

Air Pollutant Concentration Prediction over Ahmedabad Using Machine Learning

CSE523 - Machine Learning
Winter Semester 2023
Weekly Report - 8/4/2023

Yash Dahima - AU2129002

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

Kernel Ridge Model:

Ridge Regression model was developed using a periodic kernel with random search of the best parameters. 100 iterations were performed to choose the best combination of alpha, kernel length scale, and kernel periodicity. The model was fit with the dataset and evaluated using different metrics. Randomized Search provided the best values of alpha = 771.7, kernel length scale = 32.0, and kernel periodicity = 0.22 out of 100 iterations.

```
1  # -*- coding: utf-8 -*-
2  """
3
4  Kernel Ridge Regression
5
6  """
7
8  import pandas as pd, numpy as np
9
10 from sklearn.model_selection import train_test_split, RandomizedSearchCV
11 from sklearn.gaussian_process.kernels import ExpSineSquared
12 from sklearn.kernel_ridge import KernelRidge
13 from sklearn.utils.fixes import loguniform
14
15 from sklearn.metrics import explained_variance_score, mean_absolute_error, mean_squared_error,
16 r2_score, mean_absolute_percentage_error
17
18
19 # Load the dataset
20 df = pd.read_excel('C:/Users/Yash Dahima/PhD/Course Work/ML/Project/AQI/Datasets/data4.xlsx')
21 df['datetime'] = pd.to_datetime(df['datetime'])
22 df = df.set_index('datetime')
23
24 # Separate the target variable from the features
25 X = df.drop('pm2p5', axis=1)
26 y = df['pm2p5']
27
28 # Split the data into training and testing sets
29 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, shuffle=False)
30
31 kernel_ridge = KernelRidge(kernel=ExpSineSquared())
32
33 param_distributions = {
34     "alpha": loguniform(1e-2, 1e3),
35     "kernel__length_scale": loguniform(1e-2, 1e2),
36     "kernel__periodicity": loguniform(1e-2, 1e2),
37 }
```

```

38     kernel_ridge_tuned = RandomizedSearchCV(
39         kernel_ridge,
40         param_distributions=param_distributions,
41         n_iter=100,
42         random_state=0,
43     )
44
45
46     # Train the model
47     kernel_ridge_tuned.fit(X_train, y_train)
48
49     # Evaluate the model on the testing set
50     y_pred = kernel_ridge_tuned.predict(X_test)
51
52     # Compute evaluation metrics
53     #coeff = model.coef_
54     evs = explained_variance_score(y_test, y_pred)
55     mae = mean_absolute_error(y_test, y_pred)
56     mape = mean_absolute_percentage_error(y_test, y_pred)*100
57     rmse = np.sqrt(mean_squared_error(y_test, y_pred))
58     r2 = r2_score(y_test, y_pred)

```

Gaussian Process Regressor:

The Gaussian Process Regressor model was developed using a periodic kernel to address daily and seasonal periodicity, and quadratic kernel to address trend in the dataset.

```

1  # -*- coding: utf-8 -*-
2  """
3
4  Gaussian Process Regressor
5
6  """
7
8  import pandas as pd, numpy as np, time
9  from scipy.stats import pearsonr, spearmanr
10
11  from sklearn.model_selection import train_test_split
12  from sklearn.gaussian_process.kernels import ExpSineSquared, WhiteKernel, RBF, RationalQuadratic
13  from sklearn.gaussian_process import GaussianProcessRegressor
14
15  from sklearn.metrics import explained_variance_score, mean_absolute_error, mean_squared_error,
16  r2_score, mean_absolute_percentage_error
17
18
19  # Load the dataset
20  df = pd.read_excel('C:/Users/Yash Dahima/PhD/Course Work/ML/Project/AQI/Datasets/data4.xlsx')
21  df['datetime'] = pd.to_datetime(df['datetime'])
22  df = df.set_index('datetime')
23

```

```

24 # Separate the target variable from the features
25 X = df.drop('pm2p5', axis=1)
26 y = df['pm2p5']
27
28 # Split the data into training and testing sets
29 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, shuffle=False)
30
31 # Define the kernels
32 kernel_daily = ExpSineSquared(length_scale=1.0, periodicity=1.0)
33 kernel_seasonal = ExpSineSquared(length_scale=1.0, periodicity=365.0)
34 kernel_trend = RationalQuadratic(length_scale=1.0, alpha=1.0)
35
36 # Combine the kernels using addition and multiplication
37 kernel = RBF(length_scale=1.0) * (kernel_daily + kernel_seasonal) + kernel_trend
38
39 # Initialize the GaussianProcessRegressor model
40 gaussian_process = GaussianProcessRegressor(kernel=kernel)
41
42 # Fit the model to the training data
43 start_time = time.time()
44 gaussian_process.fit(X_train, y_train)
45 print(f"Time for GaussianProcessRegressor fitting: {time.time() - start_time:.3f} seconds")
46
47 # Evaluate the model on the testing set
48 y_pred = gaussian_process.predict(X_test)
49
50 # Compute evaluation metrics
51 #coeff = model.coef_
52 evs = explained_variance_score(y_test, y_pred)
53 mae = mean_absolute_error(y_test, y_pred)
54 mape = mean_absolute_percentage_error(y_test, y_pred)*100
55 rmse = np.sqrt(mean_squared_error(y_test, y_pred))
56 r2 = r2_score(y_test, y_pred)

```

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

Models	Kernel Ridge	Gaussian Process
Mean Absolute Error	40.4	32.4
Mean Absolute Percentage Error	47.4	44.5
Root Mean Squared Error	53.9	40.5
R2 Score	- 0.2	0.1

As we can see, the model performances are not satisfactory, they are still not able to capture the pattern in the data. Hence, the models such as Random Forest, SVR will be developed in the next week.