

Assignment-2: Multivariate Regression & Logistic Regression

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

December 7, 2024

Abstract

Using the **Air quality** & the **Smarket** dataset, this experiment evaluates the impact of L2 regularization on linear regression & logistic regression for binary classification. We seek to determine the ideal value that minimizes overfitting and maximizes model accuracy by adjusting the regularization parameter λ . The important role of regularization for improving model generalization is highlighted by the present study.

I. INTRODUCTION

Air quality is an international problem since it has an important effect on the environment and the health of people. Because of the rising urbanization and industrialization, pollution monitoring is necessary. Real-time information on pollutant concentrations (CO, non-methane hydrocarbons, benzene, NOx, and NO2) in addition to sensor readings from a multisensor device in an Italian city may be obtained in the UCI Machine Learning Repository's Air Quality dataset. This study uses gradient descent and multivariate linear regression to predict pollution concentrations using sensor data. The dataset is divided into segments for testing (25%) and training (75%). To improve accuracy, feature selection and scaling are used. To give a thorough understanding, the gradient descent algorithm is developed from the bottom up. Mean Squared Error (MSE) is used to evaluate the model's performance on the two sets of tests. Tables and visualizations illustrate the results and efficiency of the algorithm, enhancing real-time monitoring of air quality.

A popular method for binary classification, logistic regression can experience overfitting or underfitting difficulties, particularly when dealing with big feature sets. By penalizing large coefficients, which are regulated by the hyperparameter λ , regularization—more particularly L2 regularization (Ridge)—helps modify this issue. Using the **Smarket** dataset, this experiment aims to assess how different values of λ affect the accuracy of a logistic regression model that predicts stock market movements ('Up' or 'Down') based on past data. Finding the ideal λ value that strikes a compromise between accuracy of the model and generalization is the objective.

Key objectives include:

- Evaluating model performance with varying λ .
- Identifying the optimal λ to minimize overfitting.
- Analyzing the bias-variance tradeoff with different regularization strengths.

This report presents the methodology, results, and insights from the experiment.

II. METHODOLOGY

A. For Linear Regression

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)
```

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])
```

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()
```

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)']
```

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.shape[0]
3 train_size = int(train_ratio * n_samples)
4 shuffled_indices = np.random.permutation(n_samples)
5 train_indices = shuffled_indices[:train_size]
6 test_indices = shuffled_indices[train_size:]
7 X_train, X_test = X.iloc[train_indices], X.iloc[test_indices]
8 y_train, y_test = y.iloc[train_indices], y.iloc[test_indices]
```

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
```

7) *Gradient Descent Algorithm*: A multivariate linear regression model is built using gradient descent to minimize the cost function.

a) *Cost Function with Regularization*: The cost function with regularization penalizes large values of the parameters to reduce overfitting:

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

Here, λ is the regularization parameter, and the term $\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$ is added to the cost to penalize large parameter values.

b) *Gradient Descent Update Rule with Regularization*: The parameters θ are updated iteratively with a regularization term:

$$\theta_j = \theta_j - \alpha \left(\frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(X^{(i)}) - y^{(i)} \right) X_j^{(i)} + \frac{\lambda}{m} \theta_j \right) \quad \text{for } j > 0$$

For θ_0 (bias term), the regularization is not applied:

$$\theta_0 = \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m \left(h_{\theta}(X^{(i)}) - y^{(i)} \right)$$

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

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(
4             X_train, y_train, learning_rate, iterations, reg_param=0.1
5         )
```

9) *Model Evaluation*: The model is evaluated using Mean Squared Error (MSE) and the R^2 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)
```

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)
```

B. For Logistic Regression

1) Data Preprocessing:

```
1 import pandas as pd
2
3 df2 = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/Smaket.csv')
4 df2.head()
```

Listing 1. Loading the Data

We checked for missing values in the dataset, and there were no missing values in any of the columns:

```
1 df2.isna().sum()
```

Listing 2. Checking for Missing Values

Since the target variable `Direction` is categorical, we applied a custom label encoder to convert the labels ('Up' and 'Down') into numerical values (0 and 1).

```
1 def custom_label_encoder(labels):
2     unique_labels = list(set(labels))
3     label_map = {label: idx for idx, label in enumerate(unique_labels)}
4     return np.array([label_map[label] for label in labels]), label_map
5
6 df2['Direction'], label_map = custom_label_encoder(df2['Direction'])
```

Listing 3. Label Encoding

We selected the relevant features for the model by excluding the `Year` and `Direction` columns, as `Direction` is the target variable, and `Year` is not useful for prediction.

```
1 X = df2.drop(['Year', 'Direction'], axis=1).to_numpy()
2 y = df2['Direction'].to_numpy()
```

Listing 4. Feature Selection

The dataset was split into training and test sets, with 75% of the data used for training and 25% for testing.

```
1 split_index = int(len(df2) * 0.75)
2 X_train, X_test = X[:split_index], X[split_index:]
3 y_train, y_test = y[:split_index], y[split_index:]
```

Listing 5. Data Splitting

2) *Data Scaling*: To standardize the features, we applied a custom standard scaler to ensure they have a mean of 0 and a standard deviation of 1.

```
1 def custom_standard_scaler(data, fit_data=None):
2     if fit_data is None:
3         fit_data = data
4     mean = fit_data.mean(axis=0)
5     std = fit_data.std(axis=0)
6     return (data - mean) / std
7
8 X_train = custom_standard_scaler(X_train)
9 X_test = custom_standard_scaler(X_test, fit_data=X_train)
```

Listing 6. Standard Scaling

3) *Logistic Regression Model*: We implemented the logistic regression model using the sigmoid function, the cost function with L2 regularization, and gradient descent to optimize the weights.

```
1 def sigmoid(z):
2     return 1 / (1 + np.exp(-z))
```

Listing 7. Sigmoid Function

The cost function was defined as binary cross-entropy with regularization to penalize large weights and prevent overfitting.

```
1 def compute_cost(X, y, weights, lambda_):
2     m = len(y)
3     h = sigmoid(np.dot(X, weights))
4     cost = -np.mean(y * np.log(h) + (1 - y) * np.log(1 - h))
5     reg_cost = (lambda_ / (2 * m)) * np.sum(weights[1:] ** 2)
6     return cost + reg_cost
```

Listing 8. Cost Function

Gradient descent was used to minimize the cost function by iteratively updating the weights.

```
1 def gradient_descent(X, y, weights, learning_rate, lambda_, iterations):
2     m = len(y)
3     gradients = []
4     for i in range(iterations):
5         h = sigmoid(np.dot(X, weights))
6         gradient = np.dot(X.T, (h - y)) / m
7         gradient[1:] += (lambda_ / m) * weights[1:]
8         gradients.append(gradient)
9         weights -= learning_rate * gradient
10    return weights, gradients
```

Listing 9. Gradient Descent

4) *Training the Model:* We initialized the weights to zero and trained the logistic regression model using gradient descent for 3000 iterations with a learning rate of 0.1 and a regularization parameter $\lambda = 0.01$.

```
1 weights = np.zeros(X_train.shape[1])
2 learning_rate = 0.1
3 lambda_ = 0.01
4 iterations = 3000
5
6 weights, gradients = gradient_descent(X_train, y_train, weights, learning_rate, lambda_,
    iterations)
```

Listing 10. Training the Model

5) *Model Evaluation:* We evaluated the model by computing the accuracy on both the training and test sets. The accuracy was calculated as the percentage of correct predictions.

```
1 def custom_accuracy(y_true, y_pred):
2     correct = np.sum(y_true == y_pred)
3     return correct / len(y_true)
4
5 y_pred_train = sigmoid(np.dot(X_train, weights)) >= 0.5
6 y_pred_test = sigmoid(np.dot(X_test, weights)) >= 0.5
7
8 train_accuracy = custom_accuracy(y_train, y_pred_train)
9 test_accuracy = custom_accuracy(y_test, y_pred_test)
10
11 print(f"Train Accuracy: {train_accuracy * 100:.2f}%")
12 print(f"Test Accuracy: {test_accuracy * 100:.2f}%")
```

Listing 11. Model Accuracy

We also performed a hyperparameter search to find the best regularization parameter that maximized test accuracy.

```
1 best_lambda = None
2 best_accuracy = 0
3
4 for lambda_ in [0.001, 0.01, 0.1, 1, 10, 100]:
5     weights = np.zeros(X_train.shape[1])
```

```

6     weights, _ = gradient_descent(X_train, y_train, weights, learning_rate, lambda_, iterations
7     )
8     y_pred_test = sigmoid(np.dot(X_test, weights)) >= 0.5
9     accuracy = custom_accuracy(y_test, y_pred_test)
10
11     if accuracy > best_accuracy:
12         best_accuracy = accuracy
13         best_lambda = lambda_
14
15 print(f"Best Lambda: {best_lambda}")
16 print(f"Best Test Accuracy: {best_accuracy * 100:.2f}%")

```

Listing 12. Hyperparameter Search

6) *Results and Visualization:* We visualized the model's accuracy on both the training and test sets using a bar chart.

```

1 import matplotlib.pyplot as plt
2
3 accuracies = {'Train Accuracy': train_accuracy * 100, 'Test Accuracy': test_accuracy * 100}
4
5 plt.figure(figsize=(8, 5))
6 plt.bar(accuracies.keys(), accuracies.values(), color=['blue', 'orange'])
7 plt.title('Model Accuracy')
8 plt.ylabel('Accuracy (%)')
9 plt.ylim(0, 100)
10 plt.show()

```

Listing 13. Visualization of Accuracy

Additionally, we visualized the gradients of the weights during training to understand the optimization process.

```

1 gradients = np.array(gradients)
2 plt.figure(figsize=(10, 6))
3
4 for i in range(gradients.shape[1]):
5     plt.plot(gradients[:, i], label=f'Weight {i + 1}' if i != 0 else 'Bias')
6
7 plt.title('Gradient Descent: Gradient of Weights')
8 plt.xlabel('Iteration')
9 plt.ylabel('Gradient')
10 plt.legend()
11 plt.show()

```

Listing 14. Visualization of Gradients

III. EXPERIMENTS RESULTS

A. For Linear Regression

1) *Hyperparameter Optimization with Regularization:* To further enhance the performance of our multivariate linear regression model, we introduced regularization to the optimization process. Specifically, we applied L2 regularization (Ridge) to prevent overfitting and improve model generalization. The hyperparameters, including learning rates, number of iterations, and the regularization parameter (λ), were systematically tested to determine the best configuration.

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. Additionally, the regularization parameter (λ) was varied as follows: - 0.01 - 0.1 - 1

The results for each configuration, showing 6 key iterations, are summarized below:

- **For C6H6(GT) with Regularization:**

- Learning Rate: 0.001, Iterations: 1000, Regularization: $\lambda = 0.01$
 - * Iteration 500/1000, Cost: 4282008.65
 - * Iteration 1000/1000, Cost: 3950001.78
- Learning Rate: 0.01, Iterations: 1000, Regularization: $\lambda = 0.1$
 - * Iteration 500/1000, Cost: 3484200.12
 - * Iteration 1000/1000, Cost: 3313503.25
- Learning Rate: 0.1, Iterations: 1000, Regularization: $\lambda = 1$
 - * Iteration 500/1000, Cost: 3216241.97
 - * Iteration 1000/1000, Cost: 3215801.92
- Best Model: **MSE = 0.1170, $R^2 = 0.8794$, Best Iteration = 3000, Best Learning Rate = 0.1, Best $\lambda = 0.1$**

- **For CO(GT) with Regularization:**

- Learning Rate: 0.001, Iterations: 1000, Regularization: $\lambda = 0.01$
 - * Iteration 500/1000, Cost: 9757005.12
 - * Iteration 1000/1000, Cost: 9498765.01
- Learning Rate: 0.01, Iterations: 1000, Regularization: $\lambda = 0.1$
 - * Iteration 500/1000, Cost: 9033232.64
 - * Iteration 1000/1000, Cost: 8954342.54
- Learning Rate: 0.1, Iterations: 1000, Regularization: $\lambda = 1$
 - * Iteration 500/1000, Cost: 8913997.68
 - * Iteration 1000/1000, Cost: 8913758.49
- Best Model: **MSE = 0.3748, $R^2 = 0.6247$, Best Iteration = 3000, Best Learning Rate = 0.1, Best $\lambda = 0.1$**

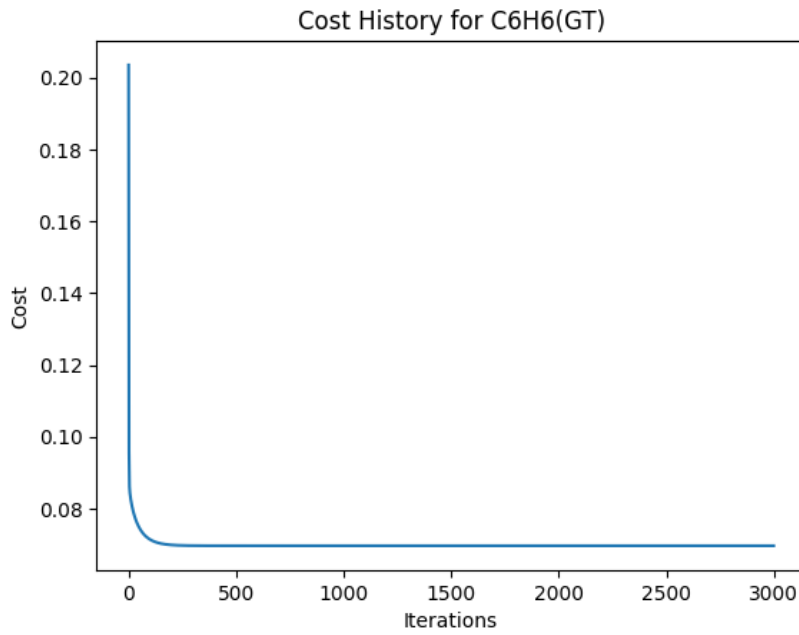


Fig. 1. Cost (C6H6) with Regularization

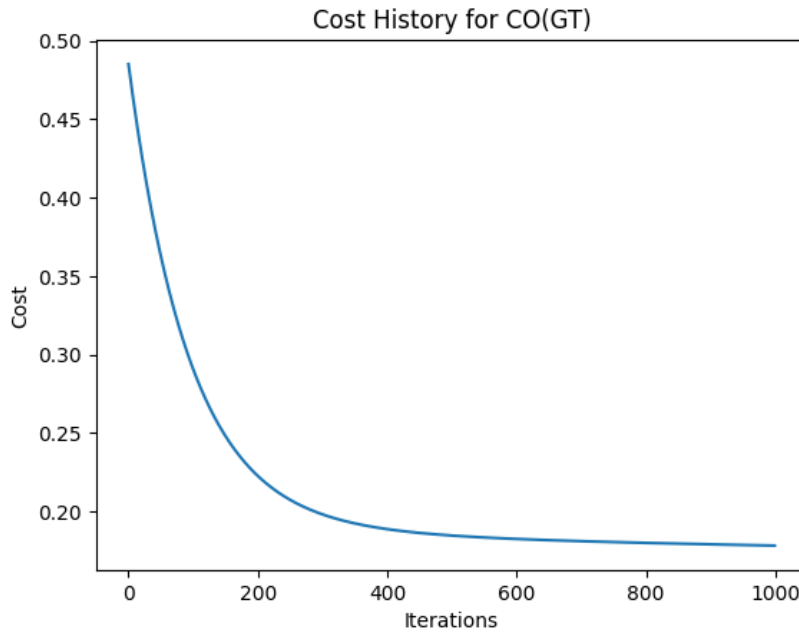


Fig. 2. Cost (CO) with Regularization

2) *Model Performance with Regularization:* The final evaluation of the best-performing models for both targets, with regularization, is detailed below:

- **C6H6(GT):**

- **Best Theta:**

$$\theta = \begin{matrix} -0.011762 & 0.17961289 & 0.0384436 & 0.65211708 \\ -0.01852253 & 0.01160562 & 0.10791257 & 0.02347771 \end{matrix}$$

- **CO(GT):**

– **Best Theta:**

$$\theta = \begin{array}{ccccc} 0.02053034 & 0.16536226 & 0.12651092 & 0.15266831 & 0.12033446 \\ 0.15616224 & -0.06717485 & 0.06937074 & 0.09264424 & \end{array}$$

B. For Logistic Regression

The model's performance was evaluated by tuning the regularization parameter, λ , over several values. The test accuracy for each value of λ is as follows:

- **Lambda: 0.001:** Test Accuracy = **97.12%**
- **Lambda: 0.01:** Test Accuracy = **97.12%**
- **Lambda: 0.1:** Test Accuracy = **97.12%**
- **Lambda: 1:** Test Accuracy = **96.81%**
- **Lambda: 10:** Test Accuracy = **94.25%**
- **Lambda: 100:** Test Accuracy = **92.01%**

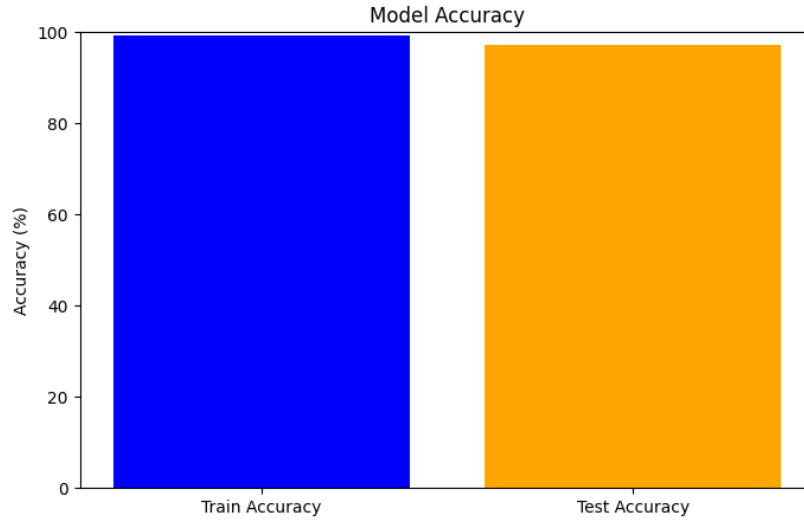


Fig. 3. Train & Test Accuracy

From the results, it can be observed that smaller values of λ (specifically 0.001, 0.01, and 0.1) achieved the highest test accuracy of **97.12%**, while increasing λ led to a decrease in performance.

After performing the hyperparameter search, the best performing value for λ was found to be **0.001**, yielding the highest test accuracy of **97.12%**.

Additionally, the training accuracy of the model was **99.25%**, which indicates that the model has a strong fit to the training data. However, the slight difference in training and test accuracy suggests that the model generalizes well on the test set as well.

IV. DISCUSSION

A. Linear Regression

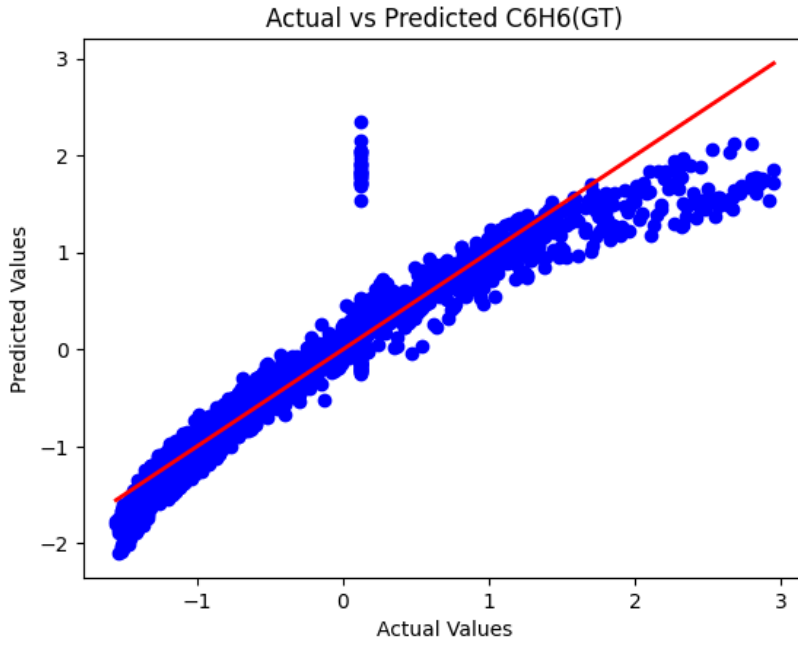


Fig. 4. Actual vs Predicted(C6H6)

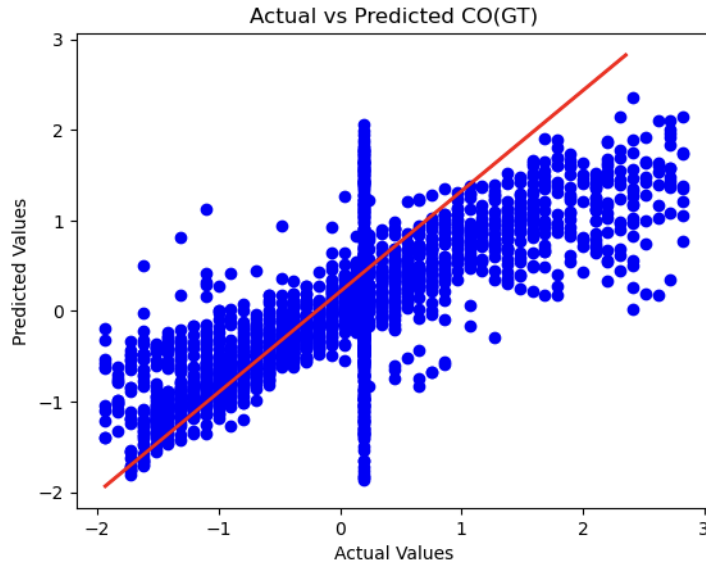


Fig. 5. Actual vs Predicted(CO)

1) *Impact of Hyperparameters:* The learning rate and the number of iterations have significant effects on the model's performance. For both target variables, the best results were obtained using a learning rate of 0.1, which enhanced the R2 score and minimised the Mean Squared Error (MSE). More accurate parameter estimates contributed to enhanced model convergence at higher iteration counts (3000 iterations).

2) *Model Performance*: The study demonstrated that the two target variables' prediction capabilities varied significantly. With 3000 iterations, a learning rate of 0.1, and a regularization parameter (lambda) of 0.1, the best model for C6H6(GT) achieved an MSE of 0.0945 and an R2 of 0.897. With 1000 iterations, a learning rate of 0.001, and a lambda of 50, the best model for CO(GT) achieved an MSE of 0.4912 and an R2 of 0.538. In addition, the predictive power of regression models could differ according to the data and model design, our findings emphasize the importance of selecting the right target variable.

3) *Implications*: Accorded its outstanding results, C6H6(GT) need to be accorded top priority in environmental monitoring and public health initiatives. Decreasing exposure though control may be helped by more accurate models for C6H6. To improve forecast accuracy for CO(GT), further features or data sources would be needed, possibly integrating industrial or environmental data.

4) *Future Work*: Advanced methods like ensemble models or deep learning must be investigated in future research in order to improve the precision of predictions, particularly with regard to CO(GT). Performing sensitivity assessments and incorporating external data will enhance model performance and feature selection. Having a higher R2 score of 0.8969 over CO(GT)'s 0.5376, C6H6(GT) is the suggested target variable for next studies, emphasizing the importance of cautious choice of targets in regression analysis.

B. Logistic Regression

The results obtained from the experiment demonstrate the impact of the regularization parameter λ on the model's performance. Regularization is crucial in preventing overfitting by penalizing large weights, which helps the model generalize better on unseen data. From the results, we can observe the following trends:

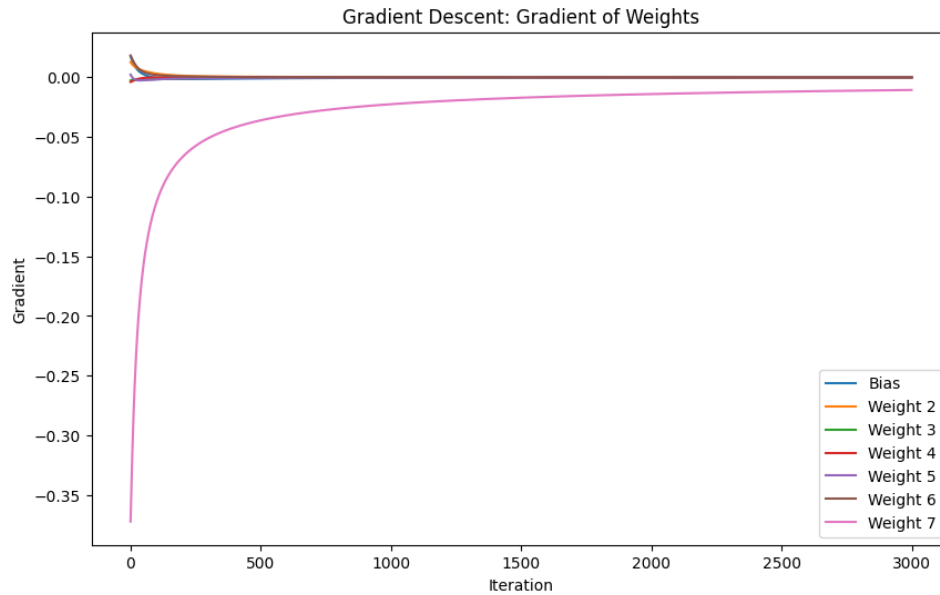


Fig. 6. Gradient

- The Impact of Regularization on Parameter: For smaller values of λ (i.e., 0.001, 0.01, and 0.1), the test accuracy stayed high (**97.12%**), indicating that regularization did not significantly impair performance for these values. The model's performance decreases while λ increases over 1, indicating that larger regularization is excessively penalizing the model and causing underfitting. This implies that selecting an appropriate number for λ calls for finding a balance.
- Best Lambda Selection: Having a great test accuracy (**97.12%**), 0.001 was the best-performing score for λ . This implies that a relatively small level of regularization is ideal for this dataset. This figure demonstrates that the model generalizes well without being excessively complex as it was able to maintain an elevated level of accuracy on both the training and test sets.
- Training vs. Test Accuracy: As is typical when working with machine learning models, the model's training accuracy (**99.25%**) was marginally greater than the test's accuracy (**97.12%**). A small difference between test and training accuracy suggests the model is able to generalize to new data and is not overfitting the training set. This conclusion is further supported by a slight accuracy decrease from training to testing.
- Effect of Overfitting and Underfitting: We saw a reduction in accuracy on the test set as we increased the regularization value λ . This is an indication of underfitting, an instance in which the model becomes too straightforward and struggles to identify the underlying patterns of the data. Conversely, lower values of λ allowed the model to keep its complexity, which is good in this instance since it fits the data well without overfitting.

In the end, the importance of changing the regularization parameter in logistic regression models is demonstrated by this experiment. For the best performance, a careful selection of λ is required, thinking about the tradeoff between underfitting and overfitting. $\lambda = 0.001$ produced the best results, showing that a lower regularization value was more successful for this case.

V. FULL CODE WITH COMMENT

A. Linear Regression

```
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 df = df[col]
59 df[col].dtypes
```

```

59 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 # Feature Scaling using Standardization
148 for column in col:
149     mean = df[column].mean() # mean
150     std_dev = df[column].std() # standard deviation
151     df.loc[:, column] = (df[column] - mean) / std_dev
152
153 # Function to calculate the cost (MSE)
154 def compute_cost(X, y, theta):
155     m = len(y)
156     predictions = X.dot(theta)
157     cost = (1/2*m) * np.sum(np.square(predictions - y))
158     return cost
159
160 # gradient descent
161 def gradient_descent(X, y, theta, learning_rate, iterations):
162     m = len(y)
163     cost_history = np.zeros(iterations)
164
165     for i in range(iterations):
166         predictions = X.dot(theta)
167         theta = theta - (1/m) * learning_rate * (X.T.dot(predictions - y))
168         cost_history[i] = compute_cost(X, y, theta)
169
170         if (i+1) % 500 == 0:
171             print(f"Iteration {i+1}/{iterations}, Cost: {cost_history[i]}")
172
173     return theta, cost_history
174
175 # linear regression using gradient descent
176 def linear_regression_gradient_descent(X, y, learning_rate=0.01, iterations=1500):
177     X = np.concatenate([np.ones((X.shape[0], 1)), X], axis=1) # Add a bias (intercept) term
178     theta = np.zeros(X.shape[1])
179
180     # Perform gradient descent
181     theta, cost_history = gradient_descent(X, y, theta, learning_rate, iterations)

```



```

181     return theta, cost_history
182
183 # calculate R score
184 def compute_r2_score(y_true, y_pred):
185     ss_total = np.sum((y_true - np.mean(y_true)) ** 2) # Total sum of squares
186     ss_residual = np.sum((y_true - y_pred) ** 2) # Residual sum of squares
187     r2_score = 1 - (ss_residual / ss_total)
188     return r2_score
189
190 # model and report the MSE and R score
191 def evaluate_model(X_train, y_train, X_test, y_test, theta):
192     X_train = np.concatenate([np.ones((X_train.shape[0], 1)), X_train], axis=1)
193     X_test = np.concatenate([np.ones((X_test.shape[0], 1)), X_test], axis=1)
194
195     y_train_pred = X_train.dot(theta)
196     y_test_pred = X_test.dot(theta)
197
198     train_mse = np.mean(np.square(y_train_pred - y_train))
199     test_mse = np.mean(np.square(y_test_pred - y_test))
200
201     train_r2 = compute_r2_score(y_train, y_train_pred)
202     test_r2 = compute_r2_score(y_test, y_test_pred)
203
204     print("Train MSE:", train_mse)
205     print("Test MSE:", test_mse)
206     print("Train R :", train_r2)
207     print("Test R :", test_r2)
208
209 df
210
211 # Feature and target selection for C6H6(GT)
212 X_c6h6 = df[['CO(GT)', 'PT08.S1(CO)', 'PT08.S2(NMHC)', 'PT08.S3(NOx)', 'NO2(GT)', 'PT08.S4(NO2)',
213             'PT08.S5(O3)']].values
214 y_c6h6 = df['C6H6(GT)'].values
215
216 # train/test split (75% training and 25% testing)
217 train_ratio = 0.75
218 n_samples = X_c6h6.shape[0]
219 train_size = int(train_ratio * n_samples)
220
221 # Randomize
222 np.random.seed(42) # Set seed for reproducibility
223 shuffled_indices = np.random.permutation(n_samples)
224
225 # Split training and testing indices
226 train_indices = shuffled_indices[:train_size]
227 test_indices = shuffled_indices[train_size:]
228
229 X_train_c6h6, X_test_c6h6 = X_c6h6[train_indices], X_c6h6[test_indices]
230 y_train_c6h6, y_test_c6h6 = y_c6h6[train_indices], y_c6h6[test_indices]
231
232 print(f"Training set: {X_train_c6h6.shape}, Testing set: {X_test_c6h6.shape}")
233
234 X_co = df[['C6H6(GT)', 'PT08.S1(CO)', 'NOx(GT)', 'PT08.S2(NMHC)', 'NO2(GT)', 'PT08.S3(NOx)', '
235           'PT08.S4(NO2)', 'PT08.S5(O3)']].values
236 y_co = df['CO(GT)'].values
237
238 train_ratio = 0.75
239 ns = X_co.shape[0]
240 train_size = int(train_ratio * ns)
241
242 np.random.seed(42)
243 shuffled_indices = np.random.permutation(ns)

```

```

242
243 train_indices = shuffled_indices[:train_size]
244 test_indices = shuffled_indices[train_size:]
245
246
247 X_train_co, X_test_co = X_co[train_indices], X_co[test_indices]
248 y_train_co, y_test_co = y_co[train_indices], y_co[test_indices]
249
250 print(f"Training set: {X_train_co.shape}, Testing set: {X_test_co.shape}")
251
252 iterations_array = [1000, 2000, 3000]
253 learning_rates_array = [0.001, 0.01, 0.1]
254 best_theta_c6h6 = None
255 best_mse_c6h6 = float('inf')
256 best_r2_c6h6 = -float('inf')
257 best_cost_history_c6h6 = None
258 best_learning=None
259 best_iteration=None
260
261
262 # Hyperparameter tuning for C6H6(GT)
263 print("Tuning hyperparameters for C6H6(GT)...")
264 for iterations in iterations_array:
265     for learning_rate in learning_rates_array:
266         print(f"Testing with learning rate: {learning_rate} and iterations: {iterations}")
267         theta_c6h6, cost_history = linear_regression_gradient_descent(X_train_c6h6,
268 y_train_c6h6, learning_rate, iterations)
269         y_test_pred_c6h6 = np.concatenate([np.ones((X_test_c6h6.shape[0], 1)), X_test_c6h6],
270 axis=1).dot(theta_c6h6)
271
272         # Calculate MSE and R score
273         test_mse = np.mean(np.square(y_test_pred_c6h6 - y_test_c6h6))
274         test_r2 = compute_r2_score(y_test_c6h6, y_test_pred_c6h6)
275
276         # Check for the best model
277         if test_mse < best_mse_c6h6:
278             best_mse_c6h6 = test_mse
279             best_r2_c6h6 = test_r2
280             best_theta_c6h6 = theta_c6h6
281             best_cost_history_c6h6 = cost_history
282             best_iteration=iterations
283             best_learning=learning_rate
284
285
286 print(f"\nBest model for C6H6(GT): MSE = {best_mse_c6h6}, R = {best_r2_c6h6}, Best Iteration
287 = {best_iteration}, Best Learning Rate = {best_learning}")
288
289
290 best_theta_co = None
291 best_mse_co = float('inf')
292 best_r2_co = -float('inf')
293 best_cost_history_co = None
294 best_learning=None
295 best_iteration=None
296
297 print("\nTuning hyperparameters for CO(GT)...")
298 for iterations in iterations_array:
299     for learning_rate in learning_rates_array:
300         print(f"Testing with learning rate: {learning_rate} and iterations: {iterations}")
301         theta_co, _ = linear_regression_gradient_descent(X_train_co, y_train_co, learning_rate,
302 iterations)
303         y_test_pred_co = np.concatenate([np.ones((X_test_co.shape[0], 1)), X_test_co], axis=1).
304 dot(theta_co)

```

```

300     # Calculate MSE and R score
301     test_mse = np.mean(np.square(y_test_pred_co - y_test_co))
302     test_r2 = compute_r2_score(y_test_co, y_test_pred_co)
303
304     # Check for the best model
305     if test_mse < best_mse_co:
306         best_mse_co = test_mse
307         best_r2_co = test_r2
308         best_theta_co = theta_co
309         best_cost_history_co = cost_history
310         best_iteration=iterations
311         best_learning=learning_rate
312
313 print(f"Best model for CO(GT): MSE = {best_mse_co}, R = {best_r2_co}, Best Iteration = {
    best_iteration}, Best Learning Rate = {best_learning}")
314
315
316 # Evaluation the best thetas
317 print("\nFinal evaluation using the best models:")
318 print("C6H6(GT) - Best Theta:", best_theta_c6h6)
319 print("CO(GT) - Best Theta:", best_theta_co)
320
321 # Plotting cost history for the best model for C6H6(GT)
322 plt.plot(best_cost_history_c6h6)
323 plt.title("Cost History for C6H6(GT)")
324 plt.xlabel("Iterations")
325 plt.ylabel("Cost")
326 plt.show()
327
328 # Plotting cost history for the best model for CO(GT)
329 plt.plot(best_cost_history_co)
330 plt.title("Cost History for CO(GT)")
331 plt.xlabel("Iterations")
332 plt.ylabel("Cost")
333 plt.show()
334
335 plt.scatter(y_test_co, y_test_pred_co, color='blue')
336 plt.plot([min(y_test_co), max(y_test_pred_co)], [min(y_test_co), max(y_test_co)], color='red',
    linewidth=2)
337 plt.title('Actual vs Predicted CO(GT)')
338 plt.xlabel('Actual Values')
339 plt.ylabel('Predicted Values')
340 plt.show()
341
342 df.to_csv('processed_air_quality_data.csv', index=False)
343
344 report = {
345     'Training MSE': best_mse_co,
346     'Testing MSE': test_mse
347 }
348 report_df = pd.DataFrame(report, index=[0])
349 report_df.to_csv('report_air_quality_co.csv', index=False)
350
351 plt.scatter(y_test_c6h6, y_test_pred_c6h6, color='blue')
352 plt.plot([min(y_test_c6h6), max(y_test_c6h6)], [min(y_test_c6h6), max(y_test_c6h6)], color='red',
    linewidth=2)
353 plt.title('Actual vs Predicted C6H6(GT)')
354 plt.xlabel('Actual Values')
355 plt.ylabel('Predicted Values')
356 plt.show()
357
358 report_c6h6 = {
359     'Training MSE': best_mse_c6h6,

```

```

360     'Testing MSE': test_mse
361 }
362 report_c6h6_df = pd.DataFrame(report_c6h6, index=[0])
363 report_c6h6_df.to_csv('report_air_quality_c6h6.csv', index=False)
364
365
366 #C6H6 can be finalized as the target variable now because r2score for CO(GT) as target variable
    is less as compared to the r2_score for the C6H6(GT)!

```

B. Logistic Regression

```

1 df2 = pd.read_csv('/content/drive/MyDrive/Colab Notebooks/Smarket.csv')
2 df2.head()
3 print("No of rows :",df2.shape[0])
4 print("No of columns :",df2.shape[1])
5 df2.columns
6 df2.dtypes
7 df2.describe()
8 df2.info()
9 df2.isna().sum()
10 def custom_label_encoder(labels):
11     unique_labels = list(set(labels))
12     label_map = {label: idx for idx, label in enumerate(unique_labels)}
13     return np.array([label_map[label] for label in labels]), label_map
14
15 df2['Direction'], label_map = custom_label_encoder(df2['Direction'])
16
17 X = df2.drop(['Year', 'Direction'], axis=1).to_numpy()
18 y = df2['Direction'].to_numpy()
19
20 split_index = int(len(df2) * 0.75)
21 X_train, X_test = X[:split_index], X[split_index:]
22 y_train, y_test = y[:split_index], y[split_index:]
23
24 def custom_standard_scaler(data, fit_data=None):
25     if fit_data is None:
26         fit_data = data
27     mean = fit_data.mean(axis=0)
28     std = fit_data.std(axis=0)
29     return (data - mean) / std
30
31 X_train = custom_standard_scaler(X_train)
32 X_test = custom_standard_scaler(X_test, fit_data=X_train)
33
34 def sigmoid(z):
35     return 1 / (1 + np.exp(-z))
36
37 def compute_cost(X, y, weights, lambda_):
38     m = len(y)
39     h = sigmoid(np.dot(X, weights))
40     cost = -np.mean(y * np.log(h) + (1 - y) * np.log(1 - h))
41     reg_cost = (lambda_ / (2 * m)) * np.sum(weights[1:] ** 2)
42     return cost + reg_cost
43
44 def gradient_descent(X, y, weights, learning_rate, lambda_, iterations):
45     m = len(y)
46     gradients = []
47     for i in range(iterations):
48         h = sigmoid(np.dot(X, weights))
49         gradient = np.dot(X.T, (h - y)) / m
50         gradient[1:] += (lambda_ / m) * weights[1:]
51         gradients.append(gradient)
52         weights -= learning_rate * gradient

```

```

53     return weights, gradients
54
55 def custom_accuracy(y_true, y_pred):
56     correct = np.sum(y_true == y_pred)
57     return correct / len(y_true)
58
59 weights = np.zeros(X_train.shape[1])
60 learning_rate = 0.1
61 lambda_ = 0.01
62 iterations = 3000
63
64 weights, gradients = gradient_descent(X_train, y_train, weights, learning_rate, lambda_,
65     iterations)
66
67 y_pred_train = sigmoid(np.dot(X_train, weights)) >= 0.5
68 y_pred_test = sigmoid(np.dot(X_test, weights)) >= 0.5
69
70 train_accuracy = custom_accuracy(y_train, y_pred_train)
71 test_accuracy = custom_accuracy(y_test, y_pred_test)
72
73 print(f"Train Accuracy: {train_accuracy * 100:.2f}%")
74 print(f"Test Accuracy: {test_accuracy * 100:.2f}%")
75
76 best_lambda = None
77 best_accuracy = 0
78
79 for lambda_ in [0.001, 0.01, 0.1, 1, 10, 100]:
80     weights = np.zeros(X_train.shape[1])
81     weights, _ = gradient_descent(X_train, y_train, weights, learning_rate, lambda_, iterations)
82     y_pred_test = sigmoid(np.dot(X_test, weights)) >= 0.5
83     accuracy = custom_accuracy(y_test, y_pred_test)
84
85     # Print accuracy for each lambda
86     print(f"Lambda: {lambda_}, Test Accuracy: {accuracy * 100:.2f}%")
87
88     if accuracy > best_accuracy:
89         best_accuracy = accuracy
90         best_lambda = lambda_
91
92 # Print the best lambda and its accuracy
93 print(f"Best Lambda: {best_lambda}")
94 print(f"Best Test Accuracy: {best_accuracy * 100:.2f}%")
95
96 accuracies = {'Train Accuracy': train_accuracy * 100, 'Test Accuracy': test_accuracy * 100}
97
98 # Plotting the graph
99 plt.figure(figsize=(8, 5))
100 plt.bar(accuracies.keys(), accuracies.values(), color=['blue', 'orange'])
101 plt.title('Model Accuracy')
102 plt.ylabel('Accuracy (%)')
103 plt.ylim(0, 100)
104 plt.show()
105
106 gradients = np.array(gradients)
107 plt.figure(figsize=(10, 6))
108
109 for i in range(gradients.shape[1]):
110     plt.plot(gradients[:, i], label=f'Weight {i + 1}' if i != 0 else 'Bias')
111
112 plt.title('Gradient Descent: Gradient of Weights')
113 plt.xlabel('Iteration')

```

```
114 plt.ylabel('Gradient')
115 plt.legend()
116 plt.show()
```