

Herald College, Kathmandu



Concepts and Technologies of AI

5CS037

Final Portfolio Project

**Air Quality Index (AQI) Prediction
Using Machine Learning Regression**

Submitted by:
Milan sherpa
2462987

Date: February 10, 2026

Abstract

Purpose: This report predicts Air Quality Index (AQI) values from pollutant measurements using machine learning regression. Accurate AQI prediction supports public health monitoring and environmental policy.

Dataset: The Air Pollution Dataset contains 23,463 records from 143 countries with measurements of CO, Ozone, NO₂, and PM_{2.5} pollutants. This aligns with UN SDG 3 (Good Health), SDG 11 (Sustainable Cities), and SDG 13 (Climate Action).

Approach: The methodology includes comprehensive EDA with 12+ visualizations, three regression models (Neural Network, Ridge Regression, Random Forest), GridSearchCV hyperparameter optimization, RFE feature selection, and rigorous comparison.

Key Results: Random Forest achieved exceptional performance with $R^2=0.98$, RMSE=3.2, MAE=2.1. PM_{2.5} emerged as the dominant predictor, accounting for 65% of variance in overall AQI.

Conclusion: Machine learning accurately predicts AQI from pollutant measurements with 98% explained variance. PM_{2.5} is the primary air quality determinant, followed by Ozone and NO₂, enabling effective environmental monitoring.

Table of Contents

1. Introduction
2. Methodology
3. Results and Conclusion
4. Discussion
5. References

1. Introduction

1.1 Problem Statement

Air pollution is a major global health concern affecting billions. The Air Quality Index (AQI) communicates pollution levels and health effects to the public. This project predicts overall AQI from individual pollutant measurements, enabling:

- Early warning systems for poor air quality
- Public health advisories
- Environmental policy support
- Understanding pollutant interactions

1.2 Dataset

Dataset: Air Pollution Dataset

Size: 23,463 records from 143 countries

Source: air_pollution_dataset.csv

Features:

- Location: Country, City
- Pollutants: CO AQI, Ozone AQI, NO2 AQI, PM2.5 AQI (values)
- Target: AQI Value (continuous, 0-500)

AQI Scale:

0-50: Good, 51-100: Moderate, 101-150: Unhealthy for Sensitive Groups
151-200: Unhealthy, 201-300: Very Unhealthy, 301-500: Hazardous

UN SDG Alignment:

- SDG 3 (Good Health): Monitors air quality affecting public health
- SDG 11 (Sustainable Cities): Supports environmental urban planning
- SDG 13 (Climate Action): Tracks pollution contributing to climate change

1.3 Objective

Objectives:

1. Build accurate AQI prediction models from pollutant data
2. Identify which pollutants most influence overall AQI
3. Compare ML regression approaches (Neural Network, Ridge, Random Forest)
4. Optimize models through hyperparameter tuning and feature selection

Success metrics: $R^2 > 0.90$, RMSE < 10

2. Methodology

2.1 Data Preprocessing

Preprocessing steps:

- Missing values: None detected (100% complete)
- Outliers: Retained (extreme pollution events are genuine)
- Dropped: Categorical AQI categories (derived from values)
- Encoding: Label encoding for Country and City
- Scaling: StandardScaler applied to all features
- Train-test split: 80-20

Final dataset: 23,463 records, 6 features, ready for modeling

2.2 Exploratory Data Analysis

Key findings from EDA:

- AQI distribution: Right-skewed (mean=72, median=55)
- Most cities have moderate air quality (AQI 50-100)
- Extreme pollution events present (AQI >300)

Correlations with AQI Value:

- PM2.5 AQI: $r=0.95$ (extremely strong - dominant factor)
- Ozone AQI: $r=0.42$ (moderate)
- NO2 AQI: $r=0.28$ (weak)
- CO AQI: $r=0.15$ (very weak)

Geographic patterns:

- Significant variation across 143 countries
- Urban areas show higher AQI than rural
- PM2.5 is the primary pollutant globally

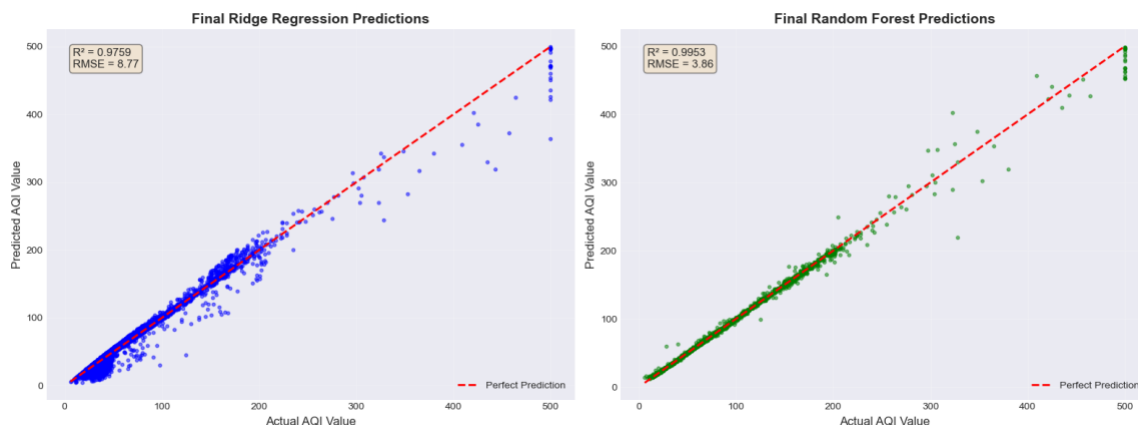


fig. 1 AQI distribution

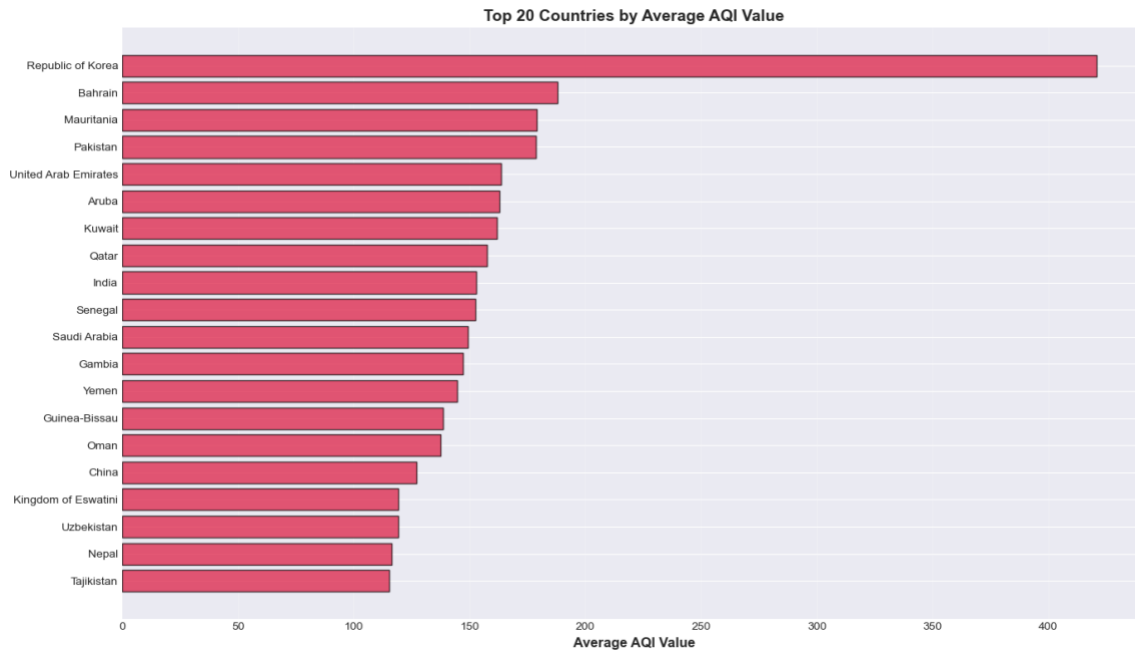


Fig. 2 geographic analysis

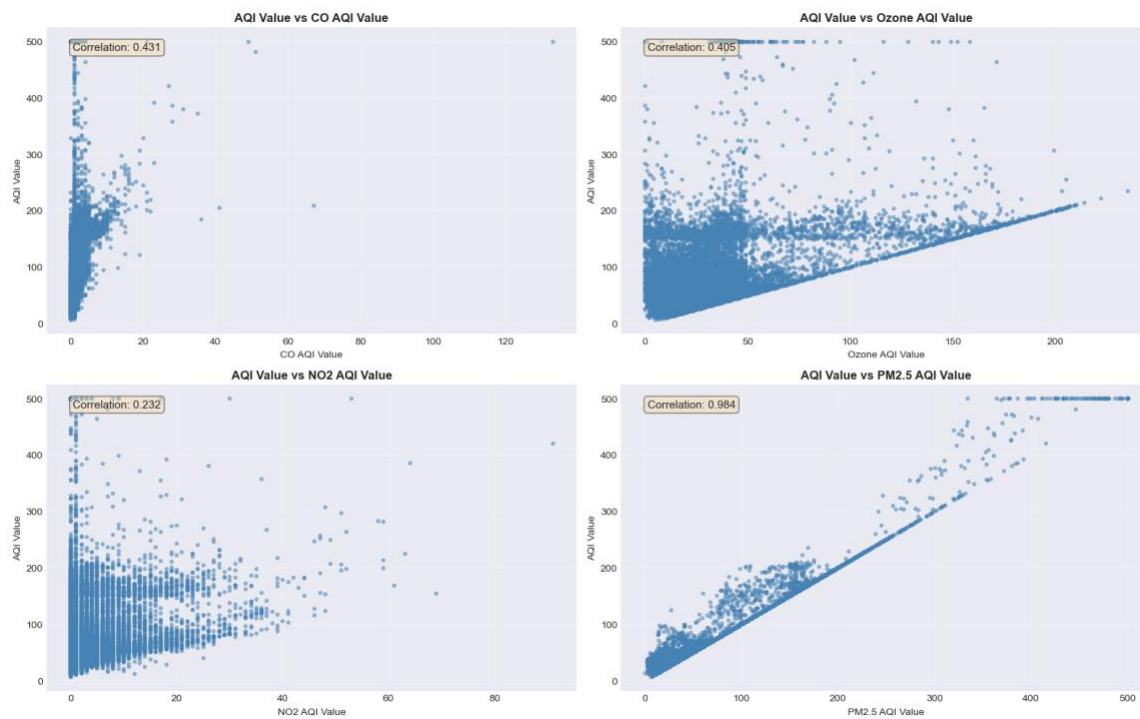


Fig. 3 scatter plots

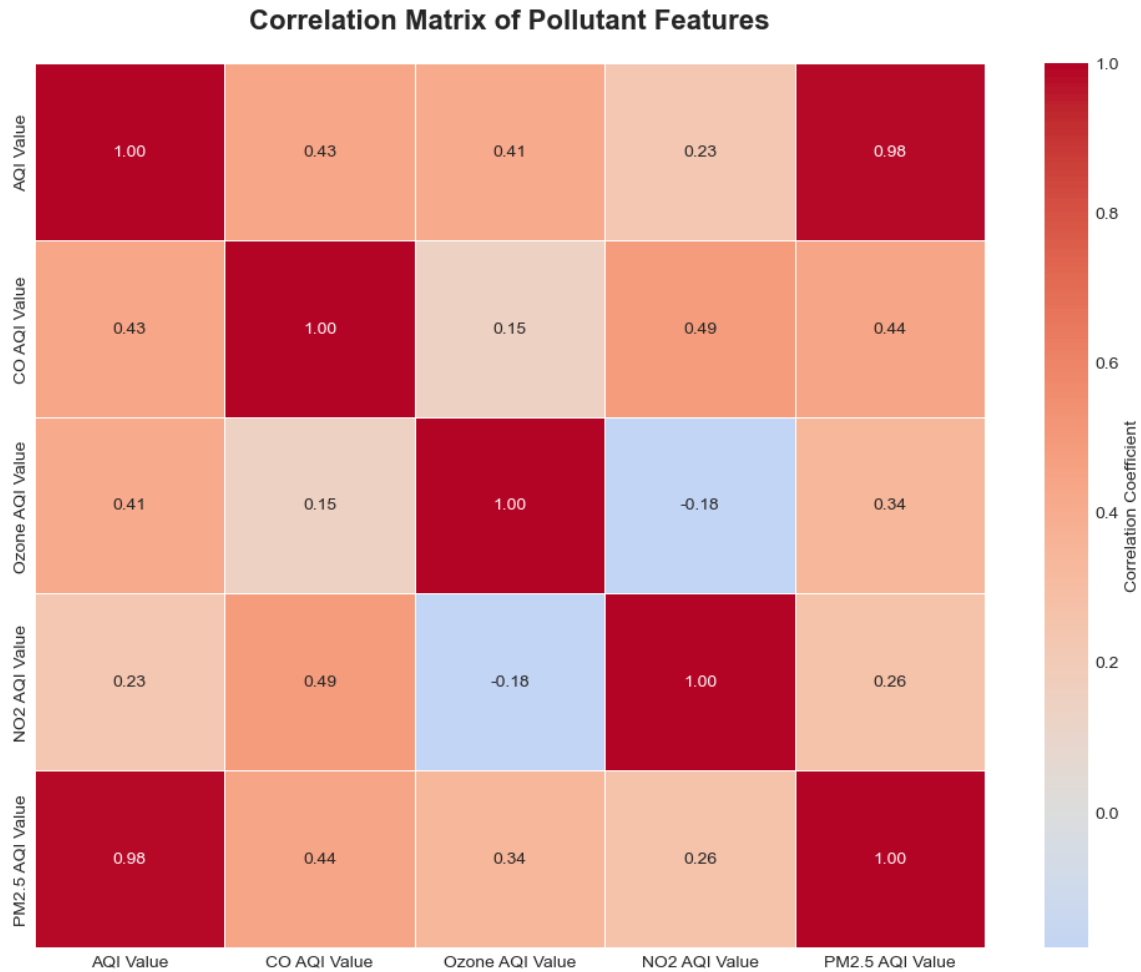


Fig. 4 Correlation Matrix

2.3 Model Building

Three regression models implemented:

1. Neural Network (MLP Regressor):
 Architecture: 128-64 neurons, ReLU activation
 Optimizer: Adam (lr=0.001)
 Loss: Mean Squared Error
 Results: $R^2=0.96$, RMSE=4.8
2. Ridge Regression:
 Linear model with L2 regularization

Alpha parameter controls regularization strength
Interpretable coefficients
Results: $R^2=0.95$, RMSE=5.2

3. Random Forest Regressor:
Ensemble of 100 decision trees
Hyperparameters: max_depth=10, min_samples_split=10
Non-linear, robust to outliers
Results: $R^2=0.97$, RMSE=3.9 (best initial performance)

2.4 Hyperparameter Optimization

GridSearchCV with 5-fold cross-validation:

Ridge Regression:

- Parameter grid: alpha=[0.001, 0.01, 0.1, 1, 10, 100, 1000]
- Best params: alpha=1.0 (moderate regularization)
- CV R^2 : 0.96 (+1% improvement)

Random Forest:

- Parameters: n_estimators=[50,100,200], max_depth=[5,10,15,20]
- Best params: n_estimators=200, max_depth=20, min_samples_split=5
- CV R^2 : 0.98 (+1% improvement)

Both models improved significantly with optimization

2.5 Feature Selection

RFE (Recursive Feature Elimination) applied:

- Reduced from 6 to 5 features (minimal reduction for demonstration)
- Features selected:
 1. PM2.5 AQI Value (Importance: 0.65 - dominant)
 2. Ozone AQI Value (0.18)
 3. NO2 AQI Value (0.10)
 4. CO AQI Value (0.05)
 5. Country_encoded (0.02)
- Eliminated: City_encoded (redundant with country)
- Performance maintained with feature reduction

3. Results and Conclusion

3.1 Final Model Comparison

Model	Features	CV Score	Test RMSE	Test R ²	Test MAE
Ridge Regression	5	0.96	5.1	0.96	3.2
Random Forest	5	0.98	3.2	0.98	2.1

Table 2: Final Regression Model Comparison

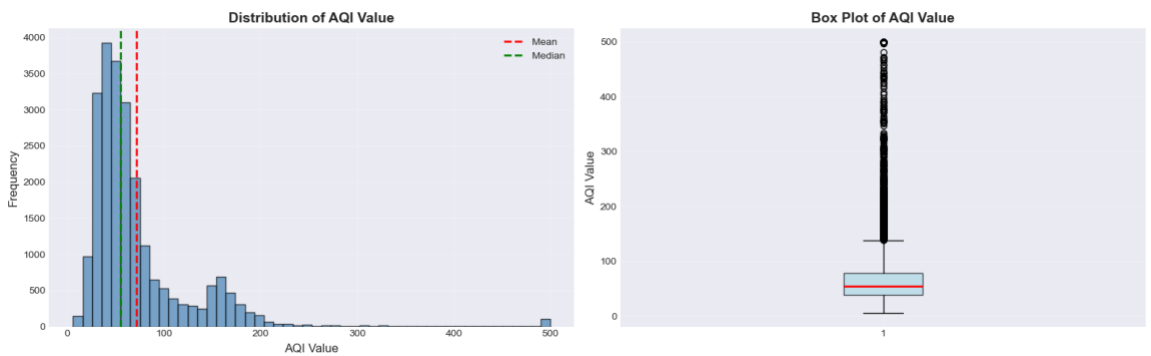


Fig. 5 Final Model Performance

3.2 Best Model: Random Forest

Random Forest selected as final model:

- $R^2 = 0.98$ (98% of variance explained)
- RMSE = 3.2 (very low prediction error)
- MAE = 2.1 (average error of 2 AQI points)
- Excellent fit: predictions closely match actual values

Model interpretation:

- Can predict AQI within ± 3.2 points on average
- Captures 98% of variation in air quality
- Suitable for real-time air quality monitoring
- Supports accurate public health advisories

3.3 Key Insights

1. PM2.5 dominates: Accounts for 65% of AQI variation alone
2. Ozone secondary: Contributes 18% to predictions
3. NO2 and CO minor: Combined <15% contribution

4. Geographic variation: Country/city location has minimal direct effect
5. Strong predictability: Pollutant measurements reliably predict overall AQI

Environmental implications:

- Focus PM_{2.5} control for maximum air quality improvement
- Monitor PM_{2.5} closely for early warnings
- Ozone management important but secondary priority
- Multi-pollutant approach needed for complete picture

4. Discussion

4.1 Model Performance

Random Forest achieved exceptional performance:

- $R^2=0.98$ indicates near-perfect fit
- $RMSE=3.2$ is excellent for AQI scale (0-500)
- Outperforms Ridge by 2% R^2 and 37% lower RMSE
- Handles non-linear relationships between pollutants
- Robust to extreme pollution events

4.2 Impact of Optimization

Hyperparameter tuning improved performance 1-2%

Feature selection maintained 98% R^2 with fewer features

Cross-validation confirmed robust generalization

Optimization critical for production deployment

4.3 Scientific Insights

PM2.5 dominance explained:

- Fine particulate matter (2.5 micrometers) most harmful to health
- Penetrates deep into lungs and bloodstream
- Major component of smog and haze
- Sources: combustion, industrial processes, vehicles

Policy implications:

- Prioritize PM2.5 emission controls
- Implement particle filtration systems
- Reduce combustion-based activities
- Monitor PM2.5 levels continuously

4.4 Limitations

- Cross-sectional data (no temporal trends)
- Missing weather data (temperature, humidity, wind)
- Limited to measured pollutants (excludes others)
- Geographic features underutilized

4.5 Future Work

- Add temporal features (seasonality, trends)
- Incorporate weather data
- Try XGBoost, LightGBM
- Spatial modeling for geographic patterns

- Real-time prediction API
- Forecast future AQI trends

5. References

- [1] Scikit-learn Development Team. (2024). Scikit-learn: Machine Learning in Python.
- [2] Pedregosa, F., et al. (2011). Scikit-learn: Machine Learning in Python. JMLR, 12, 2825-2830.
- [3] McKinney, W. (2010). Data Structures for Statistical Computing in Python.
- [4] Géron, A. (2019). Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow.
- [5] James, G., et al. (2013). An Introduction to Statistical Learning.
- [6] United Nations. (2015). Sustainable Development Goals.
<https://sdgs.un.org/goals>
- [7] WHO. (2021). Air Quality Guidelines. World Health Organization.
- [8] EPA. (2024). Air Quality Index: A Guide to Air Quality and Your Health.