

```
In [1]: import pandas as pd
import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt
sns.set_theme(color_codes=True)
pd.set_option('display.max_columns', None)
```

```
In [2]: df = pd.read_csv('Invistico_Airline.csv')
df.head()
```

Out[2]:

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	Bahá
0	0	0	2	2	4	2	3	3	0	
0	0	0	3	0	2	2	3	4	4	
0	0	0	3	2	0	2	2	3	3	
0	0	0	3	3	4	3	1	1	0	
0	0	0	3	4	3	4	2	2	0	

Data Preprocessing Part 1

```
In [3]: #Check the number of unique value from all of the object datatype
df.select_dtypes(include='object').nunique()
```

```
Out[3]: satisfaction      2
Customer Type      2
Type of Travel     2
Class              3
dtype: int64
```

Exploratory Data Analysis

```
In [4]: # Get the names of all columns with data type 'object' (categorical columns) excluding 'satisfaction'
cat_vars = df.select_dtypes(include='object').columns.tolist()

# Create a figure with subplots
num_cols = len(cat_vars)
num_rows = (num_cols + 2) // 3
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))
axs = axs.flatten()
```

```
In [4]: # Get the names of all columns with data type 'object' (categorical columns) excluding 'churn'
cat_vars = df.select_dtypes(include='object').columns.tolist()

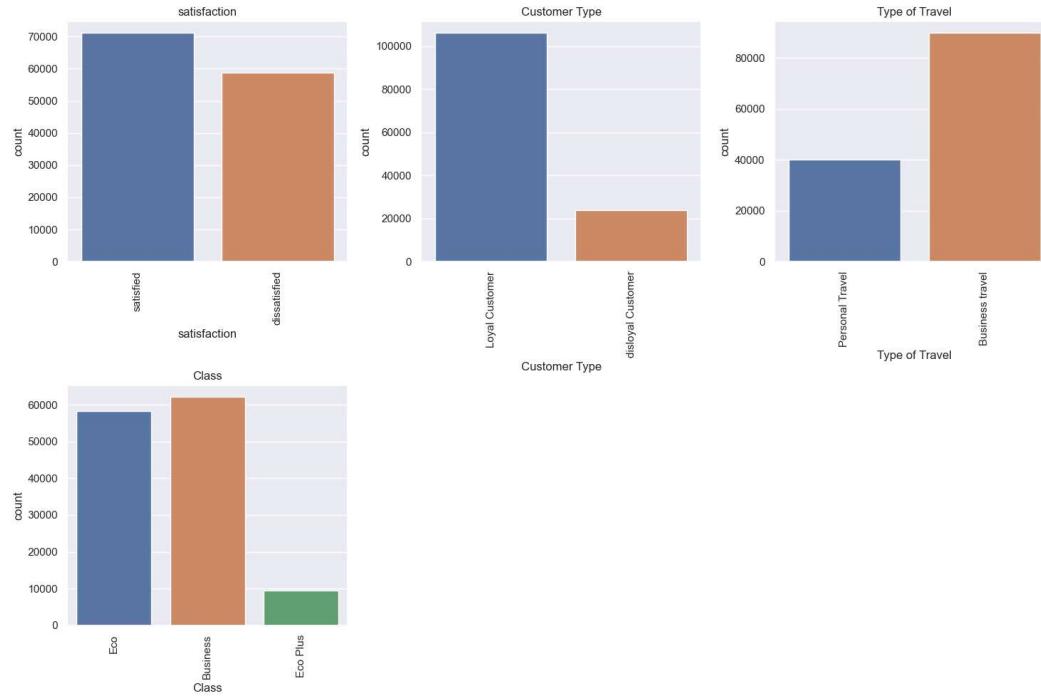
# Create a figure with subplots
num_cols = len(cat_vars)
num_rows = (num_cols + 2) // 3
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))
axs = axs.flatten()

# Create a countplot for the top 6 values of each categorical variable using Seaborn
for i, var in enumerate(cat_vars):
    top_values = df[var].value_counts().index
    filtered_df = df[df[var].isin(top_values)]
    sns.countplot(x=var, data=filtered_df, ax=axs[i])
    axs[i].set_title(var)
    axs[i].tick_params(axis='x', rotation=90)

# Remove any extra empty subplots if needed
if num_cols < len(axs):
    for i in range(num_cols, len(axs)):
        fig.delaxes(axs[i])

# Adjust spacing between subplots
fig.tight_layout()

# Show plot
plt.show()
```



```
In [5]: # Get the names of all columns with data type 'int' or 'float', excluding 'churn'
num_vars = df.select_dtypes(include=['int', 'float']).columns.tolist()

# Create a figure with subplots
num_cols = len(num_vars)
num_rows = (num_cols + 2) // 3
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))
axs = axs.flatten()
```

```
In [5]: # Get the names of all columns with data type 'int' or 'float', excluding 'churn'
num_vars = df.select_dtypes(include=['int', 'float']).columns.tolist()

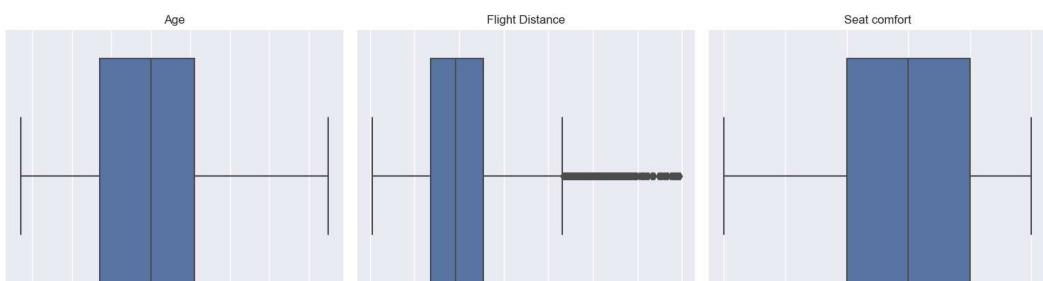
# Create a figure with subplots
num_cols = len(num_vars)
num_rows = (num_cols + 2) // 3
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))
axs = axs.flatten()

# Create a box plot for each numerical variable using Seaborn
for i, var in enumerate(num_vars):
    sns.boxplot(x=df[var], ax=axs[i])
    axs[i].set_title(var)

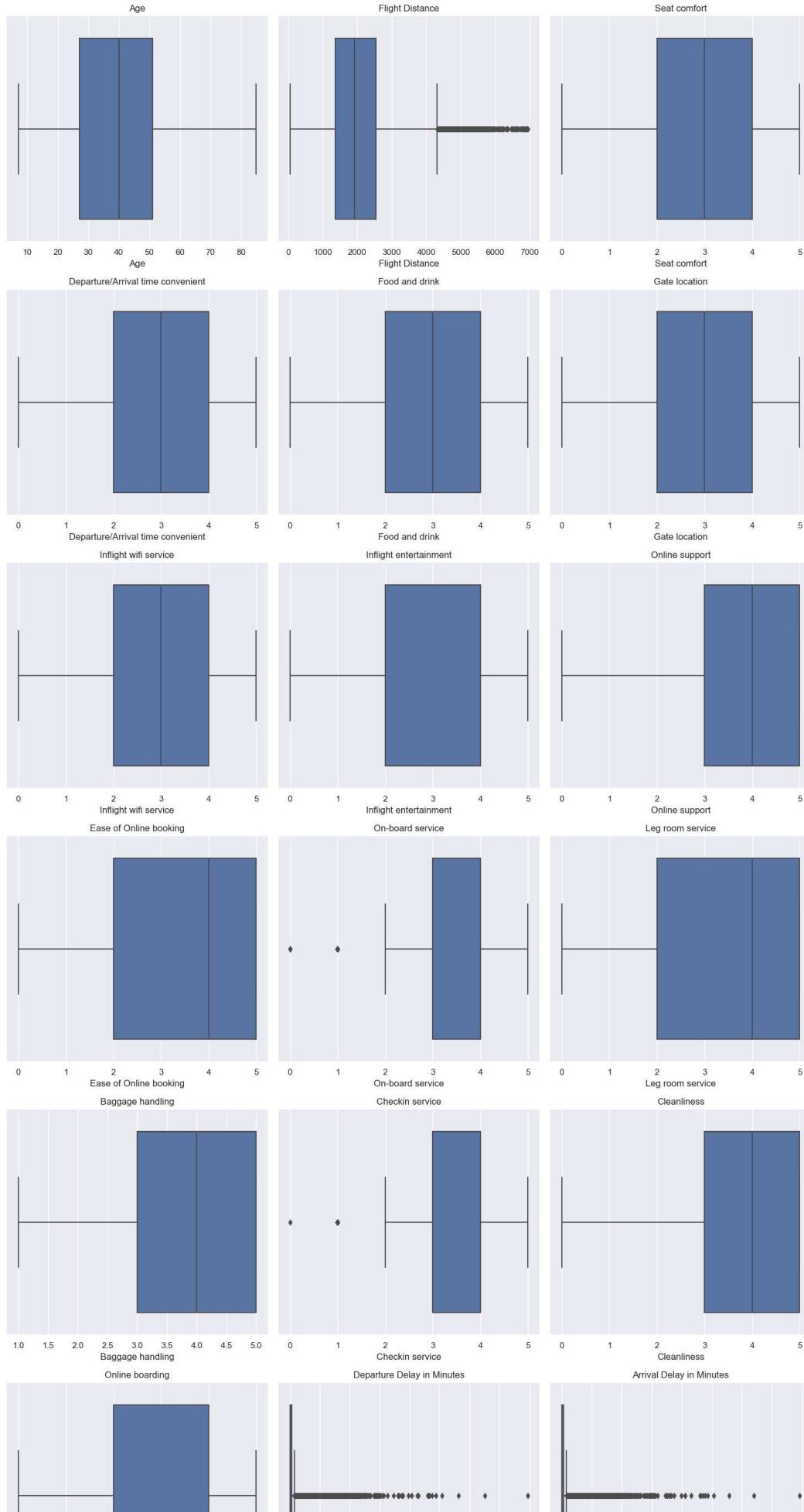
# Remove any extra empty subplots if needed
if num_cols < len(axs):
    for i in range(num_cols, len(axs)):
        fig.delaxes(axs[i])

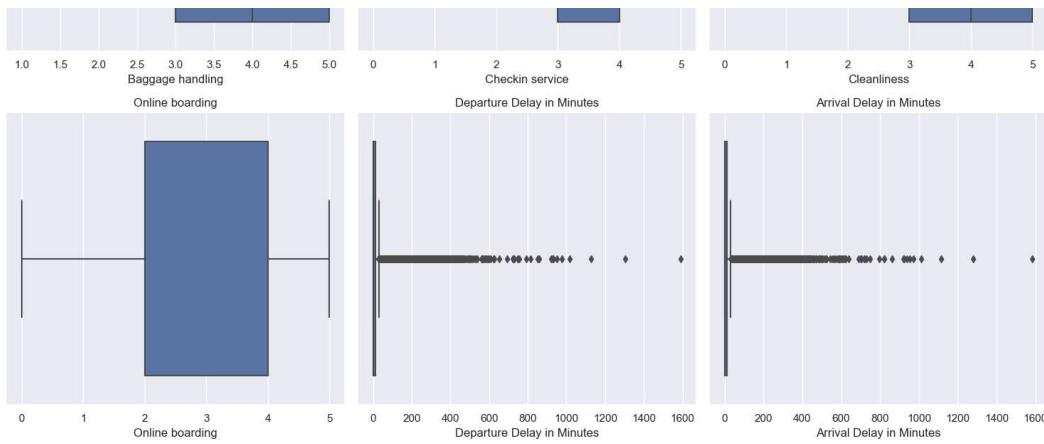
# Adjust spacing between subplots
fig.tight_layout()

# Show plot
plt.show()
```



Customer Satisfaction in Airline Prediction - Jupyter Notebook

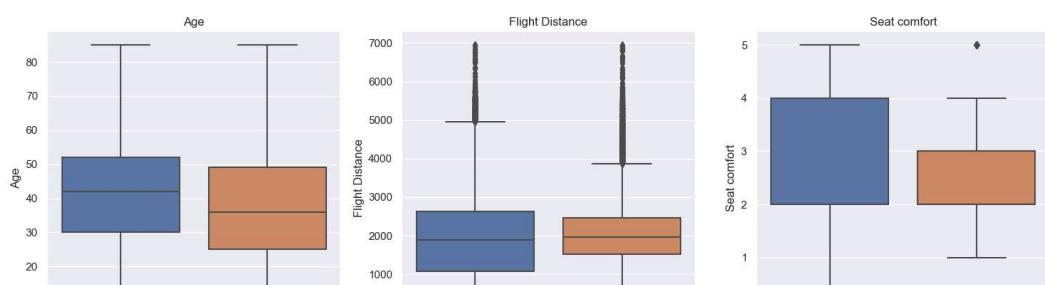




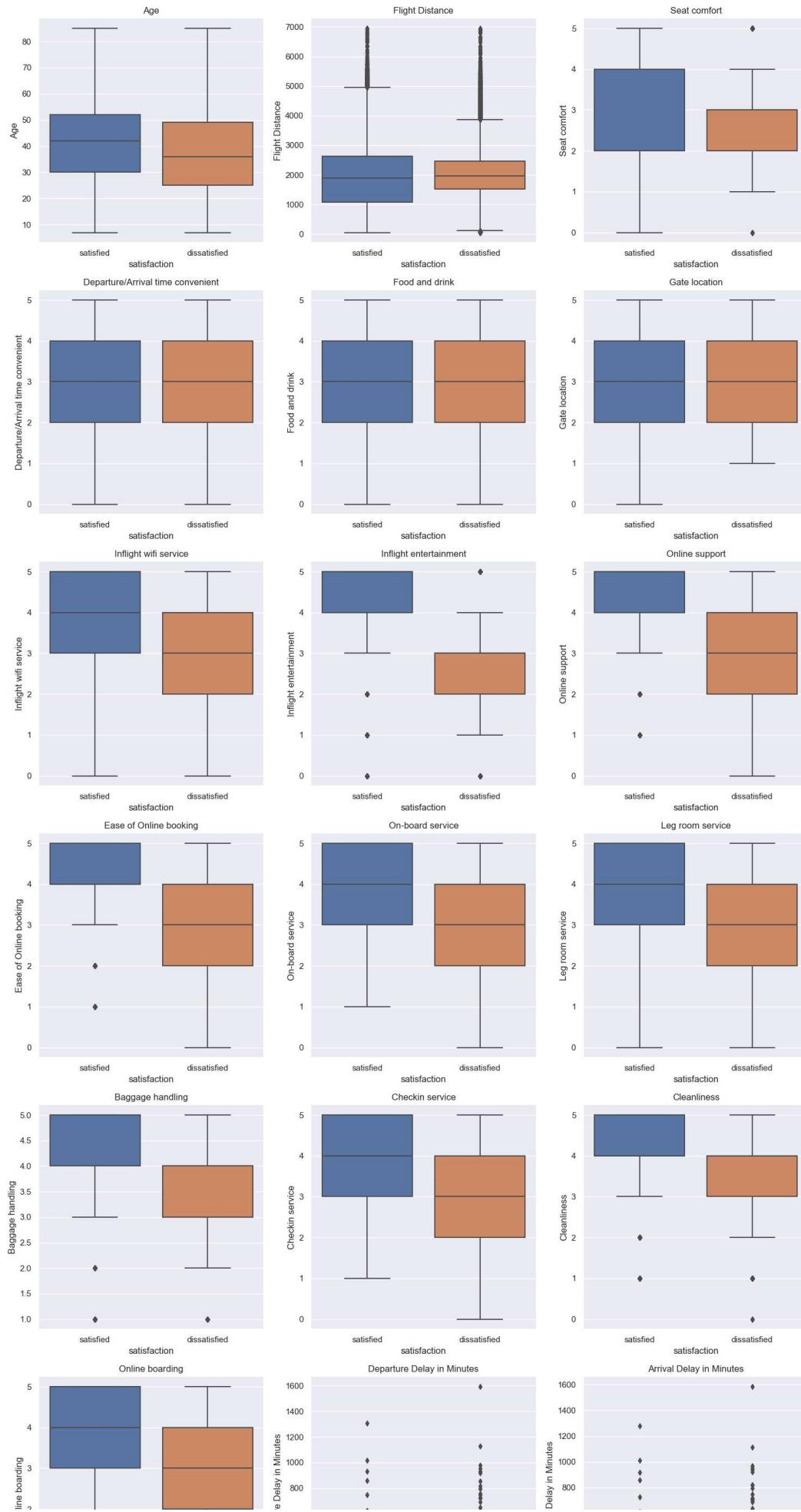
```
In [6]: # Get the names of all columns with data type 'int'
int_vars = df.select_dtypes(include=['int', 'float']).columns.tolist()
int_vars = [col for col in num_vars if col != 'satisfaction']

# Create a figure with subplots
num_cols = len(int_vars)
num_rows = (num_cols + 2) // 3 # To make sure there are enough rows for the subplots
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))
axs = axs.flatten()
```

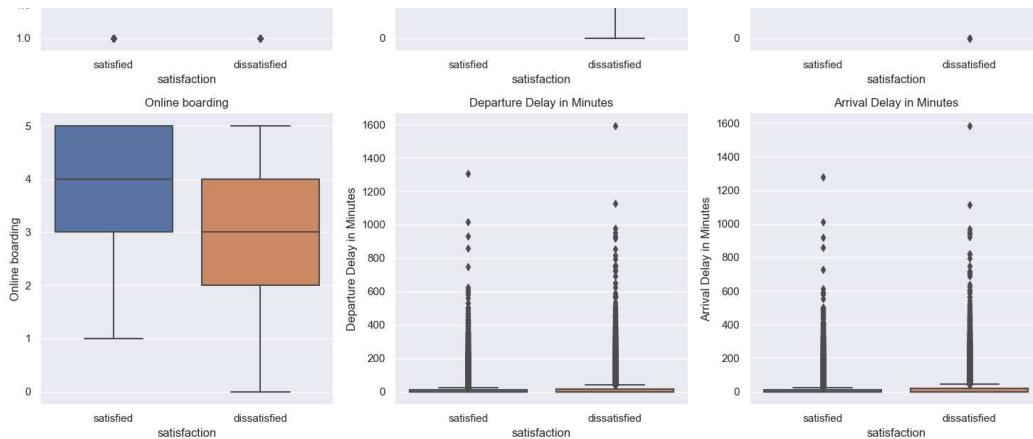
```
In [6]: # Get the names of all columns with data type 'int'  
int_vars = df.select_dtypes(include=['int', 'float']).columns.tolist()  
int_vars = [col for col in num_vars if col != 'satisfaction']  
  
# Create a figure with subplots  
num_cols = len(int_vars)  
num_rows = (num_cols + 2) // 3 # To make sure there are enough rows for the subplots  
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))  
axs = axs.flatten()  
  
# Create a box plot for each integer variable using Seaborn with hue='Attrition'  
for i, var in enumerate(int_vars):  
    sns.boxplot(y=var, x='satisfaction', data=df, ax=axs[i])  
    axs[i].set_title(var)  
  
# Remove any extra empty subplots if needed  
if num_cols < len(axs):  
    for i in range(num_cols, len(axs)):  
        fig.delaxes(axs[i])  
  
# Adjust spacing between subplots  
fig.tight_layout()  
  
# Show plot  
plt.show()
```



Customer Satisfaction in Airline Prediction - Jupyter Notebook



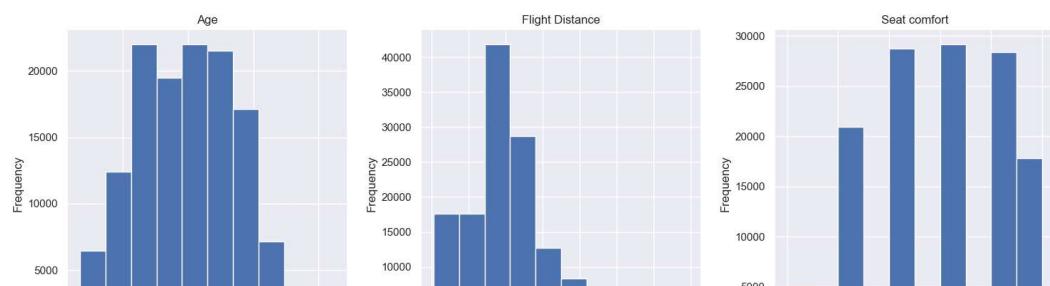
Customer Satisfaction in Airline Prediction - Jupyter Notebook



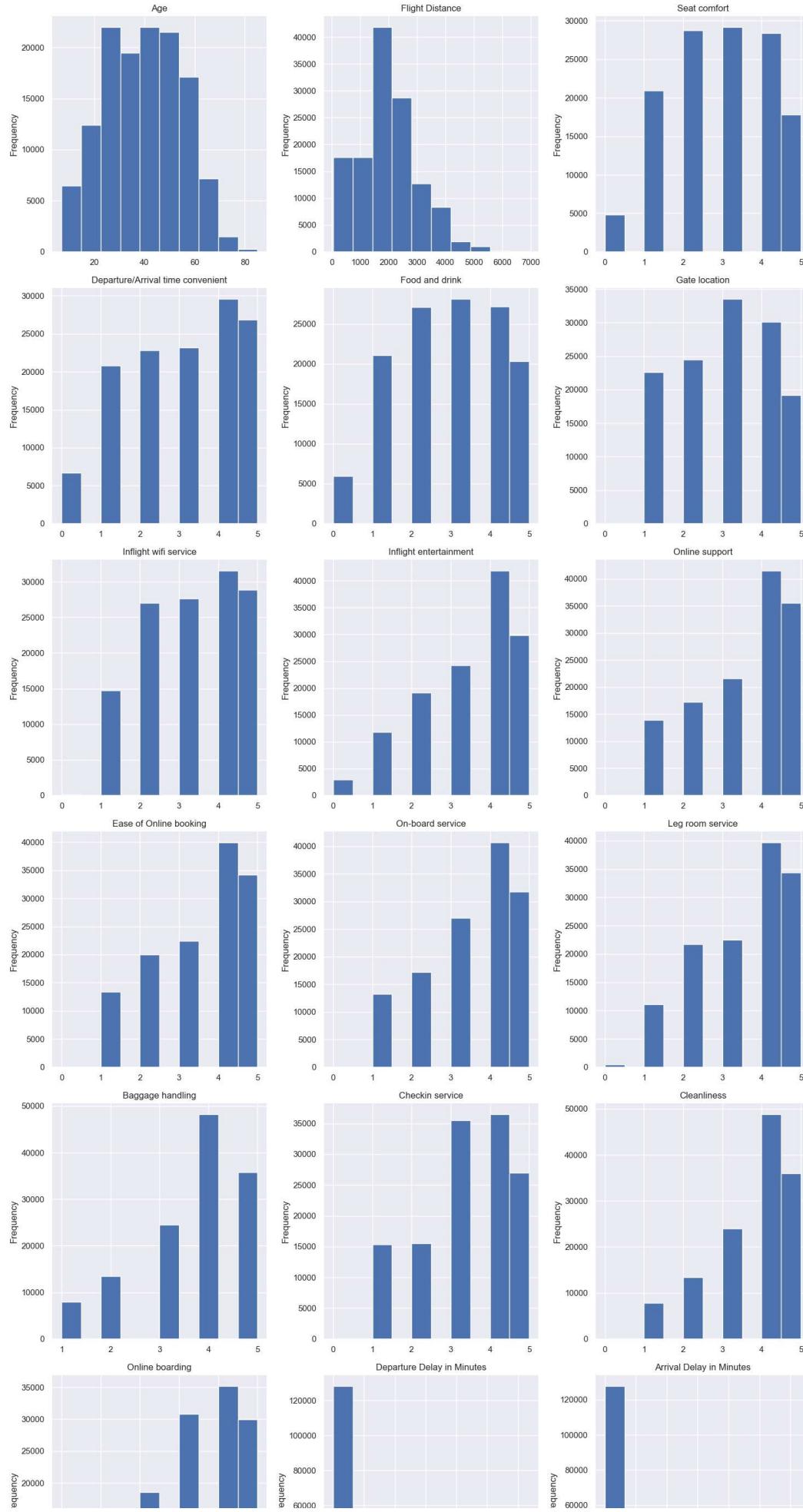
```
In [7]: # Get the names of all columns with data type 'int'
int_vars = df.select_dtypes(include=['int', 'float']).columns.tolist()

# Create a figure with subplots
num_cols = len(int_vars)
num_rows = (num_cols + 2) // 3 # To make sure there are enough rows for the subplots
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))
axs = axs.flatten()
```

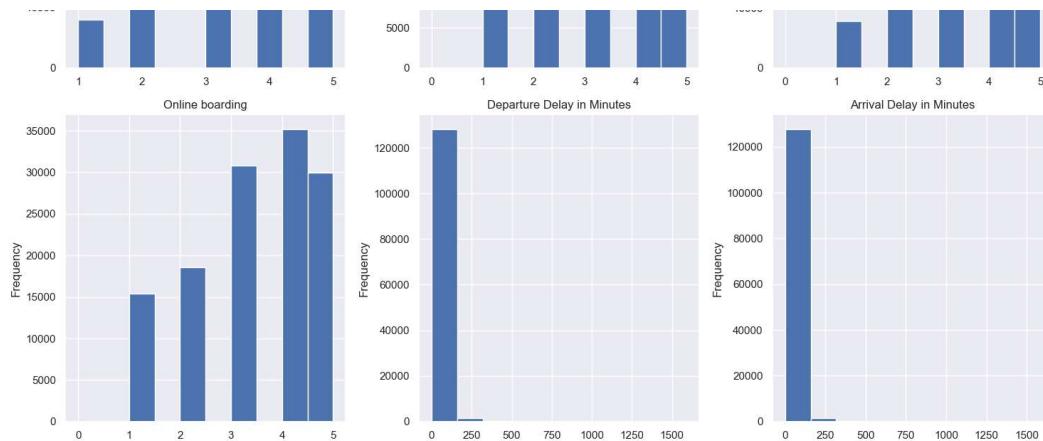
```
In [7]: # Get the names of all columns with data type 'int'  
int_vars = df.select_dtypes(include=['int', 'float']).columns.tolist()  
  
# Create a figure with subplots  
num_cols = len(int_vars)  
num_rows = (num_cols + 2) // 3 # To make sure there are enough rows for the subplots  
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))  
axs = axs.flatten()  
  
# Create a histogram for each integer variable  
for i, var in enumerate(int_vars):  
    df[var].plot.hist(ax=axs[i])  
    axs[i].set_title(var)  
  
# Remove any extra empty subplots if needed  
if num_cols < len(axs):  
    for i in range(num_cols, len(axs)):  
        fig.delaxes(axs[i])  
  
# Adjust spacing between subplots  
fig.tight_layout()  
  
# Show plot  
plt.show()
```



Customer Satisfaction in Airline Prediction - Jupyter Notebook



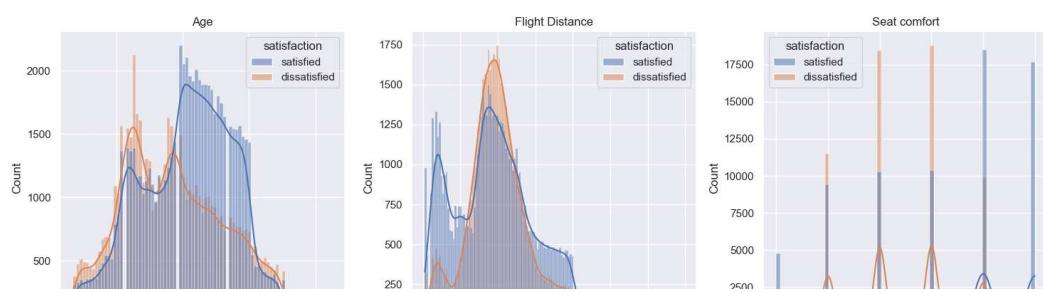
Customer Satisfaction in Airline Prediction - Jupyter Notebook



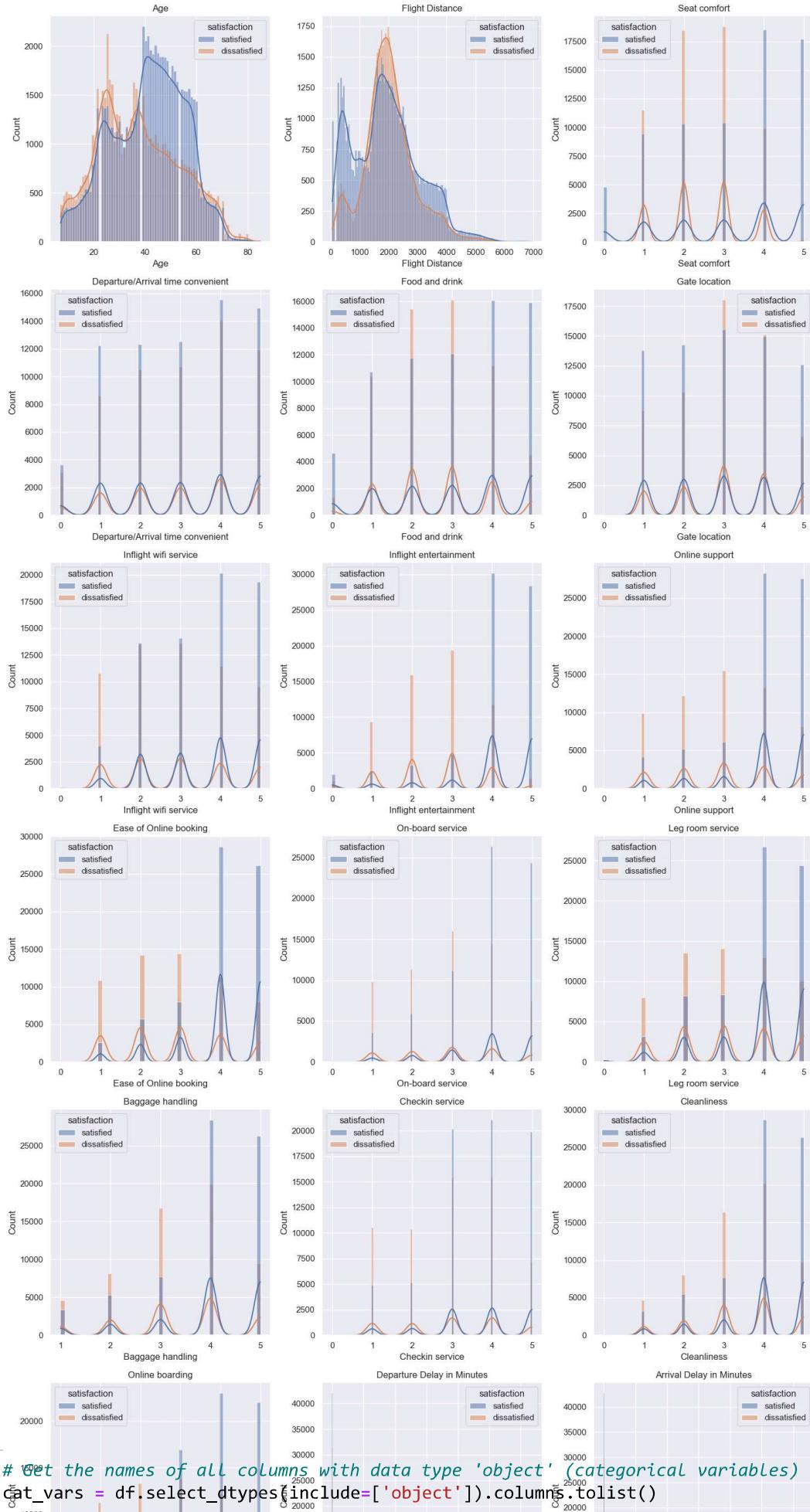
```
In [8]: # Get the names of all columns with data type 'int'
int_vars = df.select_dtypes(include=['int', 'float']).columns.tolist()
int_vars = [col for col in num_vars if col != 'satisfaction']

# Create a figure with subplots
num_cols = len(int_vars)
num_rows = (num_cols + 2) // 3 # To make sure there are enough rows for the subplots
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))
axs = axs.flatten()
```

```
In [8]: # Get the names of all columns with data type 'int'  
int_vars = df.select_dtypes(include=['int', 'float']).columns.tolist()  
int_vars = [col for col in num_vars if col != 'satisfaction']  
  
# Create a figure with subplots  
num_cols = len(int_vars)  
num_rows = (num_cols + 2) // 3 # To make sure there are enough rows for the subplots  
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))  
axs = axs.flatten()  
  
# Create a histogram for each integer variable with hue='Attrition'  
for i, var in enumerate(int_vars):  
    sns.histplot(data=df, x=var, hue='satisfaction', kde=True, ax=axs[i])  
    axs[i].set_title(var)  
  
# Remove any extra empty subplots if needed  
if num_cols < len(axs):  
    for i in range(num_cols, len(axs)):  
        fig.delaxes(axs[i])  
  
# Adjust spacing between subplots  
fig.tight_layout()  
  
# Show plot  
plt.show()
```

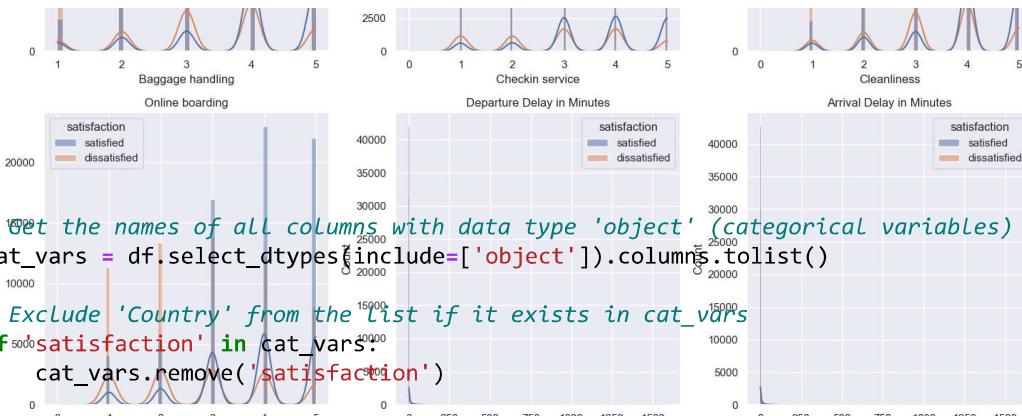


Customer Satisfaction in Airline Prediction - Jupyter Notebook



```
In [9]: # Get the names of all columns with data type 'object' (categorical variables)
cat_vars = df.select_dtypes(include=['object']).columns.tolist()
```

Customer Satisfaction in Airline Prediction - Jupyter Notebook



```
In [9]: # Get the names of all columns with data type 'object' (categorical variables)
cat_vars = df.select_dtypes(include=['object']).columns.tolist()

# Exclude 'Country' from the list if it exists in cat_vars
if 'satisfaction' in cat_vars:
    cat_vars.remove('satisfaction')

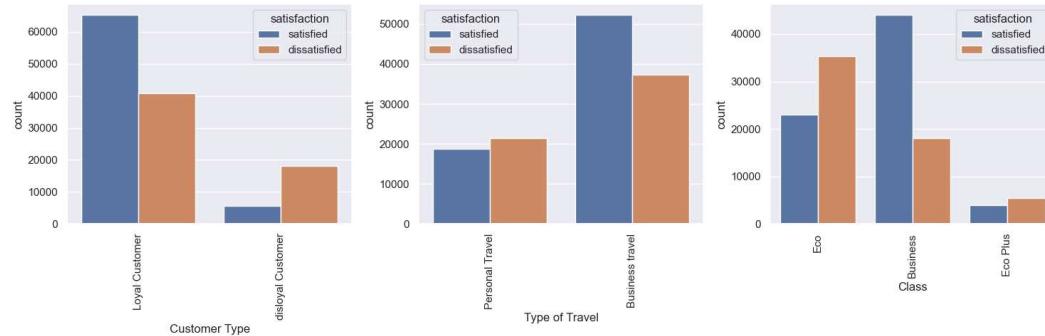
# Create a figure with subplots, but only include the required number of subplots
num_cols = len(cat_vars)
num_rows = (num_cols + 2) // 3 # To make sure there are enough rows for the subplots
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))
axs = axs.flatten()

# Create a count plot for each categorical variable
for i, var in enumerate(cat_vars):
    filtered_df = df[df[var].notnull()] # Exclude rows with NaN values in the variable
    sns.countplot(x=var, hue='satisfaction', data=filtered_df, ax=axs[i])
    axs[i].set_xticklabels(axs[i].get_xticklabels(), rotation=90)

# Remove any remaining blank subplots
for i in range(num_cols, len(axs)):
    fig.delaxes(axs[i])

# Adjust spacing between subplots
fig.tight_layout()

# Show the plot
plt.show()
```



```
In [10]: import warnings

# Get the names of all columns with data type 'object' (categorical variables)
cat_vars = df.select_dtypes(include=['object']).columns.tolist()

# Exclude 'Attrition' from the list if it exists in cat_vars
if 'satisfaction' in cat_vars:
    cat_vars.remove('satisfaction')
```

```
In [10]: import warnings

# Get the names of all columns with data type 'object' (categorical variables)
cat_vars = df.select_dtypes(include=['object']).columns.tolist()

# Exclude 'Attrition' from the list if it exists in cat_vars
if 'satisfaction' in cat_vars:
    cat_vars.remove('satisfaction')

# Create a figure with subplots, but only include the required number of subplots
num_cols = len(cat_vars)
num_rows = (num_cols + 2) // 3 # To make sure there are enough rows for the subplots
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(15, 5*num_rows))
axs = axs.flatten()

# Create a count plot for the top 6 values of each categorical variable as a density plot
for i, var in enumerate(cat_vars):
    top_values = df[var].value_counts().nlargest(6).index
    filtered_df = df[df[var].isin(top_values)]

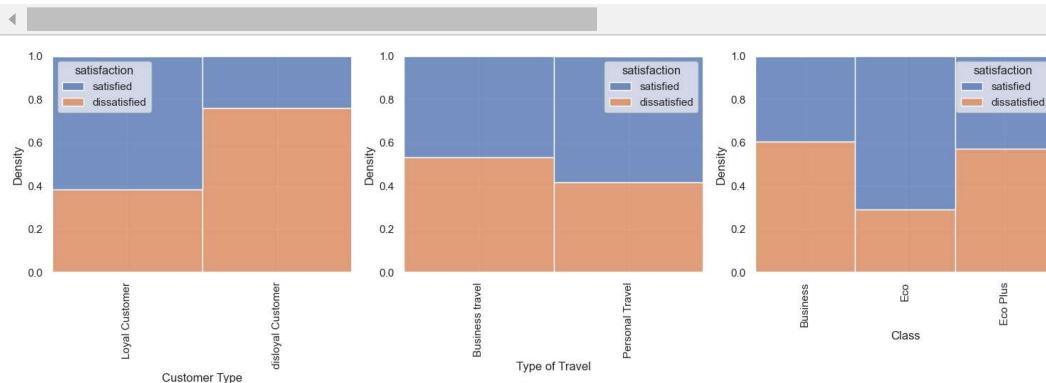
    # Set x-tick positions explicitly
    tick_positions = range(len(top_values))
    axs[i].set_xticks(tick_positions)
    axs[i].set_xticklabels(top_values, rotation=90) # Set x-tick labels

    sns.histplot(x=var, hue='satisfaction', data=filtered_df, ax=axs[i], multiple='stack')
    axs[i].set_xlabel(var)

# Remove any remaining blank subplots
for i in range(num_cols, len(axs)):
    fig.delaxes(axs[i])

# Adjust spacing between subplots
fig.tight_layout()

# Show the plot
plt.show()
```



```
In [11]: # Specify the maximum number of categories to show individually
max_categories = 5

# Filter categorical columns with 'object' data type
cat_cols = [col for col in df.columns if col != 'y' and df[col].dtype == 'object']

# Create a figure with subplots
num_cols = len(cat_cols)
num_rows = (num_cols + 2) // 3
```

```
In [11]: # Specify the maximum number of categories to show individually
max_categories = 5

# Filter categorical columns with 'object' data type
cat_cols = [col for col in df.columns if col != 'y' and df[col].dtype == 'object']

# Create a figure with subplots
num_cols = len(cat_cols)
num_rows = (num_cols + 2) // 3
fig, axs = plt.subplots(nrows=num_rows, ncols=3, figsize=(20, 5*num_rows))

# Flatten the axs array for easier indexing
axs = axs.flatten()

# Create a pie chart for each categorical column
for i, col in enumerate(cat_cols):
    if i < len(axs): # Ensure we don't exceed the number of subplots
        # Count the number of occurrences for each category
        cat_counts = df[col].value_counts()

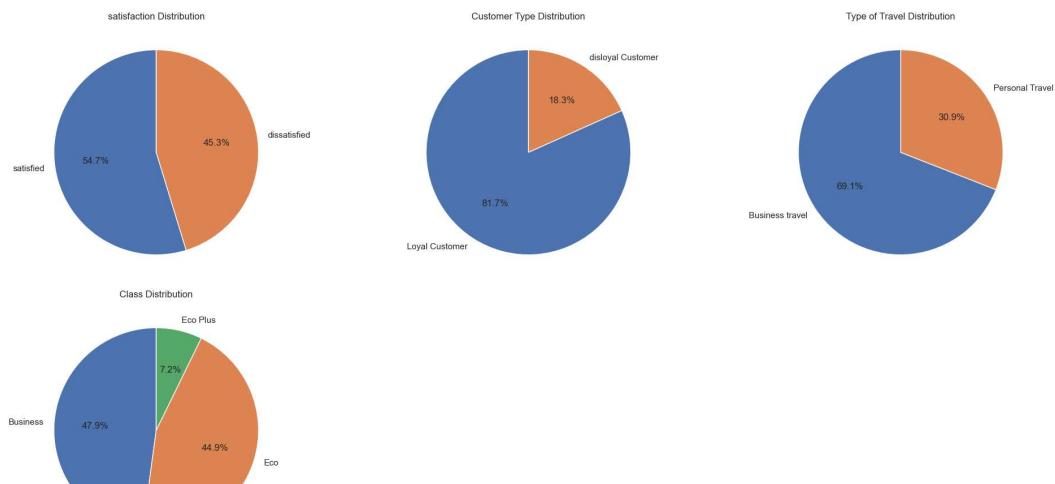
        # Group categories beyond the top max_categories as 'Other'
        if len(cat_counts) > max_categories:
            cat_counts_top = cat_counts[:max_categories]
            cat_counts_other = pd.Series(cat_counts[max_categories:]).sum(), index=[len(cat_counts_top)]
            cat_counts = cat_counts_top.append(cat_counts_other)

        # Create a pie chart
        axs[i].pie(cat_counts, labels=cat_counts.index, autopct='%1.1f%%', startangle=90)
        axs[i].set_title(f'{col} Distribution')

# Remove any extra empty subplots if needed
if num_cols < len(axs):
    for i in range(num_cols, len(axs)):
        fig.delaxes(axs[i])

# Adjust spacing between subplots
fig.tight_layout()

# Show plot
plt.show()
```



```
In [12]: # Check the amountnt of missing value
check_missing = df.isnull().sum() * 100 / df.shape[0]
check_missing[check_missing > 0].sort_values(ascending=False)
```

Out[12]: Arrival Delay in Minutes 0.302587
Type of Date

In [13]: df.shape

```
In [12]: # Check the amount of missing value
check_missing = df.isnull().sum() * 100 / df.shape[0]
check_missing[check_missing > 0].sort_values(ascending=False)
```

Out[12]: Arrival Delay in Minutes 0.302587
Bytelfloat

Data Preprocessing Part 2

```
In [13]: df.shape
```

Out[13]: (129880, 22)

```
In [14]: # Drop null value from Arrival Delay in Minutes columns because the null value is
df.dropna(subset=['Arrival Delay in Minutes'], inplace=True)
df.shape
```

Out[14]: (129487, 22)

Label Encoding for Object Datatypes

```
In [15]: # Loop over each column in the DataFrame where dtype is 'object'
for col in df.select_dtypes(include=['object']).columns:
```

```
# Print the column name and the unique values
print(f'{col}: {df[col].unique()}')
```

```
satisfaction: ['satisfied' 'dissatisfied']
Customer Type: ['Loyal Customer' 'disloyal Customer']
Type of Travel: ['Personal Travel' 'Business travel']
Class: ['Eco' 'Business' 'Eco Plus']
```

```
In [16]: from sklearn import preprocessing
```

```
# Loop over each column in the DataFrame where dtype is 'object'
for col in df.select_dtypes(include=['object']).columns:
```

```
# Initialize a LabelEncoder object
label_encoder = preprocessing.LabelEncoder()
```

```
# Fit the encoder to the unique values in the column
label_encoder.fit(df[col].unique())
```

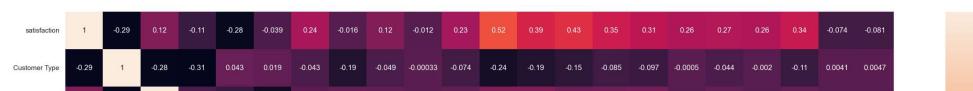
```
# Transform the column using the encoder
df[col] = label_encoder.transform(df[col])
```

```
# Print the column name and the unique encoded values
print(f'{col}: {df[col].unique()}')
```

```
satisfaction: [1 0]
Customer Type: [0 1]
Type of Travel: [1 0]
Class: [1 0 2]
```

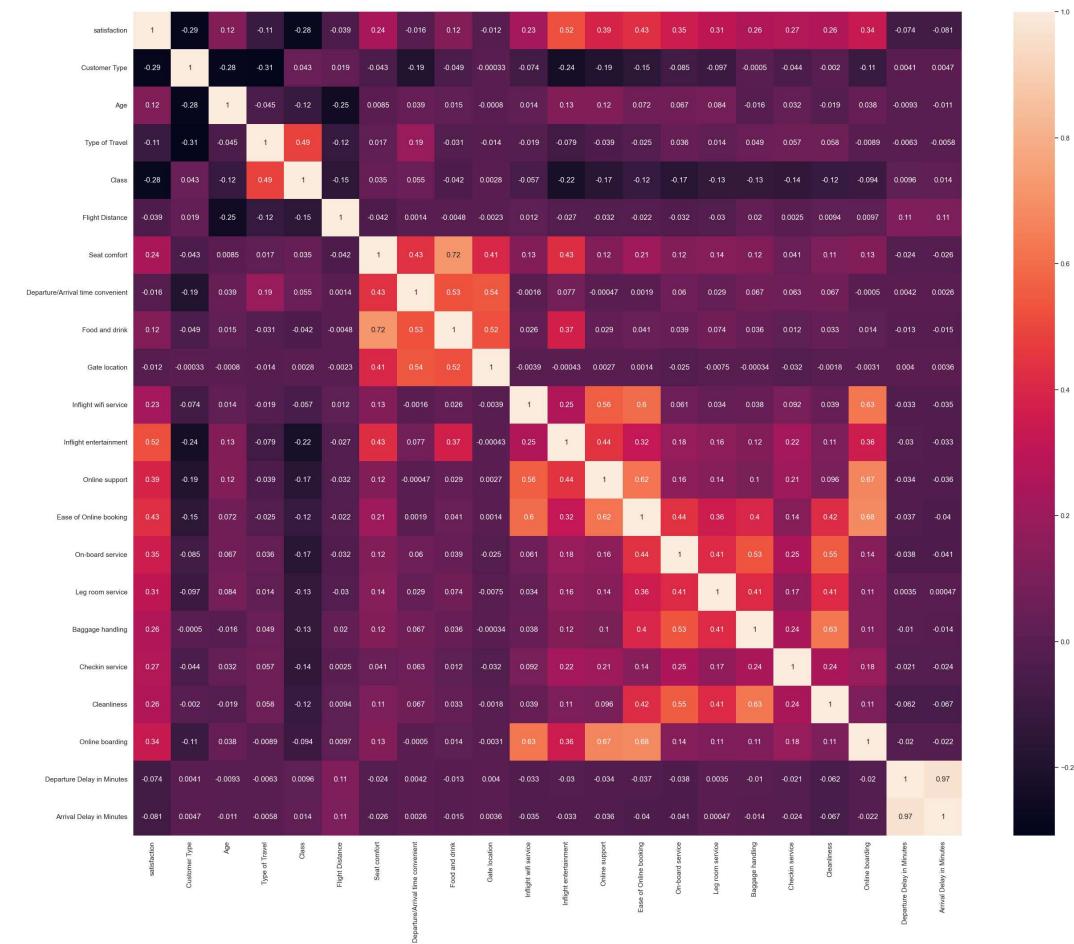
```
In [17]: # Correlation Heatmap
plt.figure(figsize=(30, 24))
sns.heatmap(df.corr(), fmt='.2g', annot=True)
```

Out[17]: <AxesSubplot:>



```
In [17]: # Correlation Heatmap
plt.figure(figsize=(30, 24))
sns.heatmap(df.corr(), fmt='.2g', annot=True)
```

Out[17]: <AxesSubplot:>



Train Test Split

```
In [18]: X = df.drop('satisfaction', axis=1)
y = df['satisfaction']
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
X_train, X_test, y_train, y_test = train_test_split(X,y, test_size=0.2,random_state=42)
```

```
In [19]: from scipy import stats
```

```
# Define the columns for which you want to remove outliers
selected_columns = ['Departure Delay in Minutes', 'Arrival Delay in Minutes']

# Calculate the Z-scores for the selected columns in the training data
z_scores = np.abs(stats.zscore(X_train[selected_columns]))

# Set a threshold value for outlier detection (e.g., 3)
threshold = 3
```

```
In [19]: from scipy import stats

# Define the columns for which you want to remove outliers
selected_columns = ['Departure Delay in Minutes', 'Arrival Delay in Minutes']

# Calculate the Z-scores for the selected columns in the training data
z_scores = np.abs(stats.zscore(X_train[selected_columns]))

# Set a threshold value for outlier detection (e.g., 3)
threshold = 3

# Find the indices of outliers based on the threshold
outlier_indices = np.where(z_scores > threshold)[0]

# Remove the outliers from the training data
X_train = X_train.drop(X_train.index[outlier_indices])
y_train = y_train.drop(y_train.index[outlier_indices])
```

Decision Tree Classifier

```
In [20]: from sklearn.tree import DecisionTreeClassifier
from sklearn.model_selection import GridSearchCV
dtree = DecisionTreeClassifier(class_weight='balanced')
param_grid = {
    'max_depth': [3, 4, 5, 6, 7, 8],
    'min_samples_split': [2, 3, 4],
    'min_samples_leaf': [1, 2, 3, 4],
    'random_state': [0, 42]
}

# Perform a grid search with cross-validation to find the best hyperparameters
grid_search = GridSearchCV(dtree, param_grid, cv=5)
grid_search.fit(X_train, y_train)

# Print the best hyperparameters
print(grid_search.best_params_)
```

{'max_depth': 8, 'min_samples_leaf': 2, 'min_samples_split': 2, 'random_state': 42}

```
In [21]: from sklearn.tree import DecisionTreeClassifier
dtree = DecisionTreeClassifier(random_state=42, max_depth=8, min_samples_leaf=2,
dtree.fit(X_train, y_train)
```

```
Out[21]: DecisionTreeClassifier(class_weight='balanced', max_depth=8, min_samples_leaf=2,
                                 random_state=42)
```

```
In [22]: from sklearn.metrics import accuracy_score
y_pred = dtree.predict(X_test)
print("Accuracy Score : ", round(accuracy_score(y_test, y_pred)*100 ,2), "%")
```

Accuracy Score : 90.61 %

```
In [23]: from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
print('F-1 Score : ',(f1_score(y_test, y_pred, average='micro')))
print('Precision Score : ',(precision_score(y_test, y_pred, average='micro')))
print('Recall Score : ',(recall_score(y_test, y_pred, average='micro')))
print('Jaccard Score : ',(jaccard_score(y_test, y_pred, average='micro')))
print('Log Loss : ',(log_loss(y_test, y_pred)))
```

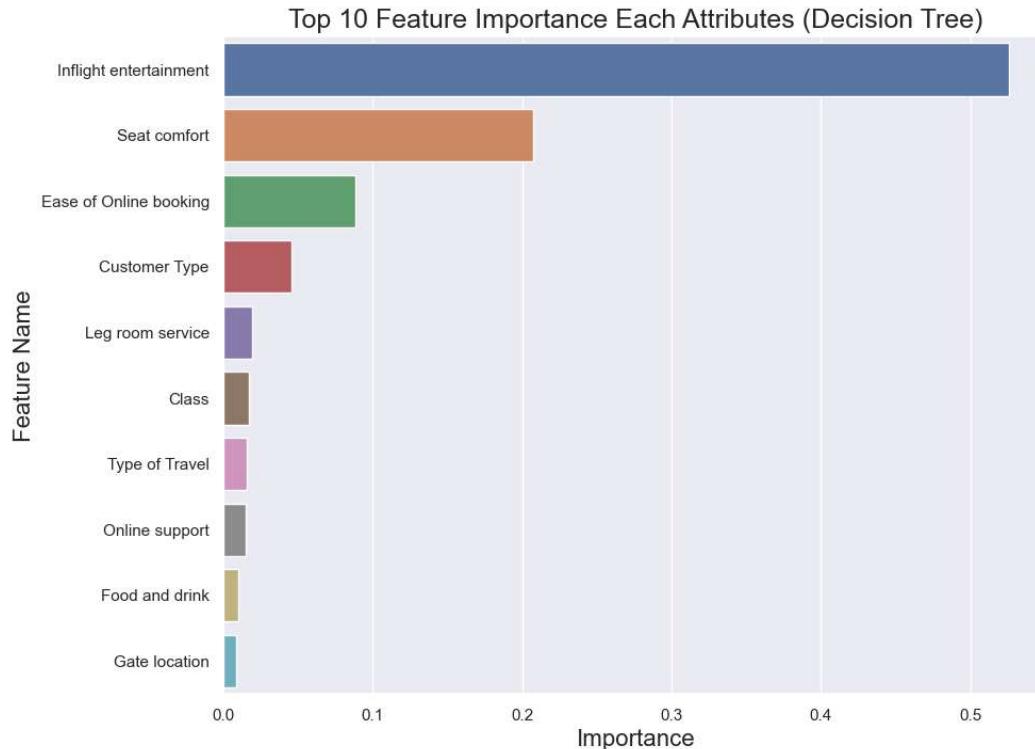
F-1 Score : 0.9060931346049887
Precision Score : 0.9060931346049887

```
In [23]: from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
print('F-1 Score :',(f1_score(y_test, y_pred, average='micro')))
print('Precision Score :',(precision_score(y_test, y_pred, average='micro')))
print('Recall Score :',(recall_score(y_test, y_pred, average='micro')))
print('Jaccard Score :',(jaccard_score(y_test, y_pred, average='micro')))
print('Log Loss :',(log_loss(y_test, y_pred)))
```

F-1 Score : 0.9060931346049887
 Precision Score : 0.9060931346049887
 Recall Score : 0.9060931346049887
 Jaccard Score : 0.8283092128485704
 Log Loss : 3.243455982335027

```
In [24]: imp_df = pd.DataFrame({
    "Feature Name": X_train.columns,
    "Importance": dtree.feature_importances_
})
fi = imp_df.sort_values(by="Importance", ascending=False)

fi2 = fi.head(10)
plt.figure(figsize=(10,8))
sns.barplot(data=fi2, x='Importance', y='Feature Name')
plt.title('Top 10 Feature Importance Each Attributes (Decision Tree)', fontsize=16)
plt.xlabel ('Importance', fontsize=16)
plt.ylabel ('Feature Name', fontsize=16)
plt.show()
```

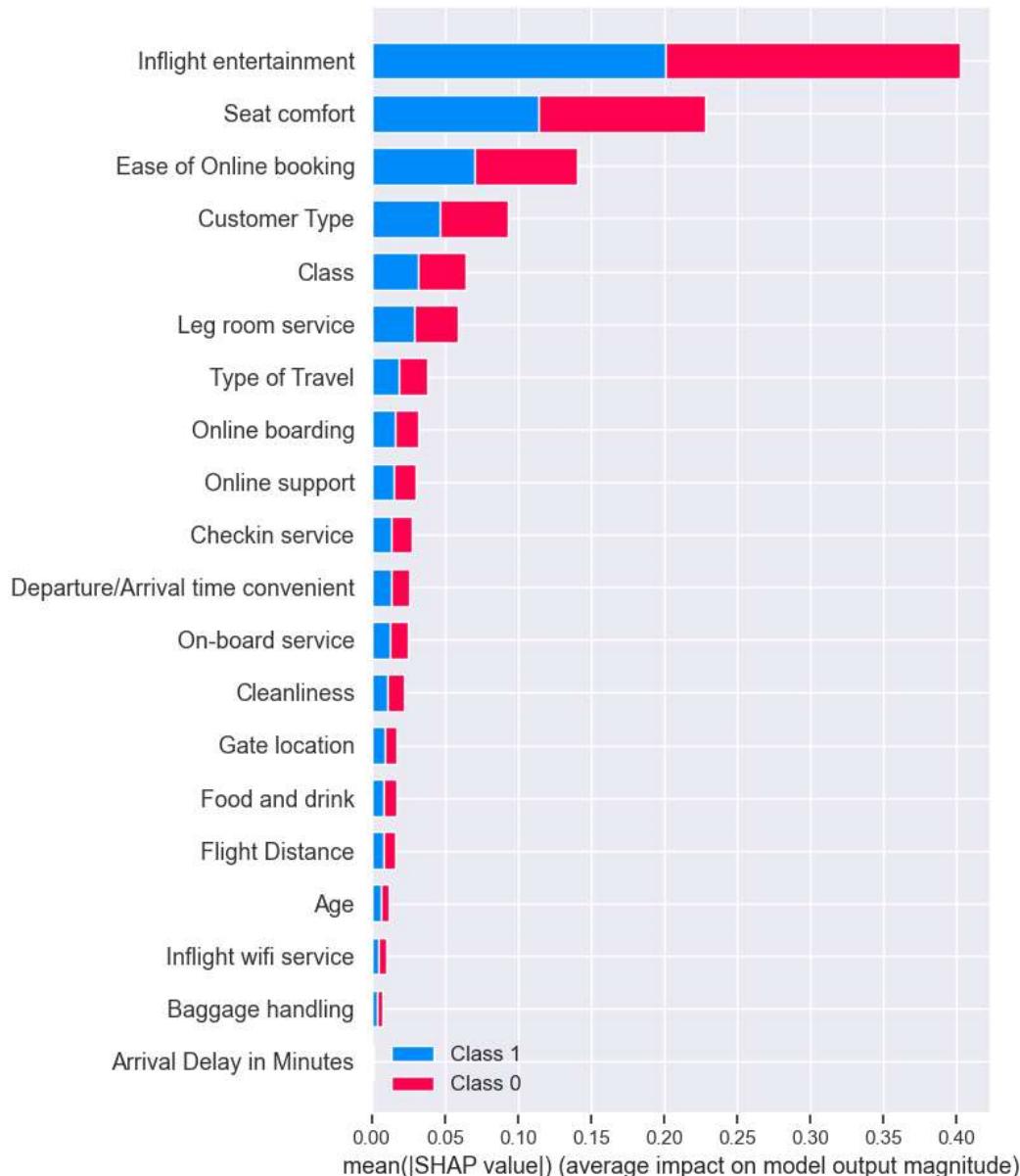


```
In [25]: import shap
explainer = shap.TreeExplainer(dtree)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values, X_test)
```

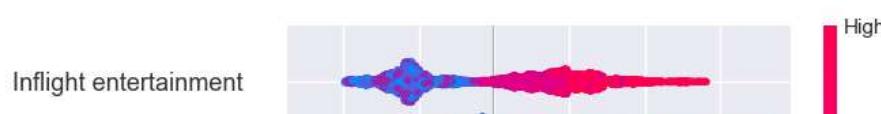
Using `tqdm.autonotebook.tqdm` in notebook mode. Use `tqdm.tqdm` instead to force console mode (e.g. in jupyter console)

```
In [25]: import shap
explainer = shap.TreeExplainer(dtree)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values, X_test)
```

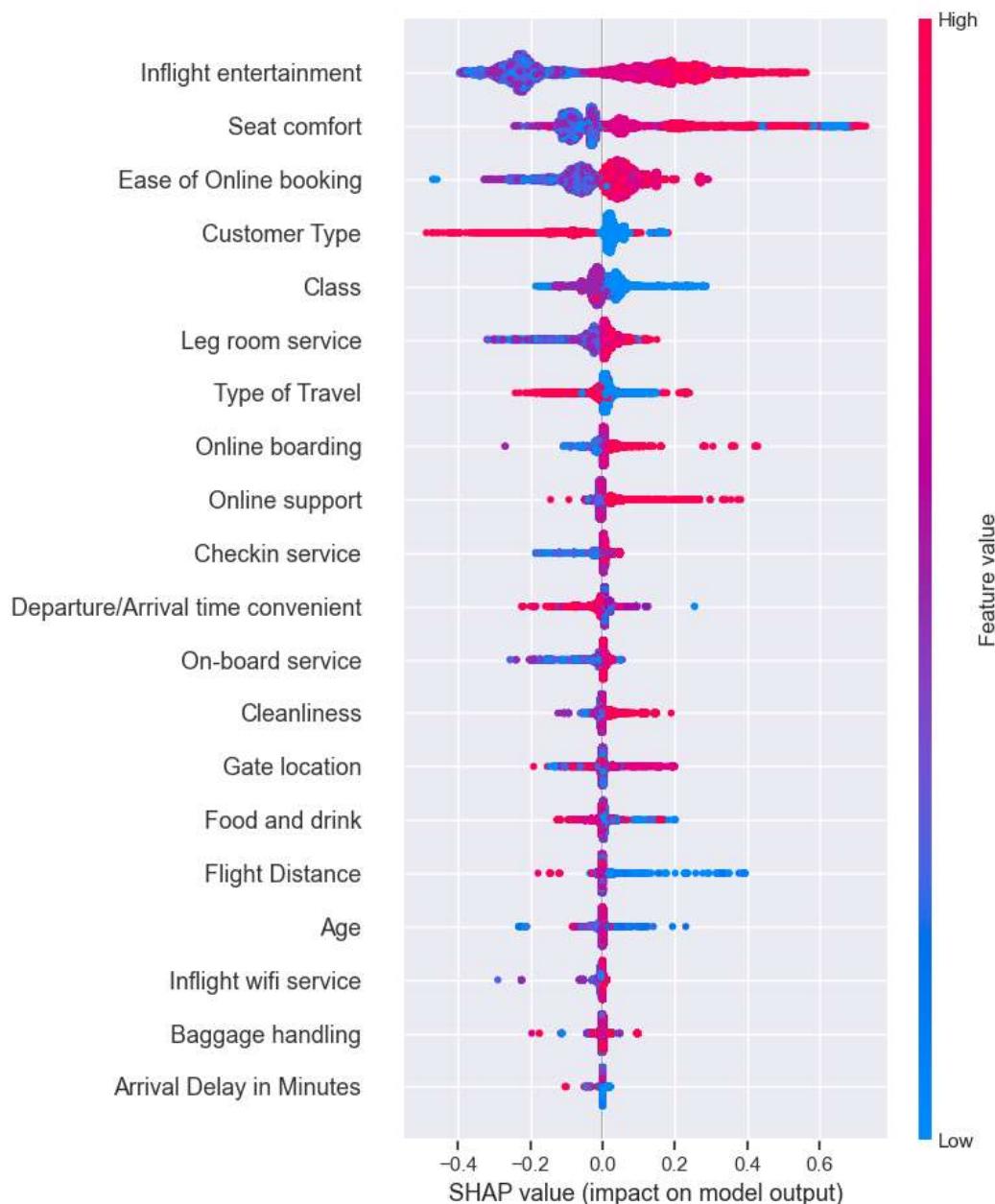
Using `tqdm.autonotebook.tqdm` in notebook mode. Use `tqdm.tqdm` instead to force console mode (e.g. in jupyter console)



```
In [26]: # compute SHAP values
explainer = shap.TreeExplainer(dtree)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values[1], X_test.values, feature_names = X_test.columns)
```



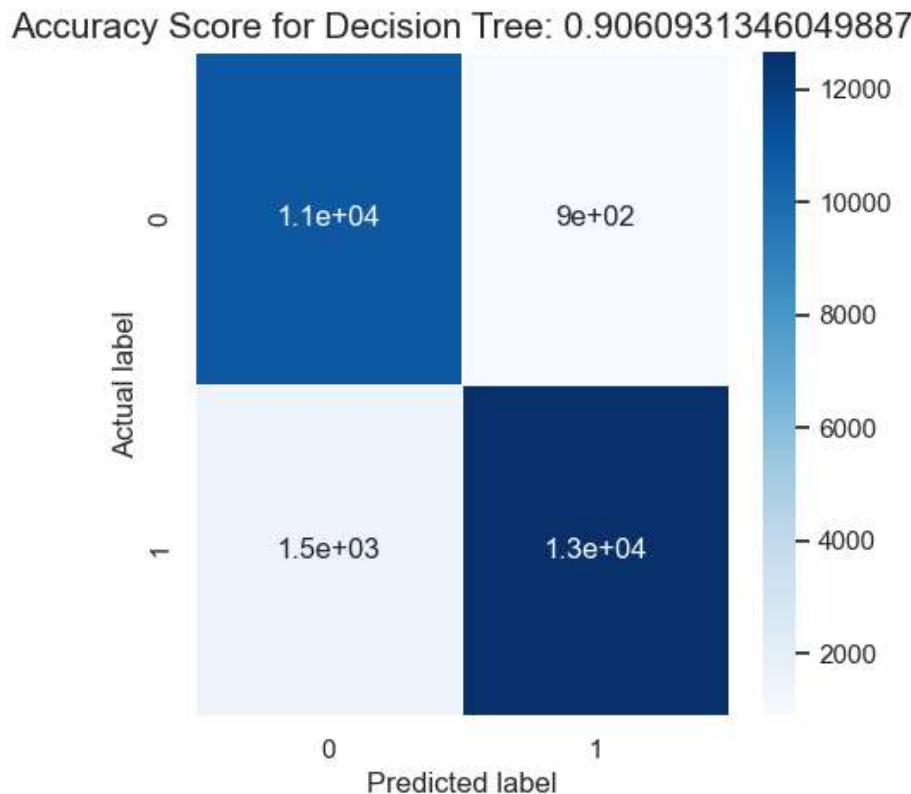
```
In [26]: # compute SHAP values
explainer = shap.TreeExplainer(dtree)
shap_values = explainer.shap_values(X_test)
shap.summary_plot(shap_values[1], X_test.values, feature_names = X_test.columns)
```



```
In [27]: from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(5,5))
sns.heatmap(data=cm, linewidths=.5, annot=True, cmap = 'Blues')
plt.ylabel('Actual label')
plt.xlabel('Predicted label')
all_sample_title = 'Accuracy Score for Decision Tree: {}'.format(dtree.score(X_1))
plt.title(all_sample_title, size = 15)
```

```
In [27]: from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
plt.figure(figsize=(5,5))
sns.heatmap(data=cm, linewidths=.5, annot=True, cmap = 'Blues')
plt.ylabel('Actual label')
plt.xlabel('Predicted label')
all_sample_title = 'Accuracy Score for Decision Tree: {0}'.format(dtree.score(X_1))
plt.title(all_sample_title, size = 15)
```

Out[27]: Text(0.5, 1.0, 'Accuracy Score for Decision Tree: 0.9060931346049887')



```
In [28]: from sklearn.metrics import roc_curve, roc_auc_score
y_pred_proba = dtree.predict_proba(X_test)[:, :, 1]

df_actual_predicted = pd.concat([pd.DataFrame(np.array(y_test), columns=['y_actual']), df_actual_predicted], axis=1)
df_actual_predicted.index = y_test.index

fpr, tpr, tr = roc_curve(df_actual_predicted['y_actual'], df_actual_predicted['y_pred'])
auc = roc_auc_score(df_actual_predicted['y_actual'], df_actual_predicted['y_pred'])
```

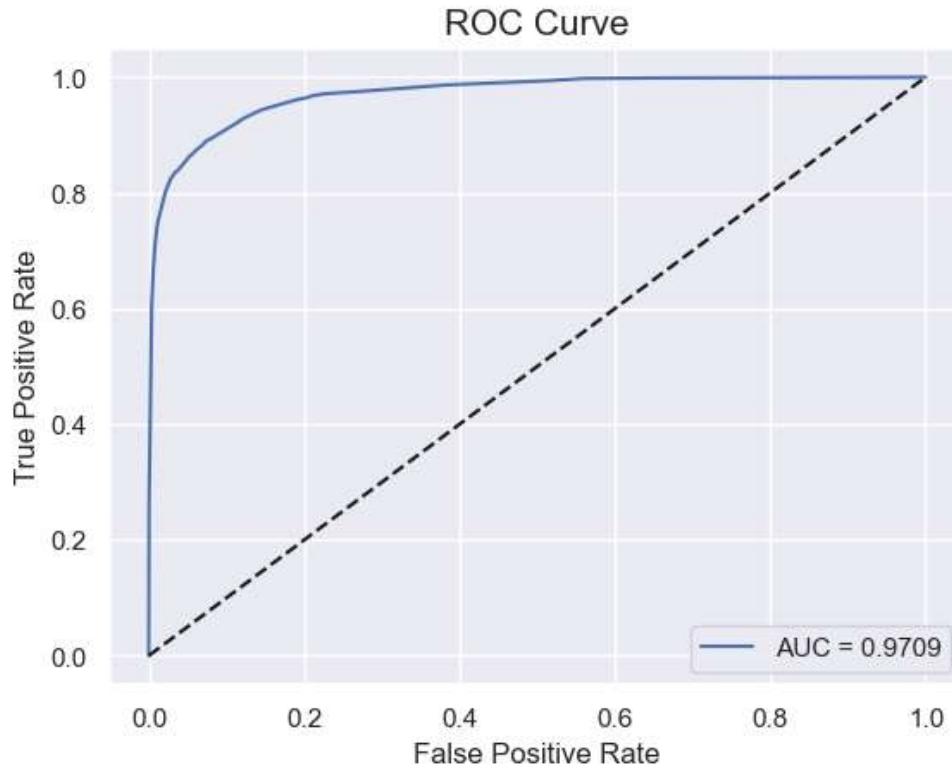
```
In [28]: from sklearn.metrics import roc_curve, roc_auc_score
y_pred_proba = dtree.predict_proba(X_test)[:,1]

df_actual_predicted = pd.concat([pd.DataFrame(np.array(y_test)), columns=['y_actual']], axis=1)
df_actual_predicted.index = y_test.index

fpr, tpr, tr = roc_curve(df_actual_predicted['y_actual'], df_actual_predicted['y_pred'])
auc = roc_auc_score(df_actual_predicted['y_actual'], df_actual_predicted['y_pred'])

plt.plot(fpr, tpr, label='AUC = %0.4f' %auc)
plt.plot(fpr, fpr, linestyle = '--', color='k')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve', size = 15)
plt.legend()
```

Out[28]: <matplotlib.legend.Legend at 0x2490475ddc0>



Random Forest Classifier

```
In [30]: from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
rfc = RandomForestClassifier(class_weight='balanced')
param_grid = {
    'n_estimators': [100, 200],
    'max_depth': [None, 5, 10],
    'max_features': ['sqrt', 'log2', None],
    'random_state': [0, 42]
}
```

```
In [30]: from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import GridSearchCV
rfc = RandomForestClassifier(class_weight='balanced')
param_grid = {
    'n_estimators': [100, 200],
    'max_depth': [None, 5, 10],
    'max_features': ['sqrt', 'log2', None],
    'random_state': [0, 42]
}

# Perform a grid search with cross-validation to find the best hyperparameters
grid_search = GridSearchCV(rfc, param_grid, cv=5)
grid_search.fit(X_train, y_train)

# Print the best hyperparameters
print(grid_search.best_params_)

{'max_depth': None, 'max_features': 'sqrt', 'n_estimators': 200, 'random_state': 42}
```

```
In [31]: from sklearn.ensemble import RandomForestClassifier
rfc = RandomForestClassifier(random_state=42, max_depth=None, max_features='sqrt')
rfc.fit(X_train, y_train)
```

```
Out[31]: RandomForestClassifier(class_weight='balanced', max_features='sqrt',
                                n_estimators=200, random_state=42)
```

```
In [32]: y_pred = rfc.predict(X_test)
print("Accuracy Score :", round(accuracy_score(y_test, y_pred)*100 ,2), "%")
```

Accuracy Score : 95.39 %

```
In [33]: from sklearn.metrics import accuracy_score, f1_score, precision_score, recall_score
print('F-1 Score :',(f1_score(y_test, y_pred, average='micro')))
print('Precision Score :',(precision_score(y_test, y_pred, average='micro')))
print('Recall Score :',(recall_score(y_test, y_pred, average='micro')))
print('Jaccard Score :',(jaccard_score(y_test, y_pred, average='micro')))
print('Log Loss :',(log_loss(y_test, y_pred)))
```

F-1 Score : 0.9538960537493243
 Precision Score : 0.9538960537493243
 Recall Score : 0.9538960537493243
 Jaccard Score : 0.9118558984201979
 Log Loss : 1.5923887104135401

```
In [34]: imp_df = pd.DataFrame({
    "Feature Name": X_train.columns,
    "Importance": rfc.feature_importances_
})
fi = imp_df.sort_values(by="Importance", ascending=False)

fi2 = fi.head(10)
plt.figure(figsize=(10,8))
sns.barplot(data=fi2, x='Importance', y='Feature Name')
```

```
In [34]: imp_df = pd.DataFrame({
    "Feature Name": X_train.columns,
    "Importance": rfc.feature_importances_
})
fi = imp_df.sort_values(by="Importance", ascending=False)

fi2 = fi.head(10)
plt.figure(figsize=(10,8))
sns.barplot(data=fi2, x='Importance', y='Feature Name')
plt.title('Top 10 Feature Importance Each Attributes (Random Forest)', fontsize=16)
plt.xlabel ('Importance', fontsize=16)
plt.ylabel ('Feature Name', fontsize=16)
plt.show()
```

