



## FLIGHT PRICE PREDICTION PROJECT

Submitted by:

Anita Thapa

## **ACKNOWLEDGMENT**

I would like to express my special thanks of gratitude to my SME Swati Ma'am for giving me an opportunity to work on the flight price prediction project. During my research on this project internet sites like [towardsdatascience](#), [medium](#), [easemytrip](#), [yatra.com](#), and [Stack Overflow](#) helped me understand the main problem statement, also to collect data and in finding their solutions.

# INTRODUCTION

- **Business Problem Framing**

Model to predict flight price so that consumers book flight tickets at the best time and save the most by taking least risk.

- **Conceptual Background of the Domain Problem**

Airline Industry has picked up a rapid pace in terms of revenue as soon as the covid-19 curbs uplifted from the country. To understand the dynamics of the flight fare changes a model needs to be built. The model will help to analyze the factors which are affecting the rise and fall of the flight fares. The prediction will be useful for consumers to book flight at the best time so that they save the most by taking the least risk.

- **Review of Literature**

The research was done on Flight prices in various online websites like easemytrip, yatra.com etc. in order to collect the data using the selenium library. Using websites like medium, towardsdatascience and Kaggle helped in modeling the dataset and in predicting the airfare. Stack overflow website helped to resolve problems faced during data scraping and modeling.

- **Motivation for the Problem Undertaken**

The main objective was to analyse the airfare changes and to model a flight price valuation prediction. In order to understand the airline industry and their price trends this project was taken up so that the consumers can book their tickets at the best time with more saving.

# Analytical Problem Framing

- Mathematical/ Analytical Modeling of the Problem

The data was present in excel format so by using pandas the excel file was read. Then the dataset details were analysed using .shape() to find the number of rows and columns, using dtypes() we found the dataset datatype like object and numeric datatype. Dataset has 4749 Rows and 12 Columns; 2 columns are numeric and rest 10 columns are object type and the Target column is Price. Dataset is a Regression model.

```
1 #Importing Libraries:
2 import numpy as np
3 import pandas as pd
```

```
1 #Reading excel file and converting it in dataframe:
2 ds=pd.read_excel("C:\\Users\\Asus\\Desktop\\flights_data.xlsx")
3 df=pd.DataFrame(ds)
4 df.head()
```

Unnamed: 0	Airline	Journey Date	Source	Destination	Route	Dep_Time	Arrival_Time	Duration	Total_Stops	Flight_no	Price	
0	0	GO FIRST	26/02/2022	Delhi	Mumbai	DEL → BOM	19:50	22:10	02h 20m	non-stop	G8-323	7740
1	1	Indigo	26/02/2022	Delhi	Mumbai	DEL → BOM	20:00	22:15	02h 15m	non-stop	6E-6261	7845
2	2	Vistara	26/02/2022	Delhi	Mumbai	DEL → BOM	19:50	22:00	02h 10m	non-stop	UK-985	8790
3	3	Air India	26/02/2022	Delhi	Mumbai	DEL → BOM	20:00	22:10	02h 10m	non-stop	AI-805	8790
4	4	Indigo	26/02/2022	Delhi	Mumbai	DEL → BOM	21:15	23:30	02h 15m	non-stop	6E-6722	9945

```
1 # Rows & Columns in Train dataset:
2 df.shape
```

```
(4749, 12)
```

```
1 4749 rows and 12 columns
```

```
1 # Datatype of dataset
2
3 df.dtypes
```

```
Unnamed: 0      int64
Airline         object
Journey Date    object
Source          object
Destination      object
Route           object
Dep_Time        object
Arrival_Time    object
Duration        object
Total_Stops     object
Flight_no       object
Price           int64
dtype: object
```

- Data Sources and their formats

Data was scraped from websites like easemytrip and yatra.com using Selenium library in Python and saved as excel file. Then data was imported in python using pandas . Information of dataset is found

using `datasetname.info()`. As per the information each column has count of 4749 rows as shown in below figure.

```
1 # Information about Dataset:
2
3 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4749 entries, 0 to 4748
Data columns (total 12 columns):
#   Column              Non-Null Count  Dtype
---  -
0   Unnamed: 0           4749 non-null   int64
1   Airline              4749 non-null   object
2   Journey Date         4749 non-null   object
3   Source               4749 non-null   object
4   Destination          4749 non-null   object
5   Route               4749 non-null   object
6   Dep_Time             4749 non-null   object
7   Arrival_Time         4749 non-null   object
8   Duration             4749 non-null   object
9   Total_Stops          4749 non-null   object
10  Flight_no            4749 non-null   object
11  Price                4749 non-null   int64
dtypes: int64(2), object(10)
memory usage: 445.3+ KB
```

- Data Pre-processing Done

Dataset had 12 columns and 4749 rows. Journey Date converted to day and month. Columns named Dep\_Time, Arrival\_Time and Duration converted into hours and minutes. Repeated values replaced in Airline, Total\_Stops and Destination column.

Converting object to datetime datatypes

```
1 # day separated from date
2 df["Journey_day"] = pd.to_datetime(df["Journey Date"], format="%d/%m/%Y").dt.day

1 # month separated from date
2 df["Journey_month"] = pd.to_datetime(df["Journey Date"], format="%d/%m/%Y").dt.month

1 df["Journey_month"].value_counts()

3    2653
4    1185
2     911
Name: Journey_month, dtype: int64
```

DEP\_TIME TO HOURS AND MINUTE

```
1 #departure hour extracted from departure time.
2 df["Depart_Hour"] = pd.to_datetime(df["Dep_Time"]).dt.hour
3
4 #departure minute extracted from departure time.
5 df["Depart_Min"] = pd.to_datetime(df["Dep_Time"]).dt.minute

1 # Arrival time is separated into arrival hour and arrival minute:
2
3 df["Arrival_hour"] = pd.to_datetime(df.Arrival_Time).dt.hour
4
5 df["Arrival_min"] = pd.to_datetime(df.Arrival_Time).dt.minute
6
```

```
1 # Separating hours and minutes from duration column from the dataset.
2 duration = list(df["Duration"])
3
4 for i in range(len(duration)):
5     if len(duration[i].split()) != 2: # Check if duration contains only hour or mins
6         if "h" in duration[i]:
7             duration[i] = duration[i].strip() + " 0m" # Adds 0 minute
8         else:
9             duration[i] = "0h " + duration[i] # Adds 0 hour
10
11 duration_hours = []
12 duration_mins = []
13 for i in range(len(duration)):
14     duration_hours.append(int(duration[i].split(sep="h")[0]))
15     duration_mins.append(int(duration[i].split(sep="m")[0].split()[-1]))
```

```
1 #replacing repeated values:
2 df["Airline"].replace({"IndiGo": "Indigo", "GO FIRST": "Go First", "Air Asia": "AirAsia"}, inplace = True)
```

```
1 df["Destination"].replace({"New Delhi": "Delhi"}, inplace = True)
```

- Data Inputs- Logic- Output Relationships

Output Column is Price column and rest other columns are Input Columns. Using Correlation we can find out the relationship between Input and output columns. Total\_Stops column is positively correlated with Price column.

Total_Stops	1.00	0.41	-0.05	0.06	-0.09	0.03	0.08	0.03	0.55	-0.02	-0.01	-0.01
Price	0.41	1.00	-0.28	-0.31	0.00	0.02	0.11	0.06	0.37	0.01	-0.09	-0.01
Journey_day	-0.05	-0.28	1.00	-0.24	-0.02	-0.02	0.01	0.01	-0.06	-0.03	0.03	0.03
Journey_month	0.06	-0.31	-0.24	1.00	-0.01	-0.03	-0.02	-0.01	0.01	0.02	-0.04	0.01
Depart_Hour	-0.09	0.00	-0.02	-0.01	1.00	0.05	-0.11	0.04	0.09	-0.01	-0.01	0.02
Depart_Min	0.03	0.02	-0.02	-0.03	0.05	1.00	0.02	0.05	0.00	-0.04	-0.02	-0.01
Arrival_hour	0.08	0.11	0.01	-0.02	-0.11	0.02	1.00	-0.02	0.02	0.01	-0.11	-0.01
Arrival_min	0.03	0.06	0.01	-0.01	0.04	0.05	-0.02	1.00	0.06	0.07	-0.04	-0.01
Duration_hours	0.55	0.37	-0.06	0.01	0.09	0.00	0.02	0.06	1.00	-0.07	-0.13	-0.11

- State the set of assumptions (if any) related to the problem under consideration

In Total\_Stops column 2stops,3 stops and 4 stops were assumed as 2+ stops

- Hardware and Software Requirements and Tools Used

Libraries used were:

1. Numpy : It is a Numerical Python library used for numerical computation like arrays in Python.
2. Panda: It was used for reading and converting excel/csv file into dataframe.
3. matplotlib: this library was used for plotting the graph in the EDA.
4. Seaborn: This library is also used for visualization as it has helped to plot different types of plots like countplot, displot, boxplot and heatmap in the notebook.
5. Label encoder: It helped to encode the object datatype into numeric datatype as the Machine learning algorithm works on numeric datatype only.

6. sklearn.preprocessing and power transform: It was used to remove the skewness from the dataset.
7. Standard Scaler : It was used to scale the feature column of the dataset before model selection.
8. sklearn: This library was used to import train\_test\_split, metrics like accuracy score, classification report and also importing Algorithms likes Linear Regression, Decision Tree Regressor, SVR, Random Forest Regressor and K-Neighbors Regressor..
9. pickle: This library was used to save the final model.

# Model/s Development and Evaluation

- Identification of possible problem-solving approaches (methods)

The Standard Scaled data was used in Model building. As the model is Regressor so Linear Regressor library along with r2 score, cross validation score, means square error metrics were imported. Then 4 Algorithms were used which were Decision Tree Regressor, SVR, Random Forest Regressor and K-Neighbors Regressor. Then after selecting best model, then the final model was deployed and the accuracy of model was 81.9% as it predicted values approximately equal to actual values.

## Model Building:

```
1 # Model selection:
2
3 from sklearn.metrics import r2_score
4 from sklearn.linear_model import LinearRegression
5 lr=LinearRegression()
6 from sklearn.metrics import mean_squared_error,mean_absolute_error
7 from sklearn.metrics import accuracy_score
8 from sklearn.model_selection import train_test_split
9
10 import warnings
11 warnings.filterwarnings('ignore')

```

```
1 for i in range(0,100):
2     train_x,test_x,train_y,test_y=train_test_split(x,y,test_size=0.2,random_state=i)
3     lr.fit(train_x,train_y)
4     pred_train=lr.predict(train_x)
5     pred_test=lr.predict(test_x)
6     print(f"At random state {i},the training accuracy is:- {r2_score(train_y,pred_train)}")
7     print(f"At random state {i},the testing accuracy is:- {r2_score(test_y,pred_test)}")
8     print("\n")

```

- Testing of Identified Approaches (Algorithms)

Algorithms used for the training and testing are:

1. XGBoost Regressor
2. Decision Tree Regressor
3. SVR
4. Random Forest Regressor
5. K-Neighbors Regressor

- Run and Evaluate selected models

Out of all the 5 algorithm used the best algorithm used is Random Forest Regressor as the accuracy score is more after using hyperparameter.



### 3. Random forest regressor

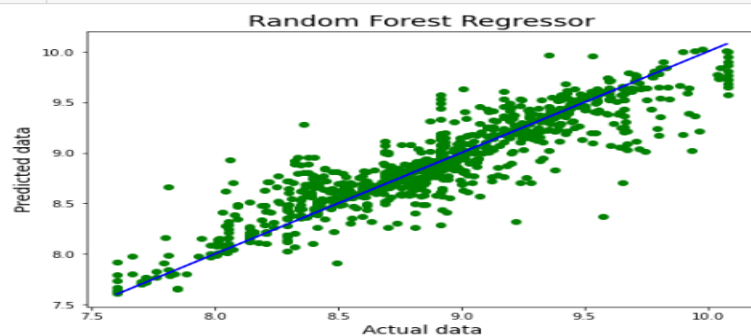
```
# Importing Libraries and Hyper parameter tuning:
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestRegressor
parameters = {'criterion':['mse', 'mae'], 'bootstrap': [True, False], 'max_depth': [5, 10, None], 'max_features': ['auto']}
rf = RandomForestRegressor()
clf = GridSearchCV(rf, parameters, cv=6, n_jobs=-1, verbose=2, refit=False)
clf.fit(train_x, train_y)
print(clf.best_params_)
```

Fitting 6 folds for each of 240 candidates, totalling 1440 fits  
{'bootstrap': True, 'criterion': 'mse', 'max\_depth': None, 'max\_features': 'auto', 'n\_estimators': 15}

```
rf = RandomForestRegressor(criterion="mse", max_features="auto", max_depth=None, bootstrap=True, n_estimators=15)
rf.fit(train_x, train_y)
rf.score(train_x, train_y)
predrf = rf.predict(test_x)
print('rf score', rf.score(train_x, train_y))
rfs = r2_score(test_y, predrf)
print('R2 Score:', rfs*100)
rfscore = cross_val_score(rf, x, y, cv=6)
rfc = rf.score.mean()
print('Cross Val Score:', rfc*100)
print('Mean squared error:', mean_squared_error(test_y, predrf))
print('Mean absolute error:', mean_absolute_error(test_y, predrf))
print('RMSE:', np.sqrt(mean_squared_error(test_y, predrf)))
```

rf score: 0.9719211732040074  
R2 Score: 81.9609580888492  
Cross Val Score: 30.15636213785588  
Mean squared error: 0.04764347710058702  
Mean absolute error: 0.14673604760568396  
RMSE: 0.21827385803294683

```
import matplotlib.pyplot as plt
plt.figure(figsize=(8,6))
plt.scatter(x=test_y, y=predrf, color='g')
plt.plot(test_y, test_y, color='b')
plt.xlabel('Actual data', fontsize=14)
plt.ylabel('Predicted data', fontsize=14)
plt.title('Random Forest Regressor', fontsize=18)
plt.show()
```

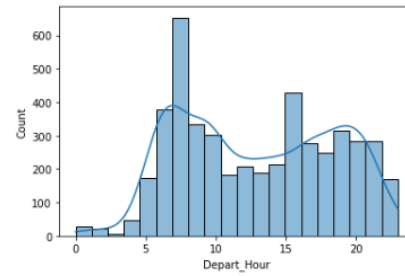
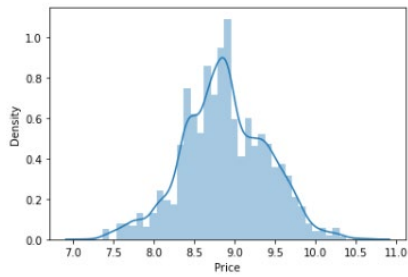


- Key Metrics for success in solving problem under consideration

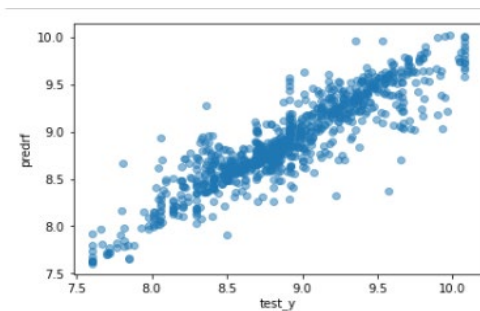
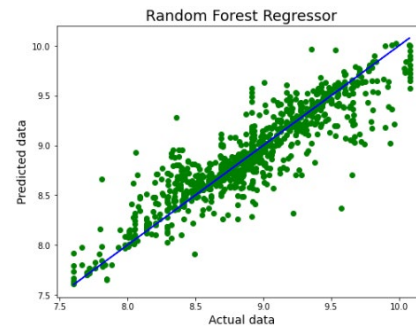
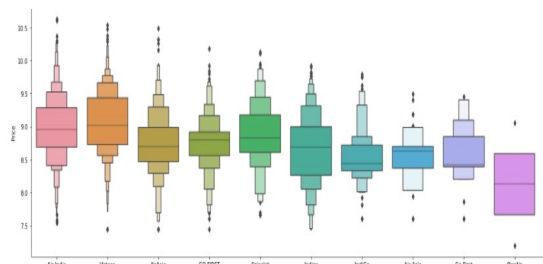
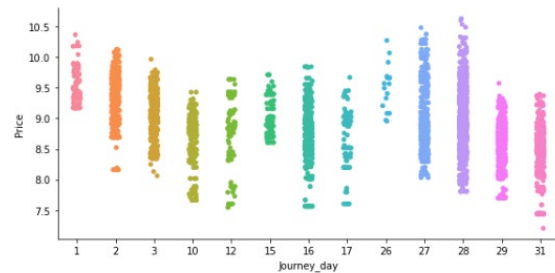
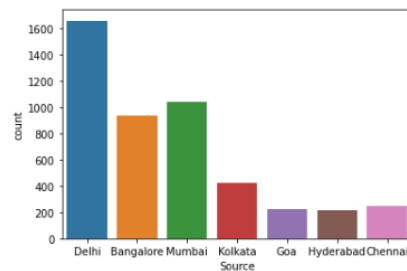
The main metric used was to find out the accuracy of the model which helped to determine the best model is comparing r2 score with the cross validation score of the algorithm. For model selection the algorithm must have small difference between cross validation score and r2 score then that algorithm is considered to be the best model.

- Visualizations

Different plots were used in the problem like displot, countplot, histplot, heatmap, boxplot and pyplot using matplotlib and seaborn library as shown in the images below.



Number of flights by source



Total_Stops	1.00	0.41	-0.05	0.06	-0.09	0.03	0.08	0.03	0.55	-0.02	-0.01	-0.0
Price	0.41	1.00	-0.28	-0.31	0.00	0.02	0.11	0.06	0.37	0.01	-0.09	-0.0
Journey_day	-0.05	-0.28	1.00	-0.24	-0.02	-0.02	0.01	0.01	-0.06	-0.03	0.03	0.0
Journey_month	0.06	-0.31	-0.24	1.00	-0.01	-0.03	-0.02	-0.01	0.01	0.02	-0.04	0.0
Depart_Hour	-0.09	0.00	-0.02	-0.01	1.00	0.05	-0.11	0.04	0.09	-0.01	-0.01	0.0
Depart_Min	0.03	0.02	-0.02	-0.03	0.05	1.00	0.02	0.05	0.00	-0.04	-0.02	-0.0
Arrival_hour	0.08	0.11	0.01	-0.02	-0.11	0.02	1.00	-0.02	0.02	0.01	-0.11	-0.0
Arrival_min	0.03	0.06	0.01	-0.01	0.04	0.05	-0.02	1.00	0.06	0.07	-0.04	-0.0
Duration_hours	0.55	0.37	-0.06	0.01	0.09	0.00	0.02	0.06	1.00	-0.07	-0.13	-0.1

## • Interpretation of the Results

From the visualizations, preprocessing and modeling of the data it was interpreted that morning flights are expensive and the Airline Air India and Vistara are expensive. If stops are more than prices also more. Consumers should purchase tickets during mid of month as they are at less price. As per the model the data can be predicted with 81.9% accuracy.

# CONCLUSION

- Key Findings and Conclusions of the Study

Key Findings of the project are:

1. Flight fare increases with increase in the number of stops.
2. Flights with duration 5 hours and above have higher prices.
3. Vistara and Air India are expensive airlines.
4. Flight fares are high during start of the month but fare is less in the mid and end last two-three days of the month.
5. Delhi source flights and Goa destination flights have higher prices.

- Learning Outcomes of the Study in respect of Data Science

Using Visualization libraries like matplotlib and seaborn it was easy to find the correlation between the dataset columns and the plotting of each column along the target column helped to analyse the problem pretty well. As visualization helped to get better insight of the data like skewness present in the data.

- Limitations of this work and Scope for Future Work

As the model accuracy is 81.9% so it can approximately predict the correct value as sometimes it might predict wrong values but chances are less.