

# modeling

July 20, 2025

## 1 Model Building: Air Quality Forecasting

This notebook builds and evaluates machine learning models to forecast pollutant concentrations.

```
[1]: # Imports
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.ensemble import RandomForestRegressor
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error, r2_score, mean_absolute_error
import joblib
%matplotlib inline
```

```
[2]: # Load the processed data
df = pd.read_csv('E:/air_quality_forecasting/data/processed_air_quality.csv',
    parse_dates=['datetime'])
df.set_index('datetime', inplace=True)
df.head()
```

```
[2]:
```

	CO(GT)	PT08.S1(CO)	NMHC(GT)	C6H6(GT)	PT08.S2(NMHC)	\
datetime						
2004-03-10 19:00:00	2.0	1292.0	112.0	9.4	955.0	
2004-03-10 20:00:00	2.2	1402.0	88.0	9.0	939.0	
2004-03-10 21:00:00	2.2	1376.0	80.0	9.2	948.0	
2004-03-10 22:00:00	1.6	1272.0	51.0	6.5	836.0	
2004-03-10 23:00:00	1.2	1197.0	38.0	4.7	750.0	

	NOx(GT)	PT08.S3(NOx)	NO2(GT)	PT08.S4(NO2)	\
datetime					
2004-03-10 19:00:00	103.0	1174.0	92.0	1559.0	
2004-03-10 20:00:00	131.0	1140.0	114.0	1555.0	
2004-03-10 21:00:00	172.0	1092.0	122.0	1584.0	
2004-03-10 22:00:00	131.0	1205.0	116.0	1490.0	
2004-03-10 23:00:00	89.0	1337.0	96.0	1393.0	

	PT08.S5(O3)	...	weekday	month	CO(GT)_rolling3	\
--	-------------	-----	---------	-------	-----------------	---

datetime	...	...	...	...
2004-03-10 19:00:00	972.0	...	2.0	3.0
2004-03-10 20:00:00	1074.0	...	2.0	3.0
2004-03-10 21:00:00	1203.0	...	2.0	3.0
2004-03-10 22:00:00	1110.0	...	2.0	3.0
2004-03-10 23:00:00	949.0	...	2.0	3.0

  

	NOx(GT)_rolling3	NO2(GT)_rolling3	C6H6(GT)_rolling3	\
datetime				
2004-03-10 19:00:00	134.500000	102.500000	10.650000	
2004-03-10 20:00:00	133.333333	106.333333	10.100000	
2004-03-10 21:00:00	135.333333	109.333333	9.200000	
2004-03-10 22:00:00	144.666667	117.333333	8.233333	
2004-03-10 23:00:00	130.666667	111.333333	6.800000	

  

	CO(GT)_lag1	NOx(GT)_lag1	NO2(GT)_lag1	C6H6(GT)_lag1
datetime				
2004-03-10 19:00:00	2.6	166.0	113.0	11.9
2004-03-10 20:00:00	2.0	103.0	92.0	9.4
2004-03-10 21:00:00	2.2	131.0	114.0	9.0
2004-03-10 22:00:00	2.2	172.0	122.0	9.2
2004-03-10 23:00:00	1.6	131.0	116.0	6.5

[5 rows x 25 columns]

## 1.1 Define Target and Features

```
[3]: # Forecasting CO(GT) as the target variable
target = 'CO(GT)'
features = df.drop(columns=[target]).select_dtypes(include=[np.number]).columns.
        tolist()

X = df[features]
y = df[target]

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
        shuffle=False)
```

## 1.2 Train Model

```
[4]: model = RandomForestRegressor(n_estimators=100, random_state=42)
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
```

### 1.3 Evaluate Model

```
[5]: print(f"MAE: {mean_absolute_error(y_test, y_pred):.3f}")  
      print(f"RMSE: {mean_squared_error(y_test, y_pred) ** 0.5:.3f}")  
      print(f"R²: {r2_score(y_test, y_pred):.3f}")
```

MAE: 0.254

RMSE: 0.333

R²: 0.943

### 1.4 Save Model

```
[6]: import os  
      os.makedirs('E:/air_quality_forecasting/models', exist_ok=True)  
      joblib.dump(model, 'E:/air_quality_forecasting/models/co_forecast_model.pkl')  
      print("Model saved successfully ")
```

Model saved successfully