

Assignment-1: Multivariate Linear Regression using Gradient Descent Algorithm

Sajan Kumer Sarker

Id: 2111131642

Email: sajan.sarker@northsouth.edu

Rafsan Jani Chowdhury

Id: 2011424642

Email: rafsan.chowdhury@northsouth.edu

Rosely Mohammad

Id: 2014219642

Email: rosely.mohammad@northsouth.edu

December 7, 2024

1. INTRODUCTION

In this assignment report, we will discuss about the implementation process of multivariate linear regression using the gradient descent algorithm. Air Quality dataset [1] was used to test the algorithm. It also involves data preprocessing, splitting it into training and testing datasets, feature selection, feature scaling, and finally, training the model using the implemented gradient descent algorithm. Then the model was tested on the testing dataset to evaluate its performance. The main objective is to predict air quality datasets parameters based on various other parameters. We train and test the model three times, each time focusing on predicting a different variable: first Temperature, then Relative Humidity (RH), and finally Absolute Humidity (AH).

2. METHODS

Necessary Imports:

```
1 import numpy as np
2 import pandas as pd
3 import math, copy
4 import matplotlib.pyplot as plt
```

We have imported all necessary libraries and modules to be used in each notebook.

Dataset Loading:

```
1 from ucimlrepo import fetch_ucirepo
2
3 # fetch dataset
4 air_quality = fetch_ucirepo(id=360)
5
6 # data (as pandas dataframes)
7 air_quality_data = air_quality.data.features
8
9 print(air_quality_data)
```

Import the dataset and print a few of the first and last samples from it.

2.1 Dataset Pre-processing Functions:

Calculate Missing Values:

```

1 # Missing values are tagged with -200 value.
2 def find_missing_values(data):
3     for column in data.columns:
4         count = 0
5         count = (data[column] == -200).sum()
6         print(f"{column:13}_has_{count:5}_occurrences.")

```

Defines a function that checks each column in a dataset for missing values tagged as -200 and prints the count of occurrences for each column.

Replace the Missing Value:

```

1 # replacing -200 with mean of each features.
2 def replace_with_mean(data):
3     for column in data.columns:
4         if pd.api.types.is_numeric_dtype(data[column]):
5             mean_value = data[data[column] != -200][column].mean()
6             data[column] = data[column].replace(-200, mean_value)
7             # print(f"Mean of {column:13} is: {mean_value:8.4f}")
8     return pd.DataFrame(data)

```

Defines a function that replaces missing values tagged as -200 in numeric columns of a dataset with the mean of that column, excluding the -200 values.

Remove Features:

```

1 # Feature remove function
2 def remove_feature(X, features_names):
3     return X.drop(columns=features_names)

```

Defines a function that removes specified features (columns) from the dataset X.

Feature Scaling (Training Data):

```

1 # Mean Normalization functions for training and test dataset
2 tr_mean = [] # global mean
3 tr_std = [] # global standard deviation
4
5 def zscore_normalize_train(X): # normalize train data only
6     m, n = X.shape
7     norm_X = []
8
9     for col in range(n):
10         col_values = X.iloc[:, col]
11         mean = sum(col_values) / m
12         tr_mean.append(mean)
13
14         sum_sq_diff = sum((val - mean)**2 for val in col_values)
15         std = (sum_sq_diff / m)**0.5
16         tr_std.append(std)
17
18         norm_col = [(val - mean) / std for val in col_values]
19         norm_X.append(norm_col)
20
21 norm_X = [list(row) for row in zip(*norm_X)]
22 norm_X = pd.DataFrame(norm_X, columns=X.columns)
23 return norm_X

```

We define a function that performs Z-score normalization on the training data by standardizing each feature (column) based on its mean and standard deviation. It returns the normalized dataset.

Feature Scaling (Test Data):

```

1 def zscore_normalize_test(X):      # normalize test data only
2     m, n = X.shape
3     norm_X = []
4
5     for col in range(n):
6         col_values = X.iloc[:, col]
7
8         norm_col = [(val - tr_mean[col]) / tr_std[col] for val in col_values]
9         norm_X.append(norm_col)
10
11 norm_X = [list(row) for row in zip(*norm_X)]
12 norm_X = pd.DataFrame(norm_X, columns=X.columns)
13 return norm_X

```

This will perform Z-score normalization on the test data using the mean and standard deviation calculated from the training data, and returns the normalized test dataset.

2.2 Data Pre-processing:

Pre-process the Data for Model:

```

1 # getting more info about the data
2 print(air_quality_data.info())
3
4 # Checking for missing Values
5 print(find_missing_values(air_quality_data))
6
7 # Delete the date and time columns, and those columns with more than 1,000 missing values
8 feature_name = ['Date', 'Time', 'CO (GT)', 'NMHC (GT)', 'NOx (GT)', 'NO2 (GT)']
9 air_quality_data = remove_feature(air_quality_data, feature_name)
10
11 # replace the missing values with mean of remaining columns
12 air_quality_data = replace_with_mean(air_quality_data)
13 print(air_quality_data.head())

```

Here we first print the general information about the dataset and check for missing values. It then removes the Date, Time, and columns with more than 1,000 missing values, followed by replacing the missing values in the remaining columns with the mean of each column. Finally, it displays the first few rows of the cleaned dataset.

Train Test Split:

```

1 # find index number of 75% data
2 train_data_index = int(len(air_quality_data) * 0.75)
3
4 train = air_quality_data[:train_data_index] # train data
5 test = air_quality_data[train_data_index:]  # test data
6
7 train_data = zscore_normalize_train(train)   # normalize train data
8 test_data = zscore_normalize_test(test)     # normalize test data
9
10
11 train_data.to_csv('./AirQualityTrainingData.csv', index=False)
12 test_data.to_csv('./AirQualityTestData.csv', index=False)
13

```

```

14 print(f"Shape_of_Whole_Dataset:{air_quality_data.shape}")
15 print(f"Shape_of_Training_Dataset:{train_data.shape}")
16 print(f"Shape_of_Test_Dataset:{test_data.shape}")
17
18 train_data = pd.read_csv('./AirQualityTrainingData.csv').to_numpy()
19 test_data = pd.read_csv('./AirQualityTestData.csv').to_numpy()

```

We Split the dataset into training and testing sets, with 75% of the data used for training. It then normalizes the training and test data using Z-score normalization. The normalized datasets are saved as CSV files, and their shapes are printed. Finally, the saved CSV files are read back into numpy arrays for further use.

Split Feature and Target and Plot the Data:

```

1 # Split features and target for train and test dataset
2 X_train = train_data[:, :8]
3 y_train = train_data[:, 8]
4 X_test = test_data[:, :8]
5 y_test = test_data[:, 8]
6
7 # visualize the train data
8 X_features = ['PT08.S1 (CO)', 'C6H6 (GT)', 'PT08.S2 (NMHC)', 'PT08.S3 (NOx)', 'PT08.S4 (NO2)', 'PT08.S5 (O3)', 'PT08.S6 (PM10)', 'PT08.S7 (PM2.5)', 'PT08.S8 (PM10.0)']
9
10 fig, ax = plt.subplots(1, 8, figsize=(30, 10), sharey=True)
11 fig.suptitle("Training_Data_Visualization", fontsize=16)
12
13 for i in range(0, 8):
14     ax[i].scatter(X_train[:, i], y_train, marker='.', c='b')
15     ax[i].set_xlabel(X_features[i])
16
17 ax[0].set_ylabel('AH')
18 plt.tight_layout()
19 plt.show()

```

Splits the training and test data into features and targets. It then visualizes the relationship between each feature and the target using scatter plots for each feature in the training dataset. The target variable AH labeled on the y-axis.

2.3 Linear Regression Algorithm Functions:

Cost Function:

```

1 # cost function
2 def compute_cost(X, y, w, b):
3     m = X.shape[0]
4     cost = 0.0
5
6     for i in range(m):
7         f_wb_i = np.dot(X[i], w) + b
8         cost = cost + (f_wb_i - y[i])**2
9
10    cost = cost / (2*m)
11    return cost

```

This function will calculate the cost (or loss) for a linear regression model. It computes the mean squared error between the predicted values \hat{y} and the actual target values y . The result is averaged over all data points and returned as the cost.

Gradient Calculation:

```

1 # compute the gradient for linear regression
2 def compute_gradient(X, y, w, b):

```

```

3  m, n = X.shape
4  dj_dw = np.zeros((n,))
5  dj_db = 0
6
7  for i in range(m):
8      err = (np.dot(X[i], w) + b) - y[i]
9
10     for j in range(n):
11         dj_dw[j] = dj_dw[j] + err * X[i, j]
12         dj_db = dj_db + err
13
14  dj_dw = dj_dw / m
15  dj_db = dj_db / m
16
17  return dj_db, dj_dw

```

It calculates the gradients of the cost function with respect to the weights and bias for linear regression. It computes the partial derivatives by iterating over the dataset and accumulating the errors. The gradients are then averaged over all data points.

Gradient Descent Algorithm:

```

1  # gradient descent algorithm
2  def gradient_descent(X, y, w_in, b_in, cost_function, gradient_function, alpha, num_iters):
3      J_history = [] # store cost J and w's at each iteration for plotting graph.
4      w = copy.deepcopy(w_in) # avoid modifying global w within function.
5      b = b_in
6
7      for i in range(num_iters):
8          dj_db, dj_dw = gradient_function(X, y, w, b)
9
10         w = w - alpha * dj_dw
11         b = b - alpha * dj_db
12
13         if i < 100000:
14             J_history.append(cost_function(X, y, w, b)) # save cost J at each iterations
15
16         if i % math.ceil(num_iters / 10) == 0:
17             print(f'Iteration_{i:4d}: cost {J_history[-1]:8.2f}')
18
19     return w, b, J_history # return final w, b and J history for plotting

```

It performs gradient descent to optimize the weights and bias for linear regression. It iteratively updates the parameters by subtracting the gradient of the cost function scaled by the learning rate. The cost is recorded at each iteration, and progress is printed. The function returns the final optimized parameters and the history of the cost function for plotting.

Calculate MSE:

```

1  # mse calculation
2  def mse(y_actual, y_pred):
3      error = 0.0
4
5      for actual, predicted in zip(y_actual, y_pred):
6          error += (actual - predicted) ** 2
7      mse = error / len(y_actual)
8      return mse

```

This function will calculate the MSE.

2.4 Running Linear Regression Gradient Descent Algorithm:

Finding Weight and Bias:

```

1 initial_w = np.zeros(8)
2 initial_b = 0
3 iterations = 1000
4 alpha = 3.0e-2
5
6 # run gradient descent
7 w_final, b_final, J_hist = gradient_descent(X_train, y_train, initial_w, initial_b, compute_co
8
9 print(f"b, w_found_by_gradient_descent: {b_final:0.2f}, {w_final}")

```

We initialize the weights as zeros, set the initial bias to 0, and specify the number of iterations (1000) and the learning rate (alpha). It then runs the gradient descent function on the training data to find the optimal weights and bias. Finally, it prints the learned bias and weights found by gradient descent.

Cost VS Iteration Plot:

```

1 plt.figure(figsize=(12, 4))
2 plt.plot(J_hist[:iterations])
3 plt.title("Cost_vs_Iteration")
4 plt.ylabel("Cost")
5 plt.xlabel("Iteration_Steps")
6 plt.show()

```

We plot the cost vs iteration figure to give a visual representation of the successful training.

Training Set MSE Calculation:

```

1 # train set predictions
2 y_train_predict = np.dot(X_train, w_final) + b_final
3
4 print(f"MSE_for_Training_Set: {mse(y_train, y_train_predict)}")

```

Calculate the MSE for the Training Set.

Plotting the Training vs Predicted Data:

```

1 # plotting actual training data vs predicted training data
2 X_features = ['PT08.S1 (CO)', 'C6H6 (GT)', 'PT08.S2 (NMHC)', 'PT08.S3 (NOx)', 'PT08.S4 (NO2)', 'PT08.S5 (
3
4 fig, ax = plt.subplots(1, 8, figsize=(30, 10), sharey=True)
5 fig.suptitle("Actual_Data_vs_Predicted_Data_(Training_Dataset)")
6
7 for i in range(8):
8     ax[i].scatter(X_train[:, i], y_train, label='Target')
9     ax[i].scatter(X_train[:, i], y_train_predict, color='FF9300', label="Prediction")
10    ax[i].set_xlabel(X_features[i])
11    ax[i].legend()
12
13 ax[0].set_ylabel('AH')
14 plt.show()

```

This code generates a set of subplots where each subplot compares the original target values with the predicted values for a specific feature in the training dataset.

Test Set MSE Calculation:

```

1 # test set predictions

```

```

2 y_test_predict = np.dot(X_test, w_final) + b_final
3
4 print(f"MSE_for_Training_Set:{mse(y_test, y_test_predict)}")

```

Calculate the MSE for the Test Set.

Plotting the Test vs Predicted Data:

```

1 # plotting actual test data vs predicted test data
2 X_features = ['PT08.S1 (CO)', 'C6H6 (GT)', 'PT08.S2 (NMHC)', 'PT08.S3 (NOx)', 'PT08.S4 (NO2)', 'PT08.S5 (CO2)', 'PT08.S6 (H2O)']
3
4 fig, ax = plt.subplots(1, 8, figsize=(30, 10), sharey=True)
5 fig.suptitle("Actual_Data_vs_Predicted_Data_(Test_Dataset)")
6
7 for i in range(8):
8     ax[i].scatter(X_test[:, i], y_test, label='Target')
9     ax[i].scatter(X_test[:, i], y_test_predict, color='FF9300', label="Prediction")
10    ax[i].set_xlabel(X_features[i])
11    ax[i].legend()
12
13 ax[0].set_ylabel('AH')
14 plt.show()

```

This code generates a set of subplots where each subplot compares the original target values with the predicted values for a specific feature in the training dataset.

3. EXPERIMENT RESULTS

3.1 Learning Rate:

Cost vs Iteration:

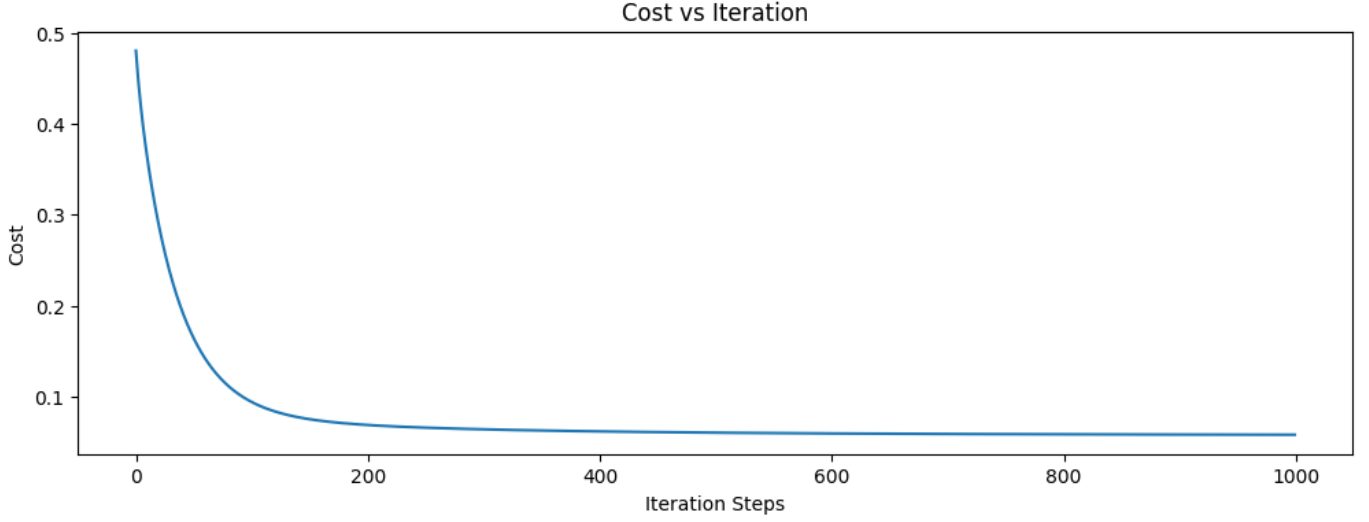


Fig. 1: Learning Rate for Alpha $3.0e-2$

The cost function was computed at each iteration during the training process with learning rate $3.0e-2$. The cost function represents the difference between the predicted values and the actual targets for the features. Here we see that the cost function decreases smoothly over the no. of iterations, gradually approaching the optimal solution and it levels out and does not decrease further. This indicates it that the optimization process is converging towards a minimum, suggesting that the algorithm has found a local or global minimum. Therefore this is our optimal learning rate

3.2 Analysis of Three Different Target Labels:

Training Data:



Fig. 2: Actual vs Predicted Data of Training Dataset for AH



Fig. 3: Actual vs Predicted Data of Training Dataset for RH



Fig. 4: Actual vs Predicted Data of Training Dataset for T

Mean Squared Error of Training Data:

| MSE for Training Dataset | | |
|--------------------------|-------|-------|
| AH | RH | T |
| 0.115 | 0.096 | 0.070 |

The table shows the Mean Squared Error (MSE) for the training dataset across three features: T, RH, and AH. The MSE values indicate that T has the lowest error (0.070), suggesting it is the most accurately predicted feature, while AH has the highest error (0.115), indicating it is the least accurately predicted feature among the three.

Test Data:

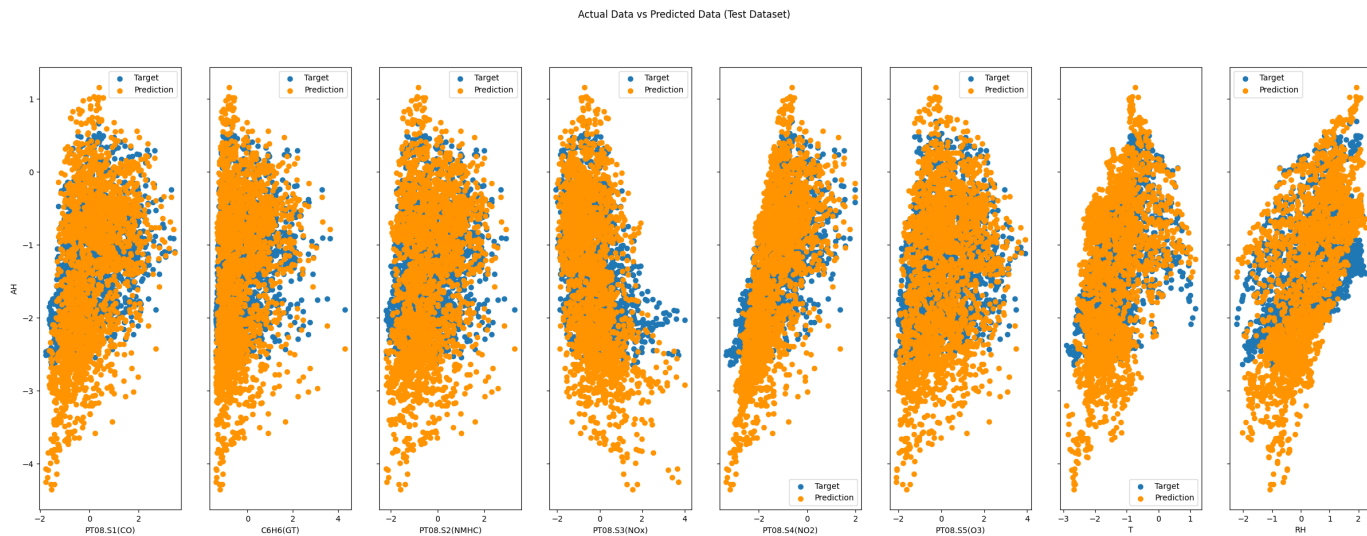


Fig. 5: Actual vs Predicted Data of Test Dataset for AH

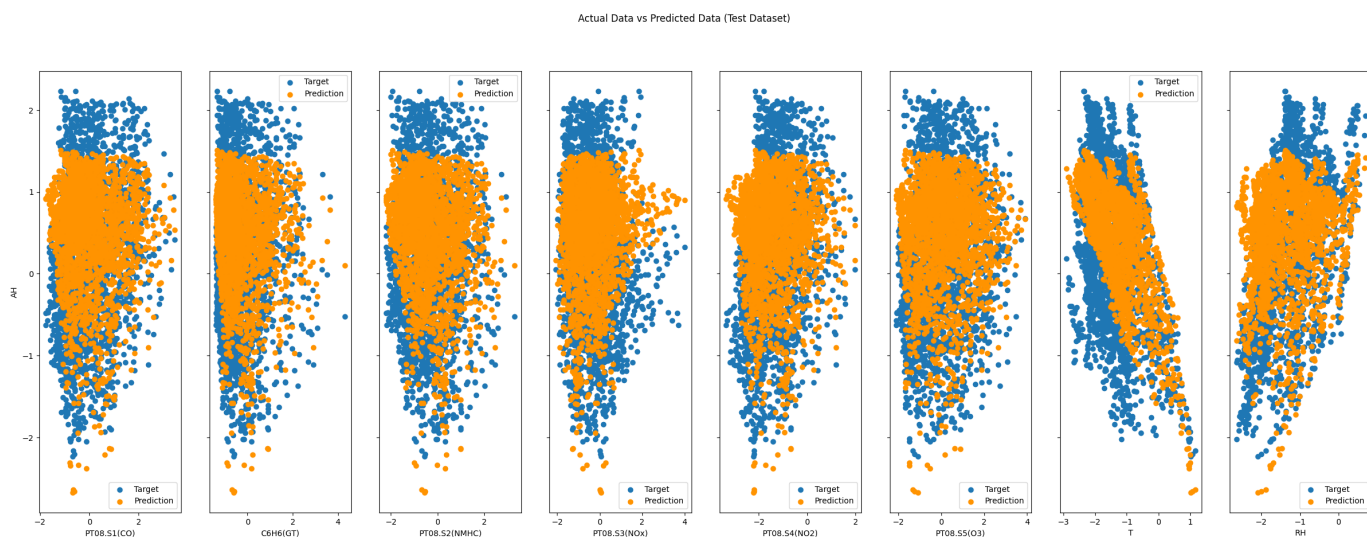


Fig. 6: Actual vs Predicted Data of Test Dataset for RH

Actual Data vs Predicted Data (Test Dataset)

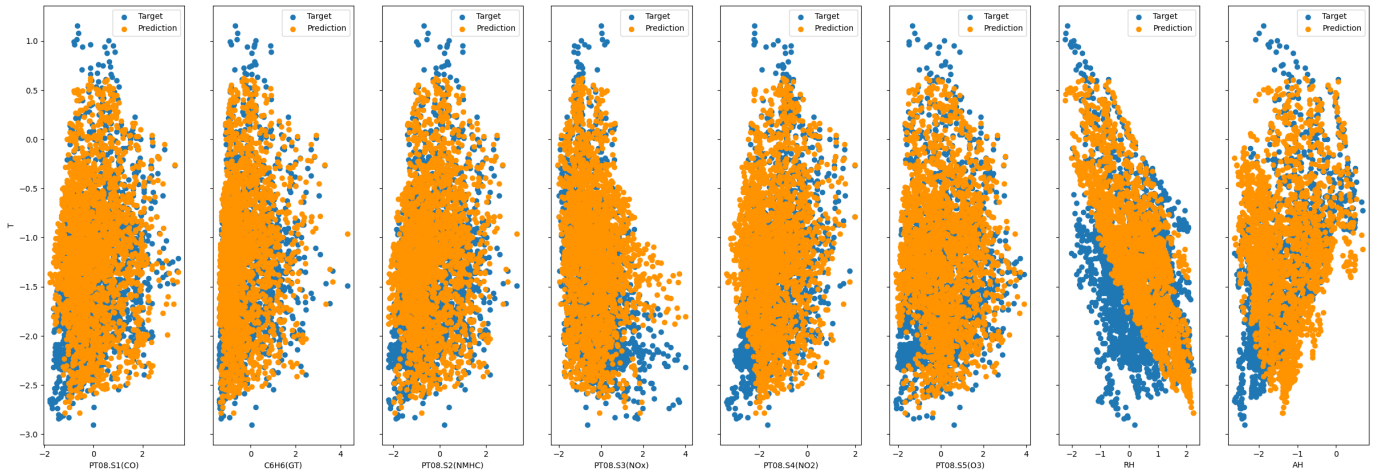


Fig. 7: Actual vs Predicted Data of Test Dataset for T

Mean Squared Error of Test Data:

| MSE for Test Dataset | | |
|----------------------|-------|-------|
| AH | RH | T |
| 0.190 | 0.299 | 0.201 |

The table shows the Mean Squared Error (MSE) for the training dataset across three features: T, RH, and AH. The MSE values indicate that AH has the lowest error (0.190), suggesting it is the most accurately predicted feature, while RH has the highest error (0.299), indicating it is the least accurately predicted feature among the three.

4. DISCUSSION

In this assignment we implemented the Multivariate Linear Regression using the Gradient Descent Algorithm and test it on the Air Quality Dataset [1]. We also learned to process dataset, apply feature engineering on dataset. We try to get the best performance for our model by experimenting with different labels (AH,RH,T) and optimize the training by choosing different learning rate and no of iterations in order to get the most accurate hypothesis equation. We choose the learning rate 3.0×10^{-2} first and it become the best learning rate so we don't need to test with other values. After training the model we check the accuracy by calculating MSE for both training and testing dataset. We found out that we got the best lowest MSE during testing process which is 0.096 while we take RH as the target.

REFERENCES

- [1] U. M. L. Repository. Air quality dataset. <https://archive.ics.uci.edu/dataset/360/air+quality>, 2008.