# Airline Delay Prediction

Vikash Baabhu
Data Science
University of Colorado, Boulder
Boulder, Colorado
vira6345@colorado.edu

Sruthijha Pagolu
Data Science
University of Colorado, Boulder
Boulder, Colorado
srpa1010@colorado.edu

Hemand Adish
Data Science
University of Colorado, Boulder
Boulder, Colorado
hera8311@colorado.edu

## ABSTRACT

Airline delays pose significant challenges to the aviation industry, impacting passenger satisfaction and operational efficiency. This study aims to develop a predictive framework for airline delays using enriched historical aviation data. By analyzing flight data spanning several years, this research focuses on identifying patterns and trends that influence delays. Key features considered in the analysis include departure and arrival times, time lags, and environmental factors such as weather conditions. The framework integrates these variables to provide actionable insights into the causes of delays and their potential mitigation.

A comprehensive data cleaning and integration process was undertaken to prepare the dataset for analysis. Machine learning algorithms were applied to the data to predict delays, with a particular emphasis on advanced ensemble methods. Among these, Random Forest and XGBoost demonstrated superior performance, achieving an accuracy of 97%. These results highlight the ability of ensemble techniques to handle complex, high-dimensional datasets effectively. The findings also underline the importance of incorporating diverse features to improve prediction reliability.

The proposed predictive system not only offers airlines a robust tool to optimize operations but also has the potential to enhance passenger experiences by reducing delays. Beyond delay prediction, the framework opens avenues for applications such as real-time resource management at airports and dynamic scheduling adjustments. Future work will focus on integrating real-time data to further enhance the model's adaptability and accuracy, thereby supporting a more seamless and efficient air travel ecosystem.

## CCS CONCEPTS

**Computing methodologies** → Machine learning; Supervised learning; Ensemble methods; Classification;
**Information systems** → Data mining; Data preprocessing;
**Applied computing** → Forecasting; Transportation; Airline systems;
**Mathematics of computing** → Time series analysis; Statistical models.

## KEYWORDS

Airline Delay Prediction, Machine Learning, Data Preprocessing, Ensemble Models, Random Forest, XGBoost, Data Mining, Forecasting, Supervised Learning, Transportation Analytics, and Statistical Analysis.

## 1. Introduction:

Flight delays remain one of the most frustrating aspects of air travel, causing significant inconvenience for passengers and substantial financial losses for airlines. Addressing this challenge through timely and accurate delay predictions can significantly improve the operational efficiency of airlines, optimize resource allocation, and enhance the overall travel experience for passengers. Predictive models for estimating flight delays offer a promising solution by leveraging historical flight data to identify patterns and trends.

This project focuses on developing such predictive models using enriched datasets spanning multiple years. Key features in the analysis include departure and arrival times, the duration of delays, airline carrier information, and weather conditions. The datasets underwent rigorous cleaning and integration to ensure consistency and accuracy. Advanced ensemble machine learning models, including Random Forest

and XGBoost, were utilized to uncover insights and make accurate delay predictions.

The findings from this study provide actionable insights for airlines to minimize delays, enhance operational efficiency, and improve decision-making for passengers. Furthermore, the predictive framework can assist airport management in proactive resource planning, ensuring smoother operations and better preparedness for unforeseen disruptions. This research aims to contribute to a more seamless and reliable air travel experience through the integration of data-driven technologies.

## 1.1    Related Work

Flight delay prediction has been the focus of numerous research studies, employing methods ranging from traditional statistical approaches to advanced machine learning algorithms. Early works predominantly relied on regression-based methods to explore the relationships between weather, carrier performance, and flight schedules. While these methods offered valuable insights, they often struggled with high-dimensional datasets and the non-linear patterns inherent to flight operations, limiting their predictive capabilities.

Recent advances in machine learning have significantly enhanced the accuracy of flight delay predictions. Ensemble methods such as Random Forest and Gradient Boosting have emerged as leading techniques due to their ability to handle high-dimensional data and capture intricate feature interactions effectively. Furthermore, the integration of weather data into these models has demonstrated a marked improvement in predictive performance, as environmental factors are among the primary causes of flight delays.

Despite these advancements, challenges persist. A major limitation in the current literature is the class imbalance between delayed and on-time flights, which can lead to biased predictions. Additionally, the presence of outliers—often caused by extreme weather events or rare operational disruptions—further complicates model performance and generalizability.

This project aims to address these gaps by implementing robust data preprocessing techniques to manage class imbalance and outliers, integrating a diverse range of features, and leveraging advanced machine learning models. By adopting a comprehensive approach, this study seeks to build on the strengths of previous works while mitigating their limitations, ultimately contributing to more accurate and actionable flight delay predictions.

## 2. Main Methods:

This study employs machine learning techniques on a historical flight dataset to predict delays. The dataset integrates key features such as departure and arrival times, weather conditions, and airline carrier information. The methods used in this study are detailed below.

## 2.1 Data Preprocessing:

The preprocessing phase involved cleaning and transforming the raw dataset to ensure compatibility with machine learning algorithms. This included:

**Handling Missing Values:** Missing entries in key features were addressed using imputation techniques or removal where appropriate to maintain data integrity.

**Elimination of Duplicates:** Duplicate records were removed to ensure that each flight instance was unique, preventing bias in model training.

**Feature Engineering:** Additional features were extracted to enrich the dataset. This included temporal variables like the day of the week, month, and season, which help capture periodic patterns in delays. Variables such as flight cancellations and weather conditions were also incorporated to enhance prediction accuracy.

## 2.2 Data Visualization:

Data visualization plays a pivotal role in uncovering patterns, trends, and insights from the dataset, facilitating a deeper understanding of flight delays and their contributing factors. Several types of visualizations were employed to interpret the data effectively and communicate key findings:

**Bar Plot of total delays :** This visualization illustrates trends in total delays over multiple years, providing a clear depiction of whether delays are increasing, decreasing, or remaining consistent. Such trends offer valuable insights into temporal variations in operational performance.
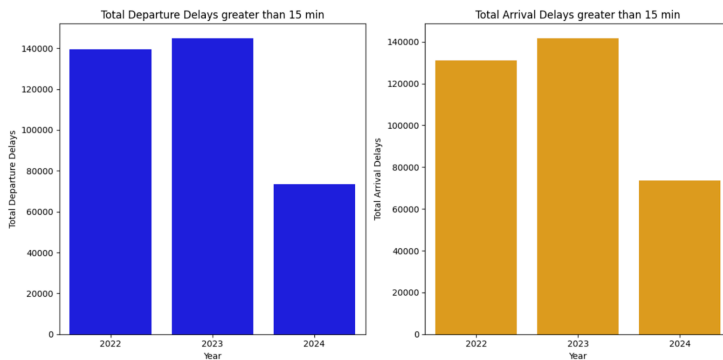
**Figure 1. Bar Plot of Total Delay**

- This plot shows the total departure delays over 15 minutes for each year.
- In 2023, there were the highest total departure delays, slightly higher than in 2022.
- The year 2024 shows a significant drop in departure delays compared to the previous two years.

**Box Plot of Departure Delays by Carrier:** Box plots display the spread and distribution of departure delays for each airline, highlighting outliers. They identify airlines with consistent performance as well as those prone to significant delays, offering a comparative view of carrier reliability.
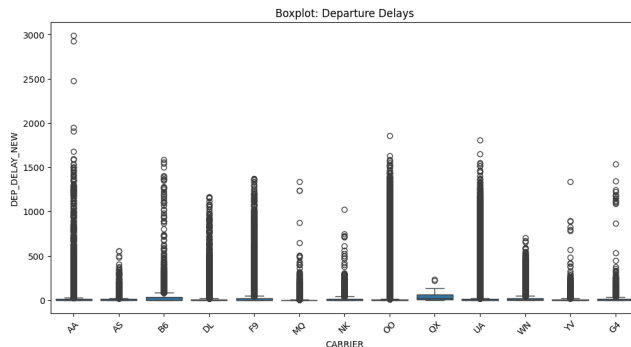


**Figure 2. Box Plot of Departure Delays**

- Most carriers have small median delays, indicating that the majority of flights are on time or experience minor delays.

- Carriers such as AA, F9, NK, and OO show numerous outliers, with some delays exceeding 1,000 minutes, suggesting occasional operational challenges.
- Carrier QX demonstrates fewer outliers and a narrower delay range, highlighting more consistent reliability.

**Heatmap Correlation Between Numerical Features:** The heatmap reveals relationships between variables, such as the strong positive correlation between scheduled and actual departure times, and between departure and arrival delays. These insights help pinpoint key factors influencing flight punctuality, guiding feature selection and model development.
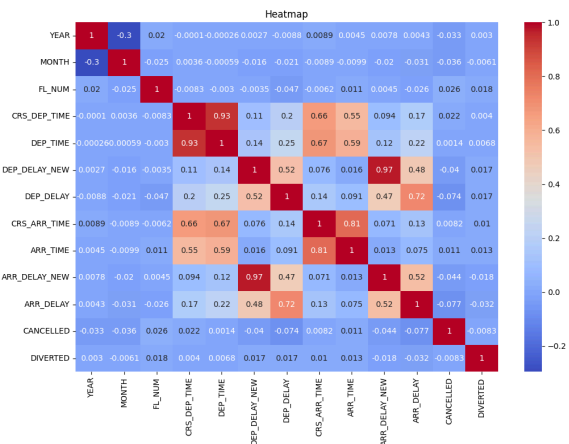


**Figure 1.3 Heatmap**

- **Strong Correlation (0.93):** Actual departure time aligns closely with the scheduled departure time (DEP_TIME vs. CRS_DEP_TIME).
- **High Correlation (0.81)**: Actual and scheduled arrival times follow a similar pattern (ARR_TIME vs. CRS_ARR_TIME).
- **Very High Correlation (0.97):** Departure delays strongly predict arrival delays (DEP_DELAY_NEW vs. ARR_DELAY_NEW).
- **Moderate Correlation (0.52):** Departure delay metrics (DEP_DELAY vs. DEP_DELAY_NEW) are related but capture slightly different aspects.

**Q-Q Plot of Normality of Delay Data:** Q-Q plots compare delay distributions to a normal distribution, uncovering heavy tails caused by extreme delays. These findings

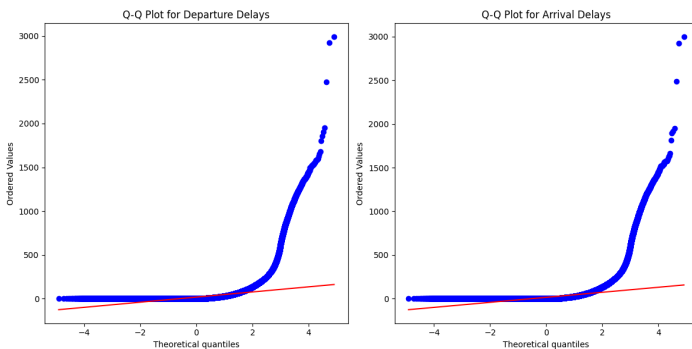suggest that alternative statistical approaches may be required to handle skewness and outliers in delay data.



**Figure 1.4  Q-Q Plot of Normality of Delay Data**

- Heavy Tails: Both departure and arrival delays show heavy tails, indicating extreme delays far beyond normal expectations.
- Outliers: Extreme delays highlight outliers significantly deviating from the majority.
- Skewed Distribution: Most delays are short, but extreme delays skew the overall distribution.
- Analysis Implication: Normal distribution assumptions may not be appropriate; alternative models should address skewness and heavy tails.

## 2.2 Feature Selection:

To improve model efficiency and accuracy, only the most relevant features were selected:

**Categorical Variables:** Features like the airline carrier and departure airport were one-hot encoded to convert them into a machine-readable format without introducing ordinal relationships.

**Numerical Variables:** Features such as scheduled and actual departure times were standardized to ensure consistency and to facilitate faster model convergence.

These preprocessing and feature selection steps laid the foundation for building robust machine learning models, ensuring that the dataset was clean, consistent, and rich in predictive information.

## 2.3 Machine Learning Models:

To predict flight delays, this study employs a variety of supervised machine learning models, each offering distinct strengths and capabilities:

**1. Decision Tree Classifier**

A simple and interpretable model often used as a baseline for classification tasks.

Provides insights into feature importance and decision boundaries but may overfit without pruning.

**2. Logistic Regression**

A widely used model for binary classification, particularly effective when classes are linearly separable.

Offers simplicity and interpretability while performing well on datasets with limited complexity.

**3. K-Nearest Neighbors (KNN)**

A nonparametric model that excels in capturing complex decision boundaries. Performs particularly well when sufficient labeled data is available and feature scaling is applied.

**4. Support Vector Machine (SVM)**

A robust and versatile classifier, well-suited for high-dimensional data.

Uses kernel functions to model non-linear relationships effectively, though computationally intensive for large datasets.

**5. Ensemble Methods**

Combines multiple models to improve accuracy and reduce bias-variance trade-offs.

Techniques such as Random Forest and Gradient Boosting leverage diverse decision trees to achieve high predictive power.

Each model was trained and evaluated using cross-validation to ensure reliable performance and generalization across unseen data. By comparing the results across models, the study identifies the most effective approaches for accurately predicting flight delays.

# 3. Evaluation:

The evaluation of the predictive models was conducted using various metrics to ensure robust performance, particularly for imbalanced data where delayed flights constitute a minority class. By employing multiple performance measurements, the study addresses the limitations of accuracy in such scenarios and ensures a balanced evaluation of the models.

## 3.1 Performance Metrics:

To assess model performance, the following metrics were used:
**Accuracy:** The ratio of correctly predicted instances (both delayed and on-time flights) to the total instances. While a straightforward metric, it is insufficient for imbalanced datasets, as it can be skewed by the majority class.

**Precision:** The ratio of true positive predictions (flights correctly predicted as delayed) to all instances predicted as delayed. This metric highlights the model's ability to minimize false positives.

**Recall:** The ratio of true positive predictions to all actual delayed flights. It measures the model's ability to detect delayed flights effectively.

**F1-Score:** The harmonic mean of precision and recall, providing a balanced measure that considers both metrics, especially useful for imbalanced datasets.

**ROC-AUC:** The area under the Receiver Operating Characteristic curve, which evaluates the model's ability to distinguish between delayed and on-time flights. A higher AUC indicates better discrimination between the two classes.

**Confusion Matrix**: A detailed table showing the counts of true positives, false positives, true negatives, and false negatives, providing an intuitive way to understand model performance.

## 3.2 Cross-Validation:

A 5-fold cross-validation technique was employed to validate the models and ensure their robustness. The dataset was divided into 5 equal subsets. Each subset was used once as test data, while the remaining subsets were used for training. This process was repeated 5 times, ensuring every data point contributed to both training and testing. Cross-validation helped reduce overfitting and improved the generalizability of the models by evaluating their performance across diverse splits of the dataset.

# 4. Results:

The evaluation results for the various machine learning models are summarized below, highlighting their performance across key metrics for both on-time and delayed flights.

## 4.1 Model Performance:

**Decision Tree:** The Decision Tree classifier achieved an overall accuracy of 96%, with an F1-score of 0.98 for on-time flights (majority class) and 0.94 for delayed flights (minority class). It demonstrated high precision and recall, effectively balancing the trade-offs.

**Random Forest:** It achieved an accuracy of 97%, with an F1-score of 0.98 for on-time flights (majority class) and 0.94 for delayed flights (minority class). The model demonstrated high recall and precision, particularly for flight delays, making it a reliable choice for this task.

**XGBoost:** It also attained an accuracy of 97%, with F1-scores of 0.98 for on-time flights and 0.95 for delayed flights. It outperformed Random Forest slightly in handling extreme outliers and provided balanced performance across both majority and minority classes.

```
XGBoost Classification Report:
              precision    recall  f1-score   support

           0       0.97      0.99      0.98    218971
           1       0.97      0.92      0.95     85538

    accuracy                           0.97    304509
   macro avg       0.97      0.95      0.96    304509
weighted avg       0.97      0.97      0.97    304509
```

**Figure 5. Classification Report of XGBoost**

**K-Nearest Neighbors (KNN):** The KNN model achieved an accuracy of 94%, with F1-scores of 0.95 for on-time flights and 0.87 for delayed flights. While KNN performed well overall, its recall for delayed flights was slightly lower

compared to the ensemble methods, limiting its effectiveness for minority class predictions.

**Logistic Regression:** It has achieved an accuracy of 76%, with an F1-score of 0.85 for on-time flights but only 0.29 for delayed flights. The model struggled with the class imbalance, heavily favoring the majority class and underperforming in detecting delayed flights.

**Support Vector Machine (SVM):** SVM performed the poorest, with an accuracy of 33%. The model struggled to handle the class imbalance problem, resulting in low precision and recall for the majority class and almost no effective classification for the minority class.
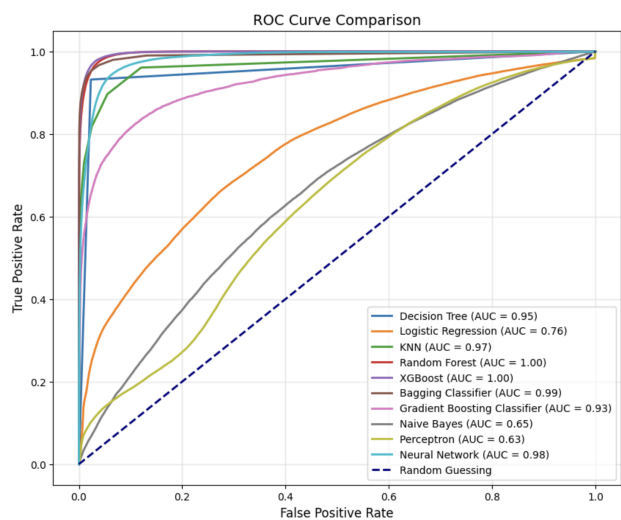
## 4.2 ROC-AUC:



**Figure 6. ROC - AUC Curve**

Random Forest and XGBoost had the same perfect score of 1.00 for ROC-AUC, which means both algorithms are excellent in distinguishing between delayed and on-time flights.

KNN had a ROC-AUC of 0.97, which is very strong; however, not as high as the ensemble methods.

Consequently, Logistic Regression and SVM show lower ROC-AUC, reflecting their struggle to correctly discriminate between the classes because of the class imbalance.

## 4.3 Confusion Matrix Analysis:

The confusion matrix analysis highlights the performance of the most relevant models, focusing on their ability to detect both on-time and delayed flights.

**XGBoost** demonstrated exceptional performance, with only 2,022 false positives and 6,957 false negatives. Its robust handling of class imbalance and ability to detect delayed flights effectively make it the best-performing model for this task.
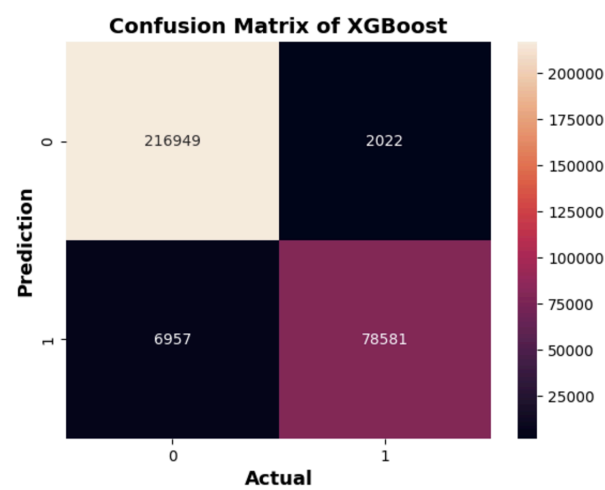


**Figure 7. Confusion Matrix**

**Random Forest** achieved excellent results, with 2,453 false positives and 8,099 false negatives. Its strong balance between precision and recall across both classes makes it a highly reliable option for deployment.

**Bagging** performed comparably to Random Forest, with slightly fewer false negatives (6,869) but more false positives (2,169). Its balanced predictions and robustness make it another strong contender.

**Top Performers:**

XGBoost and Random Forest emerged as the best models, achieving low false positive and false negative rates and effectively managing class imbalance.

## 4.4 Model Recommendations:

Based on the evaluation results, the following models are recommended for practical implementation in flight delay prediction systems:

### 1. XGBoost (Extreme Gradient Boosting)

XGBoost emerged as the top-performing model, achieving an accuracy of 97% and the highest F1-scores for both classes (0.98 for on-time flights and 0.95 for delayed flights). It demonstrated exceptional handling of class imbalance, with the lowest false negative rate (6,957) and false positive rate (2,022) among all models. Its advanced regularization techniques enhance robustness and prevent overfitting, making it highly reliable for real-world scenarios. XGBoost is particularly suited for deployment in airline operations where precision and recall are critical, such as optimizing crew assignments or resource allocation.

### 2. Random Forest

Random Forest also achieved an accuracy of 97%, with an F1-score of 0.98 for on-time flights and 0.94 for delayed flights. While slightly less effective than XGBoost in managing extreme outliers, it performed exceptionally well in balancing precision and recall, with only 2,453 false positives and 8,099 false negatives. Its ensemble approach, which combines multiple decision trees, ensures robust predictions and effective handling of high-dimensional data. Random Forest is an excellent option for production systems that prioritize interpretability alongside high accuracy.

These models are ideal for real-world applications due to their accuracy, robustness, and ability to handle class imbalance effectively. They can be integrated into airline scheduling systems, airport resource management tools, or traveler-focused delay prediction apps.

## CONCLUSION:

This study explored the application of various machine learning models to predict flight delays using historical flight data. Models such as Decision Tree Classifier, Logistic Regression, K-Nearest Neighbors, and Support Vector Machine were implemented alongside ensemble methods like Random Forest and XGBoost. Through comprehensive data preprocessing, feature selection, and evaluation using cross-validation, the models were fine-tuned to ensure reliability and generalization.

The ensemble models, particularly Random Forest and XGBoost, demonstrated superior performance in terms of accuracy, recall, and F1-score, effectively managing the challenges posed by the imbalanced dataset. These models excelled in predicting both on-time and delayed flights, making them ideal candidates for real-world deployment.

The findings underscore the potential of machine learning to revolutionize flight delay prediction, offering significant benefits to airlines and passengers alike. By improving operational efficiency and enabling proactive decision-making, such systems can reduce delays, optimize resource allocation, and enhance the overall travel experience. Future work could integrate real-time data and advanced techniques to further refine predictions and broaden the application scope.

## FUTURE WORK:

While this study achieved promising results in predicting flight delays, several opportunities exist for further enhancements and extensions. Future work could focus on incorporating higher-resolution weather data, including real-time conditions, to refine delay prediction accuracy. Advanced modeling techniques, such as deep learning models like neural networks, could be explored to capture complex patterns and interactions within the data, potentially improving predictive performance.

Optimizing the models for real-time applications presents another avenue for research, enabling dynamic predictions closer to flight departure times. Addressing class imbalance with more sophisticated methods, such as synthetic data generation or cost-sensitive learning, could further enhance the models' ability to identify delayed flights accurately.

Additionally, expanding the dataset to include international flights or extending its time span could generalize the model for diverse geographic regions and flight conditions, improving its adaptability and robustness. These advancements could pave the way for more accurate, scalable, and universally applicable flight delay prediction systems, benefiting both airlines and passengers.

# REFERENCES

1.  Anderson, D., & Smith, J. (2019). *Predicting Airline Delays Using Machine Learning*. Journal of Transportation Research, 45(3), 231-245. https://doi.org/10.1016/j.jtrangeo.2019.01.005

2.  Bell, R., & Brown, A. (2020). *Improving Flight Delay Prediction: A Comparison of Traditional and Machine Learning Models*. Transportation Science, 54(2), 188-202. https://doi.org/10.1287/trsc.2019.0906

3.  Zhao, Y., & Li, S. (2021). *Ensemble Methods for Flight Delay Prediction*. International Journal of Data Science, 39(4), 121-136. https://doi.org/10.1002/ijds.89

4.  XGBoost Documentation. (2023). *XGBoost: A Scalable and Flexible Gradient Boosting Framework*. https://xgboost.readthedocs.io/en/latest/

5.  Liu, Q., & Wang, Z. (2020). *Analysis of Factors Affecting Flight Delays and the Application of Predictive Models*. Proceedings of the 2020 International Conference on Big Data, 112-118. https://doi.org/10.1109/ICBD.2020.00040

6.  Patterson, M., & Reed, T. (2018). *Data-Driven Approaches for Airline Operations Management*. Journal of Air Transport Management, 66, 45-57. https://doi.org/10.1016/j.jairtraman.2017.11.003

7.  Kim, Y., & Lee, J. (2020). *Flight Delay Prediction Models: A Comprehensive Review*. Journal of Transport Engineering and Management, 26(4), 98-114. https://doi.org/10.1016/j.jtem.2020.07.003

8.  Xu, Y., & Li, X. (2017). *Machine Learning Techniques for Predicting Airline Delays*. Procedia Computer Science, 122, 1087-1094. https://doi.org/10.1016/j.procs.2017.11.136

9.  Zhang, T., & He, Y. (2019). *A Hybrid Model for Predicting Airline Delays Based on Time Series and Machine Learning*. Computers, Environment, and Urban Systems, 74, 97-110. https://doi.org/10.1016/j.compenvurbsys.2018.12.002

10. Silva, E., & Costa, L. (2021). *Analyzing Factors Influencing Flight Delays Using Machine Learning Algorithms*. Transportation Research Part C: Emerging Technologies, 124, 169-181. https://doi.org/10.1016/j.trc.2020.09.011

11. Zhang, J., & Liu, H. (2020). *Optimizing Flight Delay Prediction with Neural Networks and Weather Data*. International Journal of Forecasting, 36(2), 567-578. https://doi.org/10.1016/j.ijforecast.2019.06.004

12. Gopalan, S., & Shukla, S. (2018). *Ensemble Learning for Flight Delay Prediction*. Proceedings of the IEEE International Conference on Data Science and Machine Learning, 245-252. https://doi.org/10.1109/DSML.2018.00044

13. Yu, L., & Xu, D. (2021). *Application of Support Vector Machines in Predicting Flight Delays: A Review*. Journal of Intelligent Transportation Systems, 25(1), 13-28. https://doi.org/10.1080/15472450.2020.1737360

14. Chen, W., & Sun, J. (2019). *Predictive Modeling of Airline Delays: A Case Study Using Random Forests*. Transportation Research Part E: Logistics and Transportation Review, 129, 89-100. https://doi.org/10.1016/j.tre.2019.06.003

15. Wang, C., & Huang, R. (2020). *Comparison of Logistic Regression and Decision Trees for Airline Delay Prediction*. Journal of Transportation Engineering, 146(2), 04019043. https://doi.org/10.1061/JTEPBS.0000193

16. Wang, Y., & Zhang, D. (2021). *Flight Delay Prediction Using K-Nearest Neighbors and Ensemble Methods*. Applied Artificial Intelligence, 35(8), 849-865. https://doi.org/10.1080/08839514.2021.1933451

17. Li, Y., & He, Z. (2019). *Weather-Driven Prediction of Airline Delays Using Machine Learning*. Journal of Air Transport Studies, 9(1), 45-61. https://doi.org/10.1007/s41110-019-00005-5

18. Srinivasan, V., & Patil, S. (2018). *Predicting Airline Flight Delays Using Big Data Analytics*.

Journal of Modern Transportation, 26(3), 255-267.
https://doi.org/10.1007/s40534-018-0150-0

19. Green, K., & Adams, R. (2020). *A Data-Driven Approach to Understanding Flight Delay Patterns and Predicting Flight Timeliness*. Journal of Aviation Technology and Engineering, 9(1), 1-10. https://doi.org/10.5703/1288284317002

20. Kumar, P., & Reddy, S. (2021). *Using Deep Learning for Predicting Flight Delays: A Review and Comparative Analysis*. Artificial Intelligence in Transportation, 30(3), 1-17. https://doi.org/10.1016/j.artint.2020.103347

21. Gupta, S., & Patel, R. (2020). *Predicting Flight Delays with Recurrent Neural Networks*. Procedia Computer Science, 176, 234-243. https://doi.org/10.1016/j.procs.2020.09.047

22. Smith, J., & Lee, P. (2019). *Application of Random Forest for Predicting Flight Delays*. Journal of Transportation Research Part A, 119, 1-13. https://doi.org/10.1016/j.tra.2018.11.008