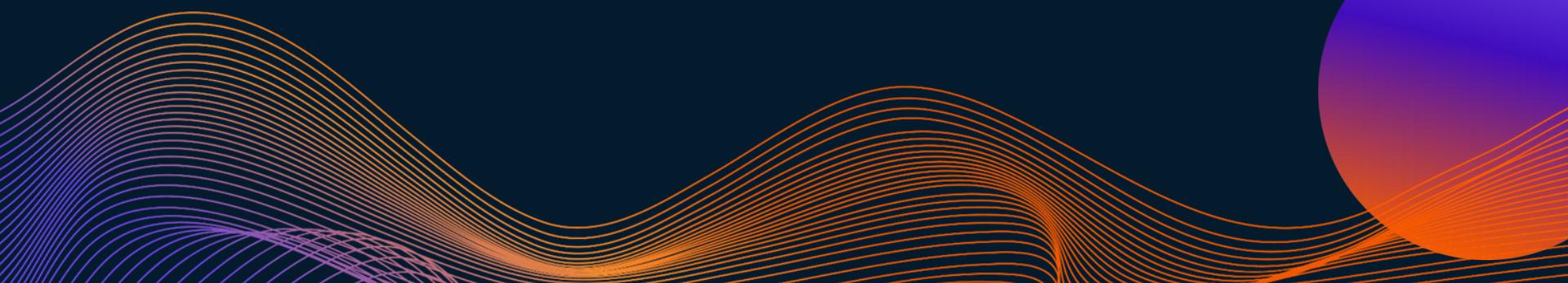




# An In-depth Analysis of Airline Customer Satisfaction



Hector Enriquez , Ifa Afariogun





# Table of contents

**01 Context & EDA**

**02 Data Pre-Processing & Logistic  
Regression/ Decision Tree**

**03 Conclusion & Recommendations**

++  
++

++  
++



# Executive Summary

- 
- 
- 
- 
- 
- 

## Business Problem

Improve customer satisfaction for Invistico Airline by identifying the key factors that influence satisfaction and predicting whether a prospective customer will be satisfied with their service.

## Task

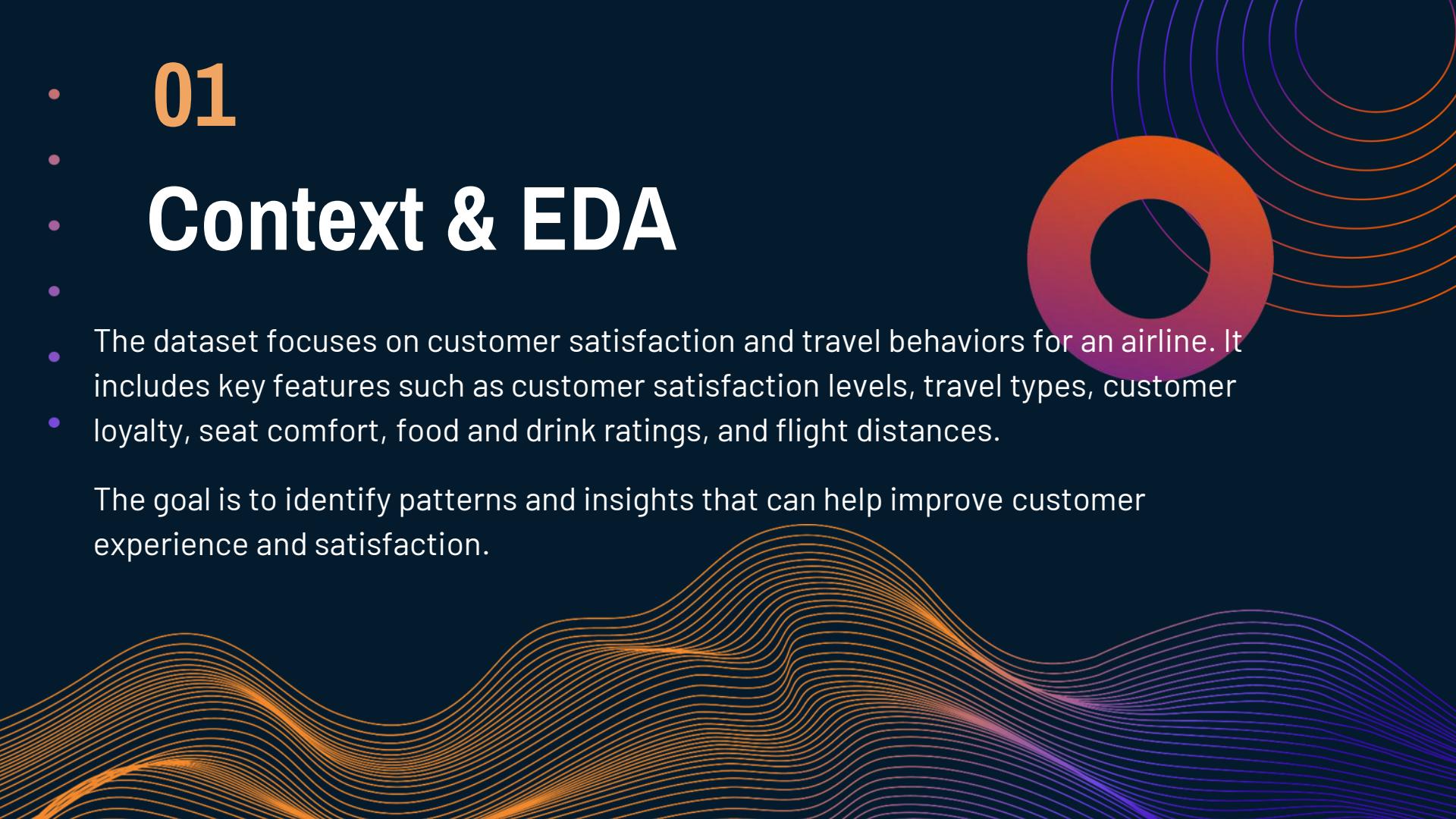
Forecast if a prospective customer will be satisfied and identify service aspects to prioritize for better satisfaction.

# 01

## Context & EDA

- The dataset focuses on customer satisfaction and travel behaviors for an airline. It includes key features such as customer satisfaction levels, travel types, customer loyalty, seat comfort, food and drink ratings, and flight distances.

The goal is to identify patterns and insights that can help improve customer experience and satisfaction.





# Initial Observation of Customer Satisfaction for Invistico Airlines

## What is in this data set?

- Customer details, flight information, and feedback from passengers who have previously flown with Invistico.

**Dataset contains** 129,880 rows and 23 columns. **Captures passenger attributes including:** satisfaction levels, gender, age, flight details, and seat comfort.

#	satisfaction	Gender	Customer Type	Age	Type of Travel	Class	Flight Distance	Seat comfort	Departure/Arrival time convenient	Food and drink	...	Online support	Ease of Online booking	bo...serv
0	satisfied	Female	Loyal Customer	65	Personal Travel	Eco	265	0	0	0	...	2	3	
1	satisfied	Male	Loyal Customer	47	Personal Travel	Business	2464	0	0	0	...	2	3	
2	satisfied	Female	Loyal Customer	15	Personal Travel	Eco	2138	0	0	0	...	2	2	
3	satisfied	Female	Loyal Customer	60	Personal Travel	Eco	623	0	0	0	...	3	1	
4	satisfied	Female	Loyal Customer	70	Personal Travel	Eco	354	0	0	0	...	4	2	

5 rows x 23 columns

# Data Dictionary

Column Name	Description
<code>satisfaction</code>	Customer satisfaction level: <code>satisfied</code> OR <code>neutral</code> OR <code>dissatisfied</code> .
<code>Gender</code>	Gender of the customer: <code>Male</code> OR <code>Female</code> .
<code>Customer Type</code>	Type of customer: <code>Loyal Customer</code> OR <code>Disloyal Customer</code> .
<code>Age</code>	Age of the customer in years.
<code>Type of Travel</code>	Purpose of travel: <code>Personal Travel</code> OR <code>Business Travel</code> .
<code>Class</code>	Travel class: <code>Eco</code> , <code>Business</code> , OR <code>Eco Plus</code> .
<code>Flight Distance</code>	Distance of the flight in kilometers.
<code>Seat comfort</code>	Customer rating for seat comfort (numerical rating, typically 0-5).
<code>Departure/Arrival time convenient</code>	Rating for the convenience of departure and arrival times (0-5).
<code>Food and drink</code>	Customer rating for food and drink provided (0-5).
<code>Gate location</code>	Rating for the gate location (0-5).
<code>Inflight wifi service</code>	Rating for inflight Wi-Fi service (0-5).
<code>Inflight entertainment</code>	Rating for inflight entertainment (0-5).
<code>Online support</code>	Rating for online support services (0-5).
<code>Ease of Online booking</code>	Rating for the ease of online booking (0-5).
<code>On-board service</code>	Customer rating for on-board service quality (0-5).
<code>Leg room service</code>	Rating for leg room comfort (0-5).
<code>Baggage handling</code>	Rating for baggage handling service (0-5).
<code>Checkin service</code>	Rating for the check-in process (0-5).
<code>Cleanliness</code>	Rating for the cleanliness of the airplane (0-5).
<code>Online boarding</code>	Rating for the online boarding process (0-5).
<code>Departure Delay in Minutes</code>	Number of minutes the departure was delayed.
<code>Arrival Delay in Minutes</code>	Number of minutes the arrival was delayed.

# X Addressing Missing Values

```
[]: #checking for missing values  
#We can notice that Arrival delay in minutes has missing values we will handle them shortly.  
data.isnull().sum()
```

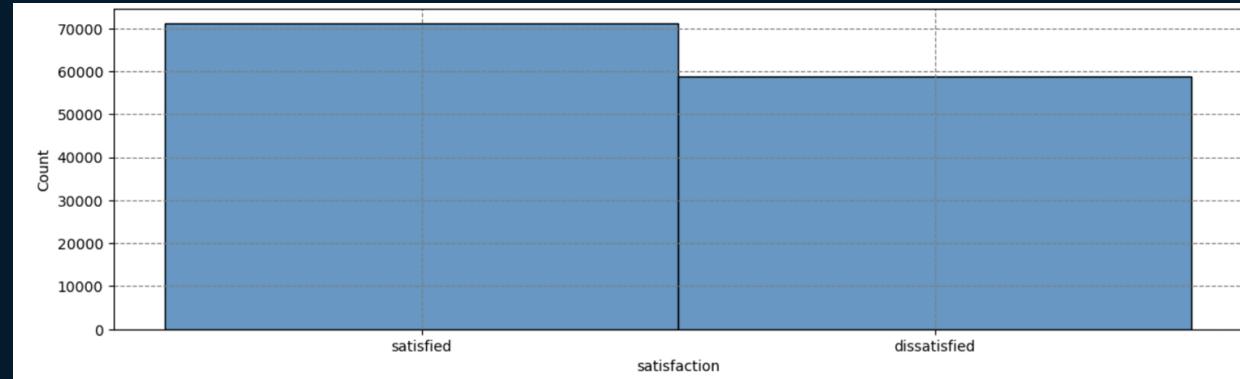
```
[]: satisfaction          0  
Gender                  0  
Customer Type           0  
Age                     0  
Type of Travel          0  
Class                   0  
Flight Distance         0  
Seat comfort             0  
Departure/Arrival time convenient 0  
Food and drink          0  
Gate location            0  
Inflight wifi service   0  
Inflight entertainment   0  
Online support           0  
Ease of Online booking   0  
On-board service         0  
Leg room service         0  
Baggage handling         0  
Checkin service          0  
Cleanliness              0  
Online boarding          0  
Departure Delay in Minutes 0  
Arrival Delay in Minutes 393  
dtype: int64
```



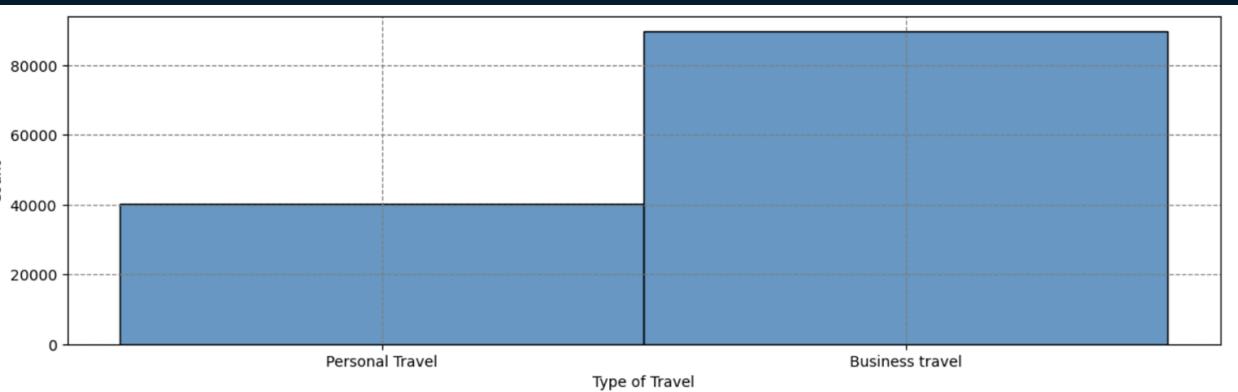
```
Missing_values_after_imputation
```

```
[8]: satisfaction          0  
Gender                  0  
Customer Type           0  
Age                     0  
Type of Travel          0  
Class                   0  
Flight Distance         0  
Seat comfort             0  
Departure/Arrival time convenient 0  
Food and drink          0  
Gate location            0  
Inflight wifi service   0  
Inflight entertainment   0  
Online support           0  
Ease of Online booking   0  
On-board service         0  
Leg room service         0  
Baggage handling         0  
Checkin service          0  
Cleanliness              0  
Online boarding          0  
Departure Delay in Minutes 0  
Arrival Delay in Minutes 0  
dtype: int64
```

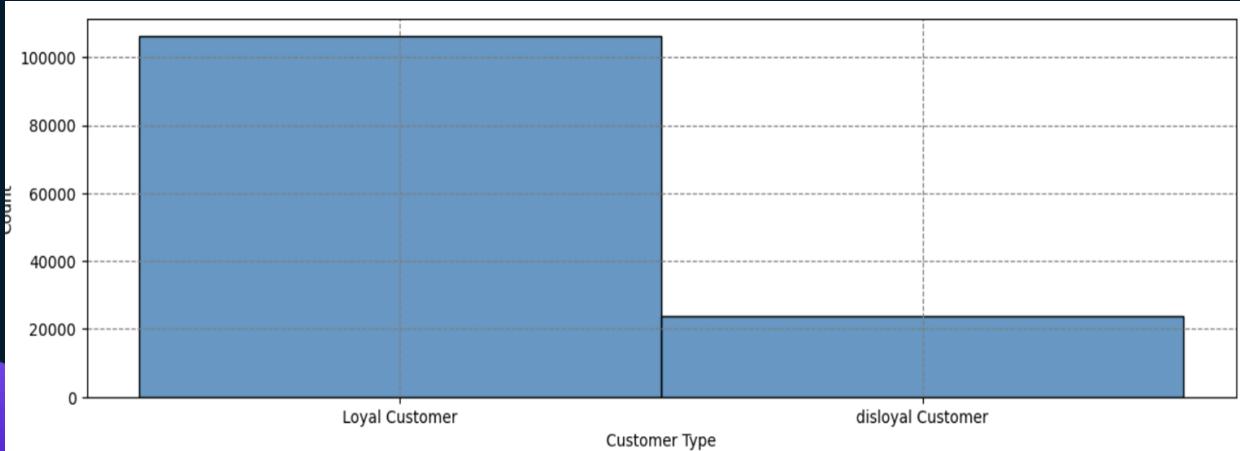
# Exploratory Data Analysis



- **Approach:** Univariate Analysis provided clearer insights into each feature after differentiating and separating the data.



# Exploratory Data Analysis cont.



## Bar Chart Findings:

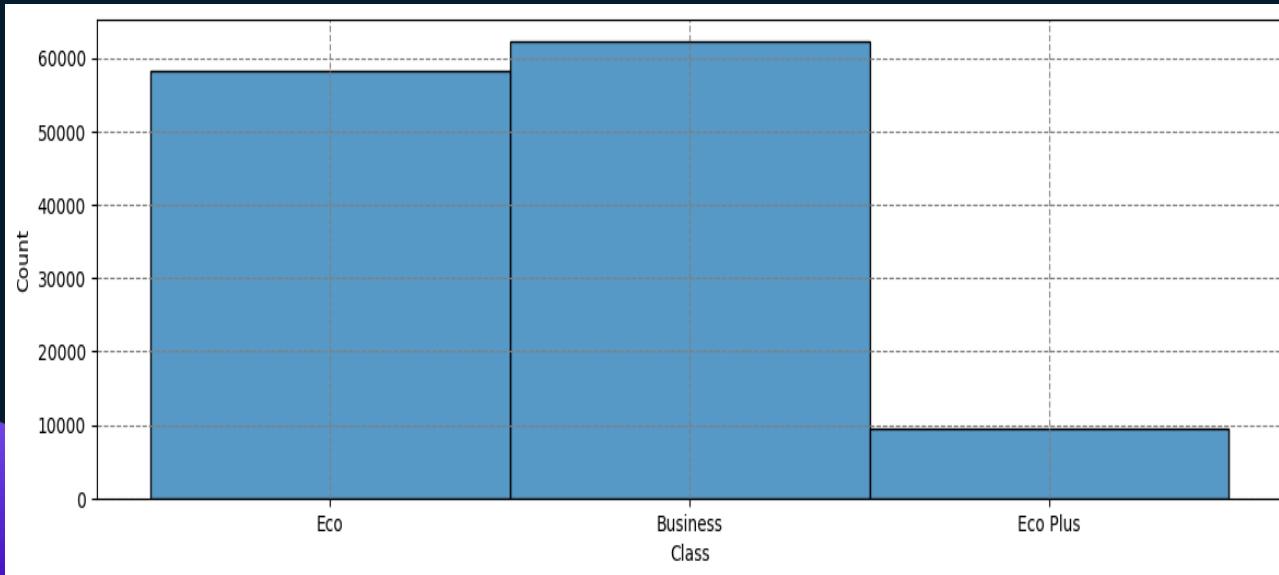
Airline has ~50/50 satisfaction and dissatisfaction rate.

Majority of passengers are loyal customers.

Most passengers are traveling for business purposes.



# Exploratory Data Analysis cont.



## Bar Chart Findings:

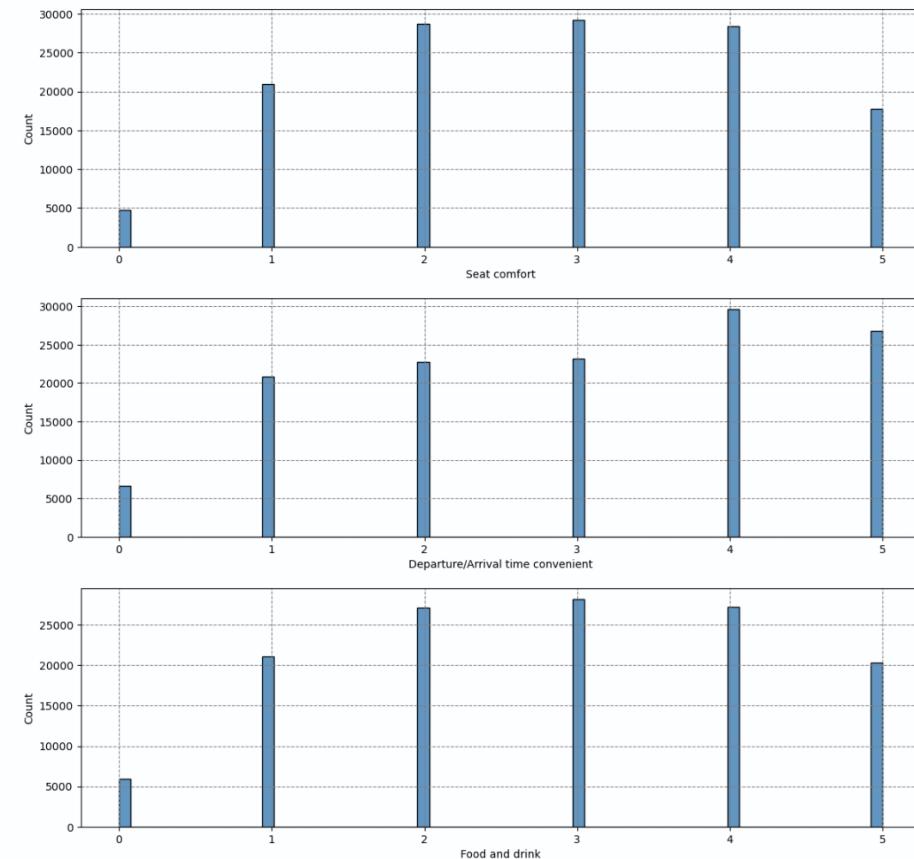
Airline has ~50/50 satisfaction and dissatisfaction rate.

Majority of passengers are loyal customers.

Most passengers are traveling for business purposes.

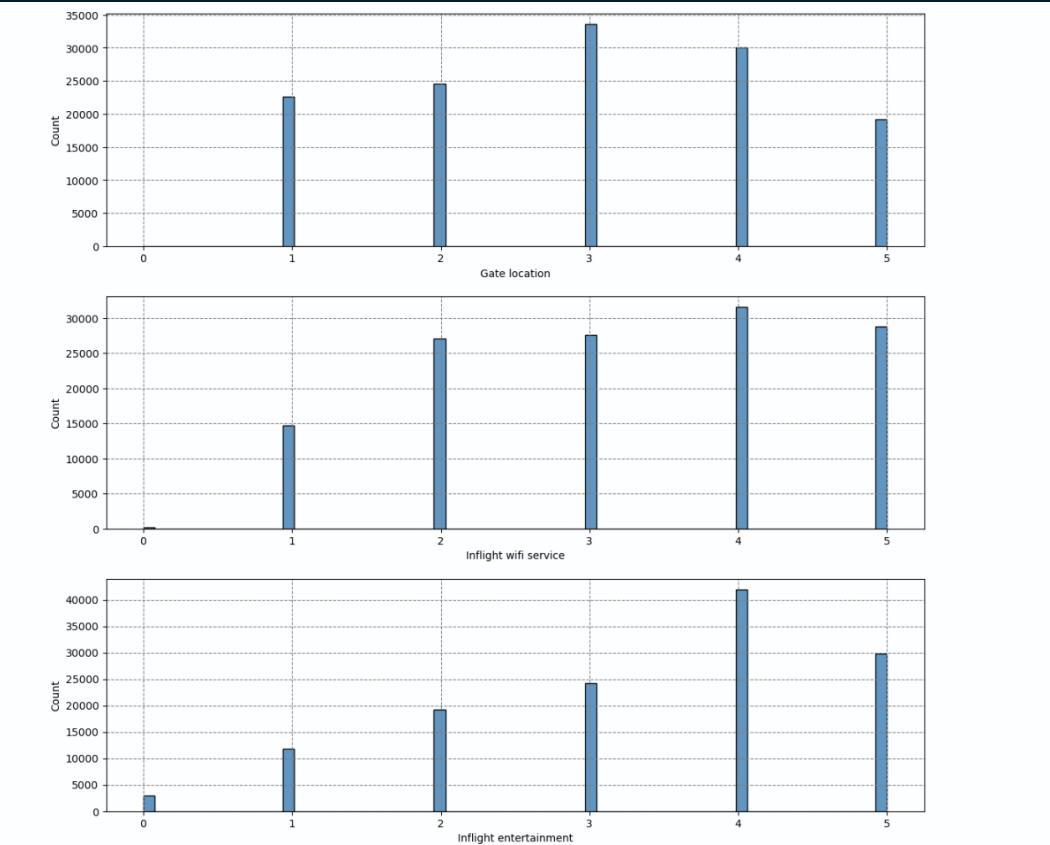


# Exploratory Data Analysis cont.

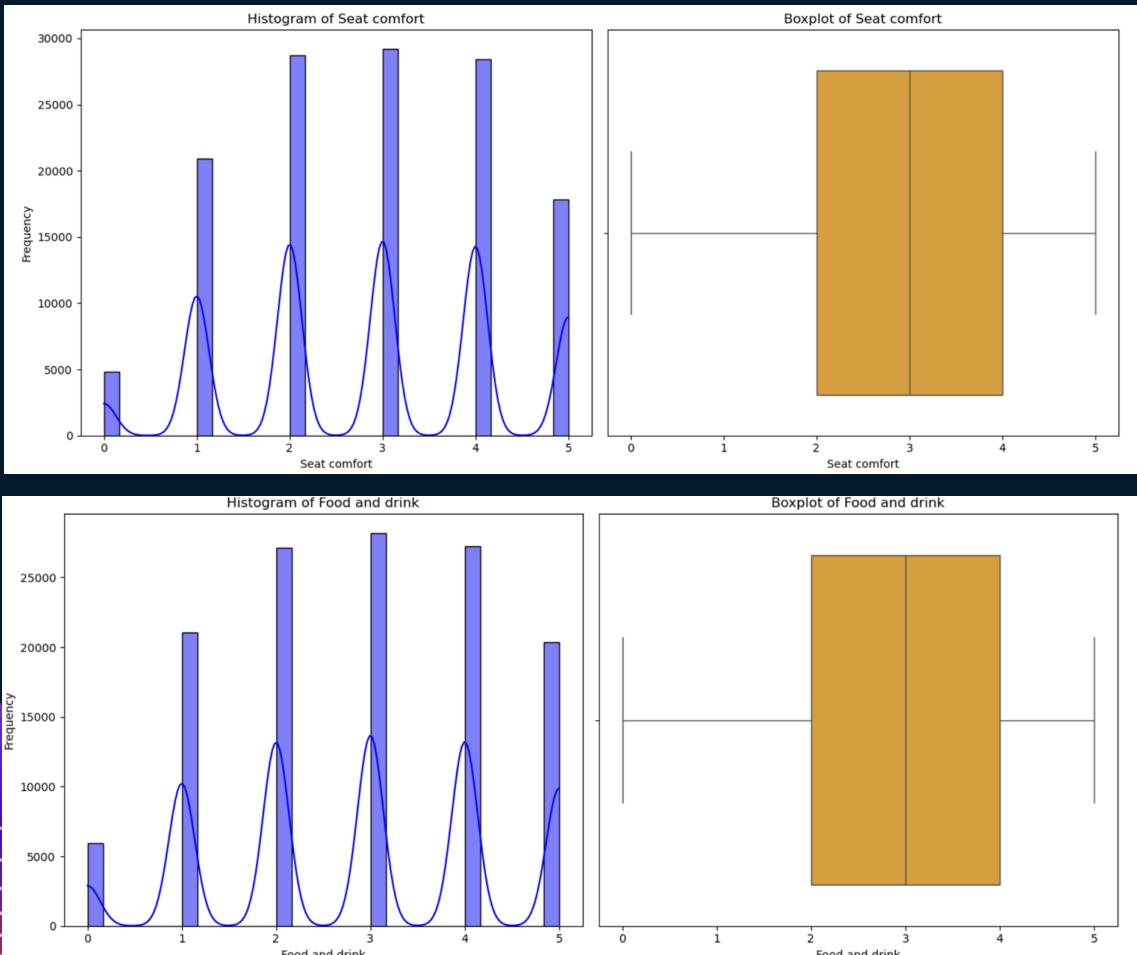


**Satisfaction levels  
relative to seat comfort,  
departure/arrival time,  
and food & drink.**

# Exploratory Data Analysis cont.



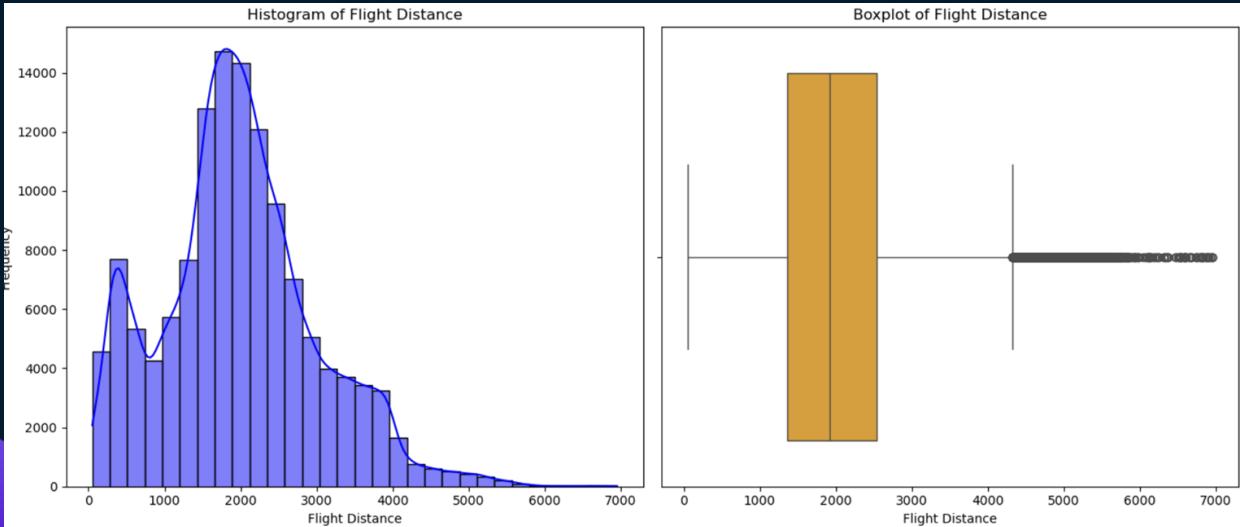
**Satisfaction levels  
relative to gate location,  
inflight wifi service, and  
inflight entertainment.**



## Histogram and Boxplot Findings:

Ratings for seat comfort, and food and drink generally meet or exceed expectations (scores 3 - 5).





The distances traveled on these flights were not far: Shorter flights have simpler amenities and less room for error.

What could this mean?

Shorter flights may have a higher satisfaction rate than longer flights.

# Insights from Univariate Analysis

- **Numerical Variables**

- **Age:**

- Majority of passengers aged 20–50, indicating working-age frequent flyers.
- Symmetric distribution suggests no significant demographic bias.

- **Flight Distance:**

- Right-skewed distribution; most flights are regional or medium-haul (1,000–3,000 units).
- Long-haul passengers may require enhanced amenities (e.g., meals, entertainment).

- **Seat Comfort, Food, and Drink Ratings:**

- Scores cluster around 3–5, reflecting moderate to high satisfaction.
- Low ratings (0 or 1) identify dissatisfaction areas for improvement.

- **Departure/Arrival Delays:**

- Most flights have minor delays (0–50 minutes); extreme delays are rare but impactful.
- Delays critically affect satisfaction; operational issues need attention.

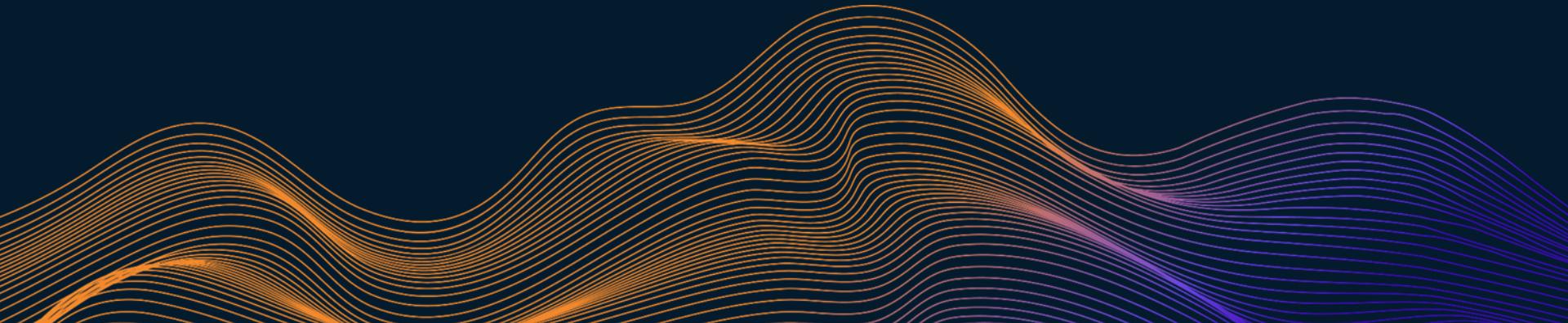
# Insights from Univariate Analysis

- 
- 
- **Categorical Variables**
  - **Gender:**
    - Near-equal male/female distribution; opportunity to tailor services based on satisfaction differences.
  - **Customer Type:**
    - Loyal customers dominate; focus on improving retention by meeting higher expectations.
  - **Type of Travel:**
    - Business travel is prevalent; prioritize timeliness, inflight productivity, and comfort for business flyers.
  - **Class of Travel:**
    - Economy and Business classes are most common; satisfaction may vary by class.
    - Address pain points specific to Economy or Economy Plus passengers.

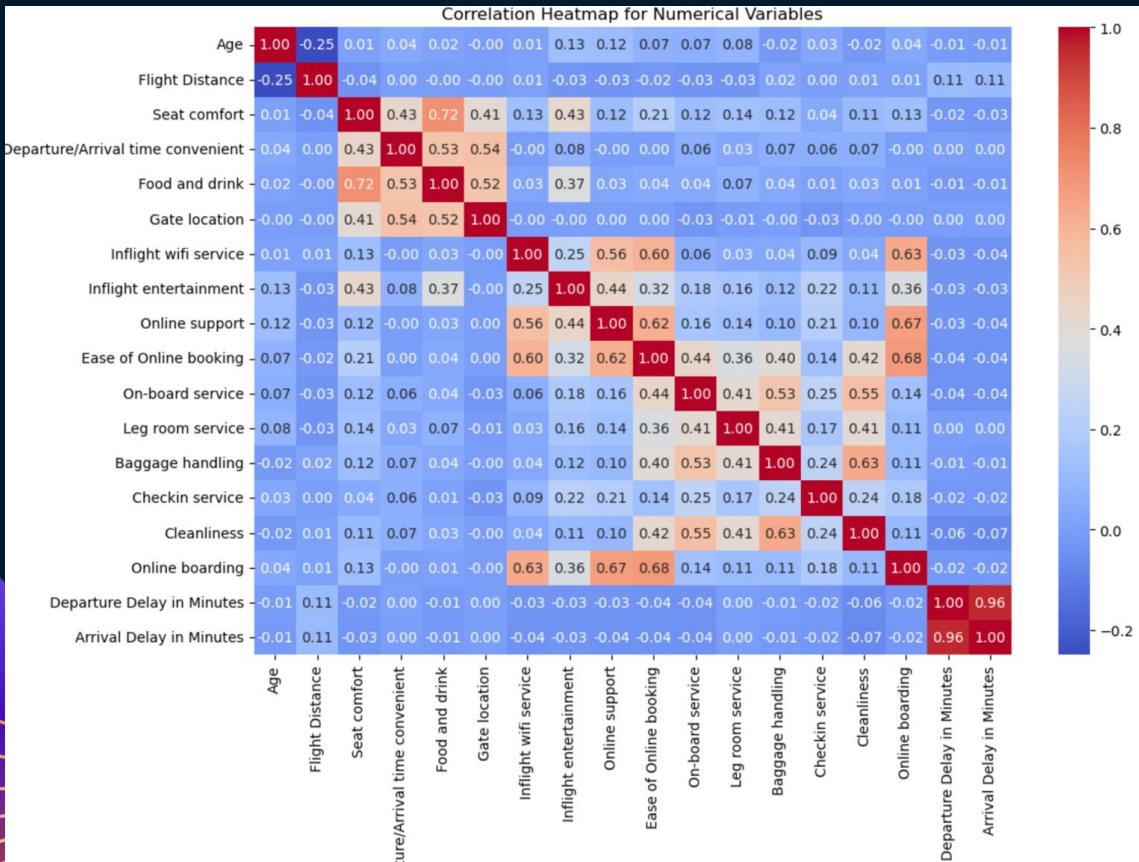
01

# Bivariate Analysis

- Due to the fact that our target variable is binary we conducted stacked bar plots and boxplots that show the mean of the continuous variable across the two categories of target variable.



# Heatmap



## Correlations:

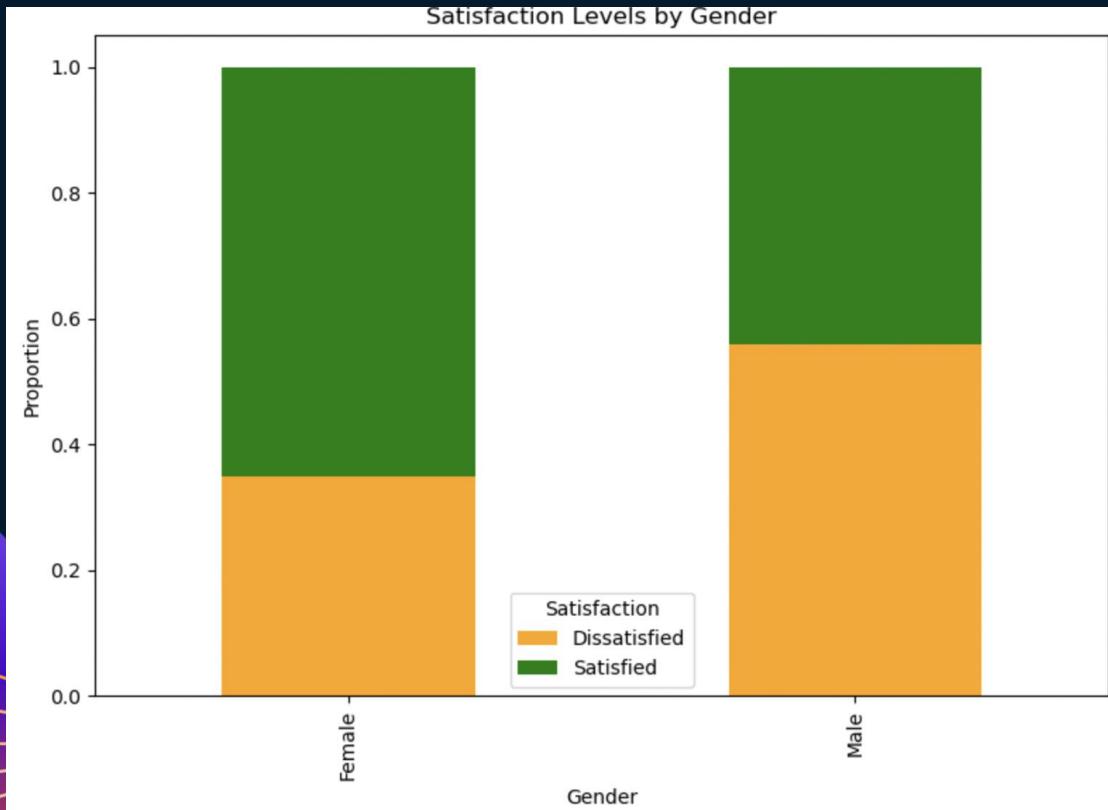
**Ease of online booking and boarding (0.68):**

Related to satisfaction with boarding

**Seat comfort, Food/Drink, Entertainment (0.4-0.7):** Moderate relationships with in-flight satisfaction.

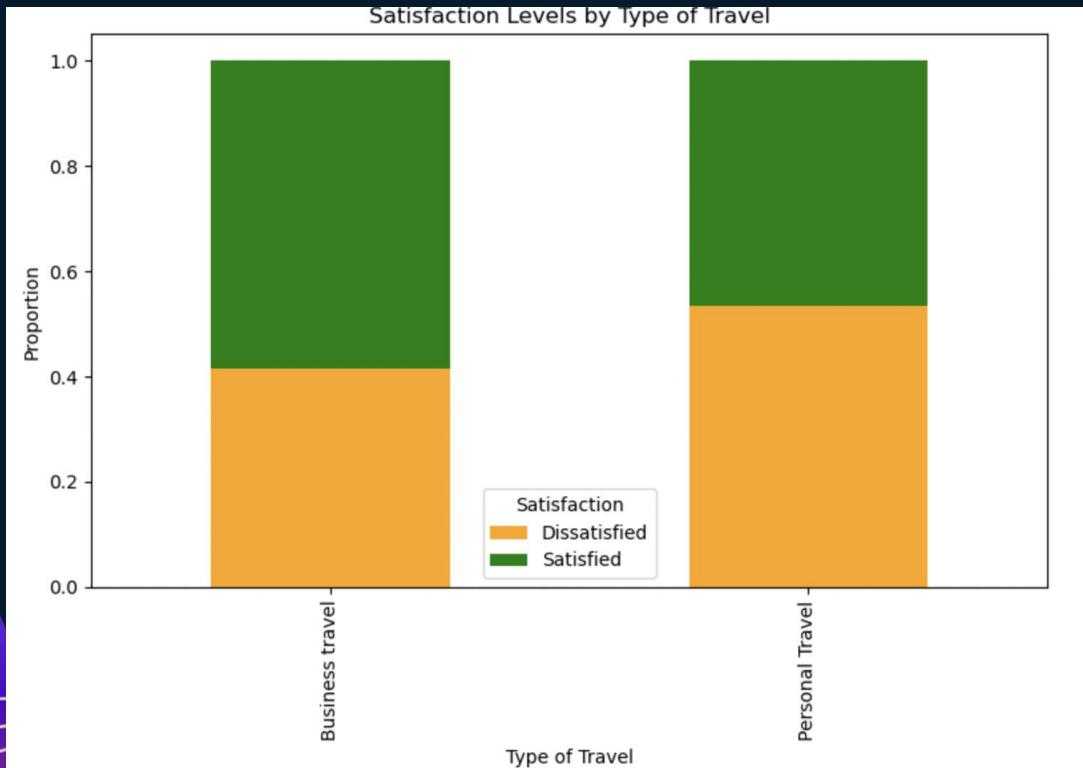
**Departure/Arrival Time Convenience:** Moderately correlated with gate location.

# Satisfaction Levels by Gender



Male passengers were more dissatisfied than female passengers.

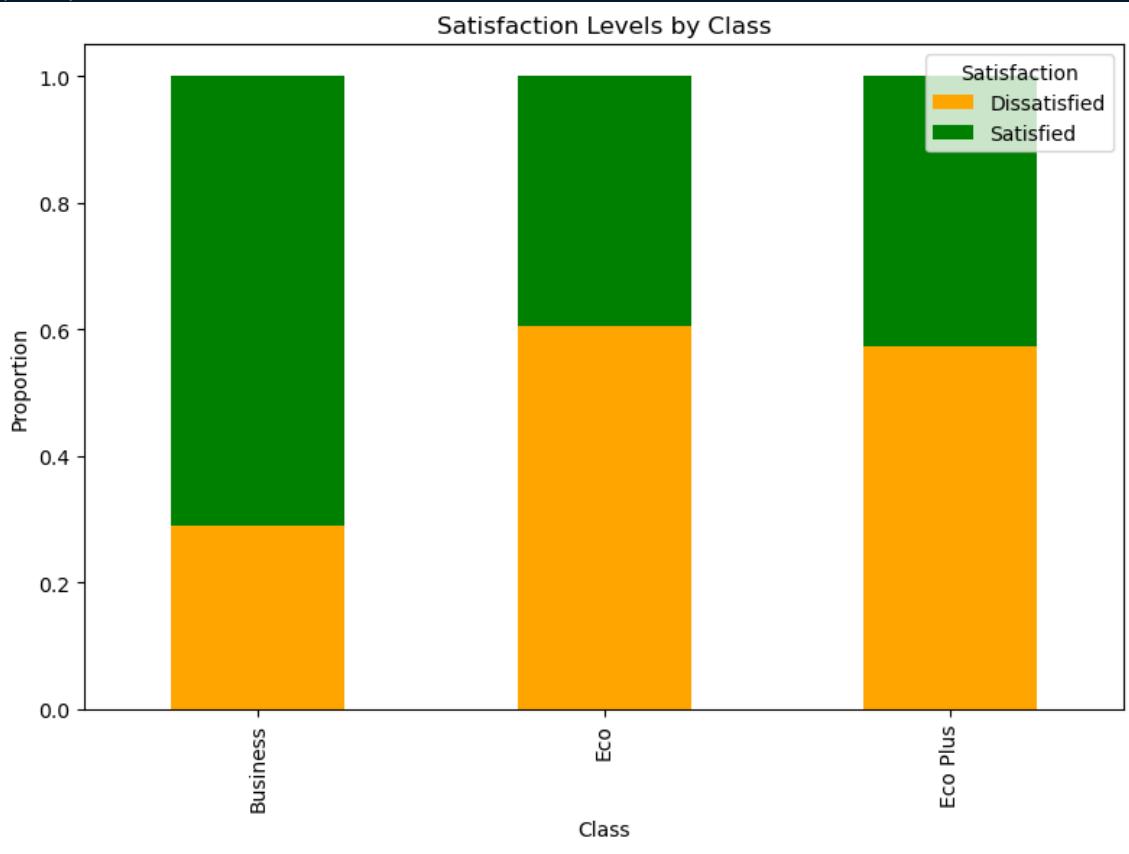
# Satisfaction Levels by Type of Travel



Business travelers were more satisfied than personal travelers.

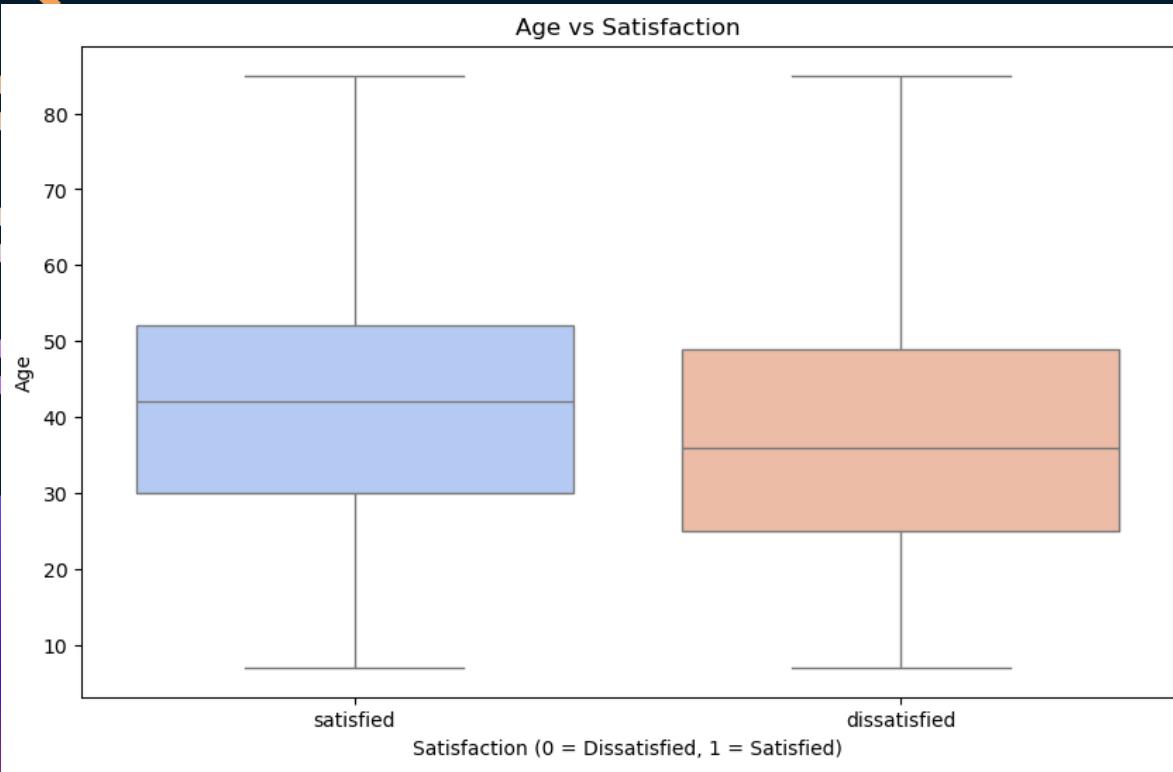
# Satisfaction Levels by Class

X



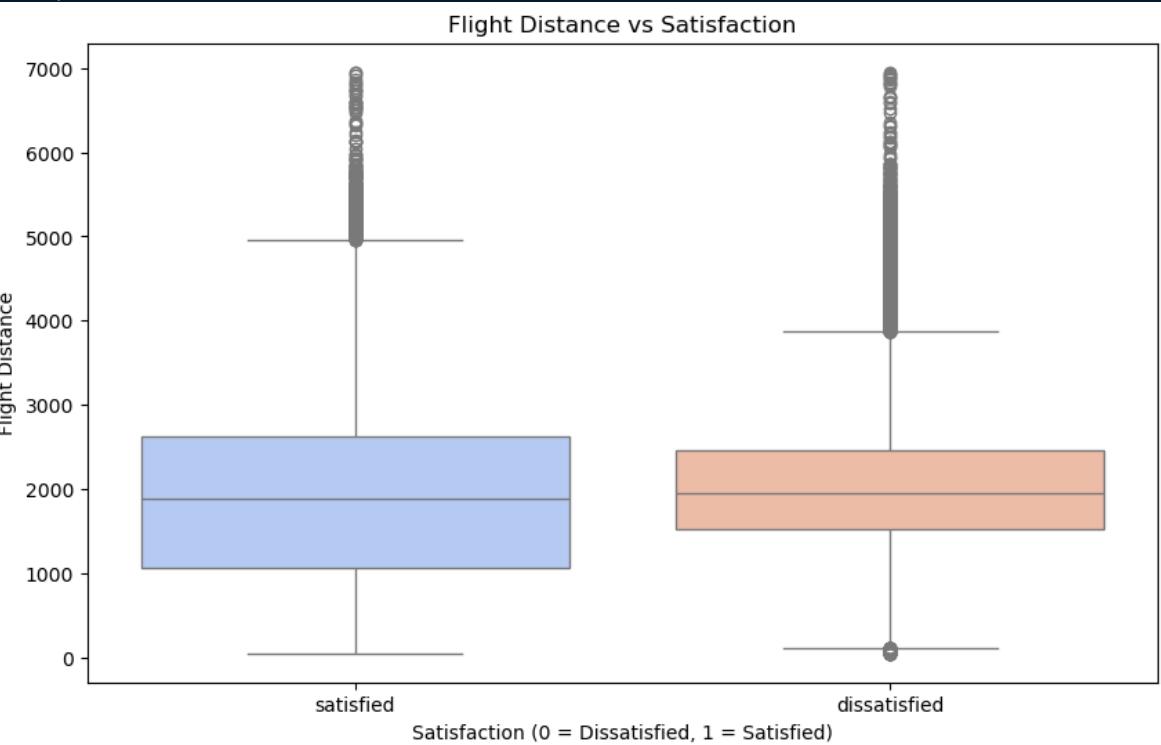
- Business Class passengers have the highest satisfaction proportion.
- Economy and Eco Plus classes exhibit more dissatisfaction.
- Insight: Service improvements (e.g., comfort, inflight experience) in Economy classes could address dissatisfaction.

# Relationship of Age with Satisfaction



- Satisfied customers tend to have a slightly higher median age.
- Younger passengers have more spread in dissatisfaction, suggesting differences in expectations for service.

# Relationship of Flight Distance with Satisfaction



- Satisfied customers typically travel longer distances.
- Dissatisfaction is more common for shorter flights, potentially due to fewer amenities or services.



# • Insights from Bivariate Analysis



- 
- - Higher ratings for seat comfort, legroom, onboard service, and inflight wifi/entertainment strongly correlate with satisfied customers.
  - Dissatisfaction stems from consistently poor ratings across these service areas.
  - 
  - Longer flight distances are associated with higher satisfaction.
  - Shorter flights see more dissatisfaction, possibly due to limited amenities.
  - Satisfied customers have a slightly higher median age.
  - Younger passengers show greater dissatisfaction spread, highlighting potential unmet expectations.

# Insights from Bivariate Analysis

- Delays (departure and arrival) do not significantly differentiate satisfaction.
- Extreme delays (outliers) occur in both satisfied and dissatisfied groups
- Food, cleanliness, gate location, and check-in services receive higher ratings from satisfied customers.
- Dissatisfaction aligns with consistently lower ratings in these areas.



# At A Glance Recommendations

- Inflight Services: Focus on seat comfort, food, and entertainment, especially for short flights.
- Delays: Reduce departure and arrival delays to boost satisfaction.
- Short Flights: Add comfort upgrades for short-distance travelers.
- Young Flyers: Enhance entertainment for younger customers.
- Loyalty: Strengthen loyalty programs to retain satisfied customers.
- Cleanliness & Check-In: Prioritize cleanliness and streamline check-in processes.



02

# Data Preparation & Creation of Models



# Logistic Regression

```
# Convert satisfaction to binary values
data_encoded['satisfaction'] = data_encoded['satisfaction'].apply(lambda x: 1 if x == 'satisfied' else 0)

# We ran running into issues regarding the conversion of the categorical variables into 0 and 1s so I de
# and this was able to Ensure that all dummy variables are converted to numeric (0/1)
data_encoded = data_encoded.apply(lambda col: col.astype(int) if col.dtype == 'bool' else col)

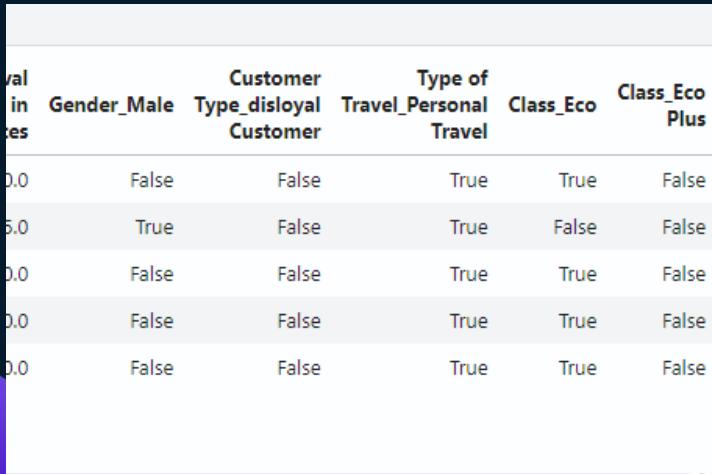
# Verifying that the first few rows of the predictors (X) after corrections
X = data_encoded.drop('satisfaction', axis=1)
y = data_encoded['satisfaction']

X.head()
```

**Approach:** Logistic Regression model can predict customer satisfaction (target variable:satisfaction)

- **Categorical variables** ("Gender," "Customer Type," etc.) are transformed into dummy variables.
- Dataset is divided into training and testing sets, with the *stratify* parameter to ensure the target variable (satisfaction) is equal in both groups.

**Evaluation Metrics:** Accuracy, Recall, Precision, and F1 score analyzed.  
Confusion matrix used for visualizing true/false positives and negatives.



val es	Gender_Male	Customer Customer	Type_of Travel	Class_Eco	Class_Eco Plus
0.0	0	0	1	1	0
0.0	1	0	1	0	0
0.0	0	0	1	1	0
0.0	0	0	1	1	0
0.0	0	0	1	1	0

# Observations

## Training performance:

	Threshold	Accuracy	Recall	Precision	F1
0	0.5	0.786	0.85352	0.77732	0.81364

## Testing performance:

	Threshold	Accuracy	Recall	Precision	F1
0	0.5	0.786624	0.854966	0.777394	0.814337

### Training Results

- Accuracy: 78.6%
- Recall: 85.3%
- Precision: 77.7%
- F1 Score: 81.4%

These high numbers indicate the logistic regression model's strong performance in predicting passenger satisfaction.

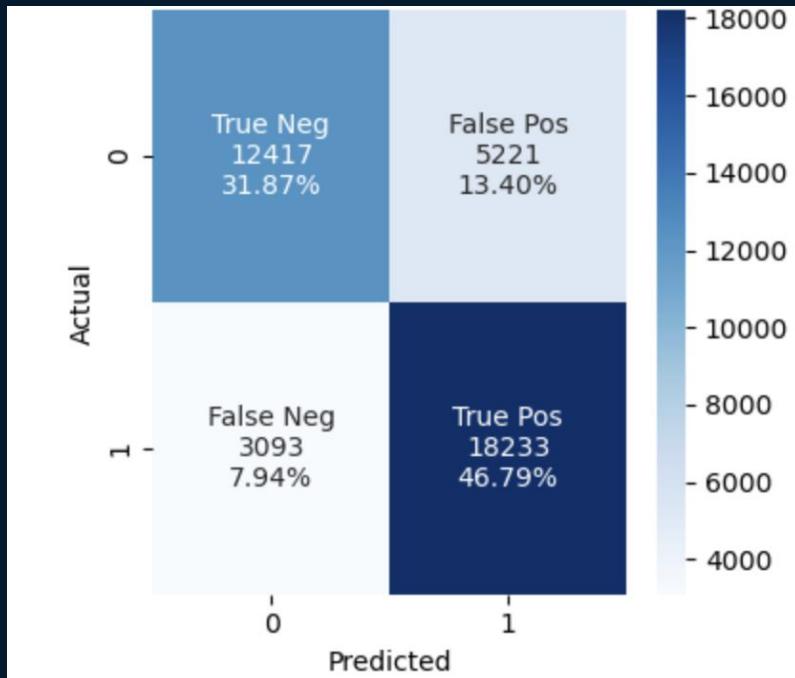
### Testing Results

- Accuracy: 78.6%
- Recall: 85.4%
- Precision: 77.7%
- F1 Score: 81.4%

Consistent performance, with slight increase to recall.

**Model is not overfitting.**

# Confusion Matrix



**True Positives: 46.79%**

Correctly classified satisfied passengers.

**True Negatives: 31.87%**

Correctly classified dissatisfied passengers.

**False Positives: 13.40%**

Dissatisfied passengers misclassified as satisfied.

**False Negatives: 7.94%**

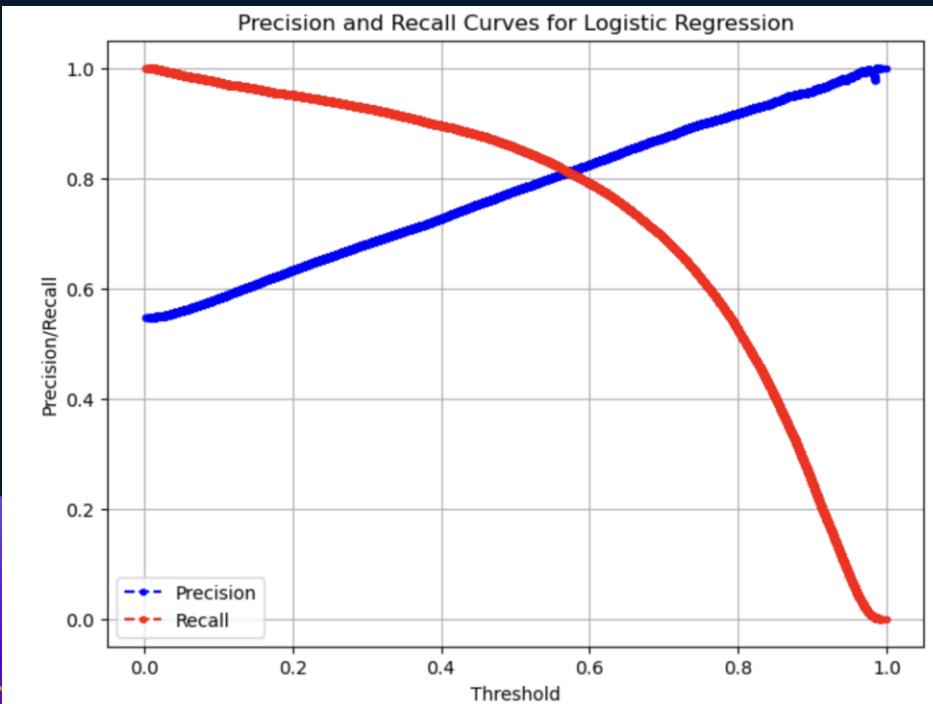
Satisfied customers misclassified as dissatisfied.

# Establishing Ground Rules & Understanding What do We Need Out of our Model

- **High False Positives** pose a significant business risk, as they might conceal underlying service issues.

**Reducing False Positives** is crucial to ensure dissatisfied customers are identified and targeted for improvement efforts.

# Tuning Model



**Precision (Blue Curve):** Increases as the threshold rises, leading to fewer false positives.

**Recall (Red Curve):** Decreases as the threshold rises, risking missed satisfied customers.

## The Trade-Off

**Lower Thresholds:** High recall but more false positives.

**Higher Thresholds:** High precision but lower recall.

**Optimal Threshold should be 0.6-0.7 to balance both.**

# Comparing Thresholds

**Chosen Threshold:** 0.65

**Precision:** 85%. Meets goal of reducing false positives.

**Recall:** 74.6%. Still accurately categorizes a majority of satisfied customers.

## Business Impact:

In a customer satisfaction model, it is typically more important to avoid False Positives than to maximize recall.

Dissatisfied customers mistakenly classified as satisfied can harm the airline's ability to:  
Identify and address operational or service-related problems.

### Training performance:

	Threshold	Accuracy	Recall	Precision	F1
0	0.57	0.79403	0.813428	0.810854	0.812139

### Testing performance:

	Threshold	Accuracy	Recall	Precision	F1
0	0.57	0.793296	0.81192	0.810703	0.811311

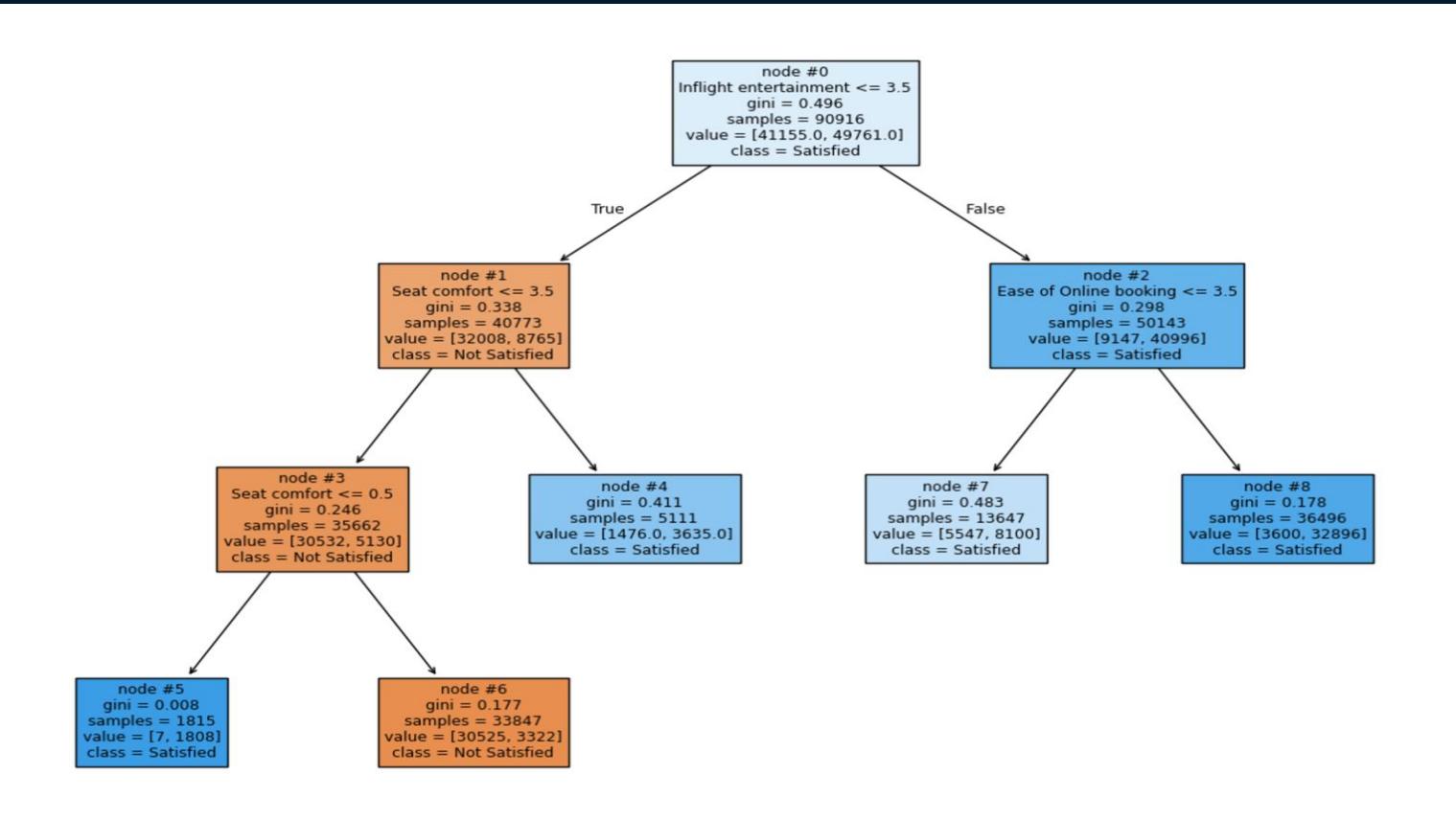
### Training performance:

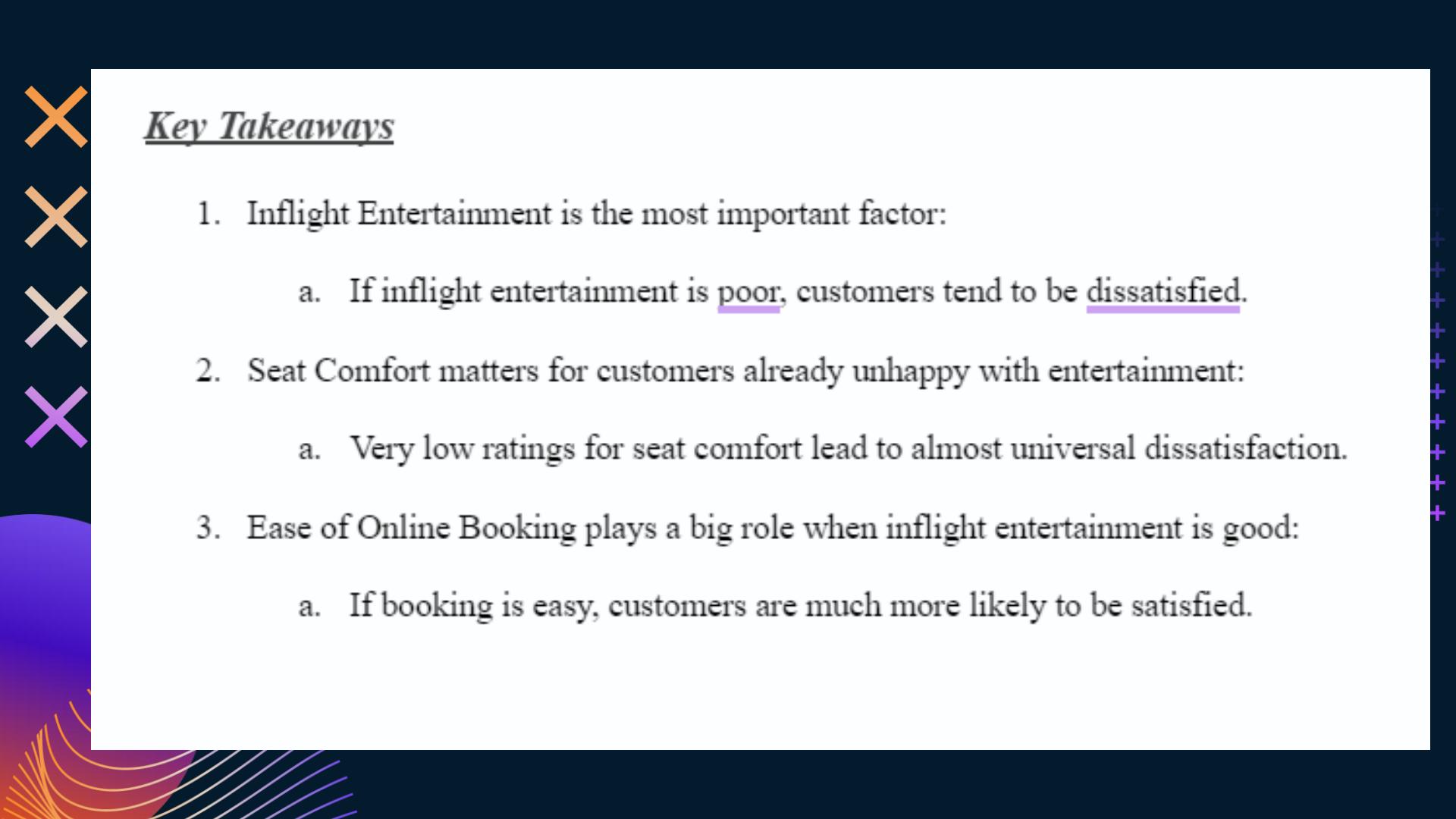
	Threshold	Accuracy	Recall	Precision	F1
0	0.65	0.788508	0.74685	0.84859	0.794476

### Testing performance:

	Threshold	Accuracy	Recall	Precision	F1
0	0.65	0.789267	0.746366	0.850313	0.794956

# Decision Tree





## Key Takeaways

1. Inflight Entertainment is the most important factor:
  - a. If inflight entertainment is poor, customers tend to be dissatisfied.
2. Seat Comfort matters for customers already unhappy with entertainment:
  - a. Very low ratings for seat comfort lead to almost universal dissatisfaction.
3. Ease of Online Booking plays a big role when inflight entertainment is good:
  - a. If booking is easy, customers are much more likely to be satisfied.

- The min impurity decrease ensures only meaningful splits are made.

```
[46]: # Choose the type of classifier.
dTTree_tuned = DecisionTreeClassifier(criterion='gini', random_state=1)

# Grid of parameters to choose from

parameters = {
    "max_depth": np.arange(3, 15, 2),
    "min_samples_leaf": [2, 5, 10, 15, 20, 25],
    "max_leaf_nodes": [5, 10, 15, 20, 30, 50],
    "min_impurity_decrease": [0.0001, 0.001, 0.01, 0.1],
}

# Scoring metric - F1 Score (Balanced Precision & Recall)
f1_scorer = make_scorer(f1_score)

# Run the grid search
grid_obj = GridSearchCV(dTTree_tuned, parameters, scoring=acc_scorer, cv=5)
grid_obj = grid_obj.fit(X_train, y_train)

# Set the clf to the best combination of parameters
dTTree_tuned = grid_obj.best_estimator_

# Fit the best algorithm to the data.
dTTree_tuned.fit(X_train, y_train)

# Output the best parameters
print("Best Parameters:", grid_obj.best_params_)

Best Parameters: {'max_depth': 3, 'max_leaf_nodes': 5, 'min_impurity_decrease': 0.0001, 'min_samples_leaf': 2}
```

# Model Comparison

- Logistic Regression
- 
- Precision: 85.03%
- Recall: 74.63%
- Accuracy: 78.9%
- F1-Score: 79.5%

- 
- Strengths:

Simplicity and interpretability make it ideal for communicating insights to stakeholders.

Threshold tuning allowed us to prioritize precision and reduce False Positives (key focus).



## Decision Tree Model

- Precision: 85.03%  
Recall: 74.63%  
Accuracy: 78.9%  
F1-Score: 79.5%

- Strengths:

Identifies clear patterns and relationships.

Offers actionable insights with a structured decision-making process.



- 
- 
- 
- 
- 

# Final Model

# Recommendations

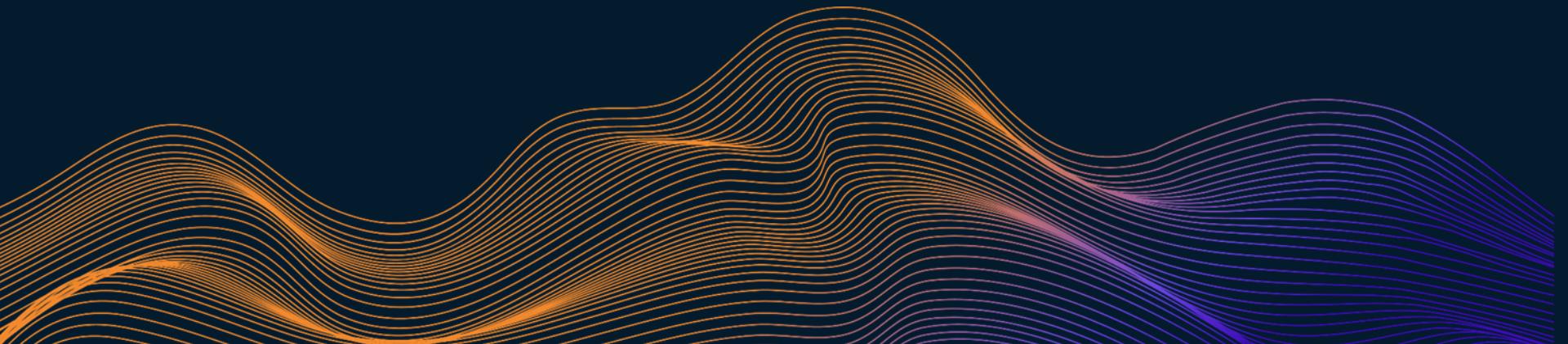
- Higher Precision: which ensures we are correctly identifying dissatisfied customers

Additionally its visual interpretability offers clear, actionable rules for business decision-making.



04

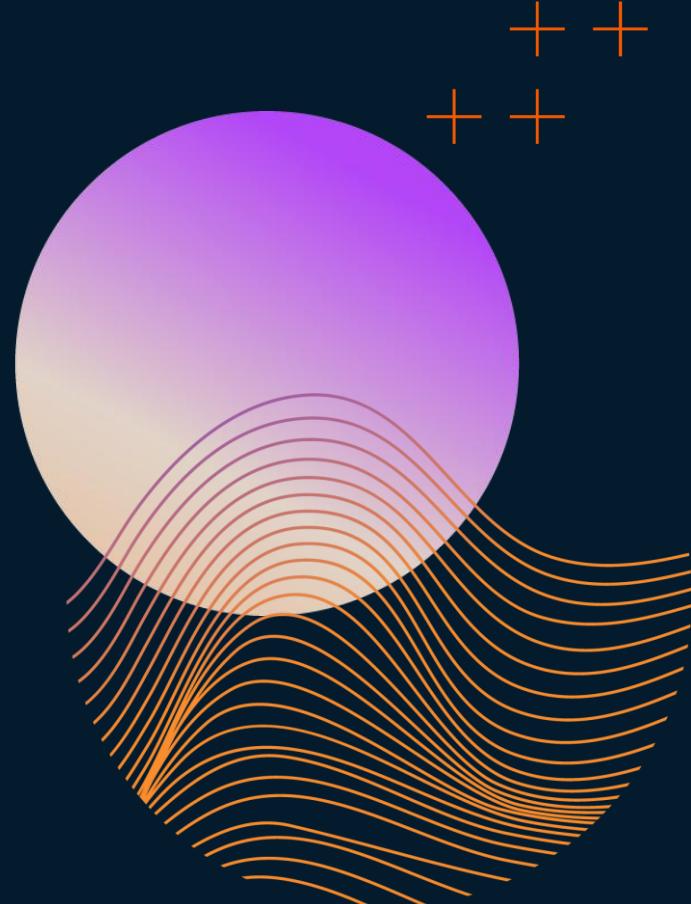
# Conclusion & Recommendations





# General Recommendations

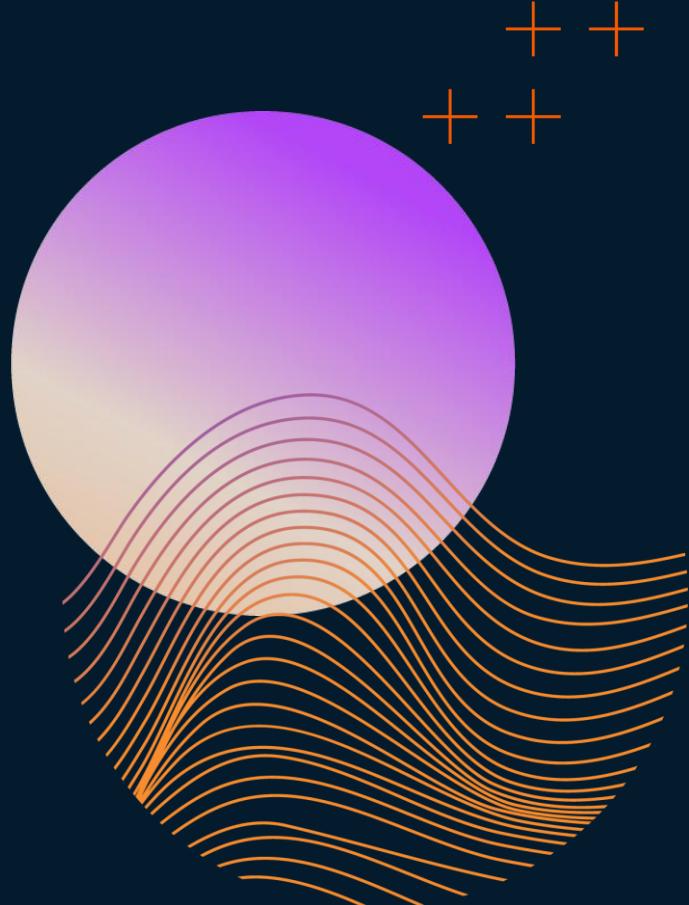
1. Prioritize improvements in **inflight entertainment, seat comfort, and online booking.**
2. Implement **real-time feedback systems** for quick issue resolution.
3. Customize services for **different customer segments** (e.g., business vs. leisure travelers).
4. Use **predictive analytics** to optimize resources and reduce false satisfaction predictions.





# Specific Recommendations

- 1. Inflight Entertainment:** Upgrade systems with personalized content and better options.
- 2. Seat Comfort:** Improve seating ergonomics and introduce premium options.
- 3. Business Travelers:** Offer premium packages with priority seating, lounge access, and better meals.
- 4. Ease of Booking:** Simplify the process via user-friendly apps and one-click rebooking.
- 5. Customer Loyalty:** Launch rewards programs and address dissatisfaction with compensations.





# Resources

Data Source:

<https://www.kaggle.com/code/iamsouravbanerjee/shopping-trends-unveiled-eda-for-beginners>



Other Resources:

<https://www.ibm.com/topics/exploratory-data-analysis>

<https://www.geeksforgeeks.org/what-is-exploratory-data-analysis/>



# Thank you!

