

# **Air Pollutant Concentration Prediction over Ahmedabad Using Machine Learning**

CSE523 - Machine Learning  
Winter Semester 2023  
Weekly Report - 15/4/2023

Yash Dahima - AU2129002

The task of model development using ensemble-based models has been performed this week.

## **Ensemble Models:**

Four models extra-trees, random forest, gradient boosting, and histogram-based gradient boosting were fit on the dataset using scikit-learn. They were run with the same random state and 10,000 number of estimators for the inter-comparison. Extra-trees performed best in this run. The models were again run with different parameters and gradient boosting performed best this time. There wasn't a big difference in the performances of these models.

```
1  # -*- coding: utf-8 -*-
2  """
3
4  Ensemble Models
5
6  """
7
8  import pandas as pd, numpy as np, time
9  from sklearn.model_selection import train_test_split
10 from sklearn.preprocessing import StandardScaler
11 from sklearn.ensemble import ExtraTreesRegressor, GradientBoostingRegressor, HistGradientBoostingRegressor,
12 RandomForestRegressor
13
14 from sklearn.metrics import explained_variance_score, mean_absolute_error, mean_squared_error, r2_score,
15 mean_absolute_percentage_error
16
17
18 # Load the dataset
19 df = pd.read_excel('C:/Users/Yash Dahima/PhD/Course Work/ML/Project/AQI/Datasets/data4.xlsx')
20 df['datetime'] = pd.to_datetime(df['datetime'])
21 df = df.set_index('datetime')
22
23 # Separate the target variable from the features
24 X = df.drop('pm2p5', axis=1)
25 y = df['pm2p5']
26
27 # Split the data into training and testing sets
28 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, shuffle=False)
29
30
31 # Create a list of models to evaluate
32 models = [ExtraTreesRegressor(n_estimators=10000, random_state=0, n_jobs=-1, criterion='friedman_mse',
33 min_samples_split=20, min_samples_leaf=5, max_features=0.3),
34           RandomForestRegressor(n_estimators=10000, random_state=0, n_jobs=-1, criterion='friedman_mse',
35 min_samples_split=20, min_samples_leaf=5, max_features=0.3),
36           GradientBoostingRegressor(n_estimators=10000, random_state=0, loss='huber', learning_rate=0.01,
37 subsample=0.9, criterion='friedman_mse', min_samples_split=20, min_samples_leaf=5, max_depth=None,
38 max_features=0.3),
39           HistGradientBoostingRegressor(max_iter=10000, random_state=0, max_leaf_nodes=None,
40 learning_rate=0.01, min_samples_leaf=5)
41 ]
```

```

37 # Create an empty dataframe to store the evaluation metrics
38 metrics_df = pd.DataFrame(columns=["Model", "evs", "mae", "mape", "rmse", "r^2 Score"])
39
40
41 # Evaluate each model using a for loop
42 for model in models:
43
44     start_time = time.time()
45
46     # Train the model
47     model.fit(X_train, y_train)
48
49     print(f"Time for fitting: {(time.time() - start_time)/60:.3f} minutes")
50
51     # Evaluate the model on the testing set
52     y_pred = model.predict(X_test)
53
54     # Compute evaluation metrics
55     evs = explained_variance_score(y_test, y_pred)
56     mae = mean_absolute_error(y_test, y_pred)
57     mape = mean_absolute_percentage_error(y_test, y_pred)*100
58     rmse = np.sqrt(mean_squared_error(y_test, y_pred))
59     r2 = r2_score(y_test, y_pred)
60
61     # Append evaluation metrics to the dataframe
62     metrics_df = pd.concat([metrics_df, pd.DataFrame({"Model": [type(model).__name__],
63                                                       "evs": [evs],
64                                                       "mae": [mae],
65                                                       "mape": [mape],
66                                                       "rmse": [rmse],
67                                                       "r^2 Score": [r2]})]),
68                               ignore_index=True)
69

```

Important model performance evaluation metrics are shown in the table below:

Models	Extra-trees	Random Forest	Gradient Boosting	Hist-Gradient Boosting
Mean Absolute Error	13.64	13.88	14.85	14.04
Mean Absolute Percentage Error	22.68	22.91	24.02	22.83
Root Mean Squared Error	18.03	18.33	19.71	18.65
R2 Score	0.82	0.81	0.78	0.80

As we can see, the model performances are far better than the previous one, they seem to be able to capture the pattern in the data.