

# Predicting flight delays using flight delay prediction

## **Problem Description :**

The objective of this project is to predict flight delays using machine learning techniques. The dataset used for this project is the Flight Delay Prediction dataset, which is available on Kaggle. The dataset contains information about airline flights departing from major airports in the United States between 2015 and 2018.

## **Dataset Description:**

The Flight Delay Prediction dataset contains a total of 7.3 million rows and 32 columns. Each row represents a single flight departing from a major airport in the United States between 2015 and 2018. The dataset includes information such as the airline, flight number, departure and arrival times, origin and destination airports, and various weather-related features.

## **Project Requirements and Deliverables:**

The requirements for this project include:

1. **Data preprocessing:** cleaning and transforming the raw data to make it suitable for machine learning models.
2. **Feature engineering:** selecting and engineering relevant features to improve the accuracy of the machine learning models.
3. **Model selection and hyperparameter tuning:** selecting appropriate machine learning models and tuning their hyperparameters to improve performance.

4. **Model training and evaluation**: training the selected machine learning models on the preprocessed data and evaluating their performance using appropriate metrics.
5. **Results interpretation and** communication: interpreting the results and communicating the findings to stakeholders in a clear and actionable way.

The deliverables for this project include a machine learning model that accurately predicts flight delays, as well as a report summarizing the data analysis and machine learning process.

# **General Framework and Steps**

The general framework for this project involves the following steps:

1. **Data collection and preprocessing:** Download and clean the Flight Delay Prediction dataset, removing any invalid or missing values and transforming the data as needed.
2. **Exploratory data analysis:** Explore the data to gain insights into the relationships between the different features and the target variable.
3. **Feature engineering:** Select and engineer relevant features that can improve the accuracy of the machine learning models.
4. **Model selection and hyperparameter tuning:** Choose appropriate machine learning models and tune their hyperparameters using cross-validation and grid search techniques.
5. **Model training and evaluation:** Train the selected machine learning models on the preprocessed data and evaluate their performance using appropriate metrics such as accuracy, precision, recall, and F1 score.
6. **Results interpretation and communication:** Interpret the results and communicate the findings to stakeholders in a clear and actionable way, including recommendations for future improvements.

By following this framework and implementing the necessary steps, we can build an accurate machine learning model that predicts flight delays with high confidence, helping airlines and passengers to plan their trips more effectively.

This framework outlines the step-by-step approach to building a machine learning model for predicting flight delays using the Flight Delay Prediction dataset.

### **Step 1: Data Collection and Preprocessing**

1.1 Download the Flight Delay Prediction dataset from Kaggle.

1.2

1.3 Clean the data by removing any invalid or missing values.

1.4

1.3 Transform the data as needed, such as converting categorical variables to numerical variables.

### **Step 2: Exploratory Data Analysis**

2.1 Explore the data to gain insights into the relationships between the different features and the target variable.

2.2 Visualize the data using various plots and graphs to identify trends and patterns.

2.3 Analyze the correlation between features and the target variable to select the most relevant features.

### **Step 3: Feature Engineering**

3.1 Select and engineer relevant features that can improve the accuracy of the machine learning models.

3.2 Extract additional features such as weather data and flight duration.

3.3 Normalize and standardize the features as needed.

## **Step 4: Model Selection and Hyperparameter Tuning**

4.1 Choose appropriate machine learning models such as Logistic Regression, Random Forest, Gradient Boosting, and Neural Networks.

4.2 Tune the hyperparameters of the selected models using cross-validation and grid search techniques.

4.3 Compare the performance of the different models and select the best-performing model.

## **Step 5: Model Training and Evaluation**

5.1 Train the selected machine learning model on the preprocessed data.

5.2

5.3 Evaluate the performance of the model using appropriate metrics such as accuracy, precision, recall, and F1 score.

5.4

5.3 Analyze the results and identify any areas for improvement.

## **Step 6: Results Interpretation and Communication**

6.1 Interpret the results and communicate the findings to stakeholders in a clear and actionable way.

6.2

6.3 Make recommendations for future improvements, such as additional features or improvements to the model architecture.

6.4

By following this framework and implementing the necessary steps, we can build an accurate machine learning model for predicting flight delays, which can be used by airlines and passengers to plan their trips more effectively.

## **Code Explanation :**

Here is the simple explanation for the code which is provided in the code.py file.

### **Section 1: Data Collection and Preprocessing**

In this section, we start by loading the Flight Delay Prediction dataset. We then remove any invalid or missing values from the dataset. After that, we convert the categorical variables into numerical variables using the factorize function. Finally, we split the dataset into features and target variable, which is like separating the chicken from the bone.

### **Section 2: Exploratory Data Analysis**

In this section, we use seaborn and matplotlib to create beautiful visualizations of the dataset. We create a heatmap of the correlation matrix, which is like a heat map of your body where you can see which parts are related to each other. We then plot a histogram of the target variable, which is like a graph of how often you are delayed in your daily life.

### **Section 3: Feature Engineering**

In this section, we extract additional features such as weather data and flight duration, which is like adding some toppings to your pizza to make it more flavorful. We then normalize and standardize the features using the StandardScaler function, which is like putting all your clothes in the washing machine to make them the same size.

### **Section 4: Model Selection and Evaluation**

In this section, we split the dataset into training and testing sets using the train\_test\_split function. We then select a machine learning model, which is like picking the right tool for the

job. We tune the hyperparameters of the selected model using the GridSearchCV function, which is like fine-tuning your car to get the best performance. Finally, we evaluate the performance of the model using the accuracy score, which is like giving your model a grade on how well it predicts flight delays.

To run this code, you need to have Python installed on your system along with the necessary libraries such as pandas, seaborn, and scikit-learn. You also need to have the Flight Delay Prediction dataset in CSV format in the same directory as the Python script. You can run the code by opening the script in your favorite Python IDE or running it on the command line using the command “python script.py”.

## **Future Work :**

### **Step 1: Data Augmentation**

One of the limitations of the current project is that the dataset only contains data for a limited time period. To overcome this limitation, we can augment the data by collecting more data for the same time period from different sources or by simulating additional data using data generation techniques.

### **Step 2: Model Selection and Evaluation**

In the current project, we used a single machine learning model to predict flight delays. However, there are many other models that we can use for this task such as neural networks, decision trees, and support vector machines. We can evaluate the performance of these models and select the best model based on the evaluation metrics.

### **Step 3: Ensemble Learning**

Another approach to improve the accuracy of the flight delay prediction is to use ensemble learning techniques such as bagging, boosting, and stacking. These techniques combine multiple models to make more accurate predictions.

### **Step 4: Time-Series Analysis**

The current project only considers the flight delays at a single point in time. However, flight delays may be correlated with other factors such as the time of day, day of the week, and season. We can use time-series analysis techniques such as ARIMA and LSTM to model these correlations and improve the accuracy of the predictions.

### **Step 5: Real-Time Prediction**



In the current project, we only predict flight delays for a historical dataset. However, it would be more useful to predict flight delays in real-time for airlines and passengers. We can use real-time data sources such as flight schedules, weather reports, and airport traffic to predict flight delays in real-time.

**To implement these future work steps, we can follow the following guide:**

1. To augment the data, we can collect additional data from different sources or simulate additional data using data generation techniques such as SMOTE or data imputation techniques.
2. To select and evaluate different machine learning models, we can use the same steps as in the current project but evaluate the performance of multiple models and select the best one.
3. To implement ensemble learning, we can combine multiple models using bagging, boosting, or stacking techniques and evaluate their performance.
4. To perform time-series analysis, we can use time-series modeling techniques such as ARIMA or LSTM and evaluate their performance on the dataset.
5. To implement real-time prediction, we can use real-time data sources such as flight schedules, weather reports, and airport traffic and build a real-time prediction model using the selected machine learning algorithm.

## **Exercise :**

**Try to answers the following questions by yourself to check your understanding for this project. If stuck, detailed answers for the questions are also provided.**

**1.What are some limitations of the current flight delay prediction project, and how can we overcome these limitations?**

Answer: One limitation of the current project is that the dataset only contains data for a limited time period. We can overcome this limitation by augmenting the data using additional sources or data generation techniques. Another limitation is that the current project only uses a single machine learning model. We can overcome this by using ensemble learning techniques to combine multiple models.

**2.What is ensemble learning, and how can it improve the accuracy of flight delay prediction?**

Answer: Ensemble learning is a technique that combines multiple machine learning models to improve the accuracy of predictions. It can be used for flight delay prediction by combining multiple models that make predictions based on different factors such as weather, time of day, and airport congestion.

**3.What is time-series analysis, and how can it improve the accuracy of flight delay prediction?**

Answer: Time-series analysis is a technique that models the correlation between data points over time. It can be used for flight delay prediction by modeling the correlation between flight delays and other factors such as time of day, day of the week, and season.

**4.What is real-time prediction, and why is it important for flight delay prediction?**

Answer: Real-time prediction is the ability to make predictions based on real-time data. It is important for flight delay prediction because flight delays can change rapidly based on factors

such as weather and airport congestion. Real-time prediction allows airlines and passengers to make more informed decisions based on the latest information.

**5.What are some real-time data sources that can be used for flight delay prediction, and how can they be incorporated into a prediction model?**

Answer: Some real-time data sources for flight delay prediction include flight schedules, weather reports, and airport traffic data. These sources can be incorporated into a prediction model by updating the model with the latest data as it becomes available, and using the updated model to make real-time predictions.