



Project: Pearls AQI Predictor
Domain: Data Sciences

Predict the Air Quality Index (AQI) in your city in the next 3 days, using a 100% serverless stack. This project involves building an end-to-end machine learning pipeline for AQI forecasting with automated data collection, feature engineering, model training, and real-time predictions through a web dashboard.

Internee: Anoosha Khalid
Computer Systems Engineering
NED UNIVERSITY OF ENG & TECH

Submitted to 10 Pearls

Abstract

Air pollution is a growing global concern, directly impacting public health and quality of life. This project, **Pearls AQI Predictor**, is a fully serverless, end-to-end machine learning system that predicts the **Air Quality Index (AQI) for the next three days** for Karachi.

The system automates **data ingestion, feature engineering, training, deployment, and real-time predictions** using modern MLOps practices. It integrates real-time weather and pollutant data, stores engineered features in a feature store (Hopsworks), retrains models automatically using GitHub Actions, and presents predictions through an interactive Streamlit web application.

1.Introduction

Air pollution causes severe health problems such as asthma, heart disease, and respiratory infections. Monitoring air quality allows people and authorities to take preventive actions.

Traditional AQI monitoring systems focus only on real-time reporting. However, **forecasting future AQI values** allows governments, healthcare systems, and citizens to prepare in advance. Pearls AQI Predictor uses **machine learning + MLOps** to predict AQI levels for the next 72 hours using real-time environmental data and historical patterns.

2.Objective

The major objective of this project was:

- Predict AQI for the next 3 days
- Automate data collection, training, and deployment
- Store reusable features in a feature store
- Provide a real-time prediction dashboard
- Follow MLOps best practices

3.System Overview

The system consists of five main components:

1. **Data Ingestion Layer** – Collects real-time weather and pollutant data from external APIs.
2. **Feature Pipeline** – Transforms raw data into meaningful features.
3. **Feature Store (Hopworks)** – Stores and versions all feature.
4. **Training Pipeline** – Trains three models and stores the best one in a registry.
5. **Prediction Dashboard** – Displays real-time predictions using Streamlit.

4.Data Sources

- **AQICN**



This provides hourly readings, which are used both for training and real-time inference.

5.Feature Engineering Pipeline

Raw data is not directly suitable for machine learning. Therefore, a feature pipeline was developed to transform raw values into meaningful indicators.

This includes:

- Time-based features (hour, day, month)
- Rolling averages of AQI
- Lag features from previous hours
- Rate of change of AQI

These features are computed automatically and stored in the feature store.

```
df['aqi'] = df['pm25'].apply(calculate_aqi_pm25)
df = df.sort_values('date')

df['aqi_lag_1'] = df['aqi'].shift(1)
df['aqi_lag_3'] = df['aqi'].shift(3)
df['aqi_change_rate'] = df['aqi'].pct_change()

df['month'] = df['date'].dt.month
df['day_of_week'] = df['date'].dt.dayofweek
df['is_weekend'] = (df['day_of_week'] >= 5).astype(int)
```

```
# List of feature names
feature_names = [f.name for f in fg.features]
print(feature_names)
```

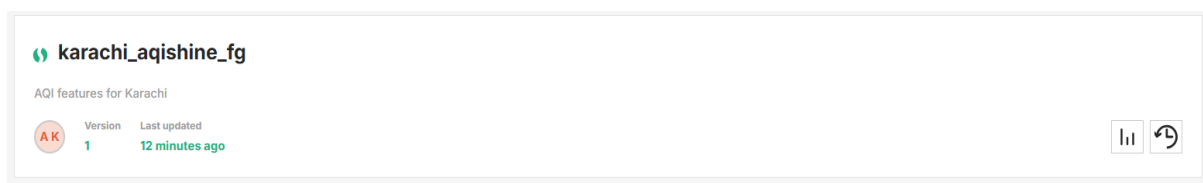
```
['date', 'pm25', 'aqi', 'aqi_lag_1', 'aqi_lag_3', 'aqi_change_rate', 'month', 'day_of_week', 'is_weekend']
```

Feature Store Using Hopsworks

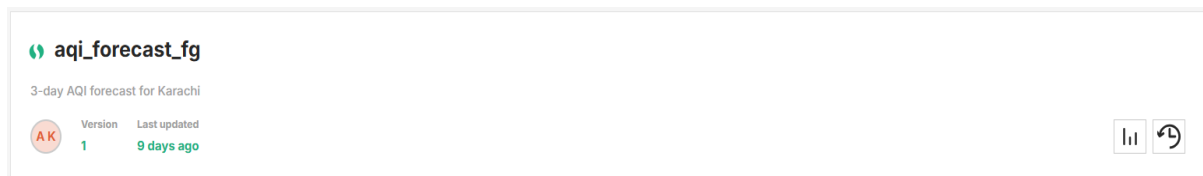
Hopsworks is used to store and manage all features. It ensures:

- Version control of features
- Consistency between training and inference
- Centralized access for pipelines

This eliminates data leakage and feature mismatches.



To store all the features



To store the prediction values, fetch by the training

Historical Data Backfilling

To train the models, historical data was generated by running the feature pipeline for previous timestamps. This created a large dataset that represents past air quality trends.

```
> karachi_aqi_last1year.csv > data
date, pm25
2024-04-16,81
2024-04-17,107
2024-04-18,70
2024-04-19,80
2024-04-20,81
2024-04-21,81
2024-04-22,91
2024-04-23,97
2024-04-25,79
2024-04-26,77
2024-04-27,79
2024-04-30,109
2024-05-02,102
2024-05-03,86
2024-05-04,87
```

Before running the feature pipeline

```
karachi_aqi_features.csv > data
date,pm25,aqi,aqi_lag_1,aqi_lag_3,aqi_change_rate,month,day_of_week,is_weekend
2024-03-07,124,222.1812434141201,242.2023182297155,224.28872497365649,-0.08266260604742237,3,3,0
2024-03-08,149,248.52476290832453,222.1812434141201,229.55742887249735,0.11856770215793211,3,4,0
2024-03-09,146,245.36354056902002,248.52476290832453,242.2023182297155,-0.012719949120203466,3,5,1
2024-03-10,130,228.5036880927292,245.36354056902002,222.1812434141201,-0.06871376422589648,3,6,1
2024-03-11,92,188.46153846153845,228.5036880927292,248.52476290832453,-0.17523633848282227,3,0,0
2024-03-12,99,195.8377239199157,188.46153846153845,245.36354056902002,0.03913894324853229,3,1,0
2024-03-13,110,207.42887249736566,195.8377239199157,228.5036880927292,0.05918751681463541,3,2,0
2024-03-14,156,300.0,207.42887249736566,188.46153846153845,0.44627889255778497,3,3,0
2024-03-15,133,231.6649104320337,300.0,195.8377239199157,-0.22778363189322093,3,4,0
2024-03-16,135,233.7723919915701,231.6649104320337,207.42887249736566,0.009097111667045743,3,5,1
```

After running the feature pipeline

6.Model Training Pipeline

The training pipeline fetches data from Hopsworks and trains multiple models, including:

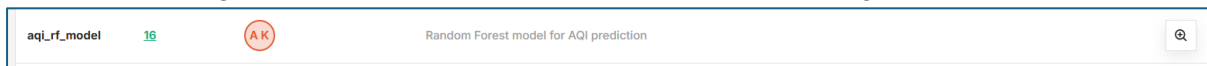
- Random Forest
- Ridge Regression
- TensorFlow Neural Network

Model Selection and Registration Strategy

After training, each model was tested on the same unseen test dataset. Their predictions were evaluated using three standard regression metrics:

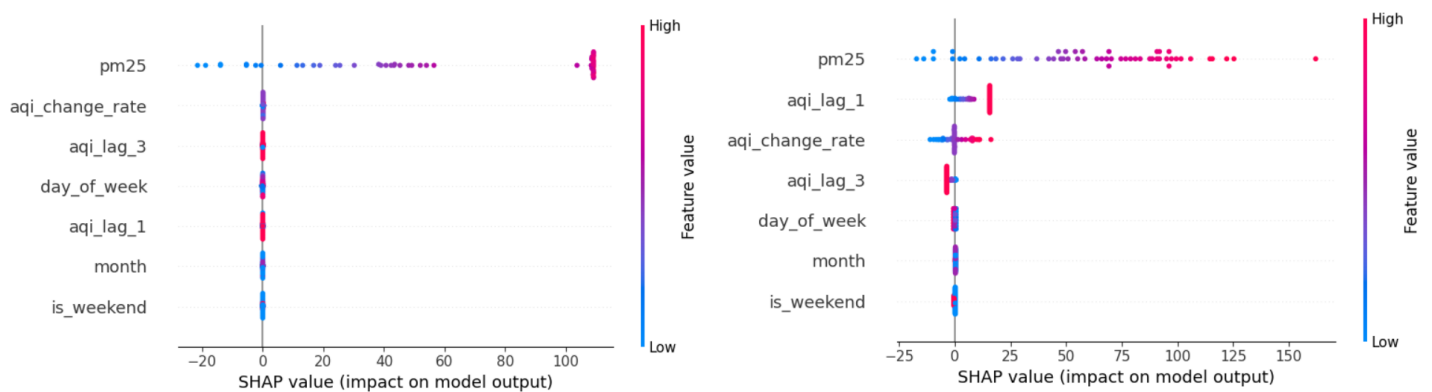
- Mean Squared Error (**MSE**) – measures the average squared difference between predicted and actual AQI values.
- Mean Absolute Error (**MAE**) – measures the average absolute prediction error.
- **R²** Score – indicates how well the model explains the variance in AQI values.

To select the best model, the system automatically identifies the model with the highest R² score, since a higher R² value means better predictive power and generalization.

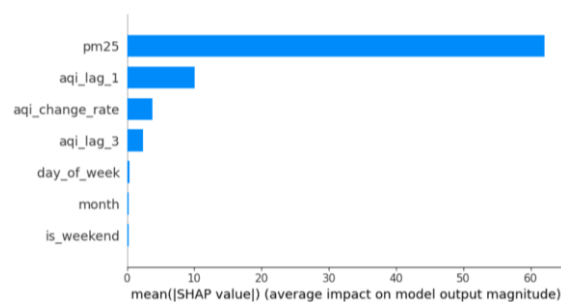


Explainability with SHAP

After selecting the best-performing model, SHAP is used to analyze how each feature contributes to the predicted AQI.



Initial : Rainforest Model & Ridge



Initial training:

Model	MSE	R2	MAE
Random Forest	0.94	1.00	0.49
Ridge Regression	263.88	0.88	11.55
Neural Network	2331.39	-0.06	38.82

7.Automation and CI/CD

All pipelines are automated using **GitHub Actions**.

- Feature pipeline runs every hour
- Training pipeline runs daily

This ensures the model stays updated without manual intervention.

✓ Feature Pipeline Feature Pipeline #220: Scheduled	main	29 minutes ago 4m 27s	...
✓ Feature Pipeline Feature Pipeline #219: Scheduled	main	Today at 12:43 AM 4m 23s	...
✓ Feature Pipeline Feature Pipeline #218: Scheduled	main	Feb 17, 11:52 PM GMT+5 4m 28s	...
✓ Feature Pipeline Feature Pipeline #217: Scheduled	main	Feb 17, 10:52 PM GMT+5 4m 40s	...

Feature Automation

✓ Training Pipeline Training Pipeline #35: Scheduled	main	Feb 17, 9:19 AM GMT+5 2m 13s	...
---	------	---------------------------------	-----

Training Automation

8.Challenges Faced

During development, several challenges were encountered:

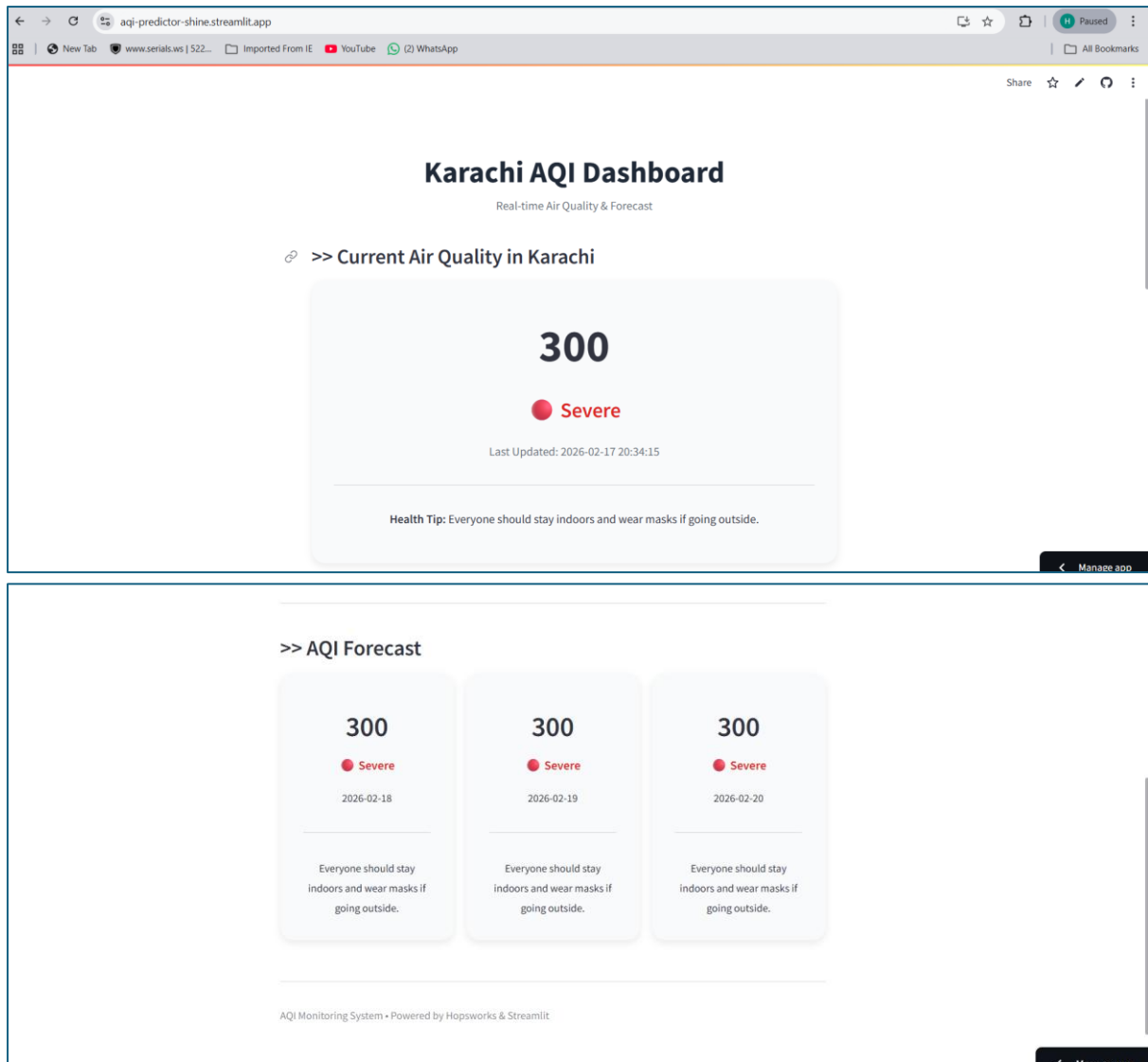
- GitHub Actions authentication failures
- Hopsworks is not compatible with VsCode
- Dependency version conflicts
- Feature store schema mismatches

These issues required debugging YAML workflows, secrets, and service permissions.

9.Web Application

The Streamlit dashboard allows users to:

- View current AQI
- See predictions for the next three days



10.Results

The system successfully predicts AQI values and updates models automatically. The automated pipelines reduce manual effort and improve reliability.

11.Limitations

- Prediction accuracy depends on API data quality
- Sudden weather changes affect forecasts
- Limited to Karachi city

Pearls AQI Predictor demonstrates a real-world machine learning system that integrates automation, feature management, and cloud deployment. It shows how predictive models can be maintained in production environments.