

AIR QUALITY INDEX (AQI) PREDICTION PROJECT REPORT

1. Project Overview

1.1 Objective

This project aimed to develop and evaluate machine learning models for accurate three-day Air Quality Index (AQI) forecasting. Three ensemble algorithms—**LightGBM, XGBoost, and Random Forest**—were compared to determine the optimal model. Reliable multi-day AQI prediction enables early warnings, pollution control, and protection of public health.

1.2 Data Description

The dataset comprises historical air quality and meteorological data stored in the Hopsworks Feature Store.

Target variables:

- target_aqi_t1: Day 1 AQI
- target_aqi_t2: Day 2 AQI
- target_aqi_t3: Day 3 AQI

Features included pollutant concentrations (PM2.5, PM10, NO₂, SO₂, CO, O₃), weather parameters (temperature, humidity, wind speed, pressure), and temporal attributes (day, month, season). Rolling averages and lag variables captured short-term and seasonal trends. Data was split chronologically to maintain forecasting realism.

1.3 Model Pipeline and Architecture

A CI/CD pipeline automated the ML workflow:

- **Data Ingestion:** Continuous retrieval of air quality and meteorological data into Hopsworks.
- **Feature Engineering:** Creation of rolling means, lag features, pollutant ratios, and time encodings stored in versioned feature groups.
- **Model Training:** Multi-output regression enabled simultaneous prediction of AQI for three days. Grid search and cross-validation optimized hyperparameters to balance bias and variance.
- **Evaluation Metrics:** R², MAE, RMSE, and accuracy ensured comprehensive model assessment.

This architecture ensured reproducibility, scalability, and seamless model updates for production use.

2. Model Results Summary

2.1 Comparative Performance

All models performed strongly for short-term predictions, with decreasing accuracy over extended horizons.

- **XGBoost:** Excellent short-term performance (Test R²: 0.9587, 0.8883, 0.8137) and lowest Day 1 errors (MAE: 2.48, RMSE: 4.60).
- **Random Forest:** Strong Day 1 accuracy (R²: 0.9545) but faster decline by Day 3 (R²: 0.7588).
- **LightGBM:** Achieved best overall balance between accuracy, generalization, and interpretability—emerging as the final winner.

2.2 LightGBM (Winner Model)

Overall	Accuracy:	91.56%
Test R ² : 0.9495 (Day 1), 0.8784 (Day 2), 0.8004 (Day 3)		

LightGBM excelled in variance explanation and low prediction error (MAE: 2.67, RMSE: 5.09). Systematic **grid search optimization** fine-tuned key parameters using 3-fold cross-validation. Training–testing gaps (4.7%, 11.5%, 18.9%) indicate strong generalization and absence of overfitting, making it robust for operational deployment.

2.3 Model Interpretability (SHAP Analysis)

SHAP (SHapley Additive Explanations) clarified how each feature influenced AQI predictions:

- **Day 1:** Dominant features—current AQI, PM2.5, and AQI change rate (`aqi_diff`).
- **Day 2:** Trend and temporal features gained importance (`aqi_diff`, `day_of_year`).
- **Day 3:** PM2.5, `aqi_diff`, and meteorological variables (wind speed, pressure, SO₂) became co-dominant.

The results confirmed that LightGBM’s predictions align with atmospheric science—reflecting pollutant concentration, temporal dynamics, and meteorological impact—without anomalous or spurious patterns.

3. Comparative Analysis and Conclusion

3.1 Model Ranking and Rationale

Winner: LightGBM (Multi-Output) — **91.56% Accuracy**

LightGBM outperformed XGBoost and Random Forest, offering the best trade-off between predictive strength, computational efficiency, and interpretability. Its gradient-boosting

framework effectively captured non-linear environmental relationships and temporal dependencies.

Key Strengths:

- Consistent accuracy across all forecast horizons.
- Scientifically explainable predictions validated by SHAP.
- Fast training and deployment within the CI/CD pipeline.
- Minimal overfitting and reliable generalization to unseen data.

3.2 Conclusion

LightGBM, after hyperparameter optimization, stands as the **optimal model** for three-day AQI forecasting with an average **accuracy of 91.34%**. It achieves high Day 1 accuracy (R^2 : 0.9495, MAE: 2.67) and maintains solid performance for Days 2–3 (R^2 : 0.8784, 0.8004).

The project demonstrates that **machine learning integrated with MLOps infrastructure** can deliver operationally reliable air quality forecasts that support public health and policy decision-making. The production-ready LightGBM model provides an accurate, interpretable, and scalable foundation for real-time AQI prediction with scope for future enhancement via ensemble extensions or feature expansion.