# Project Report



Session Fall 2022 – BSAI

## Submitted by:

- Safwan Hafeez       22I-0588

## Submitted to:

Professor Hassan Raza

**Department of Artificial Intelligence**

**National University of Computer**

**and Emerging sciences,**

**FAST**

## 1. Introduction

In this project, we aim to predict flight delays using machine learning techniques. The dataset consists of flight and weather information, with the goal of predicting whether a flight will be delayed and, if so, to what extent. We employ a series of preprocessing, feature selection, and machine learning models to achieve accurate predictions.

We use several tools, including:

- **Data Preprocessing**: Cleaning, handling missing values, and transforming data for model readiness.

- **Feature Engineering**: Mapping weather data to flight data and selecting top features based on correlation.

- **Modeling**: Using Random Forest Classifiers and Regressors for both classification (binary and multiclass) and regression tasks.

- **Hyperparameter Tuning**: Optimizing the performance of the models using RandomizedSearchCV.

- **Evaluation**: Assessing model performance through various metrics and checking Kaggle scores.

## 2. Data Acquisition and Preprocessing

### 2.1 Reading Flight and Weather Data

The first step involved reading the flight data, including both training and test datasets, and the weather data. The flight data contains key attributes such as departure and arrival times, flight delays, and other relevant details. The weather data includes environmental conditions (temperature, precipitation, etc.) at the time of the flight.
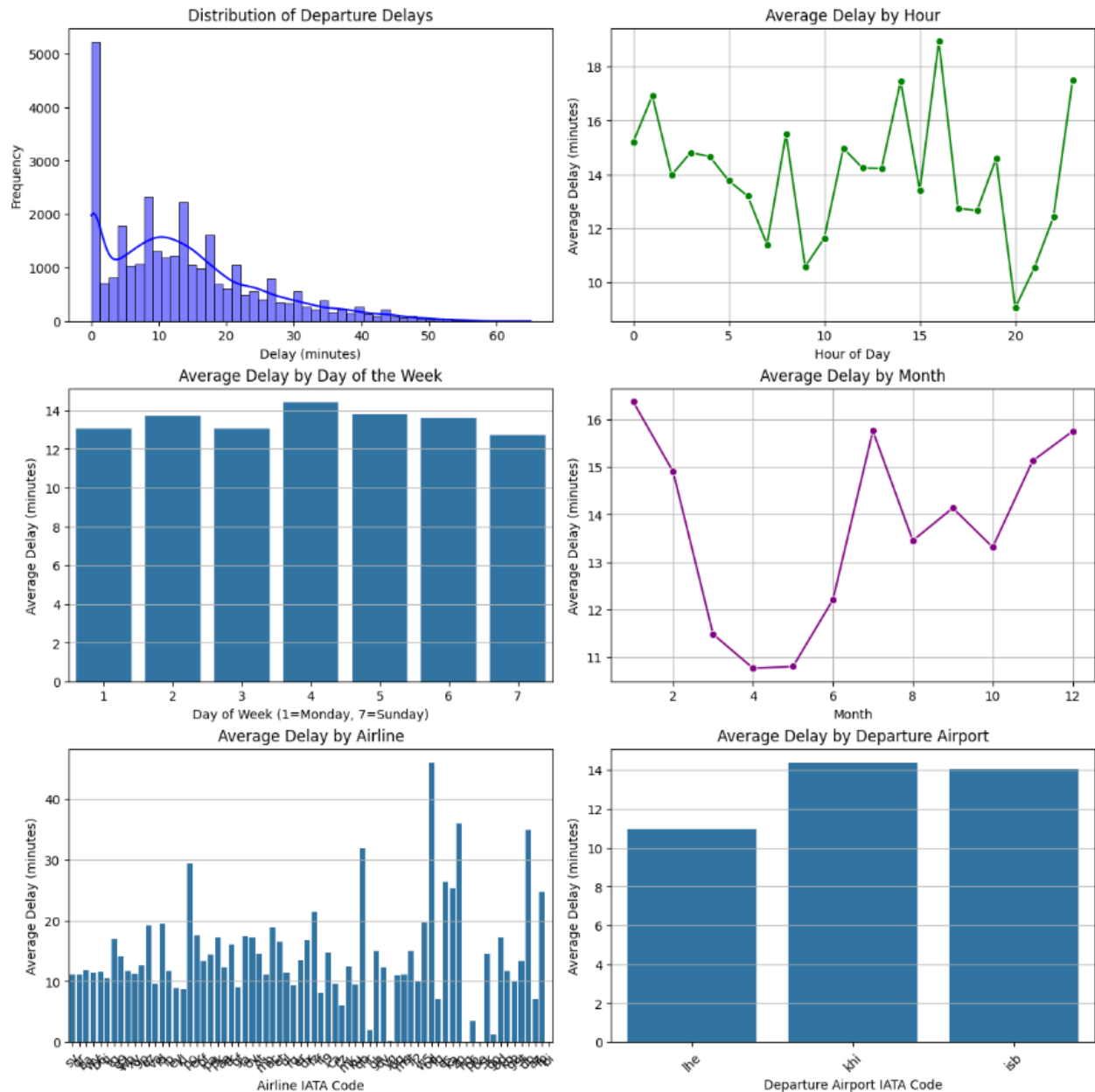
### 2.2 Data Cleaning and Preprocessing

- Missing values were handled by either imputation or removal, depending on the severity of the missing data.

- Categorical columns were encoded using One-Hot Encoding, and numerical columns were standardized using StandardScaler to ensure that all features contribute equally to the model.

- The weather data was mapped to both the training and test flight data by matching the appropriate weather attributes to the corresponding flight times.
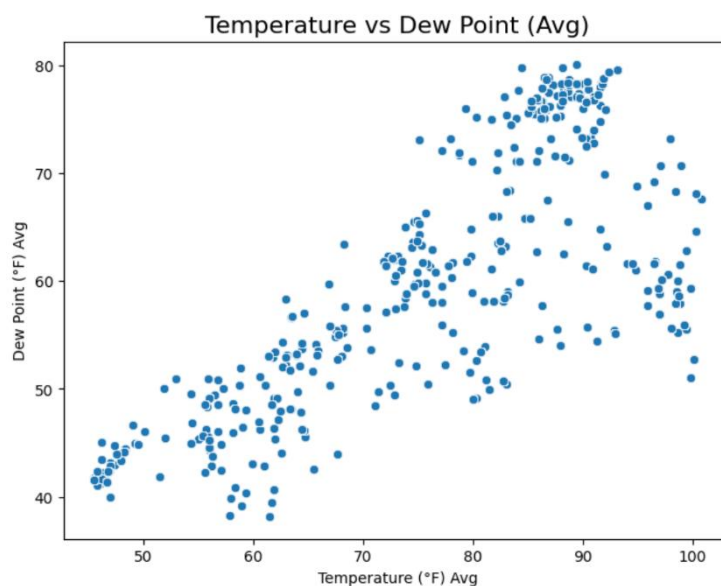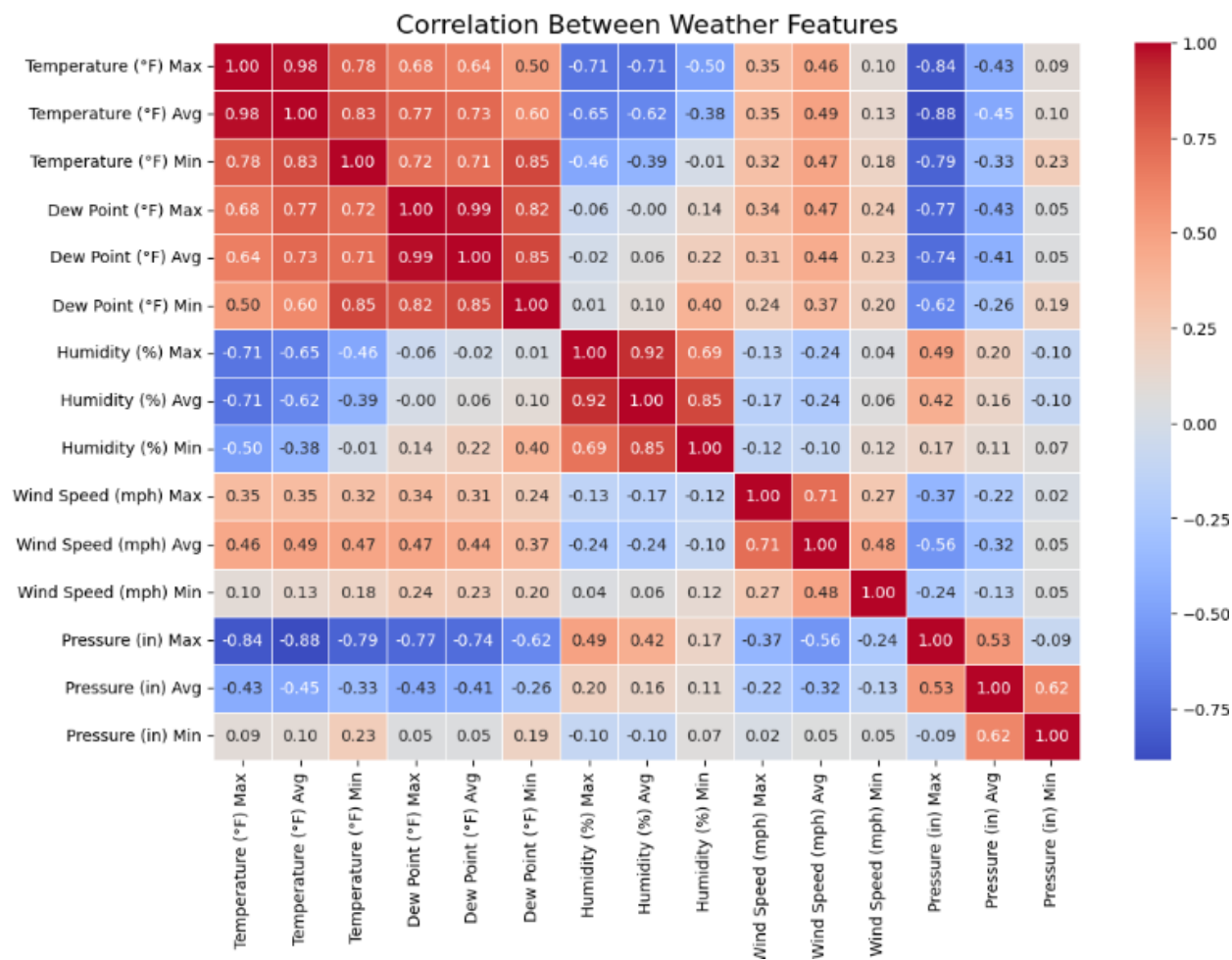
## 3. Feature Engineering

### 3.1 Selecting Relevant Features

To reduce dimensionality and enhance model performance, the top 5 most correlated features were selected based on correlation with the target variable (flight delay). Additionally, categorical variables such as day of the week, season, and weather-related features were retained for further analysis.

### 3.2 EDA and Visualization

Correlation Between Weather Features
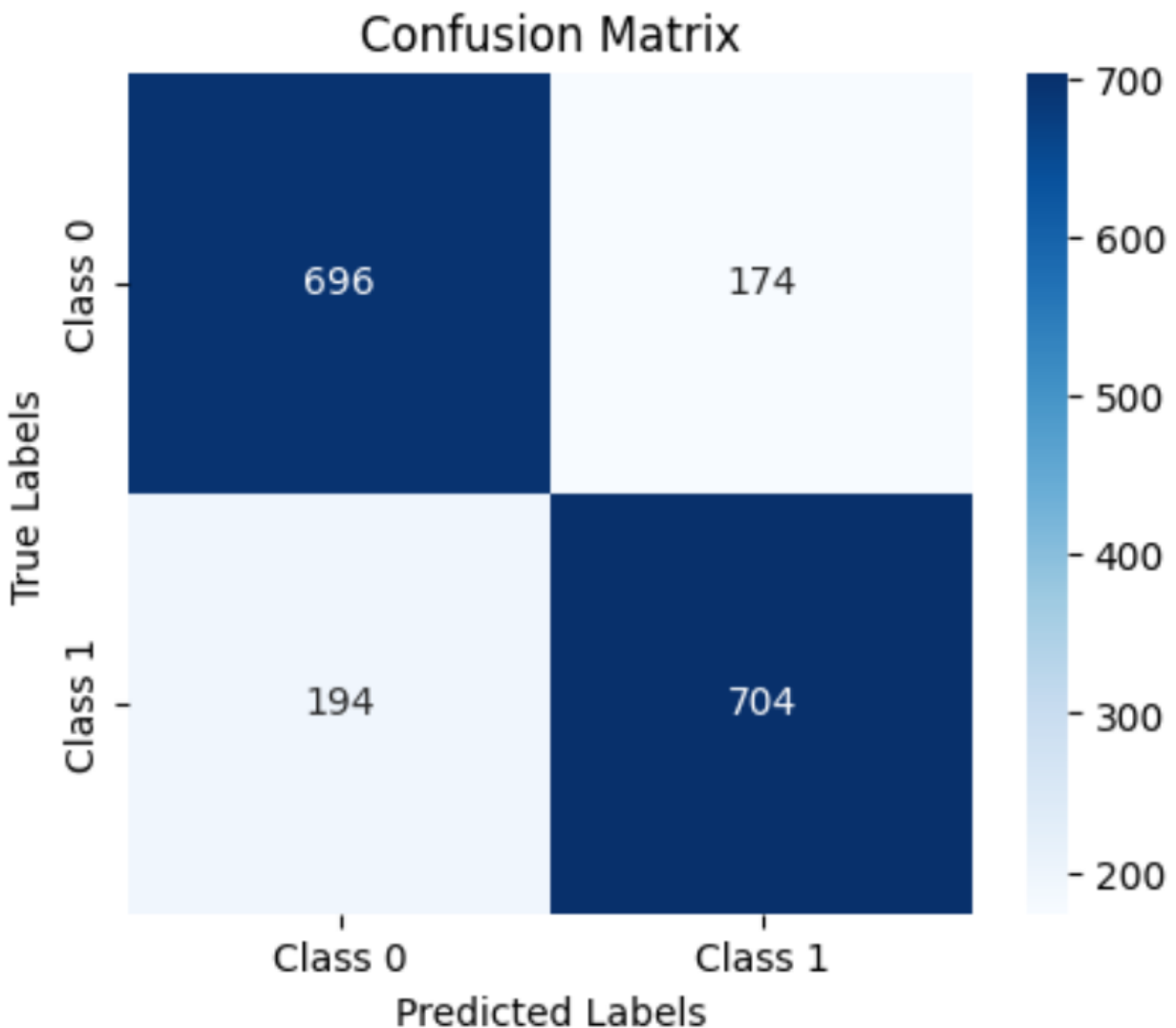


Temperature vs Dew Point (Avg)

Graphs like correlation matrices, feature distributions, and weather-related visualizations were generated to understand the relationships between different features and flight delays.

---

**4. Modeling and Evaluation**

**4.1 Random Forest Classifier for Binary Classification**

The first step in modeling was to build a binary classifier to predict whether a flight would be delayed or not. The Random Forest Classifier was chosen for this task due to its robustness and ability to handle both categorical and numerical features. The model was trained using preprocessed data, and the performance was evaluated using accuracy, precision, recall, and F1 score.
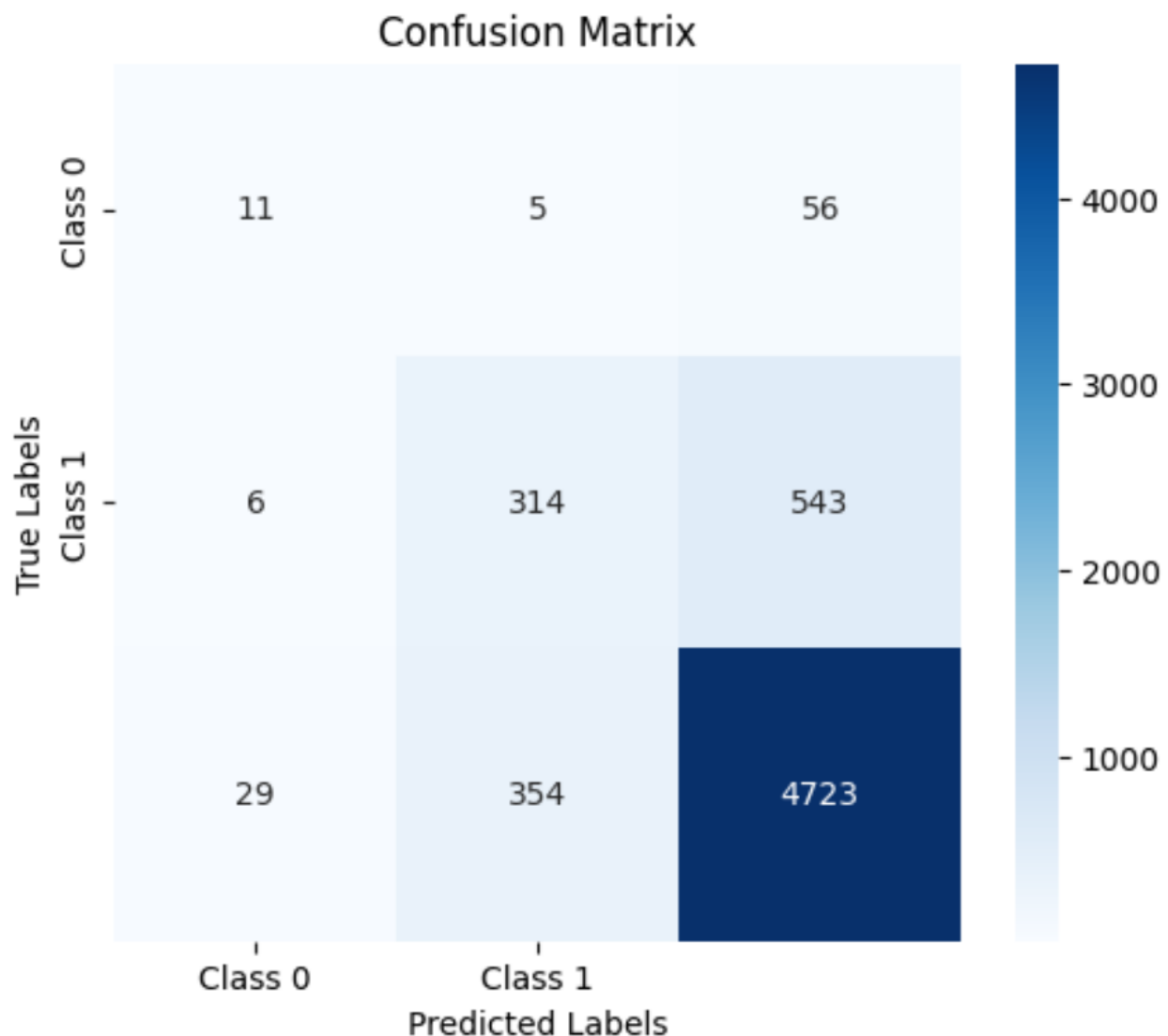
**4.2 Random Forest Classifier for Multiclass Classification**

Next, a multiclass Random Forest Classifier was trained to predict the extent of the flight delay. The delay categories were classified as:

- **No delay**

- **Short delay**

- **Moderate delay**

- **Long delay**

The model's performance was evaluated using multi-class metrics such as accuracy and classification report, which provides a detailed breakdown of precision, recall, and F1 score for each class.



Confusion Matrix

### 4.3 Random Forest Regressor for Exact Delay Prediction

A Random Forest Regressor was used to predict the exact delay time. This model helps calculate the continuous delay values rather than just the categorical class of delay. The performance was evaluated using mean squared error (MSE) and mean absolute error (MAE).

```
Mean Absolute Error (MAE): 5.855884787286874
Root Mean Squared Error (RMSE): 8.404336446981045
```

### 5. Hyperparameter Tuning

To optimize model performance, hyperparameter tuning was performed using RandomizedSearchCV. Various hyperparameters were tested, including the number of estimators, maximum depth, and minimum samples for splits. After identifying the best parameters, the models were retrained with these optimal settings.

```python
models = {
    'Logistic Regression': LogisticRegression(max_iter=1000),
    'Random Forest': RandomForestClassifier(n_estimators=100),
    'SVM': SVC(),
    'KNN': KNeighborsClassifier(n_neighbors=5),
    'Gradient Boosting': GradientBoostingClassifier(),
    'Decision Tree': DecisionTreeClassifier(),
    'Naive Bayes': GaussianNB()
}
```

### 6. Model Evaluation and Results

### 6.1 Model Performance on Test Data

After training the models with the optimal hyperparameters, the models were evaluated on the test dataset. The binary classifier, multiclass classifier, and regressor were all assessed on their performance.

The performance of the models was compared based on metrics like accuracy, precision, recall, F1 score, and mean squared error. Additionally, the Kaggle score was checked for the final model to evaluate its competitiveness in real-world applications.