Airline Passenger Satisfaction Predictor

Objective

- Determine the whether customers will be satisfied with their airline based on factors given.
- Determine correlation between satisfaction factors and/or onboard/offboard services.
- Business applicable.

Dataset Briefing

Customer satisfaction correlation to many other areas of off-flight and inflight experience

Categorical Data: Gender, Customer, Type of Travel, Class, Satisfaction

Numerical Data: Age, Flight Distance, Departure Delay in Minutes, Arrival Delay in Minutes

Numerical Data (Scale of 1-5): Inflight Wifi Service, Departure/Arrival Time Convenient, Ease of Online Booking, Gate Location, Food and Drink, Online Boarding, Seat Comfort, Inflight Entertainment, On-board service, Leg Room Service, Baggage Handling, Check-in Service, Inflight Service, Cleanliness

1 5047 Male disloyal Customer 25 Business Business 235 3 2 3 1 1 5 3 1 4 1 1 6.0 ne dissection of travel Business 1142 2 2 2 2 5 4 3 4 4 5 0 0.0 s 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1																					
2 110028 Female Loyal Customer Customer 26 Business travel 1142 2 2 2 2 2 3 4 3 4 4 4 5 0 0.0 s 3 24026 Female Loyal Customer 25 Business travel Business biranel 562 2 5 5 2 2 2 5 3 1 4 2 11 9.0 need diss 4 119299 Male Loyal Customer 61 Business travel 214 3 3 3 3 3 4 4 3 3 3 0 0.0 s 99 94171 Female disloyal Customer 23 Business travel Eco 192 2 1 2 2 3 1 4 2 3 2 3 0.0 one 99 94171 Female Loyal Customer 49 Business travel Eco 192 2 1 2 2 3 1	0	70172	Male	Loyal Customer	13	Personal Travel	Eco Plus	460		4			4			4			25	18.0	neutral or dissatisfied
3 24026 Female Loyal Customer Customer 25 Business travel 562 2 5 5 2 2 5 3 1 4 2 11 9.0 ned diss 4 119299 Male Loyal Customer 61 Business travel 214 3 3 3 3 4 4 3 3 3 0 0.0 s	1	5047	Male	disloyal Customer	25	Business travel	Business	235	3		3				3		4			6.0	neutral or dissatisfied
4 119299 Male Loyal Customer Customer 61 Business travel 214 3 3 3 3 4 4 3 3 3 0 0.0 s	2	110028	Female	Loyal Customer	26	Business travel	Business	1142	2	2	2		4		4	4	4			0.0	satisfied
<td>3</td> <td>24026</td> <td>Female</td> <td>Loyal Customer</td> <td>25</td> <td>Business travel</td> <td>Business</td> <td>562</td> <td>2</td> <td></td> <td></td> <td>2</td> <td>2</td> <td></td> <td></td> <td></td> <td>4</td> <td>2</td> <td>11</td> <td>9.0</td> <td>neutral or dissatisfied</td>	3	24026	Female	Loyal Customer	25	Business travel	Business	562	2			2	2				4	2	11	9.0	neutral or dissatisfied
00 73097 Male Loyal 49 Business Business 2347 4 4 4 5 5 5 5 5 4 0 0.0 s	4	119299	Male	Loyal Customer	61	Business travel	Business	214		3	3			4	4	3				0.0	satisfied
00 73097 Male Loyal 49 Business Business 2347 4 4 4 5 5 5 5 5 4 0 0.0 s																					
	99	94171	Female	disloyal Customer	23	Business travel	Eco	192			2	2			4			2	3	0.0	neutral or dissatisfied
01 68825 Male disloyal 30 Business Business 1995 1 1 1 1 4 3 2 4 5 5 4 7 14.0 nd diss	00	73097	Male	Loyal Customer	49	Business travel	Business	2347	4	4	4	5						4		0.0	satisfied
	01	68825	Male	disloyal Customer	30	Business travel	Business	1995				4		2	4			4		14.0	neutral or dissatisfied

Ease of Inflight
Online ... entertainment
booking

On- Leg baggage Checkin Inflight Cleanliness service service

Inflight wifi Departure/Arrival time convenient

Flight Distance

Class

Eco

Type of

22 Business

travel

27 Business Business

disloyal Customer

Loyal Customer

02 54173 Female

Male

03 62567

Travel

id Gender Customer Age

Arrival

Delay in

Minutes

satisfaction

neutral or

dissatisfied

neutral or

dissatisfied

Departure

Delay in

Minutes

Note on Datasets

```
df_train = pd.read_csv('./train.csv', index_col=[0])
df_test = pd.read_csv('./test.csv', index_col=[0])
```

Data Preprocessing

Null Values

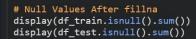
Standardization

Encoding

```
# Display the Null Values

display(df_train.isnull().sum())
display(df_test.isnull().sum())
display(df_train[pd.isnull(df_train).any(axis=1)])
```

0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
0	
310	
0	
	0 0 0 0 0 0 0 0 0 0 0 0 0



Gender	0
Customer_Type	0
Age	0
Type_of_Travel	0
Class	0
Flight_Distance	0
Inflight_wifi_service	0
Departure/Arrival_time_convenient	0
Ease_of_Online_booking	0
Gate_location	0
Food_and_drink	0
Online_boarding	0
Seat_comfort	0
Inflight_entertainment	0
On-board_service	0
Leg_room_service	0
Baggage_handling	0
Checkin_service	0
Inflight_service	0
Cleanliness	0
Departure_Delay_in_Minutes	0
Arrival_Delay_in_Minutes	0
satisfaction	0
dtype: int64	



Set the null values of Arrival Delay to Departure Delay since if one departs a certain time late, they will most likely a df_train['Arrival_Delay_in_Minutes'] = df_train['Arrival_Delay_in_Minutes'].fillna(df_train['Departure_Delay_in_Minutes']) df_test['Arrival_Delay_in_Minutes'] = df_test['Arrival_Delay_in_Minutes'].fillna(df_test['Departure_Delay_in_Minutes'])

```
# Data Standardization

scaler = StandardScaler()
numCols = df_train.select_dtypes(include=np.number).columns

df_train[numCols] = pd.DataFrame(scaler.fit_transform(df_train[numCols]), columns=df_train[numCols].columns)

df_test[numCols] = pd.DataFrame(scaler.fit_transform(df_test[numCols]), columns=df_test[numCols].columns)
```

Inflight_Entertainment	On- board_Service	Leg_Room_Service	Baggage_Handling	Checkin_Service	Inflight_Service	Cleanliness	Departure_Delay_in_Minutes	Arrival_Delay_in_Minutes
1.231704	0.479403	-0.266840	0.311769	0.549799	1.156436	1.305870	0.266393	0.070988
-1.769081	-1.849161	1.253380	-0.535045	-1.821012	0.305848	-1.742292	-0.361375	-0.238223
1.231704	0.479403	-0.266840	0.311769	0.549799	0.305848	1.305870	-0.387532	-0.392828
-1.018885	-1.072973	1.253380	-0.535045	-1.821012	0.305848	-0.980251	-0.099805	-0.160920
-0.268688	-0.296785	0.493270	0.311769	-0.240472	-0.544740	-0.218211	-0.387532	-0.392828
-1.018885	-0.296785	-1.787061	0.311769	-1.030742	-0.544740	-0.980251	-0.309061	-0.392828
1.231704	1.255590	1.253380	1.158582	1.340069	1.156436	0.543829	-0.387532	-0.392828
0.481508	-0.296785	-1.026951	0.311769	1.340069	1.156436	0.543829	-0.204433	-0.032082
-1.769081	0.479403	1.253380	-2.228672	1.340069	0.305848	-1.742292	-0.387532	-0.392828
-1.769081	-1.849161	-1.787061	0.311769	0.549799	-0.544740	-1.742292	-0.387532	-0.392828

```
# All the different values for the Categorical Columns

print(df_train['Gender'].unique())
print(df_train['Customer_Type'].unique())
print(df_train['Type_of_Travel'].unique())
print(df_train['Class'].unique())
```

['Male' 'Female']

['Loyal_Customer' 'Disloyal_Customer']
['Personal Travel' 'Business Travel']

['Eco Plus' 'Business' 'Eco']

```
categorical_columns = ['Gender', 'Customer_Type', 'Type_of_Travel']
ohe = OneHotEncoder(sparse output=False)
one hot encoded = ohe.fit_transform(df_train[categorical_columns])
one hot df = pd.DataFrame(one hot encoded, columns=ohe.get feature names out(categorical columns))
encoded_train = pd.concat([one_hot_df, df_train], axis=1)
encoded_train = encoded_train.drop(categorical_columns, axis=1)
display(encoded train)
categorical columns = ['Gender', 'Customer Type', 'Type of Travel']
ohe = OneHotEncoder(sparse_output=False)
one hot encoded = ohe.fit transform(df test[categorical columns])
one_hot_df = pd.DataFrame(one_hot_encoded, columns=ohe.get_feature_names_out(categorical_columns))
encoded_test = pd.concat([one_hot_df, df_test], axis=1)
encoded test = encoded test.drop(categorical columns, axis=1)
```

One Hot Encoding for Categorical Features

```
Satisfaction
# Ordinal Encoding for the Class where the order is Eco, Eco Plus, and then Business
                                                                                                        Class
enc = OrdinalEncoder(categories=[['Eco', 'Eco Plus', 'Business']])
                                                                                                          1.0
                                                                                                                            0.0
encoded train['Class'] = enc.fit transform(encoded train.loc[:,['Class']])
enc = OrdinalEncoder(categories=[['Neutral Or Dissatisfied', 'Satisfied']])
                                                                                                          2.0
encoded train['Satisfaction'] = enc.fit transform(encoded train.loc[:,['Satisfaction']])
                                                                                                                            0.0
                                                                                                          2.0
display(encoded train)
                                                                                                                            1.0
                                                                                                          2.0
enc = OrdinalEncoder(categories=[['Eco', 'Eco Plus', 'Business']])
                                                                                                                            0.0
encoded test['Class'] = enc.fit transform(encoded test.loc[:,['Class']])
                                                                                                          2.0
enc = OrdinalEncoder(categories=[['Neutral Or Dissatisfied', 'Satisfied']])
encoded test['Satisfaction'] = enc.fit transform(encoded test.loc[:,['Satisfaction']])
                                                                                                                            1.0
                                                                                                          0.0
                                                                                                          2.0
                                                                                                                            0.0
                                                                                                          2.0
                                                                                                                            1.0
                                                                                                          0.0
   # Standardized Encoded Data
                                                                                                                            0.0
   std encoded train = encoded train
                                                                                                          2.0
   std encoded train['Class'] = pd.Series(scaler.fit transform(encoded train[['Class']]).flatten())
   display(std encoded train)
                                                                                                                            0.0
   std encoded test = encoded test
   std encoded test['Class'] = pd.Series(scaler.fit transform(encoded test[['Class']]).flatten())
                                                                                                                            0.0
```

Exploratory Data Analysis

Satisfaction Percentage

```
fig, axes = plt.subplots(nrows=1, ncols=1)
axes.pie(df_train['Satisfaction'].value_counts(),
         labels=df_train['Satisfaction'].value_counts().index,
         explode=[0.1, 0],
         autopct='%1.2f%%',
         shadow=True,
         startangle=90,
         colors=["lightcoral", "lightgreen"])
axes.set_title('Satisfaction Percent')
legend_labels = [f'{label} ({count})' for label, count in zip(df_train['Satisfaction'].value_counts().index, df_train['Satisfaction'].value_counts())]
axes.legend(legend_labels, loc='upper left')
plt.show()
                                 Satisfaction Percent
                          Neutral_Or_Dissatisfied (58879)
                          Satisfied (45025)
                                                                 Satisfied
                                                    43.33%
                           56.67%
Neutral Or Dissatisfied
```

Age Distribution

```
sns.histplot(data=unstandardized_df_train,x='Age', color="lightskyblue")
plt.title("Age Distribution")
plt.show()
                              Age Distribution
  3000
  2500
  2000
1500 ·
  1000
   500
```

Age

Average Satisfaction on Services

```
on_board_services = ['Inflight_Wifi_Service', 'Food_and_Drink', 'Seat_Comfort', 'Inflight_Entertainment', 'On-board_Service','Leg Room_Service', 'Inflight_Service', 'Cleanliness']
average_on_ratings = unstandardized_df_train[on_board_services].mean()
average_on_ratings_df = pd.DataFrame({'On Board Service': average_on_ratings.index, 'Average Satisfaction Rating': average_on_ratings.values})
sns.barplot(x='On Board Service', y='Average Satisfaction Rating', data=average_on_ratings_df, palette='mako')
plt.xticks(rotation=45, ha='right')
plt.title('Average Satisfaction Ratings for On-board Services')
plt.show()
off_board_services = ['Departure/Arrival_Time_Convenient', 'Ease_of_Online_Booking', 'Gate_Location', 'Online_Boarding', 'Baggage_Handling', 'Checkin_Service']
average_off_ratings = unstandardized_df_train[off_board_services].mean()
average_off_ratings_df = pd.DataFrame({'Off_Board_Service': average_off_ratings.index, 'Average_Satisfaction_Rating': average_off_ratings.values})
sns.barplot(x='Off Board Service', y='Average Satisfaction Rating', data=average_off_ratings_df, palette='mako')
plt.xticks(rotation=45, ha='right')
plt.title('Average Satisfaction Ratings for Off-board Services')
plt.show()
                                                                                                    Average Satisfaction Ratings for Off-board Services
               Average Satisfaction Ratings for On-board Services
                                                                                           3.5
                                                                                           3.0
                                                                                         2.5
                                                                                         B 2.0
                                                                                            0.5
                                   On Board Service
                                                                                                                        Off Board Service
```

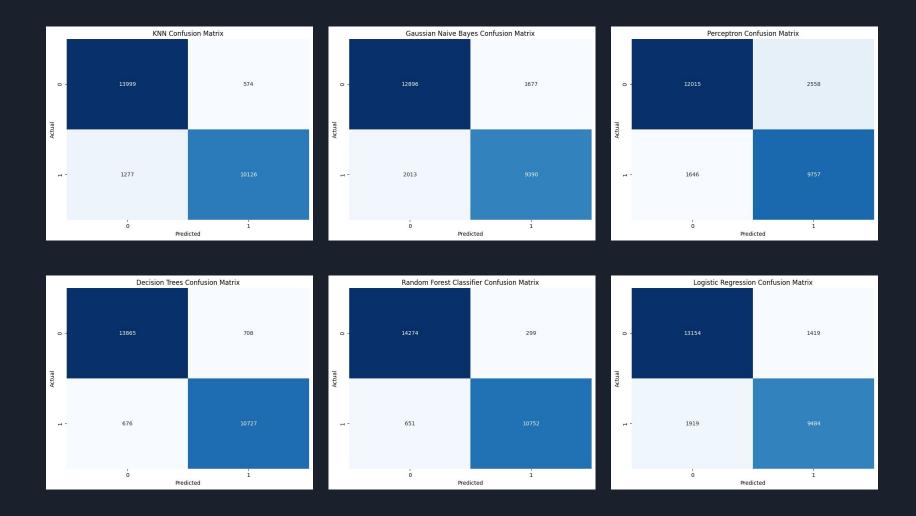
Correlation Matrix

```
plt.figure(figsize=(18,8))
corr = std_encoded_train.corr()
mask = np.triu(corr)
sns.heatmap(corr, annot=True, fmt='.2f', mask=mask, cmap="Blues")
plt.xticks(rotation=45, ha='right')
                  Gender Female -
                    Gender Male -- 1.00
 Customer Type Disloyal Customer - 0.03 -0
   Customer_Type_Loyal_Customer -- 0.03 0.03 -1.00
    Type of Travel Business Travel - 0.01 -0.01 0.31 -0.31
    Type of Travel Personal Travel -- 0.01 0.01 -0.31 0.31 -1.00
                                                                                                                                                                                               0.50
                                   0.01 0.01 -0.28 0.28 0.05 -0.0
                                   0.01 0.01 -0.11 0.11 0.55 -0.55 0.14
                                   0.01 0.01 0.23 0.23 0.27 0.27 0.10 0.45
                  Flight Distance -
              Inflight Wifi Service -
                                   0.01 0.01 0.01 0.01 0.11 0.11 0.02 0.04 0.0
                                                                                                                                                                                               - 0.25
                                   0.01 0.01 0.21 0.21 0.26 0.26 0.04 0.09 0.02 0.34
Departure/Arrival Time Convenient -
                                   0.01 0.01 -0.02 0.02 0.13 -0.13 0.02 0.11 0.07 0.72 0.44
          Ease of Online Booking -
                                   0.00 0.00 0.01 -0.01 0.03 -0.03 -0.00 0.00 0.00 0.34 0.44 0.46
                                                                                                                                                                                               - 0.00
                  Food and Drink -
                                   0.01 0.01 -0.06 0.06 0.06 -0.06 0.02 0.09 0.06 0.13 0.00 0.03 -0
                 Online Boarding - 0.04 -0.04 -0.19 0.19 0.22 -0.22 0.21 0.32 0.21 0.46 0.07 0.40 0.00 0.2
                                   0.03 -0.03 -0.16 0.16 0.12 -0.12 0.16 0.23 0.16 0.12 0.01 0.03 0.00 0.57 0.42
                                                                                                                                                                                                -0.25
                                   0.01 0.01 -0.11 0.11 0.15 -0.15 0.08 0.19 0.13 0.21 -0.00 0.05 0.00 0.62 0.29 0.61
                Baggage Handling
                                                                                                                                                                                               - -0.50
                  Inflight Service
                      Cleanliness
                                                                                                                                                                                               -0.75
      Departure_Delay_in_Minutes
         Arrival_Delay_in_Minutes
                                                                                                                                                                                              --1.00
```

Model Building Techniques Used

- K-Nearest Neighbors
- Gaussian Naive Bayes
- Perceptron
- Logistic Regression
- Decision Trees
- Random Forest Classifier

```
from sklearn.metrics import accuracy score, confusion matrix, recall score, precision score, f1 score
results = pd.DataFrame(columns=["Classifier", "Accuracy", "Recall", "Precision", "f1"])
for name, classifier in classifiers:
    Y pred = classifier.predict(X test)
    accuracy = accuracy score(Y test, Y pred)
   recall = recall score(Y test, Y pred)
    precision = precision score(Y test, Y pred)
    f1 = f1 score(Y test, Y pred)
    print("Name: ", name)
    data = pd.DataFrame([{"Classifier": name, "Accuracy": accuracy, "Recall": recall, "Precision": precision, "f1": f1}])
    confusion = confusion matrix(Y test, Y pred)
    plt.figure(figsize=(8, 6))
    sns.heatmap(confusion, annot=True, fmt="d", cmap="Blues", cbar=False)
    plt.xlabel('Predicted')
    plt.ylabel('Actual')
   plt.title(name + ' Confusion Matrix')
    plt.show()
    results = pd.concat([results, data], ignore index=True)
```



Model Comparison

```
styled_results = results.style.background_gradient(cmap='Blues', subset=["Accuracy", "Recall", "Precision", "f1"]) display(styled_results)
```

	Classifier	Accuracy	Recall	Precision	f1
0	KNN	0.928742	0.888012	0.946355	0.916256
1	Gaussian Naive Bayes	0.857946	0.823468	0.848468	0.835781
2	Perceptron	0.838158	0.855652	0.792286	0.822751
3	Logistic Regression	0.871497	0.831711	0.869852	0.850354
4	Decision Trees	0.947298	0.940805	0.939240	0.940022
5	Random Forest Classifier	0.964044	0.942559	0.974698	0.958359