

AIR QUALITY INDEX IN KARACHI FOR THE NEXT 3 DAYS

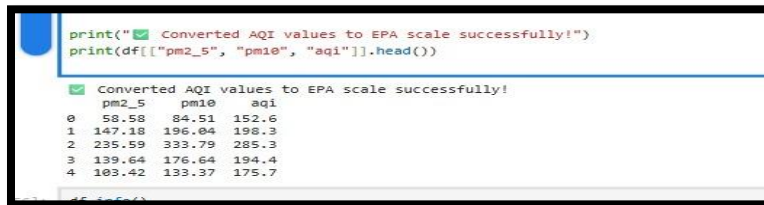
Data Collection:

We collected historical AQI and weather data from OpenWeather API over the past year. Unlike typical hourly or daily data, we chose a **5-hour interval**: Early in the data exploration, consecutive AQI values were often identical. Using such repeated values could mislead the model into memorizing these values (overfitting), reducing its predictive ability. During data collection from the OpenWeather API, the following features were retrieved for each timestamp:

Pre-Processing:

Converted to US EPA AQI Scale:

The AQI values we fetched from OpenWeather were on a simple 0–5 scale, which is too basic for real-world use and doesn't match international standards. To make the data meaningful and comparable, we converted it to the US EPA AQI scale, which ranges from 0 to 500.



```
print("Converted AQI values to EPA scale successfully!")
print(df[["pm2_5", "pm10", "aqi"]].head())
```

	pm2_5	pm10	aqi
0	58.58	84.51	152.6
1	147.18	196.04	196.3
2	235.59	333.79	285.3
3	139.64	176.64	194.4
4	103.42	133.37	175.7

Here's the approach we used:

1. Separate calculations for PM2.5 and PM10:
2. We first looked at the concentration of PM2.5 (fine particles) and PM10 (coarse particles).
3. Each concentration was mapped to its corresponding range on the EPA scale. This tells us how harmful the air quality is based on the pollutant amount.
4. Choosing the worst-case AQI
5. After converting both PM2.5 and PM10 to the EPA scale, we took the higher of the two as the final AQI.
6. This ensures that the AQI reflects the worst air quality, which is important for safety and public health.

Result: The converted AQI now ranges from 0 to 500, following the EPA standard. This scaled AQI can now be used reliably for analysis, modeling, and forecasting.

Handling Missing and Duplicate Values:

1. Missing timestamps or duplicate rows were removed.
2. Timestamps were converted to datetime objects and sorted chronologically.

Feature Engineering:

To enhance model performance, we derived both time-based and AQI-based features:

Time-based features: Hour, day, month, weekday, and whether it is a weekend.

AQI-derived features:

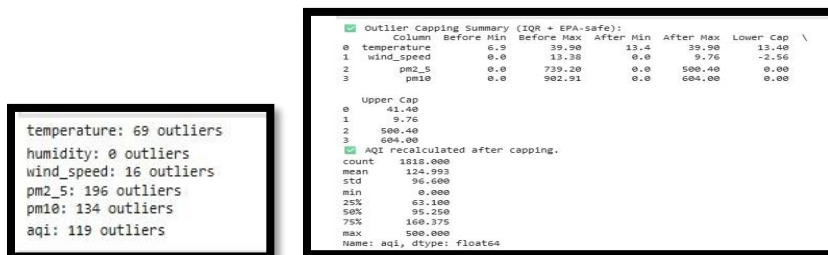
- aqi_change: Absolute change in AQI from previous measurement.
- aqi_change_rate: Relative change, capturing volatility in air quality

Correlation between AQI and features:



OUTLIER DETECTION:

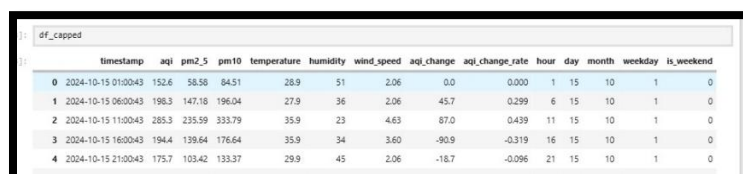
We handled outliers by checking which values were much lower or higher than most of the data, using a standard statistical method. Any extreme values that fell outside this typical range were capped to keep the dataset consistent and safe for modeling.



CAPPING Technique for Outliers:

We handled outliers by checking which values were much lower or higher than most of the data, using a standard statistical method. Any extreme values that fell outside this typical range were capped to keep the dataset consistent and safe for modeling.

Now we have processed dataframe `df_capped` that we will use for training



Feature Storage:

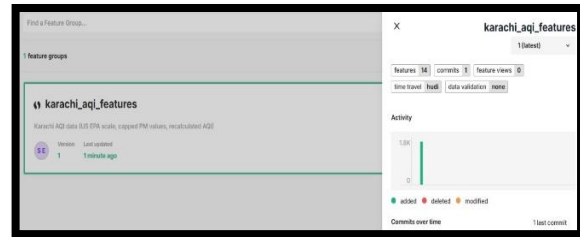
Cleaned and processed features are stored in Hopworks Feature Store for reproducibility and easy access.

```

karachi_fg.insert(df_capped)
print("New processed data inserted successfully into Feature Store.")

Uploading Dataframe: 100.00% [██████████] Rows 1818/1818 | Elapsed Time: 6
New processed data inserted successfully into Feature Store.
UserWarning: Materialization job is already running, aborting new execution

```



Model Training:

The dataset is split into 80% training and 20% testing for model evaluation.

```

Data split complete:
Training samples: 1454
Testing samples: 364

```

We experimented with both parametric and non-parametric models, as well as a neural network:

- *Non-parametric:* **Random Forest**, Extra Trees, KNN, XGBoost, SVR
- *Parametric:* Linear Regression, Ridge, Lasso, Elastic Net
- *Neural Network:* ANN

How we selected the best model out of all the algorithm's best models? To ensure models perform well and generalize, we defined a combined score:

$$\text{Combined Score} = \text{Test RMSE} + \lambda \times |\text{Train RMSE} - \text{Test RMSE}|$$

. λ weighs overfitting risk. The metric rewards low test RMSE while penalizing models that overfit training data.

Random Forest as Best Model:

Multiple hyperparameter configurations were tested. Random Forest achieved the lowest combined score and was selected as the final model. The model is saved as **best_random_forest.pkl**.*

Prediction and Deployment:

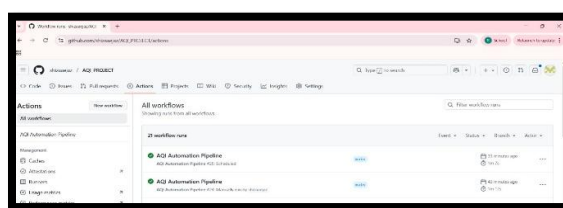
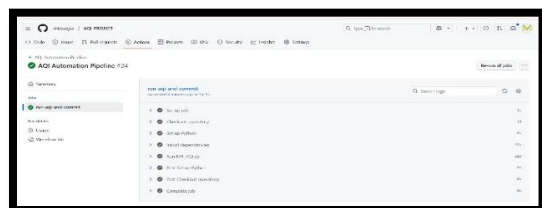
Using the best model, AQI predictions are generated:

```
forecast_data = get_forecast(df_capped, model_path="best_random_forest.pkl", save_csv=True)
```

Predictions are saved to **latest_predictions.csv** and integrated into a **Streamlit frontend**.



CI/CD via **GitHub Actions** automatically fetches new data and **updates predictions every 5 hours**.





SkyCast

Air Quality Index Prediction for Karachi



Karachi, Pakistan

03:10 PM

437

Hazardous

HEALTH RECOMMENDATIONS

- Avoid outdoor activities
- Stay indoors with air purifier
- Wear N95 if absolutely necessary

Current Weather

TEMPERATURE

23°C

Partly Cloudy

HUMIDITY

44%

Moderate

WIND SPEED

2 km/h

Light Breeze

3-Day AQI Forecast

Today

Nov 09, 2025

MIN

432

MAX

442

Tomorrow

Nov 10, 2025

MIN

432

MAX

442

Day After

Nov 11, 2025

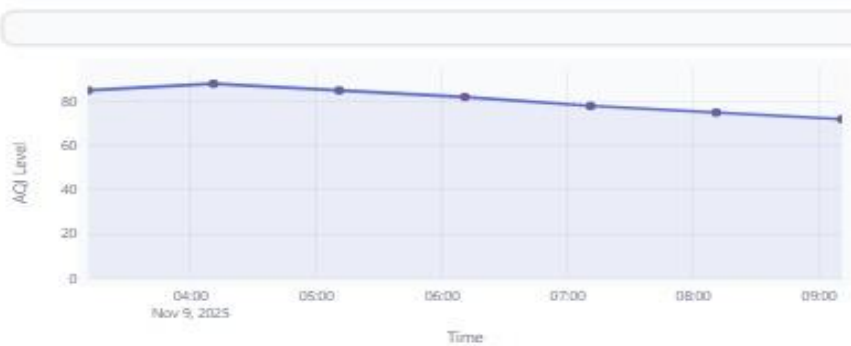
MIN

432

MAX

442

AQI Trend Over Time



SkyCast - Real-time Air Quality Predictions

Last updated: 2025-11-09 15:11:00