

FLIGHT PRICE PREDICTION PROJECT

SUBMITTED BY : **SINDHU SHREE N**

INTERNSHIP BATCH : **19**

ACKNOWLEDGMENT

With great pleasure, I sincerely express my deep sense of gratitude and heartfelt thanks to several individuals from whom I received impetus motivation during the internship project work. I am very grateful to **SHUBHAM YADAV** sir internship 19th batch guide, for the help and I do express my deep gratitude to sir for his guidance, keen interest and constant encouragement throughout the internship project work. I am thanking him for personal concern, affinity and great spirit with which he has guided the work.

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)

Machine Learning is a field of technology developing with immense abilities and applications in automating tasks, where neither human intervention is needed nor explicit programming.

The power of ML is such great that we can see its applications trending almost everywhere in our day-to-day lives. ML has solved many problems that existed earlier and have made businesses in the world progress to a great extent.

Conceptual Background of the Domain Problem

The project contains flight price dataset. And we are supposed to predict the price of the flights based on multiple features. Here we have used regression technique for price prediction. We are about to deploy an ML model for flight price prediction and analysis. This kind of system becomes handy for many people.

Review of Literature

All possible information from all the available data tables more the information, more than for EDA and feature Engineering. It's more to aggregate the data and perform all these steps in order to obtain a good result.

Motivation for the Problem Undertaken

Any kind of modifications can also be later inbuilt in this application. It is only possible to later make a facility to find out buyers. This a good idea for a great project that we can try out. The applications of Machine Learning don't end here. Similarly, there are infinite possibilities that we can explore. But for the time being, let me build the model for flight Price Prediction and its deployment process.

Analytical Problem Framing

Mathematical/ Analytical Modelling of the Problem

1. Included in 'Data-cleaning ipynb'.
2. Selecting relevant features.
3. Exploratory the data analysis and data cleaning.
4. Null value imputation
5. Handling outliers
6. Training a machine learning model
7. Hyperparameter tuning
8. Evaluate the model
9. Predictions of the model

Data Sources and their formats

We have scrapped the data using webscraping in the websites like yatra.com, Airlines.com and stored them in a excel sheet.

Creating index on all the table to be joined this will speed up and processing here we are using EDA feature Engineering, Data visualization and statics approach perform data cleaning, outlier handling, missing values build model etc.

Data Pre-processing Done

Creating index on all the table to be joined this will speed up and processing here we are using EDA feature Engineering, Data visualization and statics approach perform data cleaning, outlier handling, missing values build model etc.

Partitioning and splitting that dataset account as that dataset is skewed, stratification is used allocate the samples evenly based on sample classes so that training set and test set have similar ratio of classes.

Data Inputs- Logic- Output Relationships

Since it is a regression model where the target feature is continuous data, there is no imbalance of the data so we used correlation matrix to find the relation between the input and output features.

Set of assumptions (if any) related to the problem under consideration

I am not taken any presumptions of this problem

Hardware and Software Requirements and Tools Used

Hardware requirements:

PROCESSOR: Intel(R) Core(TM) i3 CPU

MONITOR : Any display unit

HARD DISK : 240GB SSD

RAM : 8.00GB

Software requirements:

OPERATING SYSTEM: Windows 10 Pro

FRONT END : Jupyter Notebook (Anaconda3)

BACK END : Excel 2013

Tools Used:

- 1) Pandas Library
- 2) Numpy Library
- 3) Seaborn Library
- 4) Matplotlib
- 5) Scikit_learn

Model/s Development and Evaluation

Identification of possible problem-solving approaches (methods)

- Examining the data
- Checking the basic details (null value, D type, Shape etc)
- Identifying the target and independent features and perform EDA using Data visualization and Statistical approach accordingly
- Performing data cleaning, outliers handling, missing value imputation
- Feature engineering
- Performing hyperparameter tuning
- Evaluating the model again
- Making predictions

Testing of Identified Approaches (Algorithms)

Listing down all the algorithms used for the training and testing.

- 1) XGBoost Regression
- 2) Decision Tree Regression
- 3) Gradient Boosting Regression
- 4) Random Forest Regression

Run and Evaluate selected models:

1. **XGBoosting Regression:** Extreme Gradient Boosting regressor provides an efficient and effective implementation of gradient boosting algorithm which is used for regression predictive modelling.

```
#r2_score
from sklearn.metrics import r2_score
print(r2_score(y_test, pred_xg)*100)
```

82.6211364134987

2. **Decision Tree Regression:** This is a type of regression where it splits the data instance into numerous trees and predicts the output based on it.

```
from sklearn.metrics import r2_score
print(r2_score(y_test, pred_dt)*100)
```

71.62556024024475

3. **Gradient Boosting Regression:** Gradient boosting is a machine learning technique for regression, which produces a prediction model in the form of an ensemble of weak prediction models, typically decision trees.

```
: from sklearn.metrics import r2_score
print(r2_score(y_test, pred_gbr)*100)
```

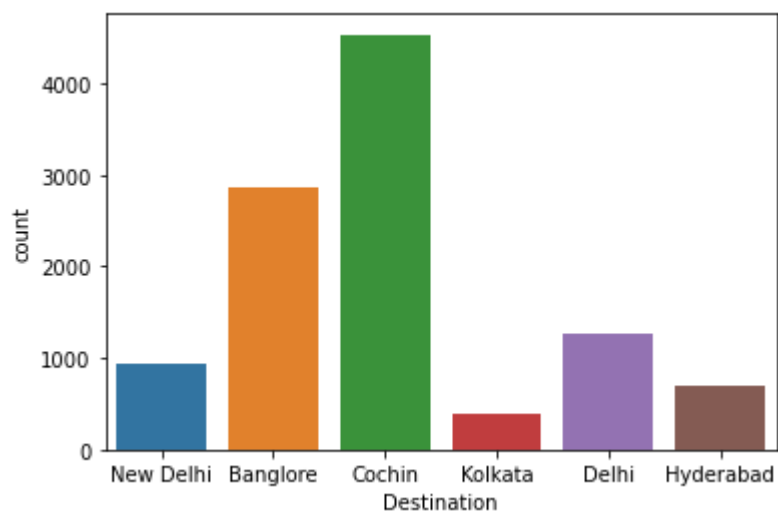
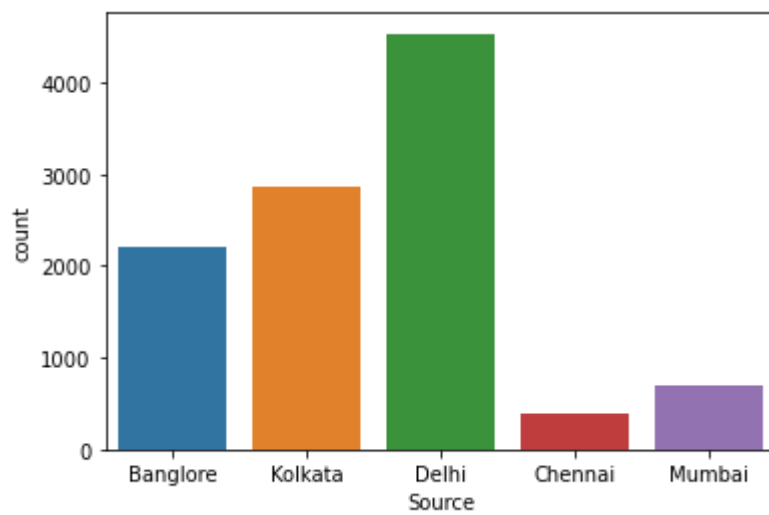
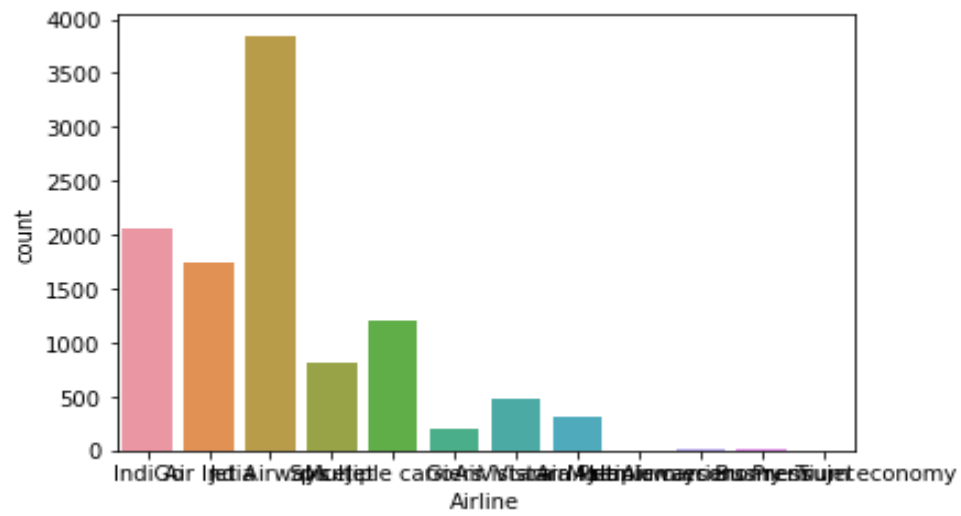
79.43001976315394

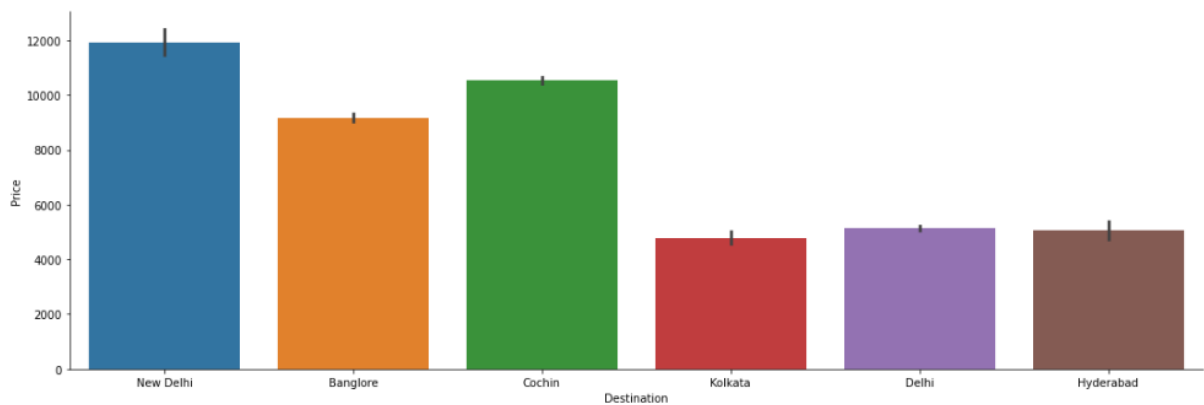
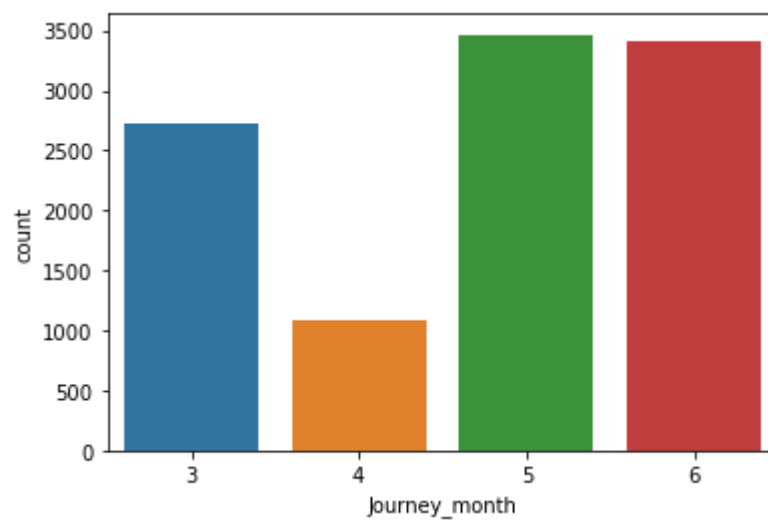
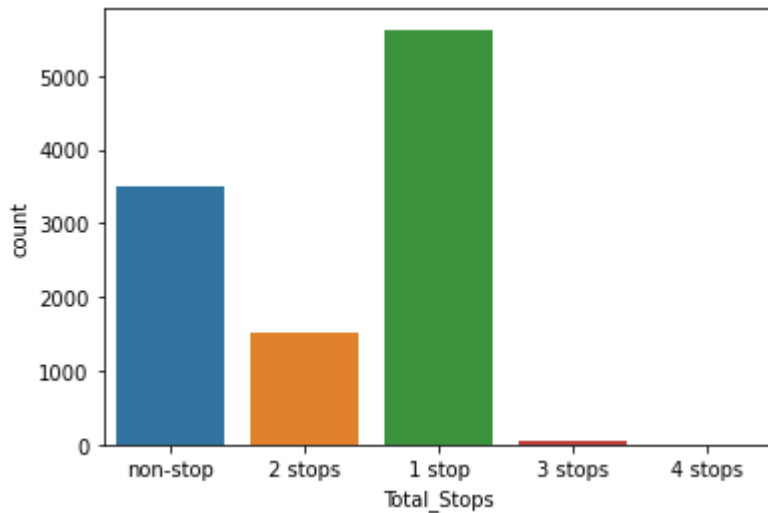
4. **Random Forest Regression:** It is a type of ensemble learning method that uses numerous decision trees to achieve higher prediction accuracy and model stability. Every tree classifies a data instance based on attributes, and the forest chooses the classification that received most instances.

```
from sklearn.metrics import r2_score
r2_score(y_test, pred_rf)*100
```

78.82659076698637

Visualizations:





Observations:

1. The count of AirIndia is more compare to other airlines.
2. Most people starts their Journey from Delhi.
3. Maximum travellers prefer 1 stop travelling followed by non-stop.
4. Most travellers are found in the months may and june.

5. For destination “NewDelhi” maximum price is fixed and followed by ‘Cochin’ and ‘Bangalore’.

Interpretation of the Results

On comparing all the results of these four models mentioned above we have observed that the xgboost, random forest and gradient boost regressors are performing well but the training set and testing set score is high almost near in xgboost regression and also has a comparative cross validation score. Thus we have choosed this model for hyperparameter tuning and there also we got a improved result. Hence **XGBoosting Regressor** is the best fit for the present dataset.

CONCLUSION

Key Findings and Conclusions of the Study:

Most Regression problems in the real world are imbalanced. Also, almost always data sets have missing values. But in this case since the target feature price is a continuous data we don't have any data imbalance and here we covered strategies to deal with missing values in the datasets. We also explored different ways of building ensembles in sklearn.

Learning Outcomes of the Study in respect of Data Science:

There is no definitive guide of which algorithms to use given any situation. What may work on some data sets may not necessarily work on others. Therefore, always evaluate methods using cross validation to get a reliable estimate.

Sometimes we may be willing to give up some improvement to the model if that would increase the complexity much more than the percentage change in the improvement to the evaluation metrics.

In some Regression problems, false negative is a lot more expensive than false positives. Therefore, we can reduce cut-off points to reduce the false negative.

Missing values sometimes add more information to the model than we might expect. One way of capturing it is to add binary features for each feature that has missing values.

Limitations of this work and Scope for Future Work

- Limited Access to information
- Time Limits
- Conflicts on biased views and personal issues
- How to structure my project research limitation correctly
- How to set my project research limitation.
- Formulation of my objectives and aims
- Implementation of my data collection methods
- Scope of discussions
- Finding my error on the codes

Concluding Thoughts

This study used different models in order to predict flight prices. However, there was a relatively small dataset for making a strong inference because of limited number of rows. Gathering more data can yield more robust predictions. Secondly, there could be more features that can be good predictors.

