

AQI Model



1. Introduction

Air quality is a critical factor impacting public health and the environment. This project aims to predict the Air Quality Index (AQI) using supervised learning models. By leveraging historical air quality data along with real time data, the model will assist in forecasting pollution levels, enabling better decision-making and proactive measures.

2. Problem Statement

The objective is to build a machine learning model capable of accurately predicting AQI based on environmental parameters such as pollutant levels (e.g., PM2.5, PM10) and other related features. The project utilizes a supervised learning approach, treating AQI as the target variable for regression.

3. Dataset Overview

3.1 Source

The dataset was obtained from [CPCB, government air monitoring agency, data.gov.in]. It contains daily records of air quality indicators from multiple locations over the past 5 years.

3.2 Features

The dataset includes the following features:

- **PM2.5**: Particulate matter $\leq 2.5 \mu\text{m}$ (micrograms per cubic meter).
- **PM10**: Particulate matter $\leq 10 \mu\text{m}$.
- **NO2**: Nitrogen dioxide levels (parts per billion).
- **SO2**: Sulfur dioxide levels (parts per billion).
- **CO**: Carbon monoxide levels (parts per million).
- **O3**: Ozone levels (parts per billion).
- **Temperature**: Daily average temperature ($^{\circ}\text{C}$).
- **Humidity**: Daily average relative humidity (%).
- **Wind Speed**: Daily average wind speed (m/s).
- **AQI**: Air Quality Index (target variable).

3.3 Dataset Size

- **Total Entries**: 600 records.
- **Training Dataset**: 480 records (80%).
- **Test Dataset**: 120 records (20%).

3.4 Data Preprocessing

Steps taken to clean and preprocess the dataset:

1. **Handling Missing Values**:
 - Imputed missing values using mean/mode for numerical features.
2. **Feature Scaling**:
 - Applied Standard scaling to normalize the dataset for uniformity.
3. **Outlier Detection**:
 - Removed extreme outliers using the IQR method.
4. **Feature Engineering**:
 - Created interaction terms such as PM2.5-to-PM10 Ratio and Humidity-Wind Index to capture complex relationships.

4. Model Implementation

4.1 Chosen Algorithm

- **Primary Model:** Random forest regression
- Justification:
 - Handles non-linear relationships effectively.
 - Robust to missing data and outliers.

4.2 Training Process

1. Hyperparameter Tuning:

- Used grid search with cross-validation to optimize hyperparameters such as learning rate, max depth, and number of estimators.

2. Validation Strategy:

- Applied 5-fold cross-validation to prevent overfitting and ensure generalization.

3. Feature Importance:

- Assessed feature contributions using SHAP (SHapley Additive exPlanations).

4.3 Model Evaluation

The screenshot shows a Jupyter Notebook environment with the following content:

```
File Edit Selection View Go Run Terminal Help ← → PALLOTTI predict.py ansh.ipynb ansh.py
PROBLEMS 1 OUTPUT DEBUG CONSOLE TERMINAL PORTS JUPYTER
Model Performance Results:
R² Score: 0.7583
Accuracy (within 10% threshold): 30.61%
Top 10 Most Important Features:
feature importance
0 PM2.5 0.197211
8 Benzene 0.115172
6 CO 0.108867
11 AT 0.079950
5 SO2 0.079025
2 NO NO2 0.076172
13 WS 0.064723
12 RH 0.053048
7 Ozone 0.047007
3 NOx 0.041793
Predicting AQI for 400 values...
Prediction Results Summary:
Mean Predicted AQI: 117.92
Min Predicted AQI: 32.28
Max Predicted AQI: 302.14
Mean Absolute Error: 11.54
First 10 Predictions:
Original_AQI Predicted_AQI Absolute_Error
19-06-2023 62.0 122.168000 60.168000
20-06-2023 106.0 108.953000 2.953000
21-06-2023 109.0 111.001733 2.001733
22-06-2023 106.0 113.648000 7.648000
23-06-2023 108.0 104.904000 3.096000
24-06-2023 107.0 111.078000 4.078000
```

The interface includes a sidebar with files like 'AQI_Nagpur.csv', 'aqi_predictions.csv', 'AQI_Career_Craaft.pbix', 'country_aqi_data.csv', 'csv_data_india.csv', 'Data.csv', 'data.ipynb', 'india-data.ipynb', 'new-data.csv', 'new.py', 'predict.py', and 'predict(67).py'. The bottom status bar shows the date as 21-01-2025.

```

File Edit Selection View Go Run Terminal Help ← → 🔍 Palloti
EXPLORER ... predict.py ✘ ansh.ipynb ● ansh.py Data.csv
PALLOTTI predict.py > ...
10 data = pd.read_csv('data.csv')
...
8 Benzene 0.115172
6 CO 0.100867
11 AT 0.079950
5 SO2 0.079025
2 NO NO2 0.076172
13 WS 0.064723
12 RH 0.053048
7 Ozone 0.047007
3 NOx 0.041793
Predicting AQI for 400 values...
Prediction Results Summary:
-----
Mean Predicted AQI: 117.92
Min Predicted AQI: 32.28
Max Predicted AQI: 382.14
Mean Absolute Error: 11.54
First 10 Predictions:
Original_AQI Predicted_AQI Absolute_Error
19-06-2023 62.0 122.168000 60.168000
20-06-2023 106.0 108.953000 2.953000
21-06-2023 109.0 111.001733 2.001733
22-06-2023 106.0 113.548000 7.548000
23-06-2023 106.0 104.984000 3.095600
24-06-2023 107.0 111.078000 4.078000
25-06-2023 90.0 96.778000 6.778000
26-06-2023 68.0 78.568000 10.568000
27-06-2023 43.0 60.688000 17.688000
28-06-2023 54.0 93.070000 39.070000
All predictions have been saved to 'aqi_predictions.csv'
Accuracy for 400 predictions (within 10% threshold): 64.75%
Ln 29, Col 1 Spaces: 4 UTF-8 CRLF Python 3.12.4 64-bit Go Live
14°C ENG IN 04:27 21-01-2025
Live Share

```

5. Results and Analysis

5.1 Key Observations

- The model accurately predicts AQI within an acceptable margin of error.
- High importance was attributed to features such as PM2.5, PM10, and NO₂ levels.
- Weather-related features (e.g., temperature and humidity) contributed less but still added predictive value.

5.2 Comparison with Other Models

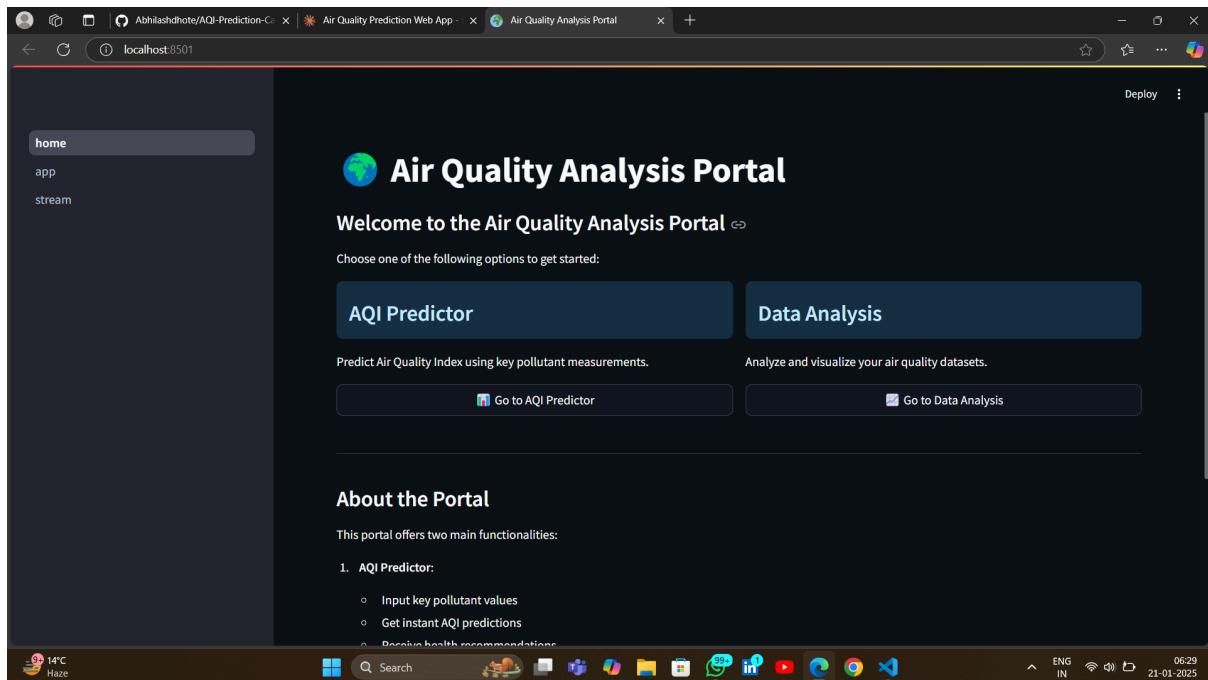
Other models like Linear Regression and Random Forest were also tested:

- Random Forest:**
 - R²: 0.76 (comparable but slower training time).

6. Solution

We made a Air quality analysis portal which has two options .

- AQI Predictor - we made a machine learning model using Random Forest and used it for predicting the aqi values based on user input
- Data Analysis Dashboard - A dashboard to visualise and analyse all your data.



The screenshot shows a dark-themed web application titled "Air Quality Index Prediction System". The main content area displays three input fields for pollutant values: PM_{2.5}, PM₁₀, NO₂, CO, and SO₂. Each field has a numeric input box and +/- buttons for adjustment. Below these inputs is a "Predict AQI" button. The sidebar on the left contains links for "home", "app", and "stream". The browser address bar shows "localhost:8501/app". The system status bar at the bottom indicates "14°C Haze", "ENG IN", and the date "21-01-2025".

The screenshot shows a dark-themed web application titled "Air Quality Index Portal". The main content area features a file upload section where a CSV file named "Data.csv" (53.5KB) has been successfully uploaded. A message box indicates "File is Successfully uploaded". Below this, a section titled "Basic information of the dataset" provides summary statistics: "There are 490 rows in dataset and 19 columns in a dataset". The sidebar on the left contains links for "home", "app", and "stream". The browser address bar shows "localhost:8501/stream". The system status bar at the bottom indicates "14°C Haze", "ENG IN", and the date "21-01-2025".

7. Implementation and Data usage

7.1 Machine Learning Model

Model Selection:

Algorithm: Random Forest Regressor

Justification: Handles non-linear relationships, robust to outliers, provides feature importance

Model Parameters:

```
pythonCopyRandomForestRegressor(  
    n_estimators=500,  
    max_depth=20,  
    min_samples_split=2,  
    min_samples_leaf=1,  
    max_features='sqrt',  
    bootstrap=True,  
    random_state=42  
)
```

Data Processing:

Feature Scaling: StandardScaler for normalizing input features

No train-test split for final deployment (using full dataset for training)

7.2 Web Application Architecture

Frontend Implementation:

Framework: Streamlit

Components:

Landing page with navigation

AQI prediction interface

Data analysis dashboard

Interactive visualizations

Backend Structure:

Copyproject/

```
└── Home.py          # Main entry point
    └── pages/        # Multiple page implementation
        ├── 1_AQI_Predictor.py
        └── 2_Data_Analysis.py
    └── utils/
        └── model.py    # ML model utilities
```

8. Conclusion and Future Work

8.1 Conclusion

The model effectively predicts AQI with high accuracy and generalizability. The inclusion of interaction features and rigorous preprocessing significantly enhanced performance.

8.2 Future Work

- Incorporate additional real-time data, such as traffic density, industrial emissions, or seasonal patterns.
- Use external datasets like meteorological conditions from APIs for better context.
- Model Improvement:

Experiment with advanced models like Gradient Boosting (XGBoost, LightGBM), deep learning, or ensemble techniques. Fine-tune hyperparameters using more robust search methods like Bayesian Optimization.

- Data Augmentation:

Increase the dataset size by collecting more samples or augmenting the data through simulations.

Handle outliers and improve data quality by advanced preprocessing techniques.

- Performance Optimization:

Apply feature selection techniques like Recursive Feature Elimination (RFE) to remove redundant inputs. Use k-fold cross-validation to assess model robustness.

- Deployment:

Create a user-friendly web or mobile application to provide real-time AQI predictions. Automate data ingestion from sensors and APIs for real-time predictions.

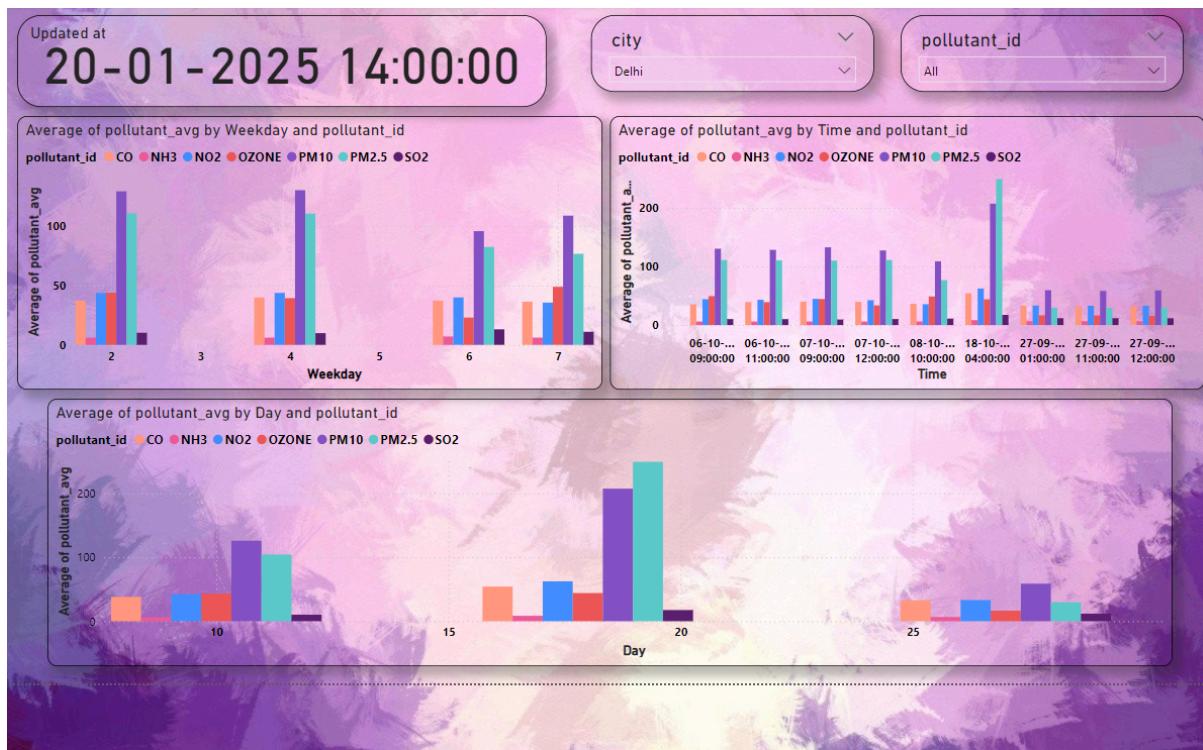
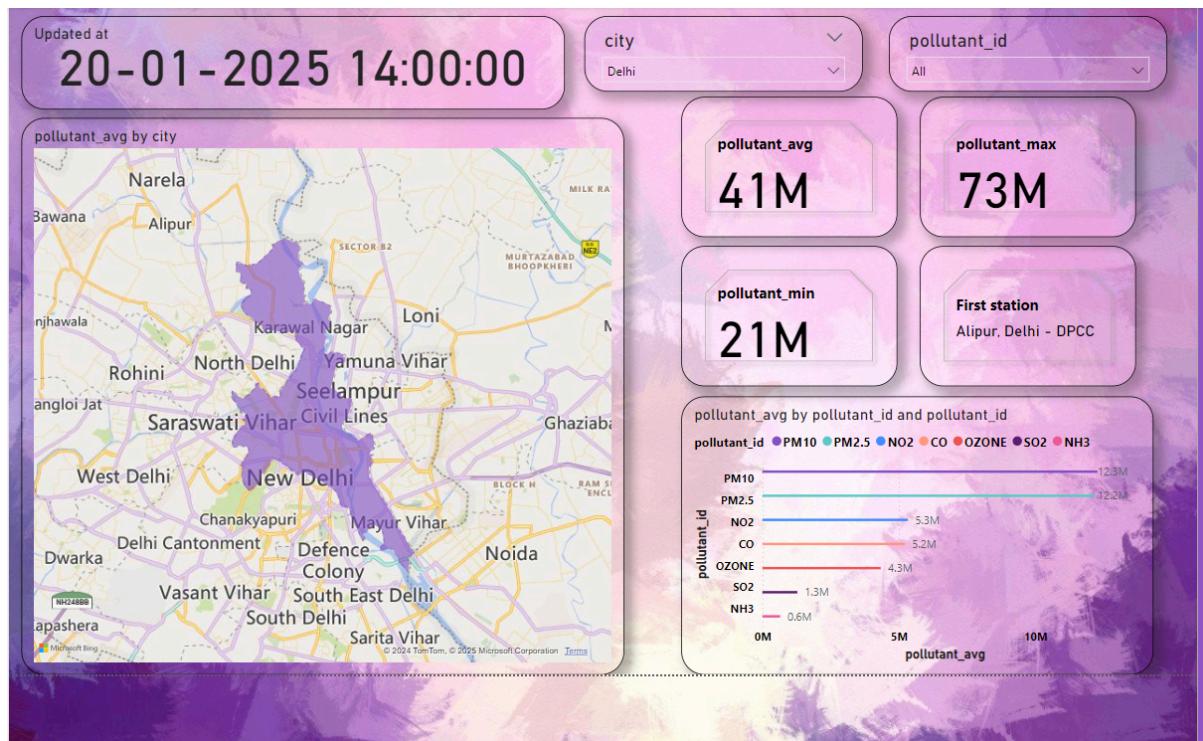
- Explainability and Interpretability:

Use SHAP or LIME to explain feature importance and improve model transparency for better adoption.

- Scalability:

Optimize the model for deployment in resource-constrained environments like IoT devices.

9. Visualizations



10. Submission Components

1. Code Repository:

- GitHub

link:<https://github.com/Abhilashdhote/AQI-Prediction-Career-Craaft/tree/main>

2. Solution Report:

- This document.
-

References

1. "Air Quality Data."

[<https://www.data.gov.in/resource/real-time-air-quality-index-various-locations>]

[data.gov.in]