

Customer Satisfaction Rating

* “Prediction” *

Overview

Created by: Md. Rafiquzzaman Rafi

This project develops a machine learning model to predict customer satisfaction in airlines using advanced techniques like **CatBoost** and **Optuna** for hyperparameter optimization.

Data Preprocessing

Exploratory Data Analysis

Feature Engineering

Model Training

Key Insights

Data Import and Basic Analysis



```
1 import numpy as np
2 import pandas as pd
3 df = pd.read_csv('data/train.csv')
4 df.head()
```

1 df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 103904 entries, 0 to 103903
Data columns (total 25 columns):
 #   Column           Non-Null Count   Dtype  
--- 
 0   Unnamed: 0        103904 non-null    int64  
 1   id               103904 non-null    int64  
 2   Gender            103904 non-null    object  
 3   Customer Type     103904 non-null    object  
 4   Age               103904 non-null    int64  
 5   Type of Travel    103904 non-null    object  
 6   Class              103904 non-null    object  
 7   Flight Distance   103904 non-null    int64  
 8   Inflight wifi service 103904 non-null    int64  
 9   Departure/Arrival time convenient 103904 non-null    int64  
 10  Ease of Online booking 103904 non-null    int64  
 11  Gate location     103904 non-null    int64  
 12  Food and drink    103904 non-null    int64  
 13  Online boarding   103904 non-null    int64  
 14  Seat comfort       103904 non-null    int64  
 15  Inflight entertainment 103904 non-null    int64  
 16  On-board service   103904 non-null    int64  
 17  Leg room service   103904 non-null    int64  
 18  Baggage handling   103904 non-null    int64  
 19  Checkin service    103904 non-null    int64  
 ...
 23  Arrival Delay in Minutes 103594 non-null    float64 
 24  satisfaction       103904 non-null    object  
dtypes: float64(1), int64(19), object(5)
memory usage: 19.8+ MB
```

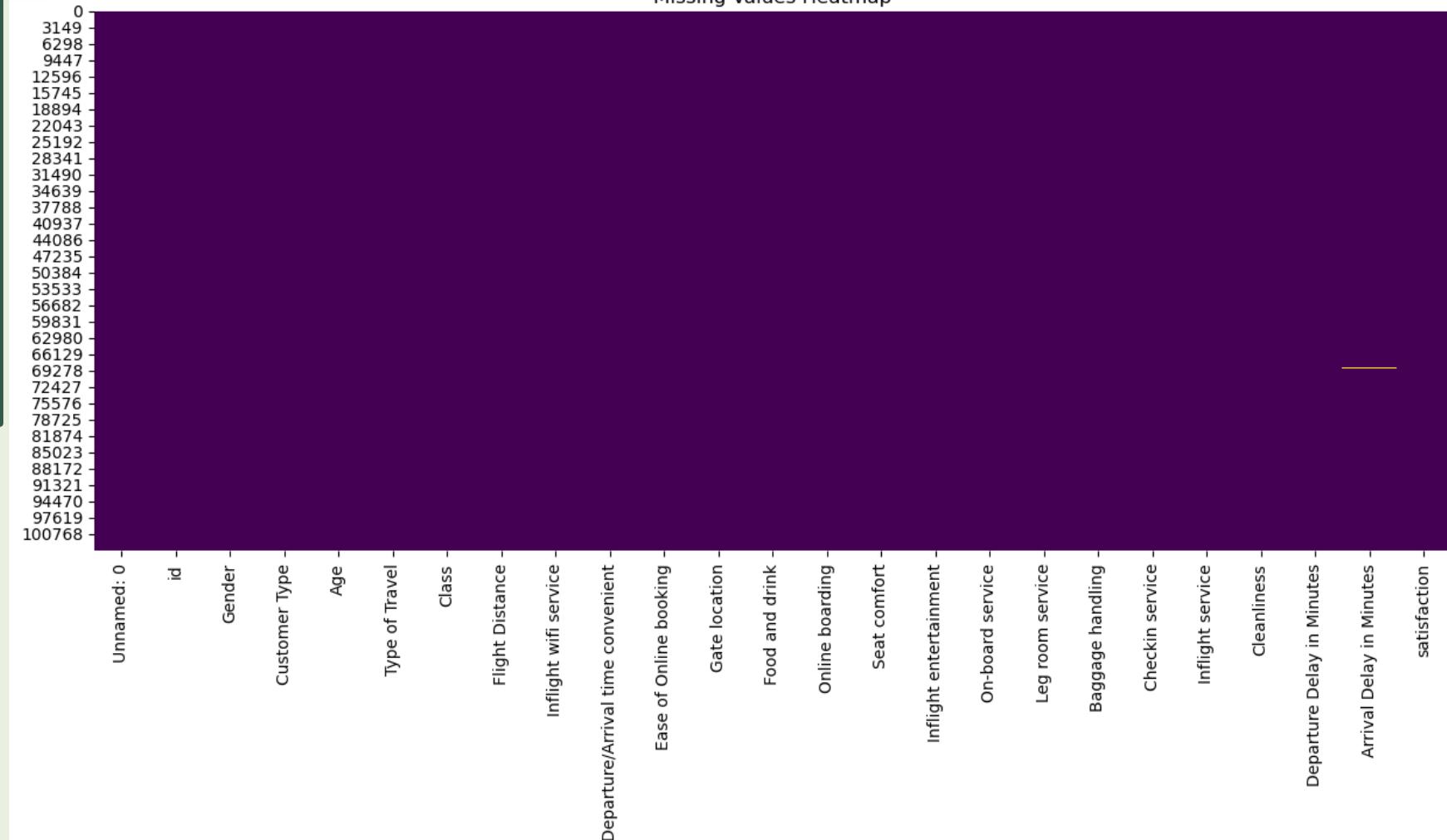
Unnamed: 0	id	Gender	Customer Type	Age	Type of Travel	Class	Flight Distance	Inflight wifi service	Departure/Arrival time convenient	...	Inflight entertainment	On-board service	Leg room service	Baggage handling	Checkin service	Inflight service	Cleanliness	Departure Delay in Minutes	Arrival Delay in Minutes	satisfaction	
0	0	70172	Male	Loyal Customer	13	Personal Travel	Eco Plus	460	3	4	...	5	4	3	4	4	5	5	25	18.0	neutral or dissatisfied
1	1	5047	Male	disloyal Customer	25	Business travel	Business	235	3	2	...	1	1	5	3	1	4	1	1	6.0	neutral or dissatisfied
2	2	110028	Female	Loyal Customer	26	Business travel	Business	1142	2	2	...	5	4	3	4	4	4	5	0	0.0	satisfied
3	3	24026	Female	Loyal Customer	25	Business travel	Business	562	2	5	...	2	2	5	3	1	4	2	11	9.0	neutral or dissatisfied
4	4	119299	Male	Loyal Customer	61	Business travel	Business	214	3	3	...	3	3	4	4	3	3	3	0	0.0	satisfied

5 rows × 25 columns

Data Preprocessing

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3 missing_values = df.isnull().sum()
4 print(missing_values)
5 plt.figure(figsize=(15, 6))
6 sns.heatmap(df.isnull(), cbar=False, cmap='viridis')
7 plt.title('Missing Values Heatmap')
8 plt.show()
```

```
Unnamed: 0          0
id                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  310
satisfaction        0
dtype: int64
```



Data Preprocessing

```
1 from sklearn.preprocessing import LabelEncoder, OneHotEncoder
2
3 cat_cols = ['Type of Travel', 'Customer Type', 'Class', 'satisfaction']
4 oh_col = ['Gender']
5
6 encoding_dict = {}
7 le = LabelEncoder()
8 oh = OneHotEncoder(sparse_output=False, drop='first')
9
10 for col in cat_cols:
11     df[col] = le.fit_transform(df[col].astype(str))
12     encoding_dict[col] = dict(zip(le.classes_, range(len(le.classes_))))
13
14 encoded_data = oh.fit_transform(df[oh_col])
15 encoded_df = pd.DataFrame(encoded_data, columns=oh.get_feature_names_out(oh_col))
16
17 df = df.drop(columns=oh_col).reset_index(drop=True)
18 df = pd.concat([df, encoded_df], axis=1)
19
20 print("\nEncoding Mappings:")
21 for col, mapping in encoding_dict.items():
22     print(f"{col}: {mapping}")
23
24 df.head()
25
```

Encoding Mappings:

Type of Travel: {'Business travel': 0, 'Personal Travel': 1}
Customer Type: {'Loyal Customer': 0, 'disloyal Customer': 1}
Class: {'Business': 0, 'Eco': 1, 'Eco Plus': 2}
satisfaction: {'neutral or dissatisfied': 0, 'satisfied': 1}

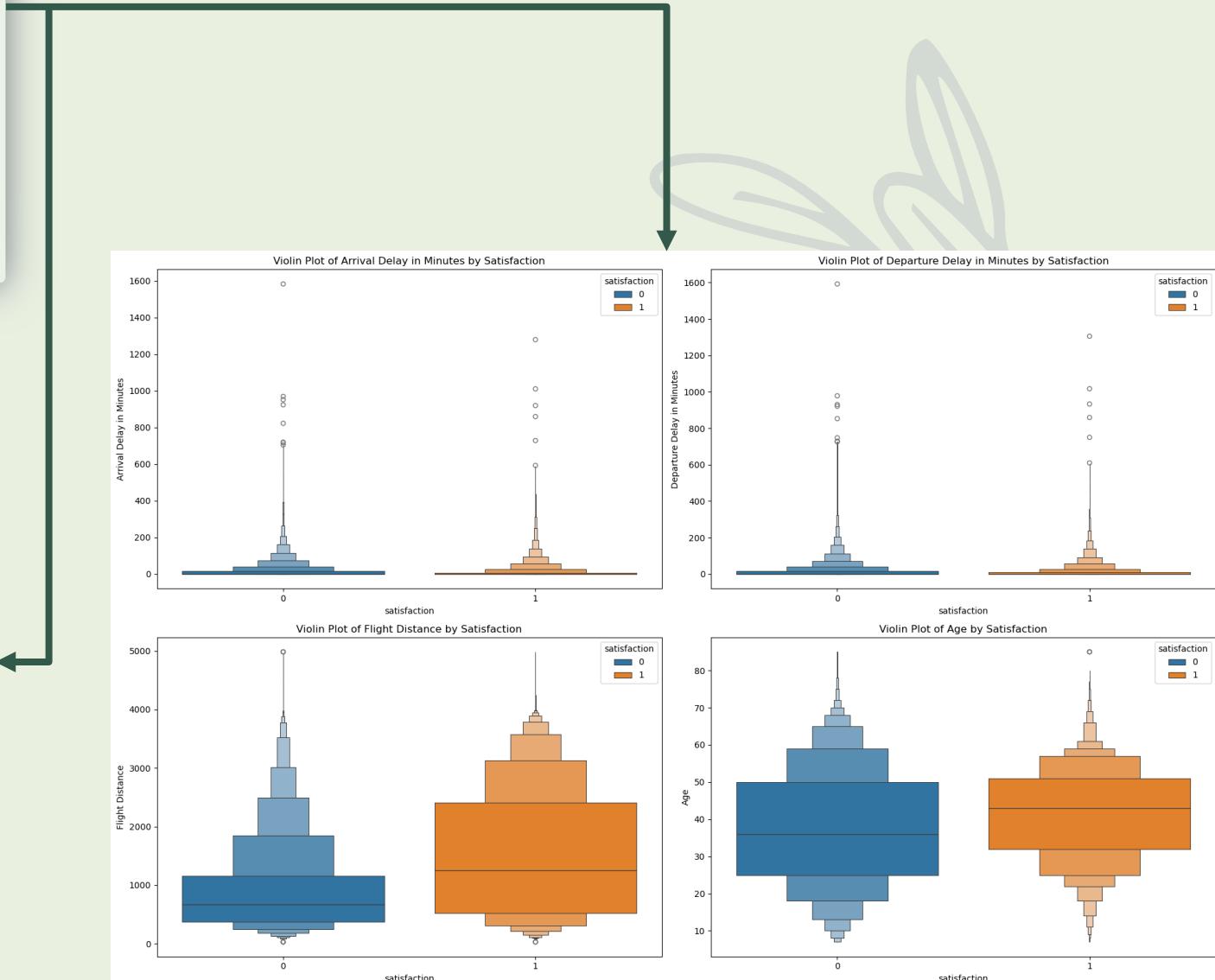
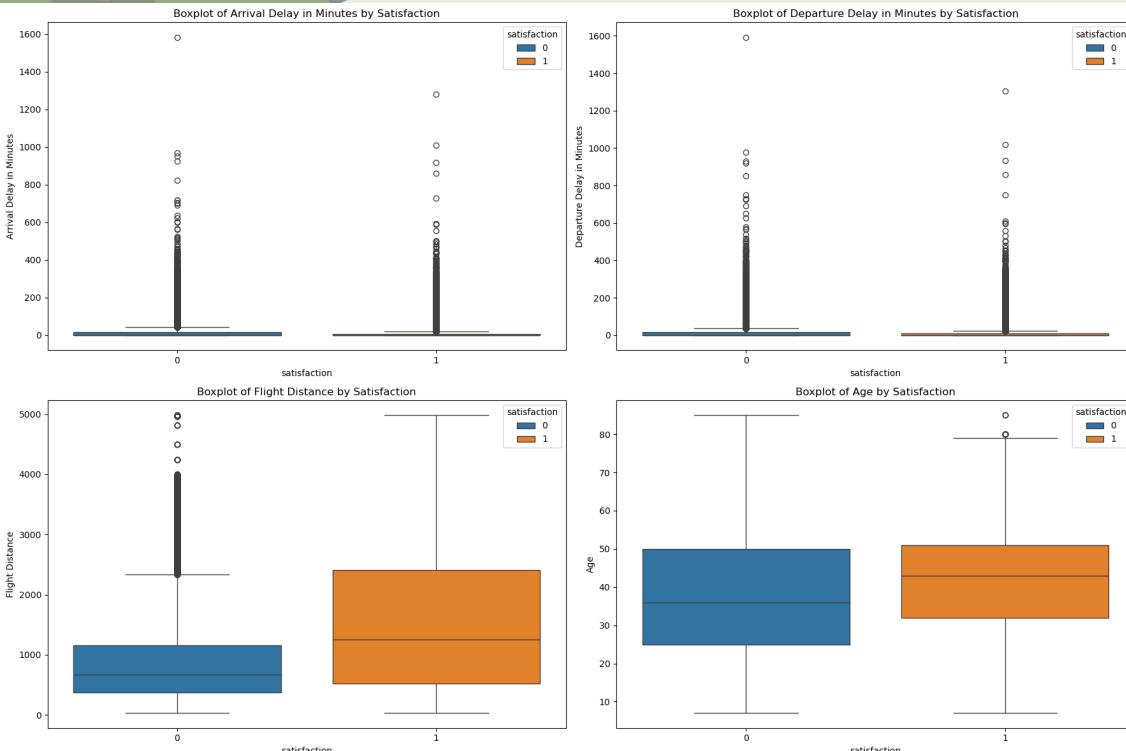
Customer Type	Age	Type of Travel	Class	Flight Distance	Inflight wifi service	Departure/Arrival time convenient	Ease of Online booking	Gate location	Food and drink	...	On-board service	Leg room service	Baggage handling	Checkin service	Inflight service	Cleanliness	Departure Delay in Minutes	Arrival Delay in Minutes	satisfaction	Gender_Male	
0	0	13	1	2	460	3	4	3	1	5	...	4	3	4	4	5	5	25	18.0	0	1.0
1	1	25	0	0	235	3	2	3	3	1	...	1	5	3	1	4	1	1	6.0	0	1.0
2	0	26	0	0	1142	2	2	2	2	5	...	4	3	4	4	4	5	0	0.0	1	0.0
3	0	25	0	0	562	2	5	5	2	...	2	5	3	1	4	2	11	9.0	0	0	0.0
4	0	61	0	0	214	3	3	3	4	...	3	4	4	3	3	3	0	0.0	1	1	1.0

5 rows × 23 columns

Exploratory Data Analysis

```
1 features = df[['Arrival Delay in Minutes', 'Departure Delay in Minutes', 'Flight Distance', 'Age', 'satisfaction']]
2 plt.figure(figsize=(18, 12))
3 for i, column in enumerate(features.columns[:-1], 1):
4     plt.subplot(2, 2, i)
5     sns.boxplot(x='satisfaction', y=column, data=features, hue='satisfaction')
6     plt.title(f'Boxplot of {column} by Satisfaction')
7 plt.tight_layout()
8 plt.show()

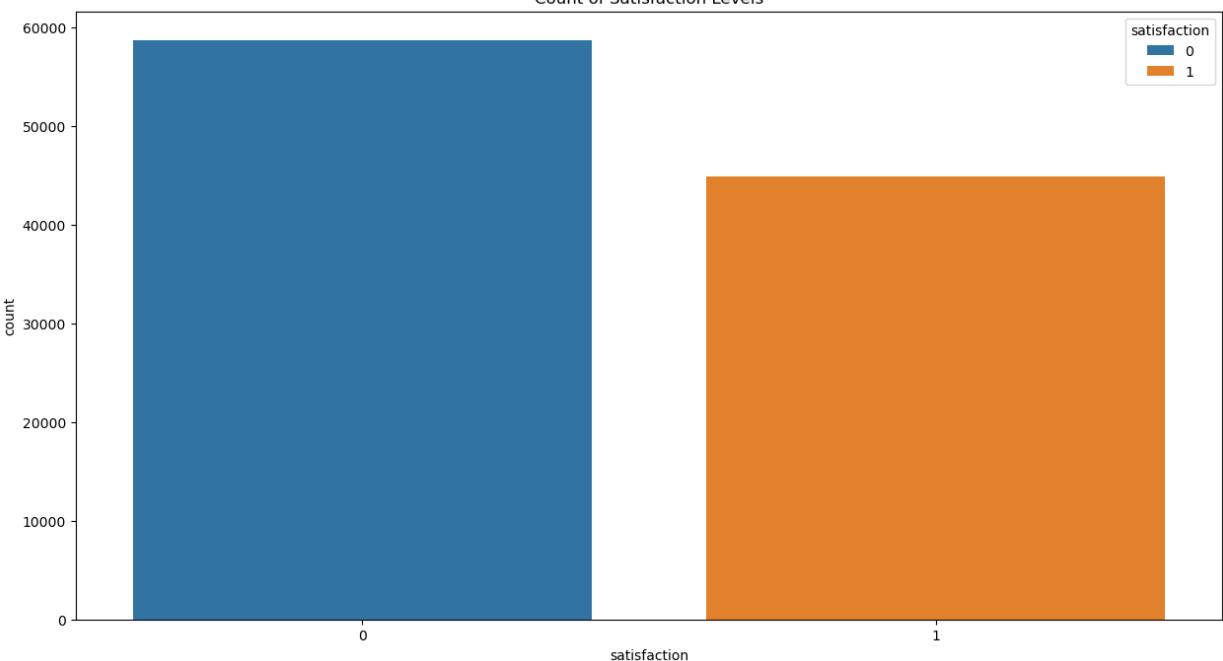
9
10 plt.figure(figsize=(18, 12))
11 for i, column in enumerate(features.columns[:-1], 1):
12     plt.subplot(2, 2, i)
13     sns.violinplot(x='satisfaction', y=column, data=features, hue='satisfaction')
14     plt.title(f'Violin Plot of {column} by Satisfaction')
15 plt.tight_layout()
16 plt.show()
17
```



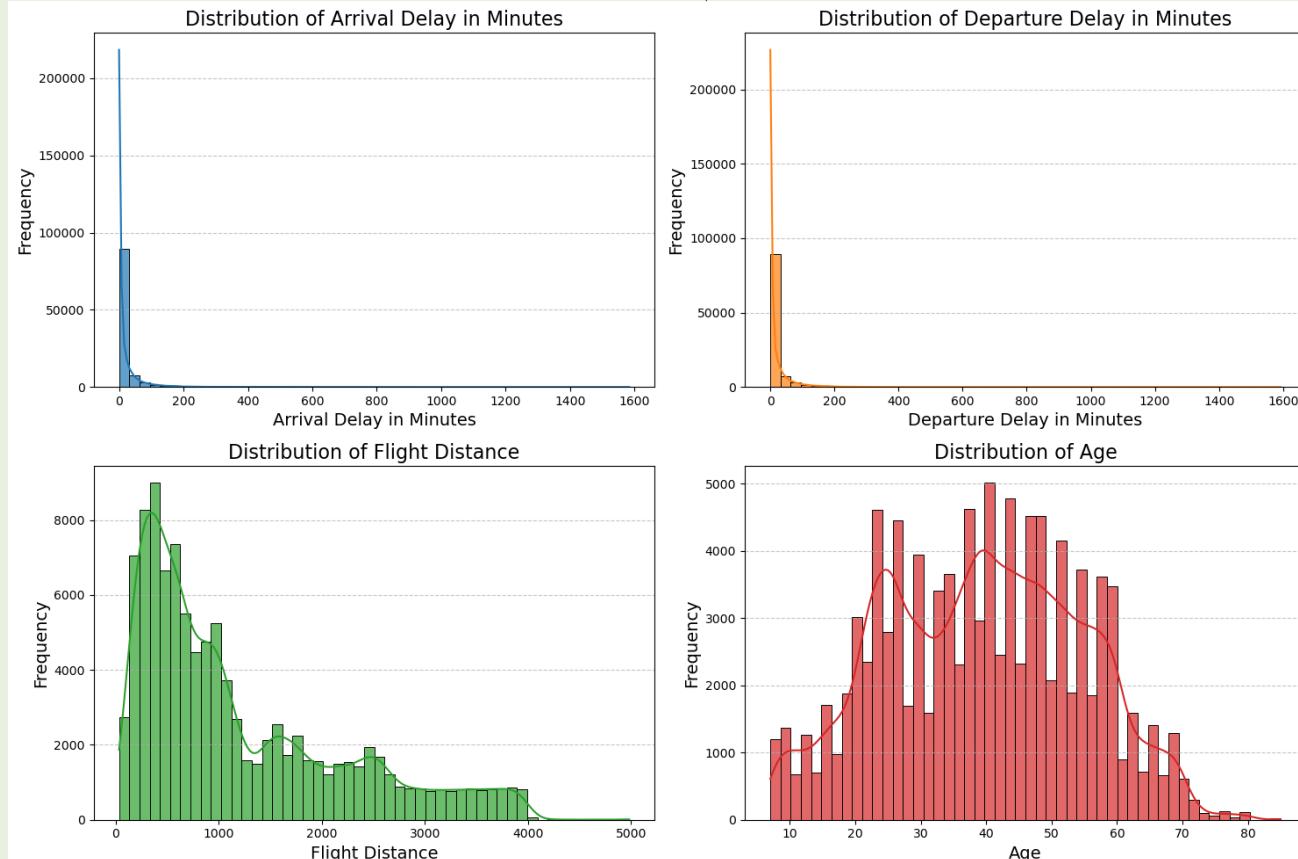
Distribution Analysis



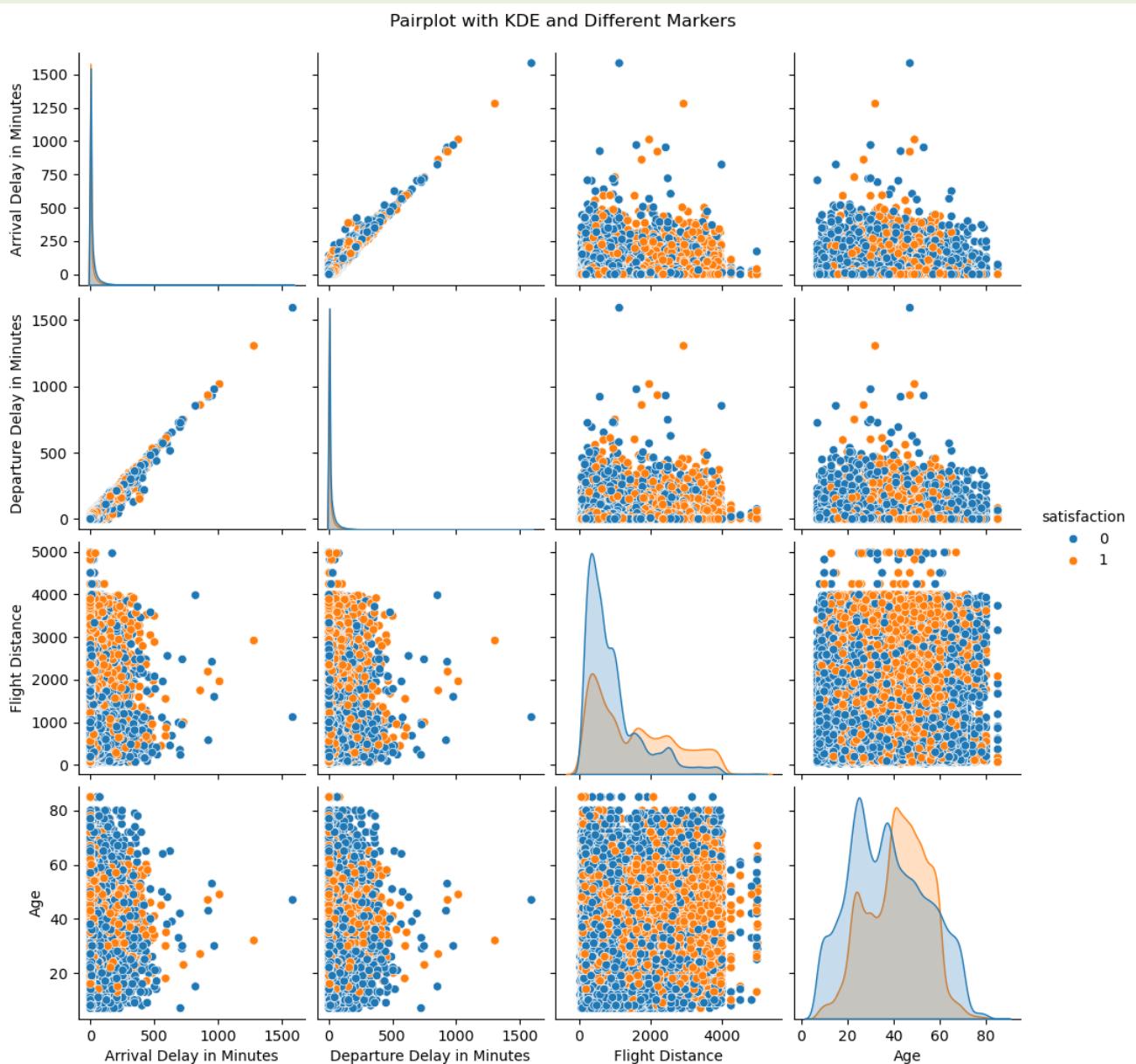
```
1 plt.figure(figsize=(15, 8))
2 sns.countplot(x='satisfaction', data=features, hue='satisfaction')
3 plt.title('Count of Satisfaction Levels')
4 plt.show()
5
```



```
1 palette = sns.color_palette(n_colors=len(features.columns)-1)
2
3 plt.figure(figsize=(15, 10))
4 for i, column in enumerate(features.columns[:-1], 1):
5     plt.subplot(2, 2, i)
6     sns.histplot(features[column], kde=True, bins=50, color=palette[i-1], alpha=0.7)
7     plt.title(f'Distribution of {column}', fontsize=16)
8     plt.xlabel(column, fontsize=14)
9     plt.ylabel('Frequency', fontsize=14)
10    plt.grid(axis='y', linestyle='--', alpha=0.7)
11
12 plt.tight_layout()
13 plt.show()
14
```



Distribution Analysis

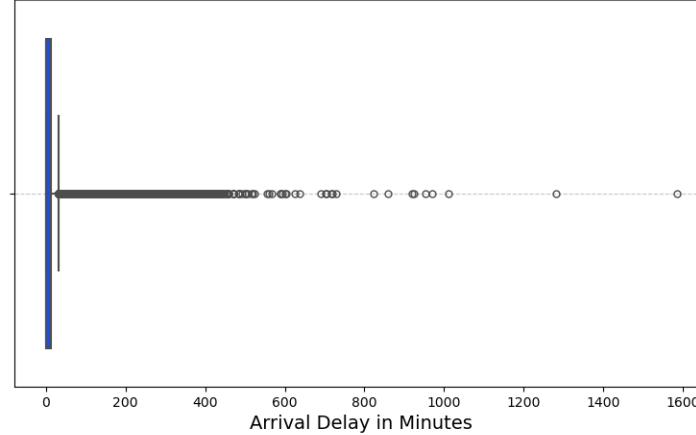


```
1 sns.pairplot(features, hue='satisfaction', diag_kind='kde', markers=[ "o", "o"])
2 plt.suptitle('Pairplot with KDE and Different Markers', y=1.02)
3 plt.show()
```

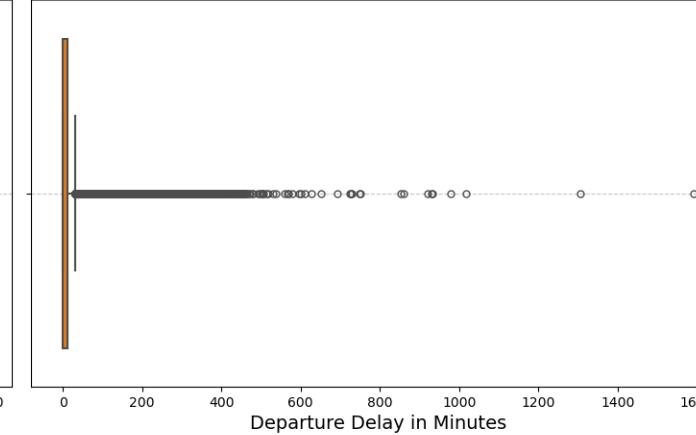
Pair Plot Analysis

Distribution Analysis

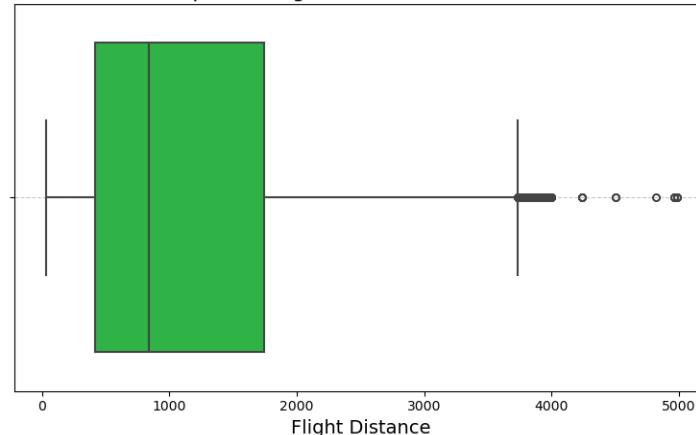
Boxplot of Arrival Delay in Minutes with Outliers



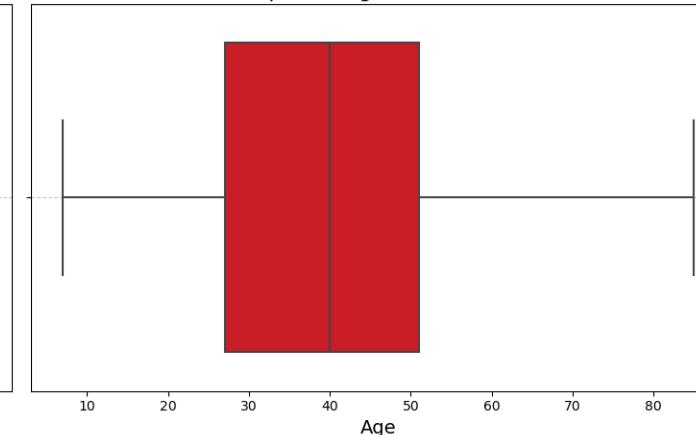
Boxplot of Departure Delay in Minutes with Outliers



Boxplot of Flight Distance with Outliers



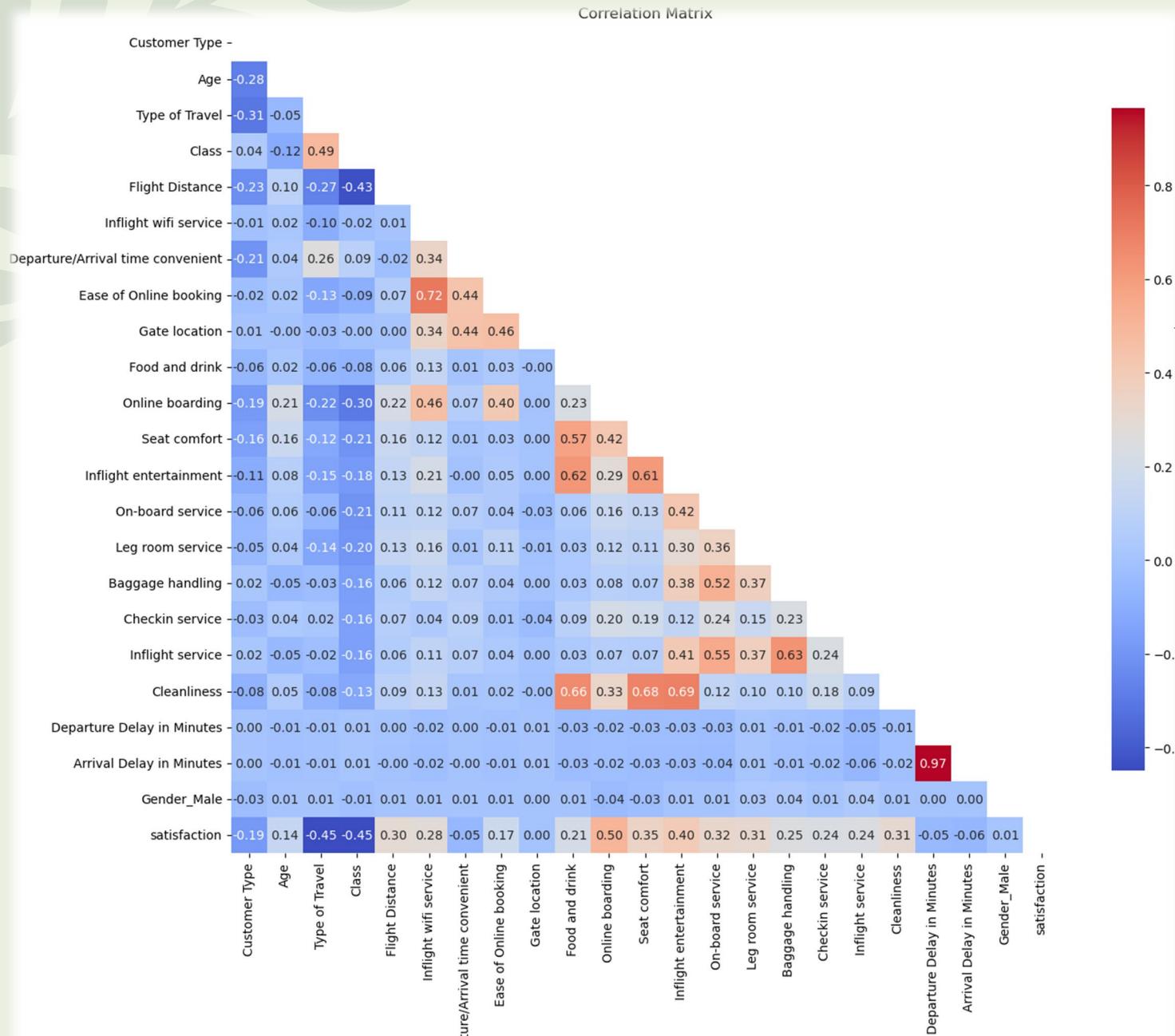
Boxplot of Age with Outliers



```
1 features = df[['Arrival Delay in Minutes', 'Departure Delay in Minutes', 'Flight Distance',
2 'Age', 'satisfaction']]
3
4 palette = sns.color_palette('bright', n_colors=len(features.columns)-1)
5
6 Q1 = features.quantile(0.25)
7 Q3 = features.quantile(0.75)
8 IQR = Q3 - Q1
9
10 outlier_condition = (features < (Q1 - 1.5 * IQR)) | (features > (Q3 + 1.5 * IQR))
11
12 plt.figure(figsize=(15, 10))
13 for i, column in enumerate(features.columns[:-1], 1):
14     plt.subplot(2, 2, i)
15     sns.boxplot(x=features[column], color=palette[i-1], fliersize=5, linewidth=1.5)
16     plt.title(f'Boxplot of {column} with Outliers', fontsize=16)
17     plt.xlabel(column, fontsize=14)
18     plt.grid(axis='y', linestyle='--', alpha=0.7)
19
20 plt.tight_layout()
21 plt.show()
```

Box Plot Analysis for
Outlier Detection

Correlation Analysis

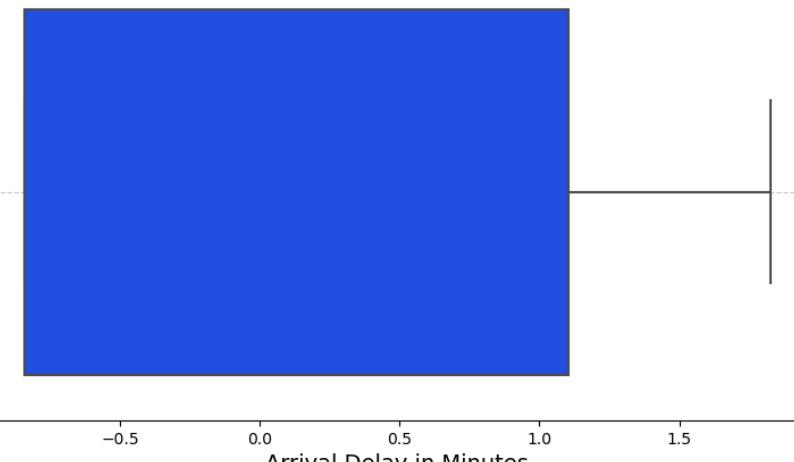


```
1 df = df[[col for col in df.columns if col !=  
'satisfaction'] + ['satisfaction']]  
2  
3 corr = df.corr()  
4 mask = np.triu(np.ones_like(corr, dtype=bool))  
5 plt.figure(figsize=(15, 12))  
6 sns.heatmap(corr, mask=mask, cmap='coolwarm',  
annot=True, fmt=".2f", square=True, cbar_kws=  
{"shrink": .8})  
7 plt.title('Correlation Matrix')  
8 plt.show()
```

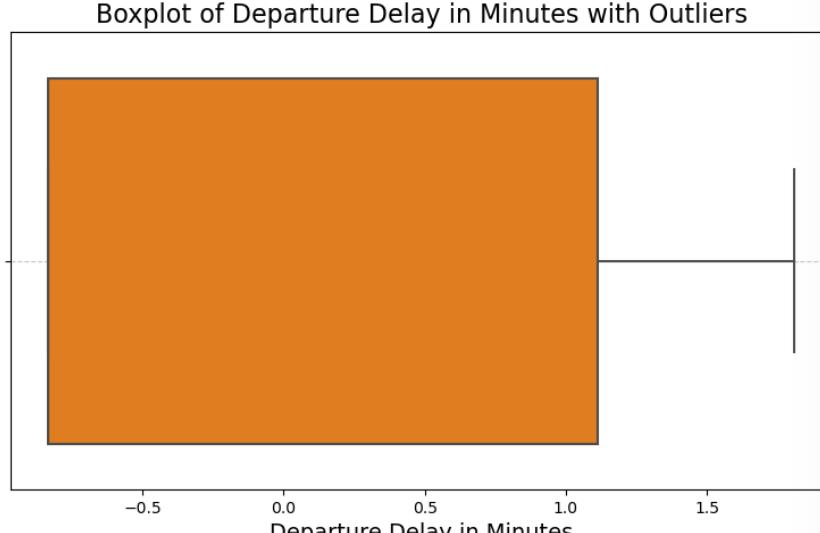
Feature Engineering

Yeo-Johnson Transformation for outlier handling

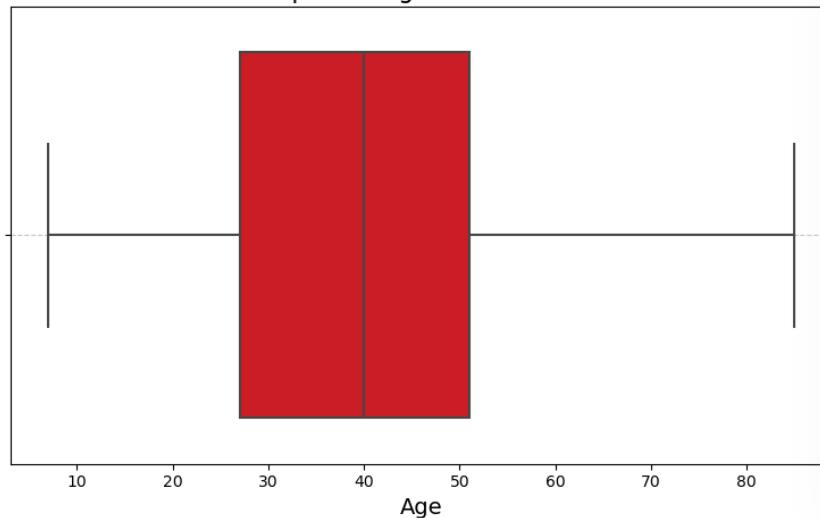
Boxplot of Arrival Delay in Minutes with Outliers



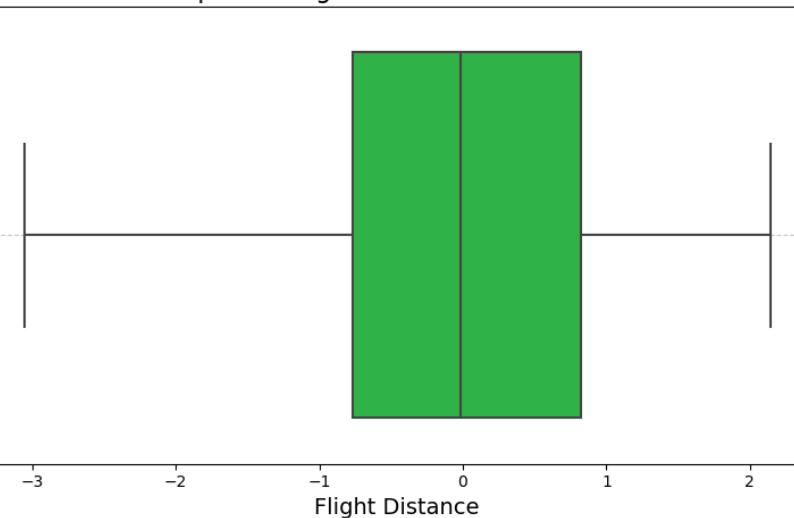
Boxplot of Departure Delay in Minutes with Outliers



Boxplot of Flight Distance with Outliers



Boxplot of Age with Outliers



```
1 features = df[['Arrival Delay in Minutes',
  'Departure Delay in Minutes', 'Flight Distance',
  'Age', 'satisfaction']]
2
3
4 palette = sns.color_palette('bright', n_colors=len(features.columns)-1)
5
6 Q1 = features.quantile(0.25)
7 Q3 = features.quantile(0.75)
8 IQR = Q3 - Q1
9
10 outlier_condition = (features < (Q1 - 1.5 * IQR)) | (features > (Q3 + 1.5 * IQR))
11
12 plt.figure(figsize=(15, 10))
13 for i, column in enumerate(features.columns[:-1], 1):
14     plt.subplot(2, 2, i)
15     sns.boxplot(x=features[column], color=palette[i-1], fliersize=5, linewidth=1.5)
16     plt.title(f'Boxplot of {column} with Outliers', fontsize=16)
17     plt.xlabel(column, fontsize=14)
18     plt.grid(axis='y', linestyle='--', alpha=0.7)
19
20 plt.tight_layout()
21 plt.show()
```

Model Training

Baseline Model Performance

Baseline Model Accuracy: 0.8759

Baseline Model ROC AUC: 0.9253

Classification Report:

	precision	recall	f1-score	support
0	0.88	0.91	0.89	11740
1	0.87	0.84	0.85	8979
accuracy			0.88	20719
macro avg	0.88	0.87	0.87	20719
weighted avg	0.88	0.88	0.88	20719

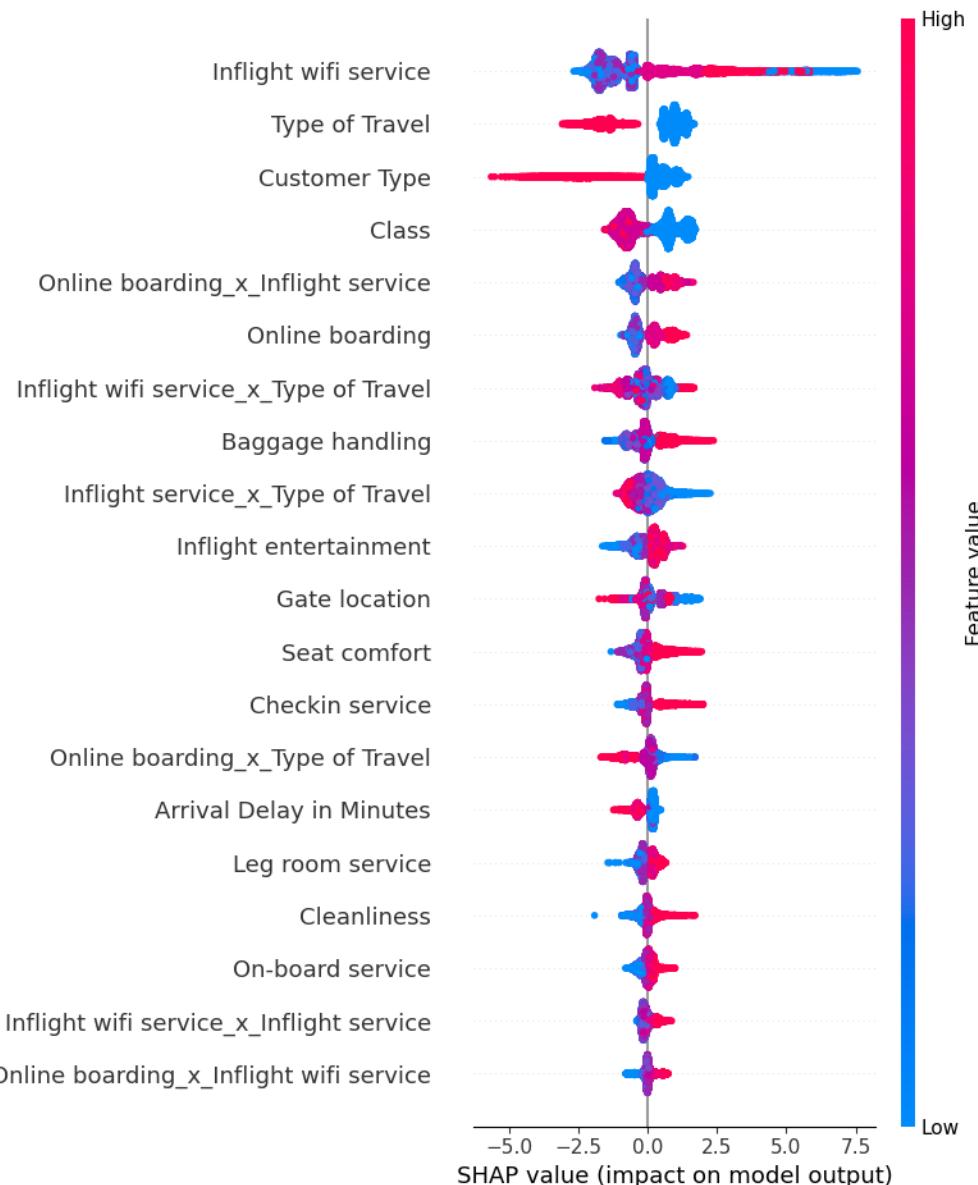


```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.preprocessing import StandardScaler
4 from sklearn.metrics import accuracy_score, classification_report, roc_auc_score
5
6 X = df.drop(columns=['satisfaction'])
7 y = df['satisfaction']
8
9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42, stratify=y)
10 scaler = StandardScaler()
11 X_train = scaler.fit_transform(X_train)
12 X_test = scaler.transform(X_test)
13
14 baseline_model = LogisticRegression(random_state=42, max_iter=100000)
15 baseline_model.fit(X_train, y_train)
16
17 y_pred = baseline_model.predict(X_test)
18
19 accuracy = accuracy_score(y_test, y_pred)
20 roc_auc = roc_auc_score(y_test, baseline_model.predict_proba(X_test)[:, 1])
21 report = classification_report(y_test, y_pred)
22
23 print(f"Baseline Model Accuracy: {accuracy:.4f}")
24 print(f"Baseline Model ROC AUC: {roc_auc:.4f}")
25 print("\nClassification Report:")
26 print(report)
```

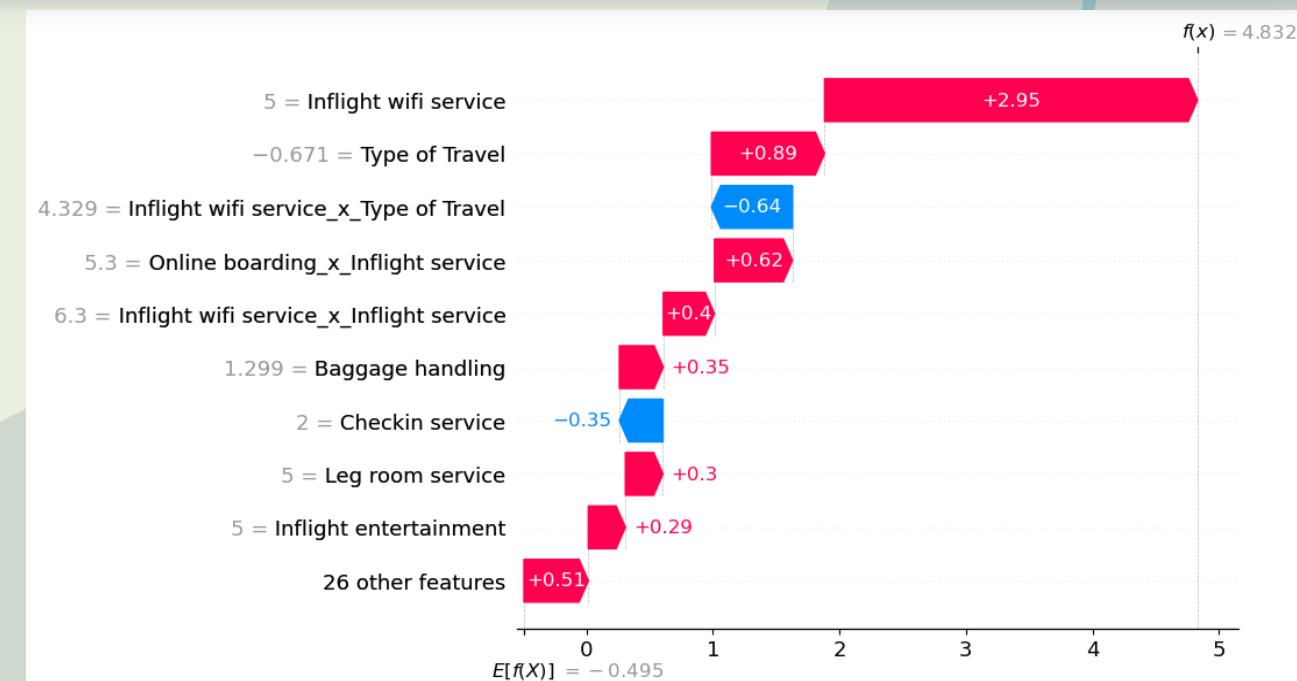
Model Comparison

Model	Accuracy	ROC AUC	R2 Score
CatBoost	0.965346	0.963122	0.858877
XGBoost	0.963801	0.961667	0.852588
LightGBM	0.963077	0.960531	0.849639
Logistic Regression	0.8759	0.9253	0.772895

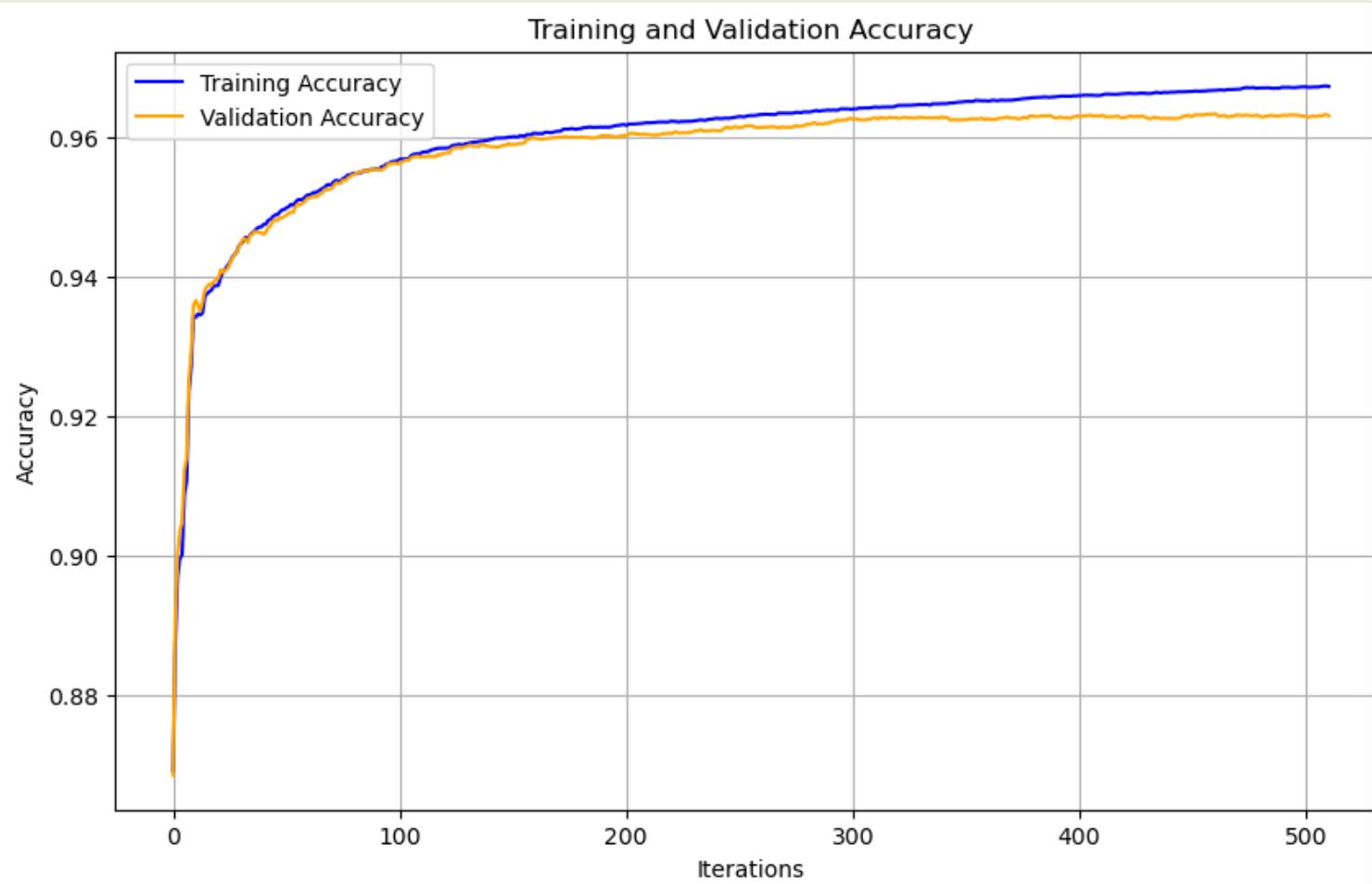
SHAP Explanation



```
1 import shap
2 from catboost import CatBoostClassifier, Pool
3
4 categorical_cols = df.select_dtypes(include=["object", "category"]).columns.tolist()
5
6 for col in categorical_cols:
7     df[col] = df[col].astype(str)
8
9 from sklearn.model_selection import train_test_split
10 X = df.drop(columns=['satisfaction'])
11 y = df['satisfaction']
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
13
14 train_pool = Pool(X_train, label=y_train, cat_features=categorical_cols)
15 test_pool = Pool(X_test, label=y_test, cat_features=categorical_cols)
16
17 model = CatBoostClassifier(iterations=500, learning_rate=0.05, depth=8, verbose=200)
18 model.fit(train_pool)
19
20 explainer = shap.TreeExplainer(model)
21 shap_values = explainer.shap_values(X_test)
22
23 shap.summary_plot(shap_values, X_test)
24
25 shap.plots.waterfall(shap.Explanation(values=shap_values[0], base_values=explainer.expected_value, data=X_test.iloc[0]))
```



Train & Validation Accuracy

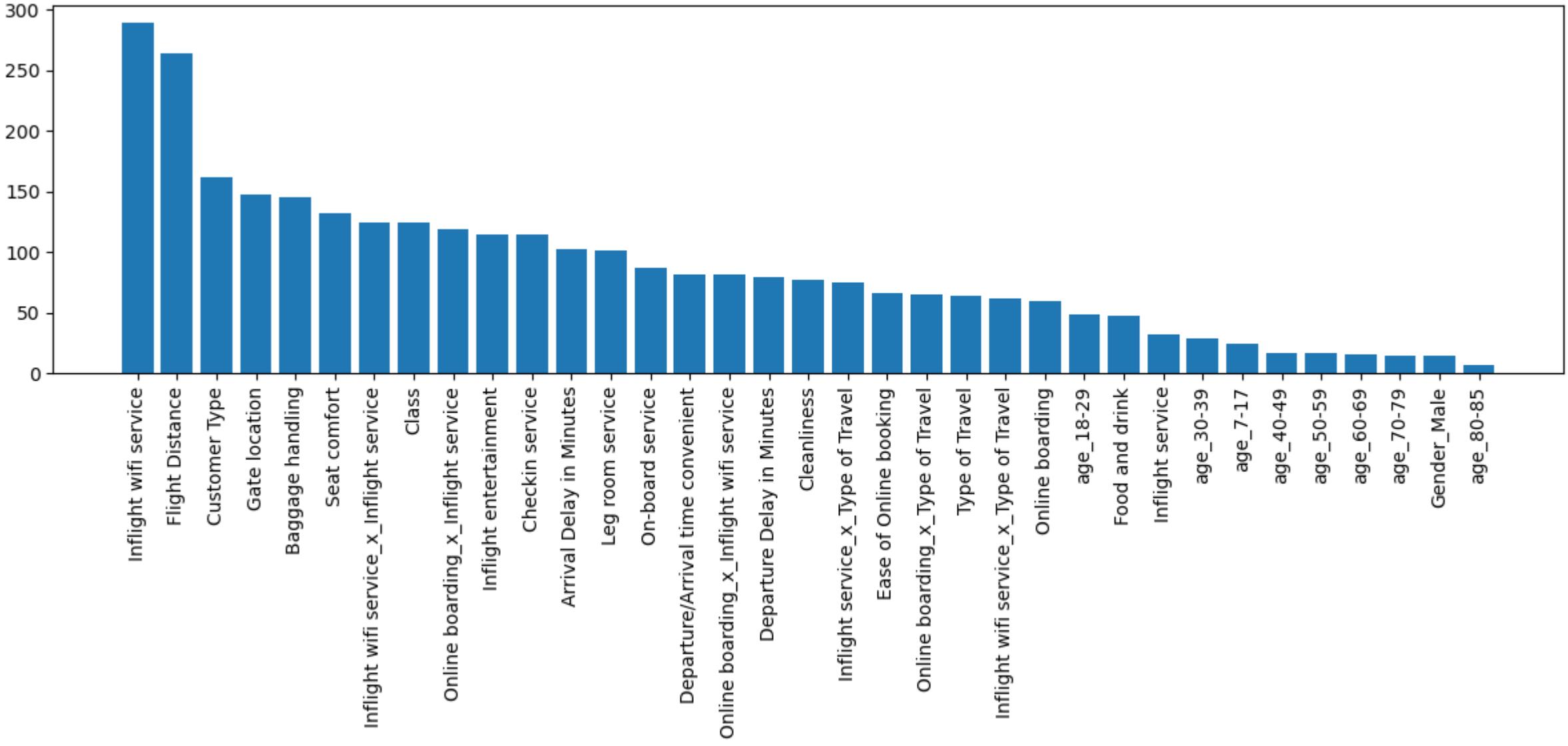


```
1  from catboost import CatBoostClassifier
2  from sklearn.metrics import accuracy_score, classification_report
3  import joblib
4
5  # Parameters found from Optuna tuning
6  cb_params = {'depth': 4, 'learning_rate': 0.17244137403542156, 'l2_le-
af_reg': 3.258820179105591, 'subsample': 0.7977884243064802}
7
8  model = CatBoostClassifier(
9      **cb_params,
10     loss_function="Logloss",
11     eval_metric="Accuracy",
12     random_seed=42,
13     verbose=200
14 )
15
16 model.fit(train_pool, eval_set=test_pool, early_stopping_rounds=50, u-
se_best_model=True)
17
18 evals_result = model.get_evals_result()
19
20 plt.figure(figsize=(10, 6))
21 plt.plot(evals_result['learn']['Accuracy'], label='Training Accurac-
y', color='blue')
22 plt.plot(evals_result['validation']['Accuracy'], label='Validation Ac-
curacy', color='orange')
23 plt.title('Training and Validation Accuracy')
24 plt.xlabel('Iterations')
25 plt.ylabel('Accuracy')
26 plt.legend()
27 plt.grid()
28 plt.show()
29
30 y_pred = model.predict(X_test)
31 print("Accuracy:", accuracy_score(y_test, y_pred))
32 print(classification_report(y_test, y_pred))
33 joblib.dump(model, 'models/cat_model.pkl')
```

After optimising using Optuna

Feature Importance

Feature Importances



Key Insights



Most Influential Features

- **Inflight wifi service:** Highest overall impact
- **Type of Travel:** Significant variability in satisfaction
- **Customer Type:** Strong influence on rating

Impact Dynamics

Positive Influences:

- ✓ Class
- ✓ Online boarding
- ✓ Inflight entertainment

Negative Influences:

- ✓ Arrival Delay (consistently reduces satisfaction)
- ✓ Inconsistent wifi service

Recommendations

1. Prioritize Wifi Quality

1. Consistent, high-speed internet
2. Uniform service across different travel types

2. Minimize Arrival Delays

1. Improve scheduling
2. Enhance operational efficiency

1. Personalized Customer Experience

1. Tailor services to different customer types
2. Segment-specific satisfaction strategies

2. Continuous Improvement

1. Monitor and enhance secondary features
2. Holistic approach to service quality

Md. Rafiquzzaman Rafi
rafiuzzamanrafil00@gmail.com



Thanks