



FLIGHT PRICE PREDICTION PROJECT

Submitted by:

SHANTANU

ACKNOWLEDGMENT

I want to express my great appreciation to my mentor for her valuable and constructive suggestions during the planning and development of this research work. Her willingness to give his time so generously has been very much appreciated. I would also like to thank the staff of the Data Trained for helping me with the problems I faced during the research work. Articles from the "Medium" platform were beneficial during the whole process. It helped me clear my concepts.

INTRODUCTION

Anyone who has booked a flight ticket knows how unexpectedly the prices vary. The cheapest available ticket on a given flight gets more and less expensive over time. This usually happens as an attempt to maximize revenue based on - 1. Time of purchase patterns (making sure last-minute purchases are expensive) 2. Keeping the flight as full as they want it (raising prices on a flight which is filling up in order to reduce sales and hold back inventory for those expensive last-minute expensive purchases) So, you have to work on a project where you collect data of flight fares with other features and work to make a model to predict fares of flights.

Objective:

This project contains two-phase-

Data Collection Phase

We scrapped more than 1500 rows of data. In this we scrapped s the data of flights from different websites (yatra.com, skyscanner.com, official websites of airlines, etc).. Generally, these columns are airline name, date of journey, source, destination, route, departure time, arrival time, duration, total stops and the target variable price.

Model Building Phase

After collecting the data, we built a machine learning model. Before model building, we did data pre-processing steps.

Followed the complete life cycle of data science. Include all the steps like.

1. Data Cleaning
2. Exploratory Data Analysis
3. Data Pre-processing
4. Model Building
5. Model Evaluation
6. Selecting the best model

Analytical Problem Framing

In the whole research process various mathematical, statistical and analytics modelling has been done. There has been reduction of the columns because few of them was not necessary for the problem solving like Id. And few of them was removed due to very less correlation with dependent variable. To fix the outliers we used z score method. After this also there was a lot of skewness in dataset so power transform has been used. To check the accuracy r^2 score was used also for cross validation `cross_val_score` is used.

DATA/ DATA PREPROCESSING:

- The dataset contains 1792 rows and 8 columns
- Fare is our dependent variable.
- We created new features from old ones.
- All columns were object data types we converted necessary ones into int and float.
- There are no null values in the dataset.
- Trimmed few columns

Hardware and Software Requirements and Tools Used

- HP 5- i5 8th generation, 8gb ram, NVidia mx130 integrated graphic,
- JupyterNotebook/Google chrome
- Libraries and packages used:
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")
from sklearn.preprocessing import LabelEncoder
from sklearn.feature_selection import VarianceThreshold
from sklearn.feature_selection import mutual_info_regression
from sklearn.feature_selection import SelectPercentile
from sklearn.preprocessing import StandardScaler
from statsmodels.stats.outliers_influence import variance_inflation_factor
from sklearn.preprocessing import power_transform
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
from sklearn.linear_model import LinearRegression
from sklearn.tree import DecisionTreeRegressor
from sklearn.neighbors import KNeighborsRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import ExtraTreesRegressor
from sklearn.svm import SVR
from sklearn.model_selection import cross_val_score
from sklearn.ensemble import BaggingRegressor
from sklearn.ensemble import AdaBoostRegressor
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.model_selection import GridSearchCV

the library used here is sklearn,numpy,matplotlib,pandas and seaborn. The matplotlib and seaborn library has been used to make charts to visualize and understand the problem, correlation, outliers and many other things, the pandas and NumPy library is used to handle dataset and perform various tasks. The seaborn library is used for model building and cross validation of the models.

Model/s Development and Evaluation

The approach to solve this problem was to get the domain knowledge to understand the data better. Which values can be the part of the data and which is not? After exploring the data, it is found that though the data has no missing value. It has extreme outliers and unrealistic value. We used Z-Score method to remove outliers. There was some skewness in the data, power transform method has been used so it dealt skewness. To check the accuracy, mean square error, mean absolute error, r2 score was used also for cross validation cross_val_score is used

Algorithm used for Training and testing:

```
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
from sklearn.linear_model import LinearRegression
from sklearn.tree import DecisionTreeRegressor
from sklearn.neighbors import KNeighborsRegressor
from sklearn.ensemble import RandomForestRegressor
from sklearn.ensemble import ExtraTreesRegressor
from sklearn.svm import SVR
from sklearn.model_selection import cross_val_score
from sklearn.ensemble import BaggingRegressor
from sklearn.ensemble import AdaBoostRegressor
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.model_selection import GridSearchCV
```

Performance of the model:

```
1 #splitting train test data
2 x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=.30, random_state=56)
3 x_train.shape,y_train.shape,x_test.shape,y_test.shape
```

((817, 32), (817,), (351, 32), (351,))

By the train test split method, 70 percent of the data has been taken for the model building while 30 percent of the data has been reserved for checking the model's performance.

```
1 #creating function
2 def model(name):
3     model=name()
4     model.fit(x_train,y_train)
5     predict=model.predict(x_test)
6     print("""mean squared error is:
7     """,mean_squared_error(y_test, predict))
8
9     print("The mean absolute error is: ", mean_absolute_error(y_test,predict))
10
11
12     print("""r2 score is:
13     """,r2_score(y_test,predict))
14
15
16     print("cross_val_score", cross_val_score(model,x,y,cv=5).mean())
```

The code above has been used to speed up the model training and its evaluation process. Here the function name model has been created which takes the name of model as argument.

#LinearRegression

```
mOdel( LinearRegression)
```

```
mean squared error is  
31é60é9.2223878595
```

```
The mean absolute error is: 1433.866794111)B72  
r2 score is'
```

```
6.25 1934025 13999397  
cross_val_score 0.016G5777862732844
```

#decisiontreeregressor

```
mOdel(DecisionTreeRegressor)
```

```
mean squared error is  
2363036.64Plé42B84
```

```
The mean absolute error is: 839.0975d48S36356  
r2 score is:
```

```
0.44691726018038  
cross_val_score -8.10927074293391207
```

#randomforestregressor

```
ma de l( RandomForestRegressor)
```

```
mean squared error is  
1783233.^960.^773 G
```

```
The mean absolute error is: 82G . 3791080769895  
r2 score is:
```

```
0.5826235105^21618  
cross_val_score 0.18439391680321432
```

#extratreesregressor

```
model(ExtraTreesRegressor)
```

```
mean squared error is:  
1836607.35 19434037
```

```
The mean absolute error is: 735.0998276676108  
r2 score is:
```

```
0.5668543757490115  
cross_val_score 0.267G53989G3^2B817
```

```
model(KNeighborsRegressor)
```

```
2B35y56.6e57507fi8  
The mean absolute error is: 991.4515580736545
```

```
0.5235191675749986  
<cross_val_score B.1064265Bb750B12^7
```

```
from sklearn.ensemble import BaggingRegressor  
from sklearn.ensemble import AdaBoostRegressor  
from sklearn.ensemble import GradientBoostingRegressor
```

```
model(BaggingRegressor)
```

```
mean squared error is:
```

```
The mean absolute error is: 855.^6B267BC54662  
r2 score is
```

```
0.5548223096347795  
cross_val_score 0.13421416Ba833171
```

```
model(AdaBoostRegressor)
```

```
The mean absolute error is: 1371.042d0BS522315  
r2 score is:
```

```
<cross_val_score B.1355557127445a
```

```
aC d' e tSook'. 'qg <q' "> so  
node1 (G' ad Tgp IB OOH I ngR ggK go à OK)
```

```
mean squared error is:  
1888658.1876762426  
The mean absolute error is: 974.2160123596385
```

```
0.516674SZ63C1f36C  
cross_val_score Bñ2778651229032538
```

```

#Doing oies! regression is the best model as Ohe RNSLE is inoixuin
#se l: i ng parameters fi-or hypei°porainetei° Sunn'i ng
parameter={
    "criterion": ["mse" "mae" $ ,
    'max_features': [ 'auto' 'sqrt' "log2" $ ,
    'min_samples_split': [2 5, GB, 15$ ,
    'min_samples_leaf': [ 1 2, 5, 16a }

#us Eng Grp dSeai°chCV fi-or Hypes° pai°aine l:er l:unn'ing
from sklearn.model_selection import GridSearchCV
qcs=G r ids ear chCV( RandomForestRegressor() parameter cv=5)

gcv.fit(x_train,y_train)

: GridSearchCV(cv=5, estimator=RandomForestRegressor()
    param_grid={ 'criterion': ['mse' 'mae'],
    'max_features': ['auto', 'sqrt', 'log2'],
    'min_samples_split': [1, 2, 5, 10],
    'min_samples_leaf': [1, 2, 5, 10] })

: #check'ing bet l: porainetei°s
gcv.best_params_

: { 'criterion': 'mae',
    'max_features': 'sqrt',
    'min_samples_leaf': 1,
    'min_samples_split': 2 }

model=RandomForestRegressor(criterion="mse",max_features="auto",min_samples_leaf=1,min_samples_split=10)
model.fit(x_train,y_train)
pred=model.predict(x_test)

print(" " "mean squared error is :
      " " " , mean_squared_error(y_test , pred) )

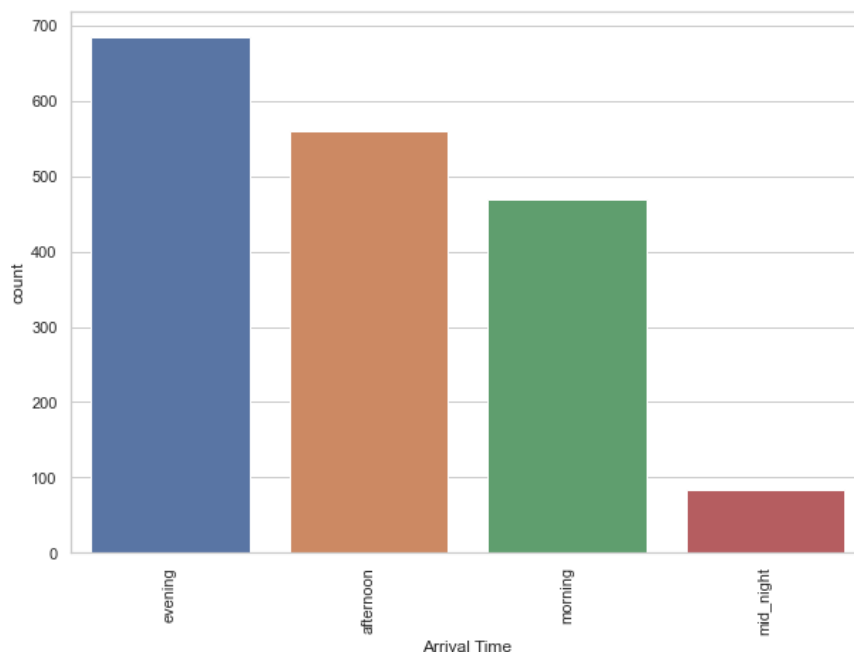
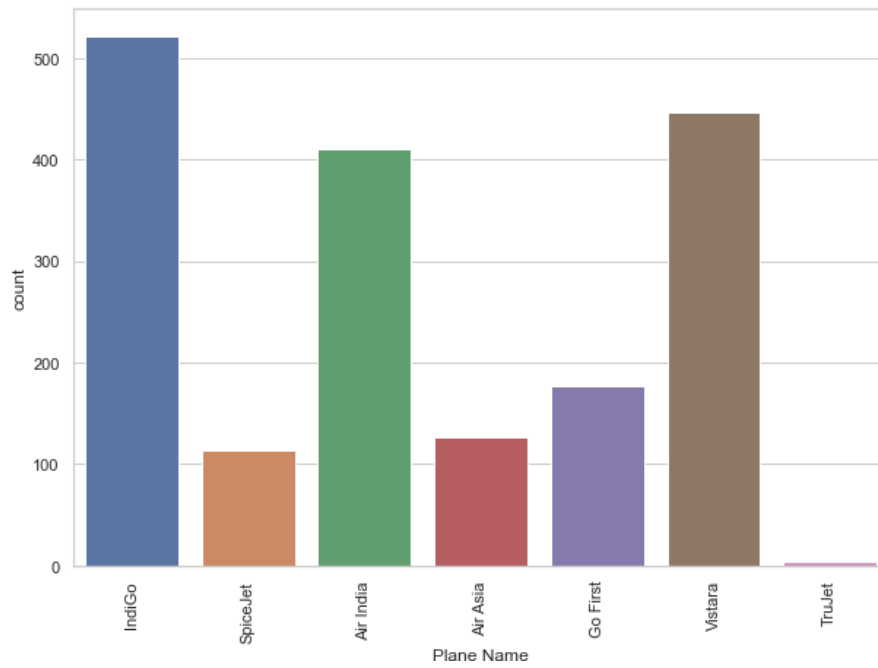
print("The mean absolute error is : " mean_absolute_error(y_test pred) )

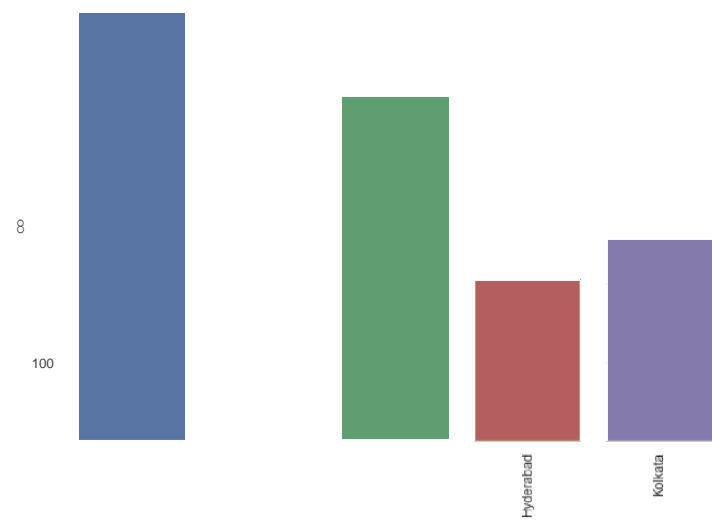
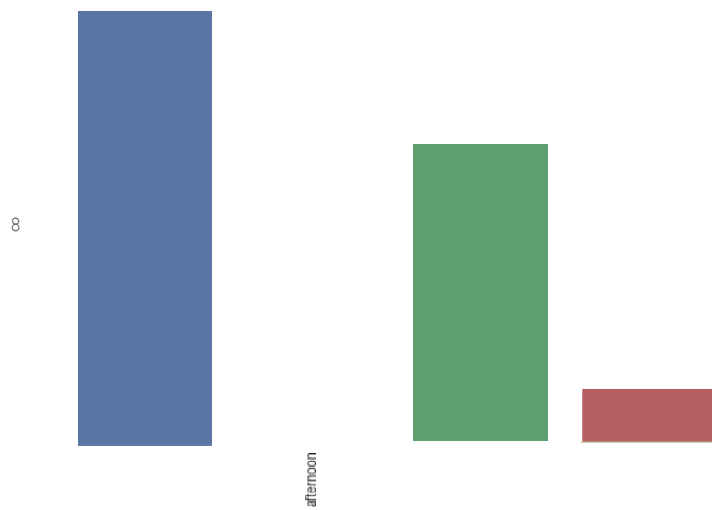
print(" r2 score is : " , r2_score(y_test pred) )

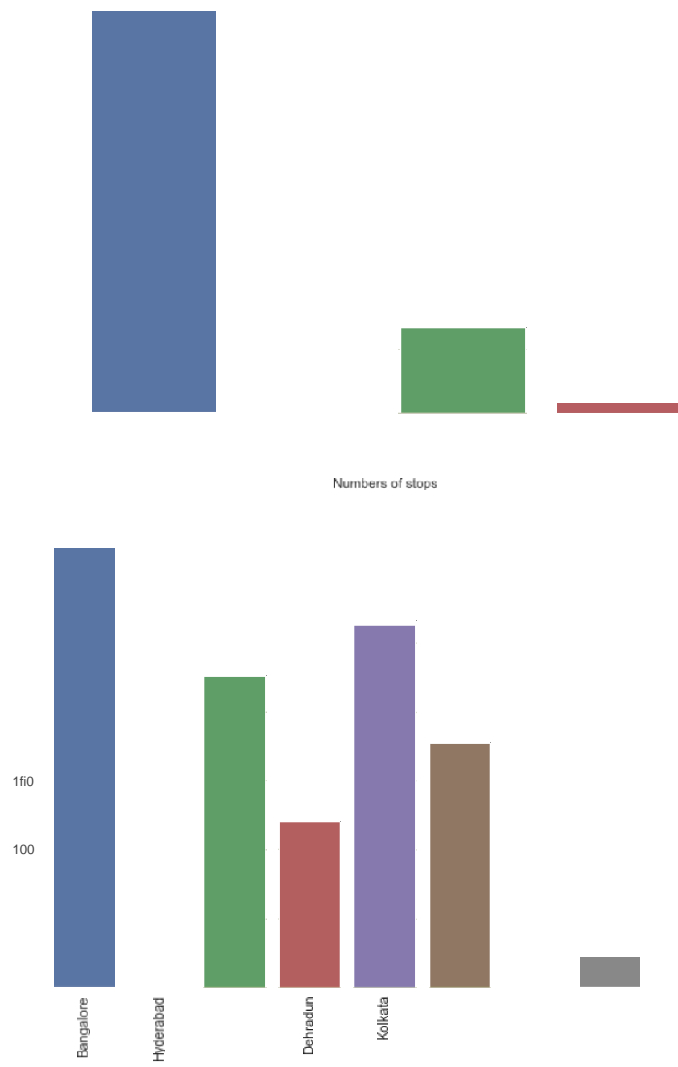
mean_squared_error is:
1W3T62. ST598T91S
The mean absolute error is : 8B1. 76751871519B5
rd score is: B . 6B12243131166367

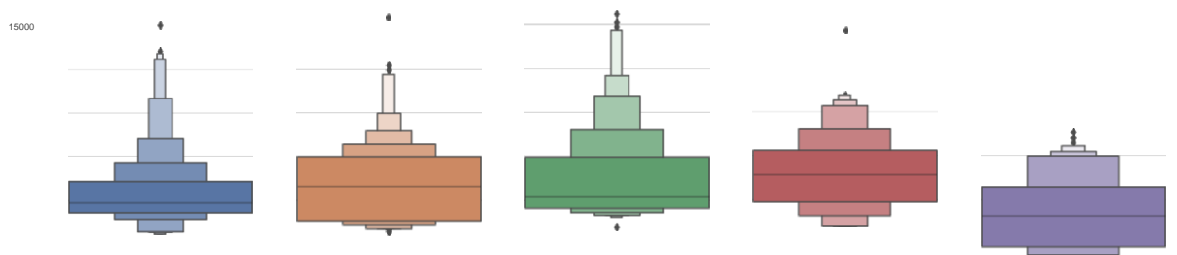
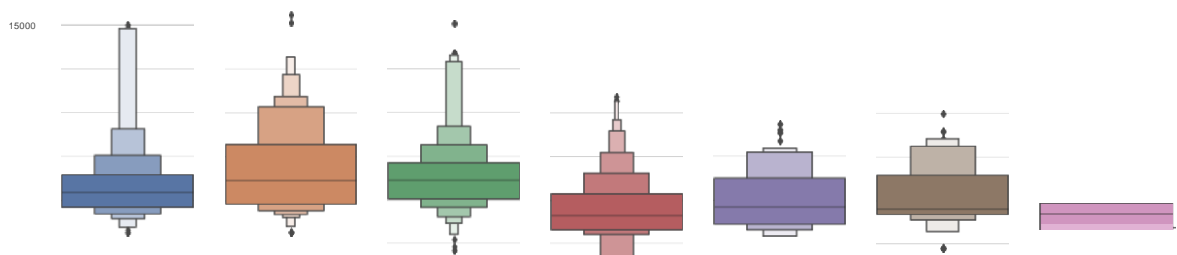
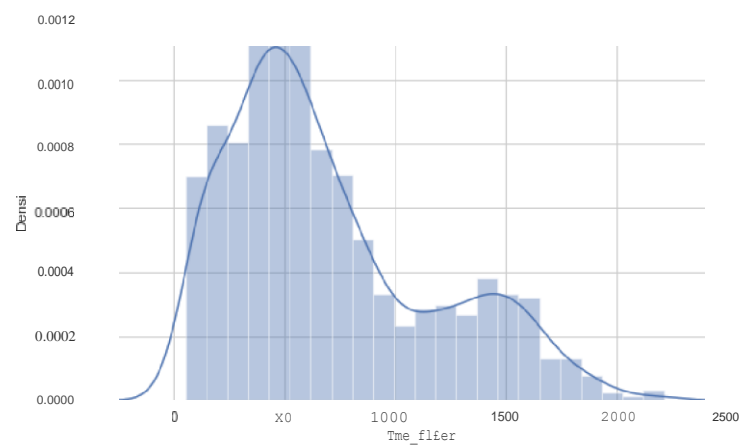
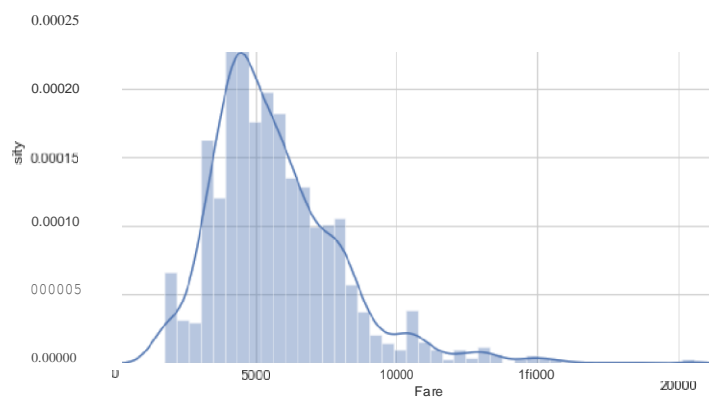
```

Visualization:

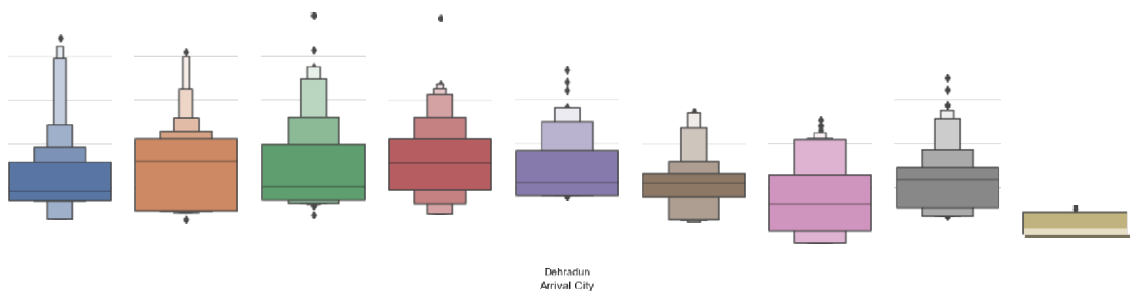




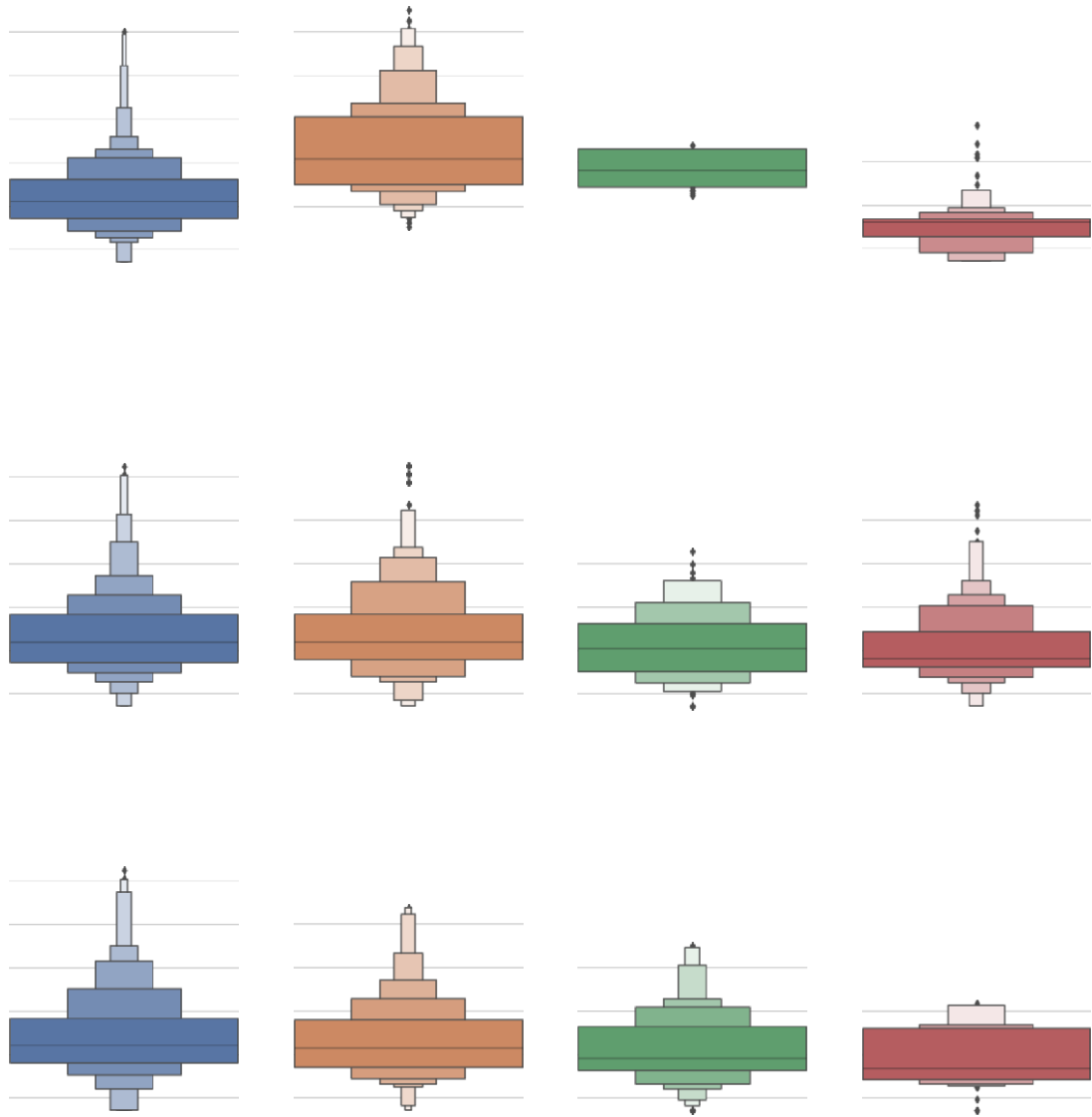


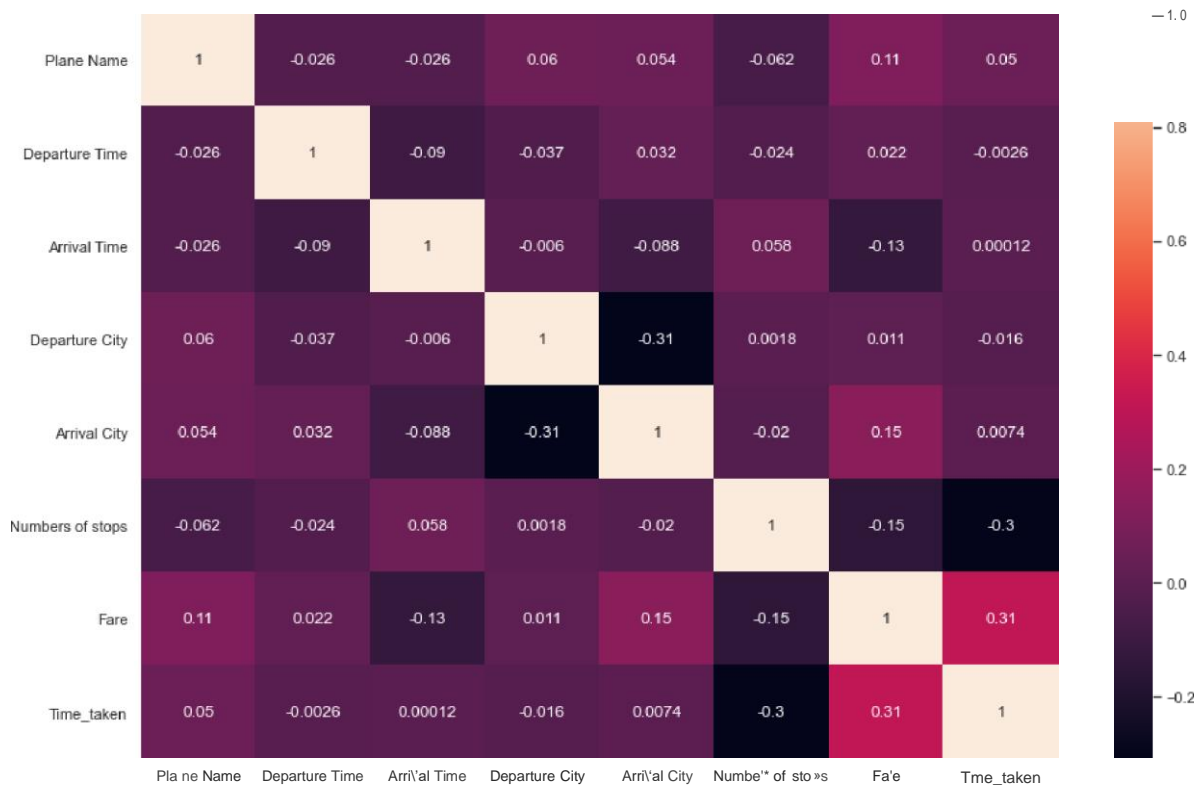
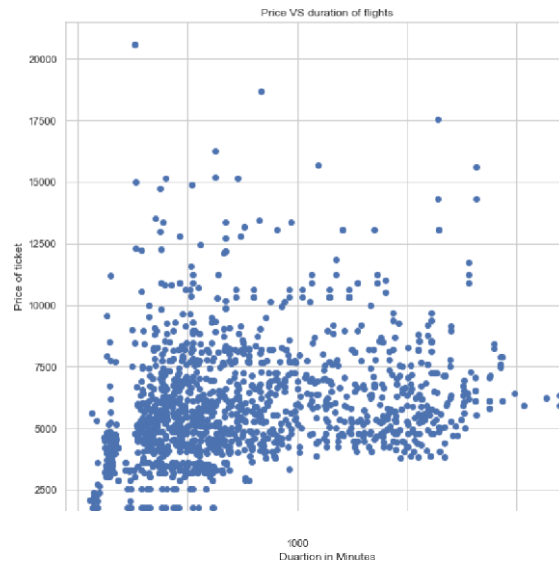


15000



15000





Observations:

#indigo provides the greatest number of services while Trujet provides least number of services.

#most flights depart in the morning

#most flight arrives in the evening

#from the collected data New Delhi is the has the most departed flight

#from the collected data New Delhi is the has the most arrived flight

#majority of flights have one stops

#the prices lie b/w 5k-7k. also there are few flights with higher price

#majority of the flights takes average 500-800 minute

#go fist has the highest fare and the lowest also

the flights that departs from Delhi has more fare than others

the flights that arrive in Bangalore has more fare than others

#the flight with 1 stop is more costly

#fights that departs in the morning costs more

##the flights that lands in afternoon has more fair

#fare shows a liner relationship with time

#fare is highly correlated with time taken

#there are some outliers in Time taken

CONCLUSION

This paper showed the model training process for the prediction of the fare Price. One of the objectives of the paper was to check the important variable for the prediction of the price and how these variables describe the price. Through model training and evaluating its performance. RandomForest proved to be as best model. As the difference between the r^2 score and cross validation score was minimum. This project has increased my understanding of the concept. During the research I came across various challenges and while solving them I learned a lot of new things. For example. How to plot different charts. For example, I learned how to plot subplot. I learned new libraries and how to use them. I explored various methods for feature selection. Also, I came to understand how can multicollinearity can cause problem during the model training. The limitation of the solution provided is that the data carried a lot of unrealistic values. Apart from that my laptop took too much time while running certain command where I lost a lot of precious time.