips
3-    the most visited area
4-    the most requested day for trips
5-    the number of trips per hour
6-    the most payment method used
7-    A color caller for the data, which shows how closely they are related to each other
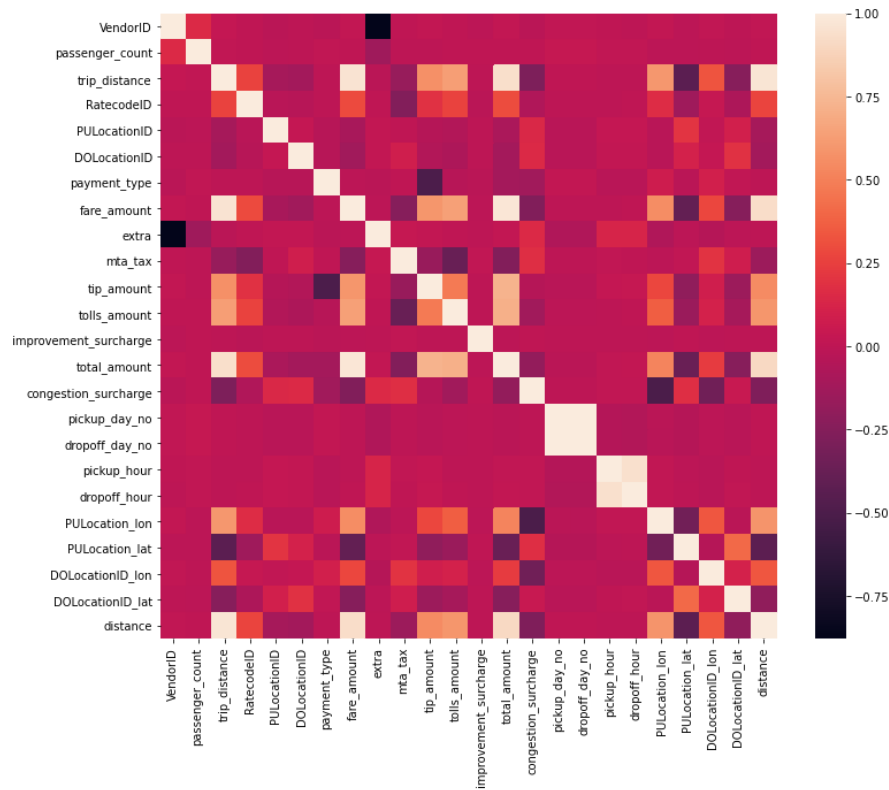


The Most used payment method



The most requested day for trips

**The most visited area**

**Histogram of Passenger Count**

**Pickup Time of Day**

**Dropoff Time of Day**

**Pickup hour**

**Dropoff hour**

The third stage is the application of machine learning to data, and in this stage there are several interconnected sequential steps
1- Apply the model to the training data and print the ratio
2- Apply the R2 and print the ratio
3- Apply the model to the verification data and print the ratio
4- Apply the R2 and print the ratio

```
In [251]: lr_model.score(X_train,y_train)

Out[251]: 0.8645609236485619

          xx,yy >> sample_val

In [253]: xx=sample_val[['dropoff_day_no','pickup_hour'

In [254]: X_val,y_val = xx, sample_val['fare_amount']

In [255]: lr_model.score(X_val,y_val)

Out[255]: 0.8560649544780876
```