# Assignment-1: Multivariate Regression

Rezwan-Ul-Alam (ID: 2011659042)
Md. Nur Alam Jowel (ID: 2012355042)
Raian Ruku (ID: 2013409642)

*Email addresses:*

`rezwan.alam1@northsouth.edu`
`alam.jowel@northsouth.edu`
`raian.ruku@northsouth.edu`

## I. INTRODUCTION

This assignment aims to develop a multi-variable linear regression model for predicting relative humidity (RH) based on the air quality factors obtained from a gas multi-sensor device deployed in an Italian city. In this context, our analysis will involve creating custom functions for computing the cost and gradient and implementing the gradient descent algorithm to optimize the model parameters. Finally, feature scaling will be applied to maximize the test performance.

TABLE I
MULTIPLE FEATURES

| feature1 | feature2 | feature3 | feature4 | target |
|----------|----------|----------|----------|--------|
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| ... | ... | ... | ... | ... |

· Notation

$$n = \text{number of features}$$

$$m = \text{number of training examples}$$

**Univariate Linear Regression:** In univariate linear regression, only one target column (dependent variable) and one feature (independent variable). The hypothesis function is represented as:

$$h_\theta(x) = \theta_0 + \theta_1 x$$

Hence, both $\theta_0$ and $\theta_1$ are scalar values

**Multivariate Linear Regression:** In contrast, multivariate linear regression involves predicting a target variable based on multiple features. If there are $n$ features, the hypothesis function is extended to:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \ldots + \theta_n x_n$$

Here, $\theta_0$ is still a scalar representing the intercept, but $\theta_1, \theta_2, \ldots, \theta_n$ form a vector. This vector ($\boldsymbol{\theta}$) contains the coefficients associated with each feature.

To express it more formally:

$$h_\theta(x) = \boldsymbol{\theta}^T \mathbf{x}$$

Where:

- $\boldsymbol{\theta} = [\theta_0, \theta_1, \theta_2, \ldots, \theta_n]$ is the parameter vector.
- $\mathbf{x} = [1, x_1, x_2, \ldots, x_n]$ is the feature vector, with $x_0 = 1$ to account for the intercept term.

In summary, while univariate linear regression deals with a single feature and scalar coefficients, multivariate linear regression extends this to multiple features, with the coefficients organized into a vector.

## II. METHOD

### A. Data Exploration and Preprocessing

- We have decided to work on **Absolute Humidity** (AH), so we will drop Temperature (T) and Relative Humidity (RH)
- **Data Exploration and Preprocessing:** Before diving into model development, exploring and preprocessing the dataset is crucial.
  - **Handling Missing Values:** In our dataset, missing values are identified as -200. So we searched on all columns for -200 and replaced the value with the mean of the column

```python
import pandas as pd
import numpy as np
df = pd.read_excel('/content/gdrive/MyDrive/data/air.xlsx')
df.head()
for column in df.columns:
    total = df[column].isnull().sum()
    print(f"Column '{column}' has {total} null values")
    print("\n")
check = -200

#checking the presence of -200 in each of the columns
count = {}
for col in df.columns:
    count[col] = (df[col] == check).sum()

# Display the presence of -200 in each column
for column, count in count.items():
    print(f"Column '{column}' has {check}, {count} times")

#replace the missing value with NaN and later replace with the mean
df.replace(-200, np.nan, inplace=True)
df.isnull().sum()

for col in df.columns:
    mean = df[col].astype(float).mean()
    df[col].replace(np.nan, mean, inplace=True)
```

  - **Dropping Unnecessary Columns:** The date and time columns do not contribute significantly to the analysis, and our target value is not temperature (T) or relative humidity (RH). While checking for missing values, the column NMHC(GT) was observed to have over 90% missing values. Therefore, we will drop these columns.

```python
df = data.drop(['Date', 'Time','T','RH'], axis=1)

for column, count in count.items():
    print(f"Column '{column}' has {check}, {count} times") #total data: 9357
#Output : Column 'NMHC(GT)' has -200, 8443 times

df = df.drop('NMHC(GT)', axis=1)
```

- **Scaling Features:** For feature scaling, we have used our own **min_max_scaler(column)** method. The code below introduces a min_max_scaling function to perform min-max scaling on numeric columns in our data frame. It calculates the scaled values, transforming X_features to a common scale between 0 and 1.

```python
1  def min_max_scaling(column):
2      minimum_value = column.min()
3      maximum_value = column.max()
4      scaled_column = (column - minimum_value) / (maximum_value - minimum_value)
5      return scaled_column
6
7  #selecting featured columns from the dataframe
8  all_columns = df2.select_dtypes(include=[float, int]).columns
9  features_columns = all_columns[:-1]
10 print(features_columns)
11
12 #applying min_max_scaler() on all the features
13 df2[features_columns] = df2[features_columns].apply(min_max_scaling)
14 df2.head()
```

## B. Cost Function : compute_cost(X, y, w, b)

- We have worked on the given solution of the univariate linear regression code. There are some changes in the code. For univariate regression, there was only one feature. As a result, $\theta_1$ (in code $w$) was a scalar. However, in multivariate regression, $\theta_1$ becomes an n-dimensional array, where $n$ is the number of features. And $X$ (ndarray$(m, n)$) becomes a matrix as multiple features are added. Therefore, we can perform **np.dot()** for calculating the value of $h_\theta(x)$. Rest of the code remains the same.

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

## C. Gradient Descent : compute_gradient(X, y, w, b)

- As with the cost function, the code has some changes. Unlike univariate, in multivariate regression, dj_dw represents the cost gradient for the parameters w. It is expected to be a 1-dimensional NumPy array of length n. Here,
  1) Initialize dj_d$\theta$ ($\theta > 0$) as a zero vector of shape $(n,)$ and dj_db ($\theta_0$) as 0.
  2) Loop over each example in the training data ($m$ times).
  3) Compute the predicted value ($f_{wb}$) for each example using the current weights ($w$) and bias ($b$).
  4) Calculate the common term, the difference between the predicted value and the actual target value (common $= f_{wb} - y[i]$).
  5) Update the gradients (dj_dw and dj_db) using the common term and the input features.
  6) After looping through all examples, normalize the gradients by dividing them by the number of examples ($m$).
  7) Return the computed gradients.
  8) **Note:** dj_dw and dj_db are nothing but $\frac{\delta}{\delta\theta_j}J(\theta)$, where $J(\theta)$ is the cost function and $\theta_j$ represents the j-th parameter.

*Repeat*

$$\theta j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

*(simultaneously update $\theta_j$ for $j = 0, 1, \dots n$)*

*}*

```
1  def compute_gradient(X, y, w, b):
2
3      m = X.shape[0]
4      n = X.shape[1]
5      dj_dw = np.zeros((n,))
6      dj_db = 0.
7
8      for i in range(m):
9          f_wb = np.dot(X[i], w) + b
10         common = f_wb - y[i]
11         for j in range(n):
12             dj_dw[j] = dj_dw[j] + common * X[i, j]
13         dj_db = dj_db + common
14     dj_dw = dj_dw / m
15     dj_db = dj_db / m
16
17     return dj_db, dj_dw
```

## D. Gradient Function : gradient_descent()

1) Initialize $w$ and $b$ with the provided initial values.
2) Loop over a specified number of iterations.
3) Compute the gradient of the cost function with respect to the parameters using the provided `gradient_function`.
4) Update the parameters ($w$ and $b$) using the gradient and the learning rate.
5) Save the cost and parameter values at each iteration in the `J_history` and `p_history` lists, respectively.
6) Print the cost, gradients, and parameter values at regular intervals during training.
7) Return the final optimized parameters ($w$ and $b$), the history of cost values, and the history of parameter values.

```
1  def gradient_descent(X,y,w_in,b_in,alpha,num_iters,
2                       cost_function,gradient_function):
3      J_history = []
4      w = copy.deepcopy(w_in)
5      p_history = []
6      b = b_in
7
8      for i in range(num_iters):
9          dj_db, dj_dw = gradient_function(X, y, w , b)
10         w = w - alpha * dj_dw
11         b = b - alpha * dj_db
12         if i<100000:
13             J_history.append( cost_function(X, y, w , b))
14             p_history.append([w,b])
15         if i% math.ceil(num_iters/10) == 0:
16           print(f"Iteration {i:4}: Cost {J_history[-1]:8.2f} ",
17                   f"dj_dw: {dj_dw}, dj_db: {dj_db: 0.3e}  ",
18                   f"w: {w}, b:{b: 0.5e}")
19
20     return w, b, J_history,p_history
```

## E. *Train Test Split*

- As per the instruction, the first 75% of the data was used for training, and the last 25% of the data was used for testing. Training values are the features, and testing values are the target values. The code is given below:

```
1  first_75 = math.ceil(0.75*len(df1))
2  data_train = df1[:first_75, :] #first 75% of the data
3  data_test = df1[first_75:, :] #last 25% of the data
4
5  x_train = data_train[:, :-1]
6  y_train = data_train[:, -1]
7  x_test = data_test[:, :-1]   # First 9 columns
8  y_test = data_test[:, -1]    # Last column
9
10 print(x_train.shape)
11 print(x_test.shape)
```
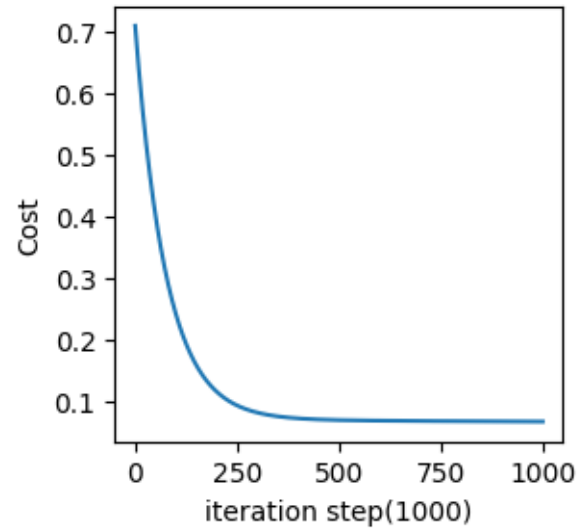
## III. EXPERIMENT RESULTS

## A. *Before Scaling*

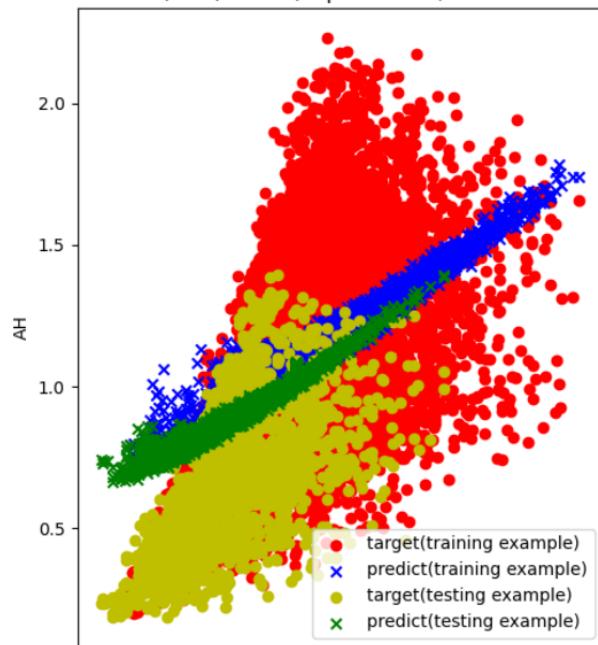- Before scling we tried to check out output for different iteration value and learning rate.

| Alpha | Iteration | Cost |
|---|---|---|
| | 100 | 0.15 |
| | 200 | 0.05 |
| 1.0e-7 | ⋮ | ⋮ |
| | 800 | 0.04 |
| | 900 | 0.04 |
| | 100 | 0.64 |
| | 200 | 0.07 |
| 1.0e-8 | ⋮ | ⋮ |
| | 800 | 0.06 |
| | 900 | 0.06 |
| | 300 | 0.64 |
| | 600 | 0.07 |
| 1.0e-8 (3000 Iterations) | ⋮ | ⋮ |
| | 2400 | 0.05 |
| | 2700 | 0.05 |
| | 100 | 0.71 |
| | 200 | 0.24 |
| 1.0e-9 | ⋮ | ⋮ |
| | 800 | 0.12 |
| | 900 | 0.07 |

| $\alpha$ | Iterations | Target Value | Predicted Value |
|---|---|---|---|
| | 1000 | 0.84 | 0.77 |
| 1.0e-7 | 1000 | 0.80 | 0.77 |
| | 1000 | 0.81 | 0.77 |
| | 1000 | 0.84 | 0.86 |
| 1.0e-8 | 1000 | 0.80 | 0.86 |
| | 1000 | 0.81 | 0.86 |
| | 3000 | 0.84 | 0.80 |
| 1.0e-8 | 3000 | 0.80 | 0.80 |
| | 3000 | 0.81 | 0.81 |
| | 1000 | 0.84 | 0.80 |
| 1.0e-9 | 1000 | 0.80 | 0.80 |
| | 1000 | 0.81 | 0.84 |

From the above tables, we can see that with different learning rates, the cost is reduced. However, when taking a relatively larger alpha value ($\alpha = 1.0e - 3$), the cost increases instead (Figure 1.1). Otherwise, in another case, the cost function decreases with iteration steps.

Cost vs. iteration(start) : alpha =0.001

Cost vs. iteration(start) : alpha =1e-09

Again, with the same alpha value, the number of iterations is affected. If we look closely at the second table, the prediction is more accurate with a larger iteration(3000).



PT08.S4(NO2) vs AH , alpha: 1e-08,iteration : 1000

PT08.S4(NO2) vs AH , alpha: 1e-08,iteration:3000

```
1  iterations = 1000 #3000
2  tmp_alpha = 1.0e-7 #1.0e-8,1.0e-9,1.0e-10
3
4  w_final, b_final, J_hist, p_hist
5          = gradient_descent(x_train ,y_train, w_init, b_init,tmp_alpha,
6               iterations,compute_cost, compute_gradient)
7
8  print(f"(w,b)␣found␣by␣gradient␣descent␣(alpha␣=␣␣␣␣␣␣␣␣{tmp_alpha:.1e}):
9  ({w_final},
10 {b_final:8.4f})")
11
12 for i in range(5):
13    print(f"alpha:␣{tmp_alpha:.1e},␣target␣value:␣{y_test[i]:0.2f},
14 ␣␣predicted␣value:␣{np.dot(x_test[i],␣w_final)␣+␣b_final:0.2f}")
15
16 # plot cost versus iteration
17 fig, (ax1) = plt.subplots(1, 1, constrained_layout=True, figsize=(3,3))
18 ax1.plot(J_hist[:iterations])
19 ax1.set_title(f"Cost␣vs.␣iteration(start)␣:␣alpha␣={tmp_alpha}");
20 ax1.set_ylabel('Cost')                 ;
21 ax1.set_xlabel(f'iteration␣step({iterations})')  ;
22 plt.show()
23
24
25 print(x_train.shape)
26 print(x_test.shape)
```
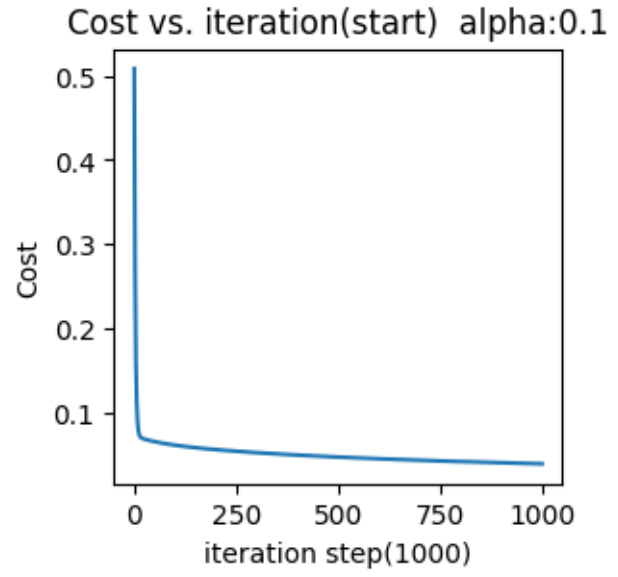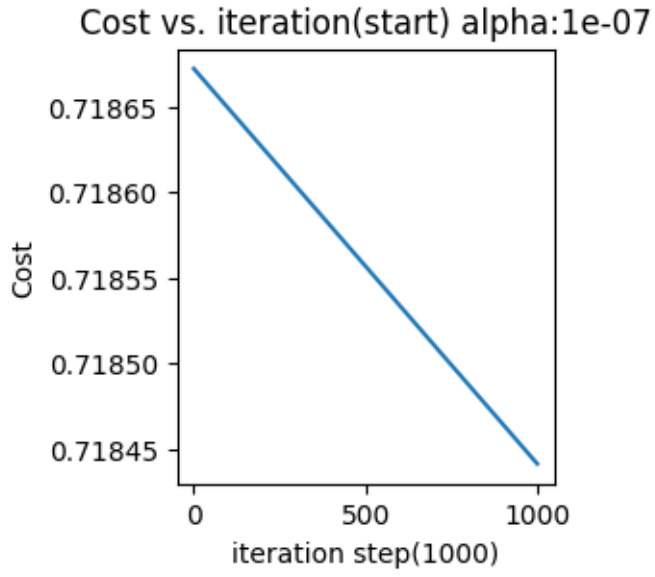
## B. *After Scaling*

- After scaling, we tried to check out the output for different iteration values and learning rates. This time, we were able to find the cost function for a relatively larger value of alpha. This time, for some larger alpha values ($\alpha = 1.0e-1, 1.0e-3$), we got a better prediction than the smaller value ($\alpha = 1.0e-10, 1.0e-7$).
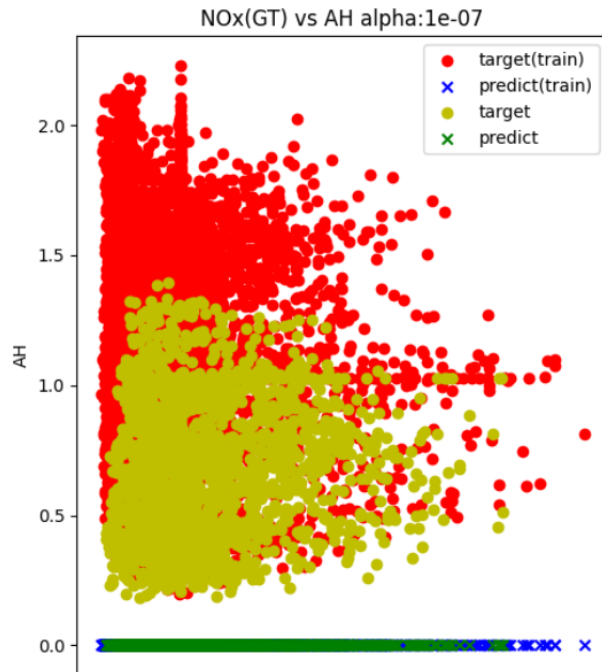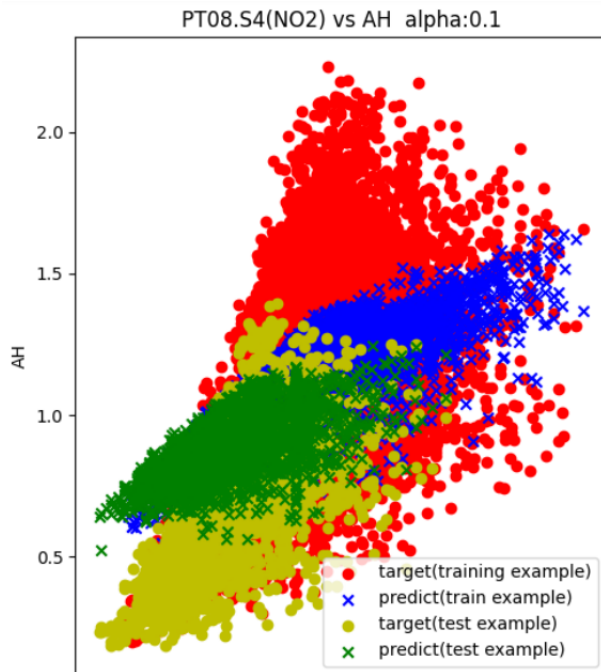
| Alpha | Iteration | Cost |
|-------|-----------|------|
|       | 100 | 0.72 |
|       | 200 | 0.52 |
| 1.0e-3 | ⋮ | ⋮ |
|       | 800 | 0.11 |
|       | 900 | 0.10 |
|       | 100 | 0.51 |
|       | 200 | 0.06 |
| 1.0e-1 | ⋮ | ⋮ |
|       | 800 | 0.04 |
|       | 900 | 0.04 |
|       | 100 | 0.72 |
|       | 200 | 0.72 |
| 1.0e-7 | ⋮ | ⋮ |
|       | 800 | 0.72 |
|       | 900 | 0.72 |
|       | 100 | 0.72 |
|       | 200 | 0.72 |
| 1.0e-10 | ⋮ | ⋮ |
|       | 800 | 0.72 |
|       | 900 | 0.72 |

| $\alpha$ | Iterations | Target Value | Predicted Value |
|----------|------------|--------------|-----------------|
|        | 1000 | 0.84 | 0.89 |
| 1.0e-1 | 1000 | 0.80 | 0.80 |
|        | 1000 | 0.81 | 0.85 |
|        | 1000 | 0.84 | 0.80 |
| 1.0e-3 | 1000 | 0.80 | 0.76 |
|        | 1000 | 0.81 | 0.78 |
|        | 3000 | 0.84 | 0.00017 |
| 1.0e-7 | 3000 | 0.80 | 0.00017 |
|        | 3000 | 0.81 | 0.00016 |
|        | 1000 | 0.84 | 0.00000 |
| 1.0e-10 | 1000 | 0.80 | 0.00000 |
|        | 1000 | 0.81 | 0.00000 |

**After scaling cost vs Iteration for** $(\alpha = 1.0e - 1, 1.0e - 7)$



**Prediction for training and testing values when** $\alpha = 1.0e - 1, 1.0e - 7$

## IV. DISCUSSION

In this section, we discuss key aspects of our analysis and implementation.

### A. Data Analysis and Pre-processing

We meticulously analyzed the dataset, employing data pre-processing techniques. This involved dropping unnecessary columns and handling missing values, crucial steps to ensure the quality and reliability of our analysis.

### B. Cost Function and Gradient Descent Implementation

Creating and implementing the cost function and gradient descent algorithm were pivotal steps in our study. We delved into the transformation of univariate to vectorized implementation, gaining a deeper understanding of the underlying mathematical principles.

### C. Effect of Learning Rate on Data

Our exploration included observing the impact of the learning rate on the dataset. Notably, we varied the learning rate and observed corresponding changes in the data. For instance, adopting a relatively low rate (e.g., $1.0 \times 10^{-9}$) without scaling yielded promising predictions on the test data. Furthermore, we discovered that iteration plays a crucial role in achieving better results with a larger iteration count (e.g., 3000) at the same learning rate.

### D. Scaling and its Influence on Model Output

Our study delved into the significance of scaling in the context of model output. Without scaling, certain learning rates (e.g., 0.1, 0.01) increased the cost function instead of decreasing it. However, after implementing scaling, we observed improved outputs even for larger learning rates like 0.001.

These findings underscore the importance of thoughtful data preprocessing, algorithm implementation, and parameter tuning in the success of machine learning models.

## REFERENCES