



Project report on
Flight Price Prediction

Submitted By
Shubham J. Chougule

ACKNOWLEDGMENT

It is my sensual gratification to present this report on FLIGHT PRICE PREDICTION project. Working on this project was a good experience that that has given me a very informative knowledge.

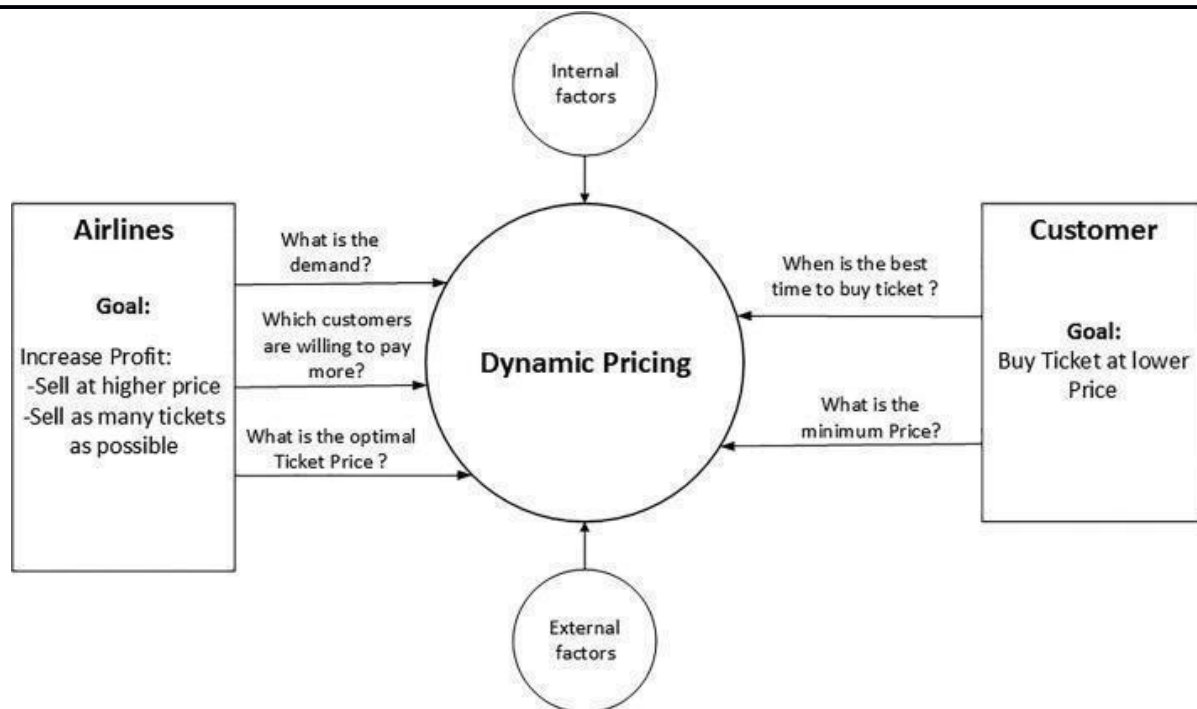
I would like to express my sincere thanks to Ms. Khushboo Garg for a regular follow up and valuable guidance provided throughout And I am also thankful to FlipRobo Technologies Bangalore for their guidance and constant supervision as well as for providing necessary information regarding the project and also for their support in completing the project.

And I am also thankful to FlipRobo Technologies Bangalore for their guidance and constant supervision as well as for providing necessary information regarding the project and also for their support in completing the project.

INTRODUCTION

Business Problem Framing

The airline industry is considered as one of the most sophisticated industry in using complex pricing strategies. Nowadays, ticket prices can vary dynamically and significantly for the same flight, even for nearby seats. The ticket price of a specific flight can change up to 7 times a day. Customers are seeking to get the lowest price for their ticket, while airline companies are trying to keep their overall revenue as high as possible and maximize their profit. However, mismatches between available seats and passenger demand usually leads to either the customer paying more or the airlines company losing revenue. Airlines companies are generally equipped with advanced tools and capabilities that enable them to control the pricing process. However, customers are also becoming more strategic with the development of various online tools to compare prices across various airline companies. In addition, competition between airlines makes the task of determining optimal pricing is hard for everyone.



Conceptual background of domain problem

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, we have to work on a project where we will collect data of flight fares with other features and work to make a model to predict fares of flights.

Analytical Problem Framing

□ Mathematical/Analytical modeling of the problem

For the given flight price prediction project I have scraped the flight prices along with some other features from a well known website that is 'yatra.com'. And this data is framed into a dataframe and saved into a .csv file. After that I have fetched this data set and performed some data processing and some EDA. The data set is having around 1645 rows and 10 columns.

Loading the data

```
In [94]: #lets import the dataset
df = pd.read_csv(r'C:\Users\hp\Dropbox\PC\Desktop\Flip Robo internship\flight price prediction\sol\Flights Details.csv')
df
```

```
Out[94]:
```

	Unnamed: 0	Name	DepTime	ArrTime	NetDuration	Stops	DepartureLoc	ArrivalLoc	Varient	Price
0	0	IndiGo	13:00	15:35	2h 35m	Non Stop	New Delhi	Bangalore	6E-5186	8,159
1	1	IndiGo	08:20	11:00	2h 40m	Non Stop	New Delhi	Bangalore	6E-308	8,159
2	2	Go First	05:50	08:35	2h 45m	Non Stop	New Delhi	Bangalore	G8-113	8,159
3	3	IndiGo	05:50	08:35	2h 45m	Non Stop	New Delhi	Bangalore	6E-6612	8,159
4	4	IndiGo	02:25	05:15	2h 50m	Non Stop	New Delhi	Bangalore	6E-5036	8,159
5	5	IndiGo	22:55	01:45	2h 50m	Non Stop	New Delhi	Bangalore	6E-6565	8,159
6	6	SpiceJet	06:05	09:00	2h 55m	Non Stop	New Delhi	Bangalore	SG-191	8,159
7	7	IndiGo	06:55	09:50	2h 55m	Non Stop	New Delhi	Bangalore	6E-5009	8,159
8	8	IndiGo	15:20	18:15	2h 55m	Non Stop	New Delhi	Bangalore	6E-781	8,159
9	9	Go First	16:45	21:35	4h 50m	1 Stop	New Delhi	Bangalore	G8-2513/242	8,159
10	10	Go First	16:00	20:55	4h 55m	1 Stop	New Delhi	Bangalore	G8-406/408	8,159

```
In [96]: #checking info about the dataset
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1645 entries, 0 to 1644
Data columns (total 10 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Unnamed: 0            1645 non-null   int64
1   Name                  1645 non-null   object
2   DepTime               1645 non-null   object
3   ArrTime               1645 non-null   object
4   NetDuration           1645 non-null   object
5   Stops                1645 non-null   object
6   DepartureLoc          1645 non-null   object
7   ArrivalLoc            1645 non-null   object
8   Varient               1645 non-null   object
9   Price                 1645 non-null   object
dtypes: int64(1), object(9)
memory usage: 128.6+ KB
```

Great, looking at the data information we can see we have one unwanted column named as 'Unnamed: 0' I will drop this column as it is not contributing to our predictions. And rest all columns are given as object data type which are required to convert into respective data types. And luckily we are not having any null values in our data set.

Data Processing

As we need to convert the data types of some feature for that I followed some of the data processing steps.

Duration column:

NetDuration Column

In [106]: *# The column NetDuration is with object datatype by doing below operations to Duration column I will extract numerical value from*

```
#Extracting numerical data using Duration
df["hour"] = df.NetDuration.str.split('h').str.get(0)
df["min"] = df.NetDuration.str.split('h').str.get(1)
df["min"] = df["min"].str.split('m').str.get(0)
df["hour"] = df["hour"].astype('float')
df["min"] = df["min"].astype('float')

df["NetDuration"] = df["hour"] + df["min"]/60
```

In [107]: *# Lets check the data set now*
df.head()

Out[107]:

	Name	DepTime	ArrTime	NetDuration	Stops	DepartureLoc	ArrivalLoc	Price	hour	min
0	IndiGo	13:00	15:35	2.583333	0	New Delhi	Bangalore	8,159	2.0	35.0
1	IndiGo	08:20	11:00	2.666667	0	New Delhi	Bangalore	8,159	2.0	40.0
2	Go First	05:50	08:35	2.750000	0	New Delhi	Bangalore	8,159	2.0	45.0
3	IndiGo	05:50	08:35	2.750000	0	New Delhi	Bangalore	8,159	2.0	45.0
4	IndiGo	02:25	05:15	2.833333	0	New Delhi	Bangalore	8,159	2.0	50.0

In [108]: *# Using hour and min column I have created Duration column with float values, now I am dropping hour and min columns*
df.drop(columns = ["hour", "min"], inplace = True)

This column is having string entries in hours and minutes; I have separated hours and minutes into two different columns by splitting and then by using both of these columns filled the entries into Duration column as given in above figure.

Departure_time & Time_of_arrival:

Time

```
In [112]: df['ArrTime']=df['ArrTime'].replace({'00:20\n+ 2 days':'00:20','00:40\n+ 2 days':'00:40'})
```

```
In [113]: #Similar to NetDuration I will extract numric data from DepTime and ArrTime columns using below codes
```

```
df["Dep_hour"] = pd.to_datetime(df.DepTime, format="%H:%M").dt.hour
df["Dep_min"] = pd.to_datetime(df.DepTime, format="%H:%M").dt.minute
df["DepTime"] = df["Dep_hour"]*60+df["Dep_min"]
df.drop(columns = ['Dep_hour','Dep_min'],inplace=True)

df["Arvl_hour"] = pd.to_datetime(df.ArrTime, format="%H:%M").dt.hour
df["arvl_min"] = pd.to_datetime(df.ArrTime, format="%H:%M").dt.minute
df["ArrTime"] = df["Arvl_hour"]*60+df["arvl_min"]
df.drop(columns = ['Arvl_hour','arvl_min'],inplace=True)
```

```
In [114]: df.head()
```

Out[114]:

	Name	DepTime	ArrTime	NetDuration	Stops	DepartureLoc	ArrivalLoc	Price
0	IndiGo	13.000000	15.583333	2.583333	0	New Delhi	Bangalore	8159.0
1	IndiGo	8.333333	11.000000	2.666667	0	New Delhi	Bangalore	8159.0
2	Go First	5.833333	8.583333	2.750000	0	New Delhi	Bangalore	8159.0
3	IndiGo	5.833333	8.583333	2.750000	0	New Delhi	Bangalore	8159.0
4	IndiGo	2.416667	5.250000	2.833333	0	New Delhi	Bangalore	8159.0

As similar to the case of duration column these two columns are also having the time but in astring format. And by using above code steps I have fetched numerical values for the time and filled to respective columns. And the extra columns from these three cases have deleted fromthe data set.

Number_of_stops:

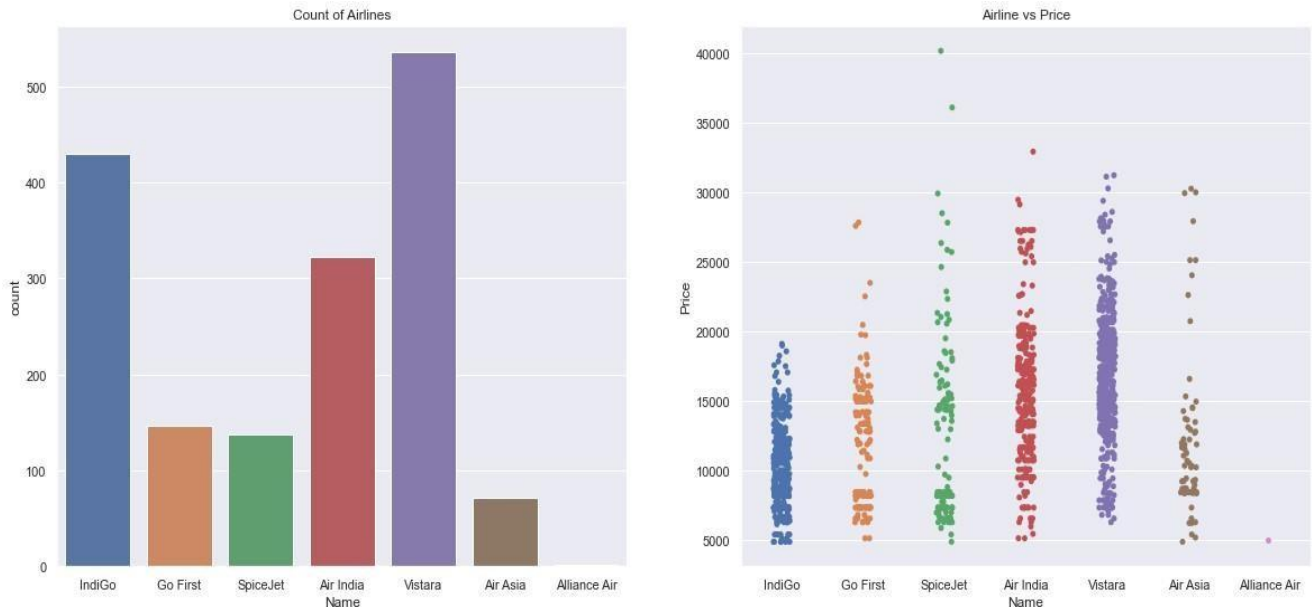
#I will replace the that categorical values from stops column to numeric data

```
df.Stops.replace({"Non Stop": 0,
                  "1 Stop": 1,
                  "2 Stop(s)": 2,
                  "3 Stop(s)": 3},
                 inplace = True)
```

This column is a categorical column and filled with string values, but we need to fill the values inordinal manner so I have replaced the entries with corresponding numeric values.

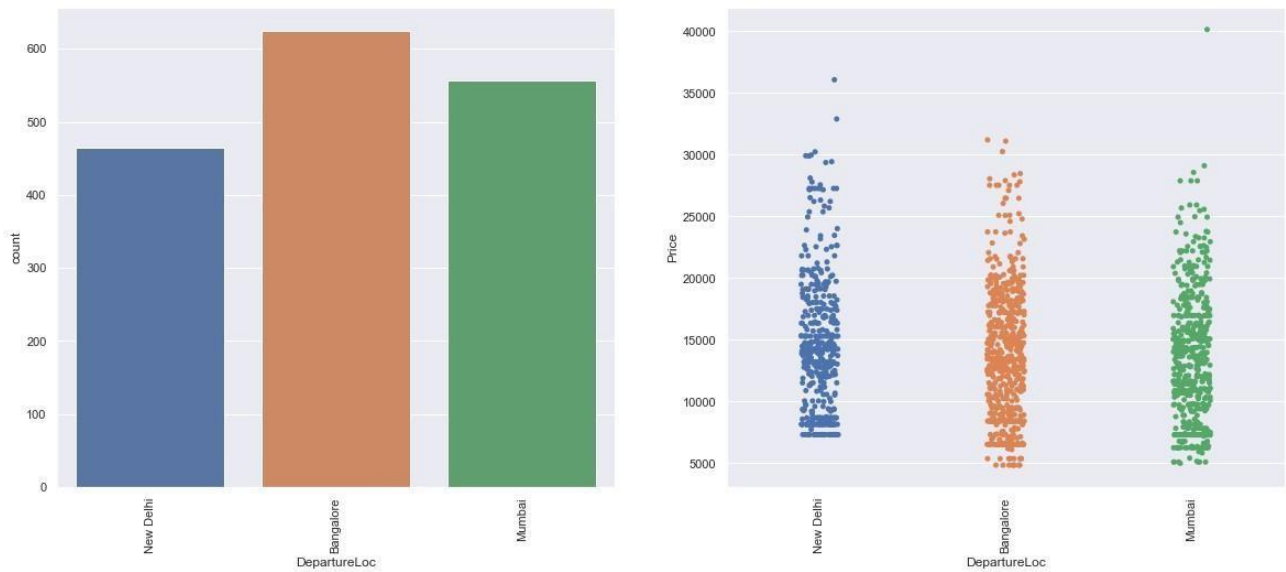
Exploratory Data Analysis

Name:



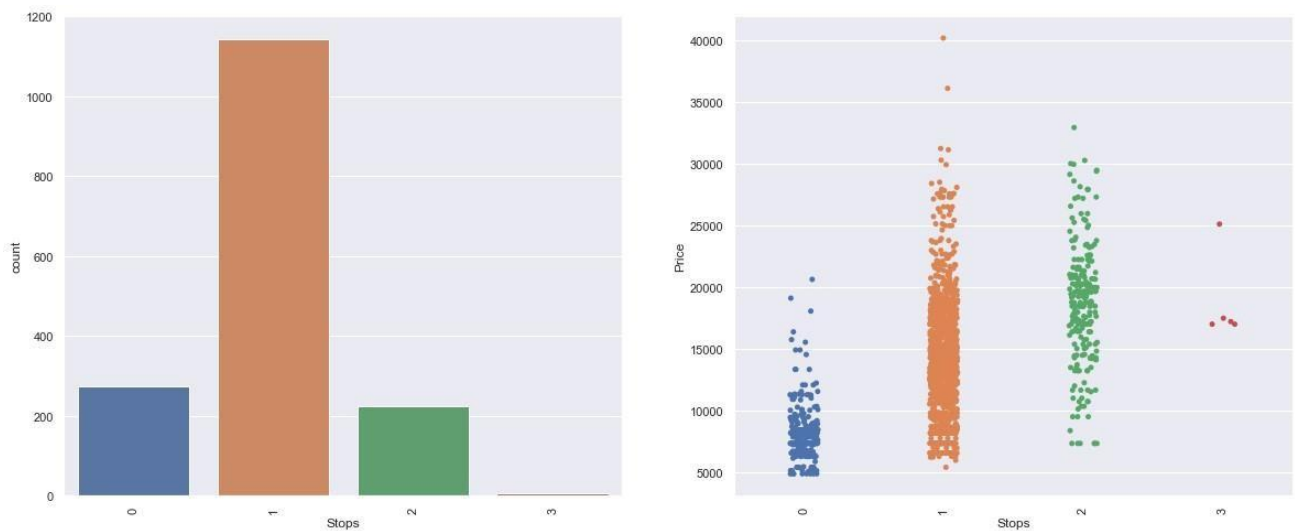
Above figure is showing two different plots one is count plot for the column Name of airline and another is a strip plot showing relation between N Airline and Prices. The count plot will tell us that there are more numbers of flights of Vistara, IndiGo and Air India than others. Flights of Airasia are very less in numbers. Strip plot will tell us that Price of Spice jet airplanes are higher than IndiGo.

Source:

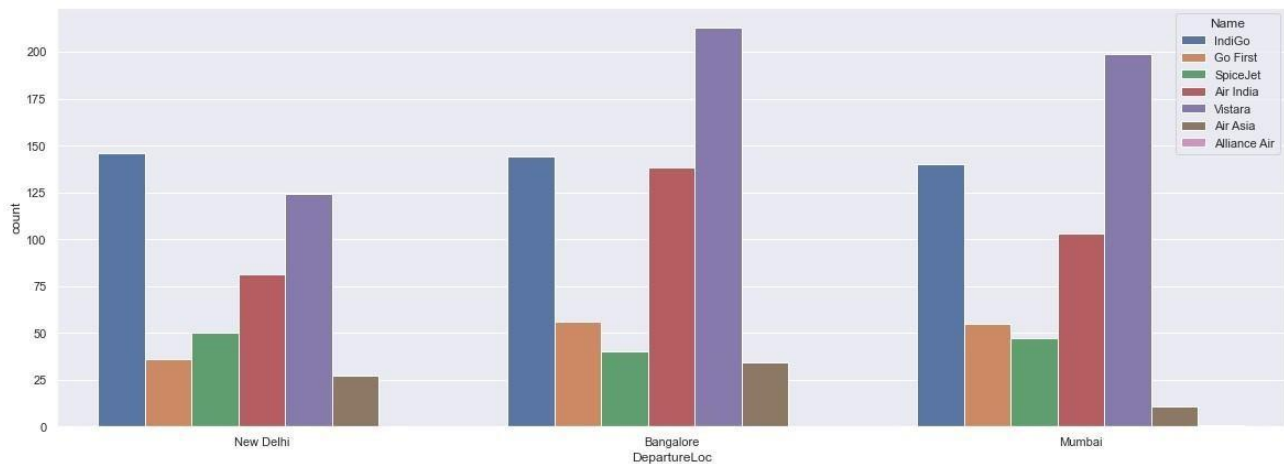


Here we are having Higher number of flights from Bangalore followed by Mumbai and New Delhi. Looking at the strip plot we can say flights from New Delhi are having some what higher prices than other cities, whereas flight from Mumbai have lowest Prices.

Number_of_stops:

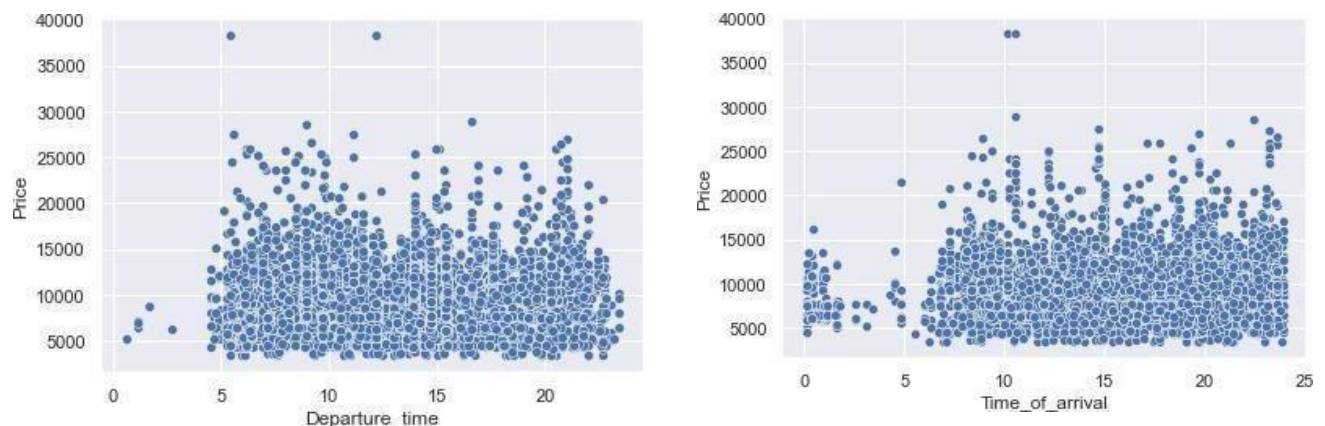


The above count plot will tell us that most of the flights are with 1 stop and very few are with 3 during the Journey. We can see that the prices are increasing with the number of stops.



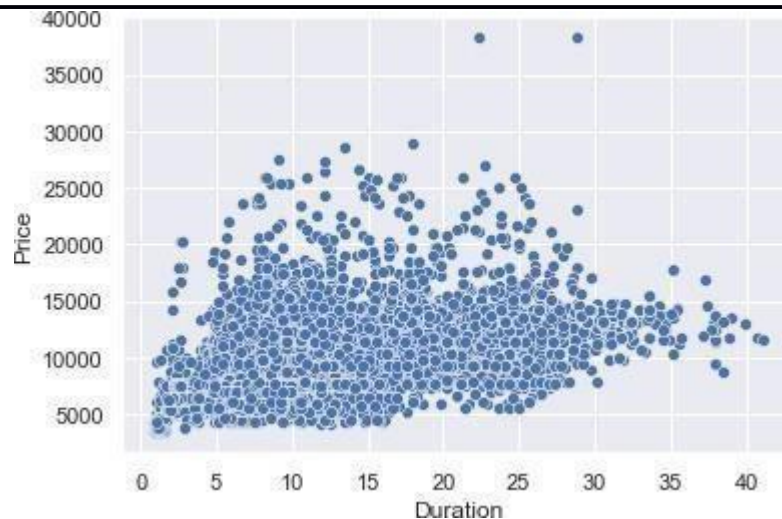
Bangalore is popular with Vistara airlines followed by Indigo. Similarly New Delhi is popular with IndiGo airlines followed by Vistara. And Mumbai is popular with Vistara airlines followed by Indigo.

Vistara and IndiGo are mostly popular in Most of the state. Air asia is least Available.



The first scatter plot is showing relationship between Departure time and flight prices. We can observe that there are very few flights departing in the early morning which are having lower price as well.

Second scatter plot is showing relation between Time of arrival and flight prices, which will tell us that very few numbers of flights are arriving in the early morning that is around 0 to 5 am. We can say the flight prices are not much dependent on the time of arrival.



The above figure is representing the scatter plot of Duration vs Price. Looking at this figure we can say that there is some linear relation between price and duration. The prices increase with duration.

Hardware and Software Requirements and Tools Used

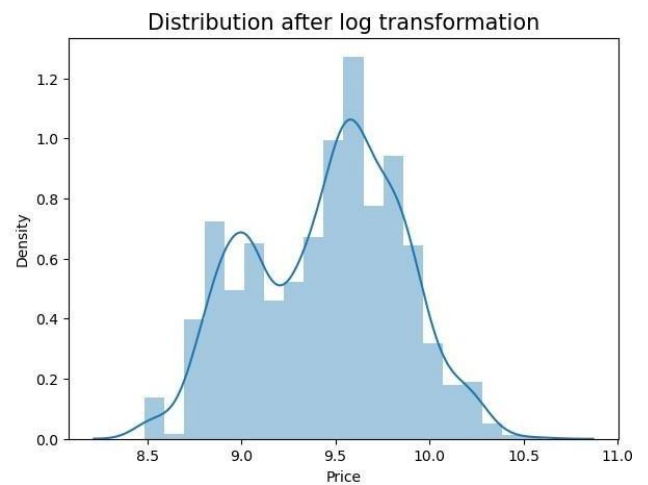
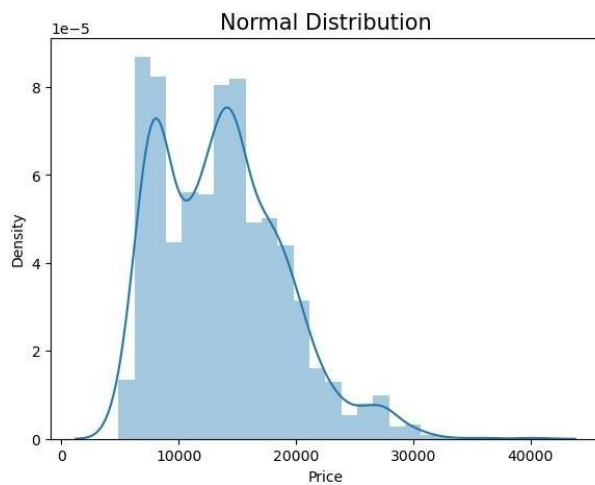
- For hardware I have used my laptop that have ryzen 5 processor and 8gb ram
- For software I have used Jupyter notebook
- For Tools I have use this following
 - library-Numpy
 - Pandas
 - Seaborn
 - Matplotli
 - bSklearn

Data Pre-processing

- ☐ Skewness treatment
- ☐ Normalizing the data(applying StandardScaler to numerical features)
- ☐ Encoding categorical features(Using OrdinalEncoder)

Model Development and Evaluation

Applying log transformation to our target variable



As we can see our target variable is right skewed that's why I will apply log transformation to it for better results.

For this project I have applied StandardScaler to numerical features for bringing it to a common scale and used ordinal encoder for categorical features. After doing such pre-processing steps I used this data for model building with the help of `train_test_split`. I have defined a function to train and evaluate our algorithms. For this task I have used many regression algorithms and selected `LGBMRegressor` as it is giving better performance than other algorithms. Linear models were giving least difference in r^2 -score and cv-score but in this case the r^2 -score was very less compared to tree based algorithms. Other than `LGBMRegressor` algorithms showing the problem of over-fitting so I selected `LGBMRegressor` which is not over-fitting much compared to others.

For this project I have used following algorithms:

- `LinearRegressionC`
- `LassoCV`
- `RidgeCV`
- `DecisionTreeRegressor`
- `RandomForestRegressor`
- `XGBRegressor`
- `ExtraTreesRegressor`
- `LGBMRegressor`

From all of these above models LGBMRegressor was giving me good performance.

KeyMetrics for success in solving problem under consideration

I have used the following metrics for evaluation:

- I have used mean absolute error which gives magnitude of difference between the prediction of an observation and the true value of that observation.
- I have used root mean square deviation is one of the most commonly used measures for evaluating the quality of predictions.
- I have used r2 score which tells us how accurate our model is.
- Also I have checked both training and testing r2-scores to check the over-fitting and under-fitting.
- Cross Validation score to check the model performance on overall data.

Hyperparameter Tuning

I did hyperparameter tuning for LGBMRegressor for the parameters like 'boosting_type', 'max_depth', 'learning_rate', 'n_estimators' using GridSearchCV

```
{'boosting_type': 'gbdt',  
 'learning_rate': 0.1,  
 'max_depth ': -1,  
 'n_estimators': 800}
```

After running the code for above mentioned parameters I got the values which are indicated in the above figure as best parametric values for our final model.

Using these parametric values I trained our final model and got good r2-score better than earlier.

Final Model

Final model

```
In [170]: #lets train and test our final model with best parameters
model = LGBMRegressor(boosting_type = 'gbdt', learning_rate = 0.1, n_estimators = 800, max_depth=-1)
model.fit(x_train,y_train)
pred = model.predict(x_test)

r2score = r2_score(y_test,pred)*100

#evaluation
mse = mean_squared_error(y_test,pred)
rmse = np.sqrt(mse)
mae = mean_absolute_error(y_test,pred)
print("MAE :", mae)
print("RMSE :", rmse)
print('-----')

# r2 score

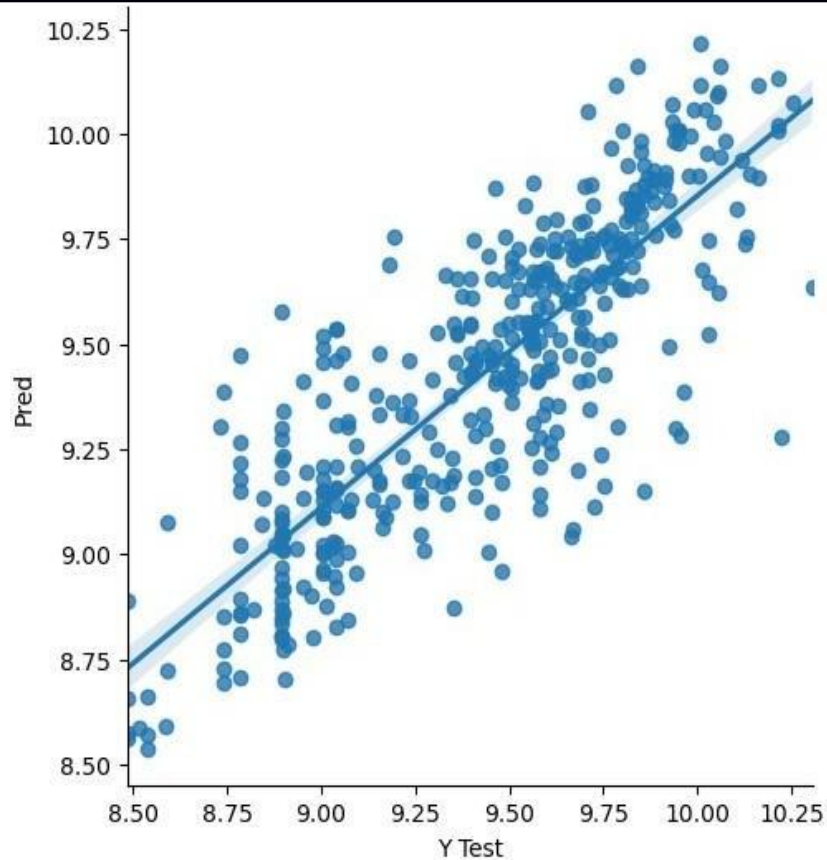
print(f" \nr2 Score:", r2score,"%")

MAE : 0.16477949722748855
RMSE : 0.22862356424245794
-----

r2 Score: 68.30705958920052 %
```

Great we have achieved a better r2 score after doing hyperparameter tuning than earlier.

Let's see the graph for actual vs predicted values



Great we have achieved a better r^2 score after doing hyperparameter tuning than earlier.

Conclusion

Key findings of the study

In this project we have scraped the flight data from yatra.com. Then the .csv file is loaded into a data frame.

Luckily we don't have any missing values in our data set.

Looking at the data set we understand that there are some features that need to be processed like converting the data types, and get the actual value from the string entries from the time-related columns.

After the data is processed I have done some EDA to understand the relation among features and the target variable.

Features like flight duration, number of stops during the journey and the availability of meals are playing a major role in predicting the prices of the flights.

Limitations of this work and scope for the future work

As looking at the features we came to know that the numbers of features are very less, due to which we are getting somewhat lower r^2 -scores.

Some algorithms are facing over-fitting problem which may be because of less number of features in our dataset.

We can get a better r^2 score than now by fetching some more features from the web scraping by that we may also reduce the over fitting problem in our models.

Thank You