



# FLIGHT TICKET PRICE PREDICTION

Submitted by:

VANISREE P G

# ACKNOWLEDGMENT

First I would like to thank the Almighty for his wonderful presence with me throughout this project and helped me to make it as a successful one.

For my internship I had the pleasure of working at FILP ROBO Was a great chance for acquired knowledge, personal and Professional development.

I extend whole hearted thanks to FILP ROBO under whom I worked and learned a lot and for enlightening me with their knowledge and experience to grow with the corporate working.

This is a great pleasure to express my deep sense of gratitude and thanks to SME for his valuable ideas, instantaneous help, effective support and continued encouragement which enabled for the successful completion of the project. I also like to thank the data trained mentors and Technical team members for helping me with technical queries.

And these are the following website which I referred for the reference

1. <https://www.kaggle.com/>
2. <https://scikit-learn.org/>
3. [www.stackoverflow.com](http://www.stackoverflow.com)
4. [www.google.com](http://www.google.com)
5. [www.geeksforgeeks.org](http://www.geeksforgeeks.org)

# INTRODUCTION

## Business Problem Framing

With respect to the season flight price also differs in the market, we have seen lot of changes in the Flight ticket price. Now some flight tickets are in demand hence making them costly and some are not in demand hence cheaper. With the change in market due to covid 19 impact, our customers are facing problems with their previous flight price valuation machine learning models. So, they are looking for new machine learning models from new data. We have to make flight ticket price valuation model.

Someone who purchase flight tickets frequently would be able to predict the right time to procure a ticket to obtain the best deal. Many airlines change ticket prices for their revenue management. The airline may increase the prices when the demand is to be expected to increase the capacity. To estimate the minimum airfare, data for a specific air route has been collected including the features like departure time, arrival time and airways over a specific period. Features are extracted from the collected data to apply Machine Learning (ML) models. This paper gives the machine learning regression methods to predict the prices at the given time.

As domestic air travel in India is becoming increasingly popular with different air ticket booking channels coming online these days, passengers are trying to understand how these airline companies make decisions over time about ticket prices. Therefore, many methods are ready to provide the proper time to do so. The customer who buys an air ticket by estimating the price of the airfare is recently proposed. The majority of these strategies make use of sophisticated Computational Intelligence Prediction Models an area of science known as Machine Learning (ML). This paper

highlights the parameters and also includes the guidelines that are important for project work to be developed that is indicated above.

Airline price ticket costs modification terribly dynamically and for a similar flight day by day. It is terribly tough for a customer to buy an air ticket within the lowest value since the value changes dynamically. We addressed the matter regarding the market section level airfare ticket cost forecasting by usage of publicly obtainable datasets and completely unique machine learning model to forecast market section level price cost of airline ticket. The purpose of this study is to raise and analyze the options that influence transportation and to develop and tune models to predict the transportation well ahead.

Planes ticket prices changes as time passes, pulling out the elements which creates the difference. Reporting the correlated and models which is used to price the flight tickets. Then, using that information, building the model which helps passengers to make pull out the ticket to buy and predicting air ticket prices which progresses in the future. Duration, Arrival time, Price, Source, Destination and much more these are the attribute used for flight price prediction.

## **Motivation for the Problem Undertaken**

To understand real world problems where Machine Learning and Data Analysis can be applied to help to predict the prices in various domains to make better decisions with the help of which they can gain profit or can be escaped from any loss which otherwise could be possible without the study of data.

The evaluating system as a result changes the charge depending on time, season, and festive days to change the header or footer on

successive pages. The ultimate aim of the airways is to earn profit whereas the customer searches for the minimum rate. Customers usually try to buy the ticket well in advance of departure date so as to avoid hike in airfare as date comes closer. But actually this is not the fact. The customer may wind up by giving more than they ought to for the same seat.

These days, domestic air travel is becoming more and more common in India. Travellers are trying to learn how these airline companies make choices over time about ticket prices with multiple air travel booking outlets coming online. For a passenger, it is a time-consuming method to search websites for deals and offers. The cost can therefore depend on various variables. This venture uses AI to show the types off light tickets after some time to estimate the costs. Both organizations have the right and the ability at any time to change their ticket prices. By reserving a ticket at the lowest cost, explorer can set aside money. People who have travelled by flight are also aware of the variations in costs. Complex revenue control policies are used by airlines for the introduction of distinctive assessment schemes.

As a result, the appraisal scheme adjusts the fee to adjust the header or footer on successive pages based on time, season, and festive days. The ultimate goal of the airways is to achieve profit ,while the customer is looking for the minimum cost. Usually consumers try to book the ticket well in advance of the departure date to prevent airfare hikes as the date gets closer. But that's not the truth, really. By giving more than they should for the same seat, the customer can finish up.

# **Analytical Problem Framing**

## **Mathematical/ Analytical Modelling of the Problem**

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. And few of them was removed due to very less correlation with dependent variable. Since the dataset contains a lot of features hence feature selection has been also done.

In machine learning, several algorithms are applied to forecast the prices of flight tickets. The algorithms are: Linear regression, Decision tree, Gradient Boosting Regression, and Random Forest Algorithm. These models have been implemented using the python library Sklearn. The parameters like MAE and MSE, RMSE are considered to check the efficiency of these models

Improving the ML structure to predict the mean plane price for the business purpose. For predicting the mean plane price with modification of R squared score, feature selection techniques were proposed in our model. Comparing the production of various ML classifiers which tells the greater plane price prediction task. Facts gathered from website that sells the planes ticket through internet apps. Authors have reported that there is limited public information access which will miss the main target attribute. Final accountable prediction model is improved by two unrelated prediction models such as Random Forest and Multilayer Perceptron. Weights for drifting with R-square value and the main estimation of the metric was used.

## Data Sources and their formats

The most critical aspect of this project is the accumulation of knowledge. To prepare the models, the distinct well springs of the data on various sites are used. Sites provide information on the different firms, hours, aircraft, and charges. For data scratching, various sources from API's to customer travel sites are available.

Data is Scrapped from a Yatra.com is an Indian online travel agency and travel search engine web portal for booked a flight ticket. The data is scraped from the Yatra.com website The data descriptions are as follow:- (2103, 10) rows and columns. To predict Flight prices using Regression. I will start by importing all the necessary libraries that we need for this task and import the dataset.

1) Importing libraries

2) Importing the dataset

They are totally 2103 rows and 10 columns in a csvfile. Our target is to find the insights of the data and to do thorough data analysis.

### 1.Uploading Data set

```
In [1]: import pandas as pd
import numpy as np
import matplotlib as mlp
import seaborn as sns
import matplotlib.pyplot as plt
import warnings
warnings.filterwarnings('ignore')
```

```
In [2]: df = pd.read_csv("Flight price data.csv")
```

```
In [2]: df = pd.read_csv("Flight price data.csv")

In [3]: df
Out[3]:
```

|      | Unnamed: 0 | Date   | Airline   | Source    | Destination | Dep_time | Arr_time | Duration | Total_Stops | Price  |
|------|------------|--------|-----------|-----------|-------------|----------|----------|----------|-------------|--------|
| 0    | 0          | 25-Jul | Vistara   | Chennai   | Mumbai      | 09:45    | 11:40    | 1h 55m   | Non Stop    | 8,666  |
| 1    | 1          | 25-Jul | Vistara   | Chennai   | Mumbai      | 12:30    | 14:30    | 2h 00m   | Non Stop    | 8,666  |
| 2    | 2          | 25-Jul | Go First  | Chennai   | Mumbai      | 13:40    | 15:40    | 2h 00m   | Non Stop    | 8,666  |
| 3    | 3          | 25-Jul | Vistara   | Chennai   | Mumbai      | 20:30    | 22:35    | 2h 05m   | Non Stop    | 8,666  |
| 4    | 4          | 25-Jul | Air India | Chennai   | Mumbai      | 15:25    | 17:10    | 1h 45m   | Non Stop    | 8,668  |
| ...  | ...        | ...    | ...       | ...       | ...         | ...      | ...      | ...      | ...         | ...    |
| 2098 | 101        | 27-Jul | Vistara   | New Delhi | Goa         | 20:35    | 12:45    | 16h 10m  | 2 Stop(s)   | 23,171 |
| 2099 | 102        | 27-Jul | Vistara   | New Delhi | Goa         | 20:35    | 12:45    | 16h 10m  | 2 Stop(s)   | 23,171 |
| 2100 | 103        | 27-Jul | SpiceJet  | New Delhi | Goa         | 19:00    | 12:55    | 17h 55m  | 1 Stop      | 24,823 |
| 2101 | 104        | 27-Jul | Vistara   | New Delhi | Goa         | 20:40    | 12:45    | 16h 05m  | 2 Stop(s)   | 26,953 |
| 2102 | 105        | 27-Jul | Vistara   | New Delhi | Goa         | 20:40    | 12:45    | 16h 05m  | 2 Stop(s)   | 26,953 |

2103 rows x 10 columns

## Data Pre-processing

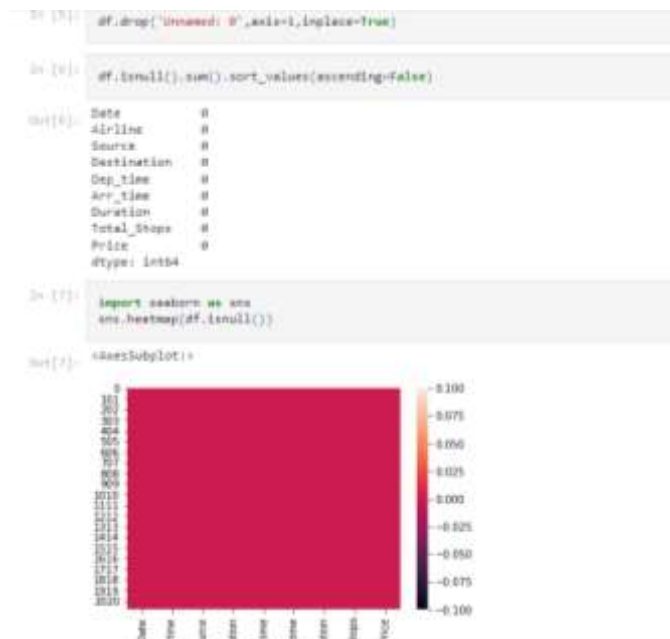
Before building model, the data should be properly pre processed and converted to quality, clean data even the resulting machine learning model will be of great quality .The data pre-processing includes three main parts that is data integration, data cleaning, data transformation. In data integration the data collected from various sources are integrated. In data cleaning process the data containing the null values, unnecessary rows with null values are being cleared. The data transformation includes the feature scaling ,categorical data, etc to set the certain range of data.

The raw data is taken and performed various steps to reduce skewness, outlier, class imbalance and scaling. There were null value was present and removed the values from the dataset. Many outlier removal and skewness removal methods are tested and best method is chosen in order to prevent data loss.

- The dataset contains 2103 rows and 10 columns
- Price is our dependent variable.
- We created new features from old ones.
- There are no null values in the dataset.
- Removed empty cells



### a) Checking missing value from the data set.



- There were no null value was present in the dataset and there is no outliers are present in the data.

### b) CORRELATION

Correlation between all the columns in the datasets. In the correlation Heat map, we have the following observations:

- Arrive time has 41 % correlation with target column which can be considered as a good bond.
- Departure time has 33 % correlation with target column which can be considered as a good bond.
- Duration has 17 % correlation with target column which can be considered as a good bond.
- Day has 16 % correlation with target column which can be considered as a good bond.

- source has 77 % correlation with target column which can be considered as a good bond.
- Airline has -0.068% correlation with target column which can be considered as a negative bond.

Correlation

```
In [77]: df.apply(lambda x: pd.factorize(x)[0]).corr(method='pearson', min_periods=1)
```

```
Out[77]:
```

|             | Airline   | Source    | Destination | Day       | Month    | Dep_time  | Arr_time | Dur_hour  | Dur_min   | Total_Stops | Price     |
|-------------|-----------|-----------|-------------|-----------|----------|-----------|----------|-----------|-----------|-------------|-----------|
| Airline     | 1.000000  | -0.041333 | 0.005968    | 0.006666  | NaN      | 0.005968  | 0.021563 | -0.328568 | 0.036875  | -0.254209   | -0.068052 |
| Source      | -0.041333 | 1.000000  | 0.422914    | 0.169687  | NaN      | 0.236819  | 0.215046 | 0.122541  | 0.037967  | 0.024925    | 0.768768  |
| Destination | 0.005968  | 0.422914  | 1.000000    | 0.178548  | NaN      | 0.198413  | 0.537545 | 0.148952  | 0.123185  | 0.136287    | 0.857347  |
| Day         | 0.006666  | 0.169687  | 0.178548    | 1.000000  | NaN      | 0.036787  | 0.117938 | -0.011818 | 0.048617  | -0.021811   | 0.185218  |
| Month       | NaN       | NaN       | NaN         | NaN       | 1.000000 | NaN       | NaN      | NaN       | NaN       | NaN         | NaN       |
| Dep_time    | 0.005968  | 0.236819  | 0.198413    | 0.036787  | NaN      | 1.000000  | 0.167180 | 0.042080  | -0.037878 | -0.081848   | 0.234821  |
| Arr_time    | 0.021563  | 0.215046  | 0.537545    | 0.117938  | NaN      | 0.167180  | 1.000000 | 0.000000  | 0.126435  | 0.025964    | 0.417605  |
| Dur_hour    | -0.328568 | 0.122541  | 0.148952    | -0.011818 | NaN      | 0.042080  | 0.000000 | 1.000000  | 0.034925  | 0.568940    | 0.172574  |
| Dur_min     | 0.036875  | 0.037967  | 0.123185    | 0.048617  | NaN      | -0.037878 | 0.126435 | 0.034925  | 1.000000  | -0.000684   | 0.021673  |
| Total_Stops | -0.254209 | 0.024925  | 0.136287    | -0.021811 | NaN      | -0.081848 | 0.025964 | 0.568940  | -0.000684 | 1.000000    | 0.140051  |
| Price       | -0.068052 | 0.768768  | 0.857347    | 0.185218  | NaN      | 0.234821  | 0.417605 | 0.172574  | 0.021673  | 0.140051    | 1.000000  |



## Data Inputs- Logic- Output Relationships

The input data contains 2103 rows and 10 columns.

Predictor variable are,

Total\_stops, Durtion, Arr\_time, Dep\_time, Destination, Source, Airline, Date.

Target variable is Price of the flight ticket

## **Hardware and Software Requirements and Tools Used**

Hardware used for doing the project is a 'Laptop' with high end specification and stable internet connection .while coming to the software part I had used 'python jupyter notebook' for do my python program and data analysis.

Excel file and Microsoft excel are required for the data handling. In jupyter notebook I had imported lot of python libraries are carried to this project.

1.Pandas-a library which is used to read the data ,visualisation and analysis of data.

2.Numpy-used for working with array and various mathematical operations in python.

3.Seaborn- visualization for plotting different type of plot.

4.Matplotlib- It provides an object-oriented API for embedding plots into applications .

## **Model/s Development and Evaluation**

### **Identification of possible problem-solving approaches (methods)**

In machine learning, several algorithms are applied to forecast the prices of flight tickets. The algorithms are: Linear regression, Decision tree, SVR, Gradient Boosting Regression, Ridge and Random Forest Algorithm. These models have been implemented using the

python library Sklearn. The parameters like MAE and MSE, RMSE are considered to check the efficiency of these models.

Regression Model with following algorithms

- Linear Regression
  - Decision Tree Regressor
  - Random forest regressor
  - SVR
  - Gradient Boosting Regressor
  - Ridge
    - Evaluation metrics
- Mean square error
  - Mean absolute error
  - R2 score
  - Root Mean Squared Error

## **Testing of Identified Approaches (Algorithms)**

Listing down all the algorithms used for the training and testing.

- LR=LinearRegression()
- DT=DecisionTreeRegressor()
- rf=RandomForestRegressor()
- svr=SVR()
- R=Ridge()
- GBR=GradientBoostingRegressor()

# Run and Evaluate selected models

Describe all the algorithms used along with the snapshot of their code and what were the results observed over different evaluation metrics. Train-test data splits were conducted. In this situation, we split the data into training and test sets, then fit candidate models on the training set, evaluate and select them on the test set.

## 1. Linear Regression

```
i) LinearRegression

In [50]:
lr=LinearRegression()
lr.fit(x_train,y_train)
print(lr.score(x_train,y_train))
lr_predict=lr.predict(x_test)

0.346582817984358

In [51]:
print('Mean Absolute Error: ',metrics.mean_absolute_error(y_test,lr_predict))
print('Mean Squared Error: ',metrics.mean_squared_error(y_test,lr_predict))
print('Root Mean Squared Error: ',np.sqrt(metrics.mean_squared_error(y_test,lr_predict)))
print('Explained Variance Score: ',metrics.explained_variance_score(y_test,lr_predict))
print('r2_score: ',r2_score(y_test,lr_predict))

Mean Absolute Error: 11.89077187920801
Mean Squared Error: 245.8485112838346
Root Mean Squared Error: 15.682775275005
Explained Variance Score: 0.4130182158791174
r2_score: 0.41525138127289485

C/V score

In [52]:
lr_cv=cross_val_score(lr,x,y, cv = 7).mean()
lr_cv

Out[52]: 0.21430167978793884
```

## 2. Random forest

```
ii) RandomForestRegressor

In [53]:
from sklearn.ensemble import RandomForestRegressor
rf=RandomForestRegressor()
rf.fit(x_train,y_train)
prediction1=rf.predict(x_test)
print(rf.score(x_train,y_train))

0.3465728712514875

In [54]:
print('Mean Absolute Error: ',metrics.mean_absolute_error(y_test,prediction1))
print('Mean Squared Error: ',metrics.mean_squared_error(y_test,prediction1))
print('Root Mean Squared Error: ',np.sqrt(metrics.mean_squared_error(y_test,prediction1)))
print('Explained Variance Score: ',metrics.explained_variance_score(y_test,prediction1))
print('r2_score: ',r2_score(y_test,prediction1))

Mean Absolute Error: 11.4217472138578962
Mean Squared Error: 84.90085565083333
Root Mean Squared Error: 9.21433114844292
Explained Variance Score: 0.7975725481981131
r2_score: 0.797454361222989

In [55]:
rf_cv=cross_val_score(rf,x,y, cv = 4).mean()
rf_cv

Out[55]: 0.2749264453295447
```

### 3. Decision Tree Regression

#### iii) DecisionTreeRegressor

```
In [76]: from sklearn.tree import DecisionTreeRegressor

DTR=DecisionTreeRegressor()
DTR.fit(x_train,y_train)
print(DTR.score(x_train,y_train))
DTR_PRED=DTR.predict(x_test)

0.9973279629636064

In [77]: print('Mean Absolute Error: ',metrics.mean_absolute_error(y_test,DTR_PRED))
print('Mean Squared Error: ',metrics.mean_squared_error(y_test,DTR_PRED))
print('Root Mean Squared Error: ',np.sqrt(metrics.mean_squared_error(y_test,DTR_PRED)))
print('Explained Variance Score: ',metrics.explained_variance_score(y_test,DTR_PRED))
print('r2_score: ',r2_score(y_test,DTR_PRED))

Mean Absolute Error: 4.89895148885491
Mean Squared Error: 385.8052362782723
Root Mean Squared Error: 19.38423791338889
Explained Variance Score: 0.7522883827808882
r2_score: 0.7522883827808882

In [78]: DTR_cv=cross_val_score(DTR,x,y, cv = 11).mean()
DTR_cv

0.15022229488744124

Out[78]:
```

#### ➤ SVR

```
In [79]: from sklearn.svm import SVR

svr=SVR()
svr.fit(x_train,y_train)
print(svr.score(x_train,y_train))
svr_predict=svr.predict(x_test)

0.25254517669953418

In [80]: print('Mean Absolute Error: ',metrics.mean_absolute_error(y_test,svr_predict))
print('Mean Squared Error: ',metrics.mean_squared_error(y_test,svr_predict))
print('Root Mean Squared Error: ',np.sqrt(metrics.mean_squared_error(y_test,svr_predict)))
print('Explained Variance Score: ',metrics.explained_variance_score(y_test,svr_predict))
print('r2_score: ',r2_score(y_test,svr_predict))

Mean Absolute Error: 12.4618287792678
Mean Squared Error: 297.3888667973557
Root Mean Squared Error: 17.244968886418572
Explained Variance Score: 0.33525204425889516
r2_score: 0.2985372922943853

In [81]: svr_cv=cross_val_score(svr,x,y, cv = 4).mean()
svr_cv

0.043016476338164555

Out[81]:
```

### 4. Gradient Boosting Regression

```
In [82]: from sklearn.ensemble import GradientBoostingRegressor

GBR=GradientBoostingRegressor()
GBR.fit(x_train,y_train)
print(GBR.score(x_train,y_train))
GBR_PRED=GBR.predict(x_test)

0.768487252847124

In [83]: print('Mean Absolute Error: ',metrics.mean_absolute_error(y_test,GBR_PRED))
print('Mean Squared Error: ',metrics.mean_squared_error(y_test,GBR_PRED))
print('Root Mean Squared Error: ',np.sqrt(metrics.mean_squared_error(y_test,GBR_PRED)))
print('Explained Variance Score: ',metrics.explained_variance_score(y_test,GBR_PRED))
print('r2_score: ',r2_score(y_test,GBR_PRED))

Mean Absolute Error: 7.59578558331815
Mean Squared Error: 106.9045288938332
Root Mean Squared Error: 10.34432822155862
Explained Variance Score: 0.7480433841418741
r2_score: 0.739888882351179

In [84]: GBR_cv=cross_val_score(GBR,x,y, cv = 11).mean()
GBR_cv

0.5284383822883088

Out[84]:
```

## 5. Ridge regression

```
In [55]: from sklearn.linear_model import Ridge

R=Ridge()
R.fit(x_train,y_train)
print(R.score(x_train,y_train))
R_predict=R.predict(x_test)

0.3464819136695814

In [56]: print('Mean Absolute Error: ',metrics.mean_absolute_error(y_test,R_predict))
print('Mean Squared Error: ',metrics.mean_squared_error(y_test,R_predict))
print('Root Mean Squared Error: ',np.sqrt(metrics.mean_squared_error(y_test,R_predict)))
print('Explained Variance Score: ',metrics.explained_variance_score(y_test,R_predict))
print('r2_score: ',r2_score(y_test,R_predict))

Mean Absolute Error: 11.888340454009262
Mean Squared Error: 245.8994878888852
Root Mean Squared Error: 15.674804237711728
Explained Variance Score: 0.41418858697458144
r2_score: 0.41384676888952641

In [57]: R_cv=cross_val_score(R,x,y, cv = 7).mean()
R_cv

Out[57]: 0.23467791885817187
```

Evaluating the model accuracy is an essential part of the process of creating machine learning models to describe how well the model is performing in its predictions. The MSE, MAE, and RMSE metrics are mainly used to evaluate the prediction error rates and model performance in regression analysis.

- MAE (Mean absolute error) represents the difference between the original and predicted values extracted by averaged the absolute difference over the data set.
- MSE (Mean Squared Error) represents the difference between the original and predicted values extracted by squared the average difference over the data set.
- RMSE (Root Mean Squared Error) is the error rate by the square root of MSE.

After evaluating the model based on MAE, MSE, RMSE, EVS, R2 SCORE the best model choose for hyper parameter tuning are RandomForestRegressor, DecisionTreeRegressor.

A. RandomForestRegressor

B. DecisionTreeRegressor

# Hyper parameter tuning

## A. RandomForestRegressor

```
In [68]: RF=RandomForestRegressor()
         param={
             'n_estimators':[100,200],
             'criterion':['mse','mae'],
             'min_samples_split':[2],
             'min_samples_leaf': [1],
         }

In [69]: from sklearn.model_selection import GridSearchCV
         RF_grid=GridSearchCV(RandomForestRegressor(),param,cv=4)

In [70]: RF_grid.fit(x_train,y_train)
         RF_grid_PXS=RF_grid.best_estimator_.predict(x_test)

In [71]: RF_grid.best_params_

Out[71]: {'criterion': 'mae',
          'min_samples_leaf': 1,
          'min_samples_split': 2,
          'n_estimators': 200}

In [72]: rf=RandomForestRegressor(criterion='mae',min_samples_leaf=1,min_samples_split=2,n_estimators=200)
         rf.fit(x_train,y_train)
         rf_predictions=rf.predict(x_test)

In [73]: import numpy as np
         from sklearn import metrics
         print('Mean Absolute Error: ',metrics.mean_absolute_error(y_test,rf_predictions))
         print('Mean Squared Error: ',metrics.mean_squared_error(y_test,rf_predictions))
         print('Root Mean Squared Error: ',np.sqrt(metrics.mean_squared_error(y_test,rf_predictions)))
         print('Explained Variance Score: ',metrics.explained_variance_score(y_test,rf_predictions))
         print('r2_score:',r2_score(y_test,rf_predictions))

Mean Absolute Error: 5.356798154843619
Mean Squared Error: 82.15186168414999
Root Mean Squared Error: 9.063766413812196
Explained Variance Score: 0.8041543358112451
r2_score: 0.8040153388822685
```

## B. Decision Tree Regression

```
In [74]: from sklearn.model_selection import GridSearchCV
         parameter = {'max_features':['auto', 'sqrt', 'log2'],
                       'criterion':['mse','friedman_mse','mae', 'poisson'],
                       'splitter':['best', 'random']}
         GCV = GridSearchCV(DecisionTreeRegressor(),parameter,cv=5)

In [75]: GCV.fit(x_train,y_train)

Out[75]: GridSearchCV(cv=5, estimator=DecisionTreeRegressor(),
                    param_grid={'criterion': ['mse', 'friedman_mse', 'mae', 'poisson'],
                                'max_features': ['auto', 'sqrt', 'log2'],
                                'splitter': ['best', 'random']})

In [76]: GCV.best_params_

Out[76]: {'criterion': 'friedman_mse', 'max_features': 'auto', 'splitter': 'best'}

In [77]: from sklearn.tree import DecisionTreeRegressor
         DTR=DecisionTreeRegressor(criterion='friedman_mse',splitter='best',max_features='auto')
         DTR.fit(x_train,y_train)
         DTR_final=DTR.predict(x_test)
```



```

In [76]: print("Score of Model is",DTR.score(x_train,x_test))
print("Mean Absolute Error", mean_absolute_error(y_test,DTR_Final))
print("Root Mean Squared Error", (mean_squared_error(y_test,DTR_Final))**0.5 )

Score of Model is 0.897317963966584
Mean Absolute Error 5.3843944533224
Root Mean Squared Error 38.9721762428284

In [77]: print("Mean Absolute Error: ",metrics.mean_absolute_error(y_test,DTR_Final))
print("Mean Squared Error: ",metrics.mean_squared_error(y_test,DTR_Final))
print("Root Mean Squared Error: ",np.sqrt(metrics.mean_squared_error(y_test,DTR_Final)))
print("Explained Variance Score: ",metrics.explained_variance_score(y_test,DTR_Final))
print("r2_score",r2_score(y_test,DTR_Final))

Mean Absolute Error: 5.3843944533224
Mean Squared Error: 150.38885142393882
Root Mean Squared Error: 38.9721762428284
Explained Variance Score: 0.7188939488322122
r2_score: 0.7127961888864142

```

The best model after hyper parameter tuning is Random Forest Regressor

```

In [80]: print("FINAL MODEL")
print("-----")
print("Mean Absolute Error: ",metrics.mean_absolute_error(y_test,rf_predictions))
print("Mean Squared Error: ",metrics.mean_squared_error(y_test,rf_predictions))
print("Root Mean Squared Error: ",np.sqrt(metrics.mean_squared_error(y_test,rf_predictions)))
print("Explained Variance Score: ",metrics.explained_variance_score(y_test,rf_predictions))
print("r2_score",r2_score(y_test,rf_predictions))

FINAL MODEL
-----
Mean Absolute Error: 5.356798154843619
Mean Squared Error: 82.15186160414899
Root Mean Squared Error: 9.063766413812196
Explained Variance Score: 0.8841543358112451
r2_score: 0.8840153308822885

```

```

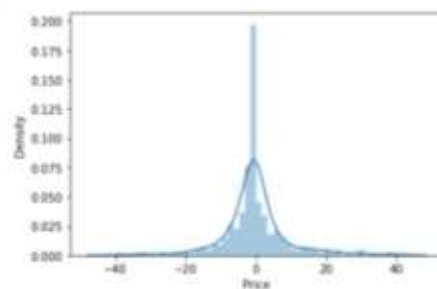
In [81]: sns.distplot(y_test,rf_predictions)

```

```

Out[81]: <AxesSubplot:xlabel='Price', ylabel='Density'>

```



```

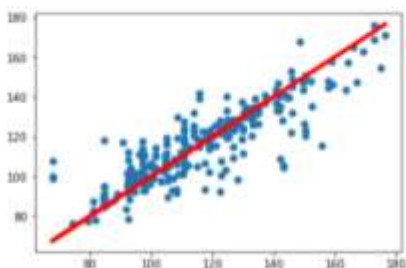
In [82]: plt.scatter(y_test,rf_predictions)
plt.plot(y_test,y_test,linewidth=4,color='red')

```

```

Out[82]: [matplotlib.lines.Line2D at 0x1fb3db13480]

```



## 7.SAVING THE MODEL

```

In [83]: import joblib
joblib.dump(rf,"final_model.pkl")

```

```

Out[83]: ['final_model.pkl']

```

In our project we had implemented various Machine Learning Algorithms such as Linear Regression, Decision Tree Regression, Random Forest Regression and compared the accuracy of results based on our test data set. Based on the various accuracy levels we find that Random Forest Regression gives the highest accuracy i.e. 80%. Therefore we selected Random Forest Regression and created User Interface based on it.

## Visualizations

After cleaning the data, we can visualize data and better understand the relationships between different variables. There are many more visualizations that you can do to learn more about your dataset, like scatterplots, histograms, boxplots.

➤ The Analysis of flight Price. The price of the ticket

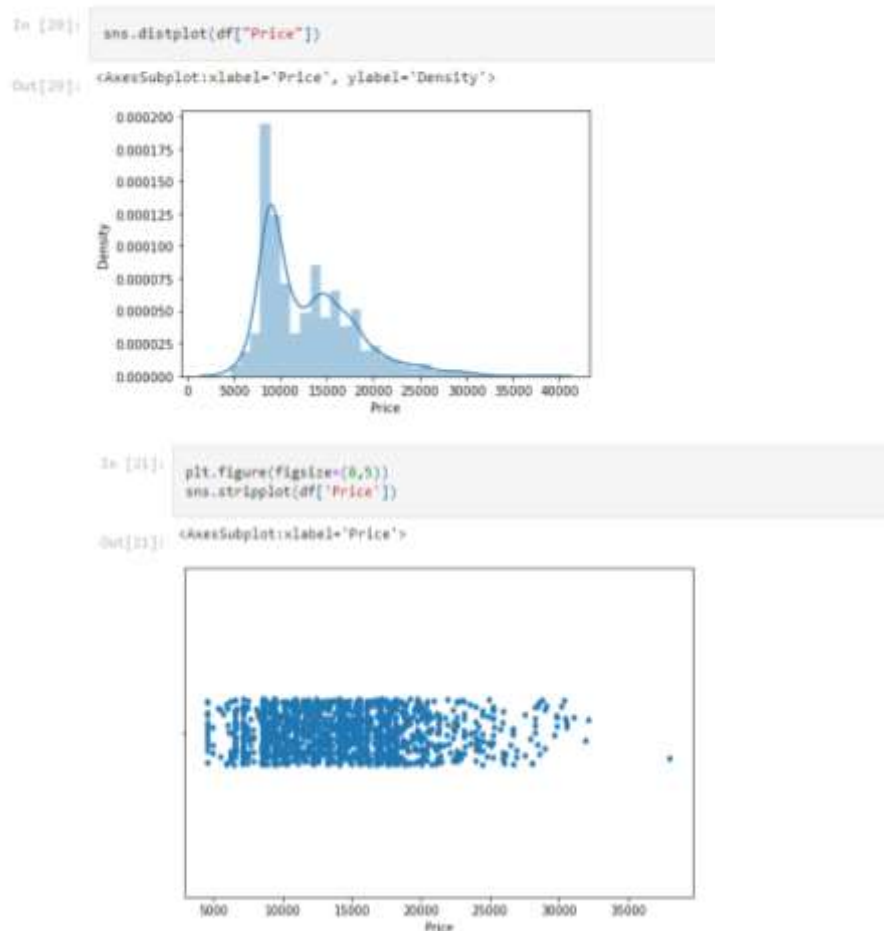


Fig. 1&2: The Analysis price of the ticket

The source from which the service begins

- In observation from New Delhi has more Flight service held compare with all other cities.
- Kochi and Pune have low flight service.



Fig. 3: The Analysis of source from which the service begins

- In observation from fig:4 Mumbai has more Flight service end point held compare with all other cities.
- Goa and Bangalore have low flight service.

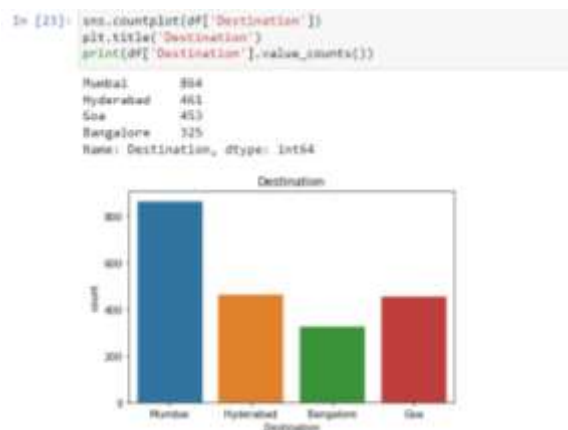


Fig. 4: The Analysis of The destination where the service ends.

- For one stop flight are higher compare with others.
- Only 20% flight are non-stop and 2 stop service

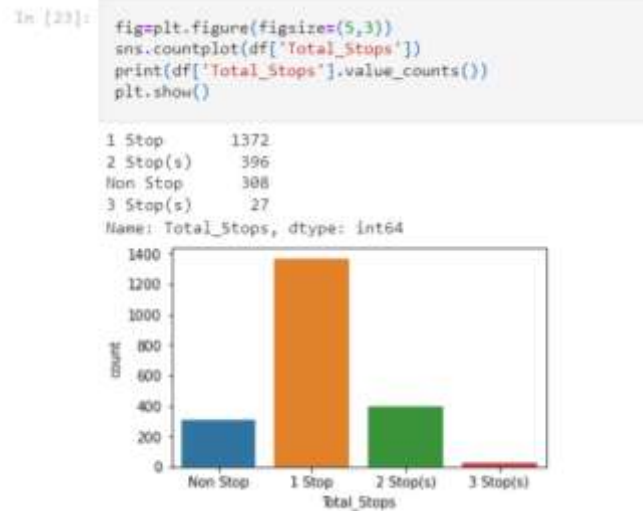


Fig. 5: The Analysis of Total stops between the source and destination.

- Vistara and IndiGo are the higher service provide.
- Air Asia has 2.3% very low service compare with other.

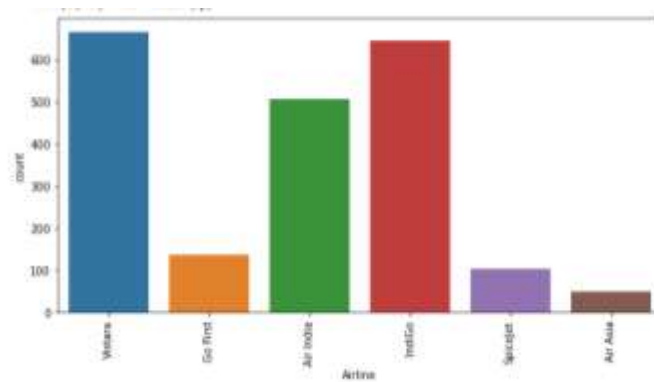


Fig. 6: The Analysis of The name of the airline.

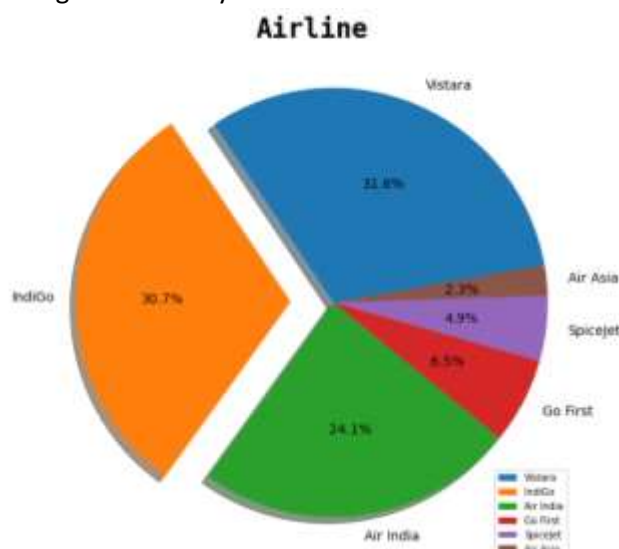


Fig.7: pie chart for The name of the airline.

➤ The source from which the service begin

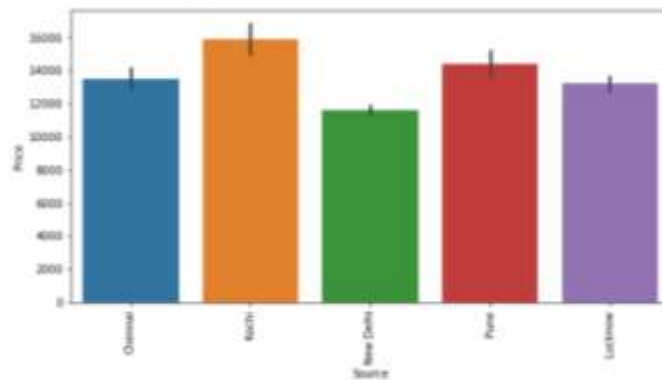


Fig. 8: The Analysis between Source and Price

- Total stops between the source and destination
- 2 stops and 3 stops have high flight price rate compare with others.
  - 1 stop and Non stop have low flight price rate.

```
In [30]: plt.figure(figsize=(10,5))
sns.barplot(df['Total_Stops'],df['Price'])
```

Out[30]: <AxesSubplot: xlabel='Total\_Stops', ylabel='Price'>

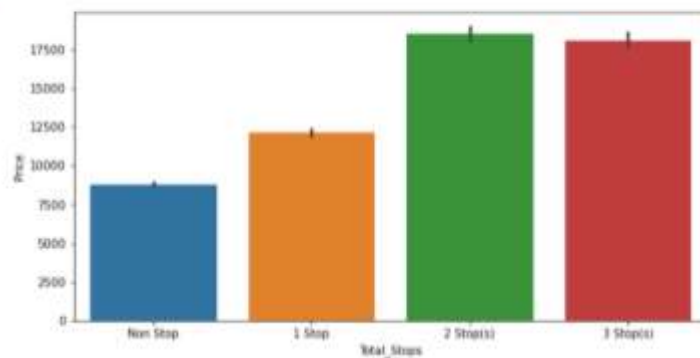


Fig. 9: The Analysis between Total stop and Price

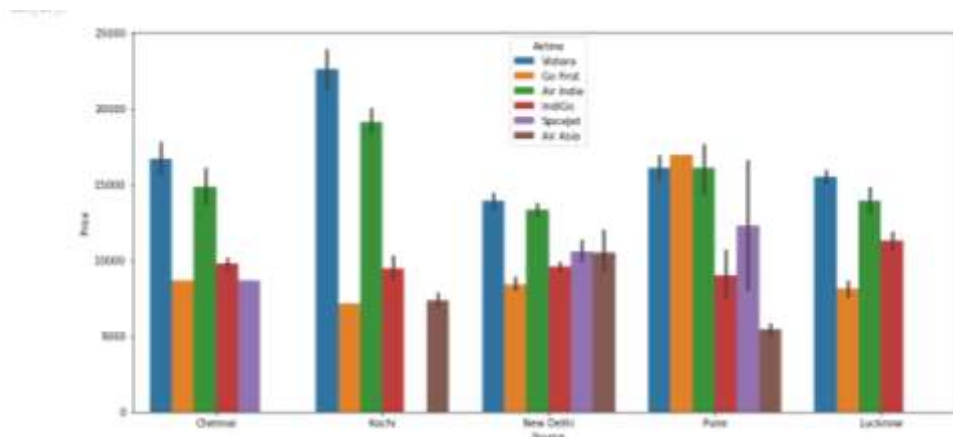


Fig. 10: The Analysis between source and Price.

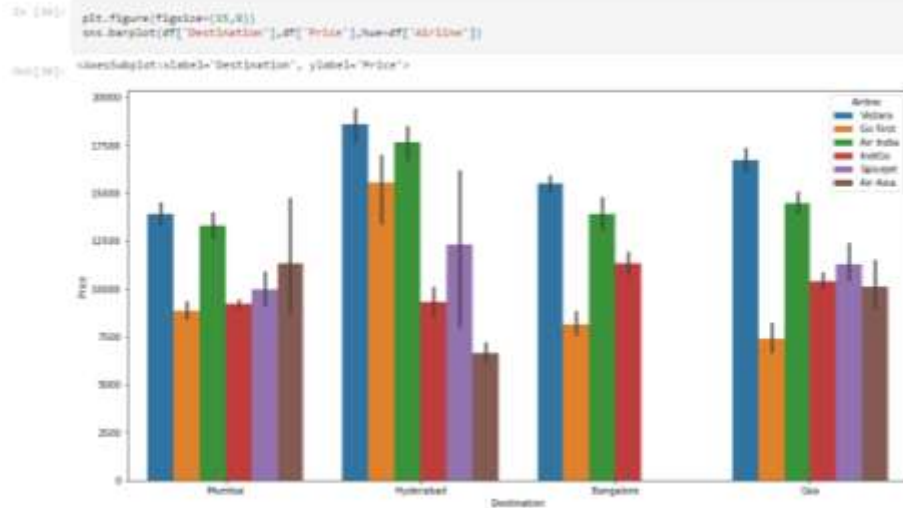


Fig. 11: The Analysis between Destination and Price.

- 24 Sunday and 25 Monday have very high price compare other day.
- Vistara and Air india are high price rated flight on everyday.
- Indigo and Go First are low price flight .

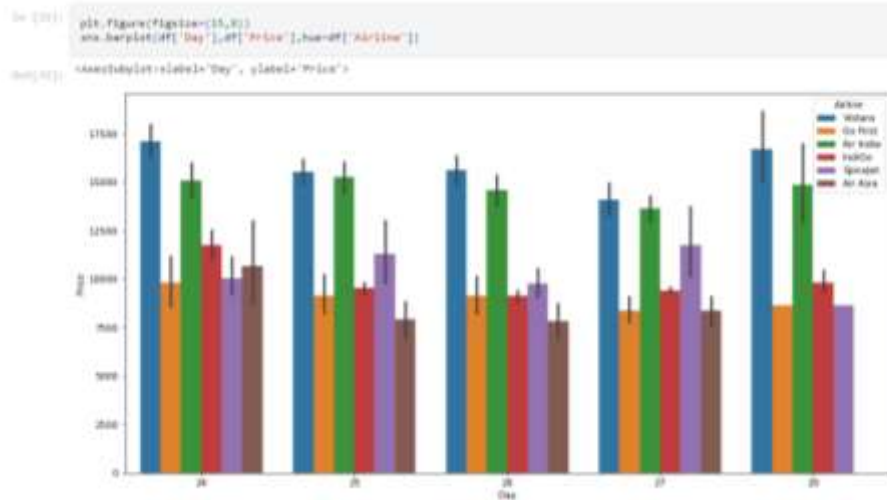


Fig. 12: The Analysis between Day and Price

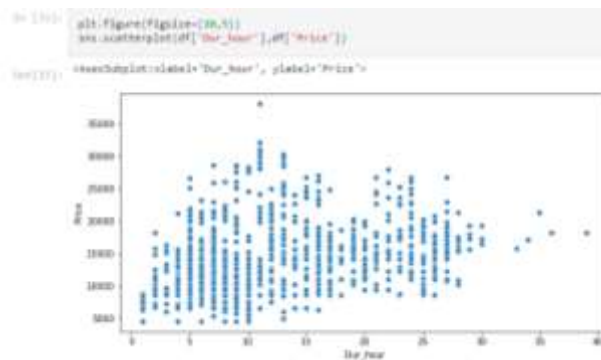


Fig. 13: The Analysis between Duration and Price.

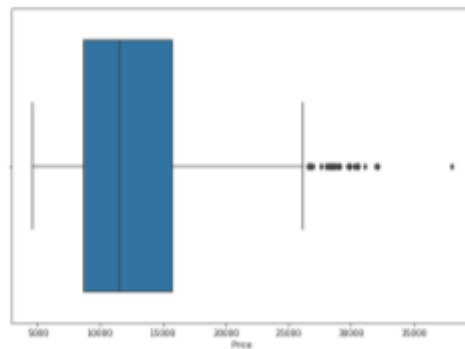
➤ The Analysis between Airline and Price.

The observation is Vistara is having the highest price of airline compared to other airlines. The second highest is Air India airline and all other airline prices are almost similar. Through this analysis the advantage is that the customers has the option to choose the various airlines with different price and comfort with their budget to travel the new places and explore the world.



Fig.14,15&16: The Analysis between Airline and Price.

## ➤ Checking the outliers



## ➤ Checking and Removing Skewness

```
In [40]: df.skew()
```

```
Out[40]: Airline    -0.125018
Source    -0.787096
Destination -0.453672
Day        0.792582
Month      0.800000
Dep_time   0.835357
Arr_time   -0.126795
Dur_hour   0.752992
Dur_min    0.175493
Total_stops 1.456016
Price      1.874889
dtype: float64
```

Using Square root transformation method to reduce Skewness

```
In [41]: #Applying the square root for data
df1 = np.sqrt(df)
df1.head(5)
```

```
Out[41]:
```

|    | Airline | Source | Destination | Day | Month   | Dep_time | Arr_time | Dur_hour | Dur_min | Total_stops | Price    |
|----|---------|--------|-------------|-----|---------|----------|----------|----------|---------|-------------|----------|
| 0  | 1.23608 | 0.0    | 1.73205     | 0.0 | 1.94771 | 7.21115  | 7.87234  | 1.08080  | 7.41616 | 1.73205     | 89.08103 |
| 1  | 1.23608 | 0.0    | 1.73205     | 0.0 | 1.94771 | 8.88794  | 9.48630  | 1.41616  | 0.00000 | 1.73205     | 95.08103 |
| 2  | 1.41421 | 0.0    | 1.73205     | 0.0 | 1.94771 | 8.42381  | 9.04876  | 1.41421  | 0.00000 | 1.73205     | 89.08103 |
| 3  | 1.23608 | 0.0    | 1.73205     | 0.0 | 1.94771 | 11.08338 | 13.07897 | 1.41421  | 2.29608 | 1.73205     | 89.08103 |
| 4  | 1.80000 | 0.0    | 1.73205     | 0.0 | 1.94771 | 11.04891 | 13.82150 | 1.00000  | 6.78236 | 1.73205     | 91.02084 |
| 5  | 1.73205 | 0.0    | 1.73205     | 0.0 | 1.94771 | 2.00000  | 4.24241  | 1.00000  | 7.07108 | 1.73205     | 91.02084 |
| 6  | 1.73205 | 0.0    | 1.73205     | 0.0 | 1.94771 | 8.00000  | 8.85935  | 1.00000  | 7.07108 | 1.73205     | 91.02084 |
| 7  | 1.73205 | 0.0    | 1.73205     | 0.0 | 1.94771 | 8.84054  | 9.00000  | 1.00000  | 7.07108 | 1.73205     | 91.02084 |
| 8  | 1.73205 | 0.0    | 1.73205     | 0.0 | 1.94771 | 8.17494  | 8.81160  | 1.00000  | 7.07108 | 1.73205     | 91.02084 |
| 9  | 1.73205 | 0.0    | 1.73205     | 0.0 | 1.94771 | 11.25717 | 12.88340 | 1.00000  | 7.07108 | 1.73205     | 91.02084 |
| 10 | 1.73205 | 0.0    | 1.73205     | 0.0 | 1.94771 | 11.81625 | 13.88920 | 1.00000  | 7.07108 | 1.73205     | 91.02084 |
| 11 | 1.80000 | 0.0    | 1.73205     | 0.0 | 1.94771 | 4.24241  | 5.36734  | 1.00000  | 7.41616 | 1.73205     | 91.02084 |
| 12 | 1.73205 | 0.0    | 1.73205     | 0.0 | 1.94771 | 11.04749 | 13.03640 | 1.41421  | 2.29608 | 1.73205     | 91.02084 |
| 13 | 1.73205 | 0.0    | 1.73205     | 0.0 | 1.94771 | 7.00000  | 7.87234  | 1.41616  | 0.00000 | 1.73205     | 91.02084 |
| 14 | 1.73205 | 0.0    | 1.73205     | 0.0 | 1.94771 | 11.78628 | 13.26489 | 1.73205  | 7.07108 | 1.80000     | 91.02084 |

```
In [42]: df1.skew()
```

```
Out[42]: Airline    -0.788851
Source    -1.386350
Destination -1.150431
Day        0.726101
Month      0.800000
Dep_time   -0.482686
Arr_time   -0.683590
Dur_hour    0.191832
Dur_min    -0.691545
Total_stops 0.949504
Price      0.660509
dtype: float64
```

Skewness is reduced



## Observations

- Airfare varies according to the day of the week of travel. It is higher for weekends and Monday and slightly lower for the other days.
- The airfare varies depending on the time of departure, making timeslot used in analysis an important parameter.
- The time duration plays the important role for making the decision to board the flight with best price. With the limited amount of time the best price can be chosen by the passengers. Everybody can afford the flight ticket with best price and best offers
- That travellers can see the highest to lowest prices to know the price differences in each airline. Passengers can share the review of the flight on the each airline websites, so that the other passengers gets the benefit out of it.
- That duration (or distance) plays a major role in affecting air ticket prices, but we see no such pattern here, as there must be pattern here and other significant factors affecting air fare like type of airline, destination of flight, date of journey of flight (higher if collides with a public holiday).
- Through this analysis the advantage is that passengers might want to reach the destination sooner, so the duration plays the important role to reach sooner with good amount of price.
- Through this analysis the advantage is that the Journey Day plays the important role because the flight charges for weekdays in might be lesser and weekends the price might be higher based on the offers and airlines chosen to board the flight.

## **ANALYSIS ON THE DATA.**

1. Do airfares change frequently?

Ans: Yes, Based on Season and Demand of the flight Tickets Price changes regularly.

2. Do they move in small increments or in large jumps?

Ans. large jumps, If you buy a ticket before few of date of journey you will get arround 44% high price.

3. Do they tend to go up or down over time?

Ans. Price of the flight goes up only down over time but its not reduces the p rice over time

4. What is the best time to buy so that the consumer can save the most by taki ng the least risk?

Ans. Consumer should buy the ticket before 30 days of date of joureney so th ey can save up to 45% to 50%

# CONCLUSION

## Key Findings and Conclusions of the Study

From this dataset I get to know that each feature play a very import role to understand the data. Data format plays a very important role in the visualization and Appling the models and algorithms .Importance of removing the skewness and outlier.

There are many systems which uses different machine learning algorithms such as Linear Regression (LR), Decision Tree, Support Vector Machine (SVM), Random Forest Algorithm, etc for predicting the price for flight ticket. In this ML based system, we are using Random Forest Algorithm which gives more accuracy in predicting the airfare. Considering the features such as departure time, the number of days left for departure and time of the day it will give the best time to buy the ticket .This system also helps the buyer to buy the flight ticket at lower price. It is easy to use and it gives more accuracy in prediction. It requires less time for prediction and it helps in reduction of over fitting. Travellers can save money if they choose to buy the ticket when its price is the lowest. It gives the customer the best time to buy a flight ticket for the desired destination and a period

Evaluating the algorithmic rule, a dataset is collected, pre-processed, performed data modelling and studied a value difference for the number of restricted days by the passengers for travelling .Machine Learning algorithms with square measure for forecasting the accurate fare of airlines and it gives accurate value of plane price ticket at limited and highest value. Information is collected from yatra websites that sell the flight tickets therefore restricting data which are often accessed .The results obtained by the random forest and decision tree algorithm has better accuracy, but best accuracy is predicted by random forest algorithm as shown is the above analysis. Accuracy of the model is also forecasted by the R-squared value.

## **Learning Outcomes of the Study in respect of Data Science**

Learnt how to process the large number of data. Tried and learnt more about distribution of the data. The power of visualization is helpful for the understanding of data into the graphical representation its help me to understand that what data is trying to say, Data cleaning is one of the most important step to remove missing value or null value fill it by mean median or by mode or by 0. Setting a good parameters is more important for the model accuracy. Finding a best random state played a vital roll in finding a better model.

Flight ticket prices can be something hard to guess, today we might see a price, check out the price of the same flight tomorrow, it will be a different story. We might have often heard travellers saying that flight ticket prices are so unpredictable. As data scientists, we are going to prove that given the right data anything can be predicted. So the collected train data should be accurate if not it may result in wrong prediction. And also it is necessary to update the train data time to time for best results.

### **Limitations of this work and Scope for Future Work**

The techniques to increase the speed of the model need to be constructed. In Upcoming days when huge amount of information is accessed as in detailed information in the dataset, the expected results in future are highly correct. For further research anyone desire to expand upon it ought to request different sources of historical data or be a lot of organized in collection knowledge manually over amount of your time to boot, a lot of different combination of plane are

going to be traversed. There is whole possibility that planes differ their execution ideas consisting characteristics of the plane. At last, it is curious to match our model accuracy with that of the business models accuracy offered nowadays.

For the prediction of the ticket prices perfectly different prediction models are tested for the better prediction accuracy. As the pricing models of the company are developed in order to maximize the revenue management .With the help of our project the travellers can find out the right time to buy their tickets at the lowest cost and also can plan accordingly. So to get result with maximum accuracy regression analysis is used. From the studies, the feature that influences the price ticket are to be considered. In future the details about number of available seats can improve the performance of the model

We might have often heard travellers saying that flight ticket prices are so unpredictable. As data scientists, we are going to prove that given the right data anything can be predicted. So the collected train data should be accurate if not it may result in wrong prediction. And also it is necessary to update the train data time to time for best results