



FLIGHT PRICE PREDICTION

Submitted by:
Ganesh kumbhar

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)

- **Conceptual Background of the Domain Problem**

Airline companies use complex algorithms to calculate flight prices given various conditions present at that particular time. These methods take financial, marketing, and various social factors into account to predict flight prices.

Nowadays, the number of people using flights has increased significantly. It is difficult for airlines to maintain prices since prices change dynamically due to different conditions. That's why we will try to use machine learning to solve this problem. This can help airlines by predicting what prices they can maintain. It can also help customers to predict future flight prices and plan their journey accordingly

- **Review of Literature**

As per the requirement of client, I have scrapped the data from online sites and based on that data I have did analysis like for based on which feature of my data prices are changing. and checked the relationship of flight price with all the feature like what flight he should choose.

- **Motivation for the Problem Undertaken**

I have worked on this on the bases of client requirements and followed all the steps till model deployment.

Analytical Problem Framing

After scrapping my data using selenium I have loaded my data into python with the help of pandas.

Home Page - Select or create a ... Cars Price Prediction - Jupyter N ... Car Price Web Scrapping - Jupy ... Untitled4 - Jupyter Notebook ... Real-Time-Flight-Price-Predictio ...

localhost:8888/notebooks/Untitled4.ipynb

Jupyter Untitled4 Last Checkpoint: Yesterday at 2:02 PM (autosaved)

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3

```
from sklearn.linear_model import LinearRegression
from sklearn.tree import DecisionTreeRegressor
from sklearn.neighbors import KNeighborsRegressor
from sklearn.svm import SVR
from sklearn.model_selection import train_test_split
from sklearn.metrics import r2_score, mean_absolute_error, mean_squared_error
from sklearn.model_selection import GridSearchCV
from sklearn.preprocessing import OrdinalEncoder, power_transform, StandardScaler, MinMaxScaler
```

In [2]: df=pd.read_csv("Web_Scraped_flight.csv",index_col=0)
df

Out[2]:

	Airline	Source	Destination	Dep_Time	Arrival_Time	Duration	Total_Stops	price	
0	IndiGo	IndiGo-137/2046	Jammu	Mumbai	11:55	19:25	7h 30m	2 Stop(s)	7,365
1	IndiGo	IndiGo-137/5041	Jammu	Mumbai	11:55	20:50	8h 55m	2 Stop(s)	7,365
2	IndiGo	IndiGo-609/6722	Jammu	Mumbai	14:00	23:30	9h 30m	2 Stop(s)	7,365
3	IndiGo	IndiGo-137/6722	Jammu	Mumbai	11:55	23:30	11h 35m	2 Stop(s)	7,365
4	SpiceJet	SpiceJet-SG-160/945	Jammu	Mumbai	11:15	16:15	5h 00m	1 Stop	7,418
...
1095	Vistara	Vistara'InUK-998/747	Pune	Kolkata	16:55	08:35n+ 1 day	15h 40m	1 Stop	10,209
1096	Vistara	Vistara'InUK-998/705	Pune	Kolkata	16:55	09:35n+ 1 day	16h 40m	1 Stop	10,209
1097	Vistara	Vistara'InUK-998/737	Pune	Kolkata	16:55	18:05n+ 1 day	25h 10m	1 Stop	10,209
1098	Vistara	Vistara'InUK-998/707	Pune	Kolkata	16:55	19:35n+ 1 day	26h 40m	1 Stop	10,388
1099	Air India	Air India'InAI-482/762	Pune	Kolkata	12:45	00:05n+ 1 day	11h 20m	2 Stop(s)	11,838

6730 rows x 8 columns

In [3]: df.isnull().sum()

Out[3]: Airline 0
Source 0

Windows Taskbar: Type here to search, 30°C Smoke, 15:50, 29-01-2022

The size of the data is 6730*8

Data Pre-processing

Checking null values

```
data.isna().sum()
```

```
Unnamed: 0      0
Unnamed: 0.1    0
Airline         0
Source          0
Destination     0
Dep_Time        0
Arrival_Time    0
Duration        0
Total_Stops     0
Price           0
dtype: int64
```

There are no missing values

as there is no null values so I can move forward

- Data Sources and their formats

I have collected data from web scrapping and I have converted it into csv format

- Data Preprocessing Done

Doing pre-processing where I am dropping some columns and filling missing values in total stops

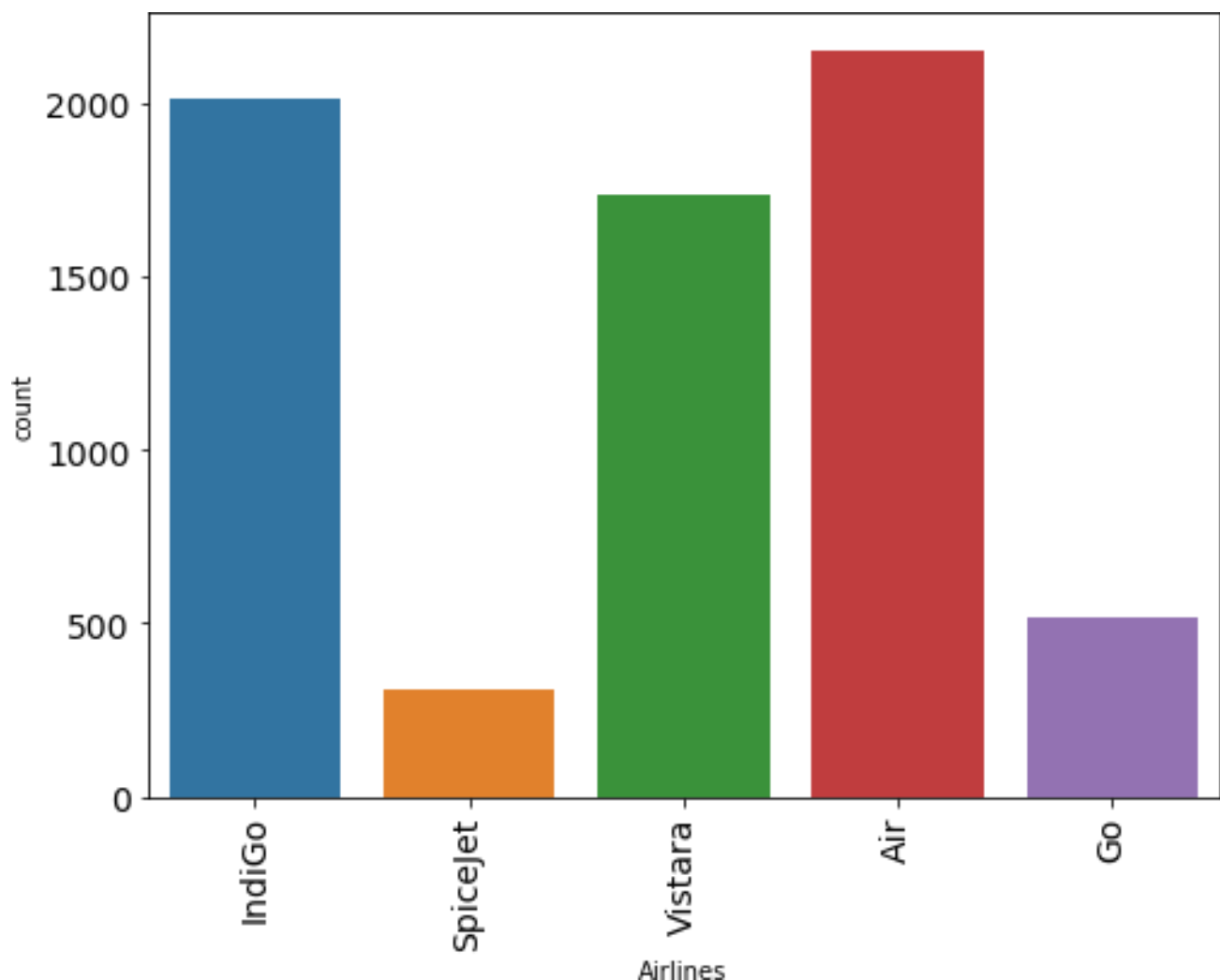
```
def preprocess1(df):
    df['Total_Stops']=df['Total_Stops'].fillna(df['Total_Stops'].mode()[0])
    df=df.drop(['Duration'],axis=1)
    return df
```

```
def preprocess2(df):
    df['Dep_hour'] = pd.to_datetime(df['Dep_Time']).dt.hour
    df['Dep_minute'] = pd.to_datetime(df['Dep_Time']).dt.minute
    df = df.drop(['Dep_Time'], axis=1)
    df['arrival_hour'] = pd.to_datetime(df['Arrival_Time']).dt.hour
    df['arrival_minute'] = pd.to_datetime(df['Arrival_Time']).dt.minute
    df = df.drop(['Arrival_Time'], axis=1)
    return df
```

Here I am converting time into hour and minute and also dropping some columns that are not useful for my model.

- Data Inputs- Logic- Output Relationships

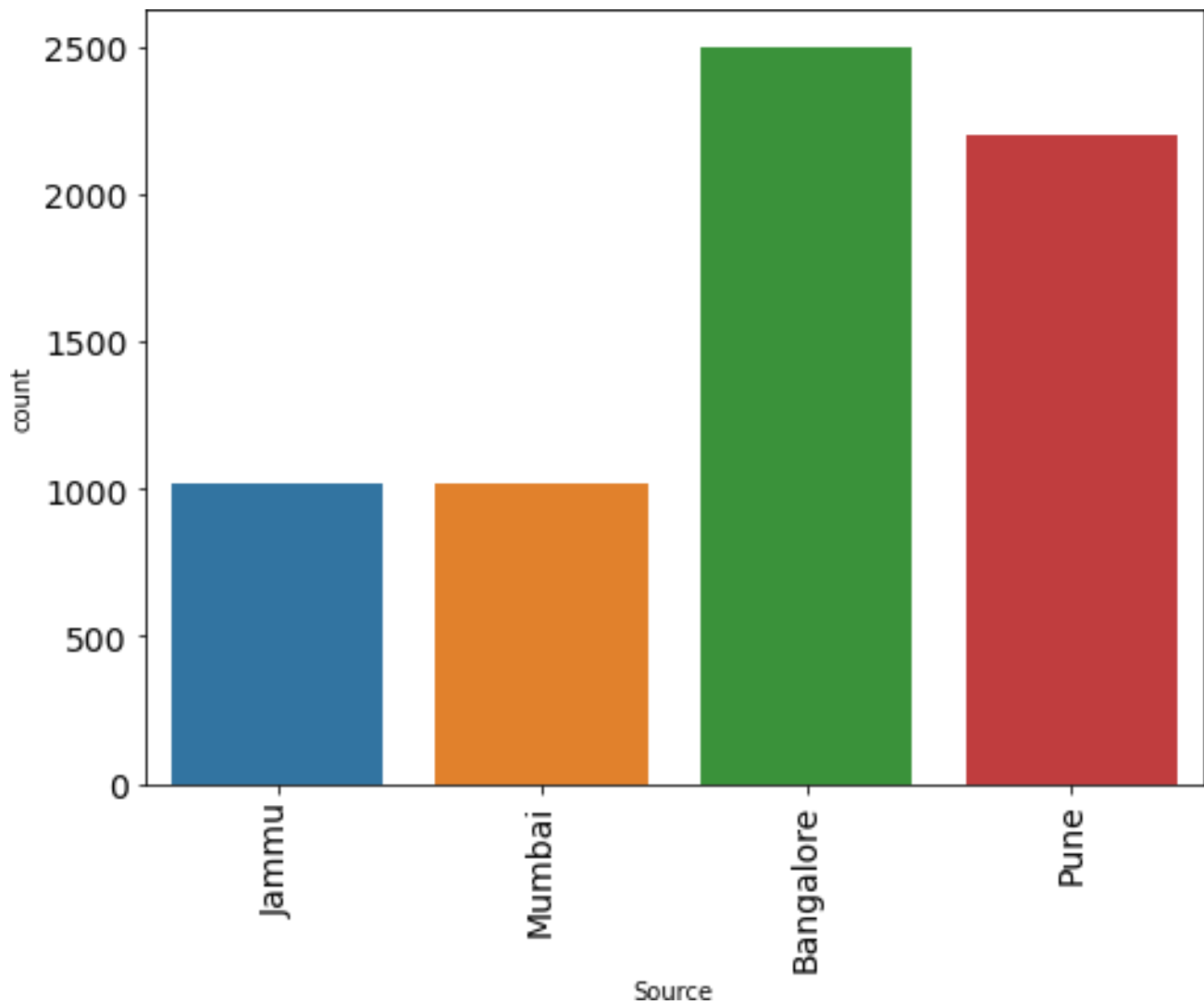
I have did EDA to understand the feature relationship.



Obseravtion

- 1-Mostly people use to travel with Air India
- 2-After air india people use to travel with IndiGo
- 3-ANd spicejet has the least count

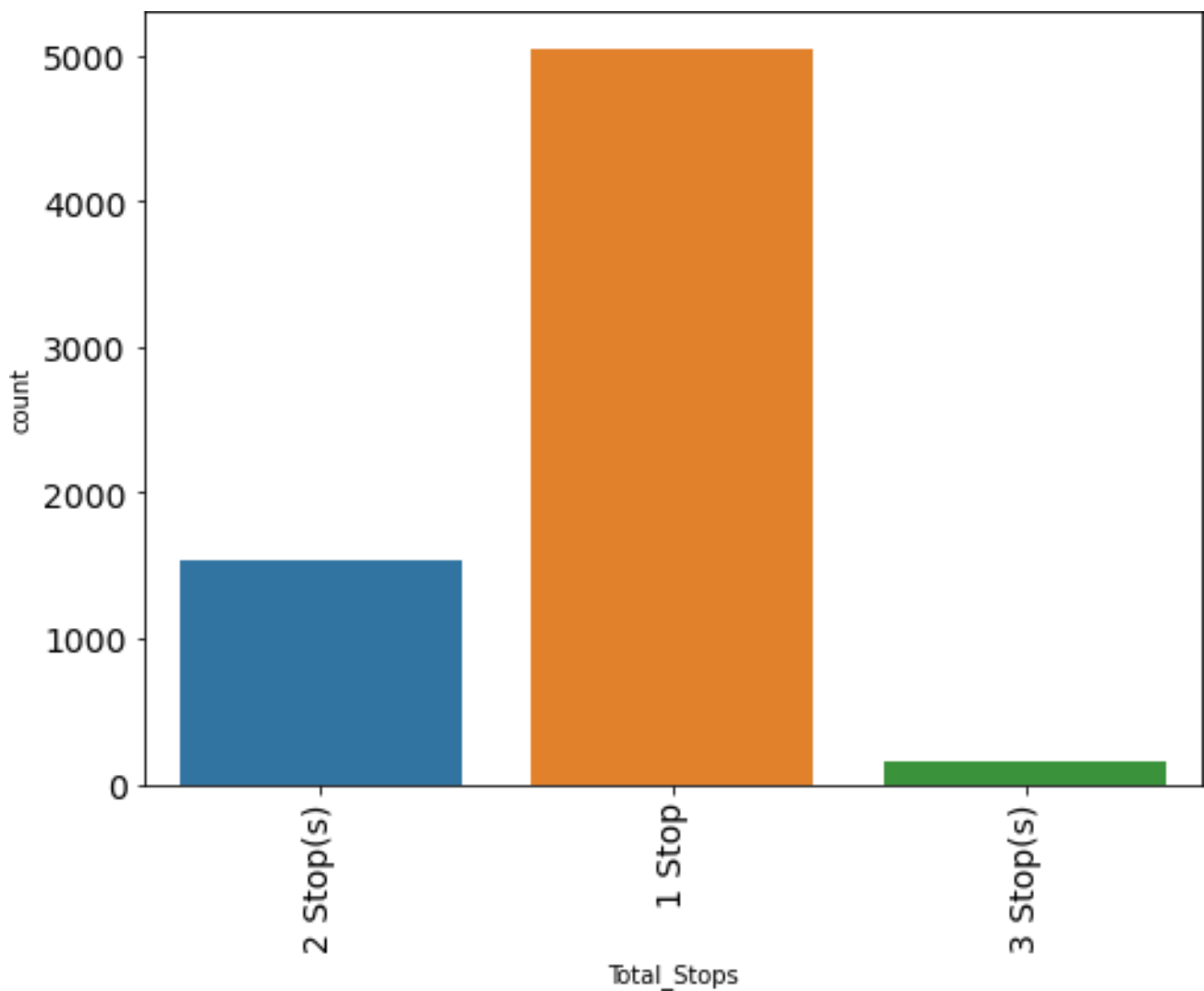
Countplot of Source



Obseravtion

- 1-Mostly Source has bangalore as high count
- 2-after bangalore pune has 2nd high count
- 3-and at least jammu and mumbai equal count

Countplot of Total Stops



Observation

1-Mostly people use to take a flight who has only one stop

3-AND only approx 1500 people use to take 2 stop flights

4-AND there are very less people who used to take 3 stops or 4 stops flight

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

Here, you can describe any presumptions taken by you.

- Hardware and Software Requirements and Tools Used

Model/s Development and Evaluation

- Identification of possible problem-solving approaches (methods)

I have did Analysis on this data to understand the value of each feature and the contribution of each feature for model creating and effect of all the feature on the prices. And considering all the point I have built a model that can predict the prices.

Testing of Identified Approaches (Algorithms)

I have trained many model and even evaluate them using all the performance metrics of regression. Here is the screenshot I have make a dataframe of all the model and metrics so here we can see the performance of every model. I have selected LIGHT

```

"Cross_Val_Score":CVS,
"R2_score":R2,
"Mean_squared_error":MSE,
"Mean_Absolute_Error":MAE,
"RMSE":RMSE
))

In [30]: models_result
Out[30]:
```

	NAME	Cross_Val_Score	R2_score	Mean_squared_error	Mean_Absolute_Error	RMSE
0	XGB Regressor	0.999999	0.999999	1.069399e-07	2.079237e-04	3.270167e-04
1	ExtraTrees Regressor	1.000000	1.000000	2.843791e-28	1.501941e-14	1.686354e-14
2	RandomForest Regressor	1.000000	1.000000	2.285704e-28	1.381318e-14	1.511854e-14
3	Linear Regression	0.311658	0.320709	5.683874e-02	1.819486e-01	2.384088e-01
4	DecisionTree Regressor	1.000000	1.000000	2.263136e-28	1.125642e-14	1.504372e-14
5	Lasso	-0.000053	-0.000011	8.367453e-02	2.198570e-01	2.892655e-01
6	LIGHT GBM	1.000000	0.999999	1.007051e-07	1.237600e-04	3.173406e-04

```

In [31]: ### I will do hyperparameter Tuning of XGB because it is giving me good CVS and 0.99% R2 Score
In [32]: XGB=XGBRegressor()
In [33]: XGB.fit(x_train,y_train)
Out[33]: XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
colsample_bynode=1, colsample_bytree=1, enable_categorical=False,
gamma=0, gpu_id=-1, importance_type=None,
interaction_constraints='', learning_rate=0.300000012,
max_delta_step=0, max_depth=6, min_child_weight=1, missing=nan,
monotone_constraints=()), n_estimators=100, n_jobs=8,
num_parallel_tree=1, predictor='auto', random_state=0, reg_alpha=0,
reg_lambda=1, scale_pos_weight=1, subsample=1, tree_method='exact',

```

GBM as a final model because of its accuracy and performance metrics.

- Key Metrics for success in solving problem under consideration

I have make a dataframe of all the model and metrics so here we can see the performance of every model. I have selected LIGHT GBM as a final model because of its accuracy and performance metrics.

- Interpretation of the Results

From the above eda we can easily understand the relationship between features and and we can even see which things are effecting the price of flights.

CONCLUSION

- **Key Findings and Conclusions of the Study**

Describe the key findings, inferences, observations from the whole problem.

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

The above research will help our client to study the latest flight price market and with the help of the model built he can easily predict the price ranges of the flight, and also will helps him to uLimitatinderstand Based on what factors the fight price is decided.

ons of this work and Scope for Future Work

The limitation of the study is that in the volatile changing market we have taken the data, to be more precise we have taken the data at the time of pandemic and recent data, so when the pandemic ends the market correction might happen slowly. So based on that again the deciding factors of the might change and we have shortlisted and taken these data from the important cities across india, if the customer is from the different city our model might fail to predict the accuracy prize of that flight.

THANK YOU !!!