**INTRODUCTION TO DATA MINING**
**FINAL PROJECT**
**AIRLINE PASSENGER SATISFACTION - CLASSIFICATION**

**GROUP MEMBERS:**
SAFA ABDUL HAI - 18625
MARIUM AFZAL -19756
MAHUM FATIMA KHAN - 19747

**Introduction:**
Recently, there has been a surge in interest in the field of data mining, where the goal is to predict correct and useful knowledge for users. We used a summary of our study and exploration of the Arline passenger satisfaction Data in this project to come up with useful, important, and intriguing data properties. The CRISP-DM data mining methodology was heavily used in the development of our project model to predict as accurate as possible statistics for our users.

**Problem Description:**
Passenger satisfaction is a high concern for airlines' management as their prime concern is whether the customer will choose their airlines for their next trip or not. Airline passenger satisfaction is affected by many factors, but at its root, this type of customer satisfaction is no different from that of any other business. Therefore, various data mining techniques are utilized in this project to build a classification model for Airlines management to determine if the customer is satisfied or not depending on various factors.

**Data Collection:**
The dataset, primarily known as "Airline Passenger Satisfaction" was found on Kaggle, https://www.kaggle.com/teejmahal20/airline-passenger-satisfaction, while we were searching for a classification data set for our Data Mining project. This data set was collected by TJ Klein in 2019 through a survey.

**Data Description:**
The data set has two csv files named Test and Train. The test.csv file contains 26k records while the train.csv file contains 104k records. There are 25 attributes in total. The following are the given attributes:

Gender: Gender of the passengers (Female, Male)

Customer Type: The customer type (Loyal customer, disloyal customer)

Age: The actual age of the passengers

Type of Travel: Purpose of the flight of the passengers (Personal Travel, Business Travel)

Class: Travel class in the plane of the passengers (Business, Eco, Eco Plus)

Flight distance: The flight distance of this journey

Inflight wifi service: Satisfaction level of the inflight wifi service (0:Not Applicable;1-5)

Departure/Arrival time convenient: Satisfaction level of Departure/Arrival time convenient

Ease of Online booking: Satisfaction level of online booking

Gate location: Satisfaction level of Gate location

Food and drink: Satisfaction level of Food and drink

Online boarding: Satisfaction level of online boarding

Seat comfort: Satisfaction level of Seat comfort

Inflight entertainment: Satisfaction level of inflight entertainment

On-board service: Satisfaction level of On-board service

Leg room service: Satisfaction level of Leg room service

Baggage handling: Satisfaction level of baggage handling

Check-in service: Satisfaction level of Check-in service

Inflight service: Satisfaction level of inflight service

Cleanliness: Satisfaction level of Cleanliness

Departure Delay in Minutes: Minutes delayed when departure

Arrival Delay in Minutes: Minutes delayed when Arrival

Satisfaction: Airline satisfaction level(Satisfaction, neutral or dissatisfaction)

**Data Preprocessing:**
The steps taken to prepare the data are as follow:

1.  **Number to String:**
    All the attributes except the continuous attributes such as Age, Flight Distance, Departure
    Delay in Minutes, Arrival Delay in Minutes are converted to String data type.



2.  **Column Filter:**
    The 4 columns ID, Gender, Age, and Flight Distance are filtered out using Column Filter
    as including these columns did not have much effect on our prediction/ accuracy.

### 3. Domain Calculator:

After applying the previous techniques on the dataset, the domain information of the data has changed. So we have used the Domain Calculator for this purpose.

**Model Building:**

After preprocessing the entire data, we used four classifiers to build our model. They are as follow:

    **1. Decision Tree:** We have used the following configuration on the Decision Tree Learner



**ROC:**

**2. Random Forest:** We have used the following configuration on the Random Forest Learner



**ROC:**

1. **Naive Bayes:** The following Configuration was used for Naïve Bayes



**ROC:**

## 2. Gradient Boost: The following configuration was used for gradient boosting



**ROC:**





Correct classified: 24,639  Wrong classified: 1,337

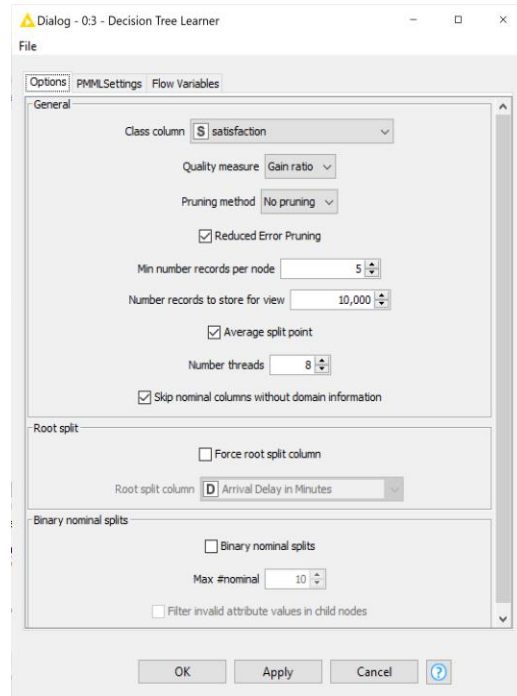Accuracy: 94.853 %  Error: 5.147 %

Cohen's kappa (κ) 0.896

| satisfactio... | satisfied | neutral or ... |
|---|---|---|
| satisfied | 10735 | 668 |
| neutral or di... | 669 | 13904 |

## Attempt 2:

## Data Preprocessing:

1. **Column Filter:** From our previous attempt we found out that the flight distance and age are important factors therefore their columns should be included and gate location was excluded since it did not make any significant difference in our prediction.

**Model Building:**

After changing the column filter's configuration, we again checked the accuracy using four classifiers. They are as follow:

   1.  **Decision Tree:**



**ROC:** The accuracy decreased a little by 0.01 percent from approximately 98% to 97%.

## 2. Random Forest:



**ROC & Scorer:** The accuracy increased.



Confusion Matrix - 0:38 - Scorer

| satisfactio... | satisfied | neutral or ... |
|---|---|---|
| satisfied | 10699 | 704 |
| neutral or di... | 236 | 14337 |

Correct classified: 25,036    Wrong classified: 940

Accuracy: 96.381 %    Error: 3.619 %

Cohen's kappa (κ) 0.926

### 3. Naive Bayes:



## ROC & Scorer:



**Confusion Matrix - 0:40 - Scorer**

| satisfactio... | satisfied | neutral or ... |
|---|---|---|
| satisfied | 9978 | 1425 |
| neutral or di... | 1440 | 13133 |

Correct classified: 23,111     Wrong classified: 2,865

Accuracy: 88.971 %     Error: 11.029 %

Cohen's kappa (κ) 0.776

P (satisfaction=satisfied) (0.9537)

## 4. Gradient Boosting:



## ROC & Scorer:

**Attempt 3:**

**Data Preprocessing:**

1. **Row Filter**

Keeping all the above conditions applied constant we applied a new filter -row filter- where we set the range of the age from 17 to 75 i.e. only the responses of people whose age lies in this range matter. People less than 18 are considered kids whose respond does not hold much value.
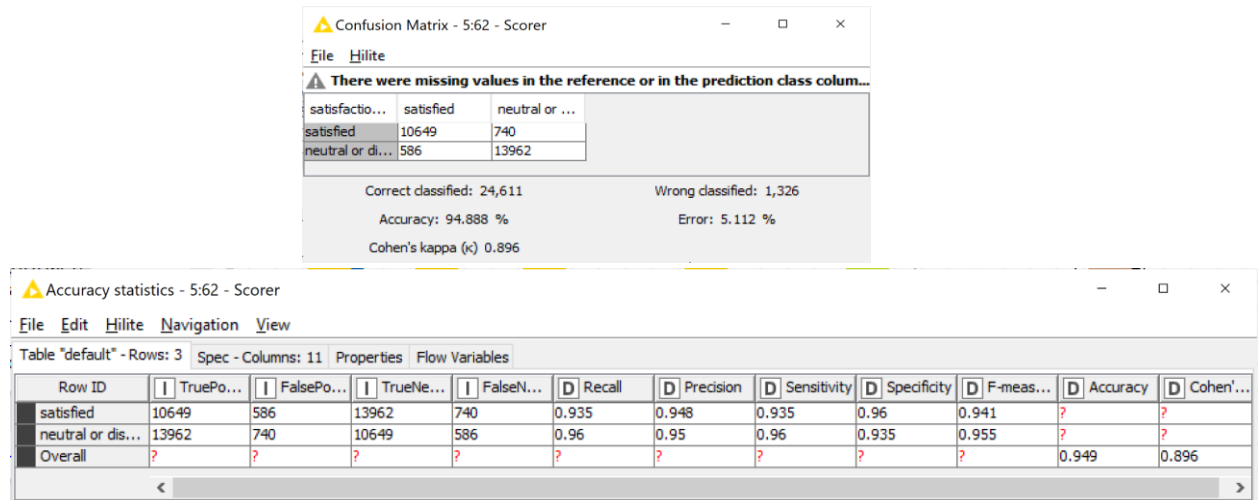


2. **Missing Value:** It was found that we had 310 missing data in arrival delay time so we used missing value filter.

**Models:** The configurations for all the models were kept constant, only a few filters mentioned above were added to our models. However the accuracy of all models were effected at most 0.01. Therefore it was concluded that applying the row filter did not have much of an effect.

1. **Decision Tree**

**Scorer**



**ROC:**

## 2. Random Forest

## ROC & Scorer:





| Row ID | TruePo... | FalsePo... | TrueNe... | FalseN... | Recall | Precision | Sensitivity | Specificity | F-meas... | Accuracy | Cohen'... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| satisfied | 10687 | 246 | 14327 | 716 | 0.937 | 0.977 | 0.937 | 0.983 | 0.957 | ? | ? |
| neutral or dis... | 14327 | 716 | 10687 | 246 | 0.983 | 0.952 | 0.983 | 0.937 | 0.968 | ? | ? |
| Overall | ? | ? | ? | ? | ? | ? | ? | ? | ? | 0.963 | 0.924 |

### 3. Naive Bayes

## ROC & Scorer:

**Confusion Matrix - 5:64 - Scorer** — □ ×

File   Hilite

| satisfactio... | satisfied | neutral or ... |
|---|---|---|
| satisfied | 9905 | 1498 |
| neutral or di... | 1432 | 13141 |

Correct classified: 23,046          Wrong classified: 2,930
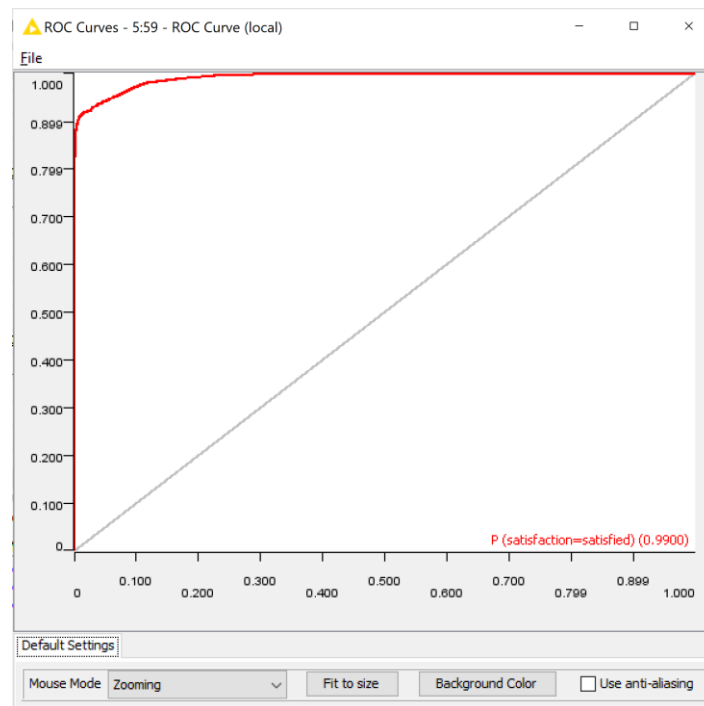
Accuracy: 88.72 %                   Error: 11.28 %

Cohen's kappa (κ) 0.771

**Accuracy statistics - 5:64 - Scorer** — □ ×

File   Edit   Hilite   Navigation   View

Table "default" - Rows: 3   Spec - Columns: 11   Properties   Flow Variables

| Row ID | TruePo... | FalsePo... | TrueNe... | FalseN... | Recall | Precision | Sensitivity | Specificity | F-meas... | Accuracy | Cohen'... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| satisfied | 9905 | 1432 | 13141 | 1498 | 0.869 | 0.874 | 0.869 | 0.902 | 0.871 | ? | ? |
| neutral or dis... | 13141 | 1498 | 9905 | 1432 | 0.902 | 0.898 | 0.902 | 0.869 | 0.9 | ? | ? |
| Overall | ? | ? | ? | ? | ? | ? | ? | ? | ? | 0.887 | 0.771 |

**ROC Curves - 5:27 - ROC Curve (local)** — □ ×

File

P (satisfaction=satisfied) (0.9530)

Default Settings

Mouse Mode   Zooming ∨   Fit to size   Background Color   ☐ Use anti-aliasing

## 4. Gradient Boosting

**ROC And Scorer:**

### Confusion Matrix - 5:65 - Scorer

File   Hilite

| satisfactio... | satisfied | neutral or ... |
|---|---|---|
| satisfied | 10719 | 684 |
| neutral or di... | 694 | 13879 |

Correct classified: 24,598          Wrong classified: 1,378

Accuracy: 94.695 %                      Error: 5.305 %

Cohen's kappa (κ) 0.892

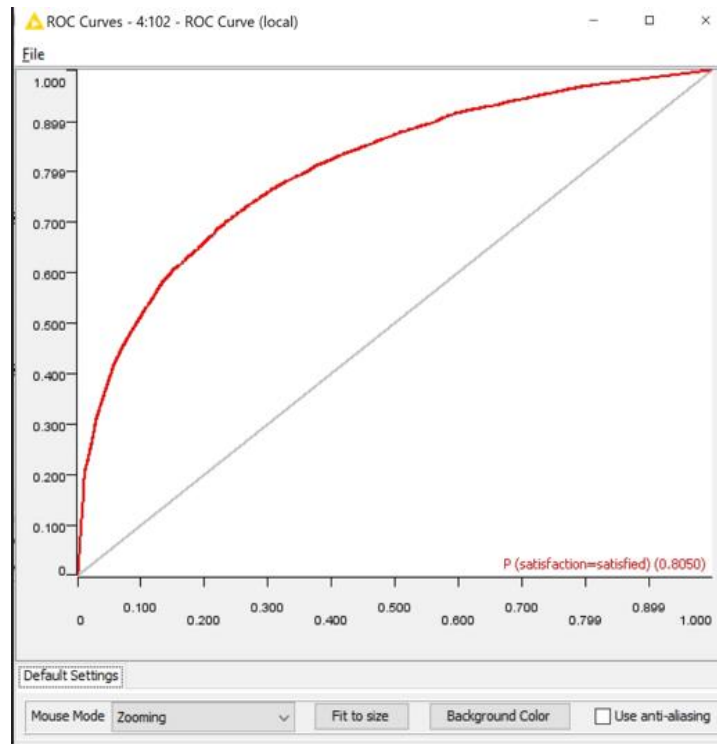### Accuracy statistics - 5:65 - Scorer

File   Edit   Hilite   Navigation   View

Table "default" - Rows: 3   Spec - Columns: 11   Properties   Flow Variables

| Row ID | TruePo... | FalsePo... | TrueNe... | FalseN... | Recall | Precision | Sensitivity | Specificity | F-meas... | Accuracy | Cohen'... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| satisfied | 10719 | 694 | 13879 | 684 | 0.94 | 0.939 | 0.94 | 0.952 | 0.94 | ? | ? |
| neutral or dis... | 13879 | 684 | 10719 | 694 | 0.952 | 0.953 | 0.952 | 0.94 | 0.953 | ? | ? |
| Overall | ? | ? | ? | ? | ? | ? | ? | ? | ? | 0.947 | 0.892 |

## Attempt 4:

### Data Preprocessing

**Low Variance Filter-**In this attempt we applied low variance filter that ignores values that are below a certain range, however in our data set there were no columns that can be filtered out using low variance filter. hence , nodes' configurations were not changed a bit. nodes' configurations were not changed a bit.





### Model Building
Since the low variance filter did not affect any of the columns therefore the accuracy was also not affected. Accuracy is the same as in attempt 3 for all the classification models.

**Attempt 5:**

**Data Preprocessing:**

1. **PCA:**
   Keeping the same configuration as Attempt 3, we added PCA to data preprocessing. Four numerical columns were reduced to two.



**Model Building**

The configuration of all the classification models is kept the same as attempt 2 and 3. We found that the accuracy of each model declined after using PCA therefore we did not use it in the next attempts.

## Attempt 6:

### Data Preprocessing

#### 1. Category to Number:

In this attempt we converted all the string columns to numerical values using category to number converter as the model used required all the values to be numerical. All the other filters used above were used with the same configuration.

## Model Building:

### K Nearest Neighbor
For this attempt we used KNN model, with the following configuration.





### Scorer
Using KNN the accuracy of our model decline from approximately 95% to 80%.



| satisfactio... | satisfied | neutral or ... |
|---|---|---|
| satisfied | 7482 | 3921 |
| neutral or di... | 2881 | 11692 |

Correct classified: 19,174    Wrong classified: 6,802

Accuracy: 73.814 %    Error: 26.186 %

Cohen's kappa (κ) 0.463

Accuracy statistics - 4:92 - Scorer

Table "default" - Rows: 3   Spec - Columns: 11   Properties   Flow Variables

| Row ID | TruePo... | FalsePo... | TrueNe... | FalseN... | Recall | Precision | Sensitivity | Specificity | F-meas... | Accuracy | Cohen'... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| satisfied | 7482 | 2881 | 11692 | 3921 | 0.656 | 0.722 | 0.656 | 0.802 | 0.687 | ? | ? |
| neutral or dis... | 11692 | 3921 | 7482 | 2881 | 0.802 | 0.749 | 0.802 | 0.656 | 0.775 | ? | ? |
| Overall | ? | ? | ? | ? | ? | ? | ? | ? | ? | 0.738 | 0.463 |

**ROC:**

## Attempt 7:

### Changing configurations in Decision Tree:
**7.1-** For this attempt, we changed the minimum number of records per node from 5 to 10 and Root split the columns from Customer Type (Loyal or Disloyal Customers). We tried root split using other attributes as well but we got the highest accuracy using Customer Type.



### ROC And Scorer:

**7.2**- Changing the Quality measure from gain ratio to gini index reduced our accuracy.



**ROC And Scorer:**

**7.3-** Changed the Quality Measure back to Gain ratio with 25 as the minimum number of records per node, force splitting the column at customer type and binary nominal splits.



**ROC and Scorer:**

**For attempts 8-10 we did not made changes in our data preprocessing but we made changes in our model building to check if it had any effect on our accuracy or not.**
**Attempt 8:**

**Changing configurations in Random Forest:**

**8.1-** I limited the number of levels (tree depth) to 10.



**ROC And Scorer:**

**8.2-** Reset the changes made in 8.1. Then tried excluding on-board service and checkin service. Just to see if it will affect the results.



**ROC And Scorer:**

**8.3-** Reset the changes made in 8.2. Increased the number of models from 100 to 2000.



**ROC And Scorer:**

## Attempt 9:

### Changing configurations in Naive Bayes Learner:

**9.1-** I changed the maximum number of unique nominal values per attribute from 20 to 10, to check whether it affects the accuracy or not.



### ROC And Scorer:

**9.2-** Reset the changes made in 9.1, realized if we go below 6 in the maximum number of unique nominal values per attribute the accuracy drops sharply.



### ROC And Scorer:

**9.3:** Changed the probability to 0.0003, minimum S.D to 0.0002 and increased the nominal values per attribute to 25. Even if we increased it more the accuracy did not change.



### ROC And Scorer:

## Attempt 10:

## Changing configurations in Gradient Boosting Learner:

**10.1:** Unticked the option of limiting the number of levels(tree depth)



## ROC And Scorer:

**10.2:** Increased the number of models to 500.



## ROC And Scorer:

**10.3:** Limited the number of levels to 15 with 200 number of models and a learning rate of 0.15



**ROC And Scorer:**

# Findings

### Finding Best Algorithm

We used various classification model for airlines to identify critical bottleneck to raise passenger satisfaction. As seen in the bar chart below, random forest gave us the best accuracy in attempt 2.



### Factors Effecting Passenger Satisfaction

To figure out which features effect our accuracy the most, we took our best model and tried excluding certain features using column filter.

From some of the simulations, it was seen that airlines should focus on improving the Services provided in the flight. For instance, airlines could develop better software to allow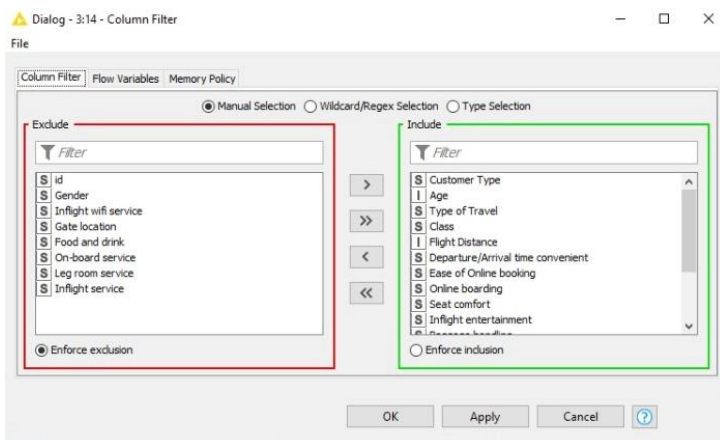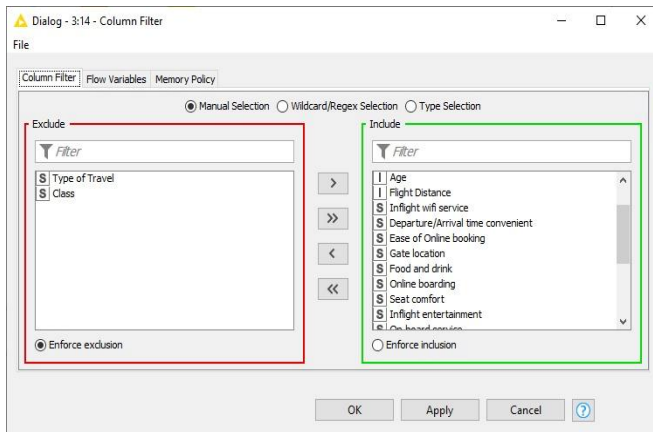 easier access to inflight wi-fi, food & etc. Excluding the "class" feature we observed how it effected the overall satisfaction of the passenger. Business class passenger are more satisfied with the flight so the airlines should initiate cheaper packages for business class.

**Concluding remarks**

The dramatic drop in demand for passenger air transport due to the COVID-19 pandemic has caused tremendous decline in the viability in the aviation industry. As a result, in order to resurrect the industry in the face of the crisis, it is critical to understand client pain points and increase their satisfaction with the services supplied. The industry should involve a third party market research company to provide them with unique, fresh, and unbiased recommendations to continuously improve their services attracting more people to choose their airlines for their trip.

There may be various options for gathering client feedback. The most straightforward and traditional method is to use the customer feedback form that is provided throughout the trip. However, the majority of passengers are uninterested in completing feedback forms. Other methods for collecting client input include the airline's internet website or online mobile applications. After the travel, the passenger can be sent an email with a link asking for feedback.

Finally, we hope that the model will provide a reference and be utilized for business value to the airlines, and will help them improve their services to allure more passengers towards them.