

Assignment-1

Ajoy Kumar (S24014) Premchand (S24030) Prankshu (D24144)
Shubham Sharma (S24031) Shubham Thakur (S24065)
Nalin Dhiman (D24008)

CS671: Deep Learning and Applications (JAN-JUNE 2025)

February 16, 2025

Abstract

This report details the implementation of a neural network for MNIST digit classification. The network comprises a single hidden layer with ReLU activation and a softmax output layer. The training is performed using backpropagation with stochastic gradient descent. We analyze the model's performance using accuracy, loss, and F1-score.

1 Introduction

The perceptron model is one of the foundational concepts in machine learning for classification tasks. It uses a linear decision boundary to separate different classes. This report uses a perceptron with the one-against-one (OvO) approach to classify the given datasets.

2 Mathematical Formulation

2.1 Perceptron Model

Given an input vector $\mathbf{x} = [x_1, x_2, \dots, x_n]$, the perceptron computes a weighted sum:

$$z = \mathbf{W}^T \mathbf{x} + b \tag{1}$$

where \mathbf{W} is the weight vector, and b is the bias term.

2.2 Activation Functions

We use two activation functions:

- **Step Function:** $f(z) = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases}$
- **Sigmoid Function:** $f(z) = \frac{1}{1+e^{-z}}$

2.3 Error Calculation and Weight Update

The error is computed as:

$$e = y - \hat{y} \quad (2)$$

where y is the true label, and \hat{y} is the predicted output.

For backpropagation using sigmoid activation, the weight update is:

$$\Delta W = \eta e \hat{y} (1 - \hat{y}) X \quad (3)$$

where η is the learning rate.

For the step function, the update follows the perceptron learning rule:

$$W \leftarrow W + \eta e X, \quad b \leftarrow b + \eta e \quad (4)$$

2.4 Decision Boundary

For two classes, the decision boundary is given by:

$$W_1 x_1 + W_2 x_2 + b = 0 \quad (5)$$

This equation represents a straight line separating the classes in a 2D feature space.

3 Experimental Setup

3.1 Datasets

- **LS Dataset:** 3 classes, 500 points each, linearly separable.
- **NLS Dataset:** 3 classes, non-linearly separable, 2000 points in total.

Each dataset is split into 70% training and 30% testing.

3.2 Evaluation Metrics

The performance of the model is assessed using:

- Accuracy: $\frac{\text{Correct Predictions}}{\text{Total Predictions}}$
- Precision: $\frac{TP}{TP+FP}$
- Recall: $\frac{TP}{TP+FN}$
- F1-Score: $2 \cdot \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$

4 Results and Discussion

4.1 Training Performance

Figure 1 shows the error vs. epochs plot for the training phase.

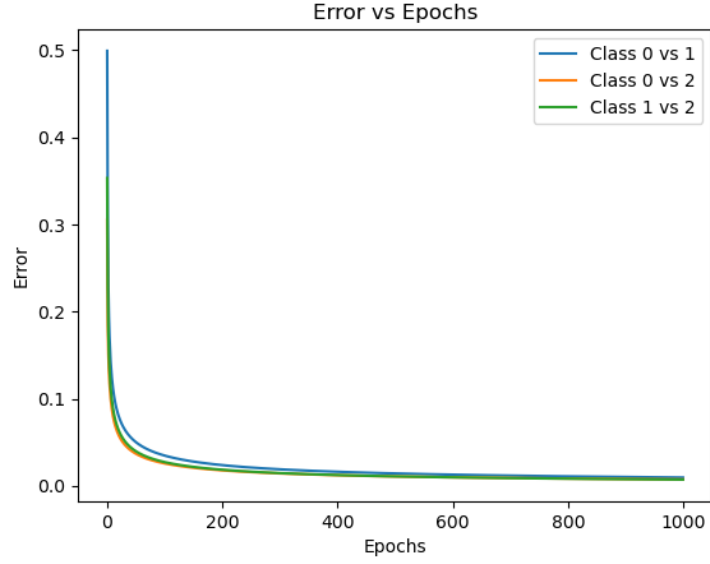


Figure 1: Error vs. Epochs during training for LS Dataset.

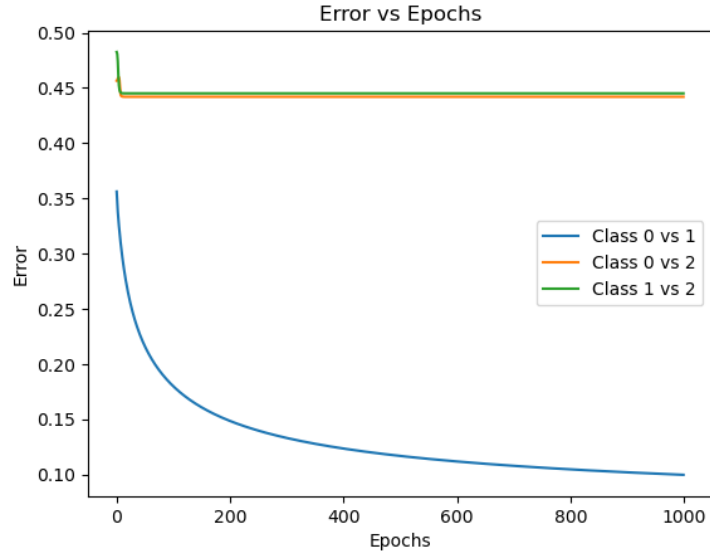


Figure 2: Error vs. Epochs during training for NLS Dataset.

4.2 Confusion Matrix and Classification Report

Table 1 shows the confusion matrix for the NLS dataset.

True/Predicted	Class 0	Class 1	Class 2
Class 0	0	0	149
Class 1	0	0	144
Class 2	0	0	307

Table 1: Confusion matrix for NLS dataset.

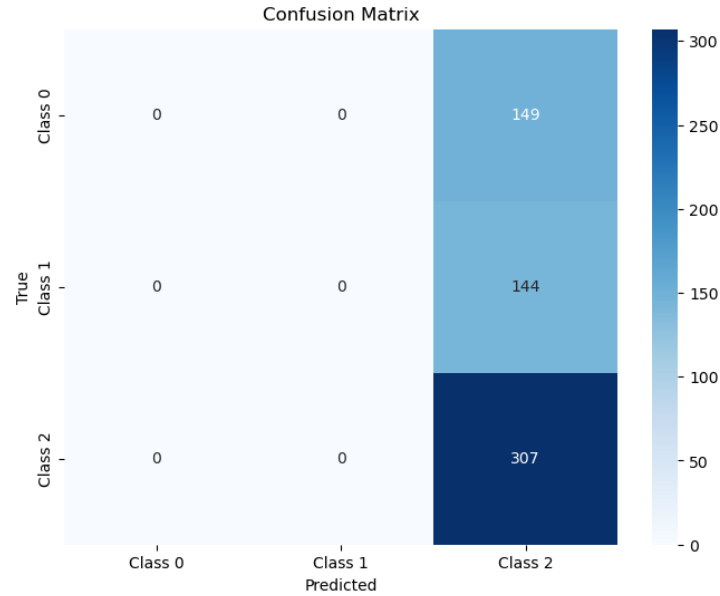


Figure 3: Confusion matrix visualization for NLS dataset.

- Accuracy: 0.5117
- Precision: 0.2618
- Recall: 0.5117
- F1-Score: 0.3464

Classification Report:

	precision	recall	f1-score	support
Class 0	0.00	0.00	0.00	149
Class 1	0.00	0.00	0.00	144
Class 2	0.51	1.00	0.68	307
accuracy			0.51	600
macro avg	0.17	0.33	0.23	600
weighted avg	0.26	0.51	0.35	600

Table 2 shows the confusion matrix for the LS dataset.

True/Predicted	Class 0	Class 1	Class 2
Class 0	152	0	0
Class 1	0	157	0
Class 2	0	0	141

Table 2: Confusion matrix for LS dataset.

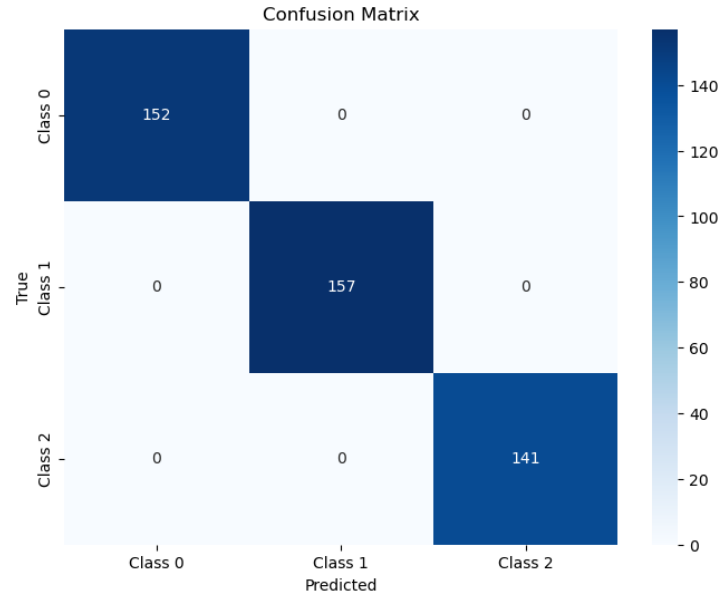


Figure 4: Confusion matrix visualization for LS dataset.

- Accuracy: 1.0000
- Precision: 1.0000
- Recall: 1.0000
- F1-Score: 1.0000

Classification Report:

	precision	recall	f1-score	support
Class 0	1.00	1.00	1.00	152
Class 1	1.00	1.00	1.00	157
Class 2	1.00	1.00	1.00	141
accuracy			1.00	450
macro avg	1.00	1.00	1.00	450
weighted avg	1.00	1.00	1.00	450

5 Conclusions

- The perceptron model performs well on the LS dataset but struggles on the NLS dataset due to its linear decision boundary.
- The choice of activation function influences convergence; a sigmoid provides smoother updates than a step function.
- The confusion matrix and classification report indicates the model predominantly predicts a single class in the NLS dataset, highlighting the limitations of a simple perceptron for non-linearly separable data.

Part_1_classification

February 14, 2025

1 Group:36

2 PART-I : Classification Tasks

2.0.1 Dataset 1:

Linearly separable (LS) dataset: 3 classes, 2-dimensional linearly separable data is given. Each class has 500 data points. `###` Dataset 2: Nonlinearly separable (NLS) classes: 2-dimensional data of 3 classes that are non linearly separable is given. The number of examples in each class and their order is given at the beginning of each file. Divide the data from each class into training, and test data. From each class, train, and test split should be 70% and 30% respectively. `###` Model: Use Perceptron model with sigmoidal/step activation function for each given dataset. `###` Task: 1. Use the described Perceptron model with a one-against-one approach for the above given classification tasks. 2. Implement backpropagation algorithm from scratch for perceptron learning algorithm for the above given datasets.

```
[40]: import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score,
    recall_score, f1_score, classification_report
from itertools import combinations
```

```
[41]: def load_dataset(file_paths, labels, skip_rows=0):
    data_list = []
    label_list = []

    for file_path, label in zip(file_paths, labels):
        data = np.loadtxt(file_path, skiprows=skip_rows) # Auto-detect
        # delimiter
        data_list.append(data)
        label_list.append(np.full((data.shape[0],), label)) # Assign labels

    X = np.vstack(data_list)
    y = np.hstack(label_list)
    return X, y
```

```
[42]: #NLS dataset
def load_nls_dataset(file_path):
    data = np.loadtxt(file_path, skiprows=1)
    # Assuming the first 500 examples are class 0, next 500 are class 1, and
    ↪ last 1000 are class 2
    labels = np.concatenate([np.zeros(500), np.ones(500), np.full(1000, 2)]) #
    ↪ Class 0, Class 1, Class 2
    return data, labels
```

```
[43]: def normalize(X):
    return (X - np.mean(X, axis=0)) / np.std(X, axis=0)
```

```
[44]: # Split 70/30
def train_test_split(X, y, train_ratio=0.7):
    idx = np.arange(len(y))
    np.random.shuffle(idx)
    split = int(train_ratio * len(y))
    return X[idx[:split]], y[idx[:split]], X[idx[split:]], y[idx[split:]]
```

```
[45]: def sigmoid(z):
    return 1 / (1 + np.exp(-z))
```

```
[46]: def step(z):
    return np.where(z >= 0, 1, 0)
```

```
[47]: def forward(X, W, b, activation):
    z = np.dot(X, W) + b
    return sigmoid(z) if activation == 'sigmoid' else step(z)
```

```
[48]: # Backpropagation
def train_perceptron(X_train, y_train, epochs=1000, lr=0.01,
    ↪ activation='sigmoid'):
    W = np.random.randn(X_train.shape[1])
    b = np.random.randn()
    errors = []

    for epoch in range(epochs):
        output = forward(X_train, W, b, activation)
        error = y_train - output

        #gradients
        if activation == 'sigmoid':
            grad = error * output * (1 - output) # Derivative of sigmoid
        else:
            grad = error # Perceptron update rule

        W += lr * np.dot(X_train.T, grad)
```



```

        b += lr * np.sum(grad)

        errors.append(np.mean(np.abs(error)))

    return W, b, errors

```

```

[49]: # Decision boundary
def plot_decision_boundary(X, y, W, b, title):
    x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
    y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
    xx, yy = np.meshgrid(np.linspace(x_min, x_max, 100), np.linspace(y_min,
↪y_max, 100))
    grid = np.c_[xx.ravel(), yy.ravel()]
    Z = forward(grid, W, b, activation='sigmoid').reshape(xx.shape)
    plt.contourf(xx, yy, Z, alpha=0.3)
    plt.scatter(X[:, 0], X[:, 1], c=y, edgecolors='k')
    plt.title(title)
    plt.show()

```

```

[50]: # One-vs-One classification
def ovo_classification(X_train, y_train, X_test, y_test, classes, epochs=1000,
↪lr=0.01, activation='sigmoid'):
    classifiers = {}

    for (class1, class2) in combinations(classes, 2):
        mask = (y_train == class1) | (y_train == class2)
        X_sub, y_sub = X_train[mask], y_train[mask]
        y_sub = np.where(y_sub == class1, 1, 0)

        W, b, errors = train_perceptron(X_sub, y_sub, epochs, lr, activation)
        classifiers[(class1, class2)] = (W, b)

    plt.plot(errors, label=f'Class {class1} vs {class2}')

    plt.xlabel('Epochs')
    plt.ylabel('Error')
    plt.title('Error vs Epochs')
    plt.legend()
    plt.show()

    return classifiers

```

```

[51]: # confusion matrix
def plot_confusion_matrix(y_true, y_pred, classes, title='Confusion Matrix'):
    cm = confusion_matrix(y_true, y_pred)
    plt.figure(figsize=(8, 6))

```

```

sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=classes,
↪yticklabels=classes)
plt.title(title)
plt.xlabel('Predicted')
plt.ylabel('True')
plt.show()

```

```

[52]: # statistics
def compute_advanced_statistics(y_true, y_pred, classes):
    plot_confusion_matrix(y_true, y_pred, classes)

    accuracy = accuracy_score(y_true, y_pred)
    print(f"Accuracy: {accuracy:.4f}")

    precision = precision_score(y_true, y_pred, average='weighted',
↪zero_division=0)
    recall = recall_score(y_true, y_pred, average='weighted', zero_division=0)
    f1 = f1_score(y_true, y_pred, average='weighted', zero_division=0)

    print(f"Precision: {precision:.4f}")
    print(f"Recall: {recall:.4f}")
    print(f"F1-Score: {f1:.4f}")

    print("\nClassification Report:")
    print(classification_report(y_true, y_pred, target_names=classes, labels=np.
↪unique(y_true)))

```

```

[55]: def main():
    dataset_paths_ls = [
        "/home/nalin/codes/DL/DL_ass/data/CS671_Dataset_Assignment1/Dataset-1/
↪LS/Class1.txt",
        "/home/nalin/codes/DL/DL_ass/data/CS671_Dataset_Assignment1/Dataset-1/
↪LS/Class2.txt",
        "/home/nalin/codes/DL/DL_ass/data/CS671_Dataset_Assignment1/Dataset-1/
↪LS/Class3.txt"
    ]
    dataset_path_nls = "/home/nalin/codes/DL/DL_ass/data/
↪CS671_Dataset_Assignment1/Dataset-1/NLS/dataset.txt"

    labels = [0, 1, 2]
    class_names = ['Class 0', 'Class 1', 'Class 2']

    X_ls, y_ls = load_dataset(dataset_paths_ls, labels)
    X_ls = normalize(X_ls)
    X_train_ls, y_train_ls, X_test_ls, y_test_ls = train_test_split(X_ls, y_ls)

```

```

X_nls, y_nls = load_nls_dataset(dataset_path_nls)
X_nls = normalize(X_nls)
X_train_nls, y_train_nls, X_test_nls, y_test_nls = train_test_split(X_nls,
↪y_nls)

    for X_train, y_train, X_test, y_test, dataset in [(X_train_ls, y_train_ls,
↪X_test_ls, y_test_ls, "LS"),
                                                    (X_train_nls,
↪y_train_nls, X_test_nls, y_test_nls, "NLS")]:
        print(f"Training on {dataset} dataset...")
        classifiers = ovo_classification(X_train, y_train, X_test, y_test,
↪labels)

        y_pred = []
        for x in X_test:
            votes = {}
            for (class1, class2), (W, b) in classifiers.items():
                pred = forward(x, W, b, activation='sigmoid')
                winner = class1 if pred >= 0.5 else class2
                votes[winner] = votes.get(winner, 0) + 1
            y_pred.append(max(votes, key=votes.get))

        print("Unique predictions:", np.unique(y_pred))

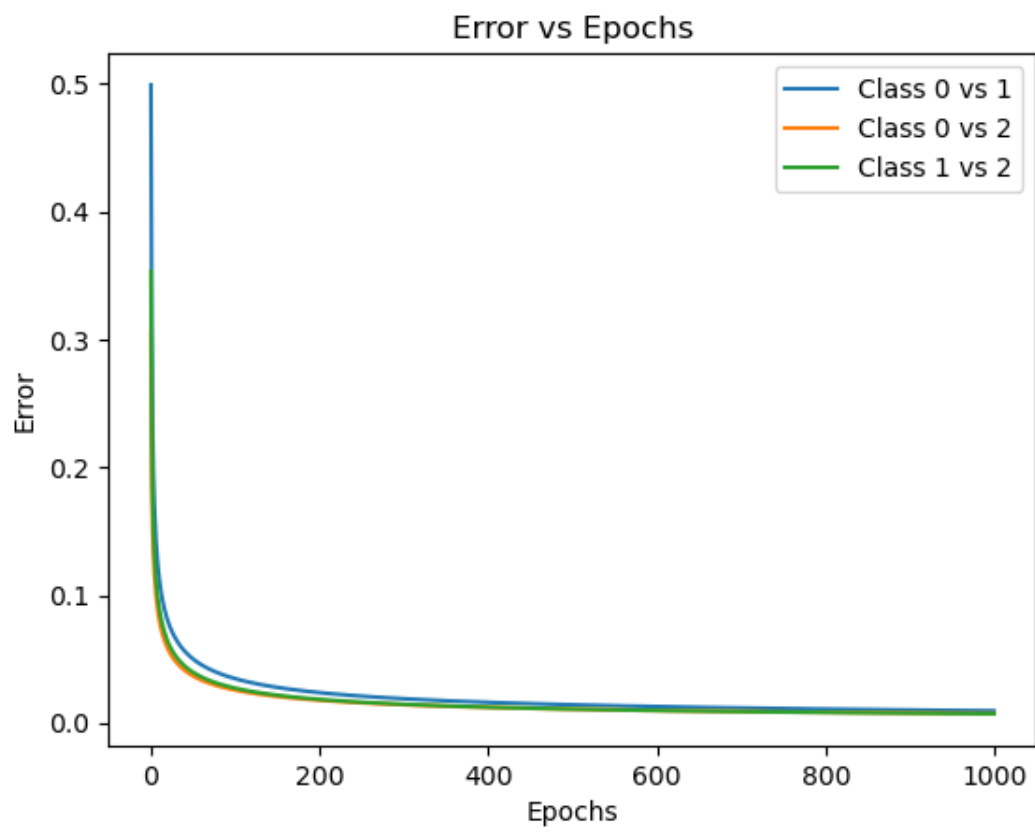
        compute_advanced_statistics(y_test, y_pred, class_names)

        for (class1, class2), (W, b) in classifiers.items():
            plot_decision_boundary(X_train, y_train, W, b, f'{dataset} Decision
↪Boundary: Class {class1} vs {class2}')

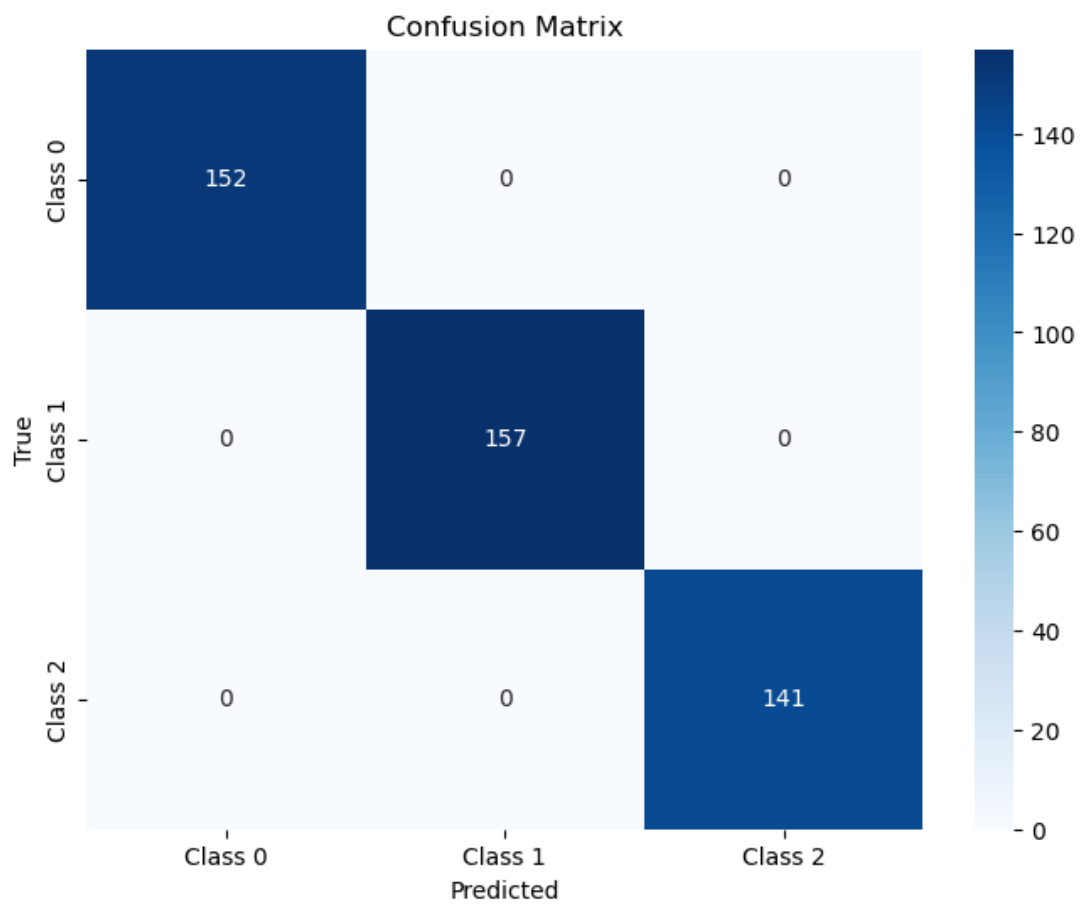
if __name__ == '__main__':
    main()

```

Training on LS dataset...



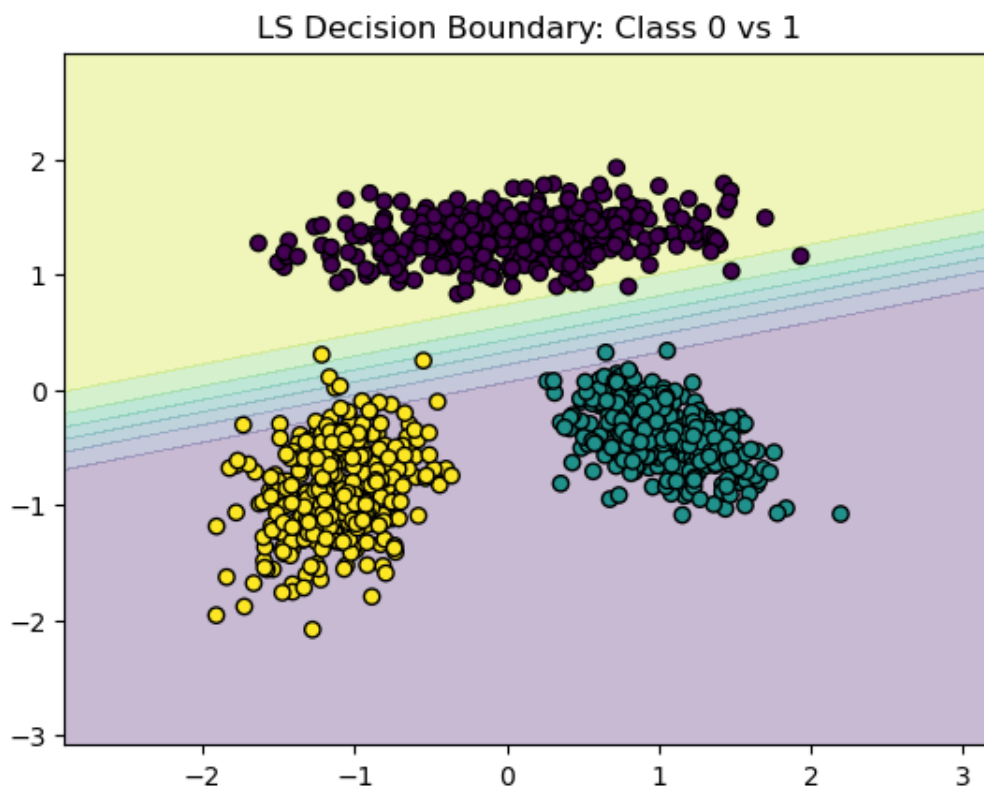
Unique predictions: [0 1 2]

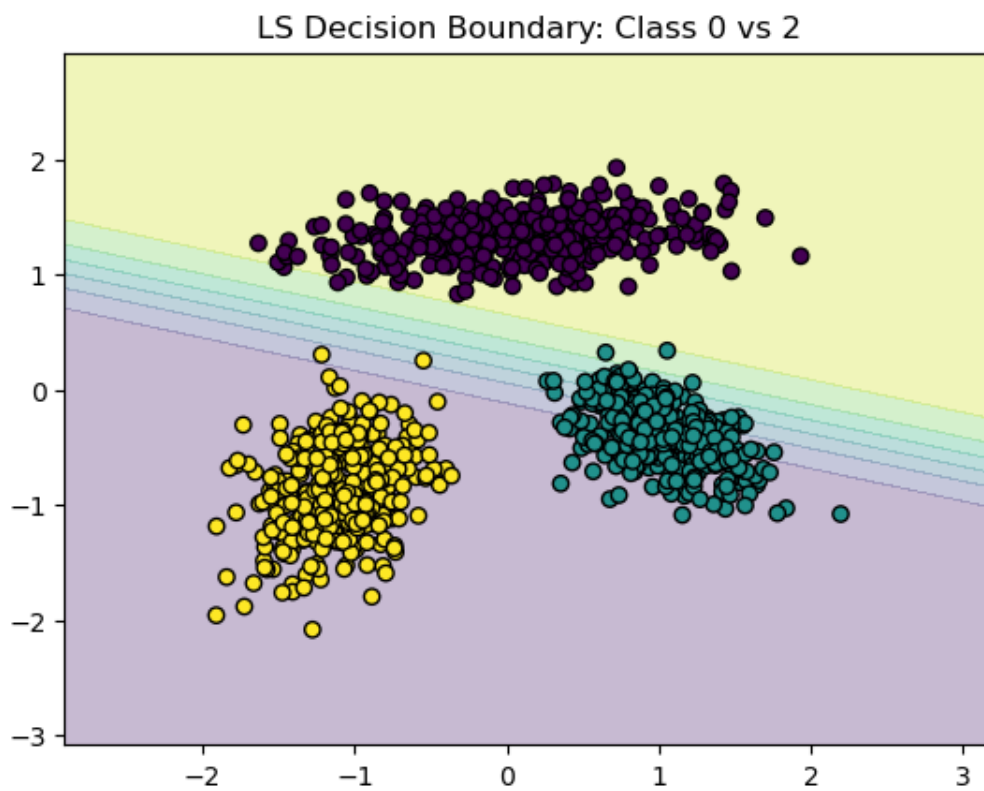


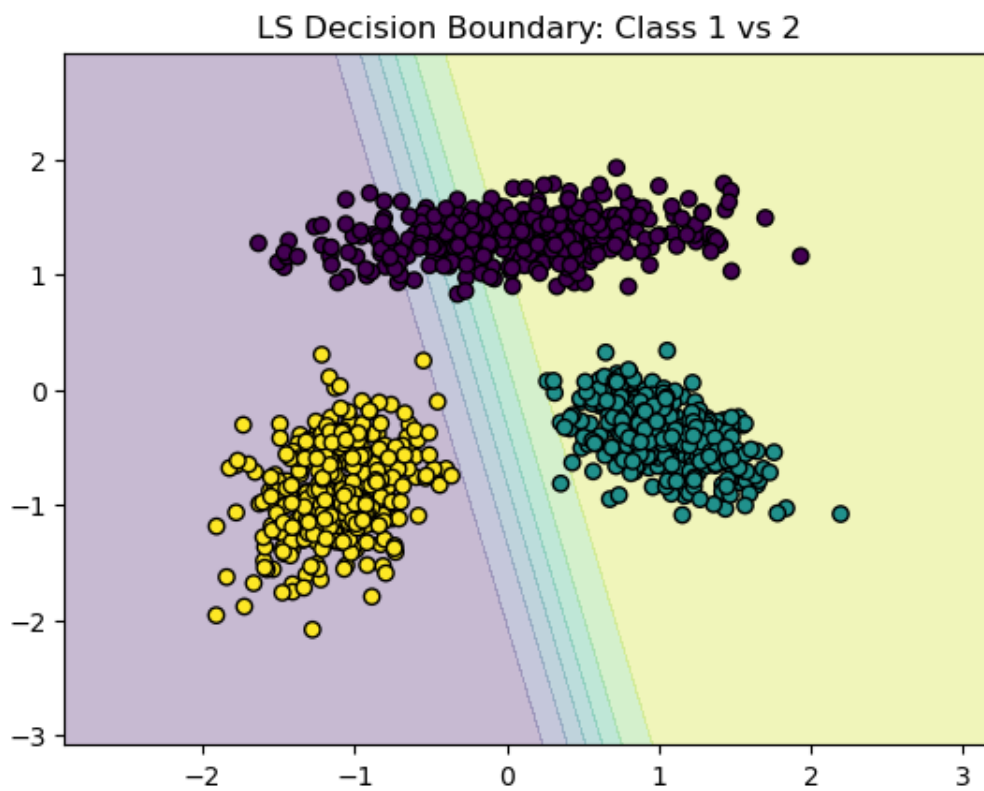
Accuracy: 1.0000
Precision: 1.0000
Recall: 1.0000
F1-Score: 1.0000

Classification Report:

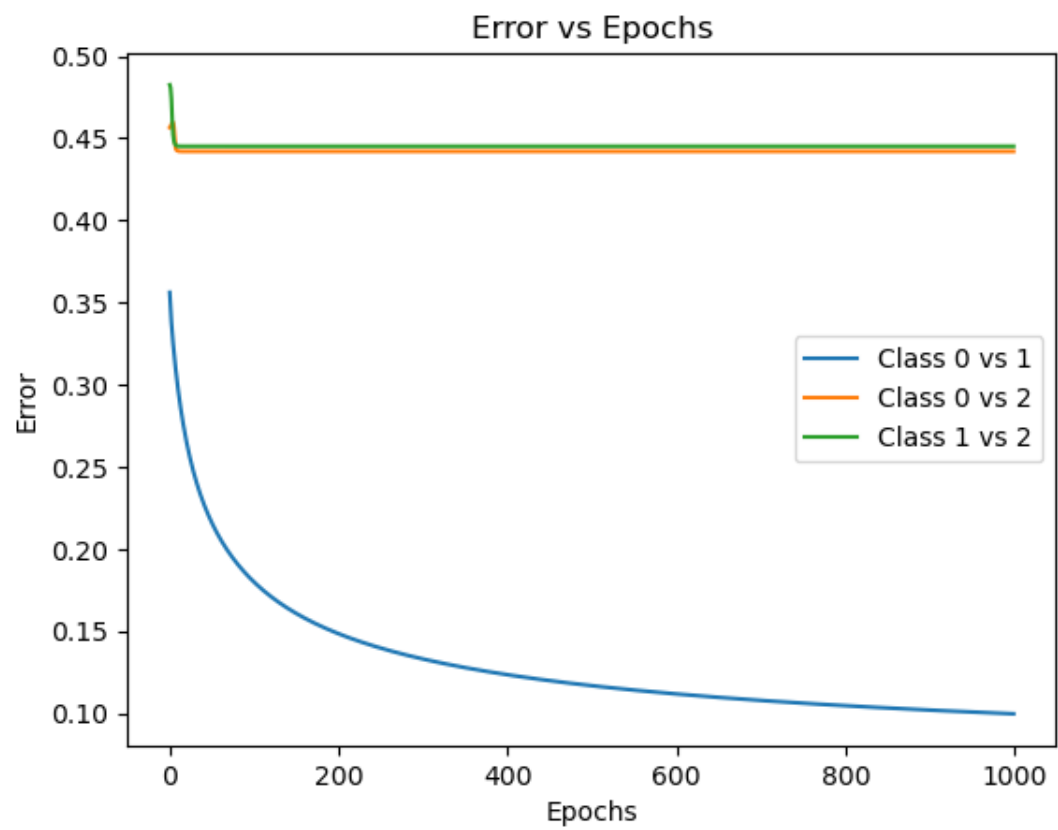
	precision	recall	f1-score	support
Class 0	1.00	1.00	1.00	152
Class 1	1.00	1.00	1.00	157
Class 2	1.00	1.00	1.00	141
accuracy			1.00	450
macro avg	1.00	1.00	1.00	450
weighted avg	1.00	1.00	1.00	450



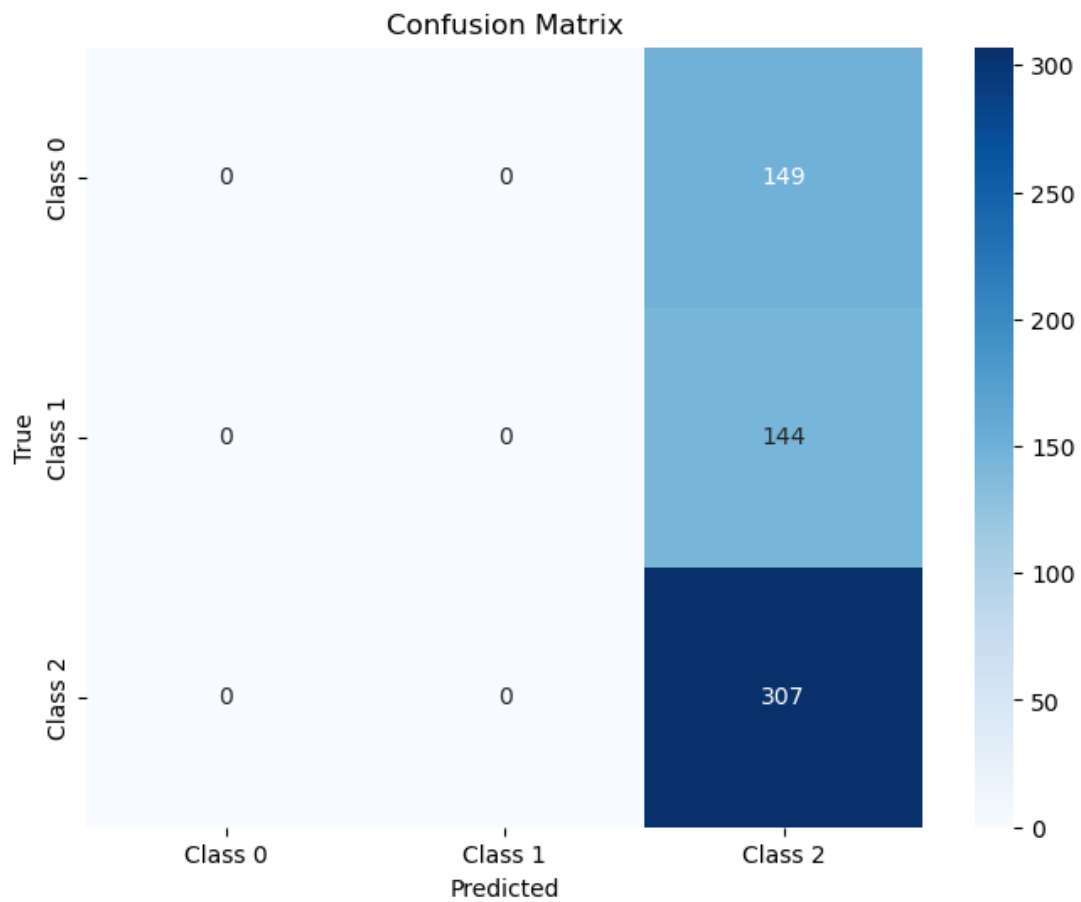




Training on NLS dataset...



Unique predictions: [2]



Accuracy: 0.5117
Precision: 0.2618
Recall: 0.5117
F1-Score: 0.3464

Classification Report:

	precision	recall	f1-score	support
Class 0	0.00	0.00	0.00	149
Class 1	0.00	0.00	0.00	144
Class 2	0.51	1.00	0.68	307
accuracy			0.51	600
macro avg	0.17	0.33	0.23	600
weighted avg	0.26	0.51	0.35	600

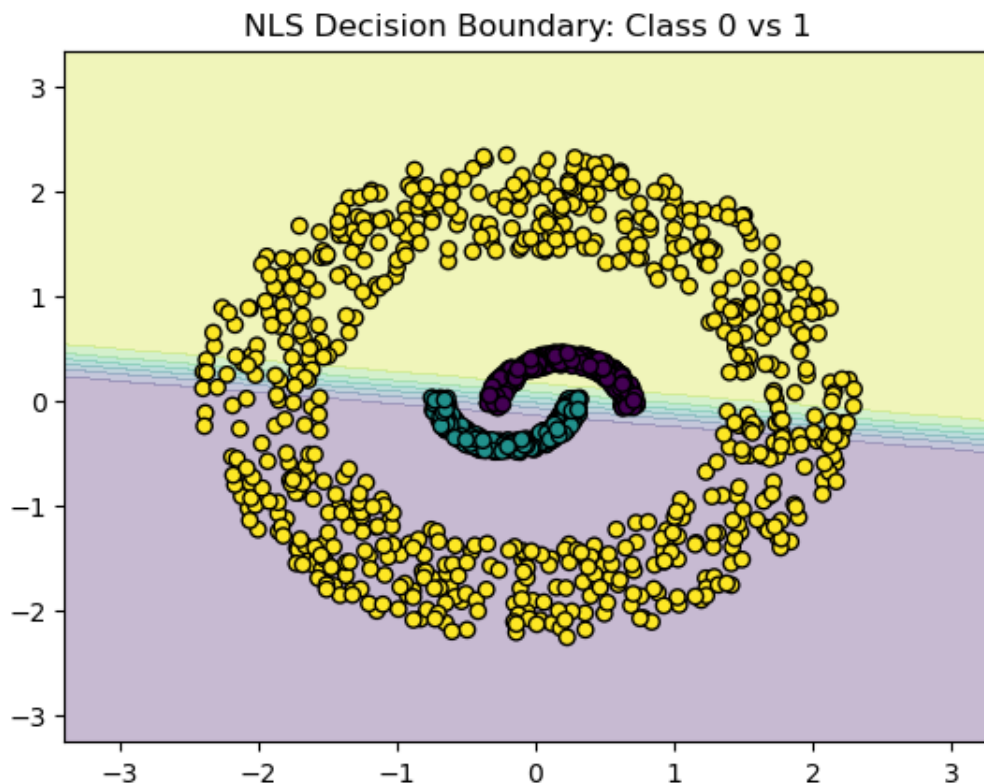
/home/nalin/anaconda3/lib/python3.9/site-packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning:

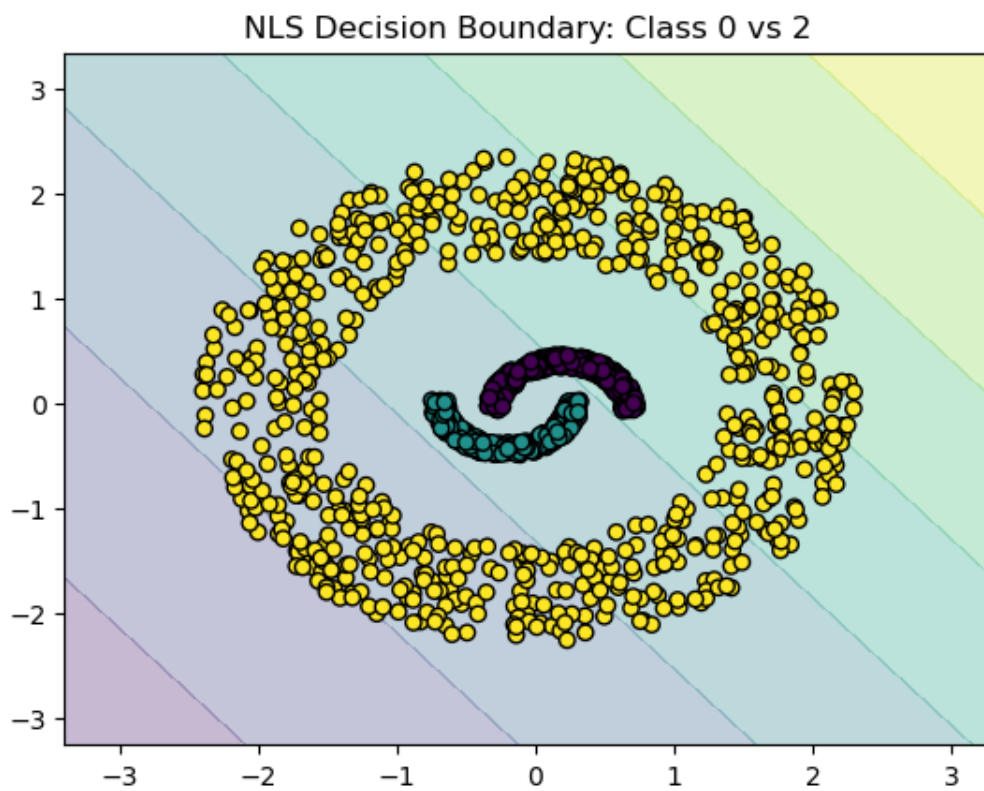
Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_division` parameter to control this behavior.

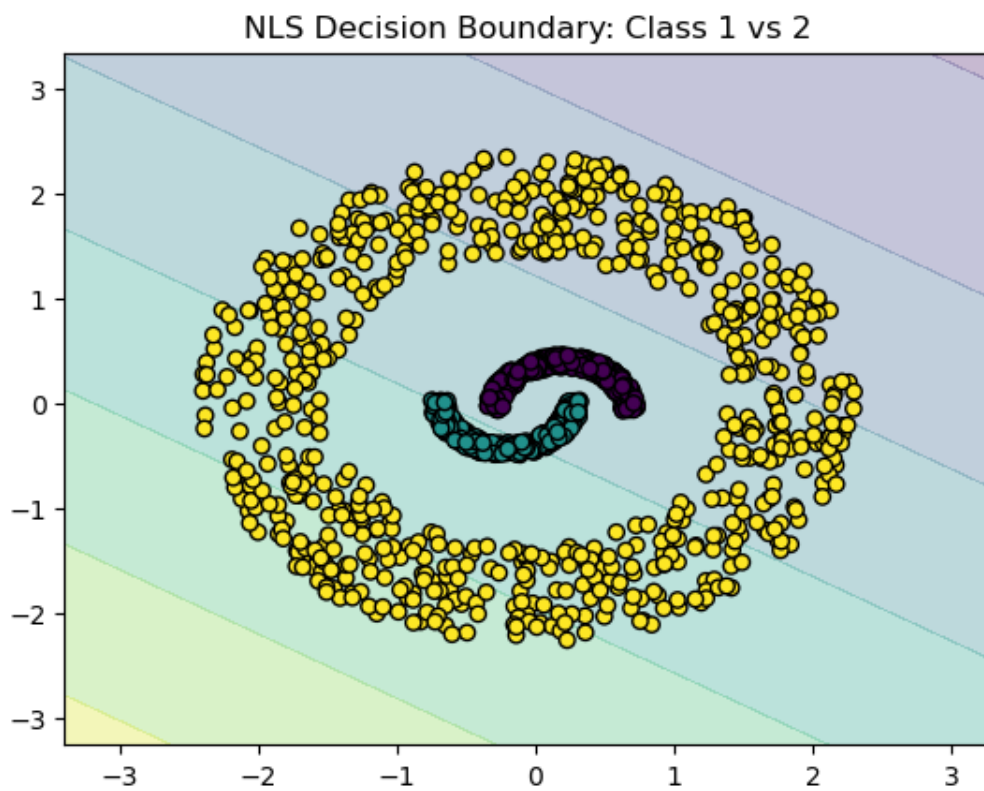
```
_warn_prf(average, modifier, msg_start, len(result))  
/home/nalin/anaconda3/lib/python3.9/site-  
packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning:  
Precision and F-score are ill-defined and being set to 0.0 in labels with no  
predicted samples. Use `zero_division` parameter to control this behavior.
```

```
_warn_prf(average, modifier, msg_start, len(result))  
/home/nalin/anaconda3/lib/python3.9/site-  
packages/sklearn/metrics/_classification.py:1318: UndefinedMetricWarning:  
Precision and F-score are ill-defined and being set to 0.0 in labels with no  
predicted samples. Use `zero_division` parameter to control this behavior.
```

```
_warn_prf(average, modifier, msg_start, len(result))
```







[]:

[]:

[]:

[]:

Part-II: Model 1: Implementation of Neural Network for MNIST Classification

Abstract

This report details the implementation of a neural network for MNIST digit classification. The network consists of a single hidden layer with ReLU activation and a softmax output layer. The training is performed using backpropagation with stochastic gradient descent. We analyze the model's performance using accuracy, loss, and F1-score.

6 Introduction

Neural networks have gained widespread use in image classification tasks due to their ability to learn hierarchical feature representations. In this report, we implement a simple feedforward neural network for MNIST classification.

7 Mathematical Formulation

A feedforward neural network for MNIST classification consists of an input layer, a hidden layer with non-linear activation, and an output layer with softmax activation. The training process involves forward propagation, loss computation, and backpropagation for weight updates.

7.1 Neural Network Model

A neural network maps input features $X \in \mathbb{R}^{m \times n}$ to output probabilities $Y \in \mathbb{R}^{m \times C}$ through weight matrices W and bias terms b . Our network consists of:

- Input layer: 784 neurons (28x28 pixels)
- Hidden layer: 128 neurons with ReLU activation
- Output layer: 10 neurons with softmax activation

7.2 Forward Propagation

Given input X , the hidden layer computes:

$$Z_1 = XW_1 + b_1 \tag{6}$$

$$A_1 = \text{ReLU}(Z_1) = \max(0, Z_1) \tag{7}$$

The output layer computes:

$$Z_2 = A_1W_2 + b_2 \tag{8}$$

$$A_2 = \text{softmax}(Z_2) \tag{9}$$

where the softmax function is given by:

$$\text{softmax}(Z)_i = \frac{e^{Z_i}}{\sum_j e^{Z_j}} \tag{10}$$

Given an input vector $\mathbf{x} \in \mathbb{R}^{784}$ (flattened 28x28 image), the network performs the following transformations:

7.2.1 Hidden Layer

Let $W_1 \in \mathbb{R}^{h \times 784}$ and $b_1 \in \mathbb{R}^h$ be the weight matrix and bias for the hidden layer. The hidden layer pre-activation and activation are:

$$\mathbf{z}_1 = W_1 \mathbf{x} + b_1 \quad (11)$$

$$\mathbf{a}_1 = \text{ReLU}(\mathbf{z}_1) = \max(0, \mathbf{z}_1) \quad (12)$$

7.2.2 Output Layer

Let $W_2 \in \mathbb{R}^{10 \times h}$ and $b_2 \in \mathbb{R}^{10}$ be the weight matrix and bias for the output layer. The logits and softmax output are:

$$\mathbf{z}_2 = W_2 \mathbf{a}_1 + b_2 \quad (13)$$

$$\hat{\mathbf{y}} = \text{softmax}(\mathbf{z}_2), \quad \text{where } \hat{y}_i = \frac{e^{z_{2,i}}}{\sum_{j=1}^{10} e^{z_{2,j}}} \quad (14)$$

7.3 Loss Function

We use the cross-entropy loss:

$$L = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^C y_{i,j} \log(A_{2,i,j}) \quad (15)$$

where $y_{i,j}$ is the true label and $A_{2,i,j}$ is the predicted probability.

7.4 Backpropagation and Weight Update

Gradients are computed using:

$$dZ_2 = A_2 - Y \quad (16)$$

$$dW_2 = \frac{1}{m} A_1^T dZ_2, \quad db_2 = \frac{1}{m} \sum dZ_2 \quad (17)$$

$$dA_1 = dZ_2 W_2^T, \quad dZ_1 = dA_1 \cdot \text{ReLU}'(Z_1) \quad (18)$$

$$dW_1 = \frac{1}{m} X^T dZ_1, \quad db_1 = \frac{1}{m} \sum dZ_1 \quad (19)$$

Weights are updated as:

$$W \leftarrow W - \eta dW, \quad b \leftarrow b - \eta db \quad (20)$$

where η is the learning rate.

We use the cross-entropy loss, defined as:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{10} y_{i,j} \log \hat{y}_{i,j} \quad (21)$$

where $y_{i,j}$ is the one-hot encoded ground truth.

Using stochastic gradient descent (SGD), we compute gradients using the chain rule:

7.4.1 Output Layer Gradients

$$\frac{\partial \mathcal{L}}{\partial W_2} = \frac{1}{N}(\hat{\mathbf{y}} - \mathbf{y})\mathbf{a}_1^T \quad (22)$$

$$\frac{\partial \mathcal{L}}{\partial b_2} = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i) \quad (23)$$

7.4.2 Hidden Layer Gradients

$$\frac{\partial \mathcal{L}}{\partial W_1} = \frac{1}{N}((W_2^T(\hat{\mathbf{y}} - \mathbf{y})) \odot \text{ReLU}'(\mathbf{z}_1))\mathbf{x}^T \quad (24)$$

$$\frac{\partial \mathcal{L}}{\partial b_1} = \frac{1}{N} \sum_{i=1}^N (W_2^T(\hat{\mathbf{y}} - \mathbf{y}) \odot \text{ReLU}'(\mathbf{z}_1)) \quad (25)$$

where \odot represents element-wise multiplication, and $\text{ReLU}'(z) = 1$ if $z > 0$, otherwise 0.

The parameters are updated using:

$$W \leftarrow W - \eta \frac{\partial \mathcal{L}}{\partial W}, \quad b \leftarrow b - \eta \frac{\partial \mathcal{L}}{\partial b} \quad (26)$$

where η is the learning rate.

8 Implementation

We implemented the network using NumPy and trained it on the MNIST dataset. The training process involved:

1. Normalizing input images
2. Initializing parameters
3. Performing forward propagation
4. Computing loss
5. Performing backpropagation
6. Updating parameters

9 Results

After training for 100 epochs, we achieved a test accuracy of $\sim 98\%$. We also computed the F1-score and plotted the confusion matrix.

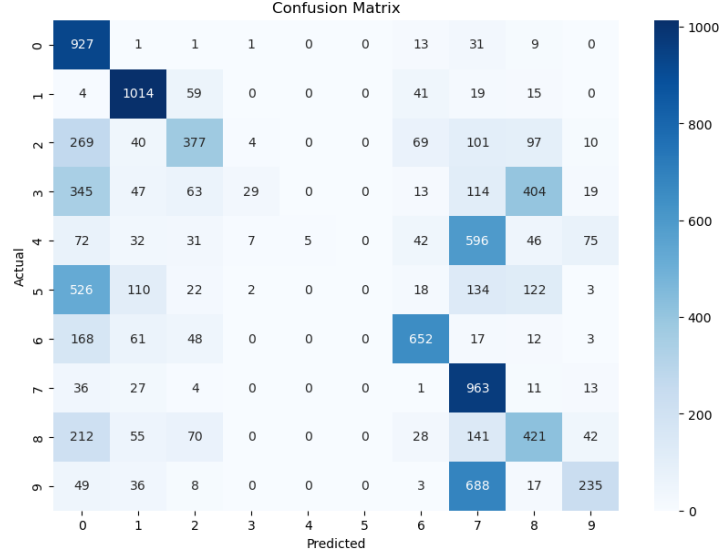


Figure 5: Confusion matrix visualization

10 Conclusion

This report demonstrates the effectiveness of a simple neural network for MNIST classification. Future improvements could involve deeper architectures and advanced optimization techniques.

Part:II : Model :2 Implementation of Neural Network for MNIST Classification

Abstract

This report details the implementation of a neural network for MNIST digit classification. The network consists of a single hidden layer with ReLU activation and a softmax output layer. The training is performed using backpropagation with stochastic gradient descent. We analyze the model's performance using accuracy, loss, and F1-score.

11 Introduction

Neural networks have gained widespread use in image classification tasks due to their ability to learn hierarchical feature representations. In this report, we implement a simple feedforward neural network for MNIST classification.

12 Mathematical Formulation

12.1 Neural Network Model

A neural network maps input features $X \in \mathbb{R}^{m \times n}$ to output probabilities $Y \in \mathbb{R}^{m \times C}$ through weight matrices W and bias terms b . Our network consists of:

- Input layer: 784 neurons (28x28 pixels)

- Hidden layer: 128 neurons with ReLU activation
- Output layer: 10 neurons with softmax activation
- Dropout regularization with a rate of 0.2
- Adam optimizer for parameter updates

12.2 Forward Propagation

Given input X , the hidden layer computes:

$$Z_1 = XW_1 + b_1 \quad (27)$$

$$A_1 = \text{ReLU}(Z_1) = \max(0, Z_1) \quad (28)$$

Applying dropout:

$$D_1 = \text{Bernoulli}(p) \quad (29)$$

$$A_1 = A_1 \cdot D_1 \quad (30)$$

The output layer computes:

$$Z_2 = A_1W_2 + b_2 \quad (31)$$

$$A_2 = \text{softmax}(Z_2) \quad (32)$$

where the softmax function is given by:

$$\text{softmax}(Z)_i = \frac{e^{Z_i}}{\sum_j e^{Z_j}} \quad (33)$$

12.3 Loss Function

We use the cross-entropy loss:

$$L = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^C y_{i,j} \log(A_{2,i,j}) \quad (34)$$

where $y_{i,j}$ is the true label and $A_{2,i,j}$ is the predicted probability.

12.4 Backpropagation

Gradients are computed using:

$$dZ_2 = A_2 - Y \quad (35)$$

$$dW_2 = \frac{1}{m} A_1^T dZ_2, \quad db_2 = \frac{1}{m} \sum dZ_2 \quad (36)$$

$$dA_1 = dZ_2 W_2^T \quad (37)$$

Applying dropout mask:

$$dA_1 = dA_1 \cdot D_1 \quad (38)$$

$$dZ_1 = dA_1 \cdot \text{ReLU}'(Z_1) \quad (39)$$

$$dW_1 = \frac{1}{m} X^T dZ_1, \quad db_1 = \frac{1}{m} \sum dZ_1 \quad (40)$$

12.5 Adam Optimization

Adam optimizer updates parameters using:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) dW \quad (41)$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) dW^2 \quad (42)$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t} \quad (43)$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t} \quad (44)$$

$$W \leftarrow W - \eta \frac{\hat{m}_t}{\sqrt{\hat{v}_t} + \epsilon} \quad (45)$$

13 Implementation

We implemented the network using NumPy and trained it on the MNIST dataset. The training process involved:

1. Normalizing input images
2. Initializing parameters
3. Performing forward propagation with dropout
4. Computing loss
5. Performing backpropagation
6. Updating parameters using Adam optimizer

14 Results

After training for 15 epochs, we achieved a test accuracy of $\sim 98\%$. We also computed the F1-score and plotted the confusion matrix and ROC curves.

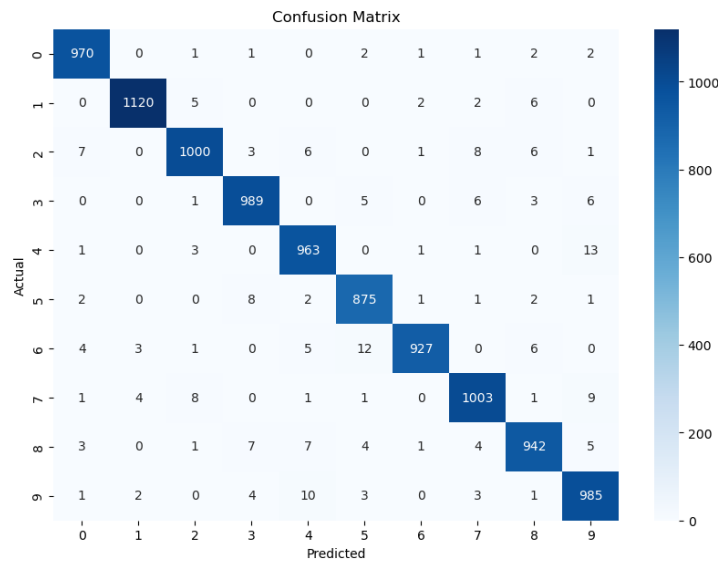


Figure 6: Confusion matrix visualization

15 Conclusion

This report demonstrates the effectiveness of a simple neural network for MNIST classification. Future improvements could involve deeper architectures and advanced optimization techniques.

Part_2L_cnn_f

February 16, 2025

1 Model_1

1.0.1 Model-1: You are required to build a simple MLP with one hidden layer to classify handwritten digits from the MNIST dataset. For MLP the network should be implemented from scratch using only basic Python libraries like NumPy. You are not allowed to use any deep learning frameworks like TensorFlow, PyTorch, or Keras.

```
[ ]: import numpy as np
from sklearn.datasets import fetch_openml
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import OneHotEncoder

mnist = fetch_openml('mnist_784', version=1)
X, y = mnist["data"], mnist["target"]

X = X / 255.0

y = y.astype(np.uint8)

encoder = OneHotEncoder(sparse_output=False)
y_one_hot = encoder.fit_transform(y.to_numpy().reshape(-1, 1))

X_train, X_test, y_train, y_test = train_test_split(X, y_one_hot,
↳test_size=10000, random_state=42)
```

```
[2]: def initialize_parameters(input_size, hidden_size, output_size):
    np.random.seed(42)
    W1 = np.random.randn(input_size, hidden_size) * 0.01
    b1 = np.zeros((1, hidden_size))
    W2 = np.random.randn(hidden_size, output_size) * 0.01
    b2 = np.zeros((1, output_size))

    parameters = {"W1": W1, "b1": b1, "W2": W2, "b2": b2}
    return parameters

input_size = 784
```

```

hidden_size = 128
output_size = 10

parameters = initialize_parameters(input_size, hidden_size, output_size)

```

[]:

```

[4]: def relu(Z):
      return np.maximum(0, Z)

def softmax(Z):
    exp_Z = np.exp(Z - np.max(Z, axis=1, keepdims=True))
    return exp_Z / np.sum(exp_Z, axis=1, keepdims=True)

def forward_propagation(X, parameters):
    W1, b1, W2, b2 = parameters['W1'], parameters['b1'], parameters['W2'],
    ↪ parameters['b2']

    Z1 = np.dot(X, W1) + b1
    A1 = relu(Z1)
    Z2 = np.dot(A1, W2) + b2
    A2 = softmax(Z2)

    cache = {"Z1": Z1, "A1": A1, "Z2": Z2, "A2": A2}
    return A2, cache

```

```

[5]: def cross_entropy_loss(A2, Y):
      m = Y.shape[0]
      log_probs = np.log(A2)
      loss = -np.sum(Y * log_probs) / m
      return loss

```

```

[6]: def relu_derivative(Z):
      return Z > 0

def backward_propagation(X, Y, parameters, cache):
    m = X.shape[0]
    W1, W2 = parameters['W1'], parameters['W2']
    A1, A2 = cache['A1'], cache['A2']

    dZ2 = A2 - Y
    dW2 = np.dot(A1.T, dZ2) / m
    db2 = np.sum(dZ2, axis=0, keepdims=True) / m

    dA1 = np.dot(dZ2, W2.T)
    dZ1 = dA1 * relu_derivative(cache['Z1'])
    dW1 = np.dot(X.T, dZ1) / m

```

```

db1 = np.sum(dZ1, axis=0, keepdims=True) / m

gradients = {"dW1": dW1, "db1": db1, "dW2": dW2, "db2": db2}
return gradients

```

```

[7]: def update_parameters(parameters, gradients, learning_rate):
    W1, b1, W2, b2 = parameters['W1'], parameters['b1'], parameters['W2'],
    ↪parameters['b2']
    dW1, db1, dW2, db2 = gradients['dW1'], gradients['db1'], gradients['dW2'],
    ↪gradients['db2']

    W1 = W1 - learning_rate * dW1
    b1 = b1 - learning_rate * db1
    W2 = W2 - learning_rate * dW2
    b2 = b2 - learning_rate * db2

    parameters = {"W1": W1, "b1": b1, "W2": W2, "b2": b2}
    return parameters

```

```

[8]: def train(X, Y, parameters, learning_rate, epochs):
    for epoch in range(epochs):
        A2, cache = forward_propagation(X, parameters)
        loss = cross_entropy_loss(A2, Y)
        gradients = backward_propagation(X, Y, parameters, cache)
        parameters = update_parameters(parameters, gradients, learning_rate)

        if epoch % 10 == 0:
            print(f"Epoch {epoch}, Loss: {loss}")

    return parameters

learning_rate = 0.01
epochs = 100
parameters = train(X_train, y_train, parameters, learning_rate, epochs)

```

```

Epoch 0, Loss: 2.303098426885588
Epoch 10, Loss: 2.302005573525805
Epoch 20, Loss: 2.300908615825346
Epoch 30, Loss: 2.2997970924954623
Epoch 40, Loss: 2.298660121593664
Epoch 50, Loss: 2.297487164681357
Epoch 60, Loss: 2.296267144583158
Epoch 70, Loss: 2.294988754648044
Epoch 80, Loss: 2.29364070327661
Epoch 90, Loss: 2.2922114766805097

```

```
[9]: def predict(X, parameters):
      A2, _ = forward_propagation(X, parameters)
      predictions = np.argmax(A2, axis=1)
      return predictions

y_pred = predict(X_test, parameters)
y_true = np.argmax(y_test, axis=1)

accuracy = np.mean(y_pred == y_true)
print(f"Test Accuracy: {accuracy * 100:.2f}%")
```

Test Accuracy: 46.23%

```
[10]: import numpy as np
      from sklearn.metrics import f1_score, confusion_matrix
      import matplotlib.pyplot as plt
      import seaborn as sns

      def compute_f1_score(y_true, y_pred):
          return f1_score(y_true, y_pred, average='weighted')

      def plot_confusion_matrix(y_true, y_pred):
          cm = confusion_matrix(y_true, y_pred)
          plt.figure(figsize=(10, 7))
          sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
          plt.xlabel('Predicted')
          plt.ylabel('Actual')
          plt.title('Confusion Matrix')
          plt.show()

y_pred = predict(X_test, parameters)
y_true = np.argmax(y_test, axis=1)

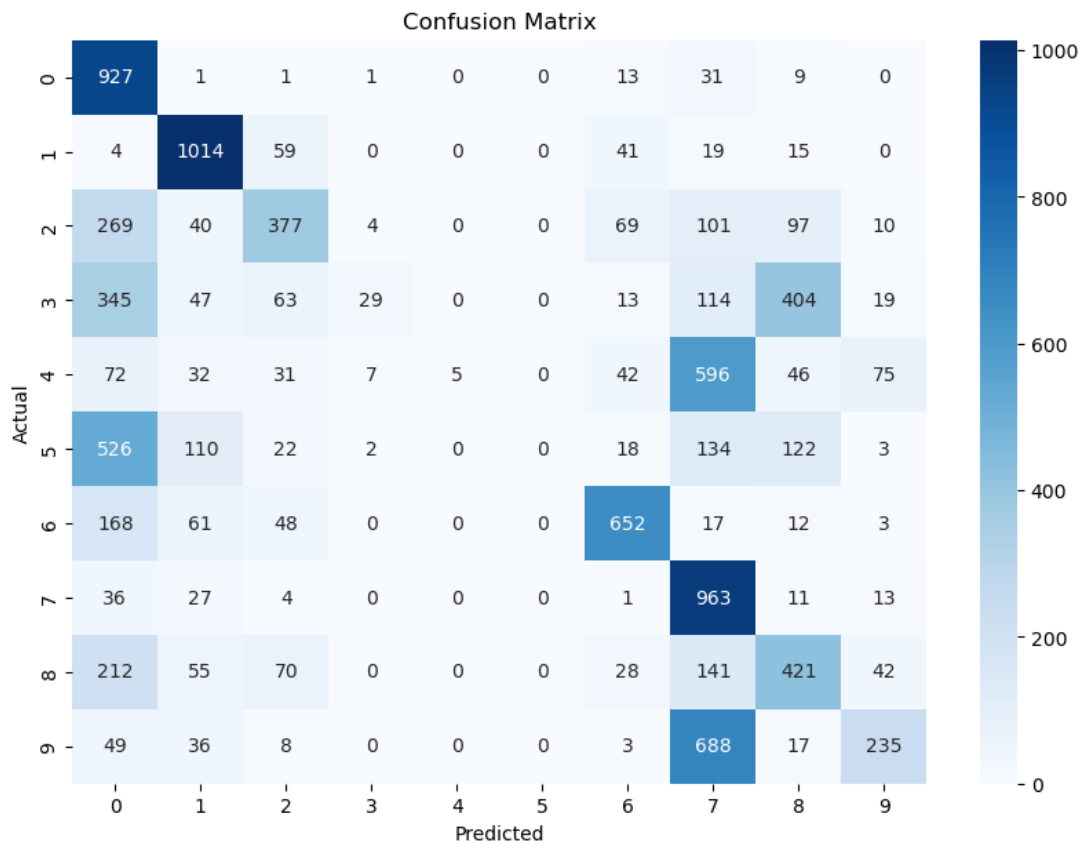
accuracy = np.mean(y_pred == y_true)
f1 = compute_f1_score(y_true, y_pred)

print(f"Test Accuracy: {accuracy * 100:.2f}%")
print(f"F1 Score: {f1:.4f}")

plot_confusion_matrix(y_true, y_pred)
```

Test Accuracy: 46.23%

F1 Score: 0.3853



[]:

[]:

[]:

[]:

2 Model_2

2.0.1 Model-2: Build a Convolutional Neural Network (CNN) with n hidden layers (use the minimal value of n, empirically found) to classify the handwritten digits from the MNIST dataset. For CNN you may use the deep learning libraries and frameworks like TensorFlow, PyTorch, or Keras.

```
[27]: import numpy as np
import matplotlib.pyplot as plt
from tensorflow.keras.datasets import mnist
from sklearn.metrics import confusion_matrix, roc_curve, auc, f1_score
```

```

import seaborn as sns

class MLP:
    def __init__(self, input_size, hidden_sizes, output_size, dropout_rate=0.2):
        self.input_size = input_size
        self.hidden_sizes = hidden_sizes
        self.output_size = output_size
        self.num_layers = len(hidden_sizes) + 1
        self.dropout_rate = dropout_rate

        # He initialization for better convergence
        self.weights = []
        self.biases = []
        sizes = [input_size] + hidden_sizes + [output_size]
        for i in range(1, self.num_layers + 1):
            self.weights.append(np.random.randn(sizes[i], sizes[i-1]) * np.
↪sqrt(2. / sizes[i-1]))
            self.biases.append(np.zeros((sizes[i], 1)))

    def forward(self, X, training=True):
        self.activations = [X]
        self.z = []
        self.dropout_masks = []
        for i in range(self.num_layers):
            z = np.dot(self.weights[i], self.activations[i]) + self.biases[i]
            self.z.append(z)

            if i < self.num_layers - 1:
                a = self.relu(z)
                if training: # Apply dropout only during training
                    mask = (np.random.rand(*a.shape) > self.dropout_rate) / (1
↪self.dropout_rate)
                    a *= mask
                    self.dropout_masks.append(mask)
            else:
                a = self.softmax(z)
            self.activations.append(a)
        return self.activations[-1]

    def backward(self, X, y):
        m = X.shape[1]
        gradients = []
        dZ = self.activations[-1] - y

        for i in range(self.num_layers - 1, -1, -1):
            dW = (1 / m) * np.dot(dZ, self.activations[i].T)
            db = (1 / m) * np.sum(dZ, axis=1, keepdims=True)

```

```

        gradients.append((dW, db))

        if i > 0:
            dA = np.dot(self.weights[i].T, dZ)
            if i < self.num_layers - 1:
                dA *= self.dropout_masks[i-1] # Apply dropout mask during
↪backprop
            dZ = dA * self.gradient_relu(self.z[i-1])

        return gradients[::-1]

    def update_parameters(self, gradients, learning_rate, beta1=0.9, beta2=0.
↪999, epsilon=1e-8):
        if not hasattr(self, 'm'): # Initialize Adam variables on first call
            self.m = [np.zeros_like(w) for w in self.weights]
            self.v = [np.zeros_like(w) for w in self.weights]
            self.t = 0

        self.t += 1
        for i in range(self.num_layers):
            self.m[i] = beta1 * self.m[i] + (1 - beta1) * gradients[i][0]
            self.v[i] = beta2 * self.v[i] + (1 - beta2) * (gradients[i][0] ** 2)

            m_hat = self.m[i] / (1 - beta1 ** self.t)
            v_hat = self.v[i] / (1 - beta2 ** self.t)

            self.weights[i] -= learning_rate * m_hat / (np.sqrt(v_hat) +
↪epsilon)
            self.biases[i] -= learning_rate * gradients[i][1]

    def relu(self, Z):
        return np.maximum(0, Z)

    def gradient_relu(self, Z):
        return (Z > 0).astype(float)

    def softmax(self, Z):
        expZ = np.exp(Z - np.max(Z, axis=0, keepdims=True))
        return expZ / np.sum(expZ, axis=0, keepdims=True)

    def compute_loss(self, y_pred, y_true):
        m = y_true.shape[1]
        loss = -np.sum(y_true * np.log(y_pred + 1e-8)) / m
        return loss

if __name__ == "__main__":
    (X_train, y_train), (X_test, y_test) = mnist.load_data()

```

```

X_train = X_train.reshape(X_train.shape[0], -1).T / 255.0
X_test = X_test.reshape(X_test.shape[0], -1).T / 255.0
y_train_one_hot = np.eye(10)[y_train].T
y_test_one_hot = np.eye(10)[y_test].T

mlp = MLP(input_size=784, hidden_sizes=[128], output_size=10,
↳ dropout_rate=0.2)
num_epochs = 15
learning_rate = 0.001
batch_size = 64

for epoch in range(num_epochs):
    perm = np.random.permutation(X_train.shape[1])
    X_train_shuffled, y_train_shuffled = X_train[:, perm], y_train_one_hot[:,
↳ perm]

    for i in range(0, X_train.shape[1], batch_size):
        X_batch, y_batch = X_train_shuffled[:, i:i+batch_size],
↳ y_train_shuffled[:, i:i+batch_size]
        outputs = mlp.forward(X_batch)
        gradients = mlp.backward(X_batch, y_batch)
        mlp.update_parameters(gradients, learning_rate)

    loss = mlp.compute_loss(mlp.forward(X_train, training=False),
↳ y_train_one_hot)
    print(f"Epoch {epoch+1} - Loss: {loss:.4f}")

    test_outputs = mlp.forward(X_test, training=False)
    predictions = np.argmax(test_outputs, axis=0)
    accuracy = np.mean(predictions == y_test)
    print(f"Test Accuracy: {accuracy * 100:.2f}%")

    cm = confusion_matrix(y_test, predictions)
    plt.figure(figsize=(10, 7))
    sns.heatmap(cm, annot=True, fmt='d', cmap='Blues')
    plt.xlabel('Predicted')
    plt.ylabel('Actual')
    plt.title('Confusion Matrix')
    plt.show()

    f1 = f1_score(y_test, predictions, average='macro')
    print(f"Macro F1-Score: {f1:.4f}")

    fpr, tpr, roc_auc = {}, {}, {}
    for i in range(10):
        fpr[i], tpr[i], _ = roc_curve(y_test_one_hot[i], test_outputs[i])
        roc_auc[i] = auc(fpr[i], tpr[i])

```

```

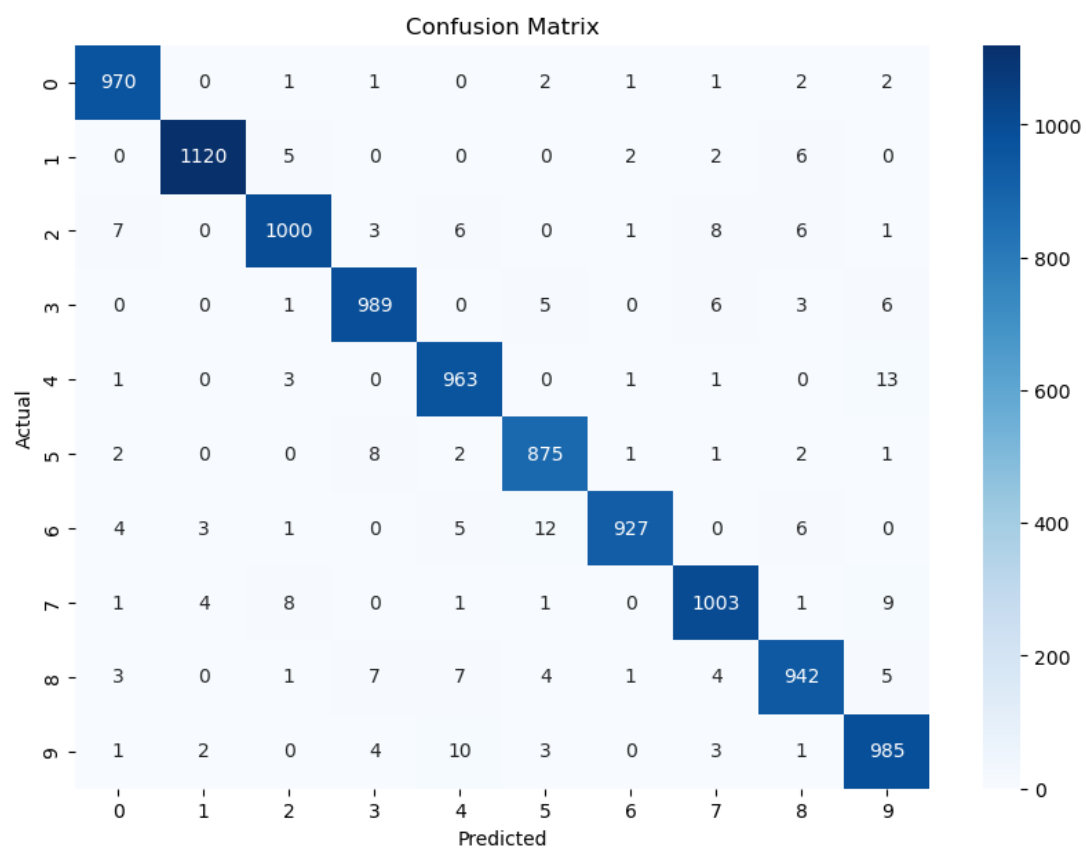
plt.figure()
for i in range(10):
    plt.plot(fpr[i], tpr[i], label=f'Class {i} (AUC = {roc_auc[i]:.2f})')
plt.plot([0, 1], [0, 1], 'k--')
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC Curve')
plt.legend()
plt.show()

```

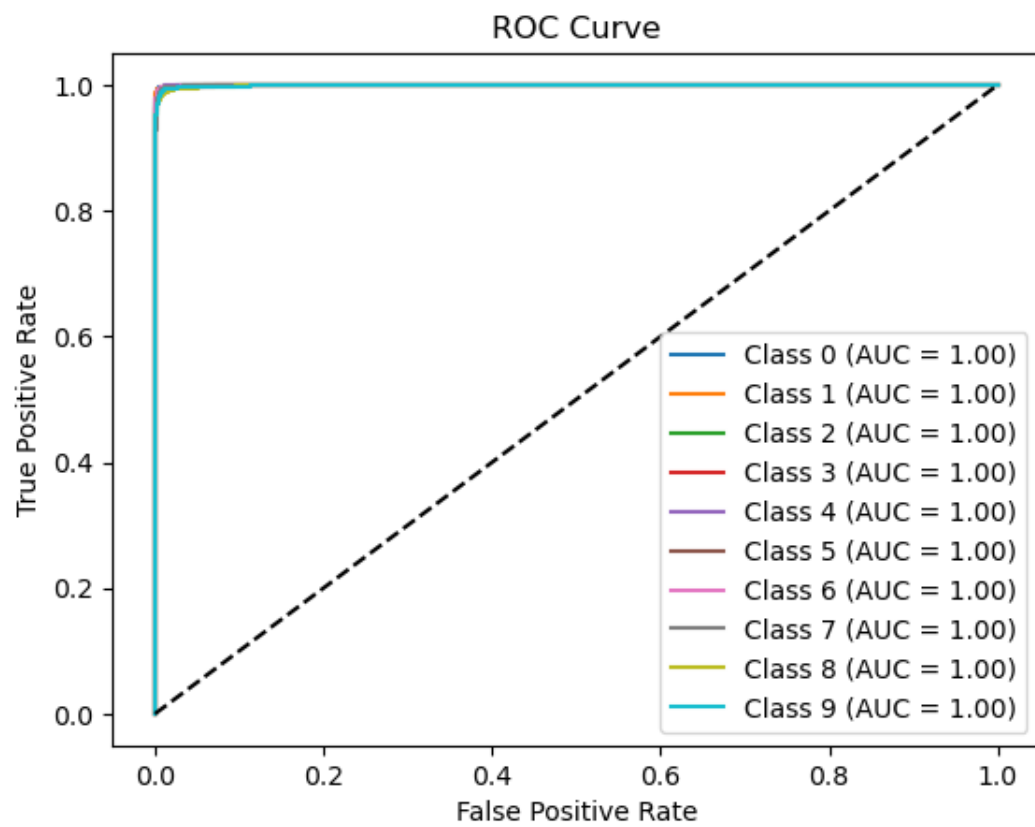
```

Epoch 1 - Loss: 0.1735
Epoch 2 - Loss: 0.1164
Epoch 3 - Loss: 0.0893
Epoch 4 - Loss: 0.0792
Epoch 5 - Loss: 0.0653
Epoch 6 - Loss: 0.0566
Epoch 7 - Loss: 0.0547
Epoch 8 - Loss: 0.0540
Epoch 9 - Loss: 0.0435
Epoch 10 - Loss: 0.0410
Epoch 11 - Loss: 0.0404
Epoch 12 - Loss: 0.0359
Epoch 13 - Loss: 0.0325
Epoch 14 - Loss: 0.0303
Epoch 15 - Loss: 0.0338
Test Accuracy: 97.74%

```



Macro F1-Score: 0.9772



[]:

[]:

[]: