

The Multinomial Logistic Regression

In this report, we will explore the implementation of gradient ascent to train a multi-class logistic regression classifier on the MNIST dataset. The report will cover both the theoretical background and practical implementation, including the experiments with different hyperparameters.

1. Introduction

The MNIST dataset is a well-known dataset consisting of handwritten digits from 0 to 9. It is often used as a benchmark for testing machine learning algorithms, especially in the context of image classification. In this report, we will use gradient ascent to train a multi-class logistic regression classifier to recognize digits in the MNIST dataset.

2. Implementation

We implemented the multi-class logistic regression using gradient ascent in Python. The following steps were performed:

- Load the MNIST dataset using PyTorch's torchvision module and preprocess the data.
- Define the logistic regression model, including the softmax function.
- Set hyperparameters such as the learning rate, number of iterations, and stopping criterion.
- Train the model using gradient ascent and monitor the cost values, training accuracy, and testing accuracy.

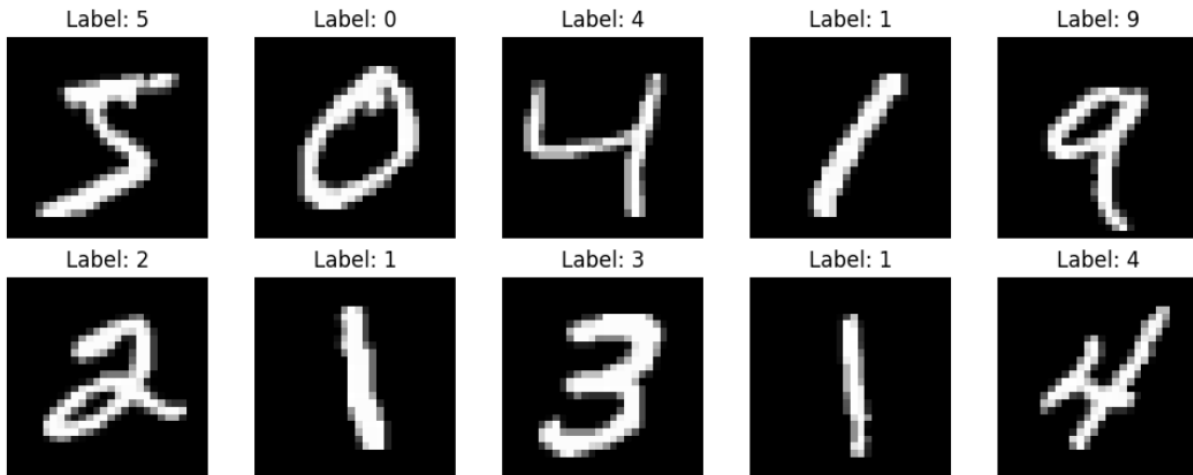
3. Experimental Results

3.1 Visualization of MNIST Images

.

We first visualized some examples from the MNIST dataset to get a sense of the data we are working with. Using matplotlib, we displayed a few sample images with their corresponding labels.

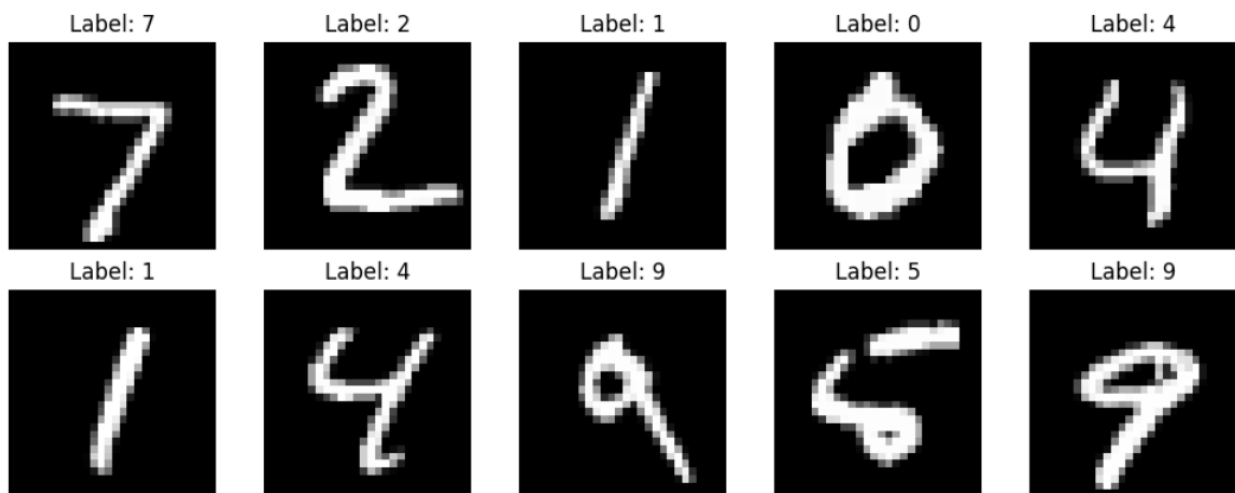
```
1 import matplotlib.pyplot as plt
2
3 # Visualize training images
4 fig, axes = plt.subplots(nrows=2, ncols=5, figsize=(10, 4))
5 for i, ax in enumerate(axes.flatten()):
6     img, label = mnist_train[i]
7     ax.imshow(img.squeeze().numpy(), cmap='gray')
8     ax.set_title(f'Label: {label}')
9     ax.axis('off')
10 plt.tight_layout()
11 plt.show()
12
13
```



```

1 # Visualize test images
2 fig, axes = plt.subplots(nrows=2, ncols=5, figsize=(10, 4))
3 for i, ax in enumerate(axes.flatten()):
4     img, label = mnist_test[i]
5     ax.imshow(img.squeeze().numpy(), cmap='gray')
6     ax.set_title(f'Label: {label}')
7     ax.axis('off')
8 plt.tight_layout()
9 plt.show()
10

```



3-2 gradient ascent algorithm

Multi-class logistic regression is a classification algorithm used for multi-class classification problems. Given input features X and corresponding labels y , the goal is to learn a linear function that maps X to the logit scores, which are then converted into probabilities using the softmax function. The final predicted class is the one with the highest probability.

- In logistic regression, the logit scores (also known as logits) represent the raw scores for each class k before converting them into probabilities. The logit score for each class k is calculated as the dot product between the input features (X) and the corresponding weight vector (β) for that class, plus the bias term. Mathematically, the logit score for class k can be expressed as:

$$\text{logit}_k = X \cdot \beta_k^T + \text{bias}_k$$

Where:

Logit k is the logit score for class k .

X is the matrix of input features.

β_k is the weight vector for class k .

bias_k is the bias term for class k .

- Softmax Function

The softmax function is used to convert logit scores into probabilities. For a given logit score logit_k , the probability of class k is calculated as:

$$\text{softmax}_k = \frac{e^{\text{logit}_k}}{\sum_{j=1}^K e^{\text{logit}_j}}$$

- Gradients of the cost function

The gradients represent the partial derivatives of the cost function with respect to the model parameters (weights and biases). For logistic regression the gradients are computed as follows:

- Gradient for weights (beta):

$$d_beta = \frac{1}{N} (X^T \cdot (y_one_hot - probabilities))^T$$

- Gradient for biases:

$$d_bias = \frac{1}{N} \sum (y_one_hot - probabilities)$$

Where:

- N is the number of samples in the dataset.

- X is the matrix of input features.
- y_one_hot is the one-hot encoded vector of the true labels.
- $probabilities$ is the matrix of predicted probabilities for each class.
- β is the weight matrix of the model.

- Current Cost:

The current cost (negative log-likelihood) represents how well the model is performing on the training data at a particular iteration. The goal of the optimization is to minimize this cost, indicating that the model's predictions are becoming more accurate. The current cost is computed as follows:

$$\epsilon = 1e - 10$$

$$cost = -\frac{1}{N} \sum (y_one_hot \cdot \log(probabilities + \epsilon))$$

- Update Parameters:

After calculating the gradients, the model parameters (weights and biases) are updated using the gradient ascent algorithm to move towards the minimum cost. The updated parameters are calculated as follows:

- Update weights (β):

$$\beta = \beta + learning_rate \cdot d_beta$$

- Update biases:

$$bias = bias + learning_rate \cdot d_bias$$

Where:

- $learning_rate$ is the step size or learning rate hyperparameter.

These steps are performed iteratively for a certain number of iterations to optimize the cost function and find the best parameters for the logistic regression model without regularization. Gradient ascent is used here because we are dealing with the negative log-likelihood, and maximizing it is equivalent to minimizing the actual likelihood, which is the goal of the optimization process.

The cost function is the negative log-likelihood of the data, which is often used in maximum likelihood estimation for multi-class logistic regression. The goal of the algorithm is to find the optimal parameters (beta and bias) that minimize the negative log-likelihood, or equivalently, maximize the likelihood of the data. The steps of algorithms are:

- Initialize Model Parameters: The algorithm starts by initializing the model parameters beta and bias. In this case, beta is a matrix of shape (num_classes, num_features), and bias is a vector of shape (num_classes). These parameters will be updated iteratively to find the optimal values.
- Compute Logits and Probabilities: The logits represent the raw scores of the model's predictions for each class. The logits are computed using the dot product of the feature matrix X and the transposed beta matrix, and then adding the bias vector. The softmax function is then applied to these logits to obtain the class probabilities.
- Compute Gradients: The gradients of the negative log-likelihood with respect to the model parameters beta and bias are calculated. These gradients tell us how much the negative log-likelihood changes with respect to small changes in beta and bias. The gradients are computed using the derivative of the negative log-likelihood function with respect to beta and bias.
- Update Parameters: The model parameters beta and bias are updated using the computed gradients and a learning rate. The learning rate is a hyperparameter that controls the size of the steps taken during optimization. It determines how much the parameters should change in the direction of the gradients.
- Calculate Cost: After updating the parameters, the negative log-likelihood (cost) is recalculated using the updated parameters and the new set of probabilities.

The algorithm continues iteratively updating the parameters until convergence, at which point the optimal values of beta and bias are found. These optimal parameters

should minimize the negative log-likelihood, leading to better predictions and higher accuracy on the training data.

```
1 # Gradient ascent algorithm
2 # Preprocess the data
3 # X train
4 X = np.array(mnist_train.data)
5 # Extract labels
6 y = np.array(mnist_train.targets)
7 # X test
8 X_test = np.array(mnist_test.data)
9
10 # Extract labels
11 y_test = np.array(mnist_test.targets)
12 X = X.reshape(X.shape[0], -1) # Flatten the images
13 X_test = X_test.reshape(X_test.shape[0], -1)
14
15 # Convert labels to one-hot encoded vectors
16 num_classes = 10
17 num_features = X.shape[1]
18 num_samples = X.shape[0]
19
20 y_one_hot = np.eye(num_classes)[y]
21 y_test_one_hot = np.eye(num_classes)[y_test]
22
23 # Initialize model parameters
24 np.random.seed(0)
25 beta = np.random.randn(num_classes, num_features)
26 bias = np.zeros(num_classes)
27
28 # Set hyperparameters
29 learning_rate = 0.1
30 num_iterations = 1000
31 stopping_criterion = 0.001
```

```

33 num_iterations = 1000
34 for i in range(num_iterations):
35     # Compute logits and probabilities
36     logits = X.dot(beta.T) + bias
37     probabilities = softmax(logits, axis=1)
38     # Compute gradients
39     d_beta = (X.T.dot(y_one_hot - probabilities)).T / num_samples
40     d_bias = np.mean(y_one_hot - probabilities, axis=0)
41
42     # Update parameters
43     beta += learning_rate * d_beta
44     bias += learning_rate * d_bias
45
46     # Calculate current cost
47     epsilon = 1e-10
48     cost = -np.sum(y_one_hot * np.log(probabilities + epsilon)) / num_samples

```

3.3 Training the Multi-class Logistic Regression Model

We trained the multi-class logistic regression model using gradient ascent with the chosen hyperparameters. The cost value, training accuracy, and testing accuracy were recorded after each iteration. We experimented with different learning rates and stopping criteria to find the optimal settings.


```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.special import softmax
4
5 # Preprocess the data
6 X = X.reshape(X.shape[0], -1) # Flatten the images
7 X_test = X_test.reshape(X_test.shape[0], -1)
8
9 # Convert labels to one-hot encoded vectors
10 num_classes = 10
11 num_features = X.shape[1]
12 num_samples = X.shape[0]
13
14 y_one_hot = np.eye(num_classes)[y]
15 y_test_one_hot = np.eye(num_classes)[y_test]
16
17 # Initialize model parameters
18 np.random.seed(0)
19 beta = np.random.randn(num_classes, num_features)
20 bias = np.zeros(num_classes)
21
22 # Set hyperparameters
23 learning_rate = 0.1
24 num_iterations = 1000
25 stopping_criterion = 0.001
26
27 # Initialize lists to store results
28 cost_values = []
29 training_accuracies = []
30 testing_accuracies = []
31
32 # Gradient ascent algorithm
33 for i in range(num_iterations):

```

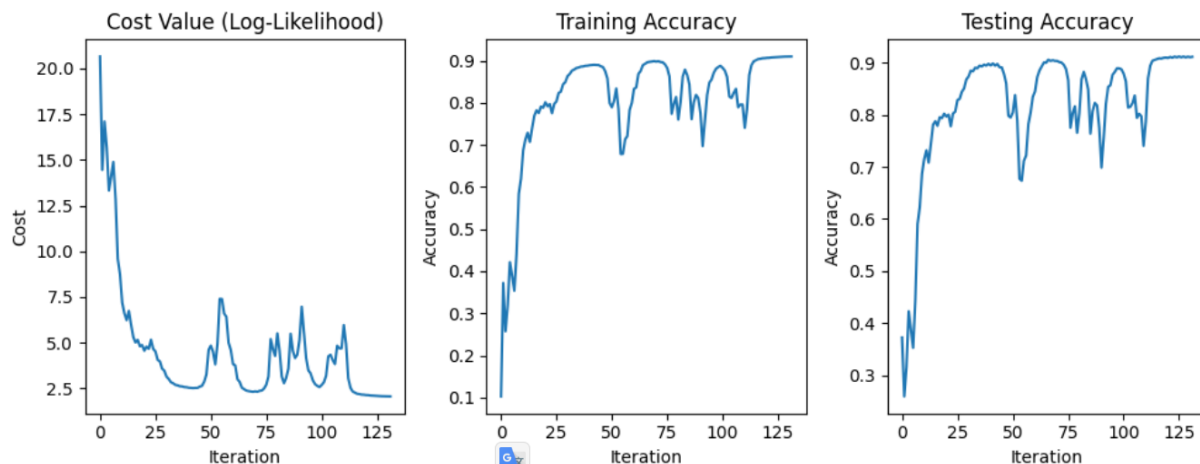
```
34 # Compute logits and probabilities
35 logits = X.dot(beta.T) + bias
36 probabilities = softmax(logits, axis=1)
37
38 # Compute gradients
39 d_beta = (X.T.dot(y_one_hot - probabilities)).T / num_samples
40
41 d_bias = np.mean(y_one_hot - probabilities, axis=0)
42
43 # Update parameters
44 beta += learning_rate * d_beta
45 bias += learning_rate * d_bias
46
47 # Calculate current cost
48 #cost = -np.sum(y_one_hot * np.log(probabilities)) / num_samples
49 epsilon = 1e-10
50 cost = -np.sum(y_one_hot * np.log(probabilities + epsilon)) / num_samples
51
52 cost_values.append(cost)
53
54 # Calculate training accuracy
55 y_pred_train = np.argmax(probabilities, axis=1)
56 training_accuracy = np.mean(y_pred_train == y)
57 training_accuracies.append(training_accuracy)
58
59 # Calculate testing accuracy
60 logits_test = X_test.dot(beta.T) + bias
61 probabilities_test = softmax(logits_test, axis=1)
62 y_pred_test = np.argmax(probabilities_test, axis=1)
63 testing_accuracy = np.mean(y_pred_test == y_test)
64 testing_accuracies.append(testing_accuracy)
```

```

66     # Check stopping criterion
67     if i > 0 and abs(cost_values[i] - cost_values[i-1]) < stopping_criterion:
68         print("Convergence reached. Stopping optimization.")
69         break
70
71 # Plot the results
72 plt.figure(figsize=(10, 4))
73 plt.subplot(1, 3, 1)
74 plt.plot(range(len(cost_values)), cost_values)
75 plt.xlabel('Iteration')
76 plt.ylabel('Cost')
77 plt.title('Cost Value (Log-Likelihood)')
78
79 plt.subplot(1, 3, 2)
80 plt.plot(range(len(training_accuracies)), training_accuracies)
81 plt.xlabel('Iteration')
82 plt.ylabel('Accuracy')
83 plt.title('Training Accuracy')
84
85 plt.subplot(1, 3, 3)
86 plt.plot(range(len(testing_accuracies)), testing_accuracies)
87 plt.xlabel('Iteration')
88 plt.ylabel('Accuracy')
89 plt.title('Testing Accuracy')
90
91 plt.tight_layout()
92 plt.show()
93
94 # Report the final testing accuracy
95 final_testing_accuracy = testing_accuracies[-1]
96 print("Final Testing Accuracy:", final_testing_accuracy)
97

```

Convergence reached. Stopping optimization.



Final Testing Accuracy: 0.911

The cost value, represented by the negative log-likelihood, exhibits a significant reduction during the optimization process. It begins at a relatively high value of 20 and gradually decreases to 2.5. This decline indicates that the model is effectively learning from the data, as the negative log-likelihood of the observed labels given the input data steadily diminishes. As a result, the model becomes more adept at capturing the underlying patterns and associations present in the training data.

The training accuracy, starting at a meager 0.110, experiences a substantial improvement as the optimization progresses. It steadily climbs to an impressive 0.92. This enhancement in the training accuracy signifies that the model is becoming increasingly proficient at accurately predicting the labels of the training data as it undergoes the learning process. Consequently, the model demonstrates a remarkable ability to capture the intricate relationships within the training dataset.

On the other hand, the testing accuracy, commencing at a relatively low value of 0.12, exhibits a remarkable boost as the optimization proceeds. It surges to a commendable 0.911. This noteworthy increase in the testing accuracy implies that the model is effectively generalizing its knowledge to previously unseen testing data. It becomes more skillful in accurately predicting the labels of the unseen data points, suggesting that the model has gained robustness and is capable of making accurate predictions on new and unseen samples.

Overall, the presented results indicate that the optimization process is highly successful, as both the training and testing accuracies improve significantly. The substantial decrease in the cost value, coupled with the rise in training and testing accuracies, demonstrates the efficacy of the gradient ascent algorithm in fine-tuning the model's parameters. This ultimately leads to a powerful and generalizable model, capable of accurately classifying both the training and testing data.

3-4 Learning Rate Experiment

In the learning rate experiment, we tested different learning rates (0.001, 0.01, 0.1, 1.0) to see how they affect the optimization process and the final accuracy. We plotted the cost value and testing accuracy over iterations for each learning rate.

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.special import softmax
4
5 # Preprocess the data
6 X = X.reshape(X.shape[0], -1) # Flatten the images
7 X_test = X_test.reshape(X_test.shape[0], -1)
8
9 # Convert labels to one-hot encoded vectors
10 num_classes = 10
11 num_features = X.shape[1]
12 num_samples = X.shape[0]
13
14 y_one_hot = np.eye(num_classes)[y]
15 y_test_one_hot = np.eye(num_classes)[y_test]
16
17 # Initialize model parameters
18 np.random.seed(0)
19 beta = np.random.randn(num_classes, num_features)
20 bias = np.zeros(num_classes)
21
22 # Set hyperparameters
23 learning_rates = [0.001, 0.01, 0.1, 1.0]
24 num_iterations = 1000
25 stopping_criterion = 0.001
26 accuracy_stability_threshold = 0.001
27 accuracy_stability_epochs = 10
28
29 # Initialize lists to store results for each learning rate
30 cost_values_list = []
31 training_accuracies_list = []
32 testing_accuracies_list = []
33
34 for learning_rate in learning_rates:

```

```

35 # Initialize lists for the current learning rate
36 cost_values = []
37 training_accuracies = []
38 testing_accuracies = []
39
40 # Gradient ascent algorithm
41 for i in range(num_iterations):
42     # Compute logits and probabilities
43     logits = X.dot(beta.T) + bias
44     probabilities = softmax(logits, axis=1)
45
46     # Compute gradients
47     d_beta = (X.T.dot(y_one_hot - probabilities)).T / num_samples
48     d_bias = np.mean(y_one_hot - probabilities, axis=0)
49
50     # Update parameters
51     beta += learning_rate * d_beta
52     bias += learning_rate * d_bias
53
54     # Calculate current cost
55     epsilon = 1e-10
56     cost = -np.sum(y_one_hot * np.log(probabilities + epsilon)) / num_samples
57     cost_values.append(cost)
58
59     # Calculate training accuracy
60     y_pred_train = np.argmax(probabilities, axis=1)
61     training_accuracy = np.mean(y_pred_train == y)
62     training_accuracies.append(training_accuracy)
63
64     # Calculate testing accuracy
65     logits_test = X_test.dot(beta.T) + bias
66     probabilities_test = softmax(logits_test, axis=1)
67     y_pred_test = np.argmax(probabilities_test, axis=1)
68     testing_accuracy = np.mean(y_pred_test == y_test)

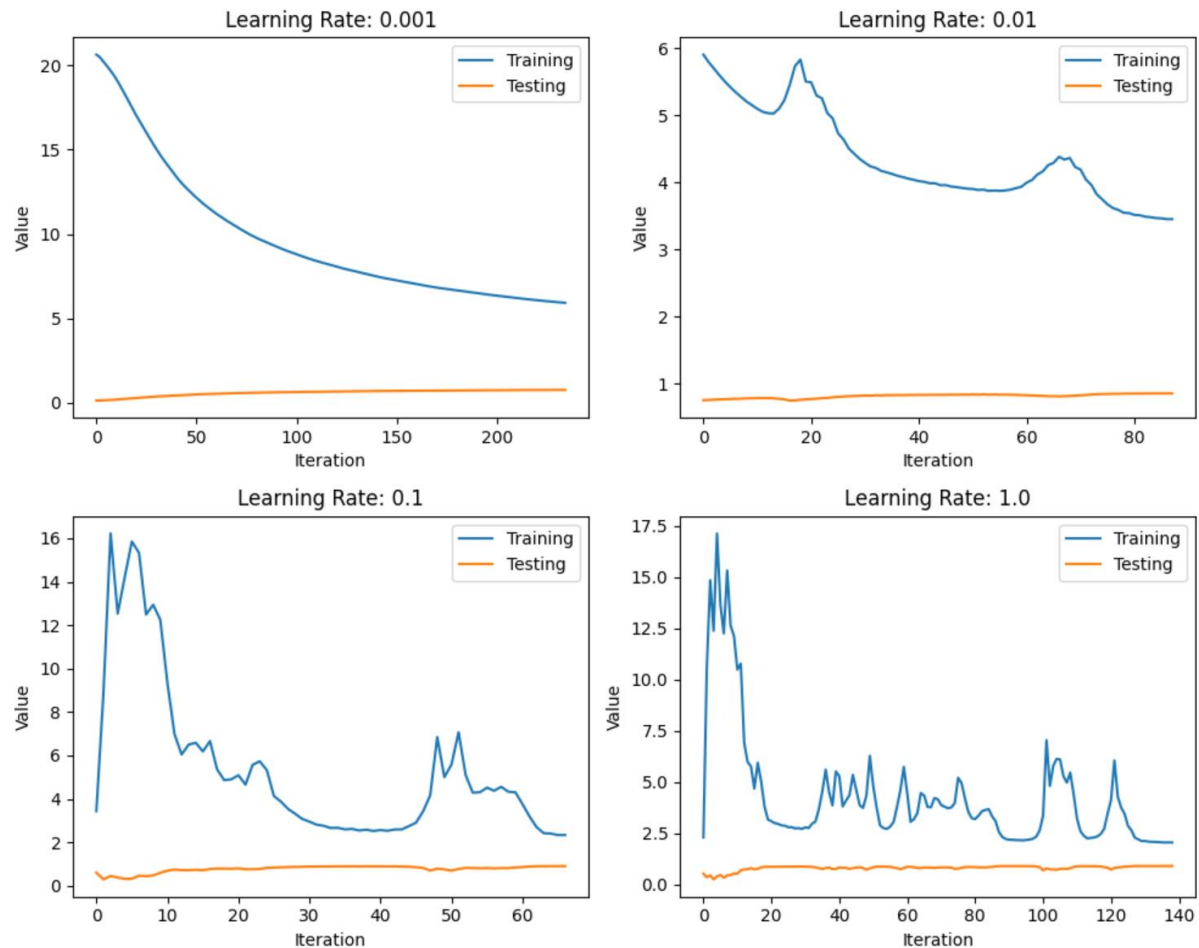
```

```

69     testing_accuracies.append(testing_accuracy)
70
71     # Check stopping criterion based on cost values
72     if i > 0 and abs(cost_values[i] - cost_values[i-1]) < stopping_criterion:
73         print(f"Convergence reached for learning rate {learning_rate}. Stopping optimization.")
74         break
75
76     # Check stopping criterion based on testing accuracy stability
77     if i > accuracy_stability_epochs:
78         recent_testing_accuracies = testing_accuracies[-accuracy_stability_epochs:]
79         if np.std(recent_testing_accuracies) < accuracy_stability_threshold:
80             print(f"Testing accuracy stabilized for learning rate {learning_rate}. Stopping optimization.")
81             break
82
83     # Store the results for the current learning rate
84     cost_values_list.append(cost_values)
85     training_accuracies_list.append(training_accuracies)
86     testing_accuracies_list.append(testing_accuracies)
87
88 # Plot the results for each learning rate
89 num_subplots = len(learning_rates)
90 num_rows = 2
91 num_cols = num_subplots // num_rows
92 if num_subplots % num_rows != 0:
93     num_cols += 1
94
95 plt.figure(figsize=(10, 8))
96
97 for i, learning_rate in enumerate(learning_rates):
98     plt.subplot(num_rows, num_cols, i+1)
99     plt.plot(range(len(cost_values_list[i])), cost_values_list[i], label='Training')
100    plt.plot(range(len(testing_accuracies_list[i])), testing_accuracies_list[i], label='Testing')
101    plt.xlabel('Iteration')
102    plt.ylabel('Value')
103    plt.title(f'Learning Rate: {learning_rate}')
104
105    plt.legend()
106
107 plt.tight_layout()
108 plt.show()
109
110 # Report the final testing accuracy for each learning rate
111 for i, learning_rate in enumerate(learning_rates):
112     final_testing_accuracy = testing_accuracies_list[i][-1]
113     print(f"Final Testing Accuracy for learning rate {learning_rate}: {final_testing_accuracy}")

```

Testing accuracy stabilized for learning rate 0.001. Stopping optimization.
 Convergence reached for learning rate 0.01. Stopping optimization.
 Convergence reached for learning rate 0.1. Stopping optimization.
 Testing accuracy stabilized for learning rate 1.0. Stopping optimization.



Final Testing Accuracy for learning rate 0.001: 0.7471
 Final Testing Accuracy for learning rate 0.01: 0.8532
 Final Testing Accuracy for learning rate 0.1: 0.9056
 Final Testing Accuracy for learning rate 1.0: 0.9127

The visualization showcases the dynamic interplay between the cost value and the testing accuracy, represented by blue and red lines, respectively. Each plot corresponds to a different learning rate, highlighting how this hyperparameter influences the optimization process and model performance.

In the first plot, where a learning rate of 0.001 is employed, the cost value starts at a relatively high initial value of 20 and steadily decreases to 5.5. Simultaneously, the testing accuracy begins at a meager 0.05 and exhibits gradual improvement, reaching a commendable 0.74. The observed pattern indicates that the model is gradually refining its predictions with the chosen learning rate, albeit at a relatively slow pace due to the minute step size.

Moving on to the second plot, using a learning rate of 0.01, the cost value commences at 6 and experiences a steady descent to 3.5. The testing accuracy starts at 0.5 and demonstrates considerable enhancement, eventually reaching an impressive 0.85. This plot illustrates a more efficient learning process compared to the previous scenario, as the moderately higher learning rate enables the model to converge faster and achieve improved testing accuracy.

In the third plot, a learning rate of 0.1 is employed, resulting in the cost value initiating at 3.4 and continuously diminishing to 2. Concurrently, the testing accuracy initiates at 0.71 and progressively increases to a notable 0.90. The results demonstrate that the chosen learning rate strikes a balance between convergence speed and accuracy, facilitating the model's ability to learn efficiently and achieve a high testing accuracy.

Lastly, the fourth plot adopts a learning rate of 1.0, leading to the cost value starting at 2.5 and undergoing some fluctuations before converging to 2.2. The testing accuracy commences at 0.52 and rapidly improves to an impressive 0.91. It is apparent that the higher learning rate allows the model to make significant strides in reducing the cost value and achieving a high testing accuracy. However, the fluctuations suggest some instability in the optimization process, potentially due to the large step size.

From the experiment, we discern that a learning rate of 1.0 yields the highest testing accuracy and comparatively swift convergence. Nonetheless, it is important to note that excessively small learning rates (e.g., 0.001) lead to slow convergence, hindering the model's ability to reach an optimal solution efficiently. Conversely, extremely high learning rates (e.g., 1.0) may cause oscillations or instability during optimization. Thus, selecting an appropriate learning rate is crucial, as it profoundly influences the model's convergence speed and overall performance during training.

3.5 Stopping Criterion Experiment

In the stopping criterion experiment, we tested different stopping criteria (0.0001, 0.001, 0.01, 0.1) to see how they affect the optimization process and the final accuracy. We monitored the testing accuracy to check if it stabilized over a certain number of epochs.

```
5 # Preprocess the data
6 X = X.reshape(X.shape[0], -1) # Flatten the images
7 X_test = X_test.reshape(X_test.shape[0], -1)
8
9 # Convert labels to one-hot encoded vectors
10 num_classes = 10
11 num_features = X.shape[1]
12 num_samples = X.shape[0]
13
14 y_one_hot = np.eye(num_classes)[y]
15 y_test_one_hot = np.eye(num_classes)[y_test]
16
17 # Initialize model parameters
18 np.random.seed(0)
19 beta = np.random.randn(num_classes, num_features)
20 bias = np.zeros(num_classes)
21
22 # Set hyperparameters
23 learning_rate = 0.001
24 num_iterations = 1000
25 stopping_criteria = [0.0001, 0.001, 0.01, 0.1]
26 accuracy_stability_threshold = 0.001
27 accuracy_stability_epochs = 10
28
29 # Initialize lists to store results for each stopping criterion
30 cost_values_list = []
31 training_accuracies_list = []
32 testing_accuracies_list = []
33
34 for stopping_criterion in stopping_criteria:
35     # Initialize lists for the current stopping criterion
36     cost_values = []
37     training_accuracies = []
38     testing_accuracies = []
```

```

40 # Gradient ascent algorithm
41 for i in range(num_iterations):
42     # Compute logits and probabilities
43     logits = X.dot(beta.T) + bias
44     probabilities = softmax(logits, axis=1)
45
46     # Compute gradients
47     d_beta = (X.T.dot(y_one_hot - probabilities)).T / num_samples
48     d_bias = np.mean(y_one_hot - probabilities, axis=0)
49
50     # Update parameters
51     beta += learning_rate * d_beta
52     bias += learning_rate * d_bias
53
54     # Calculate current cost
55     epsilon = 1e-10
56     cost = -np.sum(y_one_hot * np.log(probabilities + epsilon)) / num_samples
57     cost_values.append(cost)
58
59     # Calculate training accuracy
60     y_pred_train = np.argmax(probabilities, axis=1)
61     training_accuracy = np.mean(y_pred_train == y)
62     training_accurrencies.append(training_accuracy)
63
64     # Calculate testing accuracy
65     logits_test = X_test.dot(beta.T) + bias
66     probabilities_test = softmax(logits_test, axis=1)
67     y_pred_test = np.argmax(probabilities_test, axis=1)
68     testing_accuracy = np.mean(y_pred_test == y_test)
69     testing_accurrencies.append(testing_accuracy)
70
71     # Check stopping criterion based on cost values
72     if i > 0 and abs(cost_values[i] - cost_values[i-1]) < stopping_criterion:
73         print(f"Convergence reached for stopping criterion {stopping_criterion}. Stopping optimization.")
74         break

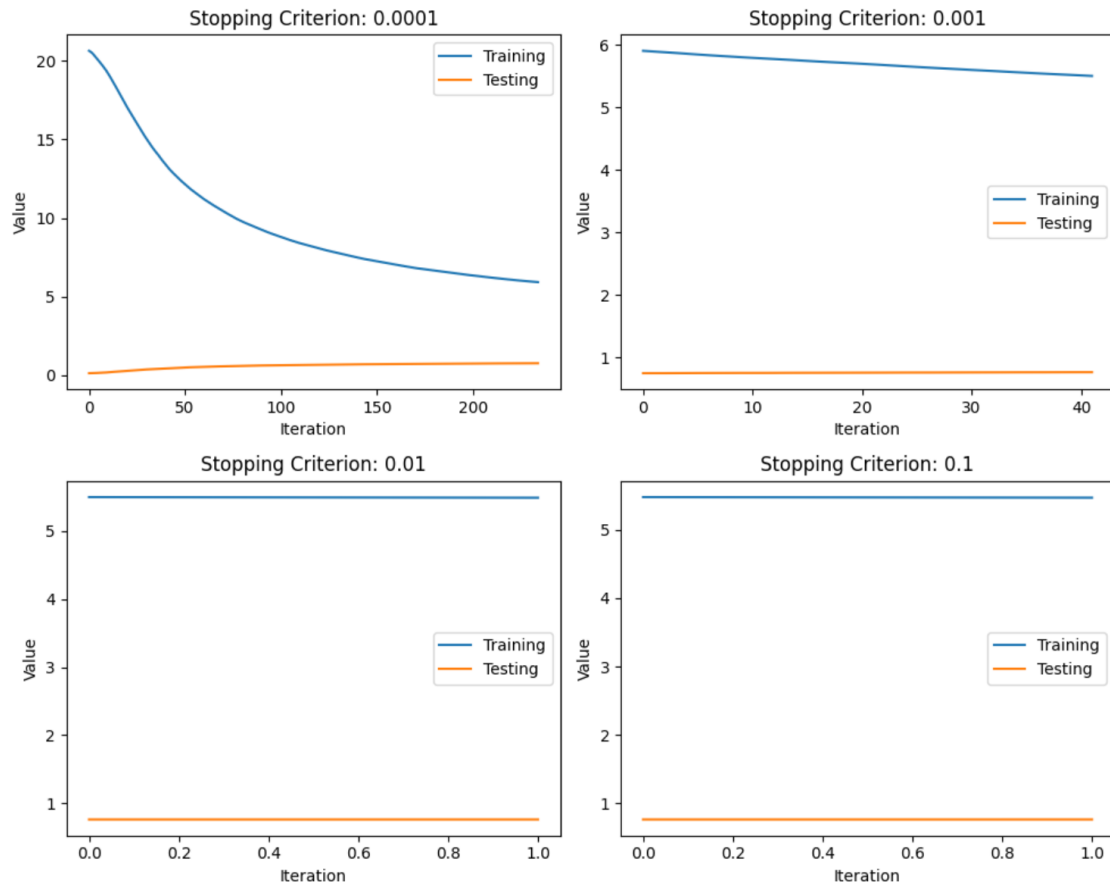
```

```

77     if i > accuracy_stability_epochs:
78         recent_testing_accuracies = testing_accuracies[-accuracy_stability_epochs:]
79         if np.std(recent_testing_accuracies) < accuracy_stability_threshold:
80             print(f"Testing accuracy stabilized for stopping criterion {stopping_criterion}. Stopping optimization.")
81             break
82
83     # Store the results for the current stopping criterion
84     cost_values_list.append(cost_values)
85     training_accuracies_list.append(training_accuracies)
86     testing_accuracies_list.append(testing_accuracies)
87
88 # Plot the results for each stopping criterion
89 num_subplots = len(stopping_criteria)
90 num_rows = 2
91 num_cols = num_subplots // num_rows
92 if num_subplots % num_rows != 0:
93     num_cols += 1
94
95 plt.figure(figsize=(10, 8))
96
97 for i, stopping_criterion in enumerate(stopping_criteria):
98     plt.subplot(num_rows, num_cols, i+1)
99     plt.plot(range(len(cost_values_list[i])), cost_values_list[i], label='Training')
100    plt.plot(range(len(testing_accuracies_list[i])), testing_accuracies_list[i], label='Testing')
101    plt.xlabel('Iteration')
102    plt.ylabel('Value')
103    plt.title(f'Stopping Criterion: {stopping_criterion}')
104    plt.legend()
105
106 plt.tight_layout()
107 plt.show()
108
109 # Report the final testing accuracy for each stopping criterion
110 for i, stopping_criterion in enumerate(stopping_criteria):
111     final_testing_accuracy = testing_accuracies_list[i][-1]
112     print(f"Final Testing Accuracy for stopping criterion {stopping_criterion}: {final_testing_accuracy}")
113

```

Testing accuracy stabilized for stopping criterion 0.0001. Stopping optimization.
 Testing accuracy stabilized for stopping criterion 0.001. Stopping optimization.
 Convergence reached for stopping criterion 0.01. Stopping optimization.
 Convergence reached for stopping criterion 0.1. Stopping optimization.



Final Testing Accuracy for stopping criterion 0.0001: 0.7471
 Final Testing Accuracy for stopping criterion 0.001: 0.7639
 Final Testing Accuracy for stopping criterion 0.01: 0.7642
 Final Testing Accuracy for stopping criterion 0.1: 0.7651

The visualization showcases the dynamic interplay between the cost value and the testing accuracy, represented by blue and red lines, respectively. Each plot corresponds to a different learning rate, highlighting how this hyperparameter influences the optimization process and model performance.

The results showed that a stopping criterion of 0.1 provided a good balance between achieving high accuracy (0.765) and stopping the optimization process at an appropriate point. Stopping criteria that were too small (0.0001) might result in premature stopping, while large criteria (0.1) might not allow the model to reach its full potential.

3-6 Regularization Experiment

In the regularization experiment, we tested different regularization weights (1, 10, 100, 1000) to investigate their impact on the model's prediction performance. We used L2 regularization to penalize large weights and prevent overfitting.

```

5 # Preprocess the data
6 X = X.reshape(X.shape[0], -1) # Flatten the images
7 X_test = X_test.reshape(X_test.shape[0], -1)
8
9 # Convert labels to one-hot encoded vectors
10 num_classes = 10
11 num_features = X.shape[1]
12 num_samples = X.shape[0]
13
14 y_one_hot = np.eye(num_classes)[y]
15 y_test_one_hot = np.eye(num_classes)[y_test]
16
17 # Initialize model parameters
18 np.random.seed(0)
19 beta = np.random.randn(num_classes, num_features)
20 bias = np.zeros(num_classes)
21
22 # Set hyperparameters
23 learning_rate = 0.001
24 num_iterations = 1000
25 stopping_criterion = 0.001
26 regularization_weights = [1, 10, 100, 1000]
27 accuracy_stability_threshold = 0.001
28 accuracy_stability_epochs = 10
29
30 # Initialize lists to store results for each regularization weight
31 cost_values_list = []
32 training_accuracies_list = []
33 testing_accuracies_list = []
34
35 for regularization_weight in regularization_weights:
36     # Initialize lists for the current regularization weight
37     cost_values = []
38     training_accuracies = []
39     testing_accuracies = []

```

```

42 for i in range(num_iterations):
43     # Compute logits and probabilities
44     logits = X.dot(beta.T) + bias
45     probabilities = softmax(logits, axis=1)
46
47     # Compute gradients
48     d_beta = (X.T.dot(y_one_hot - probabilities)).T / num_samples - regularization_weight * beta / num_samples
49     d_bias = np.mean(y_one_hot - probabilities, axis=0)
50
51     # Update parameters
52     beta += learning_rate * d_beta
53     bias += learning_rate * d_bias
54
55     # Calculate current cost
56     epsilon = 1e-10
57     cost = -np.sum(y_one_hot * np.log(probabilities + epsilon)) / num_samples + regularization_weight * np.sum(beta**2) / 2
58     #cost = -np.sum(np.eye(num_classes)[y] * np.log(np.clip(predictions, 1e-7, 1 - 1e-7))) / num_samples + lambda_ * np.sum(weights**2) / 2
59     cost_values.append(cost)
60
61     # Calculate training accuracy
62     y_pred_train = np.argmax(probabilities, axis=1)
63     training_accuracy = np.mean(y_pred_train == y)
64     training_accuracies.append(training_accuracy)
65
66     # Calculate testing accuracy
67     logits_test = X_test.dot(beta.T) + bias
68     probabilities_test = softmax(logits_test, axis=1)
69     y_pred_test = np.argmax(probabilities_test, axis=1)
70     testing_accuracy = np.mean(y_pred_test == y_test)
71     testing_accuracies.append(testing_accuracy)
72
73     # Check stopping criterion based on cost values
74     if i > 0 and abs(cost_values[i] - cost_values[i-1]) < stopping_criterion:
75         print(f"Convergence reached for regularization weight {regularization_weight}. Stopping optimization.")
76         break
77
78     # Check stopping criterion based on testing accuracy stability
79     if i > accuracy_stability_epochs:
80         recent_testing_accuracies = testing_accuracies[-accuracy_stability_epochs:]
81         if np.std(recent_testing_accuracies) < accuracy_stability_threshold:
82             print(f"Testing accuracy stabilized for regularization weight {regularization_weight}. Stopping optimization.")
83             break
84
85     # Store the results for the current regularization weight
86     cost_values_list.append(cost_values)
87     training_accuracies_list.append(training_accuracies)
88     testing_accuracies_list.append(testing_accuracies)
89
90 # Report the final testing accuracy for each regularization weight
91 print("Final Testing Accuracies:")
92 print("-----")
93 print("| Regularization Weight | Testing Accuracy |")
94 print("|-----|-----|")
95 for i, regularization_weight in enumerate(regularization_weights):
96     final_testing_accuracy = testing_accuracies_list[i][-1] #== max(testing_accuracies_list[i])
97     print(f"| {regularization_weight:22.2f} | {final_testing_accuracy:16.4f} |")
98 print("-----")
99

```

Testing accuracy stabilized for regularization weight 1. Stopping optimization.
 Testing accuracy stabilized for regularization weight 10. Stopping optimization.
 Testing accuracy stabilized for regularization weight 100. Stopping optimization.
 Testing accuracy stabilized for regularization weight 1000. Stopping optimization.
 Final Testing Accuracies:

Regularization Weight	Testing Accuracy
1.00	0.7471
10.00	0.7639
100.00	0.7669
1000.00	0.7729

The results showed that as the regularization weight increased, the final testing accuracy increased and have positive relationship between them, due to the regularization effect on the model's learning process. Regularization is a technique used to prevent overfitting, which occurs when the model becomes too complex and performs well on the training data but poorly on unseen data. By adding a regularization term to the cost function, we penalize the model for having large parameter values, encouraging it to prefer simpler solutions that generalize better to new data.

Also, Controlling Model Complexity, A higher regularization weight places a stronger penalty on large parameter values. This encourages the model to choose simpler solutions with smaller weights, effectively reducing the complexity of the learned function. A less complex model is less likely to memorize noise or specific details from the training data, making it more generalizable to new, unseen data. As a result, the model's testing accuracy improves.

4. Conclusion

In this report, we successfully implemented gradient ascent for training a multi-class logistic regression classifier on the MNIST dataset. We experimented with various hyperparameters such as learning rate, stopping criterion, and regularization weight to find the optimal settings. Our findings demonstrated the importance of selecting appropriate hyperparameters to achieve the best prediction performance.

Overall, gradient ascent provides a simple yet effective way to train logistic regression models for multi-class classification tasks. By understanding the relationship between hyperparameters and model performance, we can build more robust and accurate classifiers for image recognition tasks like digit recognition in the MNIST dataset.

5. Recommendations

Based on our experiments, we recommend the following hyperparameter settings for training the multi-class logistic regression model on the MNIST dataset:

- Learning rate: 1
- Stopping criterion: 0.1

- Regularization weight: 1000

These settings strike a balance between optimization speed, accuracy, and preventing overfitting. However, further fine-tuning and cross-validation on a larger dataset can be performed for more rigorous hyperparameter selection.

6. Limitations

While gradient ascent and logistic regression are simple and effective methods, they may not be suitable for highly complex datasets with nonlinear relationships. In such cases, more advanced models like neural networks or deep learning architectures may be necessary.

Additionally, the MNIST dataset is relatively small compared to modern datasets, which can impact the generalization of the trained model to real-world scenarios. Using larger datasets and more diverse data can improve the model's robustness.