

Flight Price Prediction Project

BY:

Akas Mandal



ACKNOWLEDGMENT

I express my sincere gratitude to Flip Robo Technologies for giving me the opportunity to work on this project on Flight Price Prediction using machine learning algorithms. I acknowledge my indebtedness to the authors of the paper titled: “Airline ticket price and demand prediction: A survey” and the online article titled: “Trying to Predict Airfares When The Unpredictable Happens” for providing me with invaluable insights and knowledge of the various factors that determine the price of Flight tickets.

Business Problem Framing

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)

Therefore, a predictive model to accurately predict Air fares is required to be made.

Conceptual Background of the Domain Problem

Predictive modelling, Regression algorithms are some of the machine learning techniques used for predicting Flight Ticket prices. Identifying various relevant attributes like Airline Brand, flight duration, source and destination etc are crucial for working on the project as they determine the valuation of air fare.

Review of Literature

A Research paper titled: “Airline ticket price and demand prediction: A survey” by Juhar Ahmed Abdella and online article titled: “Trying to Predict Airfares When The Unpredictable Happens” were reviewed and studied to gain insights into all the attributes that contribute to the pricing of flight tickets.

It is learnt that deterministic features like Airline Brand, flight number, departure dates, number of intermediate stops, week day of departure, number of competitors on route and aggregate features – which are based on collected historical data on minimum price, mean price, number of quotes on nonstop, 1-stop and multi-stoppage flights are some the most important factors that determine the pricing of Flight Tickets.

- [Airline ticket price and demand prediction: A survey - ScienceDirect](#)
- [Flight Price Predictor | American Express GBT \(amexglobalbusinessstravel.com\)](#)

INTRODUCTION

Motivation for the Problem Undertaken

With airfares fluctuating frequently, knowing when to buy and when to wait for a better deal to come along is tricky. The fluctuation in prices is frequent and one has limited time to book the cheapest ticket as the prices keep varying due to constant manipulation by Airline companies. Therefore, it is necessary to work on a predictive model based on deterministic and aggregate feature data that would predict with good accuracy the most optimal Air fare for a particular destination, route and schedule.

Analytical Problem Framing

Mathematical/ Analytical Modeling of the Problem

Various Regression analysis techniques were used to build predictive models to understand the relationships that exist between Flight ticket price and Deterministic and Aggregate features of Air travel. The Regression analysis models were used to predict the Flight ticket price value for changes in Air travel deterministic and aggregate attributes. Regression modelling techniques were used in this Problem since Air Ticket Price data distribution is continuous in nature.

In order to forecast Flight Ticket price, predictive models such as ridge regression Model, Random Forest Regression model, Decision tree Regression Model, Support Vector Machine Regression model and Extreme Gradient Boost Regression model were used to describe how the values of Flight Ticket Price depended on the independent variables of various Air Fare attributes.

Data Sources and their formats

The Dataset was compiled by scraping Data for various Air Fare attributes and Price from <https://www.yatra.com/> and <https://www.easemytrip.com/>

The data was converted into a Pandas Dataframe under various Feature and Label columns and saved as a .csv file.

```

1 fDF=pd.read_csv("FlightDataWebscrape.csv")
2 fDF.head(50)

```

	Unnamed: 0	Airline	Flight Number	Date of Departure	From	To	Duration	Total Stops	Price
0	0	Air Asia	I5-764	Tue, Feb 8	New Delhi	Mumbai	2h 10m	Non Stop	2,456
1	1	IndiGo	6E-2054	Tue, Feb 8	New Delhi	Mumbai	2h 10m	Non Stop	2,456
2	2	IndiGo	6E-5001	Tue, Feb 8	New Delhi	Mumbai	2h 10m	Non Stop	2,456
3	3	IndiGo	6E-2046	Tue, Feb 8	New Delhi	Mumbai	2h 10m	Non Stop	2,456
4	4	IndiGo	6E-5328	Tue, Feb 8	New Delhi	Mumbai	2h 15m	Non Stop	2,456
5	5	Air Asia	I5-482	Tue, Feb 8	New Delhi	Mumbai	2h 15m	Non Stop	2,456
6	6	IndiGo	6E-6278	Tue, Feb 8	New Delhi	Mumbai	2h 20m	Non Stop	2,456
7	7	IndiGo	6E-218	Tue, Feb 8	New Delhi	Mumbai	2h 20m	Non Stop	2,456
8	8	IndiGo	6E-171	Tue, Feb 8	New Delhi	Mumbai	2h 25m	Non Stop	2,456
9	9	IndiGo	6E-2081	Tue, Feb 8	New Delhi	Mumbai	2h 25m	Non Stop	2,456
10	10	Air India	AI-537	Thu, Feb 24	New Delhi	Chennai	2h 45m	Non Stop	2,343
11	11	Air India	AI-439	Thu, Feb 24	New Delhi	Chennai	2h 50m	Non Stop	2,343
12	12	Air India	AI-429	Thu, Feb 24	New Delhi	Chennai	2h 50m	Non Stop	2,343
13	13	Air India	AI-540	Thu, Feb 24	New Delhi	Chennai	2h 55m	Non Stop	2,343
14	14	IndiGo	6E-5012	Thu, Feb 24	New Delhi	Chennai	2h 45m	Non Stop	2,637

Dataset Description

The Independent Feature columns are:

- Airline: The name of the airline.
- Flight Number: Number of Flight
- Date of Departure: The date of the journey

- From: The source from which the service begins
- To: The destination where the service ends
- Duration: Total duration of the flight
- Total Stops: Total stops between the source and destination.

Target / Label Column:

☐ Price: The Price of the

Ticket

Data Pre-processing Done

- Duplicate data elements in various columns: 'Airline', 'From', 'To', which had their starting letters in upper case and lower case were converted to data elements starting with uppercase letters.
- Data in column 'Price' was converted to int64 data type.
- Columns: Unnamed: 0(just a series of numbers) was dropped since it doesn't contribute to building a good model for predicting the target variable values.
- The Date format of certain data elements in 'Date of Departure' was changed to match the general Date format of majority of the data elements of the column.

Feature Engineering:

- In order to better understand the relationships between Flight price and Air Fare attributes, 'Day', 'Date' and 'Month' columns were created based on data of existing column: 'Date of Departure'.
- The values in Column: 'Duration' were converted from Hours-Minutes format to minute format and the data type was converted to int64.

Data Inputs- Logic- Output Relationships

The Datasets consist mainly of Int and Object data type variables. The relationships between the independent variables and dependent variable were analysed.

Requirements and Tools Used

Hardware Used:

Device nameLAPTOP-S0KG8J4C

Processor Intel(R) Core (TM) i5-1035G1 CPU @ 1.00GHz 1.20 GHz

Installed RAM 8.00 GB (7.74 GB usable)

Device ID A2B7913A-CE3E-40D0-9D50-7E5599B13EAD

Product ID 00327-36296-11993-AAOEM

System type 64-bit operating system, x64-based processor

Requirements and Tools Used

Software Used:

- Windows 11 Operating System
- Anaconda Package and Environment Manager: Anaconda is a distribution of the Python and R programming languages for scientific computing, that aims to simplify package management and deployment. The distribution includes data-science packages suitable for Windows and provides a host of tools and environment for conducting Data Analytical and Scientific works. Anaconda provides all the necessary Python packages and libraries for Machine learning projects

Requirements and Tools Used

- Jupyter Notebook: The Jupyter Notebook is an open-source web application that allows data scientists to create and share documents that integrate live code, equations, computational output, visualizations, and other multimedia resources, along with explanatory text in a single document.
- Python3: It is open source, interpreted, high level language and provides great approach for object-oriented programming. It is one of the best languages used for Data Analytics And Data science projects/application. Python provides numerous libraries to deal with mathematics, statistics and scientific function.

Requirements and Tools Used

- Python Libraries used:
 - o Pandas: For carrying out Data Analysis, Data Manipulation, Data Cleaning etc o Numpy: For performing a variety of operations on the datasets.
 - o matplotlib.pyplot, Seaborn: For visualizing Data and various relationships between Feature and Label Columns
 - o Scipy: For performing operations on the datasets o Statsmodels: For performing statistical analysis

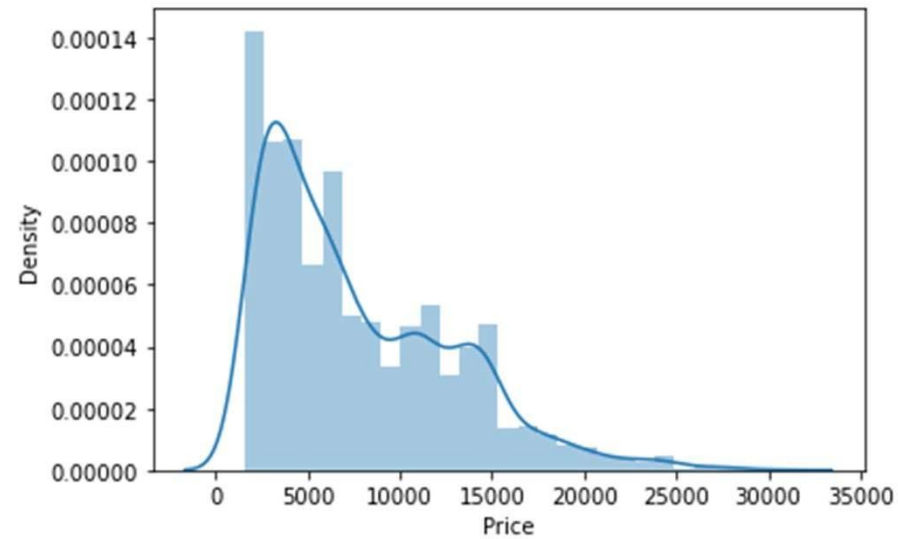
- sklearn for Modelling Machine learning algorithms, Data Encoding, Evaluation metrics, Data Transformation, Data Scaling, Component analysis, Feature selection etc.

Visualizations

Barplots, Distplots, Boxplots, Countplots, lineplots were used to visualise the data of all the columns and their relationships with Target variable.

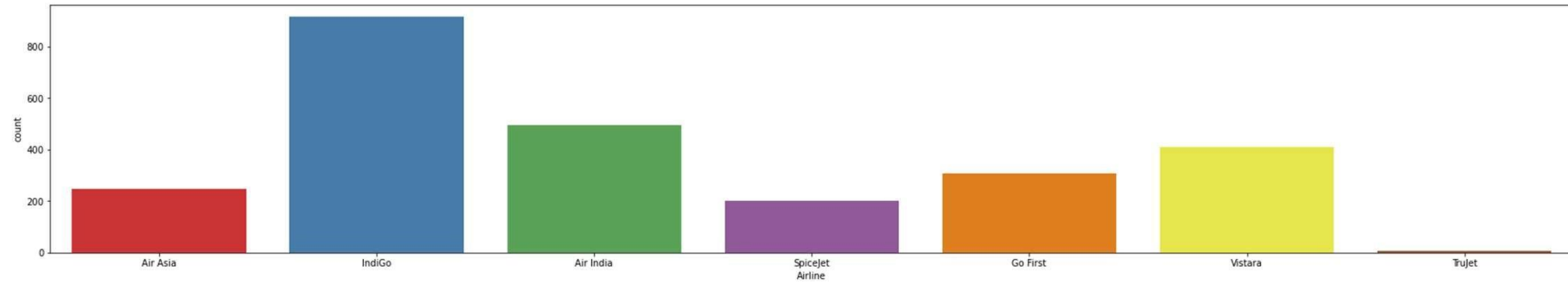
Univariate Analysis

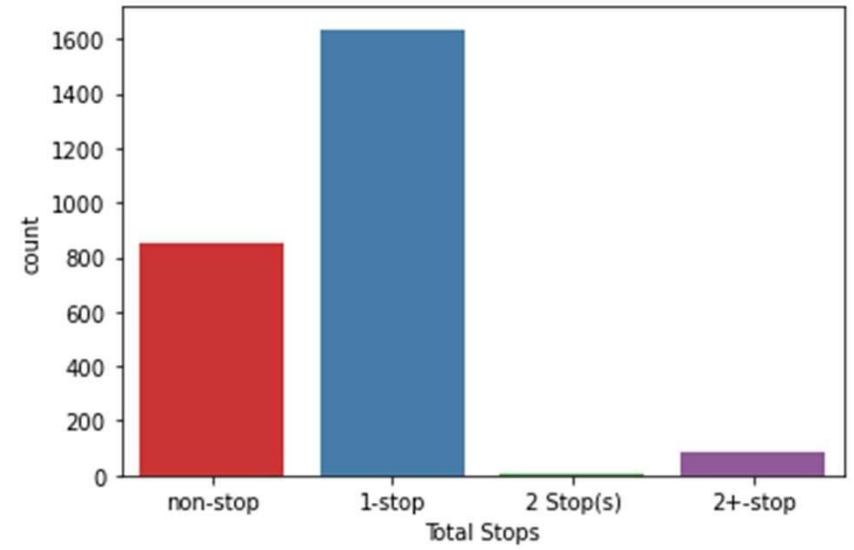
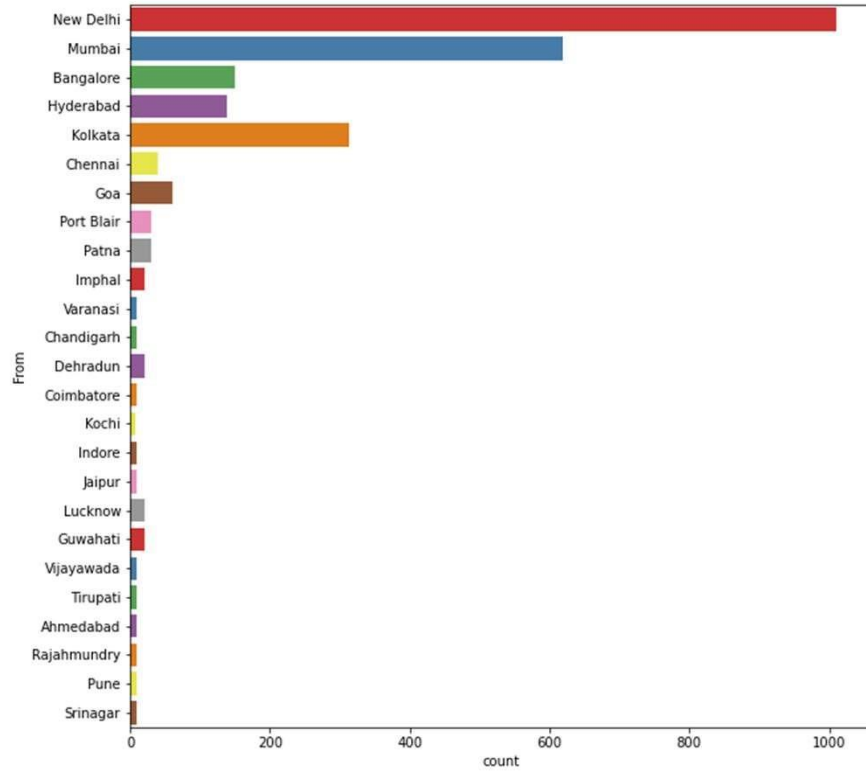
Analyzing the Target Variable

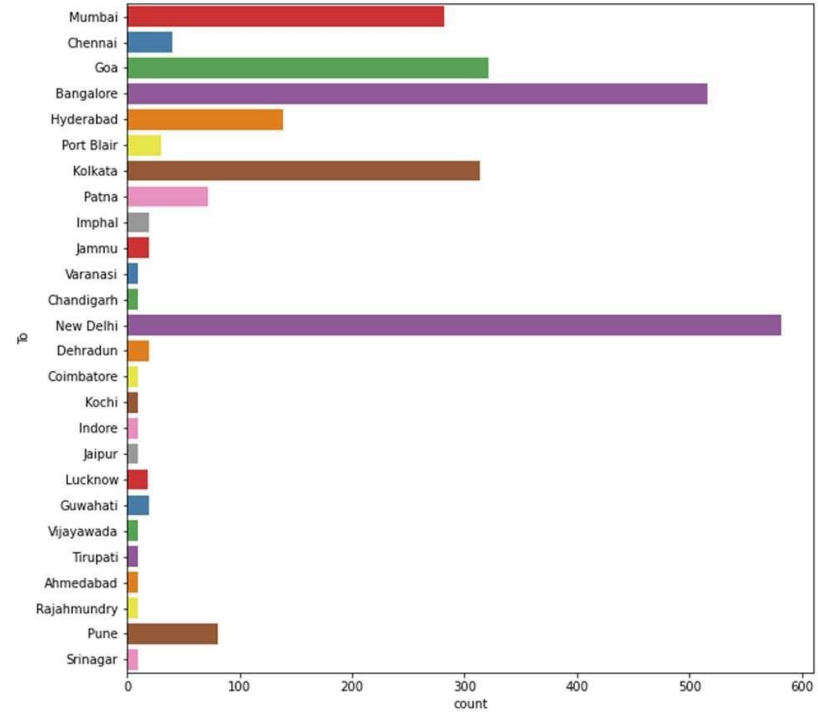


From the graph above it is observed that the Price data forms a continuous distribution with mean of 7748.33 and tails of from 15000 mark and the distribution is skewed.

Analyzing the Feature Columns







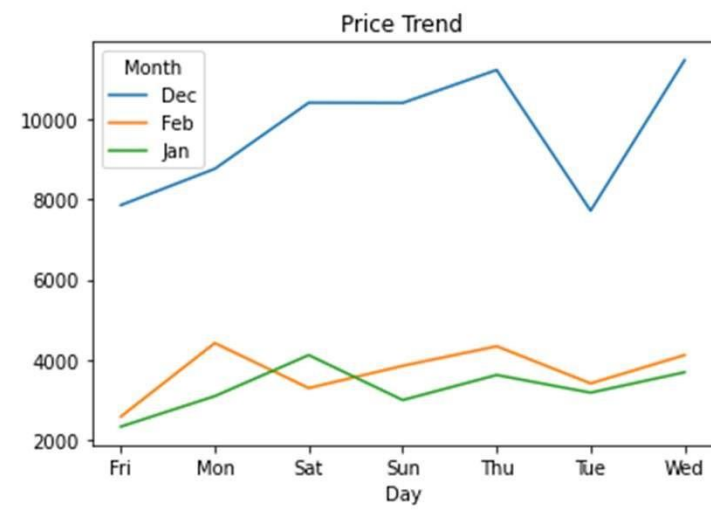
Following observations are made from graphs above:

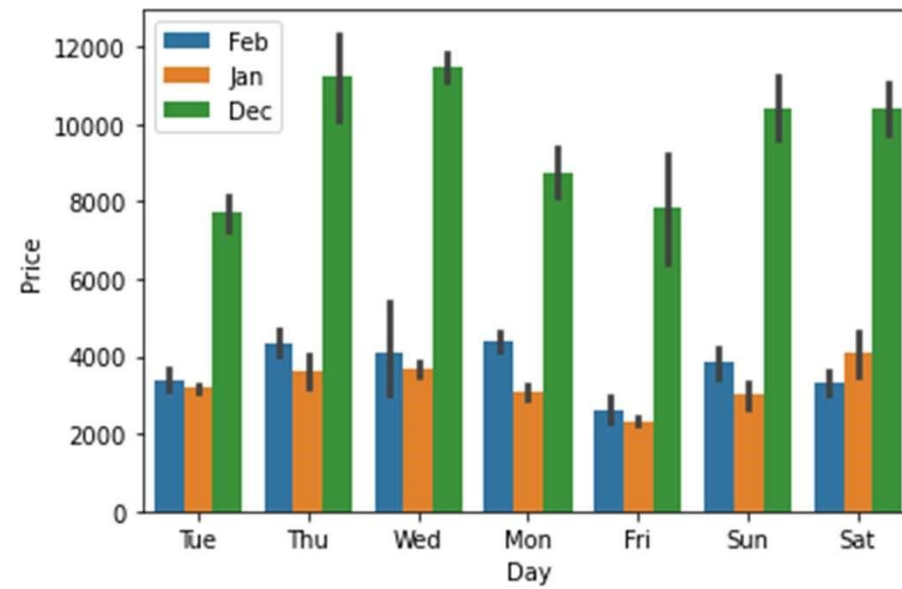
- IndiGo has the highest number of flights followed by Air India and Vistara
- Highest number of flights are from Delhi followed by Mumbai, Kolkata, Bangalore and Hyderabad
- New Delhi is the most popular destination followed by Bangalore, Goa, Kolkata and Mumbai
- Highest number of flights have only 1 stop between source and destination while 2nd highest number of flights are non stop

Bivariate Analysis

Interpreting Relationship between Dependent Variable and Independent Variable Columns

Analyzing Relationship between Day, Month columns and Price

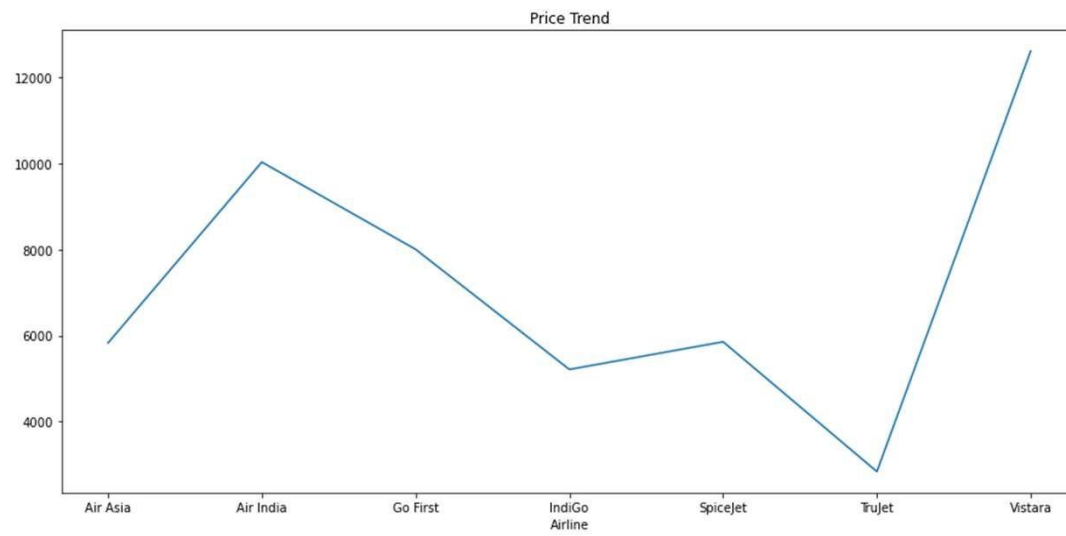


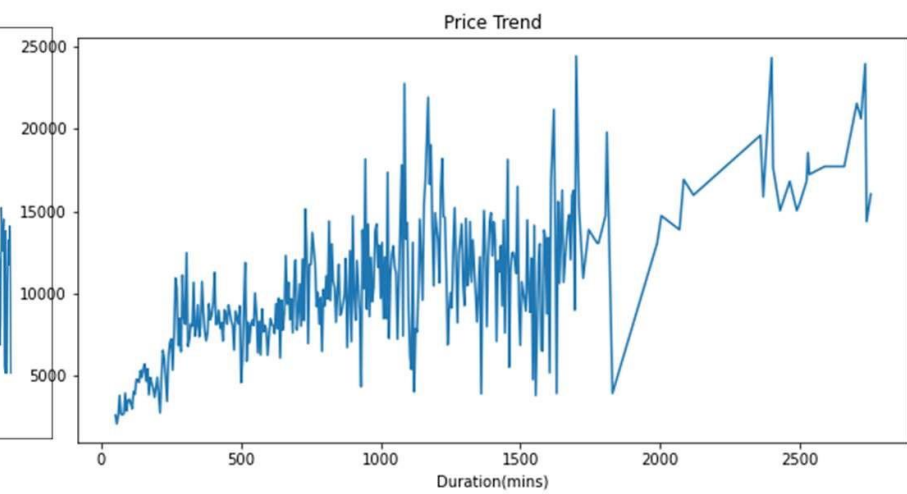
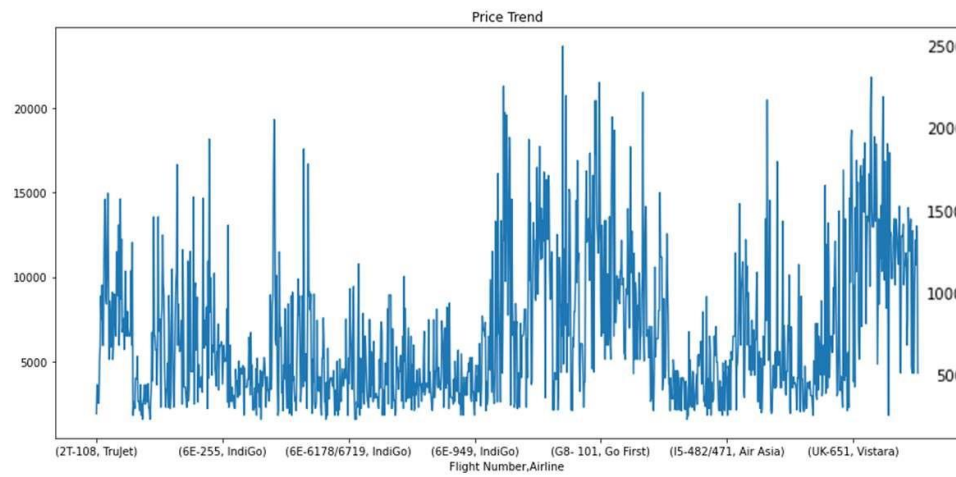


Following observations are made from graphs above:

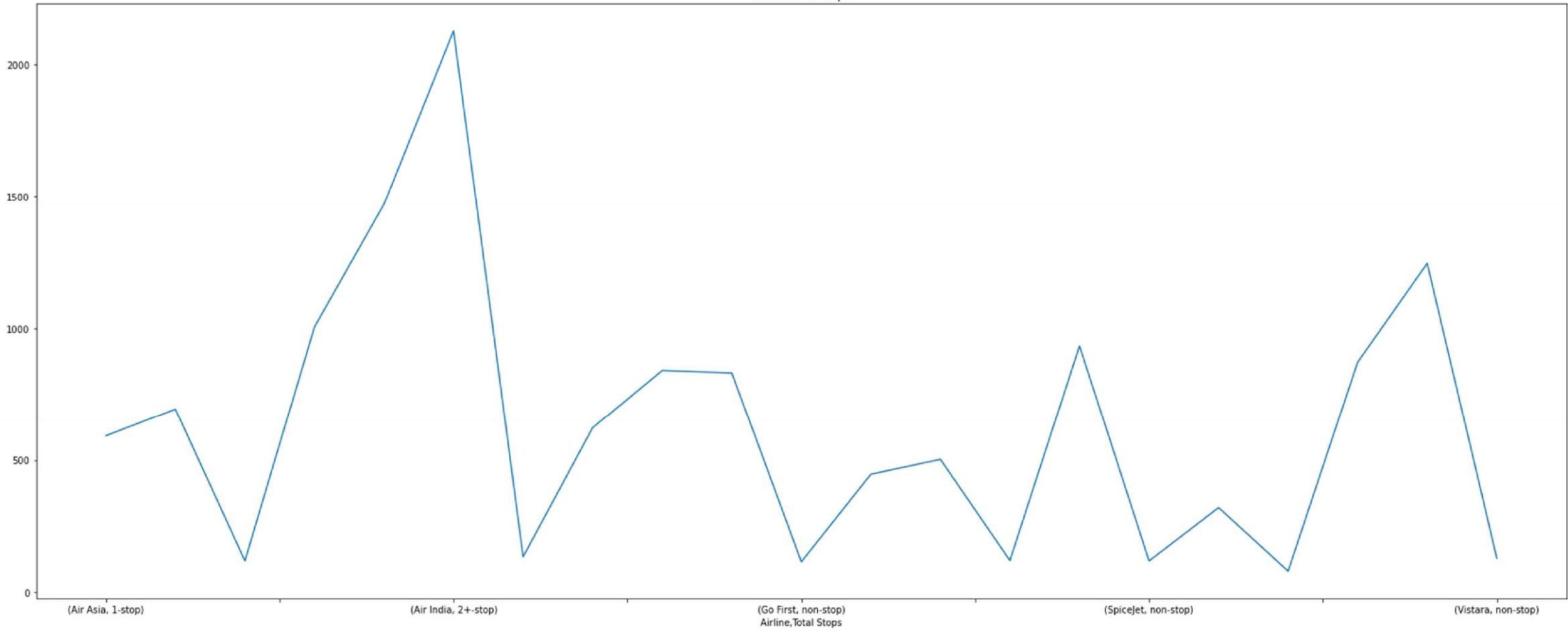
- On an average, there is a steady decline in Flight price from December to February, with the prices being lowest in January.
- Flight Prices increase on an average, as the day of departure gets nearer.
- Flight Ticket prices are the highest on Thursdays, Mondays and during the weekend on an average.

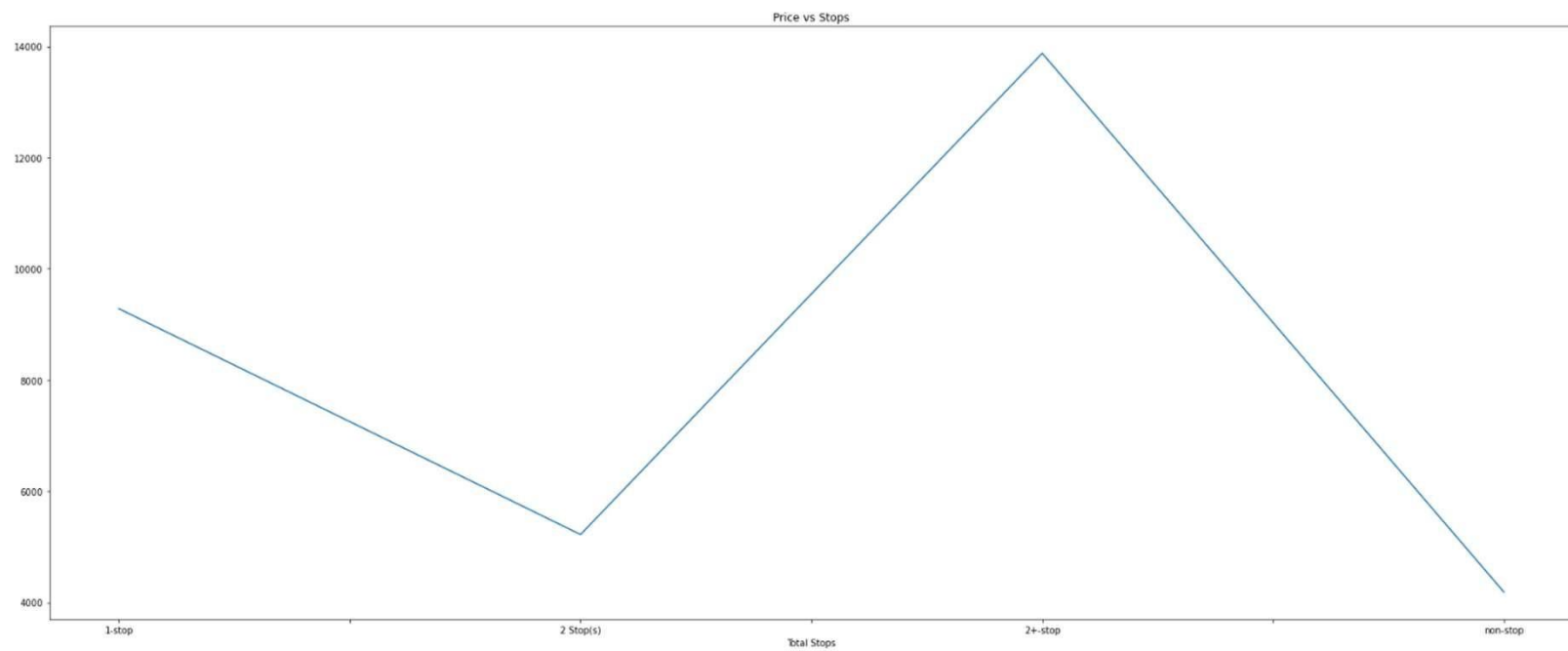
Analyzing Relationship between Airlines, Flight Duration and Price

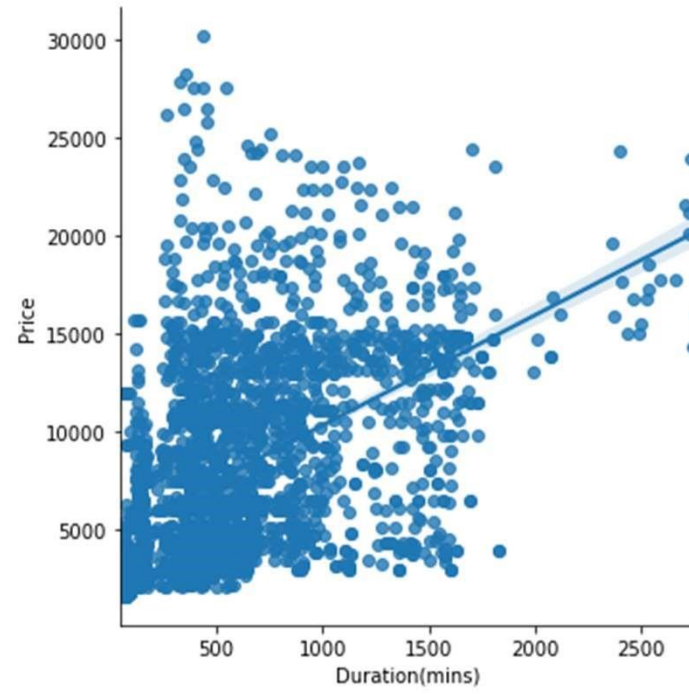




Duration vs Stops



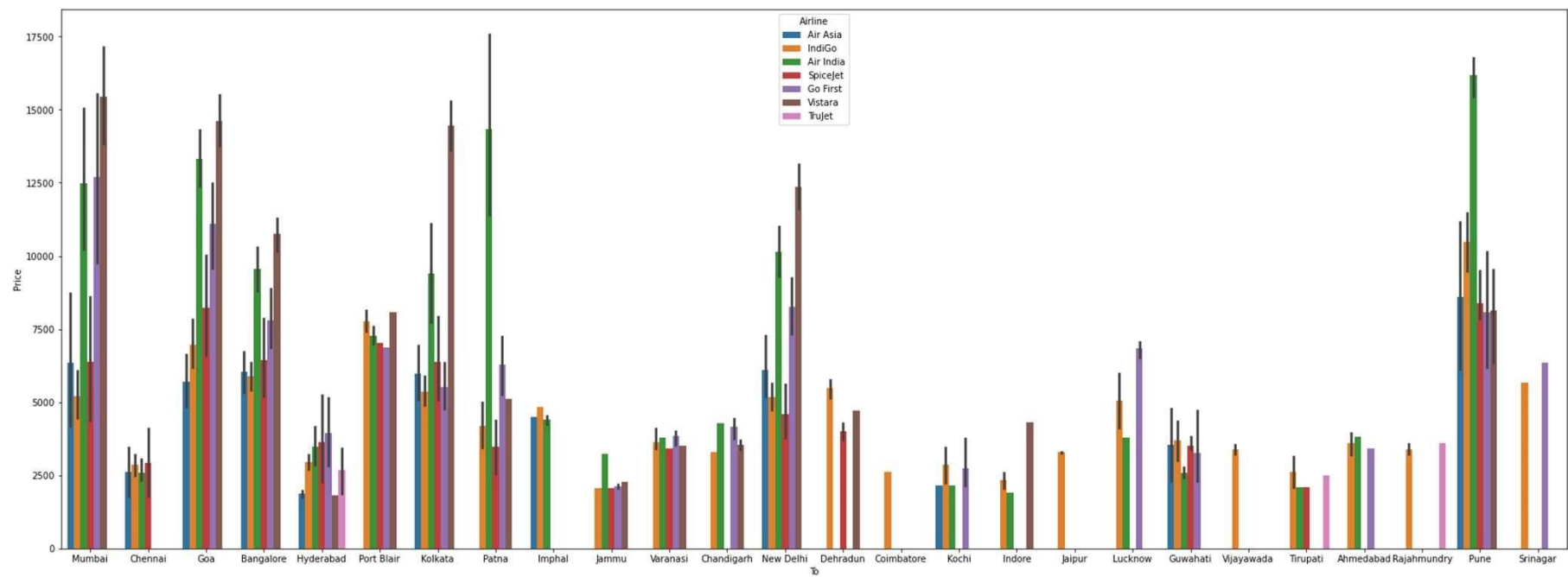




Following Observation is made from graphs above:

o Trujet, IndiGo, SpiceJet and Air Asia offer air tickets at the most affordable prices on average, whereas Vistara, Air India are the most expensive on average. o It can be observed that Number of Stops impact the travel time of Airlines. o It can be observed that Number of Stops impact the Air Ticket Pricing of Airlines. o There is a linear relationship between Price and flight duration.

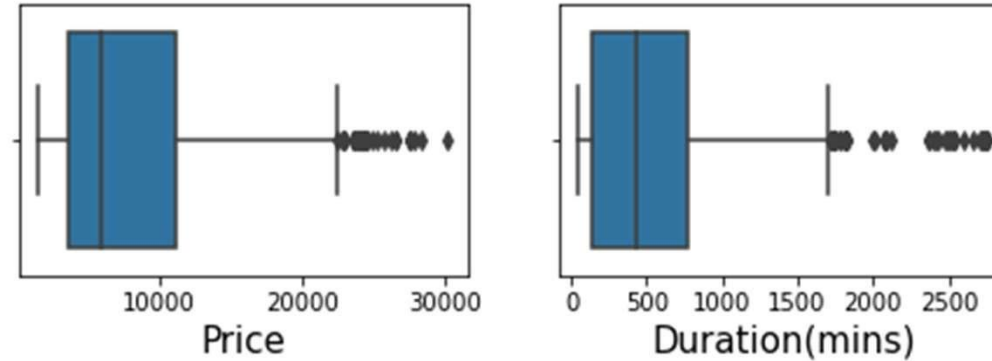
Multivariate Analysis



Following Observations are made from graphs above:

- There is a linear relationship between Price and flight duration.
- Indigo,Air Asia and Spicejet provide most affordable Airtickets to the destinations.

Checking for Outliers



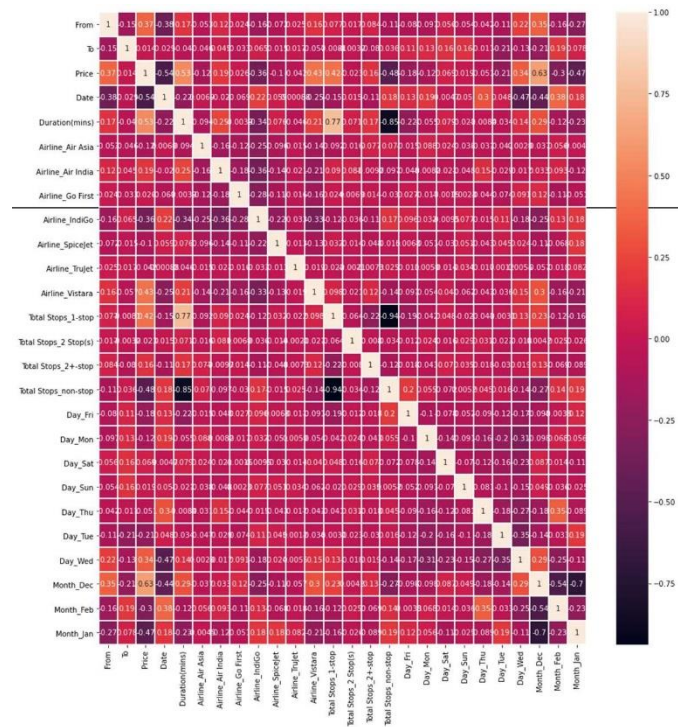
There are considerable outliers in the columns.

Outliers were Removed using Z score method which resulted in a total data loss of 1.00%, which is within acceptable range.

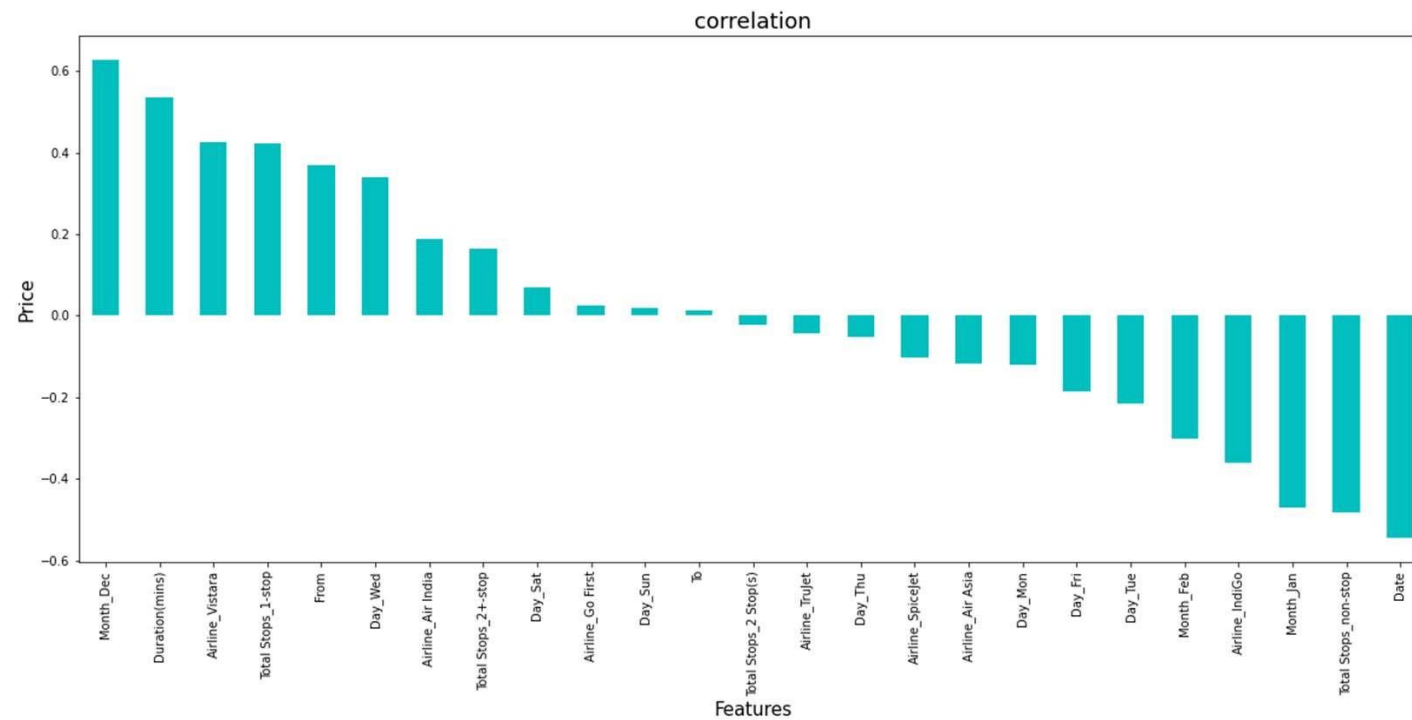
Data Normalization

Data in Column 'Duration(mins)' was normalized using Power Transformer technique.

Finding Correlation between Feature and Target columns



Exploratory Data Analysis



Exploratory Data Analysis

It is observed that Month_Dec, Duration(mins), Airline_Vistara, Total Stops_1-stop and From have the highest positive correlation with Price, while Date, Total Stops_nonstop, Month_Jan, Airline_IndiGo have the highest negative correlation with Price.

Evaluation

Feature Selection

Features were first checked for presence of multicollinearity and then based on respective ANOVA fscore values, the feature columns were selected that would best predict the Target variable, to train and test machine learning models.

Using SelectKBest and f_classif for measuring the respective ANOVA f-score values of the columns, the best features were selected. Using StandardScaler, the features were scaled by resizing the distribution values so that mean of the observed values in each feature column is 0 and standard deviation is 1.

From sklearn.model_selection's train_test_split, the data was divided into train and test data. Training data comprised 75% of total data where as test data comprised 25% based on the best random state that would result in best model accuracy.

Evaluation

The model algorithms used were as follows:

- Ridge
- DecisionTreeRegressor
- XGBRegressor
- RandomForestRegressor:
- Support Vector Regressor:

Evaluation

Regression Model Building

```
1 from sklearn.model_selection import train_test_split
```

```
1 from sklearn.metrics import r2_score
```

Finding the Best Random State

```
1 from sklearn.ensemble import RandomForestRegressor
2 maxAcc = 0
3 maxRS=0
4 for i in range(1,100):
5     x_train,x_test,y_train,y_test = train_test_split(scaled_x_best,y,test_size = .25, random_state = i)
6     modRF = RandomForestRegressor()
7     modRF.fit(x_train,y_train)
8     pred = modRF.predict(x_test)
9     acc = r2_score(y_test,pred)
10    if acc>maxAcc:
11        maxAcc=acc
12        maxRS=i
13 print(f"Best Accuracy is: {maxAcc} on random_state: {maxRS}")
```

Best Accuracy is: 0.8232413568072626 on random_state: 58

Best random state was determined to be 58

```
1 x_train,x_test,y_train,y_test = train_test_split(scaled_x_best,y,test_size = .25, random_state =58)
```

```
1 from sklearn.model_selection import GridSearchCV
2 from sklearn.linear_model import Ridge
3 from sklearn.tree import DecisionTreeRegressor
4 from sklearn.ensemble import RandomForestRegressor
5 from xgboost import XGBRegressor
6 from sklearn.svm import SVR
7
8
```

```
1 from sklearn.metrics import r2_score,mean_squared_error
```

```
1 rf = RandomForestRegressor()
2 dt = DecisionTreeRegressor()
3 xg = XGBRegressor()
4 sv = SVR()
5 r=Ridge()
```

Model/s Development and Evaluation

Training The Models

```
1 rf.fit(x_train,y_train)
2 xg.fit(x_train,y_train)
3 SV.fit(x_train,y_train)
4 r.fit(x_train,y_train)
5 dt.fit(x_train,y_train)
```

Evaluation

Analyzing Accuracy of The Models

Mean Squared Error and Root Mean Squared Error metrics were used to evaluate the Model performance. The advantage of MSE and RMSE being that it is easier to compute the gradient. As, we take square of the error, the effect of larger errors become more pronounced than smaller error, hence the model can now focus more on the larger errors.

R2 Score

```
1 r2_score(y_test,y_rf_pred)
```

```
0.8222913567368567
```

Mean Squared Error

```
1 mean_squared_error(y_test,y_rf_pred)
```

```
4444440.963715802
```

Root Mean Squared Error

```
1 np.sqrt(mean_squared_error(y_test,y_rf_pred))
```

```
2108.184281251476
```

XGB Regression Model

```
1 y_xg_pred = xg.predict(x_test)
```

R2 Score

```
1 r2_score(y_test,y_xg_pred)
```

```
0.8098289634792951
```

Mean Squared Error

```
1 mean_squared_error(y_test,y_xg_pred)
```

```
4756121.758092391
```

Ridge Regression Model ¶

```
In [99]: 1 y_r_pred = r.predict(x_test)
```

R2 Score

```
In [100]: 1 r2_score(y_test,y_r_pred)
```

```
Out[100]: 0.7058259974328809
```

Mean Squared Error

```
In [101]: 1 mean_squared_error(y_test,y_r_pred)
```

```
Out[101]: 7357205.386648202
```

Root Mean Squared Error

```
In [102]: 1 np.sqrt(mean_squared_error(y_test,y_r_pred))
```

```
Out[102]: 2712.4168902748343
```

Random Forest Regression Model

```
In [103]: 1 y_rf_pred = rf.predict(x_test)
```

Root Mean Squared Error

```
|: 1 np.sqrt(mean_squared_error(y_test,y_xg_pred))  
|: 2180.853447183554
```

Decision Tree Regression Model ¶ Support Vector Regression Model

```
1 y_dt_pred = dt.predict(x_test)
```

```
|: 1 y_svr_pred = SV.predict(x_test)
```

R2 Score

```
1 r2_score(y_test,y_dt_pred)
```

0.7247248097918293

R2 Score

```
|: 1 r2_score(y_test,y_svr_pred)
```

-0.07322019200984475

Mean Squared Error

```
1 mean_squared_error(y_test,y_dt_pred)
```

6884551.641330296

Mean Squared Error

```
|: 1 mean_squared_error(y_test,y_svr_pred)
```

26840921.72935272

Root Mean Squared Error

```
1 np.sqrt(mean_squared_error(y_test,y_dt_pred))
```

2623.842914758865

Root Mean Squared Error

```
|: 1 np.sqrt(mean_squared_error(y_test,y_svr_pred))
```

5180.822495449223

Evaluation

Using cross-validation, there are high chances that we can detect over-fitting with ease. Model Cross Validation scores were then obtained for assessing how the statistical analysis generalises to an independent data set. The models were evaluated by training several models on subsets of the available input data and evaluating them on the complementary subset of the data.

Model Cross Validation

```
1 from sklearn.model_selection import ShuffleSplit, cross_val_score
```

Ridge Regression

```
1 cross_val_score(r, scaled_x_best, y, cv=ShuffleSplit(5)).mean()
```

0.6809490952194228

Random Forest Regression

```
1 cross_val_score(rf, scaled_x_best, y, cv=ShuffleSplit(5)).mean()
```

0.767037799961446

XGB Regression

```
1 cross_val_score(xg, scaled_x_best, y, cv=ShuffleSplit(5)).mean()
```

0.7968013846376397

SV Regression

```
1 cross_val_score(SV, scaled_x_best, y, cv=ShuffleSplit(5)).mean()
```

-0.06457813564822405

Decision Tree Regression

```
1 cross_val_score(dt, scaled_x_best, y, cv=ShuffleSplit(5)).mean()
```

0.6614737269252287

Evaluation

Interpretation of the Results

Based on comparing Accuracy Score results with Cross Validation results, it is determined that Random Forest Regressor is the best model. It also has the lowest Root Mean Squared Error score.

Model/s Development and Evaluation

Hyper Parameter Tuning

- GridSearchCV was used for Hyper Parameter Tuning of the Random Forest Regressor model.

Based on the input parameter values and after fitting the train datasets The Random Forest Regressor model was further tuned based on the parameter values yielded from GridsearchCV.

- The Random Forest Regressor model displayed an accuracy of 83.46%. This model was then tested using a scaled Test Dataset.
- The model performed with good amount of accuracy.

```
1 from sklearn.model_selection import GridSearchCV
2 parameter = {'n_estimators':[30,60,80],'max_depth': [40,50,80],'min_samples_leaf':[5,10,20],'min_samples_split':[2,5,10],'criterion':['mse','mae']}
3 GridCV = GridSearchCV(RandomForestRegressor(),parameter,cv=ShuffleSplit(5),n_jobs = -1,verbose = 1)
4 GridCV.fit(x_train,y_train)
```

Fitting 5 folds for each of 486 candidates, totalling 2430 fits

```
GridSearchCV(cv=ShuffleSplit(n_splits=5, random_state=None, test_size=None, train_size=None),
             estimator=RandomForestRegressor(), n_jobs=-1,
             param_grid={'criterion': ['mse', 'mae'], 'max_depth': [40, 50, 80],
                          'max_features': ['auto', 'sqrt', 'log2'],
                          'min_samples_leaf': [5, 10, 20],
                          'min_samples_split': [2, 5, 10],
                          'n_estimators': [30, 60, 80]},
             verbose=1)
```

```
1 GridCV.best_params_
```

```
{'criterion': 'mse',
 'max_depth': 50,
 'max_features': 'auto',
 'min_samples_leaf': 5,
 'min_samples_split': 5,
 'n_estimators': 60}
```

```
1 Best_mod = RandomForestRegressor(n_estimators = 80,criterion = 'mse', max_depth= 80, max_features = 'auto',min_samples_leaf
2
3 Best_mod.fit(x_train,y_train)
```

```
RandomForestRegressor(criterion='mse', max_depth=80, min_samples_leaf=5,
                       n_estimators=80)
```

```
1 rfpred = Best_mod.predict(x_test)
2 acc = r2_score(y_test,rfpred)
3 print(acc*100)
```

83.460362103591

Random Forest Regressor has an accuracy of 83.46%


```

1 Prediction_accuracy = pd.DataFrame({'Predictions': mod.predict(scaled_x_best), 'Actual Values': y})
2 Prediction_accuracy.head(30)

```

	Predictions	Actual Values
0	2459.096777	2456
1	2452.485280	2456
2	2452.485280	2456
3	2452.485280	2456
4	2463.425381	2456
5	2470.533640	2456
6	2528.322961	2456
7	2528.322961	2456
8	2658.228352	2456
9	2658.228352	2456
10	2539.491026	2343
11	2531.843321	2343
12	2531.843321	2343
13	2539.721343	2343
14	2596.263468	2637
15	2583.768194	2637
16	2583.768194	2637
17	2588.684076	2637
18	2588.684076	2637
19	2605.416418	2637
20	3332.865026	2950
21	3550.566192	2950
22	3239.271243	3100
23	3507.979508	3100
24	3484.005072	3100
25	3555.445752	3863
26	3553.917939	3993
27	3990.100351	4018
28	4214.790966	4018

Model/s Development and Evaluation

In summary, Based on the visualizations of the feature-column relationships, it is determined that, Features like Source,month,Duration,Total Stops,Airline,Date are some of the most important features to predict the label values. Random Forest Regressor Performed the best out of all the models that were tested. It also worked well with the outlier handling.

Key Findings and Conclusions of the Study and Learning Outcomes with respect to Data Science

Based on the in-depth analysis of the Flight Price Prediction Project, The Exploratory analysis of the datasets, and the analysis of the Outputs of the models the following observations are made:

- Air Fare attributes like Date, Month, Duration, Total Stops etc play a big role in influencing the used Flight price.
- Airline Brand also has a very important role in determining the used Flight Ticket price.
- Various plots like Barplots, Countplots and Lineplots helped in visualising the Feature-label relationships which corroborated the importance of Air Fare features and attributes for estimating Flight Ticket Prices.
- Due to the Training dataset being very small, only very small amount of the outliers was removed to ensure proper training of the models.
- Therefore, Random Forest Regressor, which uses averaging to improve the predictive accuracy and controls over-fitting. performed well despite having to work on small dataset and produced good predictions that can be understood easily.

Learning Outcomes of the Study in respect of Data Science

- Data cleaning was a very important step in removing plenty of anomalous data from the huge dataset that was provided.
- Visualising data helped identify outliers and the relationships between target and feature columns as well as analysing the strength of correlation that exists between them.

Limitations of this work and Scope for Future Work

- A small dataset to work with posed a challenge in building highly accurate models. This project also relied heavily on historical data and was unable to account for various other factors that influence demand and ticket pricing like pandemic status affecting demand, government regulations on air travel, shifting in routes, weather conditions, etc.
- Most airline companies also do not publicly make available their ticket pricing strategies, which makes gathering price and air fare related data sets using web scraping the only means to build a dataset for building predicting models.
- Availability of more features and a larger dataset would help build better models.