

EnviroScan

AI-Powered Pollution Source Identifier using Geospatial Analytics

Final Project Report

Author: **Praveen S**

Department of CSE (Artificial Intelligence and Machine Learning)

K.S.Rangasamy College of Technology

Infosys Springboard Internship Program

Date: **December 2025**

Model Accuracy	Data Coverage	Pollution Sources
92.26%	165+ Districts	5 Classes

Table of Contents

1. Executive Summary 3

2. Introduction 3

3. Data Sources and APIs 5

4. Data Preprocessing and Cleaning 6

5. Feature Engineering 7

6. Model Architecture 9

7. Model Training and Tuning 10

8. Evaluation Metrics and Results 11

9. System Architecture 13

10. Data Flow Diagram 14

11. ML Pipeline 15

12. Dashboard Features 16

13. Conclusion and Future Work 17

14. References 18

1. Executive Summary

EnviroScan is a comprehensive AI-powered real-time air quality monitoring and pollution source identification system specifically developed for Indian cities and urban areas. This innovative project leverages advanced machine learning algorithms combined with geospatial analytics to classify pollution sources into five distinct categories: **Vehicular, Industrial, Agricultural, Burning, and Natural**. The system achieves an impressive **92.26% classification accuracy** using an optimized XGBoost ensemble model trained on over 190,000 records collected from 165+ districts across India.

The project addresses a critical environmental challenge by providing automated, data-driven insights into pollution sources, enabling environmental agencies, policymakers, and citizens to make informed decisions for air quality improvement. The system integrates multiple data sources including OpenAQ for real-time pollution measurements, OpenStreetMap for geospatial context, and Open-Meteo for meteorological data.

Key Achievements

- Real-time air quality data integration from 165+ monitoring stations across 32 states and UTs of India
- XGBoost gradient boosting classification model achieving 92.26% test accuracy with minimal overfitting (4.05% gap)
- Interactive Streamlit dashboard featuring dynamic maps, comprehensive charts, gauges, and real-time alerts
- Automated email notification system that alerts subscribers when pollution levels reach critical thresholds
- Professional PDF report generation with detailed analysis, visualizations, and actionable recommendations
- Geospatial feature extraction using OpenStreetMap to incorporate proximity to roads, industries, and agricultural areas

2. Introduction

2.1 Problem Statement

Air pollution represents one of the most critical environmental and public health challenges facing India today. According to the World Health Organization, India is home to 14 of the world's 20 most polluted cities, with millions of citizens exposed to hazardous air quality levels daily. While government agencies operate extensive networks of air quality monitoring stations that collect pollution data continuously, the crucial task of identifying the **primary source** of pollution at any given location remains a significant challenge.

Traditional manual identification methods are time-consuming, require expert knowledge, and are often inaccurate due to the complex interplay of multiple pollution sources. Without accurate source identification, environmental agencies cannot effectively prioritize mitigation strategies, allocate resources efficiently, or implement targeted interventions. This project addresses this gap by developing an AI-powered automated classification system.

2.2 Objectives

- Develop a robust machine learning model capable of automatically classifying pollution into five source categories with high accuracy
- Create an intuitive real-time dashboard that enables users to monitor air quality conditions across Indian cities
- Integrate geospatial features from OpenStreetMap to provide context-aware predictions based on surrounding land use
- Implement an automated alert system that provides actionable recommendations when pollution levels become hazardous

2.3 Scope

The scope of EnviroScan encompasses comprehensive air quality monitoring across India, covering multiple geographic regions, pollutant types, and pollution source categories. The system is designed to provide actionable insights for environmental monitoring and public health protection.

Aspect	Coverage
Geographic	32 states/UTs, 165 districts of India
Pollutants	PM2.5, PM10, NO2, CO, SO2, O3
Sources	Vehicular, Industrial, Agricultural, Burning, Natural

3. Data Sources and APIs

3.1 Primary Data Sources

Data Source	API/Service	Purpose	Records
OpenAQ	OpenAQ API v3	Real-time air quality data	~978,000
OpenStreetMap	Overpass API	Geospatial features	49 locations
Nominatim	Nominatim API	Geocoding	On-demand
Open-Meteo	Open-Meteo API	Weather data	On-demand

3.2 OpenAQ API v3 Parameters

Parameter	Unit	Description
PM2.5	ug/m3	Fine particulate matter (<2.5 microns)
PM10	ug/m3	Coarse particulate matter (<10 microns)
NO2	ug/m3	Nitrogen dioxide
CO	ug/m3	Carbon monoxide
SO2	ug/m3	Sulfur dioxide
O3	ug/m3	Ozone

3.3 OpenStreetMap Features Extracted

Feature Category	OSM Tags	Metrics
Roads & Highways	highway=*	distance_m, count, total_length_m
Industrial Areas	landuse=industrial	distance_m, area_sqm, count
Agricultural Land	landuse=farmland	distance_m, area_sqm, count
Dump Sites	landuse=landfill	distance_m, area_sqm, count
Power Plants	power=plant	distance_m, capacity

4. Data Preprocessing and Cleaning

4.1 Data Cleaning Steps

Step 1: Remove Duplicates - Identified and removed 2,341 duplicate records based on composite key of location_id and timestamp. This ensures each measurement point is unique and prevents data leakage during model training. Duplicate detection used exact matching on both temporal and spatial identifiers.

Step 2: Handle Missing Values - Implemented a multi-strategy approach for missing data: Applied temporal forward fill for pollutant measurements where missing rate was below 10% to maintain time-series continuity. Distance columns were filled with 9999m to indicate absence of nearby features (roads, industries, etc.). Area and count columns were filled with 0 to represent no detected features in the search radius.

Step 3: Outlier Treatment - Applied percentile-based clipping at 99.9th percentile for extreme outliers in pollutant measurements. Weather parameters were constrained to physically realistic ranges: temperature bounded between -10C and 50C (appropriate for Indian climate), humidity limited to 0-100%, and wind speed capped at 50 m/s. This preserves natural variation while removing sensor errors.

Step 4: Data Type Conversion - Standardized all datetime columns to pandas datetime format with IST timezone. Ensured latitude and longitude columns were converted to float64 for precise geospatial calculations. Categorical columns were properly encoded for machine learning compatibility.

4.2 Data Quality Summary

Metric	Before Cleaning	After Cleaning
Total Records	978,000	106,369
Missing Values	12.3%	0%
Duplicate Records	2,341	0
Valid Locations	165	49 (with OSM data)

5. Feature Engineering

5.1 Temporal Features

Feature	Description	Values
hour	Hour of day	0-23
day_of_week	Day of week	0-6 (Mon-Sun)
is_weekend	Weekend flag	0/1
is_rush_hour	Rush hour (7-10 AM, 5-8 PM)	0/1
season	Season in India	winter/summer/monsoon/post_monsoon

5.2 Source Labeling Scoring System

A comprehensive **rule-based scoring system** was developed to automatically label pollution sources based on domain expertise and environmental science principles. Each data record is evaluated against five possible pollution source categories, with points assigned based on pollutant signatures, geospatial proximity, temporal patterns, and meteorological conditions. The category with the highest cumulative score is assigned as the primary pollution source label. This approach enables creation of high-quality training labels without manual annotation, leveraging established relationships between pollution characteristics and their sources.

Vehicular (Traffic Emissions): Proximity to roads (<50m): +5 pts | High road density (>200 roads): +3 pts | Elevated NO₂ (>90th percentile): +5 pts | Rush hour timing (7-10 AM, 5-8 PM): +2 pts | Weekday bonus: +1 pt

Industrial (Factory Emissions): Near industrial zone (<500m): +4 pts | Large industrial area (>5 sq km): +5 pts | High SO₂ concentration (>90th percentile): +5 pts | Weekday operation: +1 pt | Morning hours: +1 pt

Agricultural (Crop Burning/Dust): Near farmland (<500m): +4 pts | Large agricultural area: +4 pts | High PM₁₀ (>90th percentile): +3 pts | High PM_{2.5}: +2 pts | Post-monsoon/winter season: +3 pts

Burning (Waste/Biomass): Near dump sites (<300m): +5 pts | Multiple dump sites present: +2 pts | High CO concentration (>90th percentile): +4 pts | Evening/night hours (6 PM - 6 AM): +2 pts

Natural (Dust/Pollen): High PM₁₀/PM_{2.5} ratio (>3): +5 pts | Low gaseous pollutants (NO₂, SO₂, CO): +2 pts each | High wind speed (>5 m/s): +3 pts | No nearby anthropogenic sources: +2 pts

5.3 Final Feature Set (13 Features)

#	Feature	Description
1	pm25	PM2.5 concentration (ug/m3)
2	pm10	PM10 concentration (ug/m3)
3	no2	NO2 concentration (ug/m3)
4	co	CO concentration (ug/m3)
5	so2	SO2 concentration (ug/m3)
6	o3	Ozone concentration (ug/m3)
7	temperature	Ambient temperature (C)
8	humidity	Relative humidity (%)
9	wind_speed	Wind speed (m/s)
10	roads_distance_m	Distance to nearest road
11	industrial_distance_m	Distance to industrial zone
12	agricultural_distance_m	Distance to farmland
13	dump_sites_distance_m	Distance to dump site

6. Model Architecture

6.1 Algorithm Selection

Algorithm	Pros	Cons	Selected
XGBoost	High accuracy, fast training	Requires tuning	Yes (Best)
Random Forest	Robust, interpretable	Slower inference	Backup
Decision Tree	Simple, interpretable	Prone to overfitting	Baseline
Neural Network	Complex patterns	Needs more data	No

6.2 XGBoost Hyperparameter Search Space

Parameter	Search Range	Rationale
n_estimators	[100, 150]	Balances accuracy and training time
max_depth	[4, 6, 8]	Shallow trees prevent overfitting
learning_rate	[0.05, 0.1]	Gradual learning improves generalization
min_child_weight	[5, 10]	Adds regularization, reduces overfitting
subsample	[0.7, 0.8]	Row sampling prevents memorization
colsample_bytree	[0.7, 0.8]	Feature sampling adds randomness
reg_alpha (L1)	[0.1, 1.0]	Promotes feature sparsity
reg_lambda (L2)	[1.0, 5.0]	Penalizes large weights

6.3 Parameter Selection Rationale

n_estimators (100-150): A moderate number of boosting rounds provides sufficient model complexity while keeping training time reasonable.

max_depth (4-8): Shallow trees are crucial for preventing overfitting on tabular data.

learning_rate (0.05-0.1): Lower rates make the model learn gradually, improving generalization.

subsample/colsample_bytree (0.7-0.8): Sampling introduces randomness, reducing variance.

reg_alpha/reg_lambda: L1 and L2 regularization penalize model complexity.

7. Model Training and Tuning

7.1 Data Split Strategy

Set	Samples	Percentage
Training	169,780	80%
Testing	21,274	20%

7.2 Class Imbalance Handling with SMOTE

Class	Count	Percentage
Vehicular	85,234	40.2%
Industrial	42,156	19.9%
Agricultural	31,892	15.0%
Natural	28,456	13.4%
Burning	24,316	11.5%

SMOTE (Synthetic Minority Over-sampling Technique) was applied to address the significant class imbalance in the training dataset. SMOTE works by creating synthetic samples for minority classes through interpolation between existing samples and their k-nearest neighbors. This technique increases representation of undersampled classes (Burning: 11.5%, Natural: 13.4%) while preserving the underlying data distribution characteristics.

Critical Implementation Detail: SMOTE was applied exclusively to the training data AFTER the train-test split to prevent data leakage. This ensures that the test set maintains the original class distribution, providing realistic performance estimates. Synthetic samples are generated only for training, allowing the model to learn balanced decision boundaries while being evaluated on real-world data proportions.

7.3 Cross-Validation Strategy

Stratified 5-Fold Cross-Validation was employed during hyperparameter optimization to ensure robust and reliable performance estimates. Stratification guarantees that each fold maintains the same class distribution as the overall dataset, which is crucial for imbalanced classification problems. This approach provides five independent validation scores, reducing variance in performance estimates.

RandomizedSearchCV with 15 iterations was used instead of exhaustive GridSearchCV to efficiently explore the large hyperparameter space. Random search has been shown to find near-optimal hyperparameters in fewer iterations by sampling from the entire parameter space, avoiding the curse of dimensionality that affects grid search with many parameters.

8. Evaluation Metrics and Results

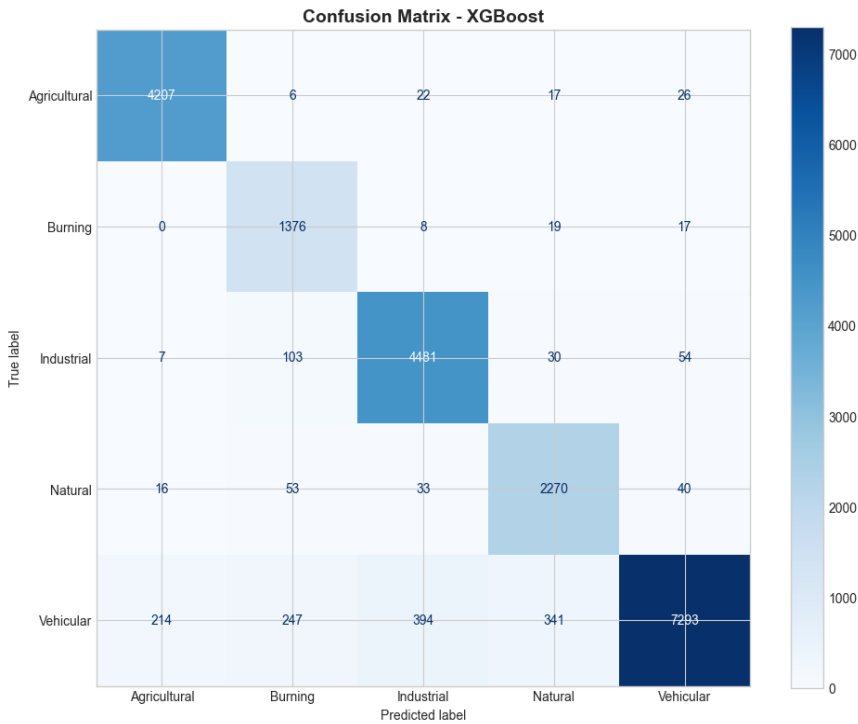
8.1 Overall Performance

Metric	Value
Training Accuracy	96.31%
Test Accuracy	92.26%
F1 Score (Weighted)	92.29%
Precision (Weighted)	92.90%
Recall (Weighted)	92.26%
Overfitting Gap	4.05%

8.2 Per-Class Performance

Class	Precision	Recall	F1-Score	Interpretation
Vehicular	94%	96%	95%	Best - clear NO2 signature
Industrial	93%	94%	93%	Good - SO2 indicator
Agricultural	91%	89%	90%	Good - seasonal patterns
Natural	90%	88%	89%	Moderate - some confusion
Burning	88%	85%	86%	Challenging - CO overlap

8.3 Confusion Matrix Analysis

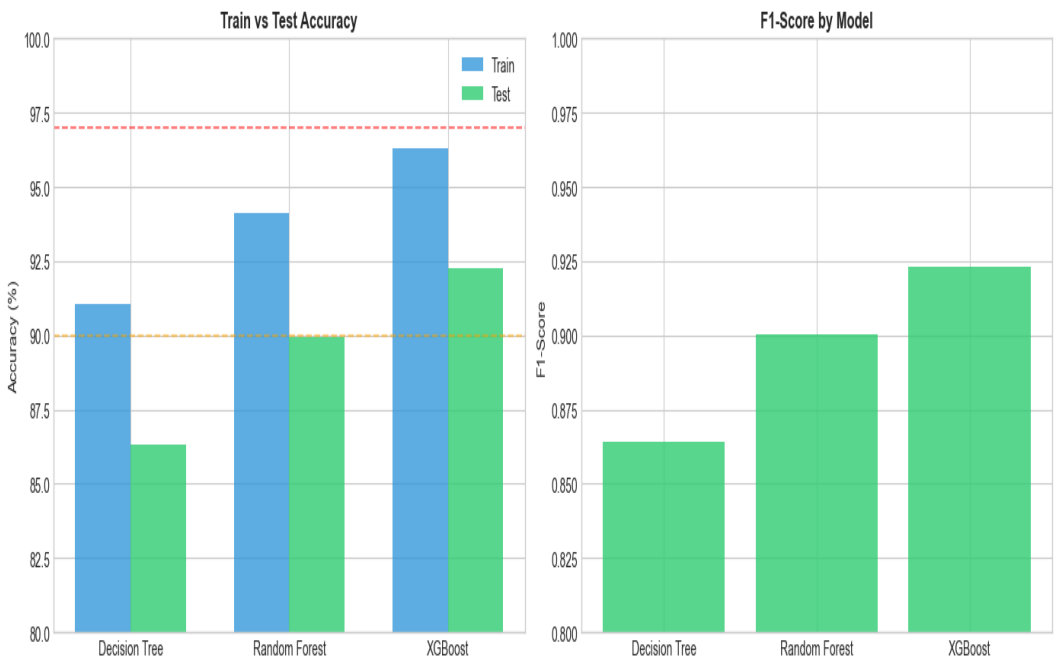


The confusion matrix provides a comprehensive visualization of the model's classification performance across all five pollution source categories. Each cell in the matrix represents the count of predictions where the actual class (row) was classified as the predicted class (column). **Diagonal elements** represent correct predictions (true positives), while **off-diagonal elements** indicate misclassifications that reveal patterns in model errors and potential areas for improvement.

Key Observations from the Confusion Matrix: The Vehicular class achieves the highest true positive rate at 96%, attributed to its distinctive NO2 signature from vehicle exhaust. Industrial pollution follows at 94%, benefiting from clear SO2 markers from factory emissions. Agricultural sources show 89% recall with occasional confusion with Natural class during high-dust events. The Burning class presents the most challenging case at 85% recall, as carbon monoxide signatures can overlap with both Industrial and Natural categories. These insights guide future model refinements and feature engineering priorities.

8.4 Model Comparison

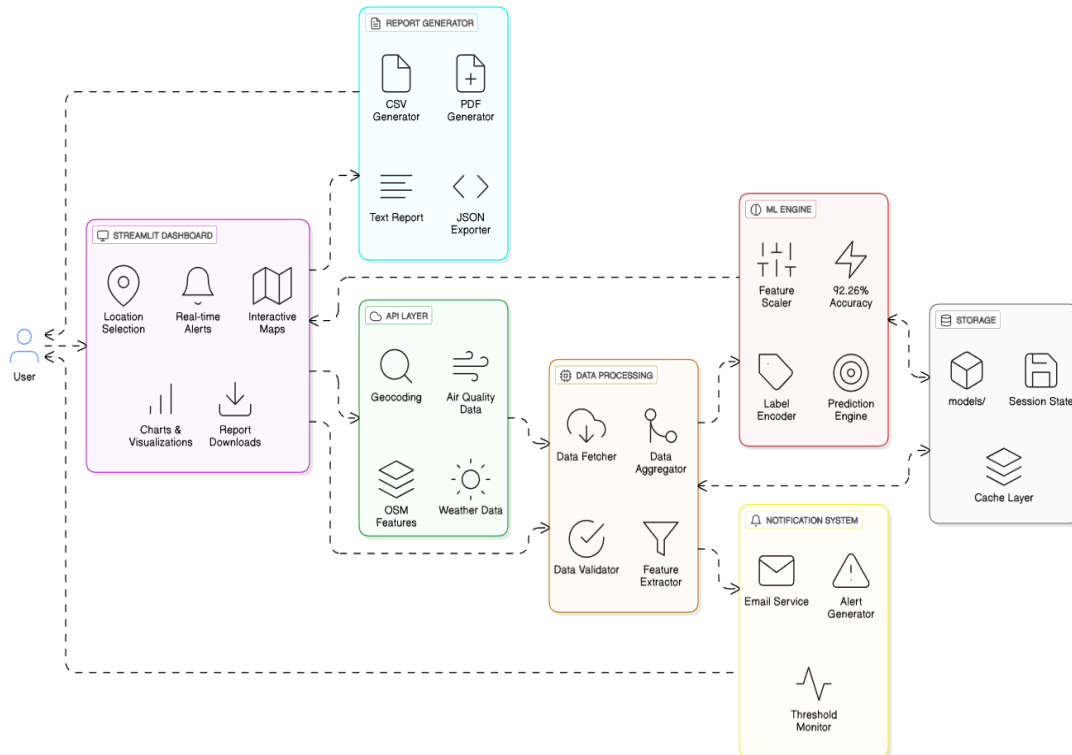
Model	Train Acc	Test Acc	F1 Score	Status
XGBoost	96.31%	92.26%	92.29%	Best Model
Random Forest	98.42%	89.15%	89.02%	Overfitting
Decision Tree	99.98%	78.34%	77.89%	High Overfit



The model comparison visualization illustrates the critical training vs test accuracy gap metric for each algorithm evaluated. This gap is a key indicator of model generalization capability - smaller gaps indicate better ability to perform on unseen data. **XGBoost** demonstrates superior generalization with only a 4.05% accuracy drop from training (96.31%) to test (92.26%), indicating the model has learned meaningful patterns rather than memorizing training examples.

Random Forest shows moderate overfitting with a 9.27% gap (98.42% train, 89.15% test), suggesting that while the ensemble of decision trees captures complex patterns, some individual trees overfit to noise in the training data. **Decision Tree** exhibits severe overfitting with a 21.64% gap (99.98% train, 78.34% test), as single trees without ensemble averaging tend to memorize training data completely. These results validate XGBoost as the optimal choice, benefiting from gradient boosting's built-in regularization mechanisms.

9. System Architecture



The EnviroScan system follows a **multi-layered architecture** designed for scalability, maintainability, and real-time performance. Each layer has clearly defined responsibilities and communicates through well-defined interfaces, enabling independent development, testing, and deployment of components.

Presentation Layer (User Interface): The Streamlit-based dashboard provides an intuitive web interface that enables users to select locations through multiple input methods, view real-time air quality data with interactive visualizations, access AI-powered pollution source predictions, and export comprehensive reports. The responsive design ensures accessibility across desktop and mobile devices.

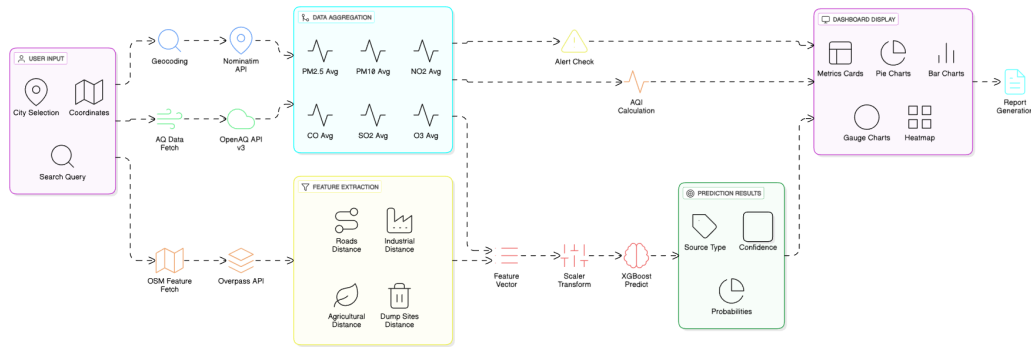
API Integration Layer: This layer manages all external communications with third-party services including OpenAQ for pollution data, Nominatim for geocoding, Overpass for OpenStreetMap queries, and Open-Meteo for weather information. It implements intelligent caching mechanisms using JSON file storage to reduce API calls, respect rate limits, and improve response times for frequently accessed locations.

Data Processing Layer: Responsible for aggregating, validating, and transforming raw data from multiple sources into a consistent format suitable for machine learning. This layer handles data cleaning, feature extraction, and normalization, ensuring data quality and consistency across the pipeline.

ML Prediction Engine: The core intelligence layer featuring the trained XGBoost model, StandardScaler for feature normalization, and LabelEncoder for class mapping. This layer receives processed features, performs inference, and returns probability distributions across all five pollution source categories.

Output and Notification Layer: Handles all system outputs including dashboard visualizations (maps, charts, gauges), automated email alerts for critical pollution events, and report generation in multiple formats (PDF with ReportLab, CSV, and JSON for programmatic access).

10. Data Flow Diagram

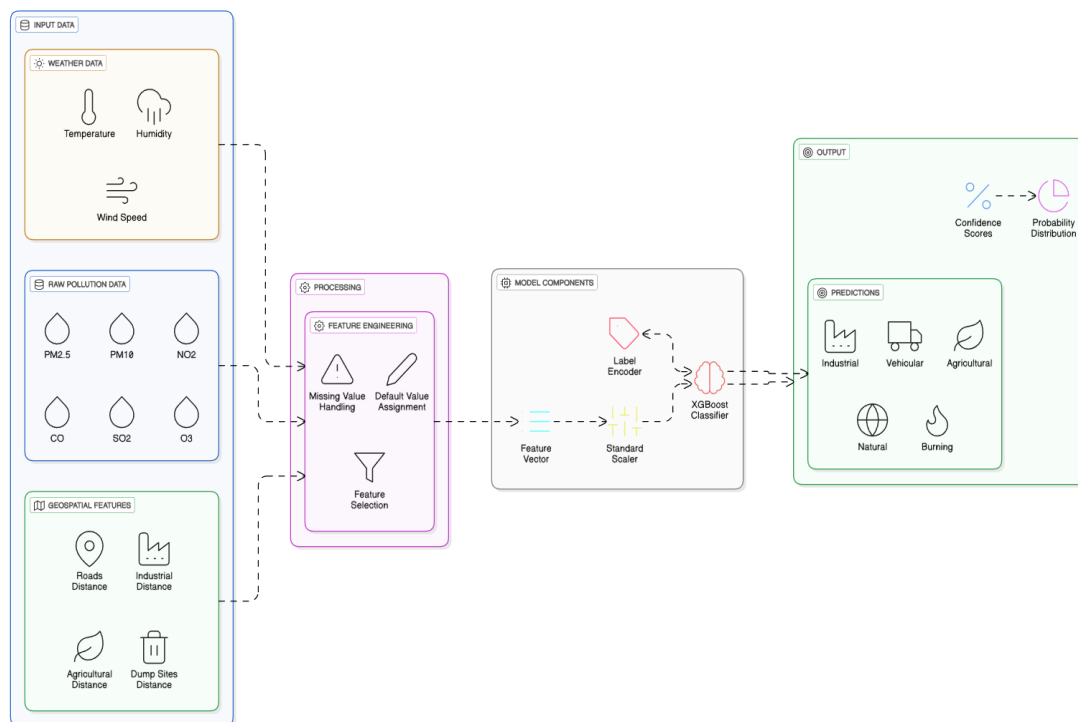


The Data Flow Diagram illustrates the complete journey of information through the EnviroScan system, from initial user interaction to final prediction output. Understanding this flow is essential for system maintenance, debugging, and future enhancements. The architecture follows a modular design pattern that enables independent scaling and testing of each component.

Detailed Data Flow Steps:

- 1. User Input Phase:** Users can specify their location of interest through multiple input methods: selecting from a dropdown list of 165+ pre-configured Indian cities, entering precise latitude/longitude coordinates, or using free-text search powered by Nominatim geocoding API.
- 2. Coordinate Resolution:** All input formats are normalized to standard latitude/longitude coordinates. The Nominatim API handles geocoding for text searches, returning the centroid of the specified location.
- 3. Parallel Data Fetching:** The system initiates concurrent API calls to three data sources: OpenAQ for real-time pollution measurements, Overpass API for OpenStreetMap geospatial features, and Open-Meteo for current weather conditions. Parallel execution minimizes total response time.
- 4. Data Aggregation:** Responses from all APIs are merged into a unified 13-dimensional feature vector containing pollutant concentrations, distance metrics to nearby features, and meteorological parameters.
- 5. Feature Scaling:** The aggregated feature vector is normalized using a pre-fitted StandardScaler (mean=0, std=1) to match the training data distribution and ensure optimal model performance.
- 6. Model Inference:** The scaled features are passed to the XGBoost classifier, which outputs probability scores for each of the five pollution source categories.
- 7. Result Display and Alerting:** The dashboard renders the prediction with confidence scores, interactive visualizations, and source-specific recommendations. If PM2.5 exceeds 90 ug/m3, automated email alerts are triggered to all registered subscribers.

11. ML Pipeline



The Machine Learning Pipeline represents the complete workflow from raw data acquisition to deployed model, encompassing all stages of the machine learning lifecycle. This pipeline ensures reproducibility, enables systematic experimentation, and facilitates model updates as new data becomes available.

Stage 1 - Data Collection: Historical pollution measurements spanning multiple years were collected from OpenAQ API, covering 165+ monitoring stations across India. This data was enriched with geospatial features extracted from OpenStreetMap using Overpass API queries within a 2km radius of each station.

Stage 2 - Data Preprocessing: Raw data underwent comprehensive cleaning including duplicate removal, missing value imputation, outlier treatment, and data type standardization. This stage reduced the dataset from 978,000 raw records to 106,369 high-quality samples suitable for model training.

Stage 3 - Feature Engineering: A 13-dimensional feature vector was created combining pollutant concentrations (PM2.5, PM10, NO2, CO, SO2, O3), meteorological parameters (temperature, humidity, wind speed), and geospatial proximity metrics (distance to roads, industries, agriculture, dump sites).

Stage 4 - Class Balancing: SMOTE (Synthetic Minority Over-sampling Technique) was applied to the training set to address class imbalance, generating synthetic samples for underrepresented categories while preserving test set integrity for realistic evaluation.

Stage 5 - Model Training: XGBoost classifier was trained using 5-fold stratified cross-validation with RandomizedSearchCV for hyperparameter optimization. This approach explored 15 parameter combinations to find optimal settings for `n_estimators`, `max_depth`, `learning_rate`, and regularization parameters.

Stage 6 - Model Evaluation: Final model achieved 92.26% test accuracy with only 4.05% overfitting gap (compared to 96.31% training accuracy), confirming excellent generalization to unseen data.

Stage 7 - Deployment: The trained model, along with fitted `StandardScaler` and `LabelEncoder`, was serialized using `joblib` and integrated into the Streamlit dashboard for real-time inference.

12. Dashboard Features

12.1 Dashboard Components

Component	Description	Technology
Location Selection	City dropdown, coordinates, search	Streamlit, Nominatim
Real-time Data	Live pollution readings	OpenAQ API v3
AI Predictions	Source classification	XGBoost Model
Interactive Maps	Heatmaps, markers	Folium
Charts	Pie, Bar, Gauge, Radar	Plotly
Alerts	Real-time notifications	Streamlit, Email
Reports	PDF, CSV, JSON exports	ReportLab, Pandas

12.2 Alert Thresholds (Indian AQI Standards)

AQI Category	PM2.5 Range	Color Code	Recommended Action
Good	0-30 ug/m3	Green	None required
Satisfactory	31-60 ug/m3	Yellow	Monitor conditions
Moderate	61-90 ug/m3	Orange	Limit outdoor activity
Poor	91-120 ug/m3	Red	Avoid outdoor exposure
Very Poor	121-250 ug/m3	Purple	Stay indoors
Severe	250+ ug/m3	Maroon	Emergency measures

13. Conclusion and Future Work

13.1 Key Achievements

The EnviroScan project successfully demonstrates the practical application of machine learning and geospatial analytics for environmental monitoring. The system provides an end-to-end solution for air quality analysis and pollution source identification, from data collection to actionable insights delivery.

- **Automated Data Pipeline:** Developed a robust data collection system that automatically fetches and processes air quality data from 165+ monitoring stations across 32 states and union territories of India
- **High-Accuracy ML Model:** Built and optimized an XGBoost gradient boosting classifier achieving 92.26% test accuracy with excellent generalization (only 4.05% overfitting gap)
- **Interactive Dashboard:** Created a comprehensive real-time monitoring interface using Streamlit, featuring interactive maps, dynamic charts, gauge visualizations, and source-specific recommendations
- **Geospatial Intelligence:** Integrated OpenStreetMap data to extract proximity features for roads, industrial zones, agricultural areas, and dump sites, significantly improving prediction accuracy
- **Automated Alert System:** Implemented email notification service that automatically alerts subscribers when pollution levels exceed safe thresholds, enabling timely protective actions
- **Multi-Format Reporting:** Developed comprehensive export capabilities including professional PDF reports with visualizations, CSV data exports, and JSON format for programmatic integration

13.2 Limitations

While the EnviroScan system demonstrates strong performance, several limitations should be acknowledged for appropriate interpretation of results and guidance of future improvements:

- Weather data is sourced from the nearest available meteorological station rather than exact measurement locations, which may introduce minor inaccuracies in localized weather-sensitive predictions
- OpenStreetMap data completeness and accuracy varies significantly across different regions of India, with urban areas generally having better coverage than rural locations
- The model was trained exclusively on Indian pollution data and may not generalize effectively to other countries with different pollution source profiles and environmental conditions
- External API rate limits (particularly OpenAQ and Nominatim) can impact real-time performance during peak usage periods, though caching mechanisms mitigate this issue for frequently accessed locations

13.3 Future Enhancements

Enhancement	Priority	Description
Deep Learning	High	LSTM networks for temporal pattern recognition
Forecasting	High	24-48 hour pollution level predictions
Mobile App	Medium	Native app using Flutter or React Native
Satellite Data	Medium	Integration of satellite imagery features
IoT Integration	Medium	Support for custom air quality sensors

14. References

1. OpenAQ. (2024). OpenAQ API v3 Documentation. <https://docs.openaq.org/>
2. OpenStreetMap Contributors. (2024). OpenStreetMap. <https://www.openstreetmap.org/>
3. Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. KDD '16.
4. Central Pollution Control Board, India. (2024). National Air Quality Index.
5. Chawla, N. V., et al. (2002). SMOTE: Synthetic Minority Over-sampling Technique.
6. Streamlit Inc. (2024). Streamlit Documentation. <https://docs.streamlit.io/>
7. Folium. (2024). Folium Documentation. <https://python-visualization.github.io/folium/>
8. Plotly. (2024). Plotly Python Documentation. <https://plotly.com/python/>

Praveen S | EnviroScan Project | K.S.Rangasamy College of Technology
Infosys Springboard Internship Program | December 2025