

# Capstone Project : Marketing-Airplane Passenger Satisfaction Prediction Using Machine Learning Techniques

Bekee Sorbarisere Yirakpoa<sup>1</sup> and Mercy Nwanyanwu<sup>2</sup>

<sup>1,2</sup>Department of Computer Science  
Captain Elechi Amadi Polytechnic, Port Harcourt

## Abstract.

Customer satisfaction questionnaires are a rich and strong source of information for companies to seek loyalty, customer and client retention, optimize resources, and repurchase products. Several advanced machine learning and statistical models have been employed to estimate the customer satisfaction score; however, there is not a single model that can yield the best result in all situations. Ensembles of regression techniques have demonstrated their effectiveness for various applications, where the success of these models lies in the construction of a set of single models. I performed an experimental study using a real dataset of 90917 samples from US airline carrier 'Falcon airlines', in order to verify the benefits of ensemble models for predicting customer satisfaction. Accordingly, in this project I evaluated the following models; Logistic Regression, Decision Tree, Bagging classifier and Random forest. The obtained results indicate that the Random forest performs better in terms of Recall and Precision.

**Keywords:** Logistic regression, customer satisfaction, Bagging classifier and Random forest.

---

## Introduction

### Problem Description

Airline businesses globally are faced with the challenge of grounding especially due to the pandemic. To stay in business, Airline operators need to determine relative parameters that can contribute to the satisfaction of passengers. This is the dilemma of a reputed US airline carrier 'Falcon airlines'. They aim to determine the relative importance of each parameter with regards to their contribution to passenger satisfaction. To achieve this aim a random sample of 90917 individuals who travelled using their flights is provided. The on-time performance of the flights along with the passengers' information is published in the Marketing Project-Flight data.csv file named 'Flight data'. These passengers were asked to provide their feedback at the end of their flights on various parameters along with their overall

experience. These collected details are available in the Marketing Project-Survey data.csv 'Survey data'.

In the survey, the respondents were requested to express their satisfaction or not with their overall flight experience and that is captured in the data of survey report under the variable labelled 'Satisfaction'. The need for this study by 'Falcon airlines' is to give themselves a competitive edge over other airline operators by identifying critical factors that lead to customer satisfaction.

### Motivation For Project

It is obvious that at the end of the pandemic, there will be an increase in demand for air travel as most persons may wish to be on vacations. Hence the need for 'Falcon airlines' to carry out this study to ascertain Passengers level of

satisfaction with their services to give themselves a competitive edge over other airline

## **Aim and Objectives**

This project was a few weeks' efforts to develop a predictive model for Airline passenger satisfaction using data from US airline carrier 'Falcon airlines'. I hope the outcome of the project will help streamline the analysis and prediction of passenger's satisfaction for US airline carrier 'Falcon airlines' and other airlines.

The objective of this project are-

1. To understand which parameters play an important role in swaying a passenger feedback towards 'satisfied'.
2. To predict whether a passenger will be satisfied or not given the rest of the details are provided, using supervised machine learning techniques.
3. To compare different classification techniques to understand which is best suitable for this application.

No doubt the development of a framework and codes that incorporate analytics and machine learning concepts studied in the program is the goal. The success of the project is predicated on the accuracy of the classification results and the extent of analysis conducted. It is my hope that

.Working steps of Machine learning technique

## **Limitations of the Study**

In this project, I evaluated the effectiveness of using specific supervised machine learning techniques to address the problem of predicting passengers satisfaction on Airline flight and survey data provided. The limitations of the methods applied in the course of this project study are as follows:

I used a pre-labelled dataset to train the algorithms. However, usually, it is difficult to find labelled data and thus applying supervised machine learning techniques may not be feasible. In such cases, the option should be to evaluate unsupervised techniques which were beyond the scope of this project.

## **Lterature Review**

This part of the project seeks to review considerable literature on the subject matter. A few researchers have also conducted literature reviews of articles published on airline passenger satisfaction and the techniques used.

the final report will serve as a benchmark for further development on this topic.

## **Research Methodology**

The typical machine learning approach was followed in this project. The identified dataset has labelled class variable, which was used as the prediction variable in machine learning models.

Through exploratory analysis, we analysed the data set in detail and identified possible predictors. Through various visualization techniques, we observed the separation between Satisfied and neutral or non-satisfied assengers. To solve the Airline passenger satisfaction prediction problem, we experimented with a few supervised machine learning techniques – Logistic Regression, Decision tree, Bagging classifier and Random Forest,

Performance measures, like Confusion Matrix and Area Under Curve (AUC), was used to compare the performance of the models.

This analysis was conducted using Python through Jupyter notebook. In-built libraries and methods were used to run the machine learning models. When needed, functions were defined to simplify specific analyses or visualizations. The diagram below shows in detail the full process that was followed in the project

This project considers marketing survey and flight data that was mined by US airline carrier 'Falcon airlines' only. I evaluated a few machine learning algorithm – Logistic Regression, Decision tree, Bagging classifier and RandomForest. Although the result of the study using these algorithms is good, it is necessary to evaluate other techniques to determine which algorithm works best for this application. Due to the large size of data, I was limited by computation capacity to explore different other techniques

Air travel is one of the most convenient way for long distance travel at both national and international level [Park et al.,2009 ]. There are many airline service providers (ASPs) around the world. The competitive world motivates the

airlines company to attract the customers. However, a traveller considers quite a few factors before deciding on any airline.

These points can be airfare, tour time, quantity of stoppages, number of baggage allowed, and existing customer feedback etc. Therefore, all ASPs are working in all these client service

areas to enhance their facility and in-flight remedy in order to attract the customers.

It is very vital to recognize the desires and remedy level of customers i.e. customer satisfaction for the duration of the flight. Therefore, client remarks are very important or any airline industry. There could be quite a possible approach to gather the customer feedback. The most easiest and regular way is the customer feedback structure available during the journey. However, most of the passengers do no longer show any activity in filling feedback forms. Another shortcoming of this strategy is that it may additionally or can also not have appropriate questionnaire and might also be biased on positive parameters i.e. the feedback form may also only have sure unique questions. Other processes for purchaser feedback collection could be via online website or on-line mobile purposes of the airlines.

After the journey, an electronic mail with a link can be despatched to the passenger to request for a feedback. However, there is no guarantee of its success. Another strategy is to send a message

on passenger's cellular cellphone and ask them to rate your provider (1 for negative and 5

for excellent) on certain parameters. All these standard strategies opted by means of the industry are limited to certain parameters only. The greater handy way for a passenger is to express their feedback, as they want.

### **Summary of related work**

Pramod., et al 2019] in their work "a literature review: customer satisfaction on airline tweet using machine learning"; presents the earlier work done by the various researchers in the field

Therefore, the most convenient way for the passengers to share their opinions is the social media instead of feedback form. Social media presents a platform the place a consumer can freely categorical his feedbacks on any issues they observed all through flight. Twitter [<https://en.wikipedia.org/wiki/Twitter>] is one of the famous platforms worldwide. The information from Twitter can be utilized to strengthen a recommender system [Abel et al., 2019]. In addition, travellers are more comfortable in sharing their views about tour experiences on Twitter.

A variety of fundamental issues influences the emotions of a passenger in air travel. These issues can be cabin crew behavior, food quality, loss of baggage, seat comfort, flight delay, airfare etc. All these issues may give upward push to each superb and bad emotion. Also, if there is a continuous trend of terrible tweets for an airline, then it might also put a negative impact to the financial growth of the airline company. Therefore, it is vital to understand the issues that provide upward e to terrible tweets so that the respective airline company can take splendid action on time. There are large number airlines operating every day to join different geographical places [[www.quora.com](http://www.quora.com)]. Therefore, we might so anticipate a large number of people journeying every day in these flights. In addition, the number of tweets by passengers for airways would be very large. Therefore, it is a difficult assignment to extract the hidden emotion at the back of a tweet. Therefore, we required some tools and techniques that are in a position to take care of such a large quantity of tweet database and can provide insights to assist airline industry

of customer satisfaction of airline tweets using machine-learning techniques such as logistic regression, SVM, KNN, random forest, Naïve Bayes classifier, AdaBoost etc as shown below

Authors	Description	Publishing Year
T.HemakalaandS.Santhoshkumar	In this research, design a framework for sentiment analysis with opinion mining for the case of airlines service feedback. Most available datasets of hotel reviews are not labelled which presents many works for researchers as far as text data pre-processing task is concerned. Twitter is a SNS that has a huge data with user posting, with this significant amount of data, it has the potential of research related to text mining and could be subjected to sentiment analysis. The airline industry is a very competitive market, which has grown rapidly in the past 2 decades. Airline companies resort to traditional customer feedback forms which in turn are very tedious and time consuming. In this work, worked on a dataset comprising of tweets for 6 major Indian Airlines and performed a multi-class sentiment analysis. This approach starts with pre-processing techniques used to clean the tweets and then representing these tweets as vectors using a deep learning concept to do a phrase-level analysis. The analysis was carried out using 7 different classification strategies: Decision Tree, Random Forest, SVM, K-Nearest Neighbors, Logistic Regression, Gaussian Naïve Bayes and Ada Boost. The outcome of the test set is the tweet sentiment.	2018
Guoning Hu et al.	Analyze the opinion of 19M Twitter users towards 62 popular industries, encompassing 12,898 enterprise and consumer brands, as well as associated subject matter topics, via sentiment analysis of 330M tweets over a period spanning a month. We find that users tend to be most positive towards manufacturing and most negative towards service industries. In addition, they tend to be more positive or negative when interacting with brands than generally on Twitter. We also find that sentiment towards brands within an industry varies greatly and we demonstrate this using two industries as use cases. In addition, we discover that there is no strong correlation between topics and sentiments of different industries, demonstrating that topics and sentiments are highly dependent on the context of the industry that they are mentioned in. We demonstrate the value of such an analysis in order to assess the impact of brands on social media. We hope that this initial study will prove valuable for both researchers and companies in understanding users' perception of industries, brands and associated topics and encourage more research in this field.	2017

Yasmin Yashodha	This study examines the extensive strategic analysis of Air AsiaBerhad that has enabled it to sustain its competitive advantage as Asia's leading low cost carrier (LCC). The study demonstrates the diverse business-level, corporate level and competitive strategies of Air AsiaBerhad, played crucial roles in the LCC to successfully penetrate the under-served market segment of the airline industry within the ASEAN region. An in-depth analysis using a wide array of academic resources, relevant financial, legal and management resources and authorized websites, including face-to-face interviews were used to provide a more consequential comprehension on the varied business and international strategies that were implemented by Air AsiaBerhad. This research exhibits critical analysis pertaining to the current macro environment of the aviation industry which includes the PESTEL framework and Porter's Industry Analysis. The competitive environment analysis for Air AsiaBerhad is thoroughly scrutinised to examine the driving determinants that attributed to the organization's competitive advantage in the industry.	2012
M. Vadivukarasial.	In Twitter, the customer of airline services can tweet their opinions about their travelled experiences in flight. So Twitter contains massive amount of data and information regarding airline services. These tweets are collected and explored the sentiments about the airline services to track customer satisfaction reports and to discover location of the customer.	2018
Janet R. McColl-Kennedy et al.	Contextualized in post purchase consumption in business-to-business settings, the authors contribute to customer experience (CX) management theory and practice in three important ways. First, by offering a novel CX conceptual framework that integrates prior CX research to better understand, manage, and improve CXs—comprised of value creation elements (resources, activities, context, interactions, and customer role), cognitive responses, and discrete emotions at touchpoints across the customer journey. Second, by demonstrating the usefulness of a longitudinal CX analytic based on the conceptual framework that combines quantitative and qualitative measures. Third, by providing a step-by-step guide for implementing the text mining approach in practice, thereby showing that CX analytics that apply big data techniques to the CX can offer significant insights that matter. The author highlights six key insights practitioners need in order to manage their customers' journey, through (1) taking a customer perspective, (2) identifying root causes, (3) uncovering at-risk segments, (4) capturing customers' emotional and cognitive responses, (5) spotting and preventing decreasing sales, and (6) prioritizing actions to improve CX. The article concludes with directions for future research.	2018

Rida Khan and Siddhaling Urolagin	Social media today is an integral part of people's daily routines and the livelihood of some. As a result, it is abundant in user opinions. The analysis of brand specific opinions can inform companies on the level of satisfaction within consumers. This research focus is on analysis of tweets related to airlines based in four regions: Europe, India, Australia and America for consumer loyalty prediction. Sentiment Analysis is carried out using Text Blob analyzer.	2018
Ankita Rane and	The airline industry is a very competitive market which has grown rapidly in the past 2 decades. Airline companies resort to traditional customer feedback	2018

nandkumar	Forms which in turn are very tedious and time consuming. This is where Twitter data serves as a good source to gather customer feedback tweets and perform a sentiment analysis. In this paper, we worked on a dataset comprising of tweets for 6 major US Airlines and perform a multi-class sentiment analysis. This approach starts off with pre-processing techniques used to clean the tweets and then representing these tweets as vectors using a deep learning concept (Doc2vec) to do a phrase-level analysis. The analysis was carried out using 7 different classification strategies: Decision Tree, Random Forest, SVM, K-Nearest Neighbors, Logistic Regression, Gaussian Naïve Bayes and Ada Boost. The classifiers were trained using 80% of the data and tested using the remaining 20% data. The outcome of the test set is the tweet sentiment (positive/negative/neutral). Based on the results obtained, the accuracies were calculated to draw a comparison between each classification approach and the overall sentiment count was visualized combining all six airlines.	
Yun Wan and Qigan Gao	In airline service industry, it is difficult to collect data about customers' feedback by questionnaires, but Twitter provides a sound data source for them to do customer sentiment analysis. However, little research has been done in the domain of Twitter sentiment classification about airline services. In this paper, an ensemble sentiment classification strategy was applied based on Majority Vote principle of multiple classification methods, including Naive Bayes, SVM, Bayesian Network, C4.5 Decision Tree and Random Forest algorithms. In our experiments, six individual classification approaches, and the proposed ensemble approach were all trained and tested using the same dataset of 12864 tweets, in which 10 fold evaluation is used to validate the classifiers. The results show that the proposed ensemble approach outperforms these individual classifiers in this airline service Twitter dataset.	2015

MustafaAltinkök	This research was conducted for the purpose of analyzing the effect of the movement education program through a 12-week-coordination on the development of basic motor movements of pre-school children. A total of 78 students of pre-school period, 38 of whom were in the experimental group and 40 of whom were in the control group, were incorporated into the study in line with their own consent after their families had also been informed.	2016
M.Vadivukarassial.	Analyzed the twitter airline dataset for finding the best and the worst airlines and also to predict the most common issues occurred during the airline services. Then the word clouds of negative tweets are created and also the location of the negatively tweeted customer is predicted and visualized using geographical analysis. Finally training and testing was done on the dataset and also compared with seven different classifiers such as Logistic Regression classifier, KNeighbors classifier, SVC, Decision Tree classifier, Random Forest classifier, AdaBoost classifier and Gaussian NB. The results of four experiments demonstrate that the Random forest approach works best in real-world practice on sentiment classification of tweet data.	2018
Bee Yee Liao and Pei Pei Tan	The purpose of this paper is to study the consumer opinion towards the low-cost airlines or low-cost carriers (LCCs) (the set water ms are used interchangeably) industry in Malaysia to better understand consumers' needs and to provide better services. Sentiment analysis is undertaken in revealing current customers' satisfaction level towards low-cost airlines. About 10,895 tweets (data collected for two and a half months) are analyzed. Text mining techniques are used during data pre-processing and a mixture of statistical techniques are used to segment the customers' opinion. Results with two different sentiment algorithms show that there is more positive than negative polarity across the different algorithms. Clustering results show that both K-	2016
	Means and spherical K-Means algorithms delivered similar results and the four main topics that are discussed by the consumer on Twitter are customer service, LCC tickets promotions, flight cancellations and delays and post-booking management.	
Xiang Ji et al.	An important task of public health officials is to keep track of health issues, such as spreading epidemics. In this paper, we are addressing the issue of spreading public concern about epidemics. Public concern about a communicable disease can be seen as a problem of its own. Keeping track of trends in concern about public health and identifying peaks of public concern are therefore crucial tasks. However, monitoring public health concerns is not only expensive with traditional surveillance systems, but also suffers from limited coverage and significant delays.	2015

## Methodology

This methodology served as the deliverables of the project. It describes the results each phase that was tried out and do a comparison between them to identify which is the best technique to address

the airline passenger satisfaction prediction problem.

Each phase of the project has an output that describes the findings in that phase. These deliverables were used in this final project are explained below

## Project Deliverables

Methodology Phases	Project Deliverables
Understanding the dataset	<ul style="list-style-type: none"> <li>• Report on the summary of the dataset and each variable it contains along with necessary visualizations</li> </ul>
Exploratory Data Analysis	<ul style="list-style-type: none"> <li>• Report on analysis conducted and critical findings with a full description of data slices considered</li> <li>• Visualizations and charts that show the differences between satisfied and neutral or non-satisfied passengers</li> <li>• Python code of the analysis performed</li> </ul>
Modeling	<ul style="list-style-type: none"> <li>• Report on the results of the different techniques tried out, iterations that were experimented with, data transformations and the detailed modeling approach</li> <li>• Python code used to build machine learning models</li> </ul>
Final Project Report	<ul style="list-style-type: none"> <li>• Final report summarizing the work done over the course of the project, highlighting the key findings, comparing different models and identifying best model for predicting airline passenger satisfaction</li> </ul>

## Tools used

This project was entirely done using Python, and the analysis was documented in a Jupyter notebook. Standard python libraries were used to conduct different analyses. These libraries are described below–

- *sklearn* – used for machine learning tasks
- *seaborn* – used to generate charts and visualizations
- *pandas* – used for reading and transforming the data
- *Gridsearch* – used for model tuning

## Dataset

The problem consists of 2 separate datasets: Flight data & Survey data. The flight data has information related to passengers and the performance of flights in which they travelled. The survey data is the aggregated data of surveys

collected post service experience. You are expected to treat both the datasets as raw data and perform any necessary cleaning/validation steps as required

## Data Information

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 90917 entries, 0 to 90916
Data columns (total 24 columns):
```

```
# Column
```

```
-----
0 CustomerId
1 Satisfaction
```

```
Non-Null Count Dtype
```

```
90917 non-null int64
90917 non-null object
```

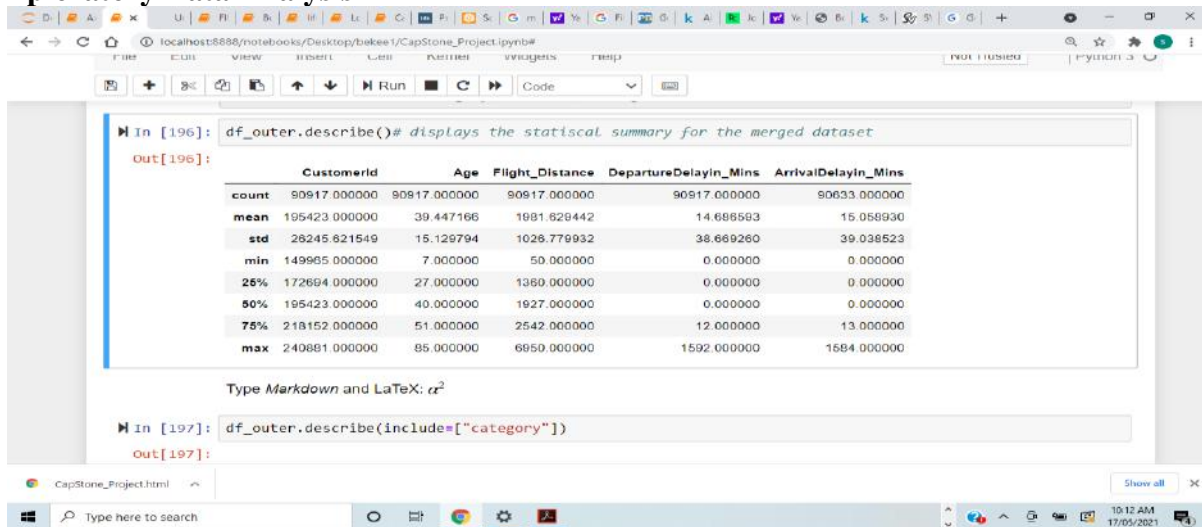


2	Seat_comfort	90917 non-null	object
3	Departure_Arrival_time_convenient	82673 non-null	object
4	Food_drink	82736 non-null	object
5	Gate_location	90917 non-null	object
6	Inflightwifi_service	90917 non-null	object
7	Inflight_entertainment	90917 non-null	object
8	Online_support	90917 non-null	object
9	Ease_of_Onlinebooking	90917 non-null	object
10	Onboard_service	83738 non-null	object
11	Leg_room_service	90917 non-null	object
12	Baggage_handling	90917 non-null	object
13	Checkin_service	90917 non-null	object
14	Cleanliness	90917 non-null	object
15	Online_boarding	90917 non-null	object
16	Gender	90917 non-null	object
17	CustomerType	81818 non-null	object
18	Age	90917 non-null	int64
19	TypeTravel	81829 non-null	object
20	Class	90917 non-null	object
21	Flight_Distance	90917 non-null	int64
22	DepartureDelayin_Mins	90917 non-null	int64
23	ArrivalDelayin_Mins	90633 non-null	float64

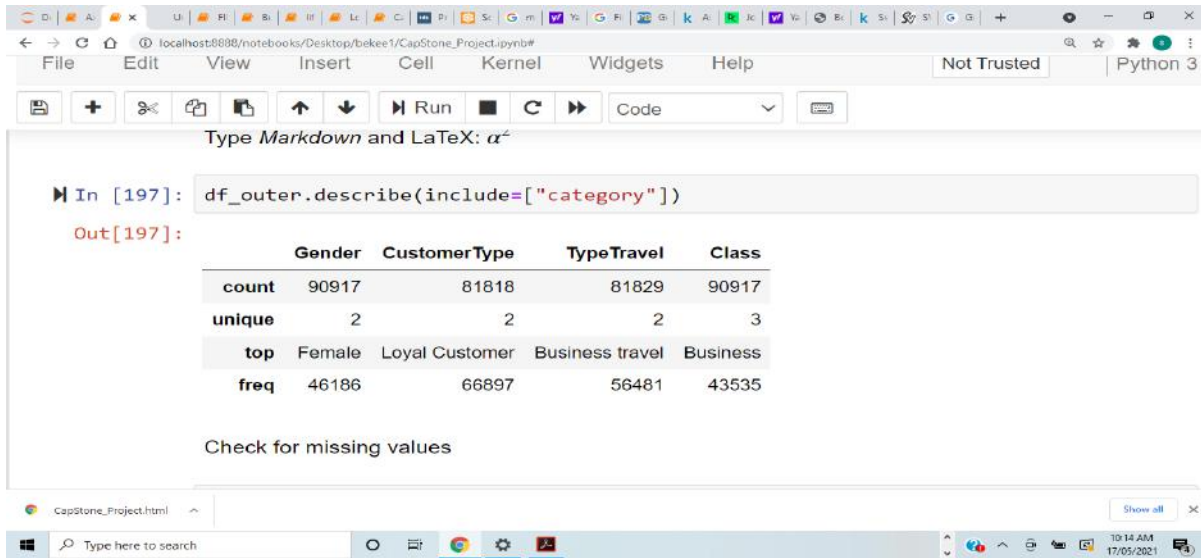
Objects will need to be converted to category as demonstrated with a few variable.

DATA TYPES: float64(1), int64(4), object(19)

### Exploratory Data Analysis



Mean Age respondents is 39.4 and Standard deviation 15.1, mean light distance 1981.6 & Standard deviation 1026.8, mean DepartureDelayin\_Mins 14.6 Std 38. ArrivalDelayin\_Mins 15.1 Std 39.0



Gender have unique count of 2, top being female, CustomerType have unique count of 2, top being Loyal Customer with count indication some value are missing same to TypeTravel and Class with 3 unque count

### Data set

The problem consists of 2 separate datasets: Flight data & Survey data. The flight data has information related to passengers and the performance of flights in which they travelled. The survey data is the

aggregated data of surveys collected post service experience. You are expected to treat both the datasets as raw data and perform any necessary cleaning/validation steps as required

### Data Information

<class 'pandas.core.frame.DataFrame'>

Int64Index: 90917 entries, 0 to 90916

Data columns (total 24 columns):

# Column Non-Null Count Dtype

```
-----
0 CustomerId                90917 non-null int64
1 Satisfaction              90917 non-null object
2 Seat_comfort              90917 non-null object
3 Departure_Arrival_time_convenient 82673 non-null object
4 Food_drink                82736 non-null object
5 Gate_location             90917 non-null object
6 Inflightwifi_service      90917 non-null object
7 Inflight_entertainment    90917 non-null object
8 Online_support            90917 non-null object
9 Ease_of_Onlinebooking     90917 non-null object
10 Onboard_service          83738 non-null object
11 Leg_room_service         90917 non-null object
12 Baggage_handling         90917 non-null object
13 Checkin_service          90917 non-null object
14 Cleanliness              90917 non-null object
15 Online_boarding          90917 non-null object
16 Gender                   90917 non-null object
17 CustomerType             81818 non-null object
18 Age                      90917 non-null int64
19 TypeTravel               81829 non-null object
```

20 Class	90917 non-null object
21 Flight_Distance	90917 non-null int64
22 DepartureDelayin_Mins	90917 non-null int64
23 ArrivalDelayin_Mins	90633 non-null float64

Objects will need to be converted to category as demonstrated with a few variable.  
 DATA TYPES: float64(1), int64(4), object(19)

### Data Preprocessing

The data set contains CustomerId variable which of type int, in my opinion I feel this variable may not add any significant value to prediction variable so I decided to drop it.

```
df_outer.drop(
  columns=["CustomerId"], inplace=True)
```

### Fixing Categorical Variables (Data Type)

Before putting our data through models, two steps that need to be performed on categorical data is encoding and dealing with missing nulls. Encoding is the process of converting text or boolean values to numerical values for processing. This approach was adopted for 19 columns that

are of type object. First the object were converted to Categories and the Ordinal values were passed as shown below;  
 categorical\_variables=df\_outer.select\_dtypes(exclude=["number","bool\_"]).columns.tolist()#  
*list of categorical variables*

forcolmnincategorical\_variables:

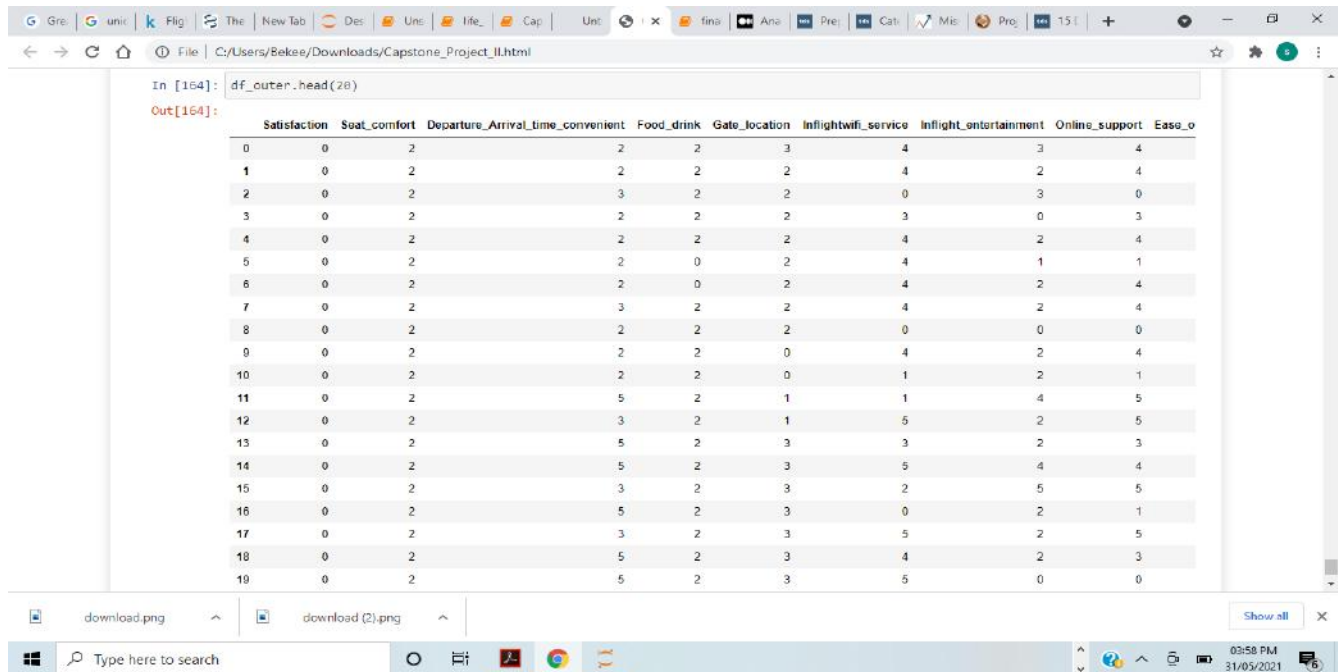
```
df_outer[colmn]=df_outer[colmn].astype('category')
```

<class 'pandas.core.frame.DataFrame'>

Int64Index: 90917 entries, 0 to 90916

Data columns (total 23 columns):

#	Column	Non-Null Count	Dtype
0	Satisfaction	90917 non-null	category
1	Seat_comfort	90917 non-null	category
2	Departure_Arrival_time_convenient	82673 non-null	category
3	Food_drink	82736 non-null	category
4	Gate_location	90917 non-null	category
5	Inflightwifi_service	90917 non-null	category
6	Inflight_entertainment	90917 non-null	category
7	Online_support	90917 non-null	category
8	Ease_of_Onlinebooking	90917 non-null	category
9	Onboard_service	83738 non-null	category
10	Leg_room_service	90917 non-null	category
11	Baggage_handling	90917 non-null	category
12	Checkin_service	90917 non-null	category
13	Cleanliness	90917 non-null	category
14	Online_boarding	90917 non-null	category
15	Gender	90917 non-null	category
16	CustomerType	81818 non-null	category
17	Age	90917 non-null	int64
18	TypeTravel	81829 non-null	category
19	Class	n-null	



### Missing Value

Missing value in a dataset is a very common phenomenon in the reality. Missing value correction was carried out on the data set to reduce bias and to produce data set suitable

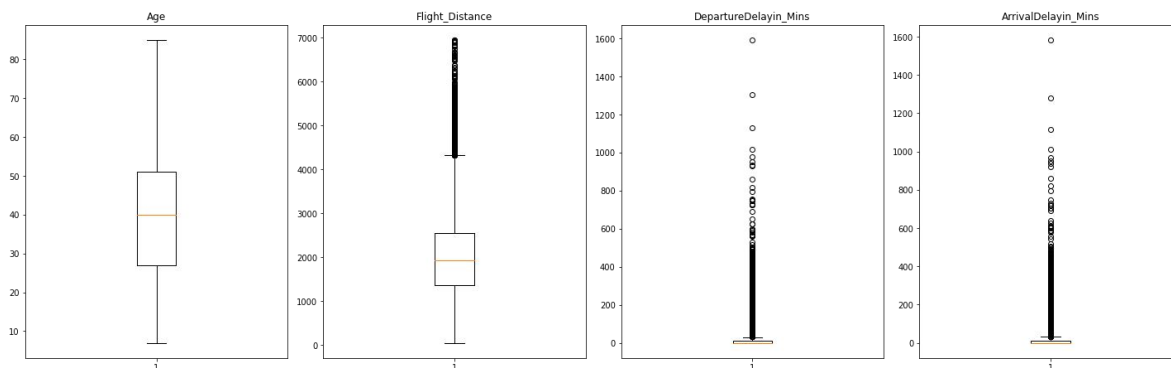
: Departure\_Arrival\_time\_convenient 8244  
 Food\_drink 8181  
 Onboard\_service 7179  
 CustomerType 9099  
 TypeTravel 9088  
 ArrivalDelayin\_Mins 284

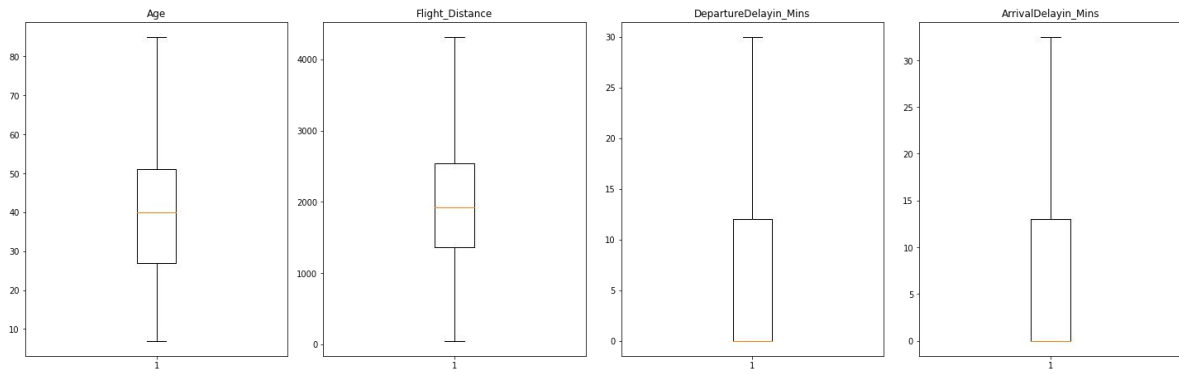
modelling. The following variables contained missing values

They were corrected with the following codes for numerical and categorical respectively

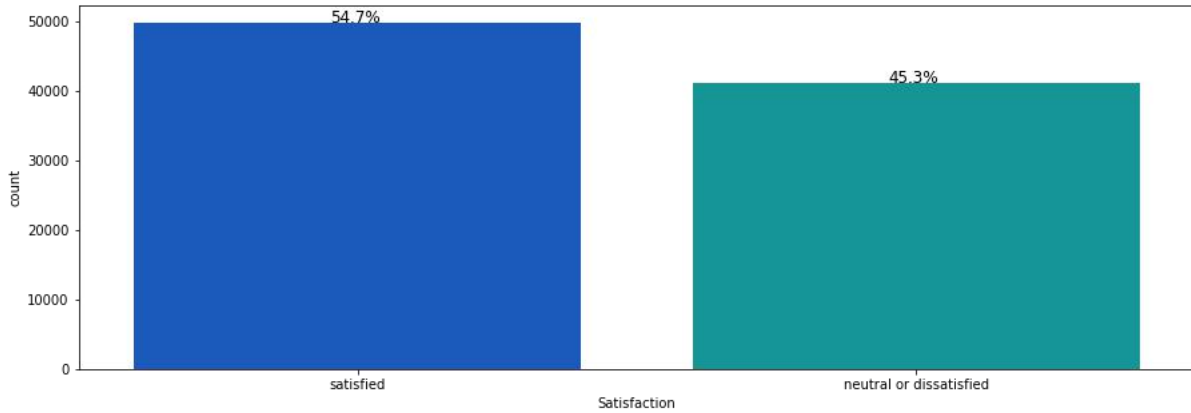
### Outlier Detection & Treatment

Outlier detection and treatment was carried out on the dataset, the outcome is shown below;

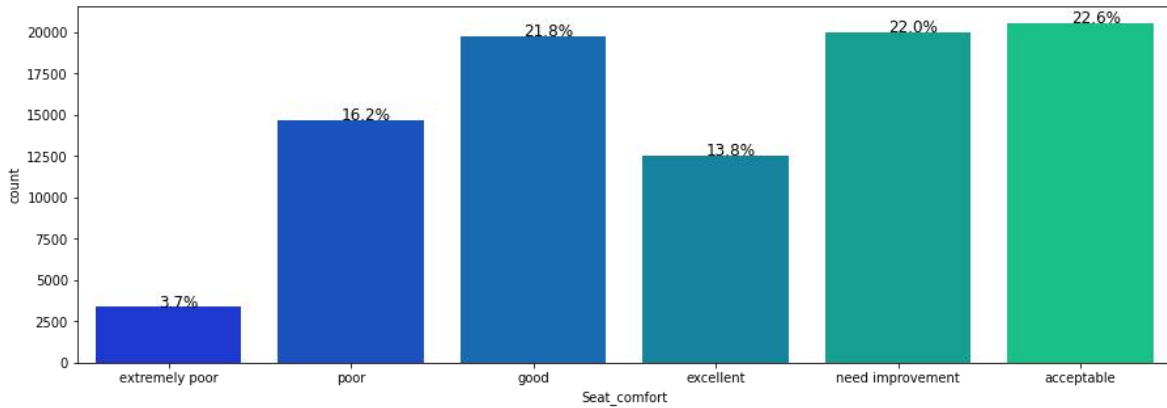




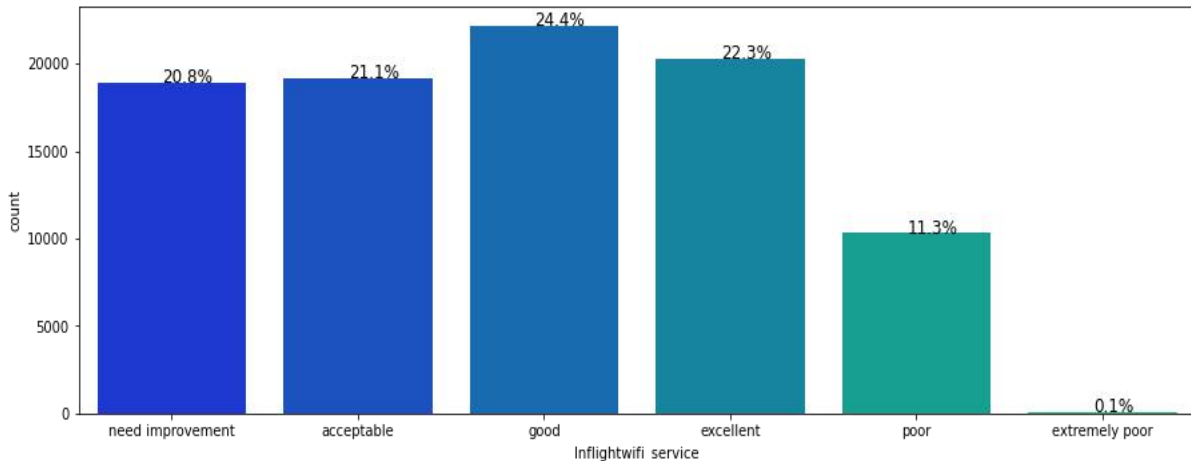
### Exploratory Data Analysis



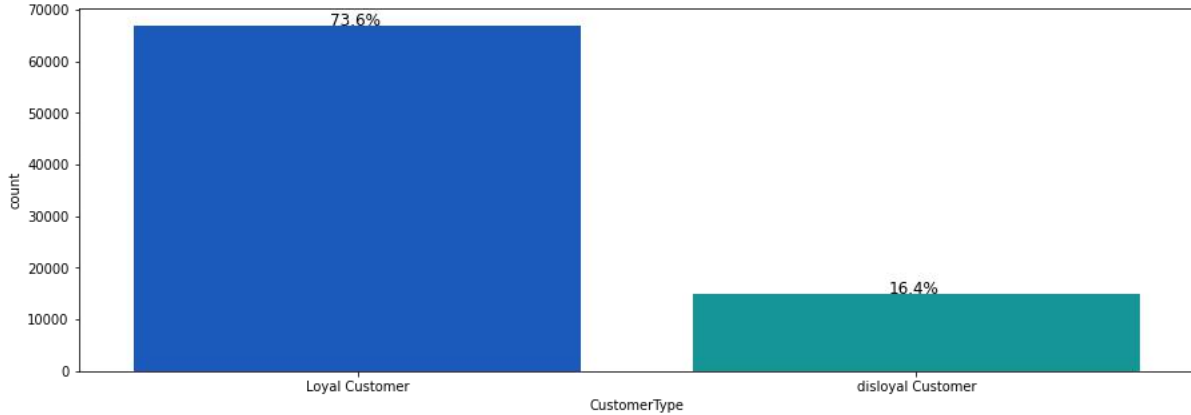
54.7 % of the passengers are satisfied while 45.3% are of the respodent fall into the group neutral or dissatisfied



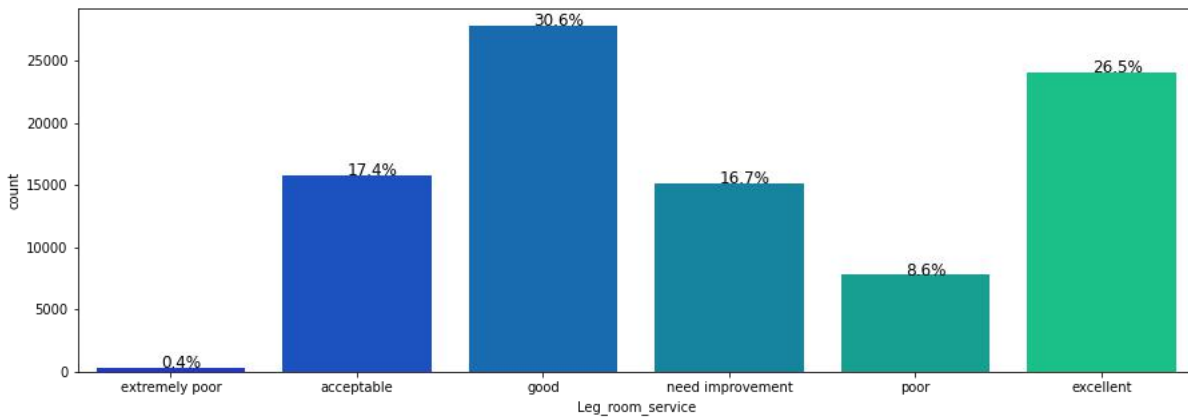
Seat\_Comfort 22.6% respodent acceptable, 22% indicate seat need improvement, 21.8% indicate good, 13.8 % excellent



inflight service has six categories with 24% good, 22.3 excellent and 21%, 20.8 acceptable & need improvement repectively



About 74% the repondent are loyal customers of the



Airline Good has the has highest percentage, followed by excellent and acceptable with 17 & 16.7 respectively

**Alternate Analytical Approach**

The two objective of this project are-  
 1.To understand which parameters play an important role in swaying a passenger feedback towards ‘satisfied’.

**Building Model**

Logistic regression is a is a type of supervised machine learning used to predict the probability of a target variable. The most common logistic regression models a binary outcome; The output of the dependent variable is represented in discrete values such as 0 and 1, true/false, yes/no, etc. In some case logistic regression can model scenarios

Security are classification problems, such as attack detection, logistic regression is a useful analytic technique.

With the above explanation, it is clear that logistic regression modelling technique suit the Airline passenger satisfaction prediction problem.

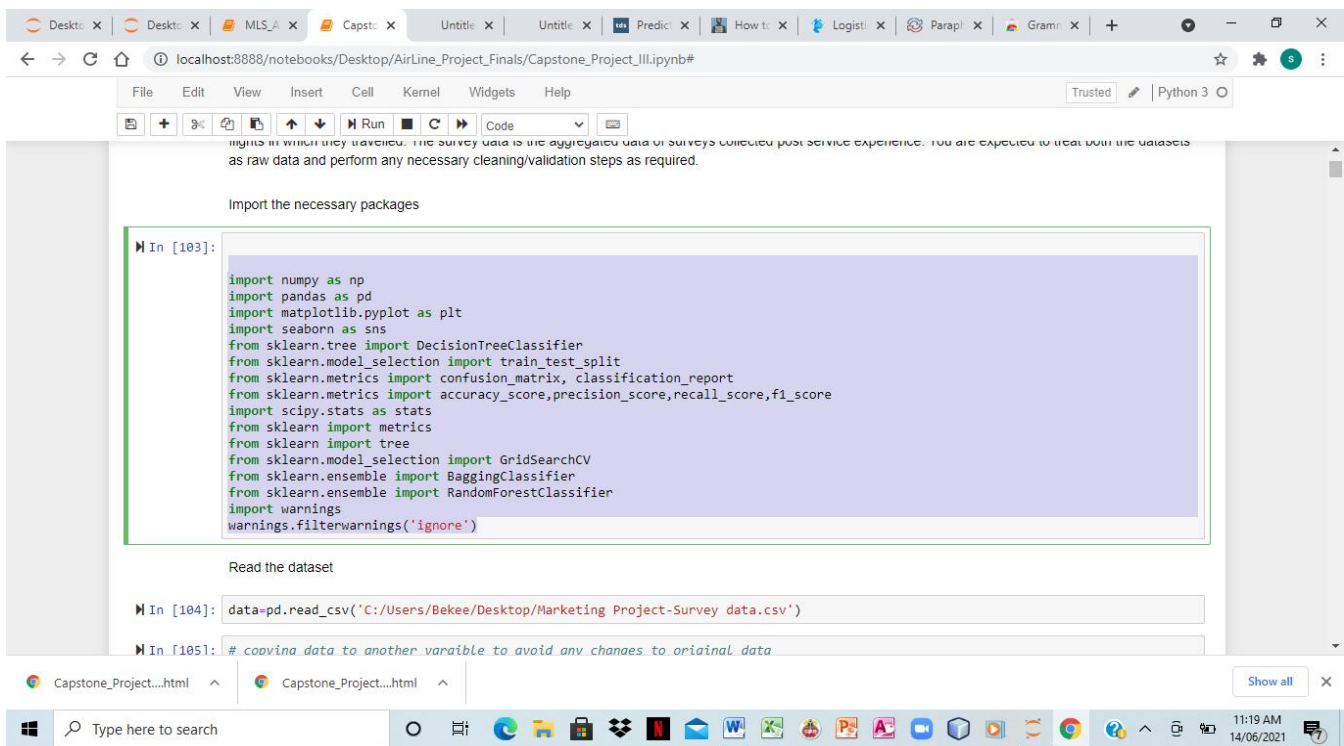
**Importing the libraries**

To build our model, the first step is to import the necessary libraries. I used the Pandas library to load in the CSV or the dataset, and Numpy to convert the data frame into arrays.

2. To predict whether a passenger will be satisfied or not given the rest of the details that are provided in the data set.

The alternate approach to this problem would be to build, test and implement a classification model.

where there are more than two possible discrete outcomes that is referred to as Multinomial logistic regression. Logistic regression is a useful analysis method for classification problems, where you are trying to determine if a new sample fits best into a category. As aspects of cyber



The second step was to define the target variable(Y) and the independent variables, and then split the data set into the training set and the test set. We will use the training set to train our logistic regression algorithm. Similarly, the test data set will be used to validate the logistic regression model.

To split the data into two sets, we will use Sklearn. The train\_split\_function can be used and we can specify the amount of data we want to set aside for training and testing

```

From sklearn.model_selection import train_test_split
# Splitting data into training and test set:
X_train, X_test, y_train, y_test =train_test_split(X, y, test_size=0.3, random_state=1,stratify=y)
print(X_train.shape, X_test.shape)
(63641, 22) (27276, 22)

```

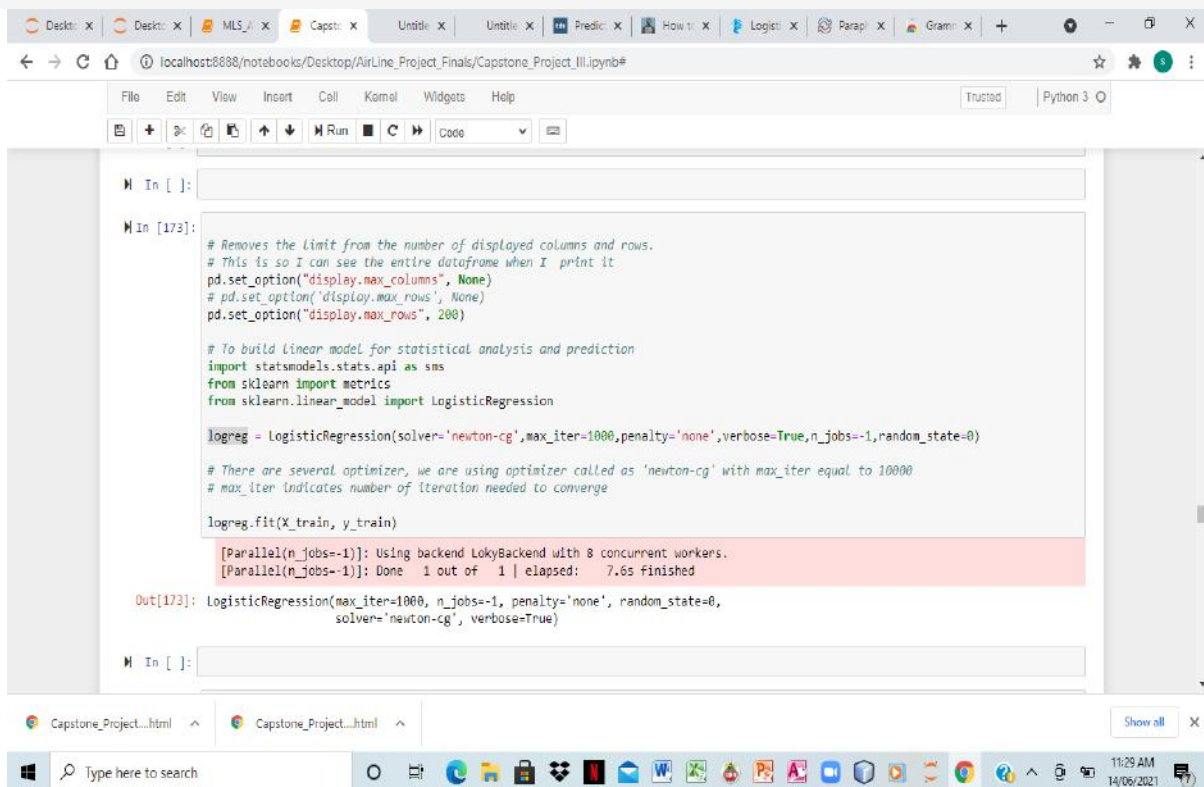
### Building The Logistic Regression Mode

Next was to build the logistic regression model and fit it to the training data set. First, we will need to import the logistic regression algorithm from Sklearn

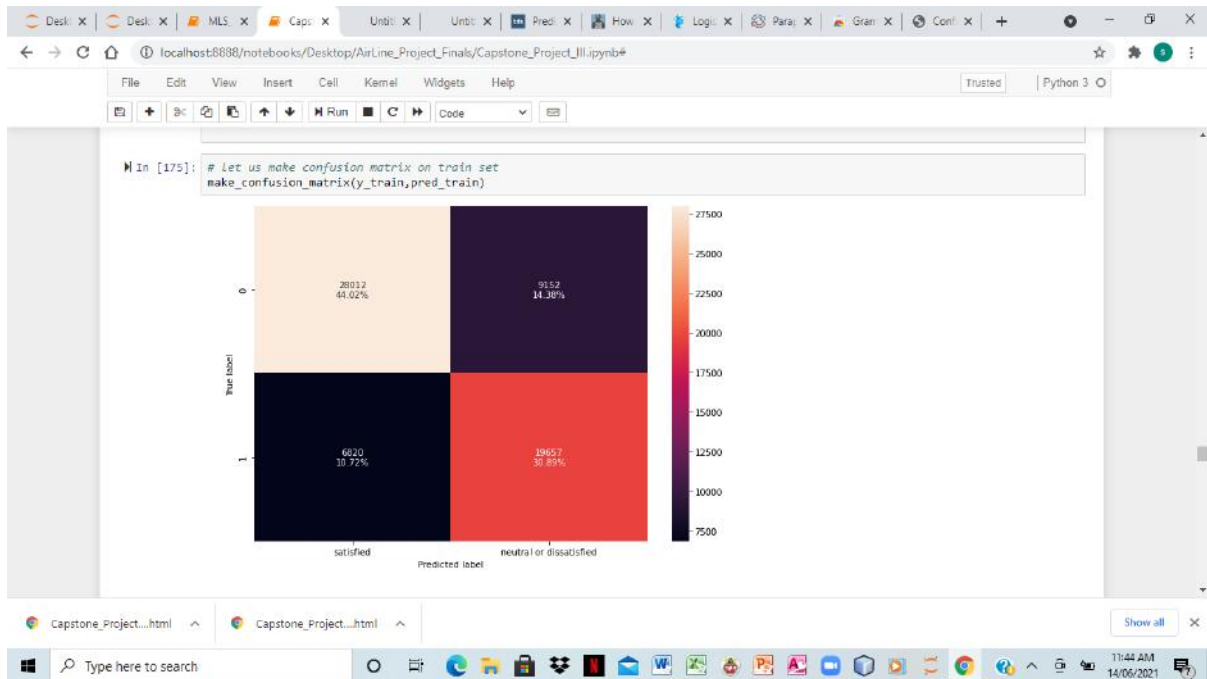
Accuracy

0.7490

$$ACC = (TP + TN) / (P + N)$$

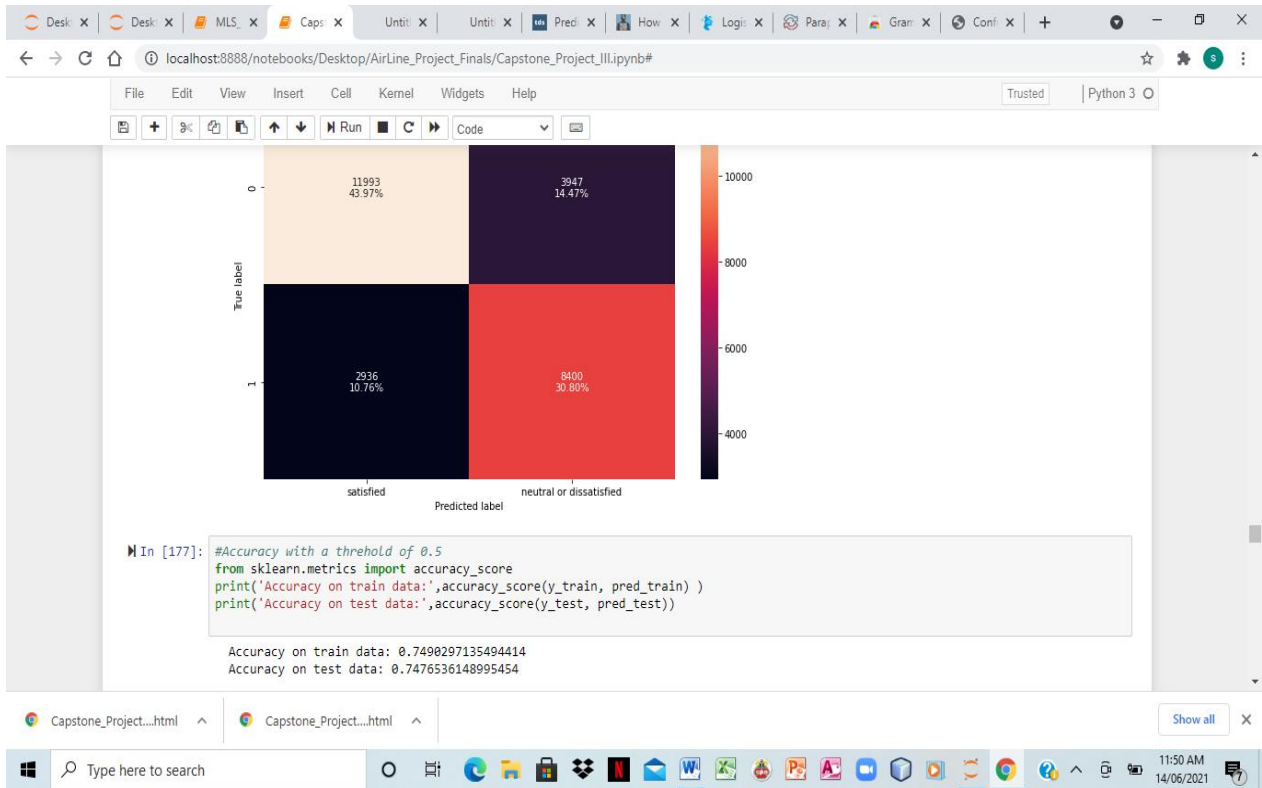


Next, was to create predictions on the test dataset and train data with the model performance as shown



low Accuracy of 74% and precision 75%

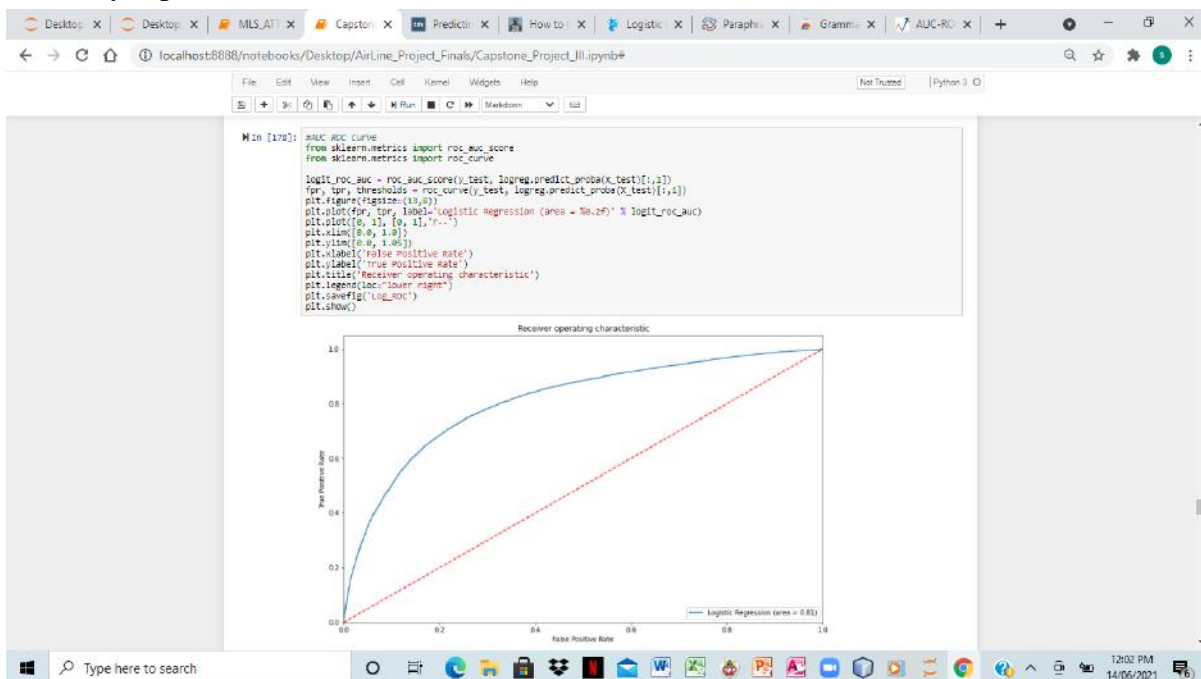




Next the AUC-ROC curve was used to visualize how well our machine learning classifier is performing.

The Receiver Operator Characteristic (ROC) curve is an evaluation metric for binary classification problems. It is a probability curve that plots the TPR against FPR at various threshold values and essentially separates the

‘signal’ from the ‘noise’. The Area Under the Curve (AUC) is the measure of the ability of a classifier to distinguish between classes and is used as a summary of the ROC curve. The result obtained is as below



The AUC of 0.81% looks quite good performance .

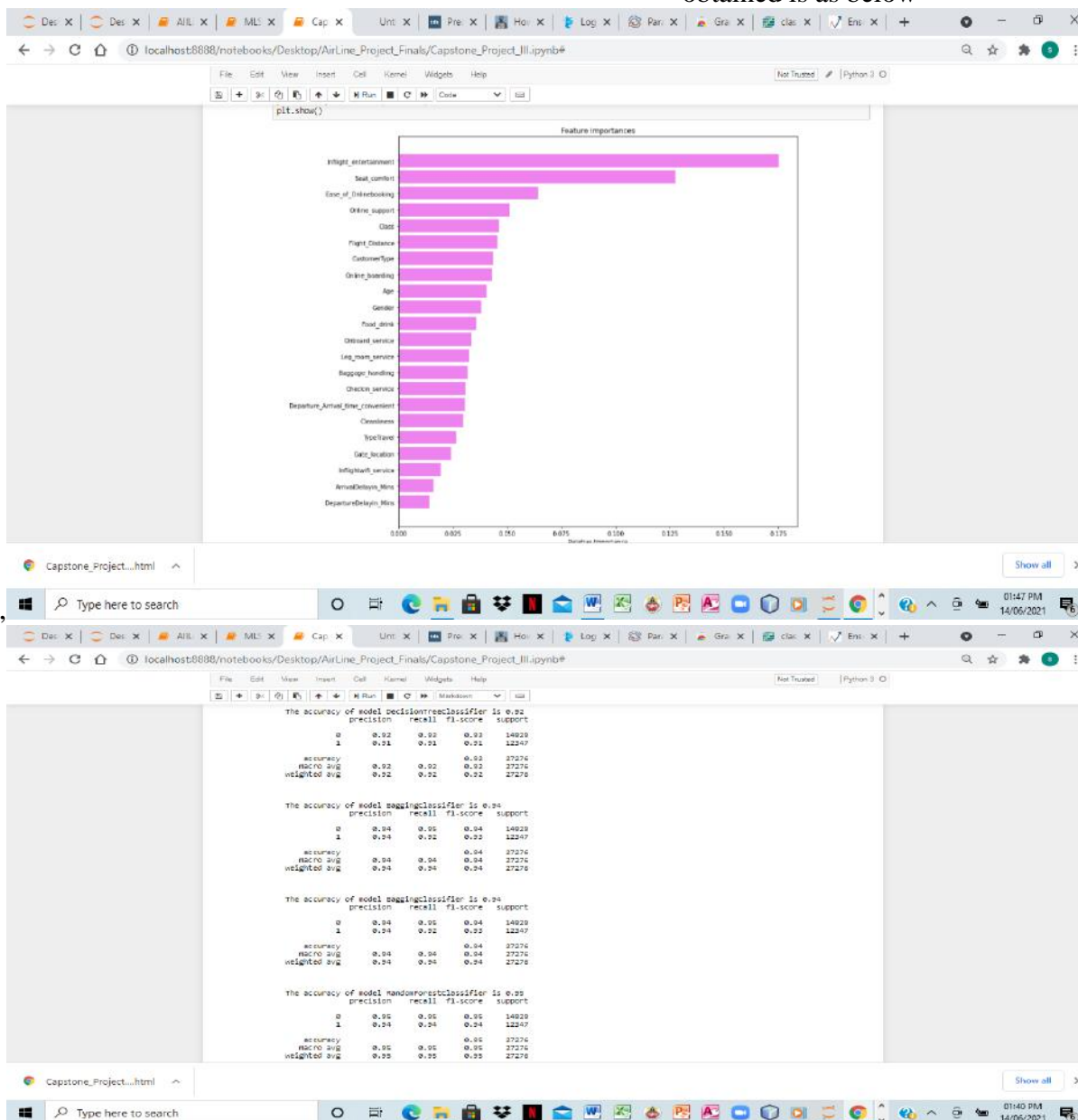
## Build Different Models

DecisionTreeClassifier, RandomForestClassifier, BaggingClassifier were used on the data set. A result of the comparison of the different models are as bellow The best model is the rf, rf\_wt

performance of 95% accurate prediction on Recall and Precession.

Feature Importance

The RandomForestClassifier was used to check for featre importance and the result obtained is as below

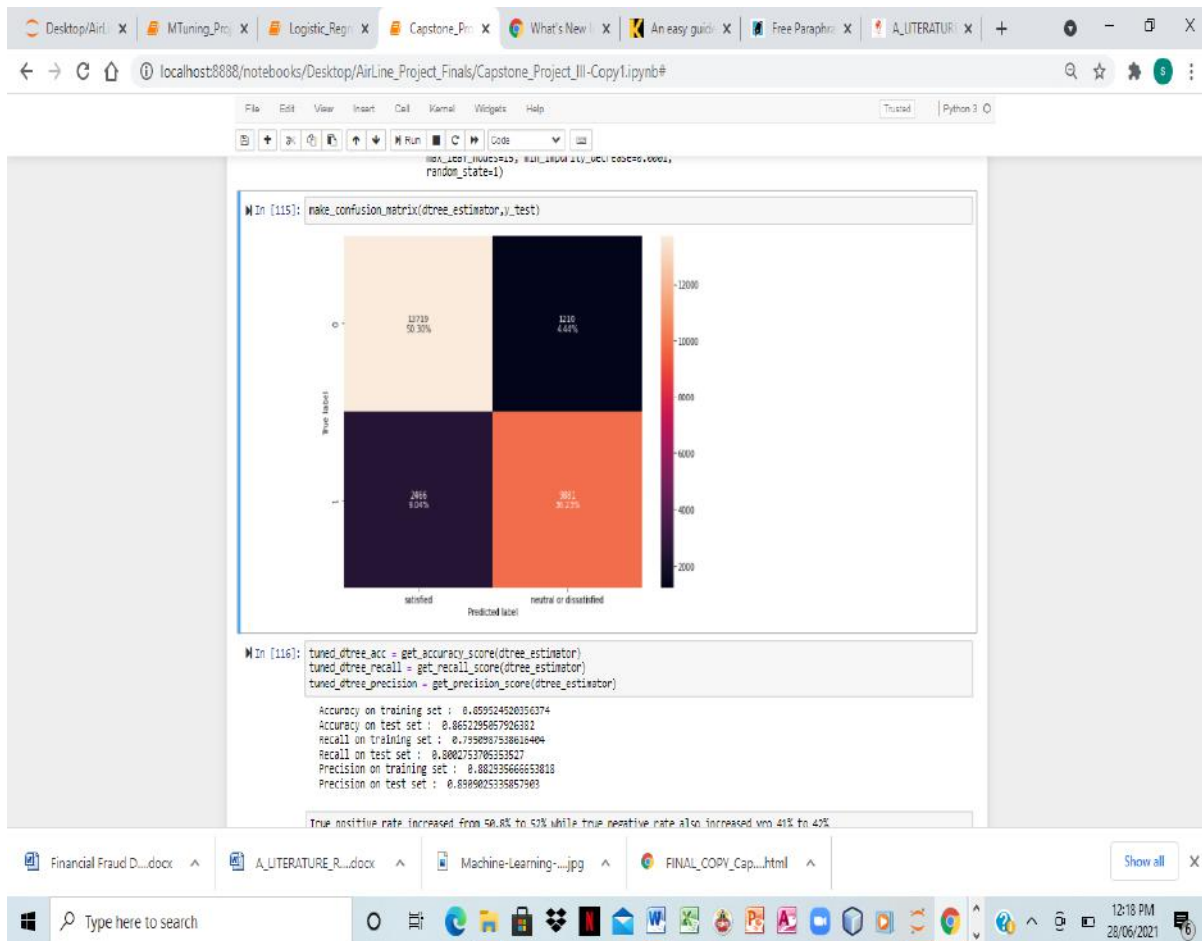


Inflight\_entertainment is the most important feature for prediction followed by Seat\_comfort ,Ease\_of\_Onlinebooking and Online\_support.

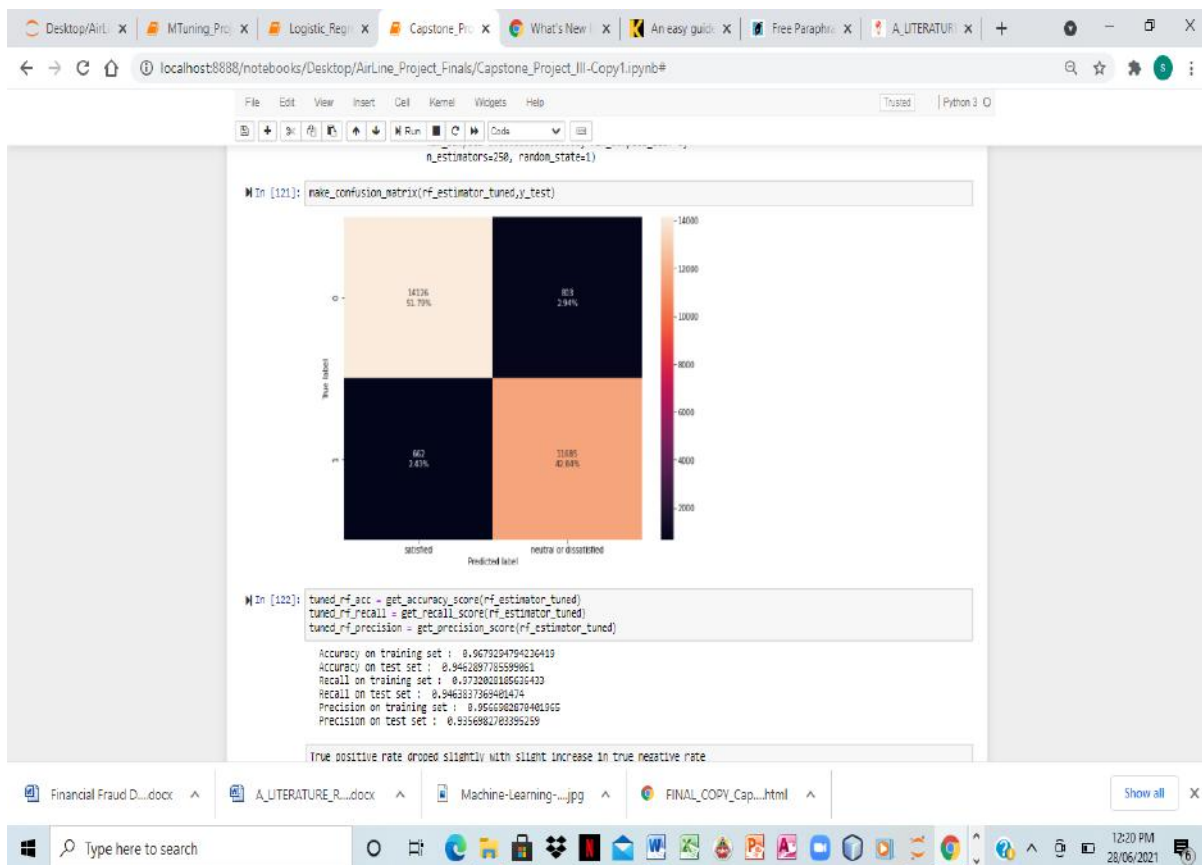
## Tuned Models Compared Alongside Initial Models

The results of the comparison of the different models are as bellow.

### Tuned Decision tree



### TUNNED RANDOM FOREST



### Model Comparison Results.

[Parallel(n\_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

[Parallel(n\_jobs=-1)]: Done 1 out of 1 | elapsed: 11.9s finished

The accuracy of model LogisticRegression is 0.75

	precision	recall	f1-score	support
0	0.75	0.80	0.78	14929
1	0.74	0.68	0.71	12347

accuracy 0.75 27276

macroavg 0.75 0.74 0.74 27276

weightedavg 0.75 0.75 0.75 27276

The accuracy of model DecisionTreeClassifier is 0.92

	precision	recall	f1-score	support
0	0.92	0.93	0.93	14929
1	0.91	0.91	0.91	12347

accuracy 0.92 27276

macroavg 0.92 0.92 0.92 27276

weightedavg 0.92 0.92 0.92 27276

The accuracy of model DecisionTreeClassifier is 0.87

	precision	recall	f1-score	support
0	0.85	0.92	0.88	14929
1	0.89	0.80	0.84	12347

accuracy 0.87 27276

macroavg 0.87 0.86 0.86 27276

weightedavg 0.87 0.87 0.86 27276

The accuracy of model BaggingClassifier is 0.94

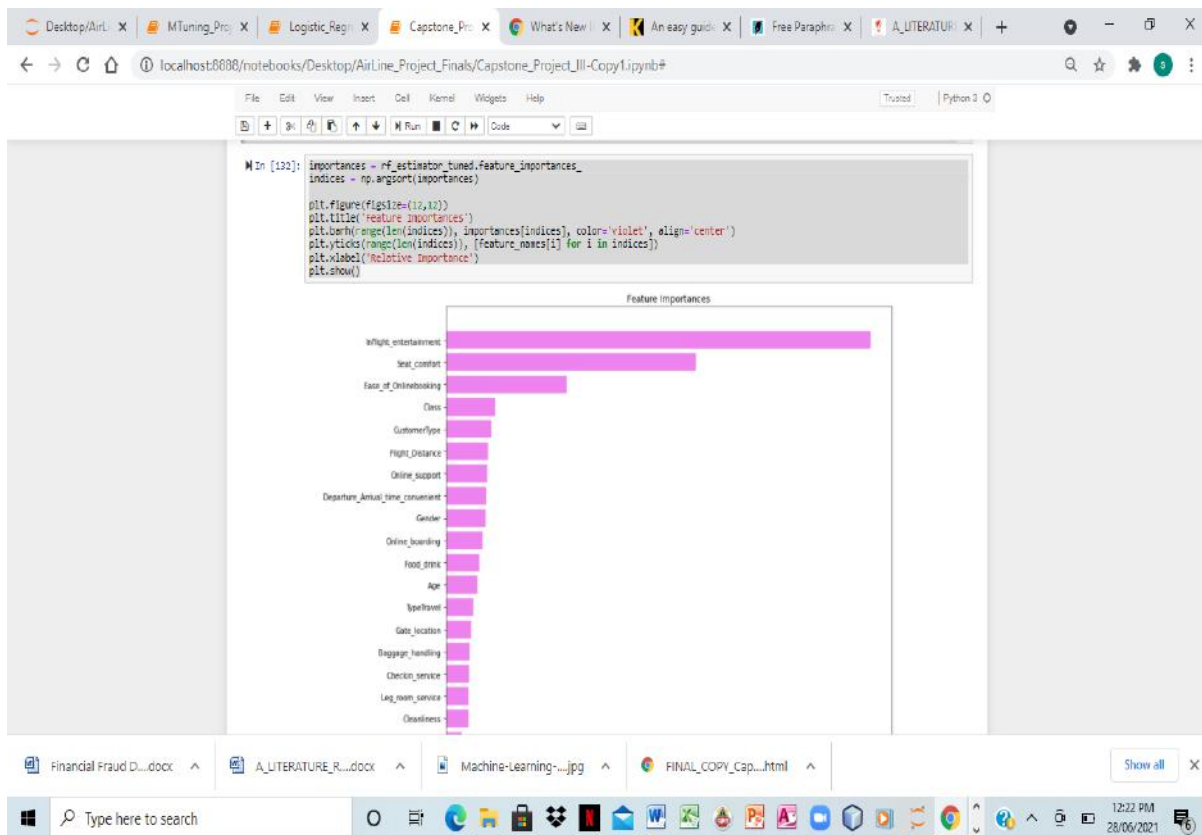
precision recall f1-score support

```

95  0.94  14929
   1  0.94  0.92  0.93  12347
accuracy                0.94  27276
macroavg  0.94  0.94  0.94  27276
weightedavg  0.94  0.94  0.94  27276
The accuracy of model BaggingClassifier is 0.95
precision recall f1-score support
   0  0.96  0.95  0.95  14929
   1  0.94  0.95  0.95  12347
accuracy                0.95  27276
macroavg  0.95  0.95  0.95  27276
weightedavg  0.95  0.95  0.95  27276
The accuracy of model RandomForestClassifier is 0.95
precision recall f1-score support
   0  0.95  0.95  0.95  14929
   1  0.94  0.94  0.94  12347
accuracy                0.95  27276
macroavg  0.95  0.95  0.95  27276
weightedavg  0.95  0.95  0.95  27276
The accuracy of model RandomForestClassifier is 0.95
precision recall f1-score support
   0  0.95  0.95  0.95  14929
   1  0.94  0.94  0.94  12347
accuracy                0.95  27276
macroavg  0.95  0.95  0.95  27276
weightedavg  0.95  0.95  0.95  27276
The accuracy of model DecisionTreeClassifier is 0.92
precision recall f1-score support
   0  0.92  0.93  0.93  14929
   1  0.91  0.91  0.91  12347
accuracy                0.92  27276
macroavg  0.92  0.92  0.92  27276
weightedavg  0.92  0.92  0.92  27276
The accuracy of model RandomForestClassifier is 0.95
precision recall f1-score support
   0  0.96  0.95  0.95  14929
   1  0.94  0.95  0.94  12347
accuracy                0.95  27276
macroavg  0.95  0.95  0.95  27276
weightedavg  0.95  0.95  0.95  27276
The accuracy of model BaggingClassifier is 0.95
The accuracy of model RandomForestClassifier is 0.95

```

### **Randomforestclassifier Important Features**



Inflight\_flight entertainment is the most important variable for predicting airline passengersatisfactionfollowed by seat\_comfort, ease\_of\_online booking, Class, customerType and Flight\_Distance.

### Conclusion

A predictive classification model has been built. Given the performance in terms of Precision and Recall, the model can be deployed to identify Passengers who may not be satisfied or are indifferent to the Airline services and shall take appropriate actions to build and improve on services that drive passengers satisfaction. Factors that drive satisfaction - Inflight\_entertainment, Seat\_comfort,Ease\_of\_Onlinebooking and Online\_support.

### Business Insights and Recommendations

I recommend that airlines should focus on improving the Inflight\_entertainment experience.

In addition, airlines should also focus on Ease of Online Booking, as business passengers prioritize on ease and convenience in their travel.Finally, I hope that the model will provide a reference for airlines and be utilized for business value

## References

[https://www.researchgate.net/figure/Machine-Learning-Model-Training-Process\\_fig1\\_336020567/download](https://www.researchgate.net/figure/Machine-Learning-Model-Training-Process_fig1_336020567/download) retrieved 28/06/2021

Park, Y., Gates, S.: Towards real-time measurement of customer satisfaction using automatically generated call transcripts. In: Proceedings of the 18th ACM Conference on Information and Knowledge Management, pp. 1387–296. ACM (2009)

<https://en.wikipedia.org/wiki/Twitter> Accessed 28 June 2021

Abel F, Gao Q, Houben G-J, Tao K (2013) Twitter-based user modeling for news recommendations. In: Rossi F (ed)

IJCAI 2013, proceedings of the 23rd international joint conference on artificial intelligence, Beijing, China, August 3–9, 2013. IJCAI/AAAI.

<https://www.quora.com/How-many-airplanes-fly-each-day-in-the-world> retrieved 28/06/2021

Jehad Ali, Rehanullah Khan, Nasir Ahmad, Imran Maqsood, "Random Forests and Decision Trees", IJCSI

International Journal of Computer Science Issues, Vol. 9, Issue 5, No 3, September 2012.

Lecture slides and Jupyter notebook on Supervised learning and Model tuning file from Greatlearning.