

**NAME OF THE PROJECT**

# **PROJECT REPORT ON FLIGHT PRICE PREDICTION USING ML TECHNIQUES**

**SUBMITTED BY:**

**MS. YASHSHREE BAVISKAR**

**FLIPROBO SME:**

**MOHD KASHIF**

## ACKNOWLEDGMENT

I would like to express my special gratitude to “Flip Robo” team, who has given me this opportunity to deal with a beautiful dataset and it has helped me to improve my analyzation skills. And I want to express my huge gratitude to Mr. Mohd Kashif (SME Flip Robo), he is the person who has helped me to get out of all the difficulties I faced while doing the project.

A huge thanks to “Data trained” who are the reason behind my Internship at Fliprobo. Last but not least my parents who have been my backbone in every step of my life.

References use in this project:

1. SCIKIT Learn Library Documentation
2. Blogs from towardsdatascience, Analytics Vidya, Medium
3. Andrew Ng Notes on Machine Learning (GitHub)
4. Data Science Projects with Python Second Edition by Packt
5. Hands on Machine learning with scikit learn and tensor flow by Aurelien Geron
6. B. Smith, J. Leimkuhler, R. Darrow, and Samuels, “Yield management at American airlines, “Interfaces, vol. 22, pp. 8–31, 1992
7. William Groves, Maria Gini, “An agent for optimizing airline ticket purchasing”, in international conference on Autonomous agents and multi-agent systems, International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC (2013)
8. Chen, Y., Cao, J., Feng, S., Tan, Y., 2015. An ensemble learning based approach for building airfare forecast service. In: 2015 IEEE International Conference on Big Data (Big Data), Santa Clara, CA, 2015, pp. 964-969.
9. Yeamduan Narangajavana, Fernando.J. Garrigos-Simon, Javier Sanchez García, Santiago Forgas-Coll, “Prices, prices and prices: A study in the airline sector”, Tourism Manage., 41 (2014), pp. 28-42
10. Bo An, Haipeng Chen, Noseong Park, V.S. Subrahmanian MAP: Frequency-Based Maximization of Airline Profits based on an Ensemble Forecasting Approach Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16), ACM, New York, NY, USA (2016), pp. 421-430
11. R. Ren, Y. Yang, and S. Yuan, “Prediction of airline ticket price,” University of Stanford, 2014.

12. T. Janssen, T. Dijkstra, S. Abbas, and A. C. van Riel, "A linear quantile mixed regression model for prediction of airline ticket prices," Radboud University, 2014.
13. K. Tziridis, T. Kalampokas, G. A. Papakostas, and K. I. Diamantaras, "Airfare prices prediction using machine learning techniques," in the 25th IEEE European signal processing conference, 2017, pp. 1036–1039.
14. C. Koopmans and R. Lieshout, "Airline cost changes: To what extent are they passed through to the passenger?" *Journal of Air Transport Management*, vol. 53, pp. 1–11, 2016.
15. G. Francis, A. Fidato, and I. Humphreys, "Airport–airline interaction: the impact of low-cost carriers on two European airports," *Journal of Air Transport Management*, vol. 9, no. 4, pp. 267–273, 2003.
16. Boruah A., Baruah K., Das B., Das M.J., Gohain N.B. (2019) "A Bayesian Approach for Flight Fare Prediction Based on Kalman Filter," [https://doi.org/10.1007/978-981-13-0224-4\\_18](https://doi.org/10.1007/978-981-13-0224-4_18)
17. G.A. Papakostas, K.I. Diamantaras and T. Papadimitriou, "Parallel pattern classification utilizing GPU-Based kernelized slackmin algorithm," doi:10.1016/j.jpdc.2016.09.001
18. T. Janssen, "A linear quantile mixed regression model for prediction of airline ticket prices," Bachelor Thesis, Radboud University, 2014.
19. Stacey Mumbower, Laurie A. Garrow, Matthew J. Higgins "Estimating flight-level price elasticities using online airline data: a first step toward *integrating* pricing, demand, and revenue optimization", *Transportation Res. Part A: Policy Practice*, 66 (2014), pp. 196-212

# CHAP 1. INTRODUCTION

## 1.1 Business Problem Framing

The Airline Companies is considered as one of the most enlightened industries using complex methods and complex strategies to allocate airline prices in a dynamic fashion. These industries are trying to keep their all-inclusive revenue as high as possible and boost their profit. 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.

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

- Time of purchase patterns (making sure last-minute purchases are expensive)
- 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, this project involves collection of data for flight fares with other features and building a model to predict fares of flights.

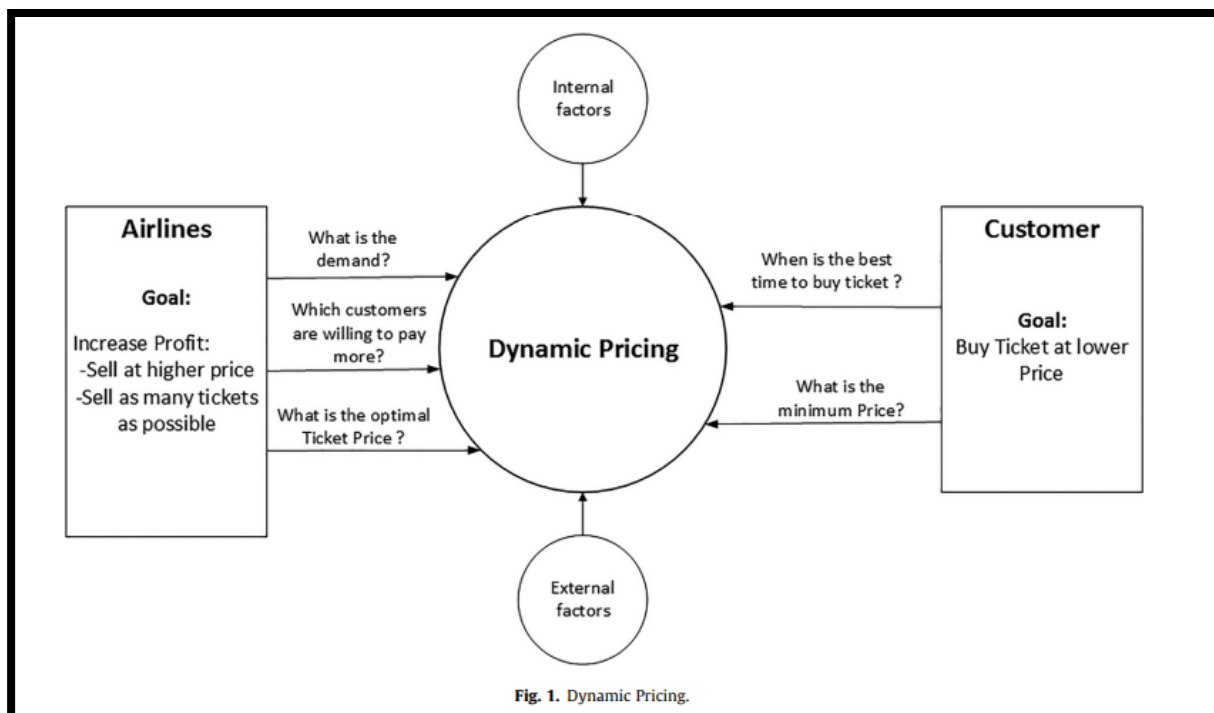
## 1.2 CONCEPTUAL BACKGROUND OF THE DOMAIN PROBLEM

A report says India's affable aeronautics industry is on a high development movement. *India is the third-biggest avionics showcase in 2020 and the biggest by 2030.* Indian air traffic is normal to cross the quantity of 100 million travellers by 2017, whereas there were just 81 million passengers in 2015. Agreeing to Google,

the expression "Cheap Air Tickets" is most sought in India. At the point when the white-collar class of India is presented to air travel, buyers searching at modest costs.

Any individual who has booked a flight ticket previously knows how dynamically costs change. Aircraft uses advanced strategies called Revenue Management to execute a distinctive valuing strategy [6]. The least expensive accessible ticket changes over a period the cost of a ticket might be high or low. This valuing method naturally modifies the toll as per the time like morning, afternoon or night. Cost may likewise change with the seasons like winter, summer and celebration seasons. The extreme goal of the carrier is to build its income yet on the opposite side purchaser is searching at the least expensive cost. Purchasers generally endeavour to purchase the ticket in advance to the take-off day.

From the customer point of view, determining the minimum price or the best time to buy a ticket is the key issue. The conception of "tickets bought in advance are cheaper" is no longer working (**William Groves and Maria Gini, 2013**) [7]. It is possible that customers who bought a ticket earlier pay more than those who bought the same ticket later. Moreover, early purchasing implies a risk of commitment to a specific schedule that may need to be changed usually for a fee. Most of the studies performed on the customer side focus on the problem of predicting optimal ticket purchase time using statistical methods. As noted by **Y. Chen et al. (2015)** [8], predicting the actual ticket price is a more difficult task than predicting an optimal ticket purchase time due to various reasons: absence of



enough datasets, external factors influencing ticket prices, dynamic behaviour of ticket pricing, competition among airlines, proprietary nature of airlines ticket pricing policies etc.

Early prediction of the demand along a given route could help an airline company pre-plan the flights and determine appropriate pricing for the route. Existing demand prediction models generally try to predict passenger demand for a single flight/route and market share of an individual airline. Price discrimination allows an airline company to categorize customers based on their willingness to pay and thus charge them different prices. Customers could be categorized into different groups based on various criteria such as business vs leisure, tourist vs normal traveller, profession etc. For example, business customers are willing to pay more as compared to leisure customers as they rather focus on service quality than price.

In a less competitive market, the market power of a given airline is stronger, and thus, it is more likely to engage in price discrimination. On the other hand, *the higher the level of competition, the weaker of the market power of an airline, and then the less likely the chance of the airline fare increases.*

### 1.3 REVIEW OF LITERATURE

On the airlines side, the main goal is increasing revenue and maximizing profit. According to **(Narangajavana et al., 2014) [9]**, airlines utilize various kinds of pricing strategies to determine optimal ticket prices: long-term pricing policies, yield pricing which describes the impact of production conditions on ticket prices, and dynamic pricing which is mainly associated with dynamic adjustment of ticket prices in response to various influencing factors.

Among the recent work performed on route demand and market share prediction is the study done by **(Bo An et al., 2016) [10]**. The authors proposed a **data mining technique designed for Maximizing Airline Profits (MAP) through prediction of total route demand and market share of an individual airline**. Unlike most other works, this work considers a broad set of routes (around 700 routes) across 13 airlines operating in those routes. The training dataset spans 10 years (40 quarters) while the testing set includes the first quarter of 2015 (a total of 9100 predictions). However, the prediction is performed quarterly and not for a short period of time which might not consider dynamic demand changes. Moreover, the routes considered are only national routes in the US.

**Ren et al. [11]** proposed using LR, Naive Bayes, Soft-max regression, and SVMs to build a prediction model and classify the ticket price into five bins (60% to 80%, 80% to 100%, 100% to 120%, and etc.) to compare the relative values with the overall average price. More than nine thousand data points, including six features (e.g., the departure week begin, price quote date, the number of stops in

the itinerary, etc.), were used to build the models. The authors reported the best training error rate close to 22.9% using LR model. **Their SVM regression model failed to produce a satisfying result.** Instead, an SVM classification model was used to classify the prices into either “higher” or “lower” than the average.

In [12], four LR models were compared to obtain the best fit model, which aims to provide an unbiased information to the passenger whether to buy the ticket or wait longer for a better price. The authors suggested using linear quantile mixed models to predict the lowest ticket prices, which are called the “real bargains”. However, this work is limited to only one class of tickets, economy, and only on one direction single leg flights from San Francisco Airport to John F. Kennedy Airport.

**Tziridis et al. [13]** applied eight machine learning models, which included ANNs, RF, SVM, and LR, to predict tickets prices and compared their performance. The best regression model achieved an accuracy of 88%. **In their comparison, Bagging Regression Tree is identified as the best model, which is robust and not affected by using different input feature sets.**

Macroeconomic data, such as crude oil price and Consumer Price Index (CPI), can also be utilized to uncover the hidden trend in airline fares. **Fuel costs can take up to 50% of the total operating cost of an airline [14].** Hence, the level of crude oil price plays an essential rule of formulating the airline’s pricing strategy. *It is a common practice for airlines to pass the cost of aviation fuel to the customer by adjusting the fare to compensate for the fluctuation of crude oil price.*

The emergence of Low-Cost Carrier (LCC) has revolutionized the entire operating model of the airline industry. *The presence of LCC in a market has had a substantial impact on the total passenger volume and the air ticket price [15].*

In detail monitoring, the passenger gets an approximation of plane price with date to choose the best blend of date and price. *The price for weekend on Sunday is not possible to calculate in this presented model,* as weekend on Sundays the most accidental price difference compared to other days in the week and needs more elements, nonlinear model for successful forecast which will be the upcoming range of study to be done for this presented technique [16]. To forecast the mean plane ticket amount on the business area, machine learning support was evolved. Selecting feature techniques authors have presented model to forecast the mean flight amount with R squared score of 80% accuracy.

The accuracy of logistic regression model is up to 70-75%. *The conclusion of the given model is that most of the plane ticket price vary from day to day.* Authors have reported that the ticket price is high for a certain period and then it gradually decreases to a certain level. **When the flight is at a difference of 2-3 days’ time the ticket price starts increasing again [17].**



Janssen [18] built up an expectation model utilizing the Linear Quantile Blended Regression strategy for San Francisco to New York course with existing every day airfares given by [www.infare.com](http://www.infare.com). **The model utilized two highlights including the number of days left until the take-off date and whether the flight date is at the end of the week or weekday.** The model predicts airfare well for the days that are a long way from the take-off date, anyway for a considerable length of time close the take-off date, the expectation isn't compelling.

***Business class flights are more inelastic as compared to leisure class as business customers have less flexibility to change or cancel their travel date*** (Mumbower et al., 2014) [19]. *In contrast, short distance flights are more elastic (more price sensitive) than long distance flights* because of the availability of other travel options (e.g., bus, train, car etc.). Airlines use price elasticity information to determine when to increase ticket prices or when to launch promotions so that the overall demand is increased

## 1.4 MOTIVATION FOR THE PROBLEM UNDERTAKEN

The project was the first provided to me by Flip Robo Technologies as a part of the internship programme. The exposure to real world data and the opportunity to deploy my skillset in solving a real time problem has been the primary motivation.

Early prediction of the demand along a given route could help an airline company pre-plan the flights and determine appropriate pricing for the route. In addition, competition between airlines makes the task of determining optimal pricing is hard for everyone. ***So prime motive is to build flight price predication system based on short range timeframe (7-14 days) data available prior to actual take-off date.***



# CHAP 2 ANALYTICAL PROBLEM FRAMING

## 1. MATHEMATICAL / ANALYTICAL MODELLING OF THE PROBLEM

First phase of problem modelling involves data scraping of flights from internet. For that purpose, flight data is scrap from [www.yatra.com](http://www.yatra.com) for timeframe of 21 June 2022 to 5th July 2022. Data is scrape for flights on route of New Delhi to Mumbai. Data is scrap for Economy class, Premium Economy class & Business class flights. Next phase is data cleaning & pre-processing for building ML Model. Our objective is to predict flight prices which can be resolve by use of regression-based algorithm. Further Hyperparameter tuning performed to build more accurate model out of best model.

## 2. Data Sources and their formats

Data is collected from [www.yatra.com](http://www.yatra.com) for timeframe of 21st June 2022 to 5th July 2022 using selenium and saved in CSV file. Data is scrape for flights on route of New Delhi to Mumbai. Data is scrap for Economy class, Premium Economy class & Business class flights. Around 3000 flights details are collected for this project.

```
print('No of Rows:',df.shape[0])
print('No of Columns:',df.shape[1])
pd.set_option('display.max_columns', None) # This will enable us to see truncated columns
df.head()
```

No of Rows: 2671  
No of Columns: 12

	Unnamed: 0	Airline	Class	Aeroplane	Date	Departure_Time	Arrival_Time	Source	Destination	Stops	Duration	Price
0	0	Vistara	Economy Class	UK-981	Tue, 21 Jun 2022	21:30	23:30	New Delhi	Mumbai	Non Stop	2h 00m	8578
1	1	Vistara	Economy Class	UK-993	Tue, 21 Jun 2022	12:50	15:00	New Delhi	Mumbai	Non Stop	2h 10m	8578
2	2	Vistara	Economy Class	UK-951	Tue, 21 Jun 2022	14:20	16:30	New Delhi	Mumbai	Non Stop	2h 10m	8578
3	3	Vistara	Economy Class	UK-933	Tue, 21 Jun 2022	15:30	17:40	New Delhi	Mumbai	Non Stop	2h 10m	8578
4	4	Vistara	Economy Class	UK-985	Tue, 21 Jun 2022	19:45	21:55	New Delhi	Mumbai	Non Stop	2h 10m	8578

Unnecessary column of index name as 'Unnamed: 0' is drop out. There are 11 features in dataset including target feature 'Price'. The data types of different features are as shown below:

```
# Lets sort columns by their datatype
df.columns.to_series().groupby(df.dtypes).groups

{int64: ['Price'], object: ['Airline', 'Class', 'Aeroplane', 'Date', 'Departure_Time', 'Arrival_Time', 'Source', 'Destination', 'Stops', 'Duration']}
```

### 3. Data Pre-processing

The dataset is large and it may contain some data error. In order to reach clean, error free data some data cleaning & data pre-processing performed data.

- Data Integrity check –

No missing values present in dataset.

duplicate entries checked and dropped the duplicate value

```
df.duplicated().sum()

110

df=df.drop_duplicates(subset=None,inplace=False)

df.shape

(2561, 11)
```

#### Conversion of Duration column from hr & Minutes format into Minutes –

By default, Duration of flights are given in format of [(hh) hours: (mm)minute] which need to convert into uniform unit of time. Here we have written code to convert duration in terms of minute. For example,

```
df['Duration'] = df['Duration'].map(lambda x : x.replace('05m','5m'))

# Conversion of Duration column from hr & Minutes format to Minutes
df['Duration'] = df['Duration'].str.replace('h','*60').str.replace(' ','+').str.replace('m','*1').apply(eval)

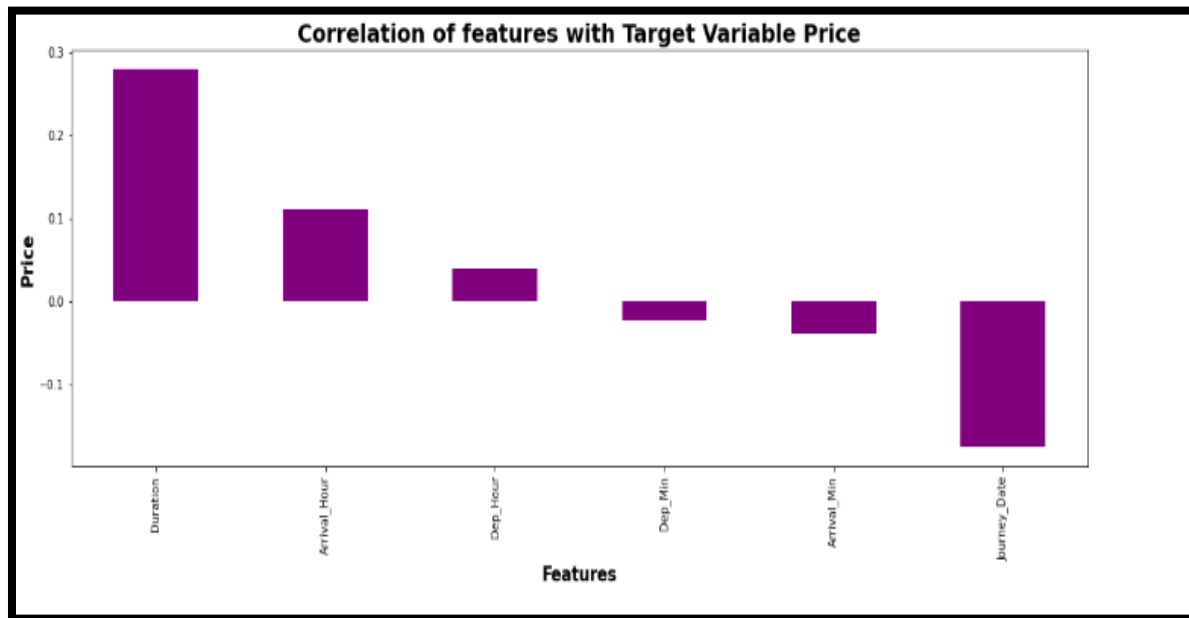
# convert this column into a numeric datatypes
df['Duration']= pd.to_numeric(df['Duration'])
```

- Create new column for day & date –

```
df['Week_Day'] = df['Date'].map(lambda x :x[:3])  
  
df['Journey_Day'] = df['Date'].map(lambda x :x[4:])  
  
df['Journey_Month'] = df['Journey_Day'].map(lambda x :x.split()[1])  
  
df['Journey_Date'] = df['Journey_Day'].map(lambda x :x.split()[0])  
  
df['Journey_Year'] = df['Journey_Day'].map(lambda x :x.split()[2])
```

New column for 'Day' & 'Date' is extracted from Date column.

## 4. DATA INPUTS- LOGIC- OUTPUT RELATIONSHIPS



Correlation heatmap is plotted to gain understanding of relationship between target features & independent features. We can see that class feature is correlated for more than -0.6 with target variable Price. Remaining feature are poorly correlated with target variable price.

## 5. Hardware & Software Requirements with Tool Used

1. Hardware Used -
  1. Processor — Intel i7 processor with 2.4GHZ
  2. RAM — 4 GB
  3. GPU — 2GB N Series Graphics card
2. Software utilised -
  1. Anaconda – Jupyter Notebook
  2. Selenium – Webscraping
  3. Google Colab – for Hyper parameter tuning

Libraries Used – General library for data wrangling & visualization

```
#Importing required libraries
import pandas as pd
import selenium
from selenium import webdriver
from selenium.common.exceptions import NoSuchElementException, StaleElementReferenceException
import time
import warnings
warnings.filterwarnings('ignore')
```

Libraries used for web scraping data from e-commerce website are

```
import pandas as pd # for data wrangling purpose
import numpy as np # Basic computation library
import seaborn as sns # For Visualization
import matplotlib.pyplot as plt # plotting package
%matplotlib inline
import warnings # Filtering warnings
warnings.filterwarnings('ignore')
```

Libraries used for machine learning model building

```
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.tree import DecisionTreeRegressor
from sklearn.ensemble import ExtraTreesRegressor
from xgboost import XGBRegressor
```

## Chap. 3 Models Development & Evaluation

### ❖ Identification Of Possible Problem-Solving Approaches (Methods)

First part of problem solving is to scrap data from [www.yatra.com](http://www.yatra.com) website which we already done. Next part of problem solving is building machine learning model to predict flight price. This problem can be solve using regression-based machine learning algorithm like linear regression. For that purpose, first task is to convert categorical variable into numerical features. Once data encoding is done then data is scaled using standard scalar. Final model is built over this scaled data. For building ML model before implementing regression algorithm, data is split in training & test data using `train_test_split` from `model_selection` module of `sklearn` library. After that model is train with various regression algorithm and 5-fold cross validation is performed. Further Hyperparameter tuning performed to build more accurate model out of best model.

### ❖ Testing of Identified Approaches (Algorithms)

Phase 1 Web Scraping Strategy employed in this project as follow:

1. Selenium will be used for web scraping data from [www.yatra.com](http://www.yatra.com)
2. Flights on route of New Delhi to Mumbai in duration of 23 Jan 2022 to 4 Feb 2022.
3. Data is scrap in three parts:
  - Economy class flight price extraction
  - Business class flight price extraction
  - Premium Economy class price extraction
4. Selecting features to be scrap from website.
5. In next part web scraping code executed for above mention details.
6. Exporting final data in Excel file.



The different regression algorithm used in this project to build ML model are as below:

- ❖ Linear Regression
- ❖ Random Forest Regressor
- ❖ Decision Tree Regressor
- ❖ XGB Regressor
- ❖ Extra Tree Regressor
- ❖ Gradient Boosting Regressor

## ❖ KEY METRICS FOR SUCCESS IN SOLVING PROBLEM UNDER CONSIDERATION

Following metrics used for evaluation:

1. Mean absolute error which gives magnitude of difference between the prediction of an observation and the true value of that observation.
2. Root mean square error is one of the most commonly used measures for evaluating the quality of predictions.
3. R2 score which tells us how accurate our model predict result, is going to important evaluation criteria along with Cross validation score.

## 4. RUN AND EVALUATE SELECTED MODELS

### 1. Linear Regression:

```
X_train, X_test, Y_train, Y_test = train_test_split(X_scale, Y, random_state= 287, test_size=0.25)
lin_reg= LinearRegression()
lin_reg.fit(X_train, Y_train)
y_pred = lin_reg.predict(X_test)
print('\033[1m+ 'Error :'+ '\033[0m')
print('Mean absolute error :', mean_absolute_error(Y_test,y_pred))
print('Mean squared error :', mean_squared_error(Y_test, y_pred))
print('Root Mean squared error :', np.sqrt(mean_squared_error(Y_test, y_pred)))
print('\033[1m+ 'R2 Score :'+ '\033[0m')
print(r2_score(Y_test,y_pred)*100)
```

**Error :**  
Mean absolute error : 4002.1921561308695  
Mean squared error : 27952391.253184095  
Root Mean squared error : 5287.002104518599  
**R2 Score :**  
49.070154943012476

### 2. Random Forest Regressor:

```
X_train, X_test, Y_train, Y_test = train_test_split(X_scale, Y, random_state= 60, test_size=0.33)
rfc = RandomForestRegressor()
rfc.fit(X_train, Y_train)
y_pred = rfc.predict(X_test)
print('\033[1m+ 'Error of Random Forest Regressor:'+ '\033[0m')
print('Mean absolute error :', mean_absolute_error(Y_test,y_pred))
print('Mean squared error :', mean_squared_error(Y_test, y_pred))
print('Root Mean squared error :', np.sqrt(mean_squared_error(Y_test, y_pred)))
print('\033[1m+ 'R2 Score of Random Forest Regressor :'+ '\033[0m')
print(r2_score(Y_test,y_pred)*100)
```

**Error of Random Forest Regressor:**  
Mean absolute error : 1147.236136224076  
Mean squared error : 4742093.943375699  
Root Mean squared error : 2177.6349426328784  
**R2 Score of Random Forest Regressor :**  
89.94015869083671

### 3. Decision Tree Regressor

```
X_train, X_test, Y_train, Y_test = train_test_split(X_scale, Y, random_state= 60, test_size=0.33)
dtc = DecisionTreeRegressor()
dtc.fit(X_train, Y_train)
y_pred = dtc.predict(X_test)
print('\033[1m+ 'Error of Decision Tree Regressor:'+ '\033[0m')
print('Mean absolute error :', mean_absolute_error(Y_test,y_pred))
print('Mean squared error :', mean_squared_error(Y_test, y_pred))
print('Root Mean squared error :', np.sqrt(mean_squared_error(Y_test, y_pred)))
print('\033[1m+ 'R2 Score of Decision Tree Regressor :'+ '\033[0m')
print(r2_score(Y_test,y_pred)*100)
```

**Error of Decision Tree Regressor:**  
Mean absolute error : 1107.147201946472  
Mean squared error : 6872305.638686132  
Root Mean squared error : 2621.508275532643  
**R2 Score of Decision Tree Regressor :**  
85.42114412350999

#### 4. Gradient Boosting Regressor

```
from sklearn.ensemble import GradientBoostingRegressor
X_train,X_test,Y_train,Y_test=train_test_split(X_scale,Y,test_size=.33,random_state=60)
GBR=GradientBoostingRegressor()
GBR.fit(X_train,Y_train)
y_pred = GBR.predict(X_test)
print('\033[1m+ 'Error of Gradient Boosting Regressor:'+' '\033[0m')
print('Mean absolute error :', mean_absolute_error(Y_test,y_pred))
print('Mean squared error :', mean_squared_error(Y_test, y_pred))
print('Root Mean squared error :', np.sqrt(mean_squared_error(Y_test, y_pred)))
print('\033[1m+'R2 Score of Gradient Boosting Regressor :'+'\033[0m')
print(r2_score(Y_test,y_pred)*100)
```

```
Error of Gradient Boosting Regressor:
Mean absolute error : 1515.589925036542
Mean squared error : 5277834.301348992
Root Mean squared error : 2297.3537606013124
R2 Score of Gradient Boosting Regressor :
88.80364325093188
```

#### 5. XGB Regressor

```
X_train, X_test, Y_train, Y_test = train_test_split(X_scale, Y, random_state= 60, test_size=0.33)
xgb = XGBRegressor()
xgb.fit(X_train, Y_train)
y_pred = xgb.predict(X_test)
print('\033[1m+ 'Error of XGB Regressor:'+' '\033[0m')
print('Mean absolute error :', mean_absolute_error(Y_test,y_pred))
print('Mean squared error :', mean_squared_error(Y_test, y_pred))
print('Root Mean squared error :', np.sqrt(mean_squared_error(Y_test, y_pred)))
print('\033[1m+'R2 Score of XGB Regressor :'+'\033[0m')
print(r2_score(Y_test,y_pred)*100)
```

```
Error of XGB Regressor:
Mean absolute error : 1258.1237008867472
Mean squared error : 5038598.295239574
Root Mean squared error : 2244.6822258929155
R2 Score of XGB Regressor :
89.31115665864509
```

**5-Fold cross validation performed over all models. We can see that Random Forest Regressor gives maximum R2 score of 89.94 and maximum cross validation score. Among all model we will select XGB Regressor as final model and we will perform hyper parameter tuning over this model to enhance its R2 Score.**

## Hyper Parameter Tuning : GridSearchCV

```
from sklearn.model_selection import GridSearchCV

X_train, X_test, Y_train, Y_test = train_test_split(X_scale, Y, random_state=287, test_size=0.25)

rfc.get_params().keys()

dict_keys(['bootstrap', 'ccp_alpha', 'criterion', 'max_depth', 'max_features', 'max_leaf_nodes', 'max_samples', 'min_impurity_decrease', 'min_samples_leaf', 'min_samples_split', 'min_weight_fraction_leaf', 'n_estimators', 'n_jobs', 'oob_score', 'random_state', 'verbose', 'warm_start'])

parameter = {'n_estimators':[90,100,125],
             'bootstrap':[True, False],
             'max_depth':[4,6,8,10, None],
             'max_features':['auto','log2','sqrt'] }

GCV = GridSearchCV(RandomForestRegressor(),parameter,verbose =10)

GCV.fit(X_train,Y_train)
```

## Final Model

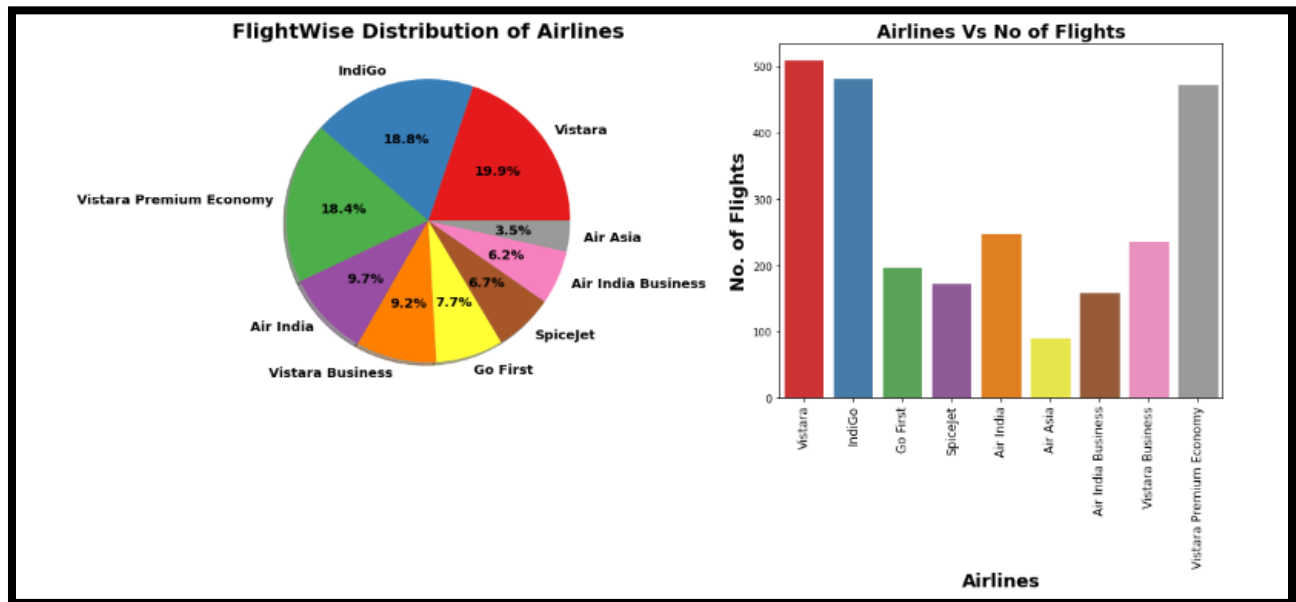
```
Final_mod=RandomForestRegressor(bootstrap=True,max_depth=None,max_features='auto',n_estimators=100)
Final_mod.fit(X_train,Y_train)
pred=Final_mod.predict(X_test)
print('R2_Score:',r2_score(Y_test,pred)*100)
print('mean_squared_error:',mean_squared_error(Y_test,pred))
print('mean_absolute_error:',mean_absolute_error(Y_test,pred))
print("RMSE value:",np.sqrt(mean_squared_error(Y_test, pred)))

R2_Score: 90.93769090668945
mean_squared_error: 4973767.525702495
mean_absolute_error: 1084.8790096463022
RMSE value: 2230.1945040068804
```

## Saving Model

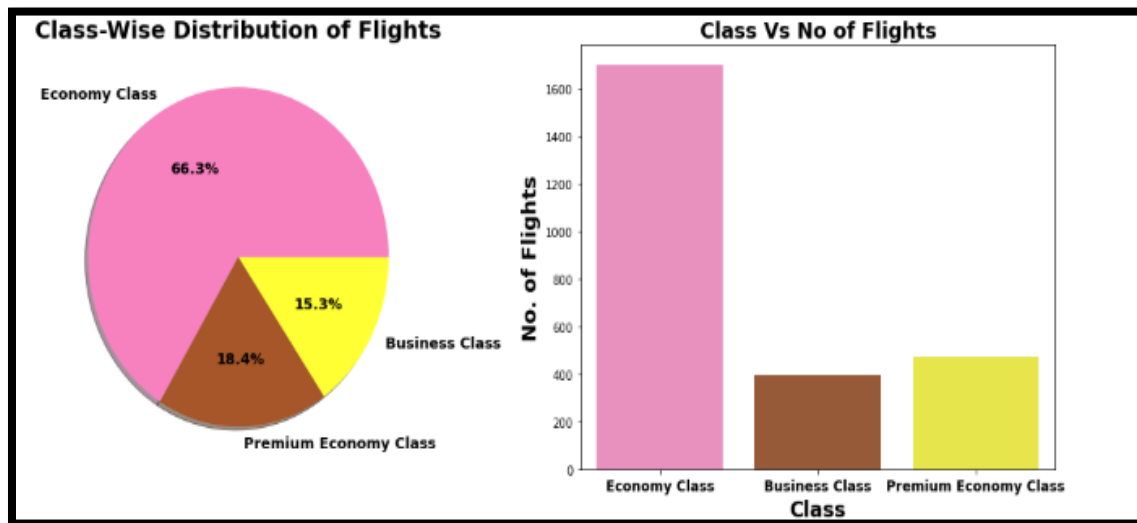
```
# Saving the model using .pkl
import joblib
joblib.dump(Final_mod,"Final_Flight_Price_Prediction.pkl")
```

## 6. VISUALIZATIONS



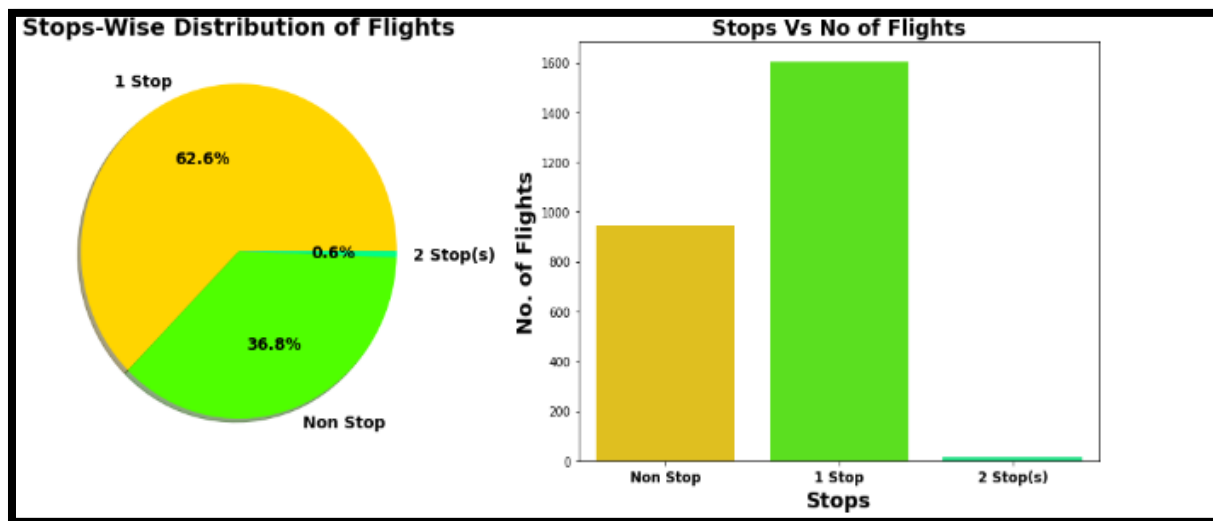
### Observation:

- We can see maximum number of flights run by Vistara Premium Economy while minimum Flights run by SpiceJet.
- Around 15% of flights of Business Class.



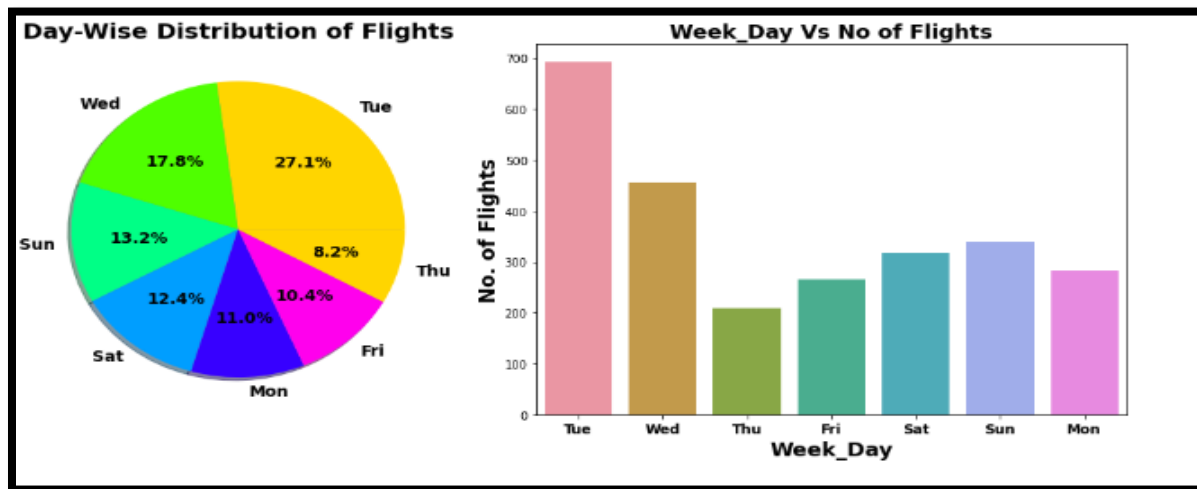
*Observation:*

- 66.3% flights are of Economy class, as they are low cost of flight & most of people prefer it.
- There are more Premium Economy flights than business class flights. It not because Business class is less costly than Premium Economy class.

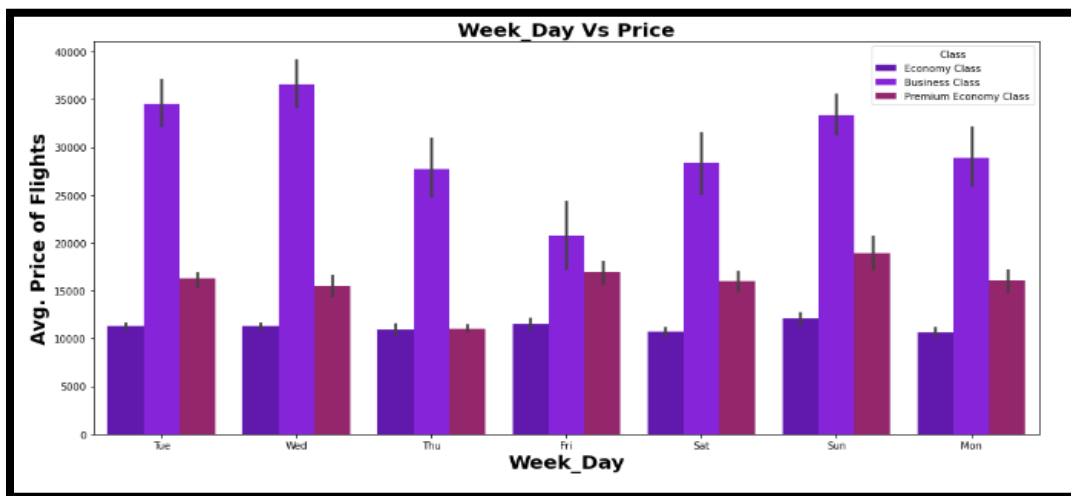


*Observation:*

- 62.6% flights take single stop in there way from New Dehli to Mumbai. It is also possible that these flights may have high flight duration compare to Non-stop Flight.
- 36.8% of flights do not have any stop in their route.



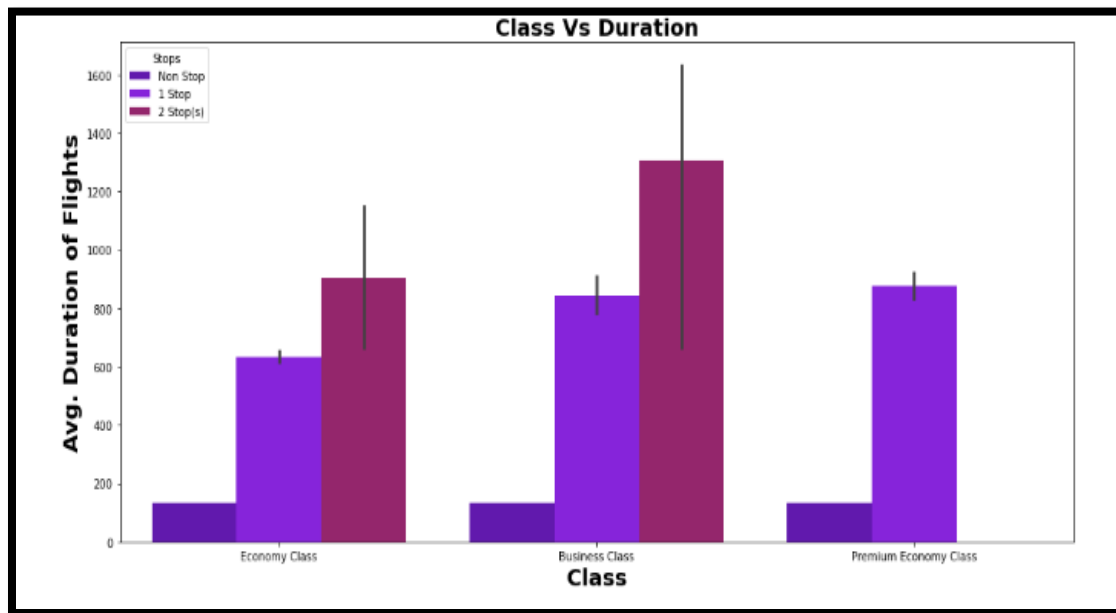
- On Tuesday Maximum flights run while on Thursday minimum flights run
- It will be interesting to investigate variation of fare as per different week days.



#### Observation:

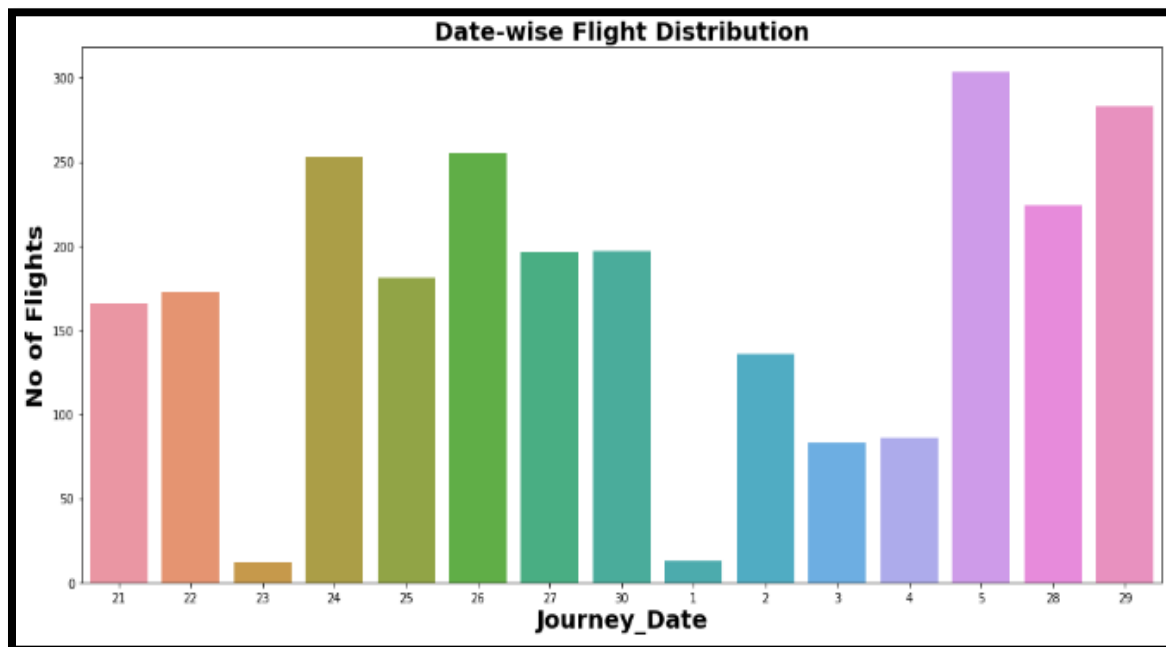
- Maximum Avg. Fare for Business Flights is on Wednesday while minimum Avg. Fare for Business flights on Thursday.
- For Economy Flights & Premium Economy Flights: Minimum Avg. Fare on Friday.
- For Economy Flights & Premium Economy Flights: Maximum Avg. Fare on Monday.





#### Observation:

- As Number of Stops increase the duration of flights increases.
- As per Class of flight Maximum Avg. Duration of flight is for Business class.



We can see those Maximum flights schedule on 5th July 2022 & Minimum flights schedule on 23 June 2022.

## Chap 4. Conclusion

### 1. Key Findings and Conclusions of the Study

Algorithm	R2 Score	CV Score
Random Forest Regressor	89.94	0.63
XGB Regressor	89.31	0.62
Linear Regression	49.07	0.63
Decision Tree Regressor	85.42	0.53
Extra Tree Regressor	84.78	0.66
Gradient Boost Regressor	88.80	0.67
Random Forest Regressor Parameter Tuned Final Model	90.93	0.82

- Random Forest Regressor giving us maximum R2 Score, so Random Forest Regressor is selected as best model.
- After hyper parameter tuning *Final Model* is giving us R2 Score of 90.93% which is slightly improved compare to earlier R2 score of 89.31%.

## 2. Limitations of this work and Scope for Future Work

- In this study we focus on flights on route of New Delhi to Mumbai, more route can incorporate in this project to extend it beyond present investigation.
- This investigation focus on short timeframe (14 days prior flights take off) which can be extended variation over larger period.
- Time series analysis can be performed over this model.