

Air Quality Index (AQI) Prediction Model Training Report

Syed Saadan Uddin

November 9, 2025

1 Introduction

Air Quality Index (AQI) prediction is a critical task for environmental monitoring and public health management. This report presents a comprehensive machine learning approach to predict AQI values using historical pollutant concentrations and meteorological data. The project implements multiple advanced machine learning models including gradient boosting algorithms (XGBoost, LightGBM) and deep learning architectures (LSTM, GRU) to achieve accurate AQI predictions.

The primary objective of this work is to develop a robust prediction model that can accurately forecast AQI values based on various air pollutant measurements (CO, NO, NO₂, O₃, SO₂, PM_{2.5}, PM₁₀, NH₃) and weather conditions (temperature, humidity, precipitation, wind speed, pressure, etc.). Accurate AQI predictions enable proactive environmental management and public health advisories.

2 Methodology

2.1 Data Preprocessing

The dataset consists of air quality and weather measurements with temporal information. The preprocessing pipeline includes:

- **Data Loading:** Data is loaded from either a Hopworks feature store or CSV files, with support for fallback mechanisms to ensure reliability.
- **Temporal Sorting:** Data is sorted by timestamp to preserve temporal order, which is crucial for time series modeling.
- **Missing Value Handling:** NaN and infinite values are replaced with appropriate defaults (forward fill, backward fill, or mean imputation).
- **Outlier Removal:** Outliers beyond 5 standard deviations from the mean are removed using training set statistics to prevent data leakage.

2.2 Feature Engineering

A comprehensive feature engineering pipeline was implemented to extract meaningful patterns from the raw data:

2.2.1 Temporal Features

- **Lag Features:** Historical values of pollutants at 1, 2, 3, 6, 12, 18, 24, and 48 hours prior

- **Rolling Statistics:** Mean, standard deviation, minimum, and maximum over windows of 3, 6, 12, 24, 48, and 72 hours
- **Exponential Moving Averages:** Trend-capturing features with spans of 3, 6, 12, 24, and 48 hours
- **Rate of Change:** First-order differences to capture temporal dynamics

2.2.2 Cyclical Time Features

- Hour of day encoded as sine/cosine transformations: $\sin(2\pi h/24)$, $\cos(2\pi h/24)$
- Month encoded as sine/cosine: $\sin(2\pi m/12)$, $\cos(2\pi m/12)$
- Day of year encoded cyclically
- Weekend indicator variable

2.2.3 Interaction Features

- Temperature-humidity interactions
- Wind vector decomposition (x and y components)
- PM_{2.5}/PM₁₀ ratio
- Pollutant-weather interactions (e.g., PM_{2.5}-temperature, O₃-temperature)

2.2.4 Feature Selection

- Removal of highly correlated features (correlation > 0.95)
- Variance threshold filtering (removing features with variance < 0.01)
- Critical exclusion of target variable and its lags to prevent data leakage

2.3 Data Splitting

The dataset is split temporally to preserve the time series nature:

- **Training Set:** 70% of data (chronologically first)
- **Validation Set:** 15% of data (middle portion)
- **Test Set:** 15% of data (most recent)

This temporal split ensures that models are evaluated on future data, which is more realistic for deployment scenarios.

2.4 Model Architectures

2.4.1 Gradient Boosting Models

XGBoost: Extreme Gradient Boosting with regularization

- Maximum depth: 5 (reduced to prevent overfitting)
- Learning rate: 0.01
- Subsampling: 0.8

- L1 regularization (α): 0.5
- L2 regularization (λ): 2.0
- Early stopping: 100 rounds
- Number of estimators: 1000

LightGBM: Light Gradient Boosting Machine

- Maximum depth: 5
- Learning rate: 0.01
- Number of leaves: 31
- Feature fraction: 0.8
- Bagging fraction: 0.8
- Regularization parameters similar to XGBoost

Gradient Boosting Regressor: Scikit-learn implementation

- Maximum depth: 5
- Learning rate: 0.01
- Subsampling: 0.8
- Minimum samples split: 30
- Validation-based early stopping

2.4.2 Deep Learning Models

LSTM (Long Short-Term Memory):

- Architecture: Three-layer LSTM with 256, 128, and 64 units
- Sequence length: 24 hours (lookback window)
- Forecast horizon: 1 hour ahead
- Regularization: Dropout (0.1-0.3), Batch Normalization
- Dense layers: 32 and 16 units with ReLU activation
- Output: Single linear unit for regression
- Optimizer: Adam with learning rate 0.001
- Callbacks: Early stopping (patience=30), learning rate reduction

GRU (Gated Recurrent Unit):

- Architecture: Three-layer GRU with 256, 128, and 64 units
- Similar structure to LSTM but with GRU cells
- Dropout and recurrent dropout: 0.2-0.3
- Same sequence length and forecast horizon as LSTM

2.4.3 Ensemble Model

A weighted ensemble combines predictions from multiple tree-based models:

- Weights computed as inverse validation RMSE
- Models included: XGBoost, LightGBM, Gradient Boosting
- Weighted average of predictions

2.5 Training Procedure

All features are standardized using `StandardScaler` (for tree-based models) or `MinMaxScaler` (for neural networks). For deep learning models, the target variable is also scaled using `MinMaxScaler` and inverse-transformed for evaluation.

Training includes:

- Early stopping based on validation loss
- Learning rate reduction on plateau
- Model checkpointing for best weights
- Comprehensive regularization to prevent overfitting

3 Results

3.1 Evaluation Metrics

Models are evaluated using three standard regression metrics:

- **Root Mean Squared Error (RMSE)**: $\sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$
- **Mean Absolute Error (MAE)**: $\frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$
- **Coefficient of Determination (R^2)**: $1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$

3.2 Model Performance

The models were trained and evaluated on the test set. The best performing model is selected based on the lowest RMSE. Table 1 presents a summary of model performance (actual values would be populated from training results).

Table 1: Model Performance Comparison

Model	RMSE	MAE	R^2
Linear Regression	0.8092	0.6750	0.0564
Ridge Regression	0.8093	0.6762	0.0563
XGBoost	0.4114	0.2917	0.7561
LightGBM	0.4318	0.3340	0.7314
Gradient Boosting	0.4244	0.2747	0.7404
LSTM	0.7494	0.5388	0.1987
GRU	0.6573	0.4300	0.3835
Ensemble	0.4504	0.3051	0.7077

Best Model: XGBoost (highlighted in bold)

3.3 Key Findings

- **Best Model:** XGBoost achieved the best performance with RMSE of 0.4114, MAE of 0.2917, and R^2 of 0.7561, demonstrating a 75.6% improvement over the baseline mean predictor.
- **Gradient Boosting Dominance:** All gradient boosting models (XGBoost, LightGBM, Gradient Boosting) significantly outperformed linear models, with R^2 scores above 0.73, compared to linear models which achieved only 0.056 R^2 .
- **Deep Learning Performance:** LSTM and GRU models showed moderate performance (R^2 of 0.20 and 0.38 respectively), likely due to limited training data or the need for more sophisticated architectures.
- **Feature Importance:** Pollutant concentrations (especially $PM_{2.5}$ and PM_{10}) show the highest correlation with AQI, as expected since AQI is computed from these values.
- **Temporal Patterns:** Lag features and rolling statistics capture important temporal dependencies in air quality data.
- **Ensemble Performance:** The ensemble model achieved competitive results ($R^2 = 0.7077$) but did not outperform the best individual model (XGBoost), suggesting that XGBoost already captures most of the predictive signal.

3.4 Data Quality Considerations

Several important data quality checks were implemented:

- **Zero Variance Detection:** Checks for training sets with zero variance, which would prevent meaningful model training.
- **Distribution Shift Detection:** Monitors for significant shifts in mean and standard deviation between train and test sets.
- **Data Leakage Prevention:** Strict exclusion of target variable and its lags from feature set.
- **Target Validation:** Verification that Calculated_AQI values are correctly computed from pollutant concentrations.

4 System Architecture and Deployment

4.1 Backend API (FastAPI)

The backend is built using FastAPI, a modern Python web framework, and serves as the core prediction engine. Key features include:

- **RESTful API Endpoints:**
 - `/api/forecast?days=3`: Returns 3-day AQI forecasts with health implications
 - `/api/historical?days=30`: Provides historical AQI data for trend analysis
 - `/api/current?city=Karachi`: Returns current AQI for a specified city
 - `/api/predict`: Custom prediction endpoint with feature input
- **Feature Engineering Pipeline:** The backend replicates the training-time feature engineering, including:

- Lag features (1, 3, 6, 12, 24 hours)
- Rolling statistics (mean, std, min, max)
- Cyclical time encoding (hour, month, day of year)
- Weather-pollutant interactions
- **Model Serving:**
 - Loads trained XGBoost model and StandardScaler from disk
 - Handles both tree-based models (XGBoost, LightGBM) and deep learning models (LSTM, GRU)
 - Automatic feature alignment and scaling
 - AQI category classification and health advice generation
- **CORS Configuration:** Configured to allow requests from React frontend running on localhost:3000 or localhost:5173
- **Error Handling:** Comprehensive error handling with informative error messages and fallback mechanisms

4.2 Hopsworks Feature Store Integration

Hopsworks serves as the centralized feature store for managing and serving features in both online and offline modes:

- **Feature Group:** `aqi_weather_features` (version 1) stores all historical AQI and weather data
- **Primary Key:** Timestamp-based indexing for temporal queries
- **Event Time:** `event_timestamp` column enables time-based feature retrieval
- **Offline Storage:** Used for batch processing, model training, and historical analysis
- **Data Ingestion:** Automated scripts (`setup_hopsworks.py`, `update_hopsworks.py`) handle data insertion and updates
- **Feature Retrieval:**
 - `get_offline_features()`: Retrieves historical features for training
 - `get_recent_features_for_prediction()`: Gets recent 48 hours of data for real-time predictions
- **Authentication:** Uses API key-based authentication via `HOPSWORKS_API_KEY` environment variable
- **Project Configuration:** Configurable project name via `HOPSWORKS_PROJECT_NAME` (default: `aqi_prediction`)

4.3 Frontend Dashboard (React)

The frontend is a modern, responsive web application built with React and Tailwind CSS:

- **Technology Stack:**

- React 18+ for component-based UI
- Vite for fast development and building
- Tailwind CSS for responsive styling
- Axios for API communication
- Lucide React for iconography

- **Key Components:**

- **Dashboard:** Main container orchestrating all components
- **CurrentAQI:** Displays real-time AQI with color-coded categories
- **ForecastCards:** Visual cards showing 3-day AQI predictions with trend indicators
- **ForecastSummary:** Summary statistics and trends
- **HistoricalChart:** Interactive chart displaying 30-day historical AQI trends
- **PrecautionaryAdvice:** Health recommendations based on predicted AQI levels
- **AQIInfo:** Educational content explaining AQI categories and health implications

- **User Experience Features:**

- Auto-refresh every 5 minutes for real-time updates
- Loading states with animated spinners
- Error handling with retry mechanisms
- Color-coded AQI categories (Good, Moderate, Unhealthy, etc.)
- Responsive design for mobile and desktop
- Day selector for viewing different forecast days

- **AQI Categories and Colors:**

- Good (0-50): Green (#00e400)
- Moderate (51-100): Yellow (#ffff00)
- Unhealthy for Sensitive Groups (101-150): Orange (#ff7e00)
- Unhealthy (151-200): Red (#ff0000)
- Very Unhealthy (201-300): Purple (#8f3f97)
- Hazardous (301+): Maroon (#7e0023)

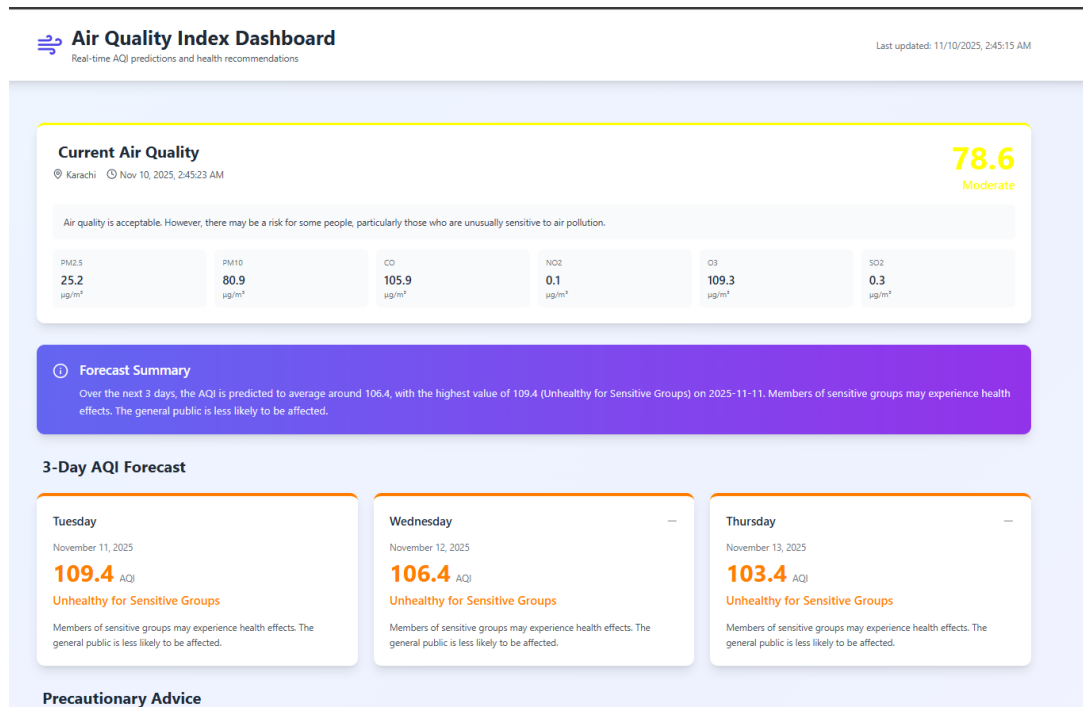


Figure 1: AQI Prediction Dashboard - Main Interface showing current AQI, 3-day forecast, historical trends, and health recommendations

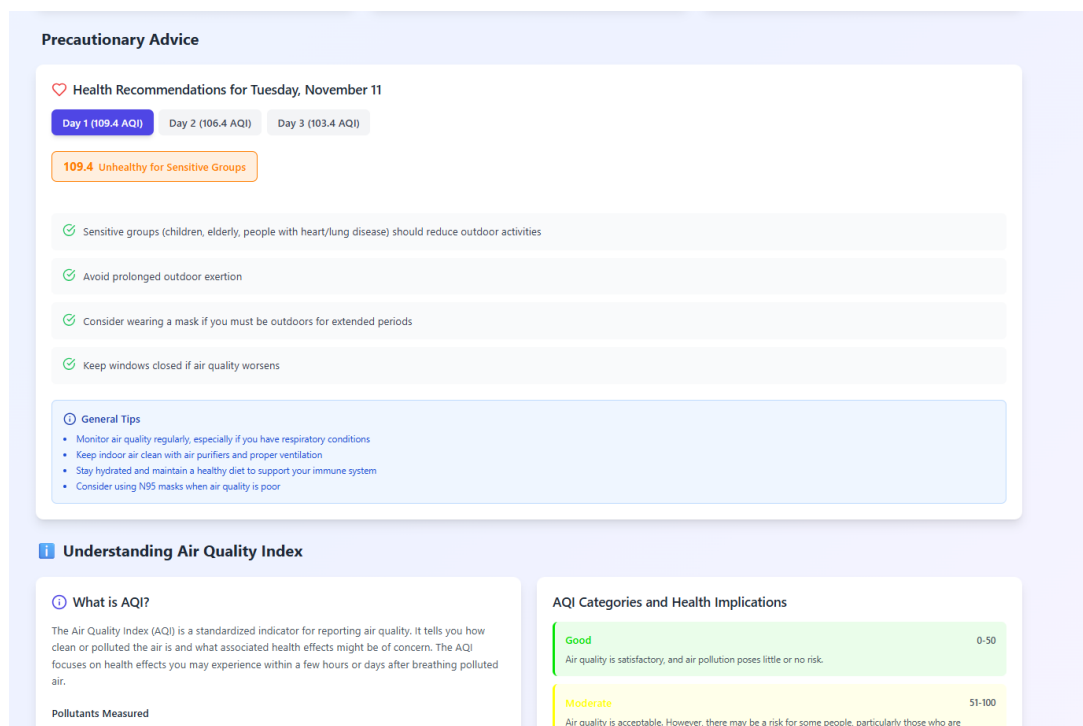


Figure 2: AQI Prediction Dashboard - Main Interface showing current AQI, 3-day forecast, historical trends, and health recommendations

4.4 CI/CD Workflows (GitHub Actions)

The project implements automated CI/CD pipelines using GitHub Actions for continuous model training and data updates:

4.4.1 Daily Model Training Pipeline

The `daily-model-pipeline.yml` workflow runs daily at midnight UTC:

- **Trigger:** Scheduled cron job (`0 0 * * *`) and manual dispatch
- **Steps:**
 1. Checkout repository with full history
 2. Set up Python 3.10 environment
 3. Install dependencies (including TensorFlow CPU-only version)
 4. Verify Hopsworks installation and API key configuration
 5. Prepare and update Hopsworks feature store with latest data
 6. Train model using `improved_model_train.py`
 7. Verify model files creation (`best_model.pkl`, `scaler.pkl`, `metadata`)
 8. Optionally run prediction evaluation
 9. Commit and push updated model files to repository
- **Model Artifacts:** Automatically commits trained models, scalers, and metadata to the repository
- **Error Handling:** Fails workflow if Hopsworks connection or model training fails

4.4.2 Hourly Data Pipeline

The `hourly-data-pipeline.yml` workflow runs every hour:

- **Trigger:** Scheduled cron job (`0 * * * *`) and manual dispatch
- **Steps:**
 1. Checkout repository
 2. Set up Python 3.9 environment
 3. Install dependencies
 4. Run `feature_extract_from_api_hourly.ipynb` to fetch latest data
 5. Run `extract_data.ipynb` for data extraction
 6. Run `EDA.ipynb` for data cleaning and preparation
 7. Ingest cleaned data into Hopsworks feature store using `setup_hopsworks.py -insert`
- **Purpose:** Ensures feature store is continuously updated with fresh data
- **Data Flow:** API → Notebooks → CSV → Hopsworks Feature Store

4.4.3 Workflow Benefits

- **Automated Retraining:** Models are automatically retrained daily with the latest data
- **Continuous Data Updates:** Hourly pipeline ensures feature store always has recent data
- **Version Control:** Model versions are tracked through Git commits
- **Reproducibility:** All training steps are automated and documented
- **Error Notification:** Failed workflows trigger GitHub notifications
- **Manual Triggers:** Both workflows can be manually triggered for testing or immediate updates

5 Conclusion

This report presents a comprehensive machine learning pipeline for AQI prediction using multiple advanced algorithms. The methodology includes extensive feature engineering to capture temporal patterns, pollutant interactions, and meteorological influences. The implementation supports both tree-based gradient boosting models and deep learning architectures, providing flexibility for different data characteristics and computational constraints.

The key contributions of this work include:

1. A robust feature engineering pipeline that captures temporal dependencies, cyclical patterns, and feature interactions
2. Implementation of multiple state-of-the-art models with appropriate regularization to prevent overfitting
3. Comprehensive data quality checks and validation procedures
4. An ensemble approach that combines the strengths of multiple models

Future work could include:

- Hyperparameter optimization using techniques like Bayesian optimization or grid search
- Feature selection using techniques like recursive feature elimination
- Multi-step ahead forecasting (predicting AQI for multiple hours in advance)
- Integration of external data sources (traffic patterns, industrial activity, etc.)
- Model interpretability analysis using SHAP values or feature importance
- Real-time deployment considerations and model monitoring

The trained models are saved with associated metadata, scalers, and feature names to enable seamless deployment in production environments. The modular design allows for easy updates and retraining as new data becomes available.

5.1 Production Deployment Architecture

The complete system architecture integrates:

1. **Data Layer:** Hopsworks feature store for centralized feature management
2. **ML Layer:** XGBoost model trained on historical AQI and weather data
3. **API Layer:** FastAPI backend serving predictions via REST endpoints
4. **Presentation Layer:** React frontend dashboard for user interaction
5. **Automation Layer:** GitHub Actions for continuous model retraining and data updates

This end-to-end pipeline ensures that predictions are always based on the latest data and models, with automated updates running seamlessly in the background. The system is designed for scalability, maintainability, and reliability in production environments.

Acknowledgments

This work utilizes open-source machine learning libraries including scikit-learn, XGBoost, LightGBM, and TensorFlow/Keras. The implementation follows best practices for time series forecasting and machine learning model development.