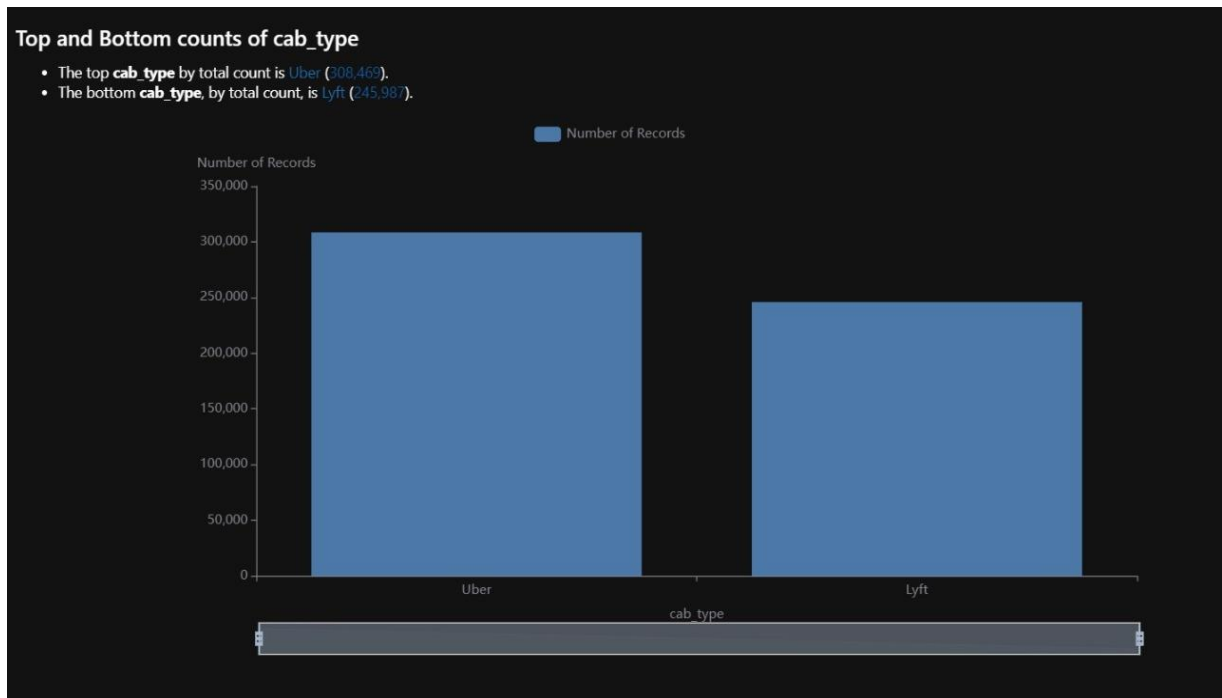# TAXI PRICE PRDICTION

Machine Learning Project

MAY 28, 2022

# Phase 1

- **Pre-processing:**
  - **"taxi_rides":**


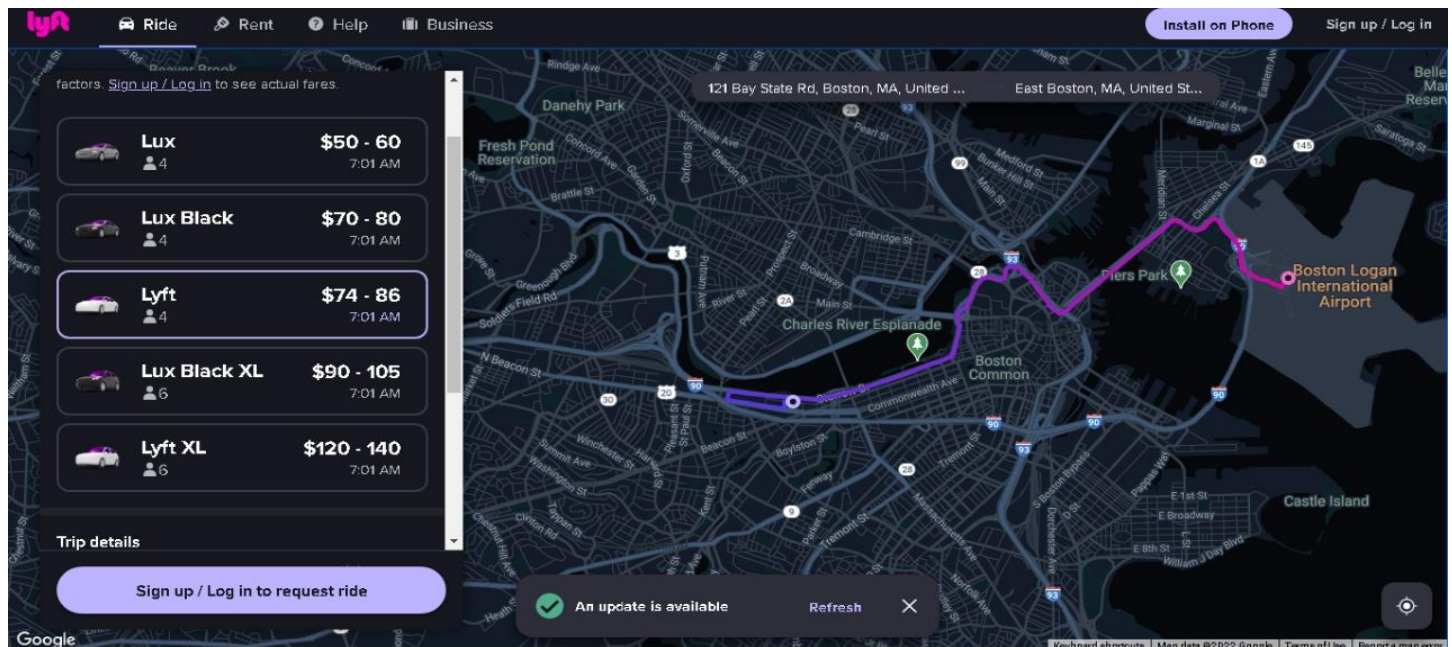
  - "cab_type" is object, it's either "uber" or "lyft", so one hot
  - "time_stamp" is float and it indicates time so it's converted to datetime.
  - "price" contains null values in taxi, so it's calculated by getting the price of the trips that are done by uberX which has the same source and destination.
  - "time_stamp" was replaced by "date" and "hour" and was dropped later.
  - "id" values are unique, so it doesn't help in prediction of the "price". It will be dropped later after merge step.
  - "lyft" has 6 product IDs, uber doesn't have any. Count of "name" is identical to count of "product_id", so a map was made to map "product_id" to "name.
  - "product_id" was dropped as "name" can replace it.

- For "lyft" cabs, names were compared to see how they affect the price, we got the description from official company website:https://help.lyft.com/hc/ru/articles/115012927427-Lyft-ride-modes-overview

  and the order is:
  - Shared: Share a car with riders headed in the same direction at a discounted price.
  - Lyft: Standard Lyft car for up to 3* riders
  - Lyft xl: SUV for up to 5* riders
  - Lux: Luxury car for up to 3* riders
  - Lux black: Premium black car service with leather seats for up to 3* riders
  - Lux black xl: Premium black SUV with leather seats for up to 5* riders

- Ordinal encoding was done based on this order on "lyft_types"

o Price differences between the Uber ride types:
The cost of your Uber ride is largely determined by the Uber service that you select.

The costs of the different services, from least expensive to most expensive: Uber Pool, Uber X, Uber Comfort, Uber XL, Uber Select, Uber Black, Uber SUV.

An example Uber from The Grove to the Century City Mall in Los Angeles
Note: This ride is 4.6 miles and 18 minutes. The price may change due to traffic, time of day, or discounts
Uber Ride Type   Est. Ride Cost (4.6 miles, 18 minutes)
  - Pool  $9-$11
  - X        $9-$12
  - Comfort    $12-$16
  - XL     $15-$20
  - Select        $24-$30
  - Black $30-$40
  - Black SUV  $42-$52

o For "uber" cabs, names were compared to see how they affect the price, the order is:

| ↓ Choose a ride | ↓ Choose a ride | ↓ Choose a ride |
|---|---|---|
| **Economy** | **Premium** | **Premium** |
| UberX Priority ▲3 $13.01<br>4:59pm 2x pts<br>Faster pickup | Black Hourly ▲3 $110.75<br>Luxury rides by the hour 2 hrs/30 miles<br>with professional drive... | **More** |
| UberX ▲3 $8.99<br>5:00pm 2x pts<br>Affordable rides, all to yourself | Black SUV Hourly ▲5 $140.75<br>4:57pm 2 hrs/30 miles<br>Luxury hourly rides for 5 with<br>professional drivers | Español ▲3 $8.99<br>5:05pm 2x pts<br>Affordable rides with Spanish-<br>speaking drivers |
| Comfort ▲3 $10.85<br>5:02pm 2x pts<br>Newer cars with extra legroom | Black ▲3 $24.86<br>4:59pm 3x pts<br>Luxury rides with professional drivers | Select ▲4 $19.25<br>4:59pm 2x pts<br>Premium rides in high-end cars |
| UberXL ▲5 $11.65<br>5:02pm 2x pts<br>Affordable rides for groups up to 5 | Black SUV ▲5 $33.10<br>4:58pm 3x pts<br>Luxury rides for 5 with professional<br>drivers | Assist ▲3 $8.99<br>5:06pm 2x pts<br>Special assistance from certified<br>drivers |
| Uber Green ▲3 $9.94<br>4:59pm 2x pts<br>Eco-Friendly | Lux ▲4 $40.55<br>5:00pm 3x pts<br>Premium rides in luxury cars | WAV ▲4 $8.99<br>5:08pm 2x pts<br>Wheelchair-accessible rides |

- UberPool
- Taxi
- UberX
- UberXl
- WAV
- Black
- Black SUV

- o Ordinal encoding was done based on this order on "uber_types"
- o "cab_type" for both "lyft" and "uber" were dropped as their values were encoded in "lyft_type" and "uber_type"
- o In "distance" there were outliers, so they were removed it .



**Average distance by source X destination**
- Overall, highest **Average distance** (5.17) was found for Financial District : Boston University. This is 136.02% higher than the overall **Average distance**(2.19).
- Overall, lowest **Average distance** (0.42) was found for Financial District : South Station.
- White cells in the heatmap denote no occurrence of data

- o **In "distance": we got all street name from {data/Boston/street_name website } using web scraping then we got the latitude and longitude of every street name using google.geoglocation api , so we can calculate distance easily between any source and destination.**

- o **"weather":**
- o "time_stamp" is float, so it's converted to datetime.
- o For "rain" 85% of the data is null, K-Nearest Neighbors imputation method was used. Normalizing data was applied in order not to generate biased replacement for the missing values.
- o Regression techniques:
  - Polynomial regression
  - Multiple regression
- o For Multiple regression:
  - Mean Square Error 16.33271939415199
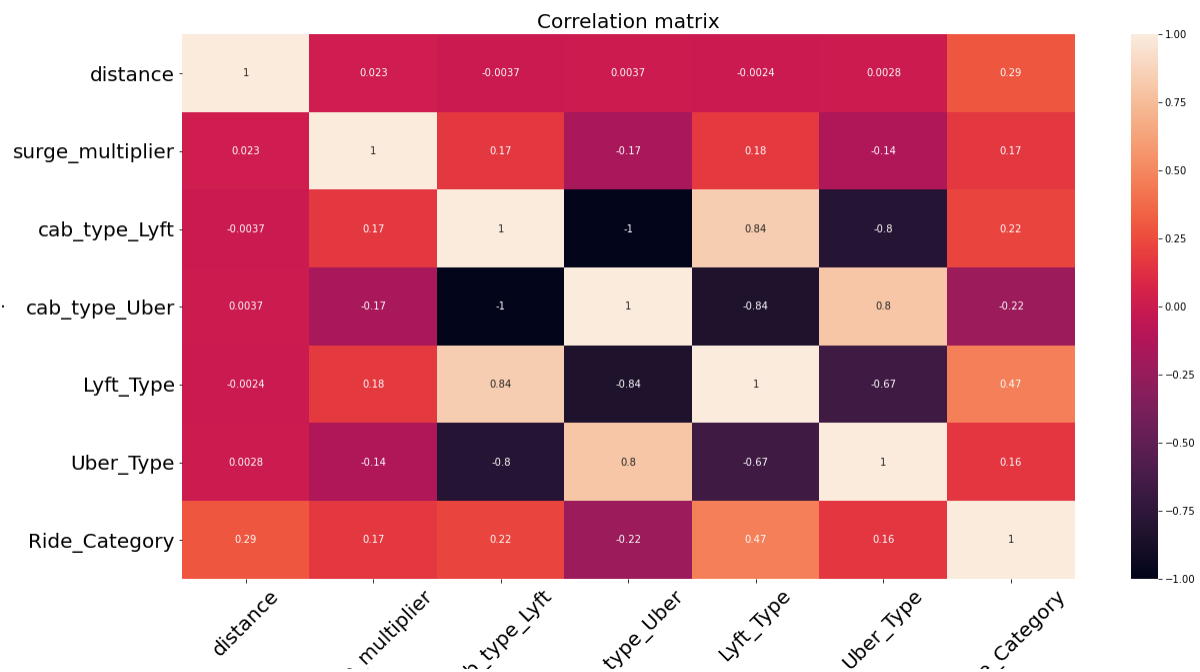  - r2_score: 75.16022240596294 %
  - Training time: 0.10199832916259766 seconds
  - For Polynomial regression: Degree : 6
  - Mean Square Error 3.203520780009517
  - r2_score: 95.94640692535074 %
  - Training time 0.12099742889404297 seconds
- o Features used for regression are:
  - Distance
  - surge_multiplier
  - Lyft_Type
  - Uber_Type
  - cab_type_Uber

- Training set size is 60% and testing set size is 40%

# How data affect each other



- **How data affect each other after feature selection**



Correlation matrix

- Data visualization for "distance" outliers:



- Conclusion:
  - o In this phase, we applied preprocessing for all features and feature selection and we have concluded that weather is not a good indicator for price and the product id (name) & distance was the main indicators for price. But for surge it was very effective in calculating the price for Lyft cabs but not for Uber cabs .

# Phase 2

- **Pre-processing:**



**We do the same preprocessing as we did in phase 1 and we noticed that Ridge Category in order of:**

       **1: unknown**

       **2: cheap**

       **3: moderate**

       **4: expensive**

       **5: very expensive**

# Feature Selection:



Correlation matrix

## Summarization:

### 1: classification accuracy



Classification Accuracy Summarize

## 2: Total training time



Total Training Time Summarize

## 3: Total test time



Total Testing Time Summarize

## 4: Mean Square Error

Mean Square Error Summarize



Hyperparameter tuning:

## RandomForestClassifier's hyper parameters.

{'bootstrap': True, 'ccp_alpha': 0.0, 'class_weight': None, 'criterion': 'gini', 'max_depth': None, 'max_features': 4, 'max_leaf_nodes': None, 'max_samples': None, 'min_impurity_decrease': 0.0, 'min_impurity_split': None, 'min_samples_leaf': 1, 'min_samples_split': 2, 'min_weight_fraction_leaf': 0.0, 'n_estimators': 3, 'n_jobs': None, 'oob_score': False, 'random_state': None, 'verbose': 0, 'warm_start': False}

# Try to Explain in detail how hyperparameter tuning affected RandomForestClassifier models'performance.

```python
from sklearn.model_selection import GridSearchCV

# Create the parameter grid based on the results of random search
param_grid = {
    'bootstrap': [True],
    'max_depth': [10,15,20],
    'max_features': [3,4],
    'min_samples_leaf': [3, 4, 5,6,7,8],
    'min_samples_split': [3,4,5,6],
    'n_estimators': [30,50]
}
# Create a based model
rf = RandomForestClassifier()
# Instantiate the grid search model
grid_search = GridSearchCV(estimator = rf, param_grid = param_grid, n_jobs = -1, verbose = 2)

# Fit the grid search to the data
grid_search.fit(x_train,y_train)
```

## Android Application

Tools:

    1: Heroku Cloud: we upload our model on Heroku cloud

So, we can use model anytime.

    2: Flask Framework: to connect model with local server    as framework.

    3: Android Studio: we use it to deploy our application.

    4: Google Places API: we use it to get any street name in all the world.

    5: Google Map API: we use it to get maps from google

    6: Google Distance API: to calculate distance between two points in map

Comparing Between our Application and the real application to predict the price between same distance

Our uber app

real uber app



| 7:20 | | | | | | | | | 0.11 KB/S | | | 94 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

TaxiPrediction

Uber

## Prediction

| Uber Pool | 8.18$-12.18$ |
|---|---|
| Shared rides are back | |

| Taxi | 9.05$-13.05$ |
|---|---|
| Local taxis at the tap of a button | |

| UberX | 12.16$-16.16$ |
|---|---|
| Affordable,everyday rides | |

| UberXL | 15.95$-19.95$ |
|---|---|
| Affordable rides for groups up to 6 | |

| WAV | 10.35$-14.35$ |
|---|---|
| Rides in wheelchair-accessible | |

| Uber Black | 22.62$-26.62$ |
|---|---|
| Premium rides in luxury cars | |

| BlackSUV | 32.62$-36.62$ |
|---|---|
| Premium rides for 6 in luxury SUVs | |

Your options

| Taxi | $16.16 |
|---|---|
| Connect | $33.35 |
| UberX | $35.31 |
| Uber Green | $35.31 |
| Comfort | $36.80 |
| Black | $37.22 |
| UberXL | $39.71 |
| Uber Pet | $41.31 |
| Black SUV | $51.20 |

## Our Lyft app

TaxiPrediction

**Prediction**

| | | |
|---|---|---|
| Shared | | Unavailable |
| Share a car with riders | | |
| Lyft | | 13.11$-15.11$ |
| Standard Lyft car for up to 3 riders | | |
| Lyft XL | | 24.13$-26.13$ |
| SUV for up to 5* riders | | |
| LUX | | 27.04$-29.04$ |
| Luxury car for up to 3 riders | | |
| Lux Black | | 30.48$-32.48$ |
| Premium black car service | | |
| Lux Black XL | | 43.92$-45.92$ |
| Premium black SUV | | |

## Real app

**Fare estimate**

Sample fares are estimates only and do not reflect variations due to discounts, traffic delays or other factors. Sign up / Log in to see actual fares.

| | | |
|---|---|---|
| Lux 4 | $21 - 24 | 7:12 AM |
| Lyft 4 | $21 - 24 | 7:12 AM |
| Lux Black 4 | $28 - 32 | 7:12 AM |
| Lyft XL 6 | $32 - 36 | 7:12 AM |
| Lux Black XL 6 | $45 - 50 | 7:12 AM |