# Project Technical Document

## INVISTICO AIRLINE CLASSIFICATION

**Contents:**

## Classification of Invistico Airline Data

## Project description:

The aim of this data science project is to analyze the customer satisfaction levels of invistico airline based on the available data. By leveraging data analysis techniques, statistical modeling and machine learning algorithms, we will extract valuable insights to determine whether the customer is satisfied or dissatisfied with the airline's service.

This project will involve data preprocessing, exploratory data analysis, feature engineering, model building, and evaluation to provide actionable recommendation for improving customer satisfaction.

## Project Technical Details:

The following diagram shows the various steps that we have followed in our project.
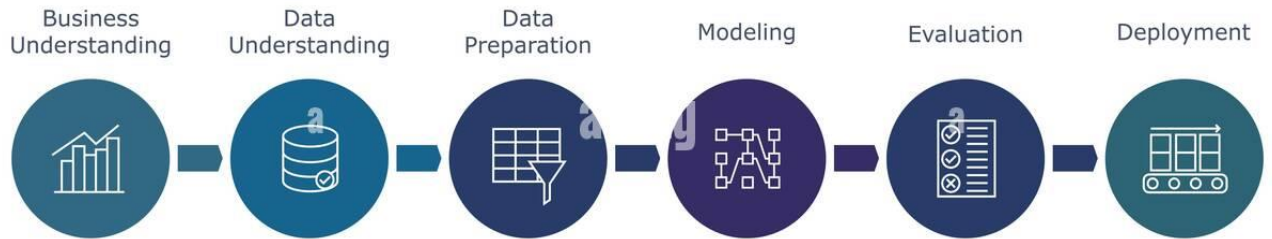
fig 1: General steps of CRISP DM process

# 1. Data collection:

Below are the raw data collected from Invistico airline excel data which includes various feature and observation. This dataset contain both categorical and numerical data.



1. There are 24 features available in the dataset and 129880 Observations.
2. Checked the null values or missing values in the dataset. Performed data imputation, filled null values with median and mode.
3. The last column named as "@" present at the last column having no data and that was irrelevant to the dataset has been removed.
4. There are 2 features having less then 30% important importance in the whole dataset, therefore it is removed.
   Departure/Arrival time convenient, Gate location
5. Considered the main 18 features such as
   Satisfaction, Gender, Customer Type, Age, Type of Travel, Class, Flight Distance,

Departure/Arrival time convenient, Food and Drink, Inflight wifi service, Online support, Ease of Online booking, On-board service, Leg room service, Baggage handling, checkin service, cleanliness, Online onboarding, Departure delay in minutes, Arrival delay in minutes.

6. The column contains 393 null values therefore we replaced that value with the help of zero.
7. Label encoding is applied to 5 columns having string values and it is important for the model.
   Those columns were gender, customer type, age, type of travel and class.
8. From the above dataset, it is clear that our target is customer satisfaction.
9. As the data present inside the dataset which is not scalable therefore it is normalized with normalization function.
10. Outliers present in the data ,explored using standard deviation,Inter Quartile Range(IQR), Median Absolute deviation(MAD), Isolation Forest, Winsorization with 95th Percentile.
11.  Separated continuous data and checked the correlation between the different features.
12. Checking the normalcy of data by using the visualization plots for all the columns.

## 2. Exploratory Data Analysis:

Analysis done on the basis of several charts which are shown below:



Fig.1 – customer satisfaction graph (left – dissatisfied, right- satisfied)

Checking data for each column is normally distributed.

## Age



## Flight Distance



## Seat comfort

## Departure/Arrival Time



## Food and Drink



## Gate location

# Wi-fi service



# Inflight Entertainment



# Online Support

# Ease of online booking



# On board service



# Leg room service

# Baggage room



# Check-in service



# Cleanliness

## Online onboarding



## Departure Delays



## Arrival Delays

Bar plot for all categorical variables in the dataset



**Bar plot for categorical values showing quantity**

**Inferences from the above charts:**

- Both genders are present approx. equal.
- Loyal customers are greater than disloyal customers.
- Airlines having customer who travel for business as compared to personal travel.
- Compared classes that customers have chosen from i.e. business, economy and economy plus.

**Correlation among all the columns.**

## 3. Feature Engineering

Following actions were performed in the feature engineering:

- Dropped features based on hypothesis testing (chi-square test)
- Selected features are:

    Satisfaction, Gender, Customer Type, Age, Type of Travel, Class, Flight Distance, Departure/Arrival time convenient, Food and Drink, Inflight wifi service, Online support, Ease of Online booking, On-board service, Leg room service, Baggage handling, checkin service, cleanliness, Online onboarding, Departure delay in minutes, Arrival delay in minutes.

- The 4 features have been dropped on Hypothesis Testing which are 'Departure/Arrival time convenient', 'Gate location'.

*Note: Hypothesis testing is performed on survey data. features are selected on the basis of sampled data. Domain knowledge has not been taken into account.*

## 4. Model building:

**Models accuracy:**

- Logistic - 76% Approx.

- Decision Tree - 54% Approx.

- KNN - 79% Approx.

- SVM - 50% Approx.

- XGB - 97% Approx.

*Selected model: XGB Classifier*

Various machine learning algorithms were explored  such as Logistic, Decision Tree, KNN, SVM, Logistic Regression, Decision Tree, ensemble technique – XGB Bagging etc. but the model that gave the highest accuracy is XGB Bagging. As compared to other models, XGB Bagging giving the best model accuracy.

**Main.py**

```python
Import pandas as pd

Import numpy as pd

from sklearn.feature_selection import SelectKBest

from sklearn.feature_selection import chi2

bestfeatures=SelectKBest(score_func=chi2,k=4)


x = data1_norm.iloc[:,1:]  #independent columns

y = data1_norm[0]  #target column is species


#apply SelectKBest class to extract top 10 best features

fit=bestfeatures.fit(x,y)

dfscores=pd.DataFrame(fit.scores_)

dfcolumns=pd.DataFrame(x.columns)


#concat two dataframes for better visualization

featureScore=pd.concat([dfcolumns,dfscores],axis=1)

featureScore.columns=['Specs','Score']

featureScore

import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.metrics import mean_squared_error

# importing machine learning models for prediction

import xgboost as xgb

# importing bagging module

from sklearn.ensemble import BaggingClassifier

# loading train data set in dataframe from train_data.csv file

from sklearn import preprocessing

model = BaggingClassifier(base_estimator=xgb.XGBClassifier())

model.fit(x_train,y_train)

pred_train=model.predict(x_train)

print(mean_squared_error(y_train, pred_train))
```

**Output:**

**Front-end**

## INVESTICO AIRLINES

| Gender : | Customer Type : |
|---|---|
| Male | Loyal Customer |

| Type of Travel : | Class : |
|---|---|
| Personal Travel | Eco |

| Age : | Flight Distance : |
|---|---|
| Enter the Age | Enter the Flight Distance |

Enter the ratings for below services ranges from 0 to 5 :

| Seat Comfort : | Departure/Arrival Time Convenient: |
|---|---|
| 0 - 5 | 0 - 5 |

| Food & Drink : | Gate Location : |
|---|---|
| 0 - 5 | 0 - 5 |

| Inflight Entertainment : | Online Support : |
|---|---|
| 0 - 5 | 0 - 5 |

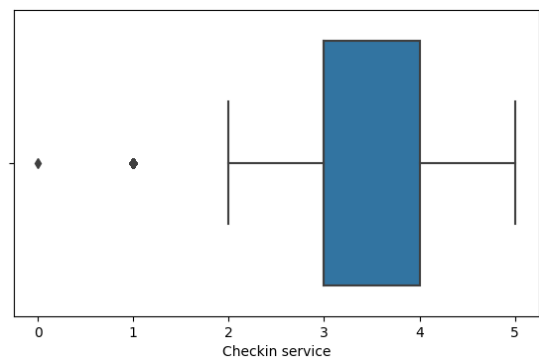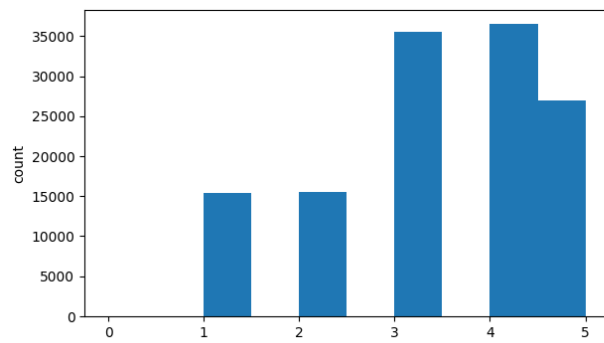| Ease Of Online Booking : | Onboard Services : |
|---|---|
| 0 - 5 | 0 - 5 |

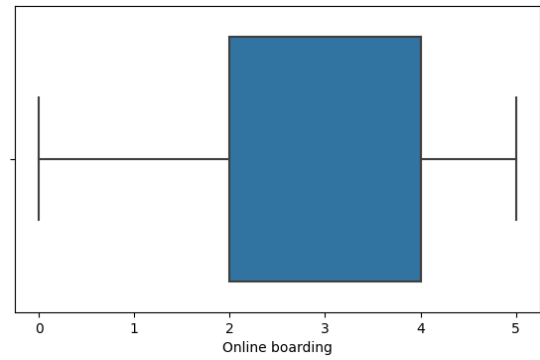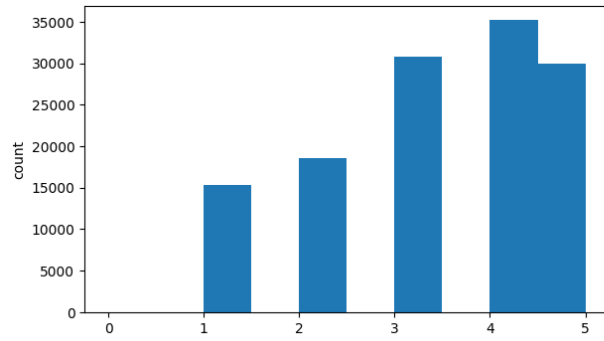| Leg Room Service : | Baggage Handling : |
|---|---|
| 0 - 5 | 0 - 5 |

| Check-in Service: | Cleanliness : |
|---|---|
| 0 - 5 | 0 - 5 |

| Online Boarding : | Department Delaying in minutes : |
|---|---|
| 0 - 5 | |

| Arrival Delaying in minutes : | Inflight Wifi Services : |
|---|---|
| | 0 - 5 |

**Satisfied or Not**

**Result:**

## INVESTICO AIRLINES

| | |
|---|---|
| Gender : | Customer Type : |
| Male | Disloyal Customer |
| Type of Travel : | Class : |
| Business Travel | Eco |
| Age : | Flight Distance : |
| 65 | 2464 |

Enter the ratings for below services ranges from 0 to 5 :

| | |
|---|---|
| Seat Comfort : | Departure/Arrival Time Convenient: |
| 2 | 2 |
| Food & Drink : | Gate Location : |
| 5 | 0 |
| Inflight Entertainment : | Online Support : |
| 0 | 2 |
| Ease Of Online Booking : | Onboard Services : |
| 2 | 3 |
| Leg Room Service : | Baggage Handling : |
| 3 | 3 |
| Check-in Service: | Cleanliness : |
| 2 | 3 |
| Online Boarding : | Department Delaying in minutes : |
| 2 | 310 |
| Arrival Delaying in minutes : | Inflight Wifi Services : |
| 305 | 0 |

Satisfied or Not

### Prediction : SATISFIED

**Invistico Airline uses the following packages and library from python:**

```python
import pandas as pd

import numpy as np

from sklearn import preprocessing

import matplotlib.pyplot as plt

import seaborn as sns

import seaborn as sns

import matplotlib.pyplot as plt

from scipy.stats import skew

from numpy import asarray

from sklearn.preprocessing import MinMaxScaler

import warnings

from sklearn.feature_selection import SelectKBest

from sklearn.feature_selection import chi2

from sklearn.model_selection import train_test_split

from sklearn.preprocessing import StandardScaler

import pandas as pd

from sklearn.model_selection import train_test_split

from sklearn.metrics import mean_squared_error

import xgboost as xgb
```

```python
from sklearn.ensemble import BaggingClassifier

from sklearn import preprocessing
```

## 5. Deployment using Flask:

Deployment process was done using flask technique.

```python
from flask import Flask, render_template, request
import pandas as pd
import pickle

app=Flask(__name__)
data=pd.read_excel('C:\VSCODE\ML TRAINING
PROJECT\Invistico_Airline.xlsx')
pklfile=pickle.load(open('C:\VSCODE\ML TRAINING
PROJECT\mainy.pkl','rb'))
@app.route('/')

def index():
    return render_template('index.html')
# Age   Type of Travel Class  Flight Distance    Seat comfort
Departure/Arrival time convenient  Food and drink Gate location
Inflight wifi service  Inflight entertainment Online support
Ease of Online booking On-board service   Leg room service
Baggage handling   Checkin service    Cleanliness    Online
boarding    Departure Delay in Minutes Arrival Delay in Minutes

@app.route('/predict',methods=['POST'])
def predict():
    Gender=request.form.get('gender')
    Customer_Type=request.form.get('customer-type')
    Age=request.form.get('age')
    Type_of_travel=request.form.get('type-of-travel')
    Class=request.form.get('class')
    Flight_Distance=request.form.get('flight-distance')
    Seat_comfort=request.form.get('seat-comfort')
    Departure_Arrival_Time_Convenient=request.form.get('datc')
    Food_Drink=request.form.get('food-drink')
    Gate_Location=request.form.get('gate-location')
    Inflight_Wifi_services=request.form.get('Iwifi')
    Inflight_entertainment=request.form.get('inflight-ent')
    Online_support=request.form.get('online-support')
    Ease_of_online_booking=request.form.get('eoob')
    On_board_services=request.form.get('onboard-services')
    Leg_room_services=request.form.get('lrs')
    Baggage_handling=request.form.get('bh')
    Checkin_services=request.form.get('cis')
    Cleanliness=request.form.get('clean')
    Online_boarding=request.form.get('online-boarding')
    Departure_delaying=request.form.get('DD-in-min')
```

```python
        Arrival_delaying=request.form.get('AA-in-min')

    if(Gender=='Male'):
        Gender=0
    else :
        Gender=1

    if(Customer_Type=='Loyal Customer'):
        Customer_Type=0
    else :
        Customer_Type=1

    if(Type_of_travel=='Personal Travel'):
        Type_of_travel=0
    else :
        Type_of_travel=1

    if(Class=='Eco'):
        Class=1
    elif(Class=='Business'):
        Class=0
    else :
        Class=2


input=pd.DataFrame([[Gender,Customer_Type,Age,Type_of_travel,Cla
ss,Flight_Distance,Seat_comfort,Departure_Arrival_Time_Convenien
t,Food_Drink,Gate_Location,Inflight_Wifi_services,Inflight_enter
tainment,Online_support,Ease_of_online_booking,On_board_services
,Leg_room_services,Baggage_handling,Checkin_services,Cleanliness
,Online_boarding,Departure_delaying,Arrival_delaying]],columns=[
'Gender','Customer Type','Age','Type of Travel','Class','Flight
Distance','Seat comfort','Departure/Arrival time
convenient','Food and Drink','Gate Location','Inflight wifi
service','Inflight entertainment','Online support','Ease of
Online booking','On-board service','Leg room service','Baggage
handling','Checkin service','Cleanliness','Online
boarding','Departure Delay in minutes','Arrival Delay in
minutes'])
    prediction=pklfile.predict(input)
    print(prediction)
    if(prediction==1):
        return "SATISFIED"
else :
        return "DISSATISFIED"

if __name__=="__main__":
```

```
    app.run(debug=True, port=5001)
```