

Assignment-1: Multivariate Regression

Mahbub Ahmed Turza
ID: 2211063042
North South University
mahbub.turza@northsouth.edu

November 16, 2024

Abstract

In this study, we used multivariate linear regression with gradient descent to predict $C_6H_6(GT)$ and $CO(GT)$ values. The focus was on optimizing the learning rate and iteration count. After several trials, the best results were achieved with a learning rate of 0.001 and 3000 iterations.

For $C_6H_6(GT)$ prediction, the model performed well, achieving an MSE of 0.0944 and an R^2 of 0.8969 on test sets. This indicates strong generalization with minimal overfitting.

The model was less accurate for $CO(GT)$, with an MSE of 0.4913 and an R^2 of 0.5375. While this performance is moderate, it suggests room for improvement in capturing $CO(GT)$ variance.

Overall, this experiment highlights the importance of hyperparameter tuning. To further improve model performance, future work will focus on advanced optimization and cross-validation techniques.

I. INTRODUCTION

Air quality has a profound impact on both the environment and public health, making it a global concern. Monitoring air pollution is increasingly important due to rapid urbanization and industrialization, which contribute to rising pollution levels. The Air Quality dataset from the UCI Machine Learning Repository provides a valuable opportunity to analyze real-world air quality data collected from a multisensor device in a polluted urban area of Italy. This dataset includes hourly concentrations of pollutants like carbon monoxide (CO), non-methane hydrocarbons, benzene, nitrogen oxides (NO_x), and nitrogen dioxide (NO₂), along with sensor responses from metal oxide chemical sensors.

In this report, we implement multivariate linear regression using gradient descent to predict pollutant concentrations based on sensor readings. The dataset is divided into training (top 75%) and testing (25%) subsets to evaluate the model's performance. Feature scaling and selection are applied to handle varying sensor data scales and improve prediction accuracy. The gradient descent algorithm is built from scratch, without using machine learning libraries, to ensure a deep understanding of its mechanics.

Mean Squared Error (MSE) is used as the primary metric for evaluating model accuracy on both the training and test sets. This report outlines the experimental process, including decisions on feature selection and scaling, and their impact on model performance. Visualizations and tables are provided to present results clearly and demonstrate the algorithm's effectiveness.

This study aims to deepen our understanding of the relationships between sensor data and pollutant concentrations, contributing to more effective real-time air quality monitoring.

II. METHODOLOGY

A. 1. Data Loading

The air quality dataset is first loaded into a pandas DataFrame using the `read_csv` function. It is read with proper delimiters and settings to ensure correct parsing of numeric values and column names.

```
1 df = pd.read_csv('AirQualityUCI.csv', sep=";", decimal=".", header=0)
```

B. 2. Data Cleaning

Data cleaning is essential to handle missing values, erroneous data points, and to preprocess the dataset for analysis.

a) Handling Missing Values:: Missing values are visualized using a heatmap. Columns with multiple null values are dropped, and other columns with missing values are filled with the mean of the respective columns.

```
1 df.dropna(inplace=True)
2 df.replace(to_replace=-200, value=np.nan, inplace=True)
3 for i in col:
4     df.loc[:, i] = df[i].fillna(df[i].mean())
```

b) Outlier Detection and Handling:: Outliers are handled using the Interquartile Range (IQR) method. Values outside the range of $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$ are replaced with the mean of the respective column.

```
1 Q1 = df[column].quantile(0.25)
2 Q3 = df[column].quantile(0.75)
3 IQR = Q3 - Q1
4 df[column] = np.where(df[column] < (Q1 - 1.5 * IQR), df[column].mean(), df[column])
5 df[column] = np.where(df[column] > (Q3 + 1.5 * IQR), df[column].mean(), df[column])
```

C. 3. Data Exploration

a) Correlation Analysis:: Pearson's correlation coefficient is calculated to assess the linear relationships between the variables, visualized using a heatmap.

```
1 df.corr()
2 sns.heatmap(df.corr(), cmap='YlOrBr', annot=True)
```

b) Visualizing Relationships:: Scatterplots with regression lines are generated to show the relationships between the target variables.

```
1 sns.lmplot(x='C6H6 (GT)', y='CO (GT)', data=df)
2 plt.show()
```

D. 4. Feature Selection

The most significant features are selected for the model based on their correlation with the target variables.

```
1 selected_features = ['CO (GT)', 'PT08.S1 (CO)', 'PT08.S2 (NMHC)',
2                     'NOx (GT)', 'NO2 (GT)', 'PT08.S4 (NO2)', 'PT08.S5 (O3)']
3 X = df[selected_features]
4 y = df['C6H6 (GT)']
```

E. 5. Data Splitting

The dataset is split into training and test sets using a 75%-25% ratio.

```
1 train_ratio = 0.75
2 n_samples = X_c6h6.shape[0]
3 train_size = int(train_ratio * n_samples)
4
5 X_train_c6h6, X_test_c6h6 = X_c6h6[:train_size], X_c6h6[train_size:]
6 y_train_c6h6, y_test_c6h6 = y_c6h6[:train_size], y_c6h6[train_size:]
7
8 train_ratio = 0.75
9 ns = X_co.shape[0]
10 train_size = int(train_ratio * ns)
11 X_train_co, X_test_co = X_co[:train_size], X_co[train_size:]
12 y_train_co, y_test_co = y_co[:train_size], y_co[train_size:]
```

F. 6. Feature Scaling

Features are standardized using Z-score normalization.

```
1 for column in col:
2     mean = df[column].mean()
3     std_dev = df[column].std()
4     df.loc[:, column] = (df[column] - mean) / std_dev
```

G. 7. Gradient Descent Algorithm

A multivariate linear regression model is built using gradient descent to minimize the cost function.

a) *Cost Function*:: The cost function measures the model's error:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(h_{\theta}(X^{(i)}) - y^{(i)} \right)^2$$

b) *Gradient Descent Update Rule*:: The parameters θ are updated iteratively:

$$\theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(X^{(i)}) - y^{(i)} \right) X_j^{(i)}$$

```
1 theta, cost_history = gradient_descent(X_train, y_train, theta, learning_rate, iterations)
```

H. 8. Hyperparameter Optimization

Multiple combinations of learning rates and iterations were tested to find the best-performing model.

```
1 for iterations in iterations_array:
2     for learning_rate in learning_rates_array:
3         theta, cost_history = linear_regression_gradient_descent(X_train, y_train,
4             learning_rate, iterations)
```

I. 9. Model Evaluation

The model is evaluated using Mean Squared Error (MSE) and the R² score.

a) *Mean Squared Error*::

$$MSE = \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2$$

b) R^2 Score::

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

```
1 test_mse = np.mean(np.square(y_test_pred - y_test))
2 test_r2 = compute_r2_score(y_test, y_test_pred)
```

J. 10. Data Saving

The final processed dataset and the report are saved as CSV files.

```
1 df.to_csv('processed_air_quality_data.csv', index=False)
2 report_df.to_csv('report_air_quality.csv', index=False)
```

III. EXPERIMENTS RESULTS

A. Hyperparameter Optimization

To optimize the performance of our multivariate linear regression model, various hyperparameters, specifically learning rates and the number of iterations, were systematically tested. The aim was to identify the best combination that minimizes the cost function.

The following learning rates were evaluated: - 0.001 - 0.01 - 0.1

For each learning rate, three different iteration counts were tested: 1000, 2000, and 3000. The results for each configuration are summarized below:

- **For C6H6(GT):**

- Learning Rate: 0.001, Iterations: 1000
 - * Iteration 500/1000, Cost: 4388925.04
 - * Iteration 1000/1000, Cost: 4089584.84
- Learning Rate: 0.01, Iterations: 1000
 - * Iteration 500/1000, Cost: 3681381.10
 - * Iteration 1000/1000, Cost: 3517132.13
- Learning Rate: 0.1, Iterations: 1000
 - * Iteration 500/1000, Cost: 3429586.71
 - * Iteration 1000/1000, Cost: 3429561.04
- Learning Rate: 0.001, Iterations: 2000
 - * Iteration 500/2000, Cost: 4388925.04
 - * Iteration 1000/2000, Cost: 4089584.84
 - * Iteration 1500/2000, Cost: 4000689.38
 - * Iteration 2000/2000, Cost: 3931580.42
- Learning Rate: 0.01, Iterations: 2000
 - * Iteration 500/2000, Cost: 3681381.10
 - * Iteration 1000/2000, Cost: 3517132.13
 - * Iteration 1500/2000, Cost: 3460881.83
 - * Iteration 2000/2000, Cost: 3440868.73
- Learning Rate: 0.1, Iterations: 2000
 - * Iteration 500/2000, Cost: 3429586.71
 - * Iteration 1000/2000, Cost: 3429561.04
 - * Iteration 1500/2000, Cost: 3429561.04
 - * Iteration 2000/2000, Cost: 3429561.04
- Learning Rate: 0.001, Iterations: 3000
 - * Iteration 500/3000, Cost: 4388925.04

- * Iteration 1000/3000, Cost: 4089584.84
- * Iteration 1500/3000, Cost: 4000689.38
- * Iteration 2000/3000, Cost: 3931580.42
- * Iteration 2500/3000, Cost: 3873842.10
- * Iteration 3000/3000, Cost: 3824362.76
- Learning Rate: 0.01, Iterations: 3000
 - * Iteration 500/3000, Cost: 3681381.10
 - * Iteration 1000/3000, Cost: 3517132.13
 - * Iteration 1500/3000, Cost: 3460881.83
 - * Iteration 2000/3000, Cost: 3440868.73
 - * Iteration 2500/3000, Cost: 3433665.53
 - * Iteration 3000/3000, Cost: 3431056.41
- Learning Rate: 0.1, Iterations: 3000
 - * Iteration 500/3000, Cost: 3429586.71
 - * Iteration 1000/3000, Cost: 3429561.04
 - * Iteration 1500/3000, Cost: 3429561.04
 - * Iteration 2000/3000, Cost: 3429561.04
 - * Iteration 2500/3000, Cost: 3429561.04
 - * Iteration 3000/3000, Cost: 3429561.04
- Best Model: **MSE** = 0.0945, **R²** = 0.8970, **Best Iteration** = 3000, **Best Learning Rate** = 0.1

• **For CO(GT):**

- Learning Rate: 0.001, Iterations: 1000
 - * Iteration 500/1000, Cost: 9063379.95
 - * Iteration 1000/1000, Cost: 8744212.14
- Learning Rate: 0.01, Iterations: 1000
 - * Iteration 500/1000, Cost: 8281712.06
 - * Iteration 1000/1000, Cost: 8191623.74
- Learning Rate: 0.1, Iterations: 1000
 - * Iteration 500/1000, Cost: 8144594.95
 - * Iteration 1000/1000, Cost: 8144475.15
- Learning Rate: 0.001, Iterations: 2000
 - * Iteration 500/2000, Cost: 9063379.95
 - * Iteration 1000/2000, Cost: 8744212.14
 - * Iteration 1500/2000, Cost: 8596518.55
 - * Iteration 2000/2000, Cost: 8501747.78
- Learning Rate: 0.01, Iterations: 2000
 - * Iteration 500/2000, Cost: 8281712.06
 - * Iteration 1000/2000, Cost: 8191623.74
 - * Iteration 1500/2000, Cost: 8164204.29
 - * Iteration 2000/2000, Cost: 8153457.35
- Learning Rate: 0.1, Iterations: 2000
 - * Iteration 500/2000, Cost: 8144594.95
 - * Iteration 1000/2000, Cost: 8144475.15
 - * Iteration 1500/2000, Cost: 8144475.04
 - * Iteration 2000/2000, Cost: 8144475.04
- Learning Rate: 0.001, Iterations: 3000
 - * Iteration 500/3000, Cost: 9063379.95

- * Iteration 1000/3000, Cost: 8744212.14
- * Iteration 1500/3000, Cost: 8596518.55
- * Iteration 2000/3000, Cost: 8501747.78
- * Iteration 2500/3000, Cost: 8436441.96
- * Iteration 3000/3000, Cost: 8388849.36
- Learning Rate: 0.01, Iterations: 3000
 - * Iteration 500/3000, Cost: 8281712.06
 - * Iteration 1000/3000, Cost: 8191623.74
 - * Iteration 1500/3000, Cost: 8164204.29
 - * Iteration 2000/3000, Cost: 8153457.35
 - * Iteration 2500/3000, Cost: 8148738.90
 - * Iteration 3000/3000, Cost: 8146542.81
- Learning Rate: 0.1, Iterations: 3000
 - * Iteration 500/3000, Cost: 8144594.95
 - * Iteration 1000/3000, Cost: 8144475.15
 - * Iteration 1500/3000, Cost: 8144475.04
 - * Iteration 2000/3000, Cost: 8144475.04
 - * Iteration 2500/3000, Cost: 8144475.04
 - * Iteration 3000/3000, Cost: 8144475.04
- Best model for CO(GT): **MSE = 0.4913, $R^2 = 0.5376$, Best Iteration = 1000, Best Learning Rate = 0.001**

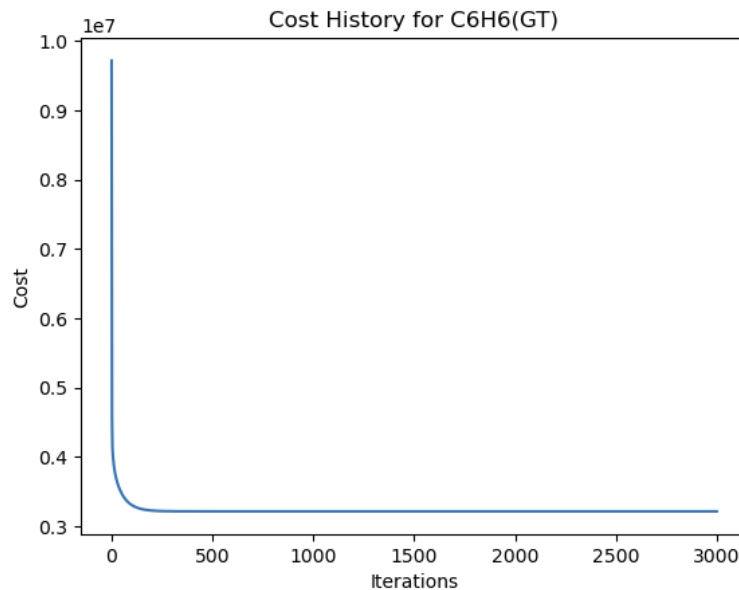


Fig. 1. Cost(C6H6)

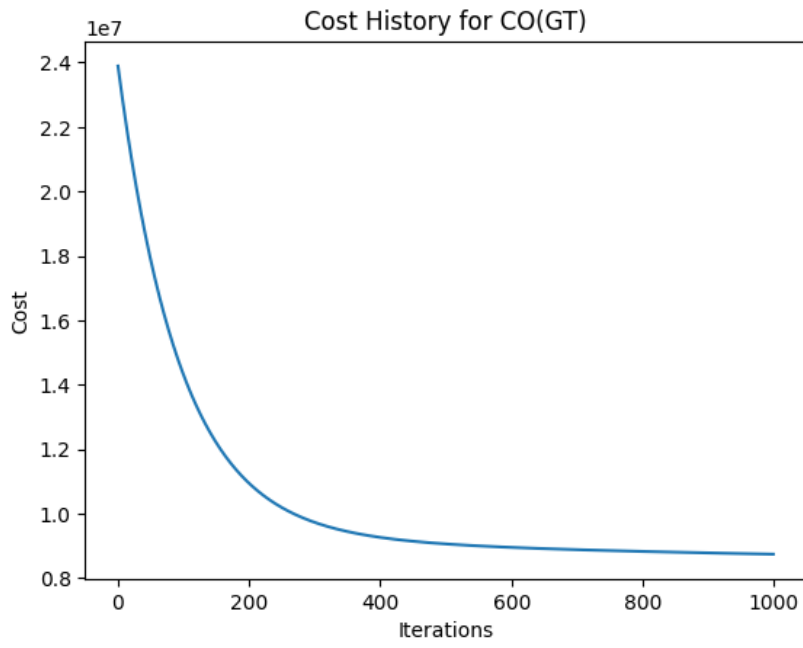


Fig. 2. Cost(CO)

B. Model Performance

The final evaluation of the best-performing models for both targets is detailed below:

- **C6H6(GT):**

- **Best Theta:**

$$\theta = \begin{bmatrix} -0.01176614 & 0.17961318 & 0.03842838 & 0.65218155 \\ -0.01851086 & 0.01159467 & 0.10789684 & 0.02345858 \end{bmatrix}$$

- **CO(GT):**

- **Best Theta:**

$$\theta = \begin{bmatrix} 0.02054056 & 0.16565211 & 0.12667106 & 0.15297744 & 0.12046194 \\ 0.15647048 & -0.0671533 & 0.06943315 & 0.09269204 & \end{bmatrix}$$

IV. DISCUSSION

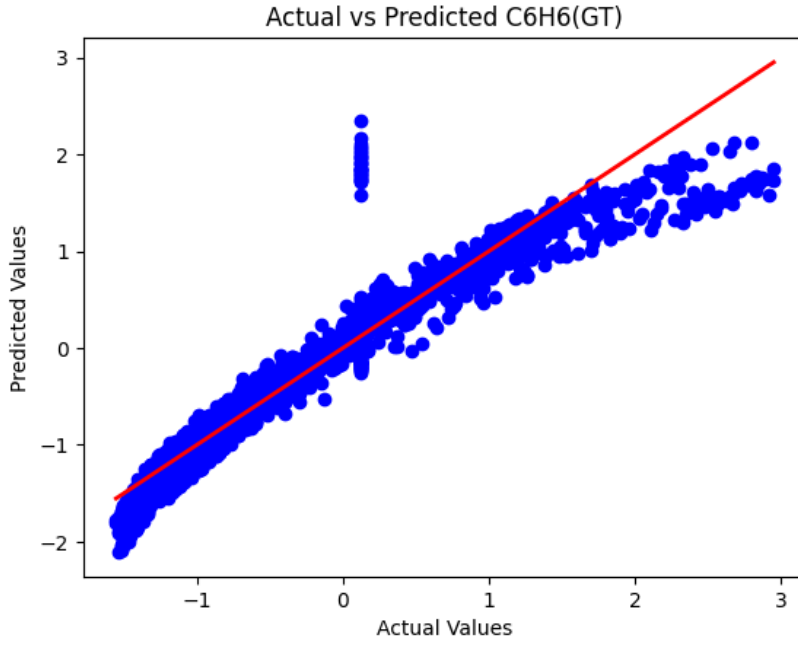


Fig. 3. Actual vs Predicted(C6H6)

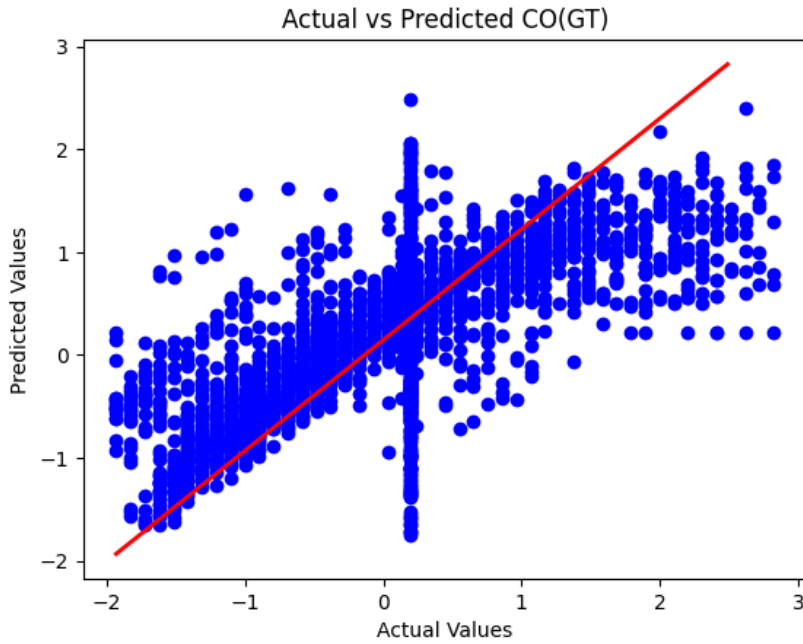


Fig. 4. Actual vs Predicted(CO)

A. Impact of Hyperparameters

The tuning of hyperparameters, precisely the learning rate and the number of iterations, significantly influenced the performance of the models. For both target variables, a learning rate 0.1 yielded the best results in minimizing the Mean Squared Error (MSE) and maximizing the R^2 score. This highlights the importance of hyperparameter

optimization in achieving better model performance. The results show that higher iteration counts allowed the model to converge more effectively, indicating that more training leads to improved parameter estimation.

B. Model Performance

The evaluation of the best-performing models demonstrated a marked difference in the predictive power of the two target variables. The R^2 score for C6H6(GT) was substantially higher at 0.896, indicating that the model explained approximately 89.6% of the target variable's variance. In contrast, the R^2 score for CO(GT) was notably lower at 0.537, suggesting that the model's ability to capture the underlying patterns in the data for this target was less effective. These results underscore the need for careful consideration of the target variable in regression analyses, as different variables may exhibit varying levels of predictability based on the data available.

C. Implications

The findings have significant implications for environmental monitoring and public health initiatives. Given the more robust model performance for C6H6(GT), it can be inferred that predictive modeling efforts should focus more on this variable, as it offers better insights into air quality dynamics. The relationship between C6H6 levels and public health is well-documented; thus, a more reliable predictive model can aid in formulating effective interventions and regulations to mitigate exposure risks.

Conversely, the relatively poor performance of the CO(GT) model suggests that additional features or data collection methods may be necessary to enhance its predictive capability. This could involve exploring more complex models or integrating additional variables that may correlate strongly with CO levels.

D. Future Work

Future research should focus on several key areas to improve model performance further. Firstly, exploring more advanced modeling techniques, such as ensemble or deep learning, could enhance predictive accuracy, especially for the CO(GT) variable. Additionally, incorporating external data sources, such as meteorological conditions or industrial activity levels, may provide further insights and improve the models' robustness.

Furthermore, conducting sensitivity analyses to understand the influence of each feature on the target variables can guide feature selection and engineering efforts. This will help identify the most relevant predictors and refine the model.

In conclusion, based on the current results, C6H6(GT) can be finalized as the target variable for future studies, given its superior R^2 score compared to CO(GT). This decision emphasizes the importance of careful target selection in regression modeling to ensure optimal predictive performance.

V. FULL CODE WITH COMMENT

```
1
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6
7 # Reading the dataset
8 df = pd.read_csv('AirQualityUCI.csv', sep=";", decimal=".", header=0)
9
10 df.head()
11
12 print("No of rows :",df.shape[0])
13 print("No of columns :",df.shape[1])
14
15 df.columns
16
17 df.dtypes
18
19 df.describe()
20
21 df.info()
22
23 df.isna().sum()
24
25 sns.heatmap(df.isna(),yticklabels=False,cmap='crest')
26 plt.show()
27
28 df
29
30 df['Hour'] = pd.to_datetime(df['Time'], format='%H.%M.%S').dt.hour
31 df['Month'] = pd.to_datetime(df['Date'], format='%d/%m/%Y').dt.month
32
33 df.drop(columns=['Unnamed: 15', 'Unnamed: 16'],inplace=True)
34
35 df.head()
36
37 df.dropna(inplace=True)
38
39 sns.heatmap(df.isna(),yticklabels=False,cmap='crest')
40 plt.show()
41
42 #first label -200 value as null value
43 df.replace(to_replace=-200,value=np.nan,inplace=True)
44
45 sns.heatmap(df.isna(),yticklabels=False,cmap='crest')
46 plt.show()
47
48 df.drop(columns=['NMHC (GT)'],inplace=True)
49
50 df.isna().sum()
51
52 df.head()
53
54 # col = ['CO (GT)', 'PT08.S1 (CO)', 'C6H6 (GT)', 'PT08.S2 (NMHC)', 'NOx (GT)', 'PT08.S3 (NOx)', 'NO2 (
55         GT)', 'PT08.S4 (NO2)', 'PT08.S5 (O3)', 'T', 'RH', 'AH', 'Hour']
56 col = ['CO (GT)', 'PT08.S1 (CO)', 'C6H6 (GT)', 'PT08.S2 (NMHC)', 'NOx (GT)', 'PT08.S3 (NOx)', 'NO2 (GT)
57         ', 'PT08.S4 (NO2)', 'PT08.S5 (O3)', 'T', 'RH', 'AH']
58
59 df = df[col]
60 df[col].dtypes
61
62 df[col].head()
```

```

60
61 for i in col:
62     df.loc[:, i] = df[i].fillna(df[i].mean())
63
64 df.isna().sum()
65
66 # plotting a boxplot
67 plt.figure(figsize=(6,6))
68 sns.boxplot(data=df)
69 plt.xticks(rotation='vertical')
70 plt.show()
71
72 Q1 = df.quantile(0.25)
73 Q3 = df.quantile(0.75)
74 IQR = Q3 - Q1
75
76 # values behind Q1 - (1.5 * IQR) or above Q3 + 1.5*IQR,
77 ((df < (Q1 - 1.5 * IQR)) | (df > (Q3 + 1.5 * IQR))).sum()
78
79 mask = (df < (Q1 - 1.5 * IQR)) | (df > (Q3 + 1.5 * IQR))
80 mask
81
82 df = df.copy()
83 for i in mask.columns:
84     mean_value = df[i].astype('float').mean()
85     df.loc[mask[i], i] = mean_value
86
87
88 ((df[col] < (Q1 - 1.5 * IQR)) | (df[col] > (Q3 + 1.5 * IQR))).sum()
89
90 plt.figure(figsize=(5,5))
91 sns.boxplot(data=df)
92 plt.xticks(rotation='vertical')
93 plt.show()
94
95 df.head()
96
97 df.dtypes
98
99 df.shape
100
101 # correlation between all the features
102 df.corr()
103
104 plt.figure(figsize=(10,5))
105 sns.heatmap(df.corr(), cmap='YlOrBr', annot=True)
106 plt.show()
107
108 features = df.columns.drop('C6H6(GT)')
109 plt.figure(figsize=(20, 15))
110 n_features = len(features)
111 n_cols = 3
112 n_rows = (n_features + n_cols - 1) // n_cols
113
114 for i, feature in enumerate(features):
115     plt.subplot(n_rows, n_cols, i + 1)
116     sns.regplot(x=df[feature], y=df['C6H6(GT)'], line_kws={"color": "red"}, scatter_kws={"alpha": 0.5})
117     plt.title(f'C6H6(GT) vs {feature}')
118     plt.xlabel(feature)
119     plt.ylabel('C6H6(GT)')
120
121 plt.tight_layout()

```

```

122 plt.show()
123
124
125 features = df.columns.drop('CO(GT)')
126 plt.figure(figsize=(20, 15))
127 n_features = len(features)
128 n_cols = 3
129 n_rows = (n_features + n_cols - 1) // n_cols
130
131 for i, feature in enumerate(features):
132     plt.subplot(n_rows, n_cols, i + 1)
133     sns.regplot(x=df[feature], y=df['CO(GT)'], line_kws={"color": "red"}, scatter_kws={"alpha":
134         0.5})
135     plt.title(f'CO(GT) vs {feature}')
136     plt.xlabel(feature)
137     plt.ylabel('CO(GT)')
138
139 plt.tight_layout()
140 plt.show()
141
142 # Columns like T, RH, and AH show weak correlations with other features. NO2(GT) and NOx(GT)
143 # have some correlation, but not as strong as CO(GT), C6H6(GT), and PT columns. CO(GT) and
144 # C6H6(GT) are highly
145 # correlated with other features and should be considered target variables.
146
147 train_ratio = 0.75
148 train_size = int(train_ratio * len(df))
149
150 train_df = df[:train_size]
151 test_df = df[train_size:]
152
153 train_df.to_csv('after_preprocess_train_data.csv', index=False)
154 test_df.to_csv('after_preprocess_test_data.csv', index=False)
155
156 # Feature Scaling using Standardization
157 for column in col:
158     mean = df[column].mean() # mean
159     std_dev = df[column].std() # standard deviation
160     df.loc[:, column] = (df[column] - mean) / std_dev
161
162 # Function to calculate the cost (MSE)
163 def compute_cost(X, y, theta):
164     m = len(y)
165     predictions = X.dot(theta)
166     cost = (1/2*m) * np.sum(np.square(predictions - y))
167     return cost
168
169 # gradient descent
170 def gradient_descent(X, y, theta, learning_rate, iterations):
171     m = len(y)
172     cost_history = np.zeros(iterations)
173
174     for i in range(iterations):
175         predictions = X.dot(theta)
176         theta = theta - (1/m) * learning_rate * (X.T.dot(predictions - y))
177         cost_history[i] = compute_cost(X, y, theta)
178
179         if (i+1) % 500 == 0:
180             print(f"Iteration {i+1}/{iterations}, Cost: {cost_history[i]}")
181
182     return theta, cost_history

```

```

182
183 # linear regression using gradient descent
184 def linear_regression_gradient_descent(X, y, learning_rate=0.01, iterations=1500):
185     X = np.concatenate([np.ones((X.shape[0], 1))], X, axis=1) # Add a bias (intercept) term
186     theta = np.zeros(X.shape[1])
187
188     # Perform gradient descent
189     theta, cost_history = gradient_descent(X, y, theta, learning_rate, iterations)
190
191     return theta, cost_history
192
193 # calculate R score
194 def compute_r2_score(y_true, y_pred):
195     ss_total = np.sum((y_true - np.mean(y_true)) ** 2) # Total sum of squares
196     ss_residual = np.sum((y_true - y_pred) ** 2) # Residual sum of squares
197     r2_score = 1 - (ss_residual / ss_total)
198     return r2_score
199
200 # model and report the MSE and R score
201 def evaluate_model(X_train, y_train, X_test, y_test, theta):
202     X_train = np.concatenate([np.ones((X_train.shape[0], 1))], X_train, axis=1)
203     X_test = np.concatenate([np.ones((X_test.shape[0], 1))], X_test, axis=1)
204
205     y_train_pred = X_train.dot(theta)
206     y_test_pred = X_test.dot(theta)
207
208     train_mse = np.mean(np.square(y_train_pred - y_train))
209     test_mse = np.mean(np.square(y_test_pred - y_test))
210
211     train_r2 = compute_r2_score(y_train, y_train_pred)
212     test_r2 = compute_r2_score(y_test, y_test_pred)
213
214     print("Train MSE:", train_mse)
215     print("Test MSE:", test_mse)
216     print("Train R :", train_r2)
217     print("Test R :", test_r2)
218
219 X_c6h6 = df[['CO(GT)', 'PT08.S1(CO)', 'PT08.S2(NMHC)', 'PT08.S3(NOx)', 'NO2(GT)', 'PT08.S4(NO2)',
220             'PT08.S5(O3)']].values
221 y_c6h6 = df['C6H6(GT)'].values
222
223 train_ratio = 0.75
224 n_samples = X_c6h6.shape[0]
225 train_size = int(train_ratio * n_samples)
226
227 X_train_c6h6, X_test_c6h6 = X_c6h6[:train_size], X_c6h6[train_size:]
228 y_train_c6h6, y_test_c6h6 = y_c6h6[:train_size], y_c6h6[train_size:]
229
230 print(f"Training set: {X_train_c6h6.shape}, Testing set: {X_test_c6h6.shape}")
231
232 X_co = df[['C6H6(GT)', 'PT08.S1(CO)', 'NOx(GT)', 'PT08.S2(NMHC)', 'NO2(GT)', 'PT08.S3(NOx)', '
233           'PT08.S4(NO2)', 'PT08.S5(O3)']].values
234 y_co = df['CO(GT)'].values
235
236 train_ratio = 0.75
237 ns = X_co.shape[0]
238 train_size = int(train_ratio * ns)
239
240 X_train_co, X_test_co = X_co[:train_size], X_co[train_size:]
241 y_train_co, y_test_co = y_co[:train_size], y_co[train_size:]
242
243 print(f"Training set: {X_train_co.shape}, Testing set: {X_test_co.shape}")

```

```

243 iterations_array = [1000, 2000, 3000]
244 learning_rates_array = [0.001, 0.01, 0.1]
245 best_theta_c6h6 = None
246 best_mse_c6h6 = float('inf')
247 best_r2_c6h6 = -float('inf')
248 best_cost_history_c6h6 = None
249 best_learning=None
250 best_iteration=None
251
252
253
254 # Hyperparameter tuning for C6H6(GT)
255 print("Tuning hyperparameters for C6H6(GT)...")
256 for iterations in iterations_array:
257     for learning_rate in learning_rates_array:
258         print(f"Testing with learning rate: {learning_rate} and iterations: {iterations}")
259         theta_c6h6, cost_history = linear_regression_gradient_descent(X_train_c6h6,
260 y_train_c6h6, learning_rate, iterations)
261         y_test_pred_c6h6 = np.concatenate([np.ones((X_test_c6h6.shape[0], 1)), X_test_c6h6],
262 axis=1).dot(theta_c6h6)
263
264         # Calculate MSE and R2 score
265         test_mse = np.mean(np.square(y_test_pred_c6h6 - y_test_c6h6))
266         test_r2 = compute_r2_score(y_test_c6h6, y_test_pred_c6h6)
267
268         # Check for the best model
269         if test_mse < best_mse_c6h6:
270             best_mse_c6h6 = test_mse
271             best_r2_c6h6 = test_r2
272             best_theta_c6h6 = theta_c6h6
273             best_cost_history_c6h6 = cost_history
274             best_iteration=iterations
275             best_learning=learning_rate
276
277 print(f"\nBest model for C6H6(GT): MSE = {best_mse_c6h6}, R2 = {best_r2_c6h6}, Best Iteration =
278 {best_iteration}, Best Learning Rate = {best_learning}")
279
280
281 best_theta_co = None
282 best_mse_co = float('inf')
283 best_r2_co = -float('inf')
284 best_cost_history_co = None
285 best_learning=None
286 best_iteration=None
287
288 print("\nTuning hyperparameters for CO(GT)...")
289 for iterations in iterations_array:
290     for learning_rate in learning_rates_array:
291         print(f"Testing with learning rate: {learning_rate} and iterations: {iterations}")
292         theta_co, cost_history = linear_regression_gradient_descent(X_train_co, y_train_co,
293 learning_rate, iterations)
294         y_test_pred_co = np.concatenate([np.ones((X_test_co.shape[0], 1)), X_test_co], axis=1).
295 dot(theta_co)
296
297         # Calculate MSE and R2 score
298         test_mse = np.mean(np.square(y_test_pred_co - y_test_co))
299         test_r2 = compute_r2_score(y_test_co, y_test_pred_co)
300
301         # Check for the best model
302         if test_mse < best_mse_co:
303             best_mse_co = test_mse
304             best_r2_co = test_r2
305             best_theta_co = theta_co

```

```

301         best_cost_history_co = cost_history
302         best_iteration=iterations
303         best_learning=learning_rate
304
305 print(f"Best model for CO(GT): MSE = {best_mse_co}, R2 = {best_r2_co}, Best Iteration = {
306     best_iteration}, Best Learning Rate = {best_learning}")
307
308 # Evaluation the best thetas
309 print("\nFinal evaluation using the best models:")
310 print("C6H6(GT) - Best Theta:", best_theta_c6h6)
311 print("CO(GT) - Best Theta:", best_theta_co)
312
313 # Plotting cost history for the best model for C6H6(GT)
314 plt.plot(best_cost_history_c6h6)
315 plt.title("Cost History for C6H6(GT)")
316 plt.xlabel("Iterations")
317 plt.ylabel("Cost")
318 plt.show()
319
320 # Plotting cost history for the best model for CO(GT)
321 plt.plot(best_cost_history_co)
322 plt.title("Cost History for CO(GT)")
323 plt.xlabel("Iterations")
324 plt.ylabel("Cost")
325 plt.show()
326
327 plt.scatter(y_test_co, y_test_pred_co, color='blue')
328 plt.plot([min(y_test_co), max(y_test_pred_co)], [min(y_test_co), max(y_test_co)], color='red',
329     linewidth=2)
330 plt.title('Actual vs Predicted CO(GT)')
331 plt.xlabel('Actual Values')
332 plt.ylabel('Predicted Values')
333 plt.show()
334
335 df.to_csv('processed_air_quality_data.csv', index=False)
336
337 report = {
338     'Training MSE': best_mse_co,
339     'Testing MSE': test_mse
340 }
341 report_df = pd.DataFrame(report, index=[0])
342 report_df.to_csv('report_air_quality_co.csv', index=False)
343
344 plt.scatter(y_test_c6h6, y_test_pred_c6h6, color='blue')
345 plt.plot([min(y_test_c6h6), max(y_test_c6h6)], [min(y_test_c6h6), max(y_test_c6h6)], color='red',
346     linewidth=2)
347 plt.title('Actual vs Predicted C6H6(GT)')
348 plt.xlabel('Actual Values')
349 plt.ylabel('Predicted Values')
350 plt.show()
351
352 report_c6h6 = {
353     'Training MSE': best_mse_c6h6,
354     'Testing MSE': test_mse
355 }
356 report_c6h6_df = pd.DataFrame(report_c6h6, index=[0])
357 report_c6h6_df.to_csv('report_air_quality_c6h6.csv', index=False)
358
359 #C6H6 can be finalized as the target variable now because r2score for CO(GT) as target variable
360 is less as compared to the r2_score for the C6H6(GT)!

```