# CSE 422: ARTIFICIAL INTELLIGENCE

Fall 2024

# Project Report

**AIR QUALITY PREDICTOR**

Student Information:

Name : B M Rauf

Student ID : 22201782

Department of Computer Science and Engineering

Brac University

# TABLE OF CONTENTS

# Introduction

Air quality prediction is a critical task in environmental science and public health management. This project aims to classify air quality into different categories based on various environmental factors using machine learning models. The dataset contains attributes such as temperature, humidity, pollutant levels (PM2.5, PM10, NO2, SO2, CO), proximity to industrial areas, and population density, which serve as predictors for the classification task.

The goal is to preprocess the data, scale the features, split the dataset, and train multiple models to identify the most effective one for the classification task. With increasing urbanization and industrial activities, air quality monitoring is essential to ensure public safety and mitigate health risks associated with poor air quality.
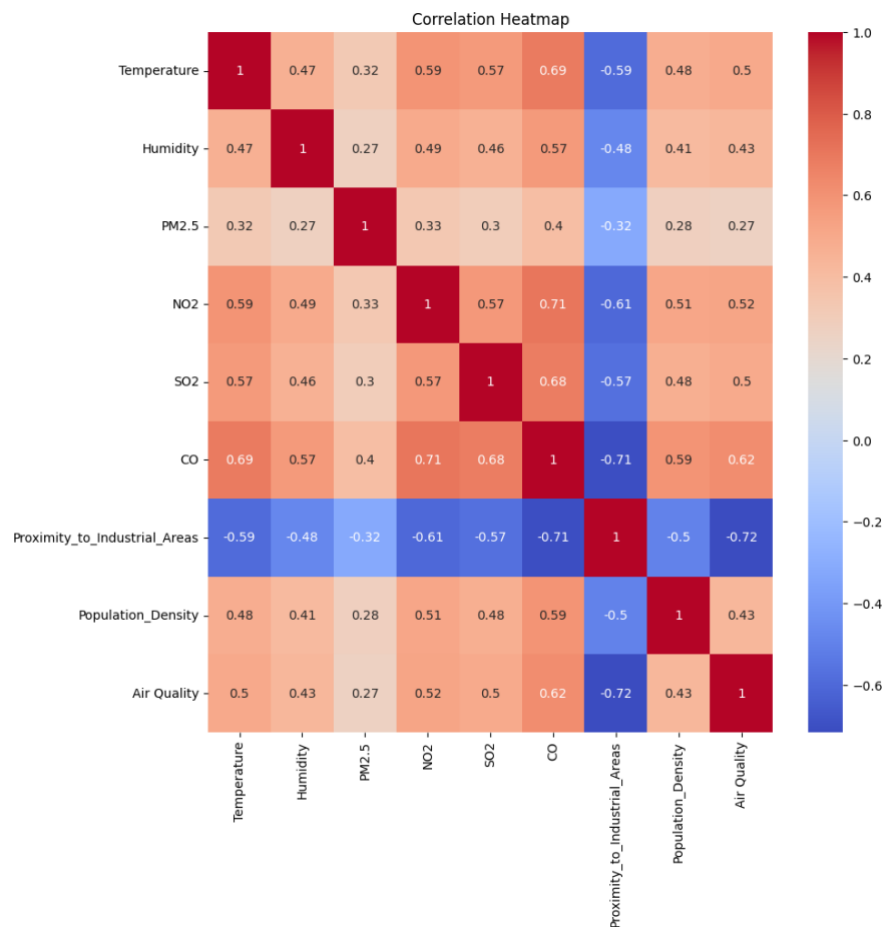
# Dataset Description

- **Source:**

  The dataset was sourced from Kaggle, containing information related to air pollution and its contributing factors.

  - Link: Air Quality Dataset
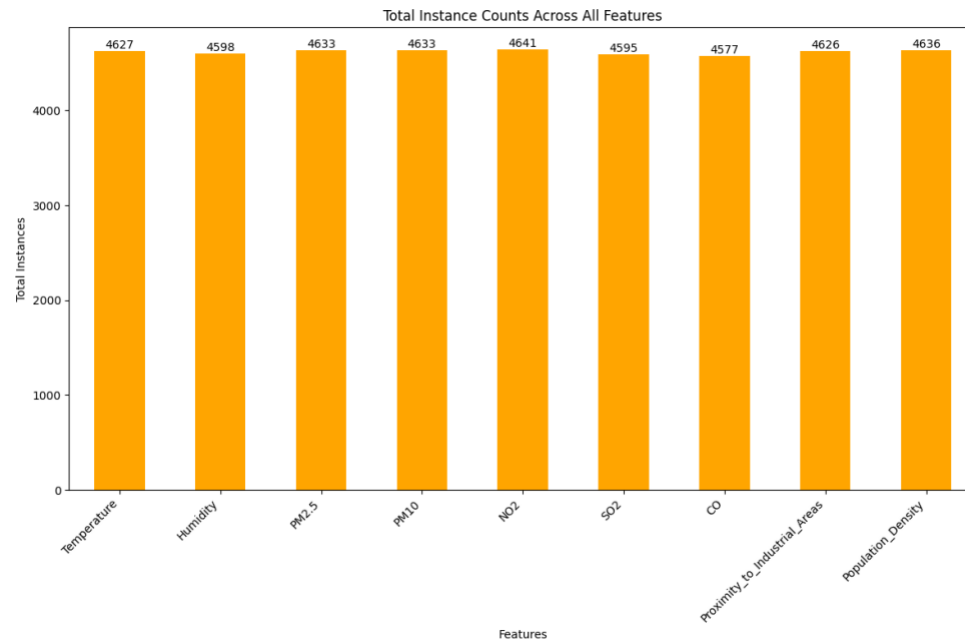
- **Dataset Description:**

  - **Number of Features:** 9 features, including temperature, humidity, pollutant levels, and population density.

  - **Problem Type:** Classification problem, as the target is categorical (Good, Moderate, Poor, Hazardous).

  - **Number of Data Points:** 5,000 samples.

o **Feature Types:** All features are numerical except for the target variable, which was encoded into numerical values using LabelEncoder.

o **Correlation Analysis:** A heatmap was generated to visualize correlations, highlighting highly correlated features that might impact the model.



Correlation Heatmap

o **Unbalanced Dataset:**

The dataset is not perfectly balanced because the number of instances varies across features. Here, the instance counts range from **4577 to 4641**, meaning some features have slightly

MORE OR FEWER VALUES THAN OTHERS.

Total Instance Counts Across All Features



# DATASET PREPROCESSING

- **FAULTS:**

  o **NULL VALUES:** THE DATASET DID NOT CONTAIN ANY NULL VALUES INITIALLY. HOWEVER, WE RANDOMLY INSERTED 5–10% NULL VALUES FOR TESTING PURPOSES.

  o **CATEGORICAL VALUES:** THE TARGET VARIABLE ("AIR QUALITY") WAS CATEGORICAL.

- **SOLUTIONS:**

  o **DROPPED COLUMNS:** THE '**PM10**' COLUMN WAS REMOVED DUE TO ITS HIGH CORRELATION WITH '**PM2.5**' TO REDUCE REDUNDANCY. MISSING VALUES IN CATEGORICAL DATA, SPECIFICALLY THE '**AIR QUALITY**' COLUMN, WERE FILLED USING THE **MODE** (MOST FREQUENT VALUE), WHILE MISSING VALUES IN NUMERICAL COLUMNS WERE REPLACED WITH THE **MEDIAN** (MIDDLE VALUE) TO PRESERVE DATA DISTRIBUTION.

FINALLY, ANY REMAINING ROWS WITH NULL VALUES AFTER THIS PROCESS WERE DROPPED TO ENSURE DATA COMPLETENESS.

- o **ENCODING:** THE TARGET VARIABLE ("AIR QUALITY") WAS ENCODED USING LABELENCODER TO CONVERT CATEGORICAL VALUES INTO NUMERIC REPRESENTATIONS.

## FEATURE SCALING

FEATURE SCALING IS ESSENTIAL FOR ENSURING THAT MACHINE LEARNING ALGORITHMS TREAT ALL FEATURES EQUALLY. THE DATASET WAS SCALED USING THE MINMAXSCALER FROM SCIKIT-LEARN, WHICH NORMALIZES THE VALUES BETWEEN 0 AND 1. THIS STEP ENSURES THE DATASET IS SUITABLE FOR ALGORITHMS SENSITIVE TO MAGNITUDE DIFFERENCES IN FEATURES.

## DATASET SPLITTING

THE DATASET WAS SPLIT INTO TRAINING AND TESTING SUBSETS TO EVALUATE MODEL PERFORMANCE:

- SPLITTING TECHNIQUE: STRATIFIED RANDOM SAMPLING WAS USED TO ENSURE THAT THE DISTRIBUTION OF THE TARGET VARIABLE IN THE TRAINING AND TESTING SUBSETS MATCHES THE ORIGINAL DATASET.
- TRAINING SET: 70% OF THE DATA.
- TESTING SET: 30% OF THE DATA.

THE SPLITTING ENSURED AN UNBIASED EVALUATION OF THE MODEL'S PERFORMANCE ON UNSEEN DATA WHILE MAINTAINING THE REPRESENTATIVENESS OF THE TARGET VARIABLE DISTRIBUTION. THE RANDOM_STATE PARAMETER WAS USED TO ENSURE REPRODUCIBILITY OF THE SPLITS.

# MODEL TRAINING & TESTING

MULTIPLE MACHINE LEARNING MODELS WERE TRAINED TO CLASSIFY AIR QUALITY, INCLUDING:

1. **LOGISTIC REGRESSION**

   - **DESCRIPTION:**

     LOGISTIC REGRESSION IS A STATISTICAL MODEL USED FOR BINARY AND MULTICLASS CLASSIFICATION TASKS. IT ESTIMATES PROBABILITIES USING A LOGISTIC FUNCTION, WHICH MAPS THE INPUT FEATURES TO A PROBABILITY VALUE BETWEEN 0 AND 1. THE OUTPUT IS THEN CATEGORIZED INTO CLASSES BASED ON A THRESHOLD (COMMONLY 0.5).

   - **HOW IT WORKS:**
     - **LOGISTIC FUNCTION:** THE ALGORITHM USES THE LOGISTIC

       $$P(y = 1|X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + ... + \beta_n X_n)}}$$

       (SIGMOID) FUNCTION:

       WHERE $X_1, X_2, ..., X_n$ ARE THE INPUT FEATURES, AND $\beta_0, \beta_1, ..., \beta_n$ ARE THE COEFFICIENTS THAT THE MODEL LEARNS DURING TRAINING.

     - **TRAINING:** DURING TRAINING, THE MODEL OPTIMIZES THE COEFFICIENTS ($\beta$) TO MINIMIZE THE LOSS FUNCTION, TYPICALLY THE LOG-LOSS (CROSS-ENTROPY LOSS), WHICH MEASURES HOW WELL THE PREDICTED PROBABILITIES MATCH THE ACTUAL CLASS LABELS.

     - **PREDICTION:** FOR A NEW DATA POINT, THE MODEL CALCULATES THE PROBABILITY OF BELONGING TO EACH CLASS. IF THE PROBABILITY

EXCEEDS THE THRESHOLD (E.G., 0.5), THE DATA POINT IS CLASSIfiED AS ONE CLASS; OTHERWISE, IT IS CLASSIfiED AS ANOTHER.

- o **MULTICLASS CLASSIfiCATION:** FOR MULTICLASS PROBLEMS, LOGISTIC REGRESSION CAN USE TECHNIQUES LIKE ONE-VS-REST (OVR) OR SOFTMAX REGRESSION TO HANDLE MULTIPLE CLASSES.

## 2. NAIVE BAYES

- **DESCRIPTION:**

NAIVE BAYES IS A PROBABILISTIC CLASSIfiER BASED ON BAYES' THEOREM , ASSUMING INDEPENDENCE BETWEEN FEATURES. DESPITE ITS SIMPLICITY AND "NAIVE" ASSUMPTION OF FEATURE INDEPENDENCE, IT PERFORMS SURPRISINGLY WELL IN MANY REAL-WORLD SCENARIOS, ESPECIALLY IN TEXT CLASSIfiCATION AND SPAM fiLTERING.

- **HOW IT WORKS:**

$$P(y|X) = \frac{P(X|y) \cdot P(y)}{P(X)}$$

- o **BAYES' THEOREM:**

WHERE:

- - P(Y|X): POSTERIOR PROBABILITY OF CLASS Y GIVEN THE FEATURES X.
  - P(X|Y): LIKELIHOOD OF OBSERVING THE FEATURES X GIVEN THE CLASS Y.
  - P(Y): PRIOR PROBABILITY OF CLASS Y.

- P(X): Probability of observing the features X (normalizing constant).

  o **Training:**

  - The model calculates the prior probabilities P(Y) for each class and the likelihoods P(X|Y) for each feature given the class.

  - Since Naive Bayes assumes feature independence, the joint probability of all features given the class is calculated as the product of individual feature probabilities:

    - $P(X|Y) = P(X_1|Y) \cdot P(X_2|Y) \cdot ... \cdot P(X_N|Y)$

  o **Prediction:**

  - For a new data point, the model calculates the posterior probability P(Y|X) for each class using Bayes' theorem.

  - The class with the highest posterior probability is chosen as the predicted class.

3. **Random Forest Classifier**
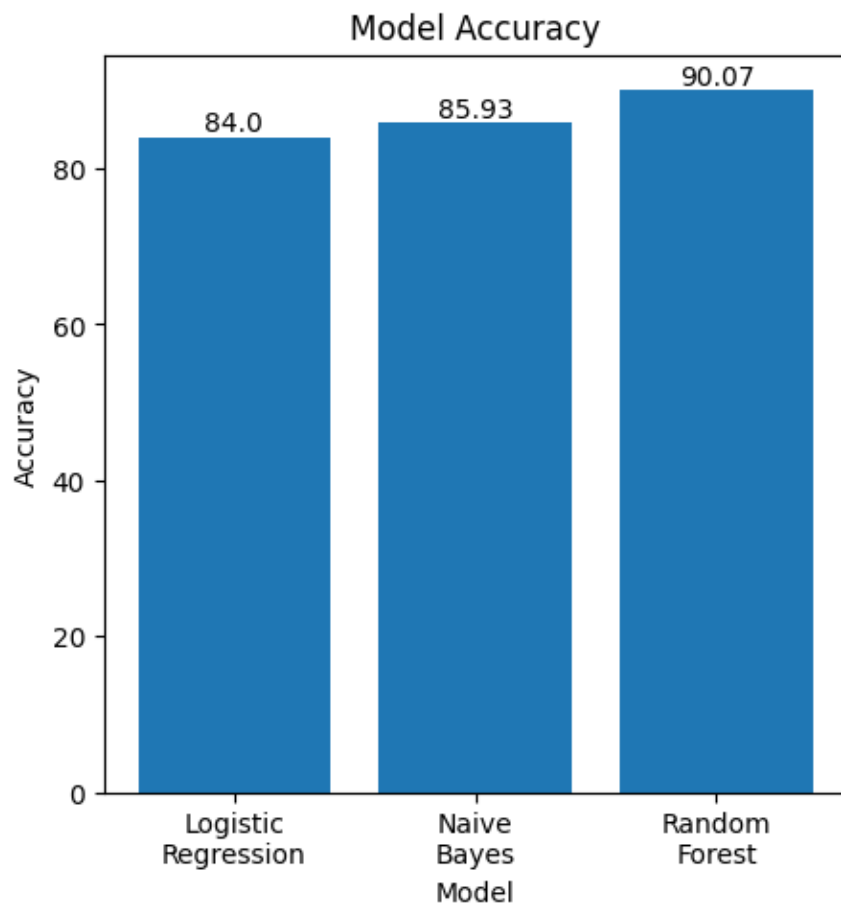
- **Description:**

  Random Forest is an ensemble learning method that builds multiple decision trees during training and aggregates their predictions (e.g., majority voting for classification) to improve accuracy and reduce overfitting. It is robust to noise and works well on large datasets.

- **How It Works:**

- o **Bootstrap Sampling:** Random Forest creates multiple decision trees by sampling the dataset with replacement (bootstrap sampling). Each tree is trained on a different subset of the data.

- o **Random Feature Selection:** At each split in a tree, only a random subset of features is considered. This introduces diversity among the trees and reduces the risk of overfitting.

- o **Tree Construction:** Each tree is grown fully without pruning, ensuring maximum depth and complexity.

- o **Aggregation:** For classification tasks, the final prediction is made by majority voting among all the trees. Each tree "votes" for a class, and the class with the most votes is selected as the final prediction.

- o **Advantages:**
  - ▪ Reduces overfitting compared to individual decision trees.
  - ▪ Handles high-dimensional data well.
  - ▪ Provides feature importance scores, helping identify the most influential features.

# Model Selection/Comparison Analysis

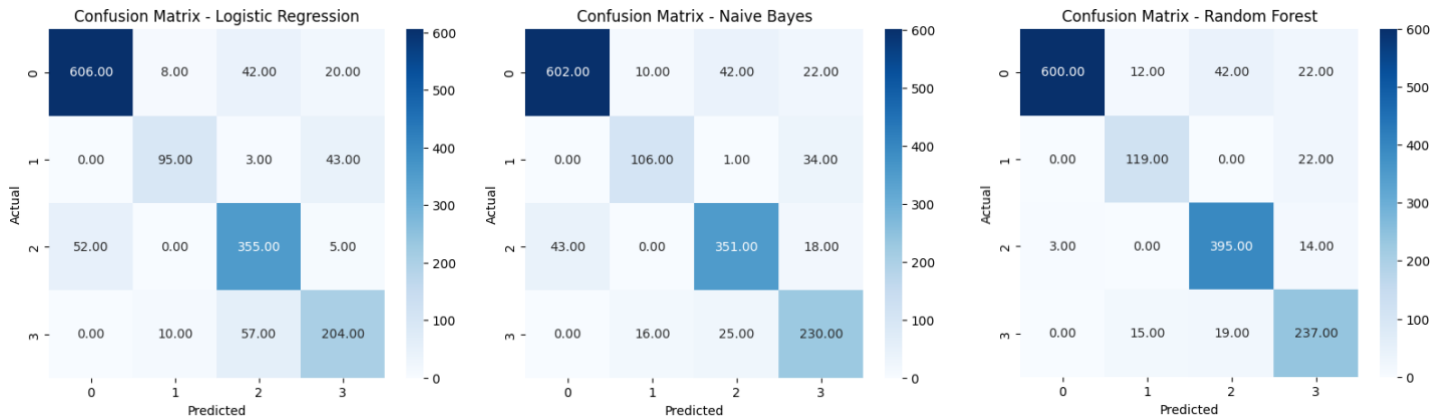- • **Bar Chart:** Showcased the prediction accuracy of all models.

Model Accuracy

- PRECISION AND RECALL:
  - PRECISION: THE PROPORTION OF CORRECTLY PREDICTED POSITIVE INSTANCES OUT OF ALL PREDICTED POSITIVES. IT MEASURES ACCURACY IN POSITIVE PREDICTIONS.
  - RECALL: THE PROPORTION OF CORRECTLY PREDICTED POSITIVE INSTANCES OUT OF ALL ACTUAL POSITIVES. IT MEASURES THE ABILITY TO fiND ALL POSITIVE INSTANCES.

```
                 Model  Precision     Recall
0  Logistic Regression   0.842941   0.840000
1          Naive Bayes   0.862886   0.859333
2        Random Forest   0.908111   0.900667
```

o CONFUSION MATRIX: ANALYZED FOR EACH MODEL TO EVALUATE CLASSIFICATION PERFORMANCE IN DETAIL.



# CONCLUSION

THE AIR QUALITY PREDICTOR PROJECT SUCCESSFULLY USED MACHINE LEARNING TO CLASSIFY AIR QUALITY INTO CATEGORIES SUCH AS GOOD, MODERATE, POOR, AND HAZARDOUS. THROUGH PROPER PREPROCESSING, SCALING, AND SPLITTING OF THE DATASET, THE MODELS — LOGISTIC REGRESSION, NAIVE BAYES, AND RANDOM FOREST — ACHIEVED ACCURACIES OF 84.0%, 85.93%, AND 90.07%, RESPECTIVELY. METRICS LIKE PRECISION, RECALL, AND CONFUSION MATRIX CONFIRMED THEIR EFFECTIVENESS. THIS PROJECT DEMONSTRATES HOW MACHINE LEARNING CAN BE A RELIABLE TOOL TO MONITOR AIR QUALITY AND PREVENT HEALTH RISKS ASSOCIATED WITH POLLUTION.