

# A Comprehensive Analysis of Airline Passenger's Satisfaction through Classification Model

# Sejal Jethwa<sup>1</sup>, Neha Vora<sup>2</sup>

<sup>1</sup>Student, <sup>2</sup>Assistant Professor, SVKM's Usha Pravin Gandhi College of Arts, Science and Commerce, Mumbai, India

Email: 1sejaljethwa2201@gmail.com, 2nehavora2501@gmail.com

#### **Abstract**

Passengers' satisfaction with airlines is one of the important determinants of competitive advantage for airlines with globalization and customer-oriented aviation industries of today. Customer satisfaction directly impacts brand loyalty, customer retention, and profitability, becoming an important key performance indicator for airline operators. The research examines a dataset of over 120,000 passengers, evaluating various aspects of their flying experience such as inflight services, comfort, punctuality, and overall satisfaction. An in-depth analysis of the key factors for customer satisfaction is done based on several data analyses and using different machine learning techniques like Logistic Regression, K-Nearest Neighbors, Decision Tree, and Random Forest. Through these techniques, our results show that, among others, service quality, punctuality, and comfort are the essential requirements for customer satisfaction. Conclusively, this study offers actionable recommendations for airlines that are vital in improving passenger satisfaction, which is an important factor in developing effective customer loyalty in the highly competitive airline business. The K-Nearest Neighbors (KNN) model performed best with an F1 score of 0.93, excelling in balancing precision and recall. Other models like Decision Tree and Random Forest were also used, with Random Forest providing robustness due to its ability to handle large datasets without overfitting, while Logistic Regression gave interpretable results but had the lowest F1 score at 0.89.

**Keywords:** Classification, Airline, Satisfaction, Predicting, Machine Learning Model.

#### 1. Introduction

Customer satisfaction in the airline industry has been widely researched due to its direct impact on an airline's success. With growing competition, particularly from low-cost carriers, airlines need to highlight service quality, reduce delays, and generally enhance the travel experience. Research has demonstrated that passengers' perceptions of the flight experience include comfort, in-flight entertainment, food and beverage services, and punctuality. These key drivers make a significant difference in enhancing customer loyalty, thereby improving profitability.

The current research attempts to pinpoint the important determinants of airline passenger satisfaction based on a dataset containing 120,000 records. The dataset covers passengers' demographic information, details of flights, and passengers' evaluation of inflight Wi-Fi service, food and beverage service, seat comfort, and inflight entertainment. Through statistical analysis and machine learning methods, this study emphasizes the most influential factors that drive passenger satisfaction and provides recommendations for service quality improvement. In the past decade, as market competition has intensified, airlines have placed greater emphasis on enhancing service quality to increase passenger satisfaction. Studies indicate a strong connection between service quality, passenger satisfaction, and airline profitability.

In today's highly competitive airline industry, offering high-quality services is essential for airlines to stay profitable and grow. Many studies suggest that passengers view service quality as a complex, multi-faceted concept. In-flight services, such as the performance of flight attendants and the amenities available on the plane, are especially important because they directly affect passengers' experiences and their overall satisfaction with the airline.

#### 2. Related Work

Research regarding this aspect has been performed extensively, as airline passenger satisfaction is fundamental in customer retention brand loyalty, and, finally, profitability [1]. Service quality has always ranked high as the main determinant of passenger satisfaction for airlines. Research by [2] recently used classification algorithms to predict passenger satisfaction, and this study shows that service quality is the most important variable in creating a customer experience and loyalty.[3] used Indian Airlines as an example to prove how

"different factors like inflight services, food and beverage, and the performance of flight attendants" are the most considerable factors for customer satisfaction in general.

According to [4], comfort is one of the airline loyalty drivers. Comfort is physical, and it even becomes a critical factor that may tell passengers who get a comfortable feel during their flight will express overall satisfaction more likely. The attributes of service quality, punctuality, and comfort on board have been identified as the crucial drivers in most studies conducted to establish the determinants of passenger satisfaction. [5]. While satisfaction is a direct response to consumption, service quality is defined as the overall impression of the customer about the service delivered [6].

If the received service is as expected, the service quality will be satisfying, and if the service received is more than their expectations, then customers will be delighted and will consider the quality of service to be very good. It is important to conceptualize the characteristics of an airline's service to estimate them accurately. Airline service is a concept that represents all types of services provided by airlines. It is essential to clarify and then apply the concept of airline service [7]. Their study investigates service quality from the perspective of domestic airline passengers in Tamil Nadu, offering insights into how regional preferences can shape customer satisfaction metrics[8]. Their work evaluates how airline service quality influences customer satisfaction and loyalty in the Indian market, providing insights into the critical aspects of airline operations that affect passenger loyalty [9].

This author focuses on evaluating passenger satisfaction with airline service quality using consecutive methods, providing a comprehensive look at how different service dimensions influence overall satisfaction [10]. Their research examines how expectations of service quality vary among airline passengers across different regions, offering valuable insights into global variations in customer satisfaction [11]. The study explores passengers' satisfaction with airline services through the Classification and Regression Tree (CART) approach together with the Importance-Performance Analysis [12].

# 3. Methodology

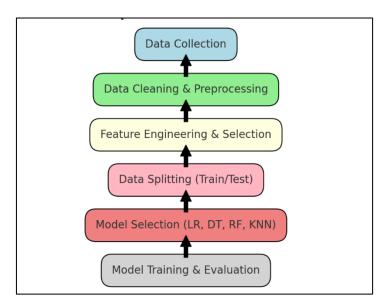


Figure 1. Proposed Work Flowchart

A research process (Figure 1) includes developing and restructuring questions, gathering, preparing, arranging, and analyzing data, drawing conclusions from it, and deducing from it. Metrics such as accuracy, precision, recall, and F1-score were used for the model's performance. The understanding gathered from it would be of immense use to the board of aircraft regarding enhancing the quality of service and the overall customer experience.

For this study, we used a secondary dataset from Kaggle. This dataset was collected by airlines through a survey to measure their passengers' satisfaction across various factors [13]. This data is secondary because it was originally gathered for a different purpose and is for use in our analysis. The dataset comprises 120,000 records. Degree levels (satisfied vs. neutral/dissatisfied) appear to be the classes, though no class distribution is indicated.

### 3.1 Dataset Preparation

Table 1 outlines key variables identified from a study related to airline passenger experiences, based on variables including demographics of passengers like ID, gender, and age; and information about their travel, such as customer type, travel purpose, class of travel, distance, delays, and satisfaction ratings for various services: check-in, boarding, in-flight services, and baggage handling. Each variable is further classified with its data type, providing an example of values when appropriate.

Table 1. Distribution Table

| Variable                           | Description                         | Data Type   | Example<br>Values     |
|------------------------------------|-------------------------------------|-------------|-----------------------|
| ID                                 | Unique ID for each passenger        | Integer     | 1, 2, 3               |
| Gender                             | Gender of the passenger             | Categorical | Male, Female          |
| Age                                | Age of the passenger                | Integer     | 25, 34, 45            |
| Customer Type                      | Type of customer                    | Categorical | First-time,<br>Return |
| Type of Travel                     | Purpose of travel                   | Categorical | Personal,<br>Business |
| Class                              | Travel class on the airline         | Categorical | Eco, Eco Plus,<br>Bus |
| Flight Distance                    | Distance in miles                   | Integer     | 500, 1200             |
| Departure Delay                    | Minutes of delay at departure       | Integer     | 0, 15, 45             |
| Arrival Delay                      | Minutes of delay at arrival         | Float       | 0, 20, 60             |
| Departure & Arrival<br>Convenience | Satisfaction for timing convenience | Integer     | 1, 3, 5               |
| Online Booking Ease                | Ease of online booking              | Integer     | 1, 4, 5               |
| Check-in Service                   | Satisfaction for check-in service   | Integer     | 2, 3, 5               |
| Online Boarding                    | Online boarding process             | Integer     | 1, 2, 5               |
| Gate Location                      | Satisfaction for gate location      | Integer     | 3, 4, 5               |
| On-board Service                   | Satisfaction for onboard service    | Integer     | 2, 4, 5               |
| Seat Comfort                       | Seat comfort level                  | Integer     | 1, 3, 5               |
| Leg Room Service                   | Leg room sen 👃                      | Integer     | 2, 4, 5               |

# 3.2 Preprocessing

**Data Cleaning:** Preprocessing techniques involve data cleaning. Several rows were removed from the dataset because their missing values exceeded the numbers needed in the respective fields. Missing values, particularly in the Arrival Delay feature, were handled using the most frequent strategy with the SimpleImputer.

**IQR** (Interquartile Range): A statistical method for discovering and correcting outliers from skewed distribution and extreme values was used to eliminate outliers in fields like "Flight Distance" and "Departure Delay." The Flight Distance values that were over the upper bound, were classified as an outlier, and the departure delay and arrival delay were removed. The boxplot method was used to eliminate outliers by employing IQR.

Feature Engineering: It transforms raw data into meaningful features that improve the performance of machine learning models. The missing values were imputed using the most frequent strategy with SimpleImputer. These are numerical features that are scaled using Robust Scaling to deal with extreme value and make the model robust. The pandas libraries were utilized for data manipulation, while Scikit-learn libraries such as SimpleImputer,

RobustScaler, and get\_dummies were utilized for categorical variable encoding and numerical feature scaling.

**Feature Selection:** After feature engineering, the features were used to train the machine learning models. These features encompass both numerical and categorical data that are important indicators of passenger satisfaction, making them relevant for predicting whether a passenger is satisfied or dissatisfied with their flight experience.

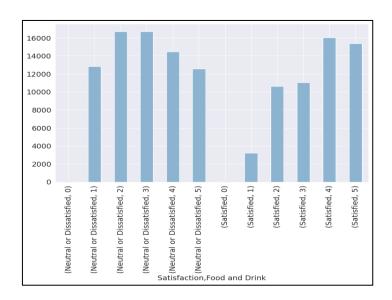
**Feature Importance:** In some models like Decision Tree and Random Forest, feature importance was computed, showing that features such as Flight Distance, Departure Delay, and Seat Comfort had a significant impact on the prediction of passenger satisfaction.

# 3.3 Experimental Setup

This experiment utilized Python libraries like Pandas, NumPy, and Scikit-learn for machine learning models, Matplotlib and Seaborn for data visualization, and a Python Notebook for development.

# 3.4 Data Splitting

The split set separates the training set (80%) and test set (20%) to assess model generalization to new data. This step allows machine learning models to learn from the data without overfitting, ensuring independent validation of model performance.



**Figure 2.** Relation between the Food and Drink Ratings and the Satisfaction of the Customer

Figure 2 shows the relation of customer satisfaction with their rating of food and drink. This chart, distributed based on customer satisfaction levels by food and drink, is categorized into various degrees of satisfaction. Neutral or Dissatisfied: N/A and Dissatisfied, with subcategories from 0 to 5 Satisfied: has the same range of subcategories from 0 to 5. All these numbers probably are rating scales, wherein "0" states least satisfaction and "5" states utmost satisfaction. The y-axis is plotted as the number of respondents for every satisfaction level (on a scale ranging from 0 to 16,000).

Neutral or Unsatisfied (0-5): There is a peak at levels 2 and 3, which indicates that most of the respondents rated their level of satisfaction as mid-range Satisfied (0-5): The number of respondents for levels 4 and 5 are the highest. It may indicate that customers rated themselves as 4 or 5, meaning highly satisfied with food and drink. This chart is probably a survey or customer feedback that suggests there is indeed a large number of customers who are either neutral, unhappy, or dissatisfied with the food and drinks, while at the same time remaining a significant percentage of very satisfied customers.

# 3.5 Data Visualization

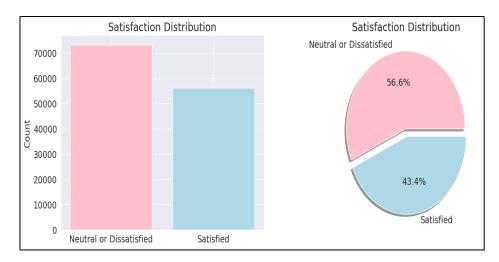


Figure 3. Visualization of Satisfaction

Above Figure 3. is a graph of distribution in satisfaction levels in two different formats: as a bar graph and as a pie graph. The bar graph contains two categories "Satisfied" and "Neutral or Dissatisfied. "For the category "Satisfied," the number of people is approximately 48,000. For the category "Neutral or Dissatisfied," the number of people is approximately 68,000. This graph reflects that more people were either Neutral or Dissatisfied than those who reported they were Satisfied.

The pie chart reiterates the information on the bar chart. Blue represents the "Satisfied" category, which is 41.4%. Red represents the "Neutral or Dissatisfied" category, which is 58.6%. The graph states that while the majority of the respondents are classified as either neutral or dissatisfied at 58.6%, the minority is satisfied at 41.4%.

#### 3.6 Model Used

**Logistic Regression:** This is a simple but powerful model that generally can be used for binary classifications. The reason chosen is that it can deliver very interpretable results by the probability estimates. It helps further the understanding of the features, and flight experience variables, in this case, related to customer satisfaction.

**Decision Tree:** Decision trees are very interpretable models. Airlines can directly see through the decision trees the path leading the customer to either be satisfied or dissatisfied with their flight. While competitive in their accuracy, decision trees tend to overfit the training data. Thus, the model may generalize worse when faced with new unseen data.

**Random Forest:** It is an ensemble of decision trees, hence more robust, with less overfitting than a single Decision Tree. It was chosen because it can handle large sizes of datasets and lots of features, which would be especially important in this problem to determine which aspects of flight experience. The model also identifies the importance of ranking features for airlines to focus on areas of priority for improvement.

**K-Nearest Neighbors:** KNN was chosen for its capability of capturing nonlinear relationships in the dataset. KNN predicts by comparing a given profile of a passenger with the most similar ones within the training set. This model had the best performance, reaching an F1-Score of 0.93, showing that it is good at balancing precision and recall. It is simple yet shows a very strong performance, which makes it a suitable fit for this dataset, even though it can be computationally expensive concerning large datasets.

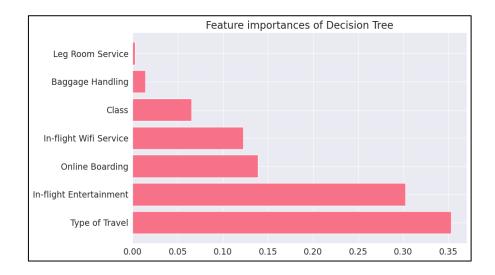


Figure 4. Feature Importance of Decision Tree Model

Figure 4 illustrates the feature importance in the Decision Tree model. Those are the factors contributing the most to its decisions within the model. "Type of Travel" is the most important feature, having the highest impact on the model's decisions. "In-flight Entertainment" and "Online Boarding" also contribute significantly to the decision from this model. "Leg Room Service" and "Baggage Handling" contributed the least influence on model output.

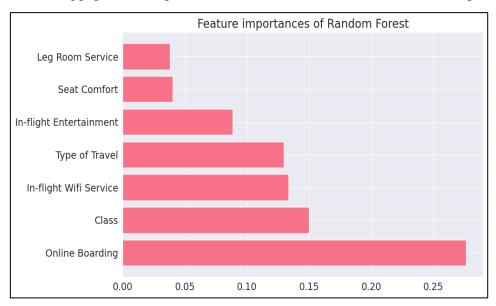


Figure 5. Feature Importance of Random Forest Model

Figure 5 explains the importance of the different variables involved in predicting with a Random Forest Model. Every bar represents a different variable (such as Class, Leg Room Service, etc.) The height of that bar depicts how important that variable was for coming up

with the predictions. The taller the bar is, the more important this variable is. The shorter the bar is, the less important this variable is.

#### 4. Results and Discussion

The machine learning models (Logistic Regression, Decision Tree, K-nearest Neighbors (KNN), and Random Forest) were evaluated using these four metrics to measure their performance on a classification task. Each model was likely trained and tested on the same dataset, with performance metrics calculated to compare the models' effectiveness. The evaluation metrics like accuracy, precision, recall, and F1 score were calculated. Table 2 below shows the overall results.

 Table 2. Overall Result

| Model                           | Accuracy | Precision | Recall | F1-Score |
|---------------------------------|----------|-----------|--------|----------|
| Logistic<br>Regression          | 0.87     | 0.88      | 0.90   | 0.89     |
| Decision<br>Tree                | 0.92     | 0.91      | 0.96   | 0.91     |
| K-nearest<br>neighbors<br>(KNN) | 0.92     | 0.91      | 0.96   | 0.93     |
| Random<br>Forest                | 0.91     | 0,91      | 0.94   | 0.92     |

However, in respect of those models, there were some model-specific hyperparameters, which are as follows: K-Nearest Neighbors (KNN): The number of neighbors was defined with n\_neighbors=7. Decision Tree: The parameters used were max features = 6 and max depth = 4. Random Forest: Some hyperparameters were defined, including maximum features: max\_features=5 and maximum depth: max\_depth=4.

The model's strengths and weaknesses for each model in the results section: K-Nearest Neighbors gave the best F1-score of 0.93. Precision equaled its recall but resulted in high computation with large data. Decision Trees are easier to interpret, but they tend to overfit the model. Random Forest helps in dealing with big data and avoids overfitting; however, it is less interpretable compared to others. The simple and interpretable model is Logistic Regression. The lowest F1 score was 0.89. Unlike the logistic regression model, it just performed poorly.

It has the lowest F1-Score of any model in this comparison, at 0.89. In general, KNN proved to be the best match for this dataset, yielding the best measures in terms of accuracy, precision, recall, and F1 score. Conversely, both the Decision Tree and the Random Forest models represent competitive alternatives as well, depending on the specific needs regarding interpretability and computational efficiency in practice.

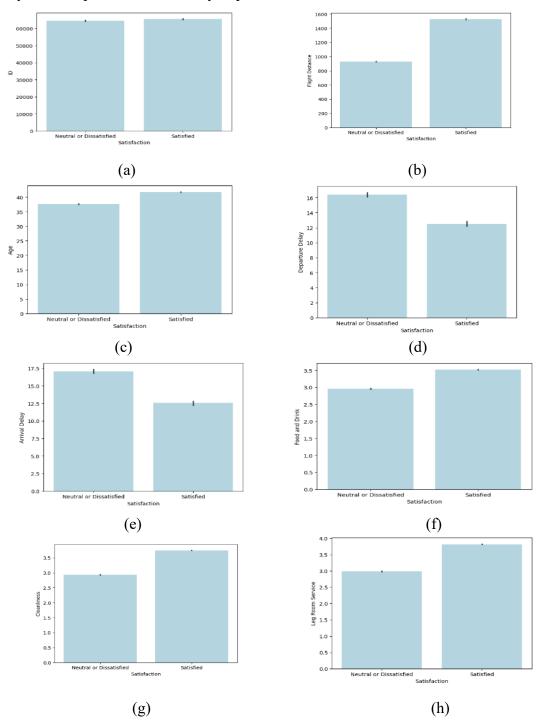


Figure 6. Comparison of Passenger Satisfaction for Airline Services Offered

Figure 6: graphically represents passenger satisfaction by examining the relationship that may occur between delays (arrival and departure) and customer satisfaction on an aggregate basis.

**Focus on Delays:** The figure mostly focuses on how both departures and arrival delays affect the satisfaction of air passengers. According to the study, delays at the beginning or the end of the journey are major contributors to dissatisfaction from airline passengers.

**Unsatisfaction Caused by Delay:** The graph indicates that most of the delayed passengers exhibited dissatisfaction with their flight experience. This, therefore, unmistakably points out that there is a link between rescheduling and dissatisfied feedback from customers.

Comparison Across Different Delay Types: The above figure compares both the forms of delays, which are the departure and arrival ones but concludes only that any form of delay negatively affects satisfaction and then does not indicate which of these is more adverse; it lumps them together as two of the key reasons for dissatisfaction. Thus Figure 6 indicates that the reduction of delay should be placed on the 'must to do' list by airlines if they want to have a chance at sustaining the benefits experienced in terms of customer satisfaction, as delays are the leading cause of dissatisfaction according to this analysis.

### 5. Conclusion and Future Scope

In conclusion, airlines must be able to understand and enhance airline passenger satisfaction if they want to be competitive and successful. If the passengers are unhappy and unsatisfied, the company cannot grow. The limitation of the airline passenger satisfaction feature could have been much better if neutral/unsatisfied were separate categories. The data could have been divided into 3 groups as satisfied, neutral, and dissatisfied passengers. This can offer a meaningful prediction as it's hard to divide between neutral and unsatisfied passengers. In the future, making airline passengers happier will involve a few key changes. Airlines will use advanced technology to offer more personalized experiences, like tailoring in-flight services to each passenger's preferences. Technology will also improve, with better entertainment options, easier self-service check-ins, and faster boarding processes.

#### References

- [1] Usha, P., and E. Kusuma. "A study on survive quality and passenger satisfaction on Air India services." International Journal of Advance Research, Ideas And Innovations In Technology 3, no. 4 (2017): 534-546.
- [2] B.Herawan Hayadi, Jin-Mook Kim, Khodijah Hulliyah, Husni Teja Sukmana. "Predicting Airline Passenger Satisfaction with Classification Algorithms." International Journal of Informatics and Information Systems 4.1 (2021): 82-94.
- [3] Dolnicar, Sara, Klaus Grabler, Bettina Grün, and Anna Kulnig. "Key drivers of airline loyalty." Tourism Management 32, no. 5 (2011): 1020-1026.
- [4] Archana, R., and M. V. Subha. "A study on service quality and passenger satisfaction on Indian airlines." International Journal of Multidisciplinary Research 2, no. 2 (2012): 50-63.
- [5] Watkins, Marley W. A step-by-step guide to exploratory factor analysis with SPSS. Routledge, 2021.
- [6] Malhotra, Naresh K. Marketing research: an applied prientation. pearson, 2020.
- [7] Bhat, Vasanthakumar N. "A multivariate analysis of airline flight delays." International Journal of Quality & Reliability Management 12, no. 2 (1995): 54-59.
- [8] Bhuvaneswaran, R., & Vijayanagar, D.(2013)."Service quality towards domestic airlines passenger perception in Tamilnadu". Asia Pacific Journal of Research, 2(9), 133–140.
- [9] Shashikala, P. (2020). Advanced Business Analytics. Cengage Learning India Private Limited
- [10] Agarwal, Ira, and Kavitha R. Gowda. "The effect of airline service quality on customer satisfaction and loyalty in India." Materials Today: Proceedings 37 (2021): 1341-1348.
- [11] Tahanisaz, Sahar. "Evaluation of passenger satisfaction with service quality: A consecutive method applied to the airline industry." Journal of Air Transport Management 83 (2020): 101764.

- [12] Punel, Aymeric, Lama Al Hajj Hassan, and Alireza Ermagun. "Variations in airline passenger expectation of service quality across the globe." Tourism Management 75 (2019): 491-508.
- [13] https://www.kaggle.com/datasets/mysarahmadbhat/airline-passenger-satisfaction