

# Airfare Price Detection

## Abstract

Flight fares are notoriously dynamic and influenced by various factors, making it challenging for travellers to find affordable tickets. This research proposes a machine learning approach to predict flight fares using historical data. The study utilizes a Random Forest Regressor model to capture complex relationships between features like airline, source, destination, date of journey, and more. The model's performance is evaluated using metrics such as R-squared, Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE). The results demonstrate the effectiveness of the model in accurately predicting flight fares, providing valuable insights for travellers and the travel industry. This research contributes to the growing body of knowledge on flight fare prediction and offers practical applications for cost-effective travel planning.

## 1. Introduction

### I. Significance and Contribution

Emphasize the potential benefits and contributions of your research. Highlight how the developed model can assist travellers in making informed decisions, saving costs, and enhancing travel planning. Discuss any novel aspects or improvements your research brings compared to existing work in the field.

- **Example**

"This research has the potential to significantly benefit travellers by providing an accurate and reliable flight fare prediction tool. The developed model can assist in making informed booking decisions, identifying cost-saving opportunities, and optimizing travel plans. This study contributes to the existing body of knowledge on flight fare prediction by [mentioning specific contributions, like improved accuracy, novel feature engineering techniques, or addressing specific limitations of previous work]."

### II. Paper Organization

Conclude the introduction by briefly outlining the structure of your research paper. Mention the subsequent sections and their respective content.

- **Example**

"The remainder of this paper is organized as follows: Section 2 provides a detailed description of the dataset and data preprocessing techniques. Section 3 describes the feature engineering process. Section 4 explains the model selection and training procedures. Section 5 presents the model evaluation results and analysis. Finally, Section 6 concludes the paper with a discussion of findings, limitations, and future research directions."

### III. Problem Statement and Motivation

Begin by introducing the real-world problem of fluctuating flight fares and its impact on travellers. Highlight the difficulty in predicting fares due to various dynamic factors like demand, seasonality, and competition. Emphasize the need for an accurate and reliable flight fare prediction system to empower travellers with informed decision-making and cost savings.

- **Example**

"Fluctuating flight fares pose a significant challenge for travellers seeking affordable air travel. The dynamic nature of pricing, influenced by factors like demand, seasonality, and competition, makes it difficult to predict ticket costs accurately. This unpredictability can lead to missed opportunities for savings and increased travel expenses. Therefore, developing an accurate and reliable flight fare prediction system is crucial to empower travellers with informed choices and cost-effective travel planning."

### IV. Research Objectives and Scope

Clearly state the main objectives of your research, focusing on developing a predictive model for flight fares. Define the scope of your project, including the specific airlines, routes, or timeframes considered. Mention any limitations or constraints of your study.

- **Example**

"This research aims to develop a machine learning-based predictive model for flight fare prediction. The model will leverage historical flight data to identify patterns and relationships between various factors and ticket prices. The scope of this study includes domestic flights within [specify region or country] operated by [mention specific airlines, if applicable]. While the model strives for generalizability, it's essential to acknowledge limitations regarding external factors like unforeseen events or sudden market shifts."

### V. Methodology Overview

Provide a concise overview of the methodology employed in your research. Briefly mention the key steps involved, such as data collection, cleaning, feature engineering, model selection, training, and evaluation. Highlight the use of machine learning algorithms, such as the Random Forest Regressor, and any specific techniques like hyperparameter tuning.

- **Example**

"This research employs a machine learning approach to flight fare prediction. The methodology involves collecting historical flight data, cleaning and preprocessing the data, engineering relevant features, selecting an appropriate machine learning model, training the model on the prepared data, and evaluating its performance using various metrics. The Random Forest Regressor algorithm will be utilized for prediction, and hyperparameter tuning techniques will be applied to optimize model accuracy."

## 2. Methodology

### I. Data Collection

The assortment of data is the very first step in machine learning projects. There are various sources of data available on numerous websites that are deployed to construct the models. These sites supply a huge variety of data regarding different airlines, routes, times, and tolls. In this part, data gathered from the various available sources are studied. For the execution of this, information is brought from a site called Kaggle. For the assortment of the data and to execute the model's Python is utilized. The dataset collected contains information about different airlines in India. It consists of various factors which affect the price of a flight ticket including the price for a particular flight. It contains 10683 rows of data. The features present in the dataset are the name of companies, Date of travelling, Origin, terminus, path of travelling, Time of Departure, Time of Arrival, Travelling Hours, Total Stoppage, Additional Info, and Price

### A. Dataset Description

The dataset likely contained the following features:

**Airline:** The name of the airline operating the flight.

**Date\_of\_Journey:** The date of the flight.

**Source:** The origin city of the flight.

**Destination:** The destination city of the flight.

**Route:** The route taken by the flight.

**Dep\_Time:** The departure time of the flight.

**Arrival\_Time:** The arrival time of the flight.

**Duration:** The total duration of the flight.

**Total\_Stops:** The number of stops during the flight.

**Additional\_Info:** Any additional information about the flight.

**Price:** The price of the flight ticket (target variable).

### B. Data Collection Process

The data collection process would typically involve the following steps:

**Identifying a suitable data source:** Searching for publicly available datasets containing flight fare information.

**Downloading the dataset:** Obtaining the dataset in a suitable format, such as a CSV or Excel file.

**Inspecting the data:** Examining the dataset to understand its structure, features, and data types.

**Cleaning and preprocessing the data:** Handling missing values, transforming data types, and engineering new features as needed.

### C. Ethical Considerations

When collecting data for a research project, it's important to consider ethical implications, such as data privacy and potential biases in the data. If the dataset contains personally identifiable information, it should be anonymized or aggregated to protect individual privacy. Additionally, any potential biases in the data should be acknowledged and addressed to ensure fairness and accuracy in the model's predictions.

## II. Cleaning and Preparing of Data

Data cleaning and preparation are crucial steps in any machine learning project, ensuring the quality and suitability of the data for model training. In this flight fare prediction project, several data cleaning and preparation techniques were applied to enhance the model's performance.

### A. Handling Missing Values

**Identification:** The dataset was examined for missing values using the `isnull().sum()` function.

**Removal:** Rows with missing values were removed using the `dropna()` function to prevent errors during model training. This approach was chosen as the number of missing values was relatively small, minimizing the impact on the overall dataset.

### B. Feature Engineering

**Date and Time Features:** The 'Date\_of\_Journey' feature was split into separate 'Journey\_Day' and 'Journey\_Month' features to capture the day and month information. Similarly, 'Dep\_Time' and 'Arrival\_Time' were transformed into 'Dep\_hour', 'Dep\_minute', 'Arrival\_hour', and 'Arrival\_minute' to represent departure and arrival times in hours and minutes. The 'Duration' feature was converted into 'Duration\_Hours' and 'Duration\_Mins' to represent the flight duration in numerical form.

**Categorical Data Handling:** Categorical features like 'Airline', 'Source', and 'Destination' were transformed using one-hot encoding. This technique creates binary features for each category, allowing the model to learn the impact of each category without imposing any ordinal relationship between them. The 'Total\_Stops' feature was converted into numerical form by replacing string values with corresponding numerical values (e.g., 'non-stop' to 0, '1 stop' to 1).

### C. Data Scaling

**No Scaling for Random Forest:** Scaling was not performed in this project as the chosen model, Random Forest Regressor, is not sensitive to the scale of features. Random Forest models are tree-based and do not require feature scaling for optimal performance.

## III. Feature Engineering

Feature engineering is a crucial step in machine learning that involves transforming raw data into features that are more informative for the model. In this flight fare prediction project, several feature engineering techniques were applied to enhance the model's performance.

### A. Date and Time Features

**Date\_of\_Journey:** This feature was split into two separate features, 'Journey\_Day' and 'Journey\_Month', to capture the day and month of the journey. This allows the model to learn potential patterns based on the time of year and day of the week.

**Dep\_Time and Arrival\_Time:** These features were transformed into 'Dep\_hour', 'Dep\_minute', 'Arrival\_hour', and 'Arrival\_minute' to represent the departure and arrival times in hours and minutes. This helps the model to understand the impact of departure and arrival times on the fare.

**Duration:** The 'Duration' feature, which was originally in a string format (e.g., '2h 50m'), was converted into two numerical features, 'Duration\_Hours' and 'Duration\_Mins', to represent the duration of the flight in hours and minutes. This allows the model to directly utilize the duration information.

### B. Categorical Data Handling

**Airline, Source, and Destination:** These categorical features were transformed using one-hot encoding. This technique creates binary features for each category, allowing the model to learn the impact of each category on the fare without imposing any ordinal relationship between them.

**Total\_Stops:** This ordinal feature was converted into numerical form by replacing the string values ('non-stop', '1 stop', '2 stops', etc.) with corresponding numerical values (0, 1, 2, etc.). This allows the model to understand the relationship between the number of stops and the fare.

## C. Feature Selection

**Correlation Analysis:** A heatmap was generated to visualize the correlation between the features and the target variable (Price). This helps to identify features that have a strong relationship with the fare.

## D. Feature Importance

**Random Forest Feature Importance:** The feature importance scores obtained from the Random Forest model provide insights into the relative importance of each feature in predicting the fare. Features with higher importance scores have a greater impact on the model's predictions.

# IV. Predictive Modelling

Predictive modelling, a core aspect of machine learning, involves building a model that can accurately predict future outcomes based on historical data. In this flight fare prediction project, predictive modelling is employed to forecast airfare prices based on various factors like airline, source, destination, date of journey, and more.

## A. Training and Evaluation

**Data Splitting:** The dataset was divided into training and testing sets to evaluate the model's performance on unseen data. This helps to assess the model's generalization ability.

**Model Training:** The Random Forest Regressor was trained on the training data, learning the patterns and relationships between the features and the fare.

**Model Evaluation:** The trained model was evaluated using the testing data, comparing its predictions to the actual fares. Evaluation metrics such as R-squared score, Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE) were used to assess the model's accuracy.

## B. Hyperparameter Tuning

**Randomized Search:** A randomized search approach was used to optimize the model's hyperparameters, improving its predictive performance. This involves exploring different combinations of hyperparameter values to find the settings that yield the best results.

## C. Model Deployment

**Pickling:** The trained model was saved using pickling, allowing it to be reused for future predictions without retraining. This enables the model to be deployed in a real-world setting for practical use.

# V. Model Evaluation

Model evaluation is a critical step in the machine learning process, ensuring that the developed model accurately predicts outcomes and generalizes well to unseen data. In this flight fare prediction project, rigorous model evaluation was performed to assess the performance and reliability of the predictive model.

## A. Evaluation Metrics

Several metrics were used to evaluate the model's performance:

**R-squared (R2) Score:** Measures the proportion of variance in the target variable (Price) explained by the model. A higher R2 score indicates a better fit.

**Mean Absolute Error (MAE):** Calculates the average absolute difference between the predicted and actual fares. It represents the average prediction error in the original units of the target variable.

**Mean Squared Error (MSE):** Computes the average squared difference between the predicted and actual fares. It gives more weight to larger errors.

**Root Mean Squared Error (RMSE):** The square root of MSE, providing an error measure in the same units as the target variable. It is often preferred over MSE as it is more interpretable.

## B. Evaluation Procedure

**Data Splitting:** The dataset was divided into training and testing sets (80% for training, 20% for testing). The model was trained on the training set and evaluated on the unseen testing set to assess its generalization ability.

**Prediction:** The trained model was used to predict fares for the testing set.

**Metric Calculation:** The evaluation metrics (R2, MAE, MSE, RMSE) were calculated by comparing the predicted fares to the actual fares in the testing set.

**Interpretation:** The results of the evaluation metrics were analyzed to understand the model's accuracy and potential areas for improvement.

## C. Visualization

**Distribution Plot:** A distribution plot of the residuals (difference between predicted and actual fares) was used to visualize the error distribution. A normal distribution of residuals indicates a well-performing model.

**Scatter Plot:** A scatter plot of predicted fares against actual fares was generated to visualize the relationship between the two. A strong positive correlation indicates good predictive accuracy.

# 3. Machine Learning Techniques

This project leverages machine learning techniques to build a predictive model for flight fares. The core technique employed is the **Random Forest Regressor**, an ensemble learning method known for its robustness and accuracy in regression tasks. Let's delve deeper into the specifics:

## I. Random Forest Regressor

**Ensemble Learning:** Random Forest belongs to the ensemble learning family, where multiple decision trees are combined to create a more powerful and accurate model.

**Decision Trees:** Each decision tree in the forest is built on a random subset of the data and features, introducing diversity and reducing overfitting.

**Bagging:** This technique, also known as bootstrap aggregating, involves creating multiple subsets of the training data with replacement. Each subset is used to train a separate decision tree.

**Feature Randomness:** Random Forest further enhances diversity by randomly selecting a subset of features at each node of the decision tree, preventing over-reliance on any single feature.

**Aggregation:** Predictions from individual trees are aggregated, typically by averaging, to produce the final prediction. This reduces variance and improves generalization.

## ❖ Why Random Forest for Flight Fare Prediction?

**Handling Complex Relationships:** Flight fares are influenced by numerous factors with complex interactions. Random Forest effectively captures these non-linear relationships.

**Robustness to Outliers:** The averaging mechanism in Random Forest makes it less sensitive to outliers in the data.

**Feature Importance:** Random Forest provides insights into the relative importance of different features in predicting fares, aiding in feature selection and understanding the problem domain.

## II. Hyperparameter Tuning

**Randomized Search:** This technique was employed to optimize the hyperparameters of the Random Forest model. Randomized search explores a range of hyperparameter values randomly, offering a balance between efficiency and performance.

**Cross-Validation:** During hyperparameter tuning, the training data is further divided into folds, and the model is trained and evaluated on different fold combinations. This helps to assess the model's performance across different data subsets and prevents overfitting to the training data.

## III. Model Evaluation

**Metrics:** Several metrics, including R-squared score, Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE), were used to evaluate the model's predictive performance.

**Visualization:** Distribution plots of residuals and scatter plots of predicted vs. actual fares were used to visualize the model's accuracy and error distribution.

# 4. Machine Learning Algorithm Used

This project utilizes the Random Forest Regressor as the primary machine learning algorithm for predicting flight fares. Random Forest is a powerful ensemble learning method renowned for its robustness, accuracy, and ability to handle complex datasets. Here's a detailed explanation suitable for your research paper:

## I. Ensemble Learning

**Concept:** Random Forest belongs to the ensemble learning family, where multiple individual models (decision trees, in this case) are combined to create a stronger and more accurate predictive model.

**Advantages:** Ensemble methods generally outperform single models by reducing overfitting, improving generalization, and handling noise in the data.

## II. Decision Trees

**Building Blocks:** The fundamental units within a Random Forest are decision trees. These are tree-like structures that recursively partition the data based on feature values to make predictions.

**Structure:** Each decision tree consists of nodes representing features, branches representing decision rules, and leaf nodes representing predicted outcomes.

### III. Random Forest Construction

**Bagging (Bootstrap Aggregating):** Random Forest employs bagging, where multiple subsets of the training data are created by randomly sampling with replacement. Each subset is used to train a separate decision tree.

**Feature Randomness:** To further enhance diversity and prevent overfitting, Random Forest randomly selects a subset of features at each node of the decision tree. This prevents any single feature from dominating the model.

**Aggregation:** Predictions from individual decision trees are aggregated, typically by averaging, to produce the final prediction. This reduces variance and improves the overall accuracy of the model.

### IV. Algorithm in Action for Flight Fare Prediction

**Training:** The Random Forest Regressor is trained on a labelled dataset containing historical flight data, including features like airline, source, destination, date of journey, and more. The algorithm learns patterns and relationships between these features and the target variable (flight fare).

**Prediction:** Once trained, the model can predict the fare for a new flight by considering its features and traversing the decision trees in the forest. The final prediction is an aggregation of the predictions from individual trees.

### V. Advantages of Random Forest for Flight Fare Prediction

**Handling Complex Relationships:** Flight fares are influenced by numerous factors with complex interactions. Random Forest effectively captures these non-linear relationships, leading to accurate predictions.

**Robustness to Outliers:** The averaging mechanism in Random Forest makes it less sensitive to outliers in the data, improving its reliability.

**Feature Importance:** Random Forest provides insights into the relative importance of different features in predicting fares. This information can be valuable for feature selection and understanding the problem domain.

## 5. Results

This section presents the results of the flight fare prediction model, focusing on its performance evaluation and analysis. The primary goal is to demonstrate the model's accuracy, reliability, and effectiveness in predicting flight fares.

### I. Evaluation Metrics

**R-squared (R2) Score:** Report the R2 score achieved by the model on the testing dataset. This metric represents the proportion of variance in the target variable (flight fare) explained by the model. A higher R2 score indicates a better fit and predictive power.

**Mean Absolute Error (MAE):** Present the MAE value, which quantifies the average absolute difference between the predicted and actual fares. It provides a measure of the model's prediction error in the original units of the target variable.

**Mean Squared Error (MSE):** Report the MSE value, which calculates the average squared difference between the predicted and actual fares. It gives more weight to larger errors and is useful for assessing the overall model performance.

**Root Mean Squared Error (RMSE):** Present the RMSE value, which is the square root of MSE and provides an error measure in the same units as the target variable. It is often preferred over MSE for its interpretability.



## II. Visualization of Results

**Distribution Plot of Residuals:** Include a distribution plot of the residuals (difference between predicted and actual fares). A normal distribution of residuals indicates a well-performing model. Discuss any deviations from normality and their potential implications.

**Scatter Plot of Predicted vs. Actual Fares:** Present a scatter plot with predicted fares on the x-axis and actual fares on the y-axis. A strong positive correlation and points clustered around the diagonal line indicate good predictive accuracy.

## III. Interpretation and Analysis

**Model Performance:** Analyze the evaluation metrics and visualizations to assess the overall performance of the model. Discuss the strengths and weaknesses of the model's predictions.

**Feature Importance:** If available, present the feature importance scores obtained from the Random Forest model. This provides insights into the relative contribution of each feature in predicting fares and can be used to interpret the model's behaviour.

**Comparison with Baseline or Other Models:** If applicable, compare the performance of your model with a baseline model or other machine learning models. This highlights the improvements achieved by your approach.

## IV. Statistical Significance (if applicable)

If you conducted statistical tests to compare model performance, report the p-values and discuss their significance. This adds rigor to your analysis and supports your claims about the model's effectiveness.

## V. Limitations and Future Work

Acknowledge any limitations of your study, such as data limitations, assumptions made, or potential biases.

Suggest directions for future research to address these limitations and further improve the model's performance.

### Example

"The Random Forest Regressor achieved an  $R^2$  score of 0.85 on the testing dataset, indicating that the model explains 85% of the variance in flight fares. The MAE was found to be \$50, suggesting an average prediction error of \$50, suggesting an average prediction error of \$50. The RMSE was calculated as \$75, providing a measure of the typical prediction error in the same units as the target variable. The distribution plot of residuals exhibited a near-normal distribution, suggesting a well-performing model. The scatter plot of predicted vs. actual fares showed a strong positive correlation, further supporting the model's accuracy. Feature importance analysis revealed that [mention key features and their importance]. Future work could explore incorporating external data sources, such as weather or fuel prices, to potentially improve the model's predictive power."

# 6. Discussion

The discussion section is where you interpret your findings, discuss their implications, and relate them to existing research in the field. It's also an opportunity to acknowledge limitations and propose future research directions. Here's a breakdown of how to structure this section for your flight fare prediction project:

## I. Summary of Key Findings

Begin by summarizing the main findings from your results section. Briefly reiterate the model's performance using key metrics like R-squared, MAE, and RMSE. Highlight any noteworthy patterns or insights revealed by the feature importance analysis.

## Example

"The results demonstrate that the Random Forest Regressor effectively predicts flight fares, achieving an R-squared score of [value] and an MAE of [value]. This indicates that the model accurately captures the relationships between various factors and ticket prices. Feature importance analysis revealed that [mention key features and their contribution], suggesting their significant influence on fare determination."

## II. Interpretation and Implications

Discuss the implications of your findings in the broader context of flight fare prediction and travel planning. Explain how your model's insights can be applied to real-world scenarios. Consider the potential benefits for travellers, airlines, and the travel industry.

## Example

"The model's ability to accurately predict fares empowers travellers with valuable information for making informed booking decisions. By understanding the factors that influence prices, travellers can optimize their travel plans and potentially save costs. Airlines can leverage these insights to adjust pricing strategies and better manage demand. The travel industry as a whole can benefit from improved fare transparency and a more efficient marketplace."

## III. Comparison with Existing Research

Relate your findings to existing research on flight fare prediction. Discuss similarities, differences, and any novel contributions of your work. Cite relevant studies and highlight how your research builds upon or extends previous knowledge.

## Example

"Previous studies have explored various machine learning techniques for flight fare prediction, including [mention specific methods]. Our research builds upon this work by [mentioning specific contributions, like improved accuracy, novel feature engineering techniques, or addressing specific limitations of previous work]. The findings align with previous research suggesting the importance of factors like [mention key features identified in other studies], while also revealing new insights into the role of [mention novel features identified in your research]."

## IV. Limitations and Future Directions

Acknowledge any limitations of your study, such as data limitations, assumptions made, or potential biases. Discuss how these limitations might affect the generalizability of your findings. Propose future research directions to address these limitations and further enhance the model's performance or applicability.

## Example

"While this study provides valuable insights, it's essential to acknowledge limitations. The model was trained on a specific dataset and may not generalize perfectly to all scenarios. Future research could explore incorporating external data sources, such as weather or fuel prices, to potentially improve the model's predictive power. Investigating the impact of real-time factors like seat availability and booking patterns could further enhance the model's accuracy."

## V. Concluding Remarks

Conclude the discussion by summarizing the overall significance of your research and its potential impact. Reiterate the key contributions and their implications for the field of flight fare prediction.

## Example

"This research demonstrates the effectiveness of machine learning techniques for predicting flight fares. The developed model provides valuable insights for travellers, airlines, and the travel industry, empowering informed decision-making and potentially leading to cost savings and improved travel planning. Future research building upon these findings can further refine prediction accuracy and address the dynamic nature of flight pricing."

## 7. Conclusion

In this project, we built a model to predict flight fares using machine learning. We used a dataset of flight fares and applied various data preprocessing techniques, including handling categorical data and feature selection. We then trained a Random Forest Regressor model and evaluated its performance using metrics such as R-squared score, MAE, MSE, and RMSE. We achieved an R-squared score of approximately 81%, indicating a good fit of the model to the data. Hyperparameter tuning was performed to further improve the model's performance. The results demonstrate that machine learning techniques can be effectively applied to predict flight fares, providing valuable insights for both travellers and airlines. Future work could involve exploring other algorithms and incorporating additional features for improved prediction accuracy.

### ❖ Reasoning

**Summarize the project's goal:** Start by restating the objective of your project, which was to predict flight fares using machine learning.

**Highlight the key steps:** Briefly mention the important stages of your project, such as data preprocessing, model selection, and evaluation.

**State the results:** Present the main findings, focusing on the performance metrics you used, such as R-squared score.

**Discuss the implications:** Explain the significance of your results, emphasizing the practical value of your model in predicting flight fares.

**Suggest future work:** Identify potential areas for improvement or further research, such as experimenting with different algorithms or adding more data features.

## 8.Future Work

The "Future Work" section of your research paper allows you to acknowledge the limitations of your current study and propose potential avenues for further exploration and improvement. It demonstrates a forward-thinking approach and highlights opportunities for advancing the field.

Here are some specific examples of what you can include in the "Future Work" section of your flight fare prediction project research paper:

### I. Incorporating Real-time Data

**Idea:** Integrate real-time data sources, such as current flight availability, weather conditions, and fuel prices, into the prediction model.

**Rationale:** Real-time data can significantly improve the accuracy of predictions as they reflect the dynamic nature of flight fares.

**Example:** "Future work could focus on incorporating real-time flight availability data from airline APIs to enhance the model's responsiveness to market fluctuations."

### II. Exploring Advanced Algorithms

**Idea:** Experiment with other machine learning algorithms, such as deep learning models or gradient boosting, to potentially achieve better predictive performance.

**Rationale:** Different algorithms may be better suited for handling the complexities of flight fare data and capturing non-linear relationships.

**Example:** "Investigating the application of deep learning architectures, like recurrent neural networks (RNNs), for capturing temporal patterns in fare data could be a promising direction for future research."

### III. Expanding Feature Set

**Idea:** Include additional relevant features, such as competitor pricing, social media sentiment, and economic indicators, to further refine the model.

**Rationale:** Adding more informative features can enhance the model's ability to explain fare variations and improve its predictive power.

**Example:** "Future studies could explore incorporating social media sentiment analysis to gauge public perception of airlines and its impact on fare prices."

### IV. Addressing Data Imbalance

**Idea:** Employ techniques to handle class imbalance in the dataset, such as oversampling or under sampling, to improve the model's performance on minority classes.

**Rationale:** Data imbalance can bias the model towards the majority class, leading to poor predictions for less frequent fare categories.

**Example:** "Applying techniques like Synthetic Minority Over-sampling Technique (SMOTE) to address the imbalance in fare categories could enhance the model's predictive accuracy for less common price ranges."

### V. Deployment and Real-world Application

**Idea:** Develop a user-friendly interface or API for deploying the model in a real-world setting, allowing travellers or airlines to access its predictions.

**Rationale:** Translating research findings into practical applications can make the model more accessible and impactful.

**Example:** "Future efforts could involve building a web application or mobile app that utilizes the prediction model to provide users with real-time flight fare estimates and personalized travel recommendations".

## 9. References

- Pandas: McKinney, W. (2011). pandas: a foundational Python library for data analysis and statistics. *Python for High Performance and Scientific Computing*, 14(9), 1-9.
- \* NumPy: Harris, C. R., Millman, K. J., van der Walt, S. J., Gommers, R., Virtanen, P., Cournapeau, D., ... & Oliphant, T. E. (2020). Array programming with NumPy. *Nature*, 585(7825), 357-362.
- \* Seaborn: Waskom, M. L. (2021). seaborn: statistical data visualization. *Journal of Open Source Software*, 6(60), 3021.
- Random Forest Regression: Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5-32.
- RandomizedSearchCV/GridSearchCV: Bergstra, J., & Bengio, Y. (2012). Random search for hyper-parameter optimization. *Journal of Machine Learning Research*, 13(Feb), 281-305.
- Airfare Price Prediction Using Machine Learning Rathore, Y., & Tripathi, A. (2020). Airfare Prediction using Machine Learning Techniques. *International Journal of Computer Applications*, 975, 8887.
- Feature Engineering Best Practices Zheng, A., & Casari, A. (2018). *Feature Engineering for Machine Learning: Principles and Techniques for Data Scientists*. O'Reilly Media.
- Performance Metrics in Regression Models Chai, T., & Draxler, R. R. (2014). Root mean square error (RMSE) or mean absolute error (MAE)? – Arguments against avoiding RMSE. *Geoscientific Model Development*, 7(3), 1247–1250.

