



# **FLIGHT PRICE PREDICTION**

Submitted by: Naveen

# **ACKNOWLEDGEMENT**

The flight price prediction project has been completed through the aid of many factors, we have searched over the web for appropriate websites in order to collect the latest data, In that process we landed up in finding website yathra.com where we got almost 1500 data which belongs to 6 major aviation companies of India, The data collection part was challenging because the price of the ticket depends upon several factors like class, distance, no of travels etc., We webscrapped the data using selenium ,We took the help of software like tableau, Microsoft power bi for data analysis and Jupyter for model building.

In the process of making an informative, productive project we traversed through many websites out of which towardsdatascience, geeksforgeeks, stackoverflow are significant one, we also got required aid by flirobo and datatrained company.

# **INTRODUCTION**

- **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

1. Time of purchase patterns (making sure last-minute purchases are expensive)
2. Keeping the flight as full as they want it (raising prices on a flight which is filling up in order to reduce sales and hold back inventory for those expensive last-minute expensive purchases) So, you have to work on a project where you collect data of flight fares with other features and work to make a model to predict fares of flights.

- **Conceptual Background of Domain Problem**

First of all we should be knowing how to scrap the required data from various websites through the techniques of webscrapping, In order to do that we should extract all the urls of each ticket using relevant xpath's , we should be knowing how to build a linear regression model through collected data.

- **Literature Review**

To get the enough amount of relevant data we have to crawl through google to get perfect website, In that process we ended up landing in yathra.com, We have used many other websites like stackoverflow, towardsdatascience,geeksforgeeks for many queries on the project.

- **Motivation for the Problem Undertaken**

The motivation for this project is to provide our client with new model to excel in their business and also to upgrade myself by solving the real world problem through virtual environment.

# Analytical Problem Framing

- **Mathematical and Analytical understanding of the problem**

The data we collected was in terms of rows and columns, there are 10 columns and 1587 rows out of which 3 columns were of object type and remaining are integer type, so in order to make the machine understand this we have to convert the categorical value to numerical values we accomplish it through label encoding or one hot encoding.

We did scaling of the dataset in order to have values of variables within certain limits so that machine can perform better, for this purpose we used different types of scaling namely standard scaling, min max scaling and Robust scaler, later we reduce the skewness using various techniques.

## Importing the dataset ¶

```
In [1]: import pandas as pd
df=pd.read_csv('flight_tableau.csv')
df.head()
```

```
Out[1]:
```

	Unnamed: 0	Company Name	No of stops	No of days in advanced booked	class	where to where	Route Value	Departure Time	Duration	Price	distance in km
0	0	Air Asia	1	1	Economy	Delhi - Mumbai	3	8.00	6.58	5953	1148
1	1	Air Asia	1	1	Economy	Delhi - Mumbai	3	9.42	6.58	5953	1148
2	2	Air Asia	1	1	Economy	Delhi - Mumbai	3	12.67	7.58	5953	1148
3	3	Air Asia	1	1	Economy	Delhi - Mumbai	3	11.92	8.33	5953	1148
4	4	Air Asia	1	1	Economy	Delhi - Mumbai	3	8.00	8.58	5953	1148

- **Data Source and Their Format**

The major portion of the data we collected from car trade.com , there are 10 major parameters and rows out which 3 are of categorical type column and remainings are of categorical type columns.

```
: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1587 entries, 0 to 1586
Data columns (total 10 columns):
 #   Column                                Non-Null Count  Dtype  
---  -
 0   Company Name                         1587 non-null   object  
 1   No of stops                          1587 non-null   int64   
 2   No of days in advanced booked        1587 non-null   int64   
 3   class                               1587 non-null   object  
 4   where to where                       1587 non-null   object  
 5   Route Value                          1587 non-null   int64   
 6   Departure Time                       1587 non-null   float64  
 7   Duration                             1587 non-null   float64  
 8   Price                                1587 non-null   int64   
 9   distance in km                       1587 non-null   int64   
dtypes: float64(2), int64(5), object(3)
memory usage: 124.1+ KB
```

## • Data Processing

There were enough amount of anomalies in the data, we looked column by column for anomalies and we cleaned each column separately , through the aid of pandas, loops and regular expressions we achieved it.

1.We imported the collected data to a separate notebook which was meant for data preprocessing

2.The ticket price column contains the data which were of strings type and we adopted regular expressions to clean it out.

```
In [1]: import pandas as pd
import numpy as np
```

```
In [2]: # Creating the dataframes
dataframes=[]
for i in range(1,73):
    k1='df'
    k2=str(i)
    k=k1+k2
    dataframes.append(k)
```

```
In [3]: # Creating the datafiles name
flight=[]
for i in range(1,73):
    k1='flight_data_'
    k2=str(i)
    k3='.csv'
    k=k1+k2+k3
    flight.append(k)
```

```
In [4]: #importing all the dataframe
k=0
for i in dataframes:
    locals()[i]=pd.read_csv(flight[k])
    k=k+1
```

```
In [5]: df_list=[]
for i in dataframes:
    df_list.append(locals()[i])
```

```
In [6]: df=pd.DataFrame()
df = df.append(df_list)
```

3.The Name of Company column contained a string out of which we have to extract only name of the just name of the company, it was a challenging one and we accomplished it through regular expressions.

- **Data Input-Logic-Output Relationships**

**Data Input :** It was a precise dataset, which has been scaled having low skewness and very minimal outliers.

**Logic:** The logic here is linear regression algorithm which predict the response variable, the linear regression algorithm we used in this case are

- 1.Linear Regression
- 2,Lasso Regression
- 3.Ridge Regression
- 4.ElasticNet
- 5.Ransac Regressor
- 6.Support vector Regressor
- 7.Random Forest Regressor

**Data output :**We got a model predicting the price of second hand car

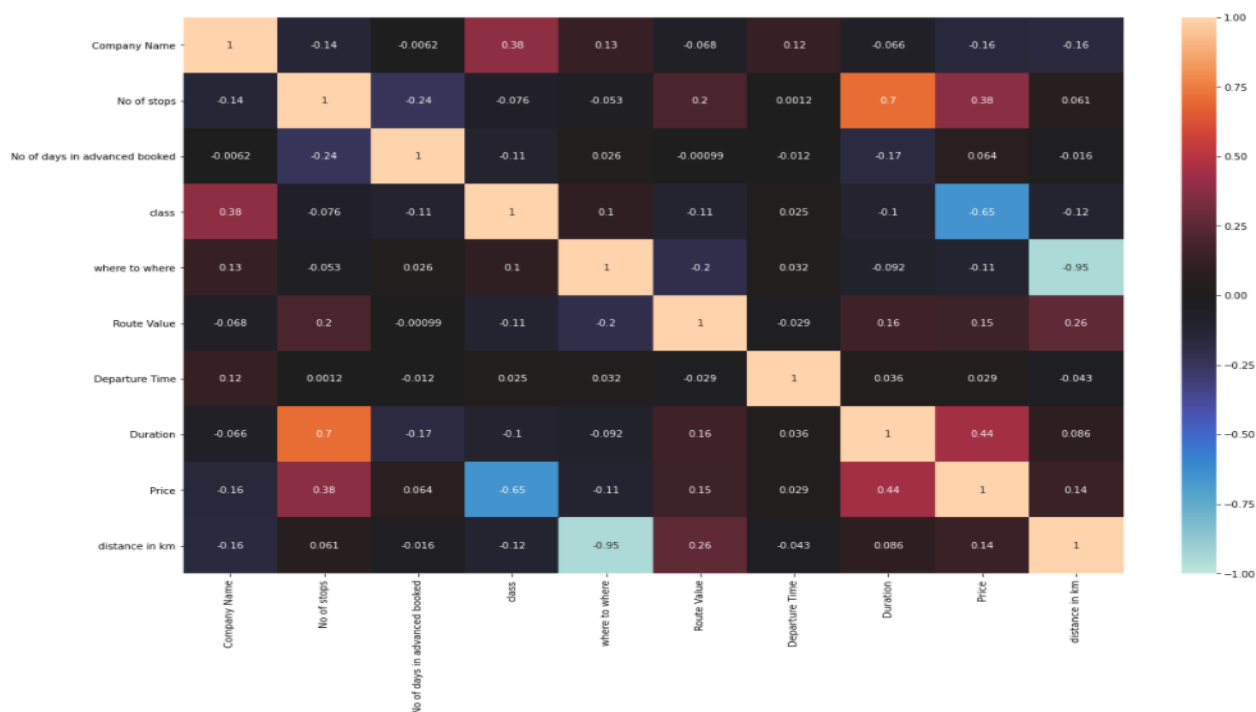
- **Hardware and software tools used**

**Hardware :** We used the hardware of 8GB RAM,1Tb ROM and i5 processor.

**Software :** For better visualization we used Tableau Public, Jupyter notebook from anaconda navigator for coding and webscrapping and Microsoft word and Power Point Presentation for creating the report and making the presentation respectively.

## Exploratory Data Analysis

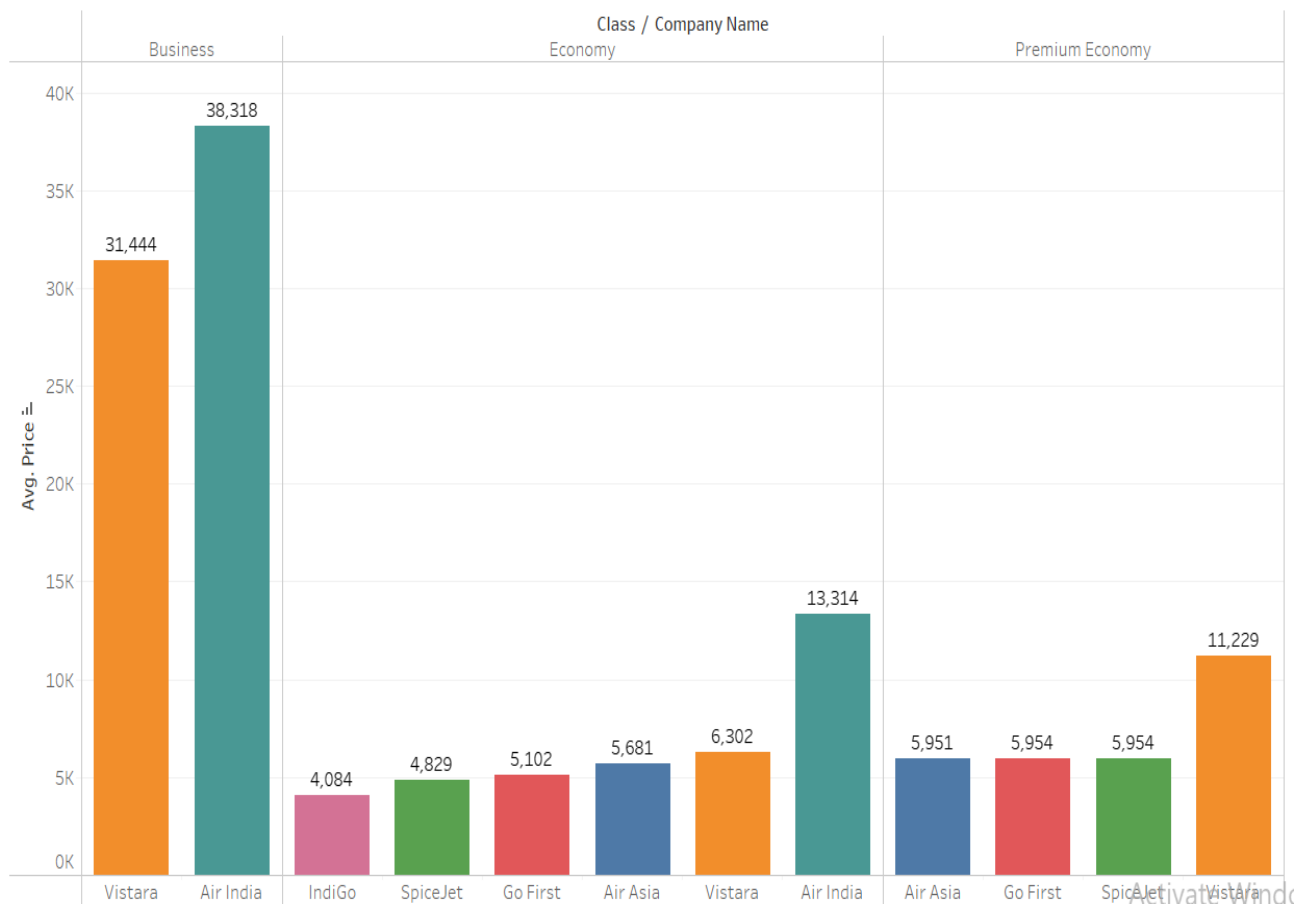
The following are the interpretation we got from data analysis, For analyzing the data we took the help of data visualization libraries like seaborn, Matplotlib, plotly,Tableau software, First we drew the correlation and heatmaps from which we came to know which are the significant variables which decides the flight price.



From the above heatmap we came to know about those independent variable which has got significant effect over the response variable, But here we have a problem in analyzing suppose say we want to know how the advanced booking decides the price but there are factors due to which we can't find directly because we have to give class as hue as well as destination, but providing two hue values and analyzing is bit difficult in python so we did whole analysis part in Tableau.

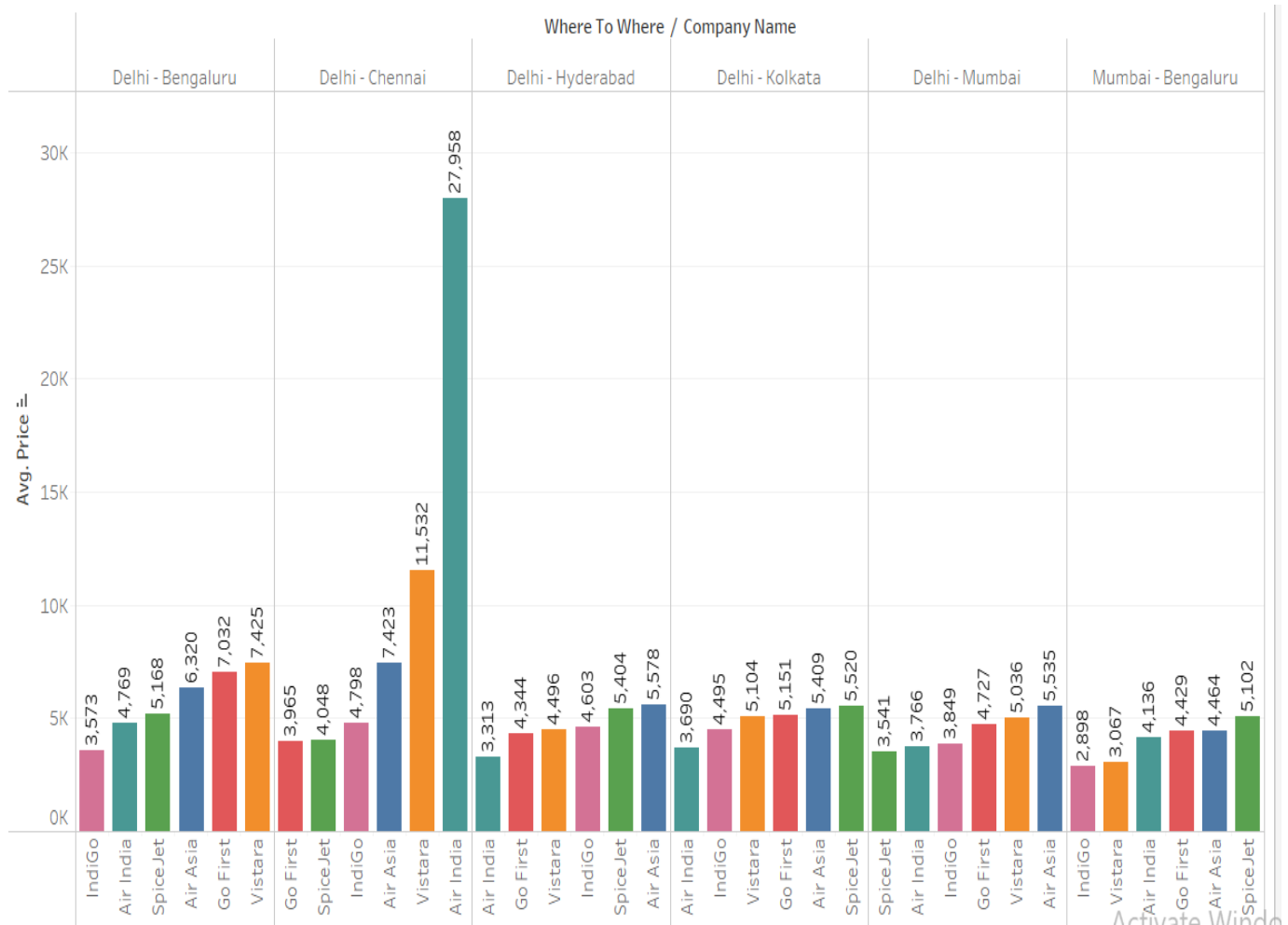
To analyze in tableau first we have prepared set of questions and we tried to find answers for those in tableau.

1. Average of each class price based on company name



The above graphs clearly indicates that if you want to travel in Business class choose vistara it has least average price,and for Economy class choose IndiGo, and for Premium Economy go for Air Asia.

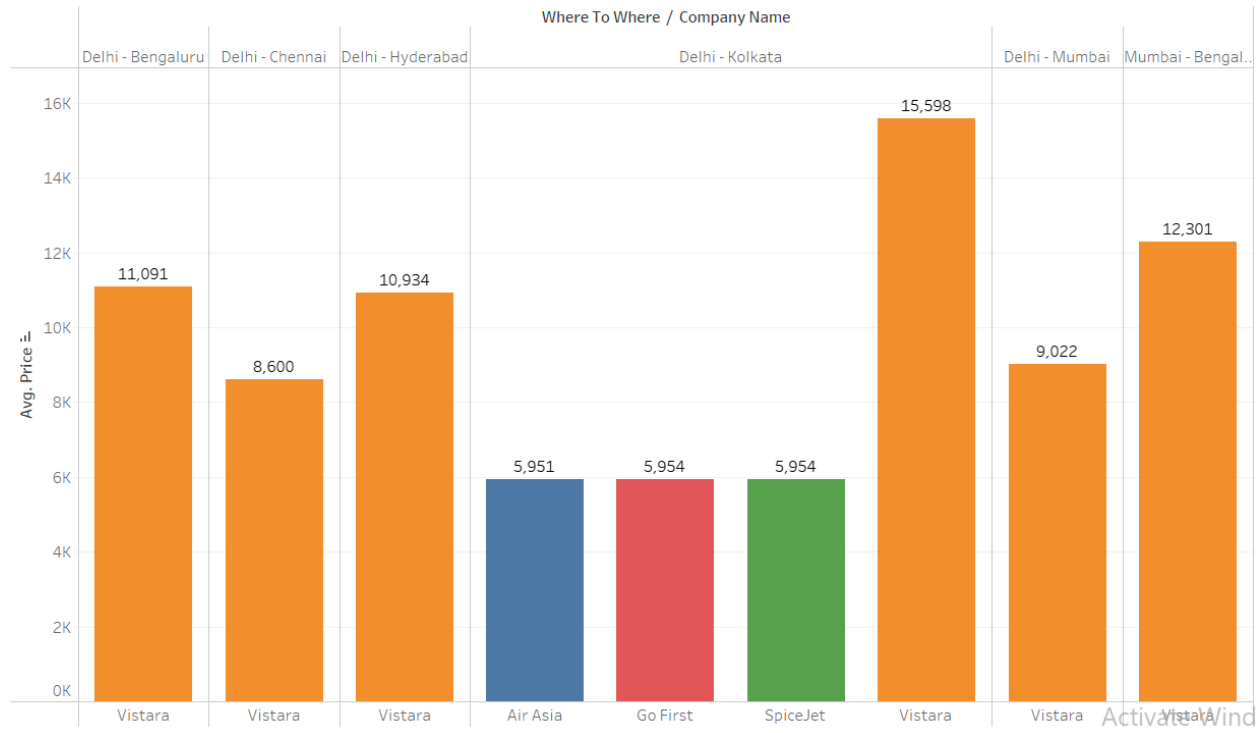
## 2. What's the least price to travel from Delhi and Mumbai to other locations in Economy class



- If you want to travel from Delhi to Bengaluru in Economy class then IndiGo will offer least price followed by Air India and SpaceJet and Vistara is the costly one.
- If you want to travel from Delhi to Chennai in Economy class then Go First will offer least price followed by SpaceJet and IndiGo and Air India is the costly one.
- If you want to travel from Delhi to Hyderabad in Economy class then Air India will offer least price followed by GoFirst and Vistara and Air Asia is the costly one.
- If you want to travel from Delhi to Kolkata in Economy class then Air India will offer least price followed by IndiGO and Vistara and SpaceJet is the costly one.
- If you want to travel from Delhi to Mumbai in Economy class then SpaceJet will offer least price followed by Air India and IndiGo and Air Asia is the costly one.
- If you want to travel from Mumbai to Bengaluru in Economy class then IndiGo will offer least price followed by Vistara and Air India and SpaceJet is the costly one.

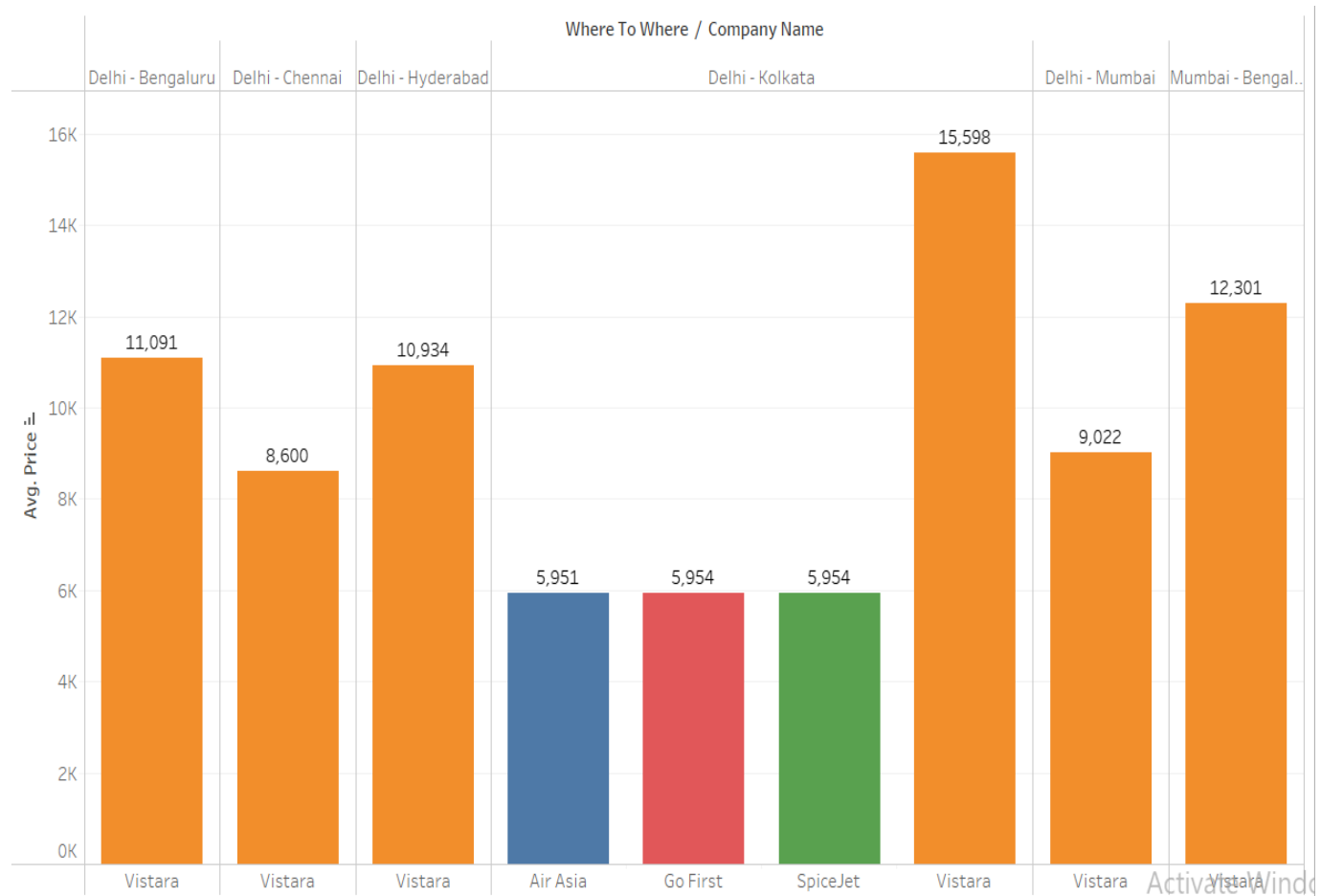


### 3.What's the least price to travel from Delhi and Mumbai to other locations in Premium Economy class



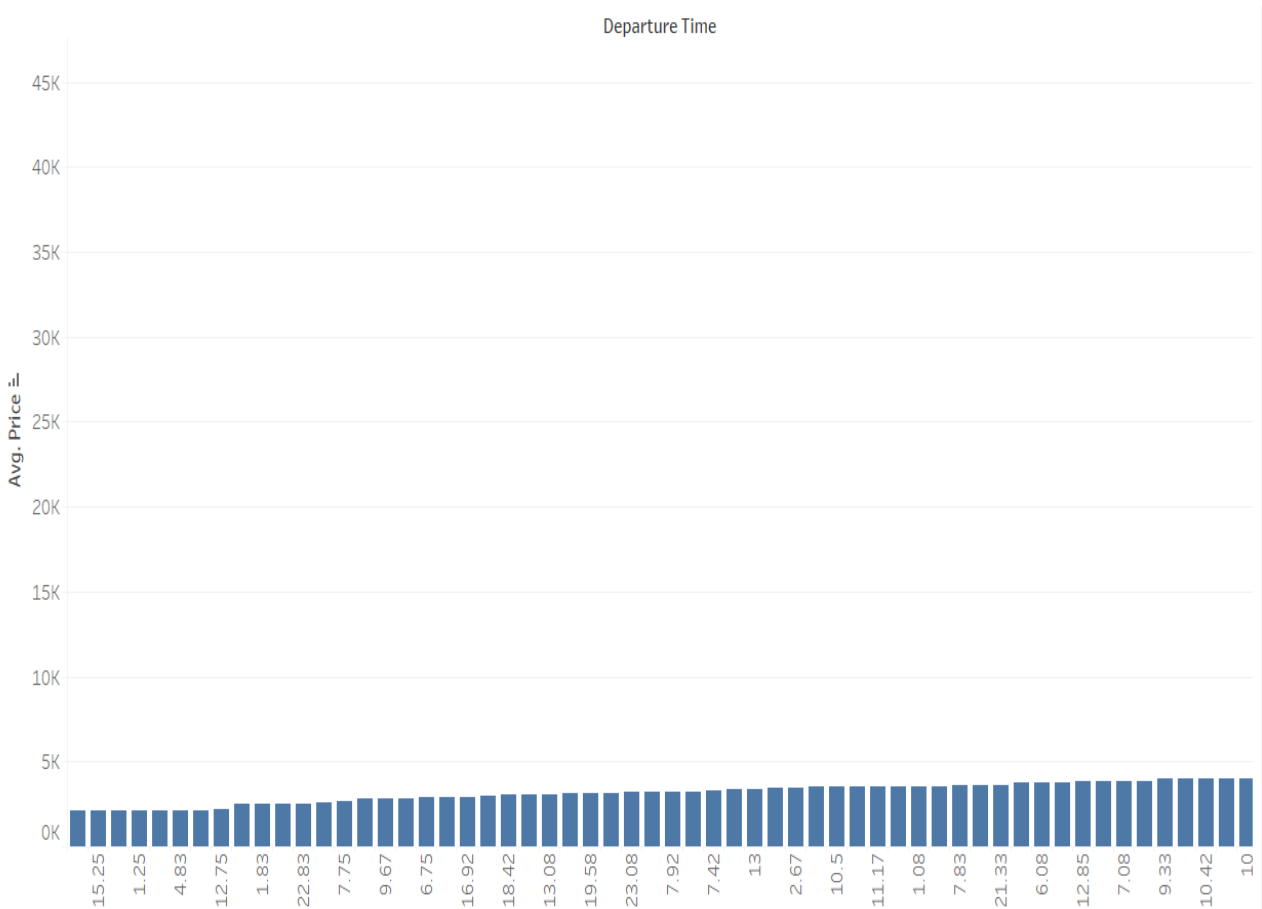
- If you want to travel from Delhi to Bengaluru in PremiumEconomy class then Vistara is the best choice
- If you want to travel from Delhi to Chennai in Premium Economy class then Vistara is the best choice
- If you want to travel from Delhi to Hyderabad in Premium Economy class then Vistara
- If you want to travel from Delhi to Kolkata in Premium Economy class then Air Asia is the cheapest one followed by GoFirst and SpiceJet
- If you want to travel from Delhi to Mumbai and Mumbai to Bengaluru in Premium Economy class then Vistara is the best choice.

#### 4.What's the least price to travel from Delhi and Mumbai to other locations in Business class



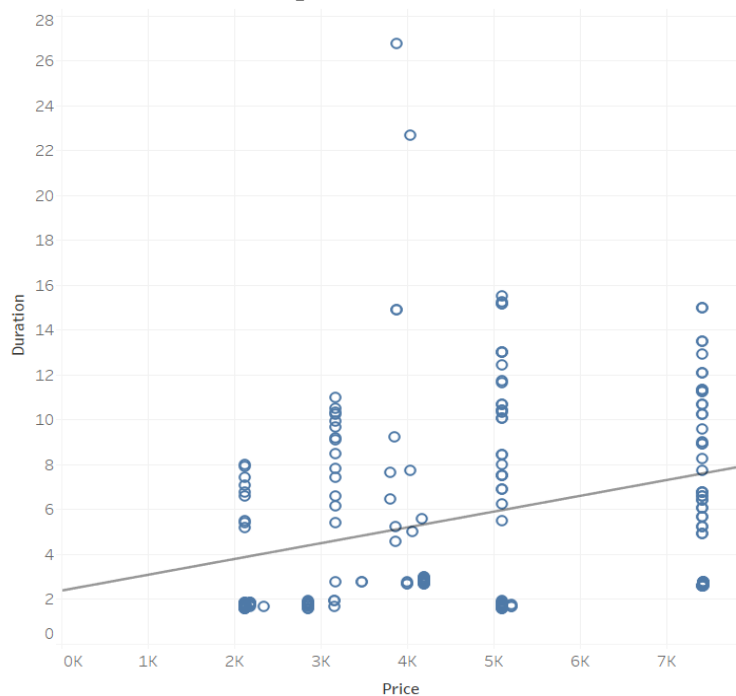
- If you want to go from Delhi to Bengaluru Chennai or Hyderabad or Mumbai or Mumbai to Bengaluru then vistara is the only available as per this data at this moment, It may vary on coming days.
- If you want to go from Delhi to Kolkata in Business class then Air Asia is the cheapest one followed by GoFirst and SpaceJet and Vistara is the costliest one according to this data.

## 5. What's the general departure time at which Economy class will have low price comparatively?



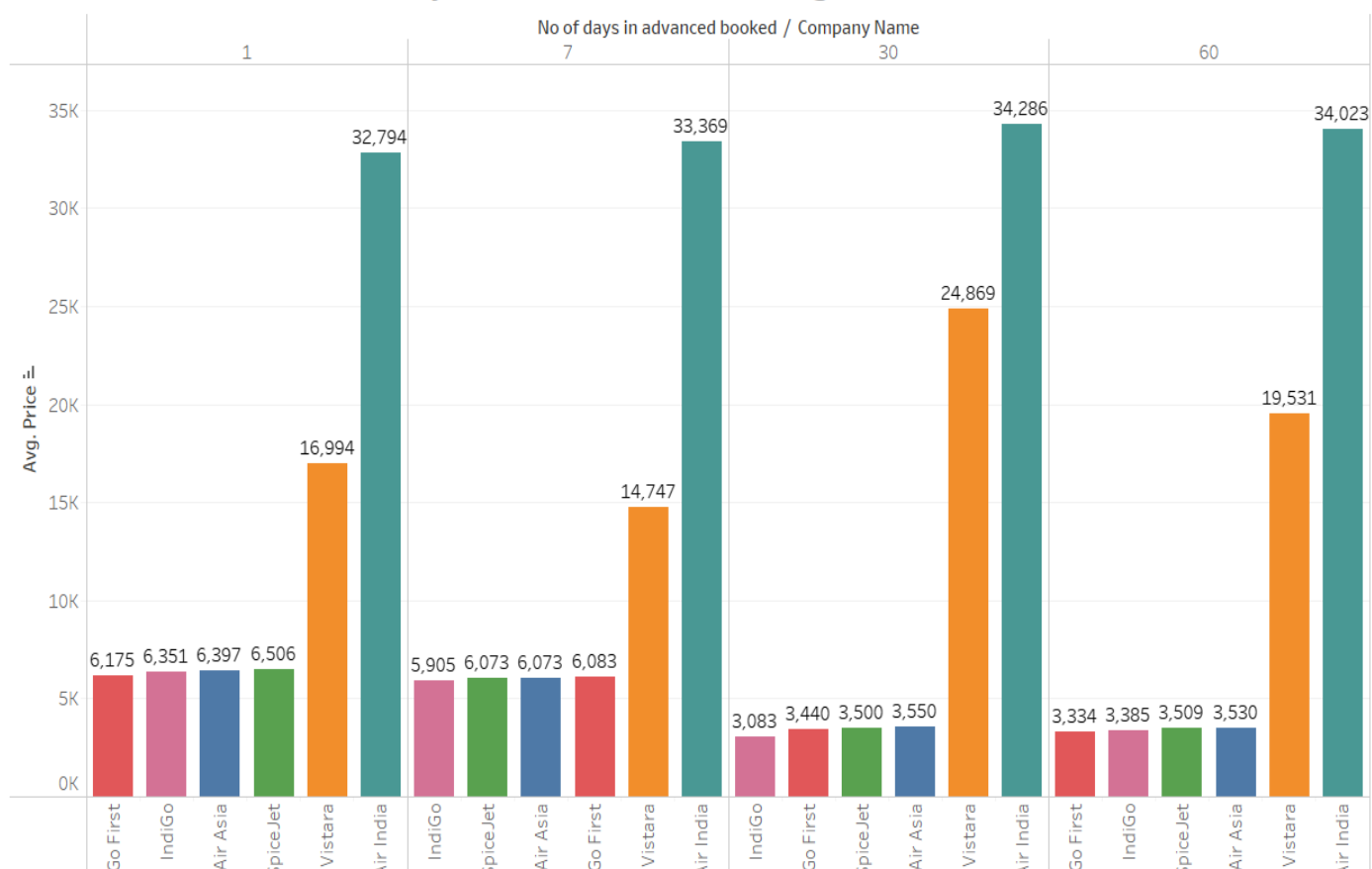
- The above graph depicts that if you want to book for Economy class then if you choose 4am-5am in morning you have maximum chance that your booking price will be lowest in that whole day.
- The highest price is generally observed for Economy class is around 3-5pm.
- For premium economy class the least price observed around 11am to 12pm and costly price is observed between 3pm -4 pm
- For Business class the least price observed around 11am to 12pm and costly price is observed between 3pm -4 pm

## 6. How duration and price are related based on class and destination?



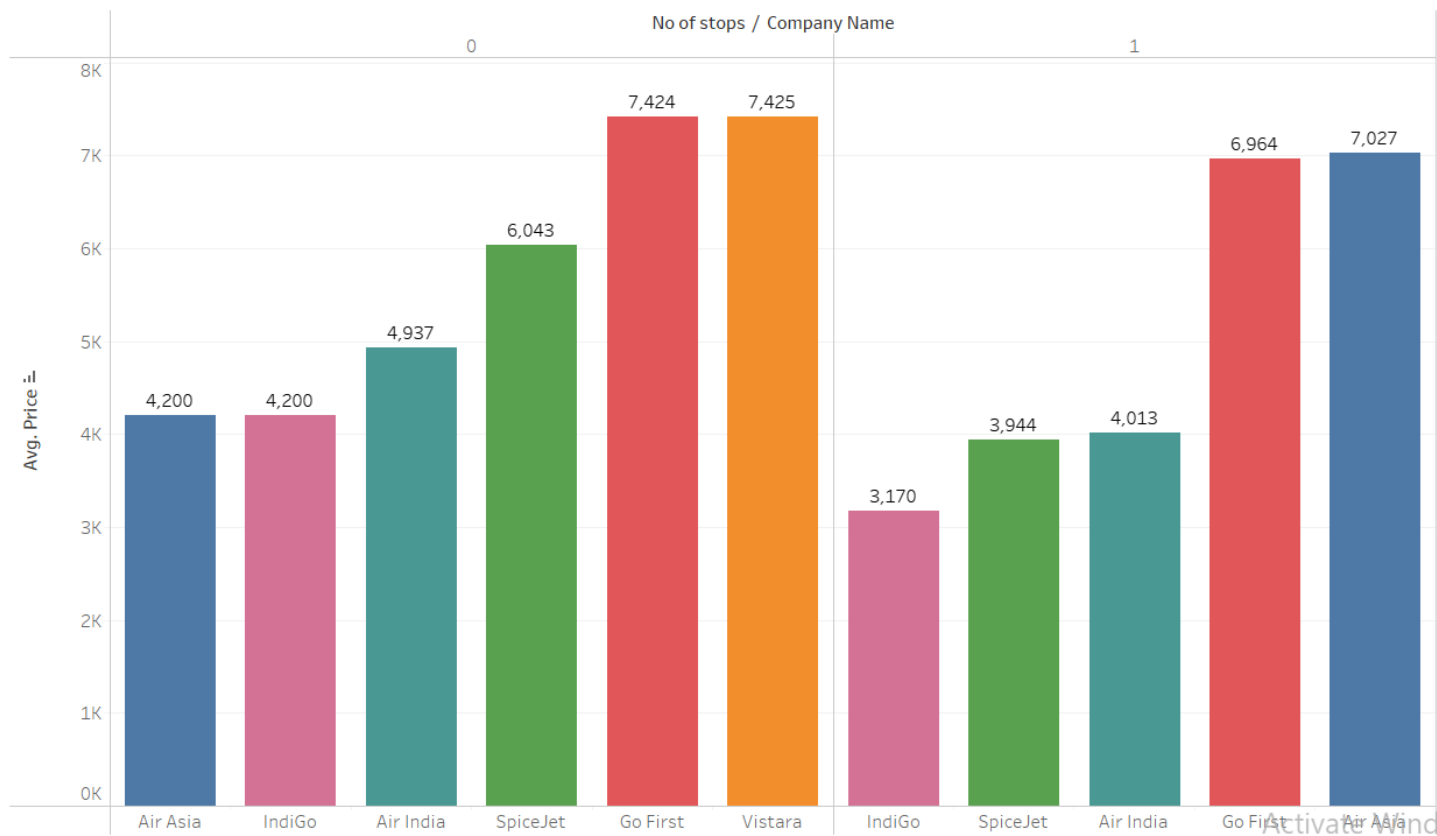
- From the above graph we can clearly see that as duration increases the price is also going to be increased.

## 7. How advanced booking affects the price of the ticket?



- If you book any ticket in an advance of 7 days you will be profited by 4.37% of cost of booking one day in advance
- If you book any ticket in an advance of 30 days you will be profited by 51% of cost of booking one day in advance
- There won't be much difference if you book 30 or 60 days in advanced.

8. How number of stops varies the price of the ticket to the same destination?



- On an average if there is one stops the ticket price will drop by 24% of the cost of ticket without any stop.

# Model Building

We have used several linear regression models to evaluate and finalize the best models, The major models we have used as follows.

## 1.Linear Regression

```
: #Linear model
ln=LinearRegression()
ln.fit(x_train,y_train)
predln=ln.predict(x_test)
print('r2 score is :',round((r2_score(y_test,predln)),3))
print('RMSE:',round(np.sqrt(mean_squared_error(y_test,predln)),3))
print('mean absolute error:',round((mean_absolute_error(y_test,predln)),3))

r2 score is : 0.813
RMSE: 7111.926
mean absolute error: 5103.609
```

## 2.Lasso Regression model

```
: #Lasso model
ls=Lasso(alpha=9)
ls.fit(x_train,y_train)
predls=ls.predict(x_test)
print('r2 score is :',round((r2_score(y_test,predls)),3))
print('RMSE:',round(np.sqrt(mean_squared_error(y_test,predls)),3))
print('mean absolute error:',round((mean_absolute_error(y_test,predls)),3))

r2 score is : 0.813
RMSE: 7104.568
mean absolute error: 5078.614
```

## 3.Ridge Regression

```
: #Ridge model
rd=Ridge(alpha=16)
rd.fit(x_train,y_train)
predrd=rd.predict(x_test)
print('r2 score is :',round((r2_score(y_test,predrd)),3))
print('RMSE:',round(np.sqrt(mean_squared_error(y_test,predrd)),3))
print('mean absolute error:',round((mean_absolute_error(y_test,predrd)),3))

r2 score is : 0.814
RMSE: 7099.817
mean absolute error: 5045.134
```

## 4.Elasticnet Regression

```
#ElasticNet model
enr=ElasticNet(alpha=0.001)
enr.fit(x_train,y_train)
predenr=enr.predict(x_test)
print('r2 score is :',round((r2_score(y_test,predenr)),3))
print('RMSE:',round(np.sqrt(mean_squared_error(y_test,predenr)),3))
print('mean absolute error:',round((mean_absolute_error(y_test,predenr)),3))

r2 score is : 0.813
RMSE: 7110.851
mean absolute error: 5100.703
```

## 5.Ransac Regressor

```
ran = RANSACRegressor(base_estimator=LinearRegression(), max_trials=100)
ran.fit(x_train, y_train)
predran=ran.predict(x_test)
print('r2 score is :',round((r2_score(y_test,predran)),3))
print('RMSE:',round(np.sqrt(mean_squared_error(y_test,predran)),3))
print('mean absolute error:',round((mean_absolute_error(y_test,predran)),3))
```

```
r2 score is : 0.762
RMSE: 8019.396
mean absolute error: 4890.685
```

---

## 6.Support Vector Regressor

```
: svr=SVR()
svr.fit(x_train, y_train)
predpoly=svr.predict(x_test)
print('r2 score is :',round((r2_score(y_test,predpoly)),3))
print('RMSE:',round(np.sqrt(mean_squared_error(y_test,predpoly)),3))
print('mean absolute error:',round((mean_absolute_error(y_test,predpoly)),3))
```

```
r2 score is : -0.169
RMSE: 17787.638
mean absolute error: 13066.372
```

## 7.Random Forest Regressor

```
rf = RandomForestRegressor(n_estimators=100)
rf.fit(x_train, y_train)
predrf=rf.predict(x_test)
print('r2 score is :',round((r2_score(y_test,predrf)),3))
print('RMSE:',round(np.sqrt(mean_squared_error(y_test,predrf)),3))
print('mean absolute error:',round((mean_absolute_error(y_test,predrf)),3))
```

```
r2 score is : 0.976
RMSE: 2537.227
mean absolute error: 1190.071
```

## 8.AdaBoost Regressor

```
: ada = AdaBoostRegressor()
ada.fit(x_train, y_train)
predada=ada.predict(x_test)
print('r2 score is :',round((r2_score(y_test,predada)),3))
print('RMSE:',round(np.sqrt(mean_squared_error(y_test,predada)),3))
print('mean absolute error:',round((mean_absolute_error(y_test,predada)),3))
```

```
r2 score is : 0.847
RMSE: 6424.581
mean absolute error: 5423.27
```

# Cross Validation Score

```
models=[ln,ls,rd,enr,nan,svr,rf,ada]
for m in models:

    score=cross_val_score(m,x,y,cv=5)
    print(m,'score is:')
    print(round((score.mean()),3))
    print('\n')
```

LinearRegression() score is:  
0.672

Lasso(alpha=9) score is:  
0.676

Ridge(alpha=16) score is:  
0.688

ElasticNet(alpha=0.001) score is:  
0.673

RANSACRegressor(base\_estimator=LinearRegression()) score is:  
0.688

SVR() score is:  
-0.241

RandomForestRegressor() score is:  
0.772

AdaBoostRegressor() score is:  
0.687

---

**The difference between accuracy and cross validation is less for random forest regressor, so it is the best model**



# Hyper Parameter Tuning

```
rf=RandomForestRegressor()
grid_param={
    'criterion':['mse','mae'],

    'max_depth':[10,20,30,40,50],
    'max_features':['auto', 'sqrt', 'log2'],
    'min_samples_split':[2,5,10,15,20],
    'bootstrap':[True,False]
}

gd_sr=GridSearchCV(estimator=rf,
                    param_grid=grid_param,
                    scoring='r2',
                    cv=5)

gd_sr.fit(x,y)

best_parameters=gd_sr.best_params_
print(best_parameters)
best_result=gd_sr.best_score_
print(best_result)
```

```
{'bootstrap': True, 'criterion': 'mse', 'max_depth': 40, 'max_features': 'auto', 'min_samples_split': 20}
0.7991795802183211
```

---

## Final Accuracy

```
rf = RandomForestRegressor(n_estimators=100)
rf.fit(x_train, y_train)
predrf=rf.predict(x_test)
print('r2 score is :',round((r2_score(y_test,predrf)),3))
print('RMSE:',round(np.sqrt(mean_squared_error(y_test,predrf)),3))
print('mean absolute error:',round((mean_absolute_error(y_test,predrf)),3))
```

```
r2 score is : 0.976
RMSE: 2537.227
mean absolute error: 1190.071
```

# **Conclusions**

- **Key findings and the conclusions of the study**

Based on the above key findings we would like to give some prescription for our client if you want to travel in Business class choose vistara it has least average price, and for Economy class choose IndiGo, and for Premium Economy go for Air Asia.

- If you want to travel from Delhi to Hyderabad in Economy class then Air India will offer least price followed by GoFirst and Vistara and Air Asia is the costly one.
- If you want to travel from Delhi to Kolkata in Economy class then Air India will offer least price followed by IndiGo and Vistara and SpaceJet is the costly one.
- If you want to travel from Delhi to Mumbai in Economy class then SpaceJet will offer least price followed by Air India and IndiGo and Air Asia is the costly one.
- If you want to travel from Mumbai to Bengaluru in Economy class then IndiGo will offer least price followed by Vistara and Air India and SpaceJet is the costly one.
- If you want to travel from Delhi to Bengaluru in PremiumEconomy class then Vistara is the best choice
- If you want to travel from Delhi to Chennai in Premium Economy class then Vistara is the best choice
- If you want to travel from Delhi to Hyderabad in Premium Economy class then Vistara
- If you want to travel from Delhi to Kolkata in Premium Economy class then Air Asia is the cheapest one followed by GoFirst and SpiceJet
- If you want to travel from Delhi to Mumbai and Mumbai to Bengaluru in Premium Economy class then Vistara is the best choice.
- If you want to go from Delhi to Bengaluru Chennai or Hyderabad or Mumbai or Mumbai to Bengaluru then vistara is the only available as per this data at this moment, It may vary on coming days.
- If you want to go from Delhi to Kolkata in Business class then Air Asia is the cheapest one followed by GoFirst and SpaceJet and Vistara is the costliest one according to this data.
- The above graph depicts that if you want to book for Economy class then if you choose 4am-5am in morning you have maximum chance that your booking price will be lowest in that whole day.
- The highest price is generally observed for Economy class is around 3-5pm.
- For premium economy class the least price observed around 11am to 12pm and costly price is observed between 3pm -4 pm
- For Business class the least price observed around 11am to 12pm and costly price is observed between 3pm -4 pm
- From the above graph we can clearly see that as duration increases the price is also going to be increased.

- **Learning outcomes of the study in respect to data science**

Learnt many things from this project, following are the significant one

1. How to deal with xpaths and getting the urls.
2. Before we step into the project we should have a foresight of its complete picture and completed version.
3. Learnt that whenever we are checking whether a variable is of object or integer type we shouldn't include int in double inverted comma because it's a keyword.
4. Overall experience got enriched from this project.

**!! Thank You !!**