

Mini-Project: Deep Learning from scratch

Tamer, Abu-shakra
319147260

Tomer, Varsanno
206727471

Naor, Guetta
316017359

April 21, 2024

Contents

1 Part 1.1 :	2
1.1 Softmax Regression	2
1.2 Softmax-Regression-Loss : cross-entropy function	2
1.3 Gradients by W and X for SM	2
2 Part 1.2 :	5
2.1 SGD on Least Squares	5
2.2 Stochastic gradient descent with Momentum	6
3 Part 1.3:	8
3.1 Stochastic gradient descent with Momentum on SoftMax	8
3.2 Batch size	9
3.3 Learning rate	10
4 Part 2.1 :	11
4.1 activation functions	11
4.2 One step forward for the standard neural network	11
4.3 Verification of Derivatives via Jacobian Test for standard network	12
5 Part 2.2 :	14
5.1 One step forward for the residual neural network	14
5.2 Verification of Derivatives via Jacobian Test for residual network	15
6 Part 2.3 :	18
6.1 Initialize network	18
6.2 Forward pass	18
6.3 Back propagation	19
6.4 Gradient test for the complete network	20
7 Part 2.4 :	21
7.1 train network	21
7.2 data validation test	21
7.3 Optimal Network Configurations: Evaluating Different Lengths	21
8 Part 2.5	31
8.1 Optimizing Neural Network Efficiency - "Swiss Roll"	31
9 Part 2.5	33
9.1 Optimizing Neural Network Efficiency - "GMM"	33
10 Code implementation here	34

Part 1 : The classifier and optimizer

1 Part 1.1 :

1.1 Softmax Regression

We implemented the Soft-Max-Regression function to transform the neural network's final layer into the probability of an instance belonging to a particular class.

$$\text{Softmax } (Z) = \frac{\exp(xw_j - \eta)}{\sum_{i=1}^N \exp(xw_i - \eta)}$$

We added the η as recommended to ensure safe computation.

Z : The weights matrix (W) time Data matrix (X)

X : Input data matrix where each row represents a feature and each column represents a sample.

Shape: (number of features, number of samples)

W : Weight matrix representing the parameters of the model.

Shape : (number of classes, number of features)

Returns: Output of the softmax regression, representing the predicted probabilities for each class.

Shape: (number of classes, number of samples)

Code implementation :

$$Z = X \cdot W$$

$$\text{Exp} = \exp(Z - \max(Z))$$

$$\text{Softmax Output}_i = \frac{\exp(Z_i)}{\sum_{j=1}^N \exp(Z_j)}$$

Here is an example for a 3 data features and 4 data samples:

$$\begin{bmatrix} a_{11} & \color{red}{a_{12}} & a_{13} \\ a_{21} & \color{red}{a_{22}} & a_{23} \\ \color{blue}{a_{31}} & a_{32} & \color{blue}{a_{33}} \\ a_{41} & \color{red}{a_{42}} & a_{43} \end{bmatrix} \quad \begin{array}{l} \text{row = Feature} \\ \text{column = Data sample} \end{array}$$

1.2 Softmax-Regression-Loss : cross-entropy function

$$\text{softmax regression loss } (X, W, C) = -\frac{1}{m} \sum_{k=1}^l c_k \log \left(\text{diag}(\sum_{j=1}^{\ell} \exp(XW_j))^{-1} \cdot \exp(XW_k) \right)$$

C : One-hot encoded class labels for each sample.

Shape: (number of samples, number of classes)

Returns : Softmax regression loss, calculated as the negative log-likelihood of the true class probabilities.

Code implementation :

$$\text{Probabilities} = \text{soft_max_regression}(X, W)$$

$$\text{Log Probs} = -\log(\text{Probabilities})$$

$$\text{Loss} = \sum_{i=1}^m \text{Log Probs}_i$$

1.3 Gradients by W and X for SM

Softmax regression gradient by W

$$\text{softmax regression gradient by W}(X, W, C) = \frac{1}{m} \left[\text{diag}(\sum_j \exp(W_j X))^{-1} \cdot \exp(W_k X) - C_p \right] X^T$$

Returns : Gradient by W of the softmax regression loss function.

Shape: (number of classes, number of features)

Code implementation :

$$\text{Probabilities} = \text{soft max regression}(X, W)$$

$$\text{Probabilities} = \text{Probabilities} - 1$$

$$\text{Gradient} = \frac{1}{m} \cdot \text{Probabilities} \cdot X^T$$

```

def gradient_test_byW(function, gradient):
    # Gradient Check Setup
    batchSize = 2
    EPSILON = 0.1
    outputnum = 2

    X , C =pickRandomBatch(data_training_X,data_training_C,batchSize)
    X = np.vstack([X, np.ones((1, X.shape[1]))])

    W = np.random.randn(outputnum, X.shape[0])

    d = np.random.randn(*W.shape)

    def F(W):
        Z = W.dot(X)
        return function(Z, C)

    def g_F(W):
        return gradient(X, W, C)

    F0 = F(W)
    g0 = g_F(W)
    y0 = np.zeros(8)
    y1 = np.zeros(8)

    print("k\terror order 1\t\terror order 2")

    for k in range(1, 9):
        epsk = EPSILON * (0.5 ** k)
        Fk = F(W + epsk * d)
        F1 = F0 + epsk * np.sum(g0 * d)
        y0[k-1] = abs(Fk - F0)
        y1[k-1] = abs(Fk - F1)
        print(f"\"{k}\\"{y0[k-1]}\"{y1[k-1]}\"")

    # Plotting
    plt.semilogy(range(1, 9), y0, label="Zero order approx")
    plt.semilogy(range(1, 9), y1, label="First order approx")
    plt.legend()
    plt.title("Gradient Test for Softmax Regression By W")
    plt.xlabel("k")
    plt.ylabel("Error")
    plt.show()

gradient_test_byW(compute_loss_for_testing,softmax_regression_gradient_by_W)

```

Figure 1: SM gradient test by W code

```

def softmax_regression_gradient_by_W(X ,W, C):
    Z = W.dot(X)
    m = Z.shape[1]
    probabilities = softmax(Z)
    # for i in range(m):
    #     probabilities[:, i] -= C[:, i]
    probabilities -= C
    gradient = np.dot(probabilities, X.T) / m
    return gradient

```

k	error order 1	error order 2
1	0.09386261850254052	0.002451807682035234
2	0.047550851886846934	0.000660631612054408862
3	0.02392785072508503	0.00015075582113888153
4	0.012001719664543503	3.758420852850719e-05
5	0.006610268629983472	9.383006732477384e-06
6	0.0030074816578324892	2.3441604355411982e-06
7	0.0015043270472882009	5.858618458143283e-07
8	0.0007523100010101702	1.4645355683740036e-07

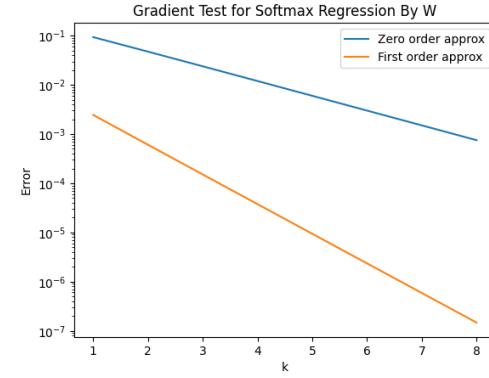


Figure 2: SM gradient by W code + results

Softmax regression gradient by X

softmax regression gradient by $X(X, W, C) = \frac{1}{m} W^T \left[\exp(W_k X) \oslash \text{diag}(\sum_j \exp(W_j X))^{-1} - C \right]$

Returns : Gradient by X of the softmax regression loss function.

Shape: (number of features, number of samples)

Code implementation :

Probabilities = soft max regression(X, W)

Probabilities = Probabilities - C

$$\text{Gradient} = \frac{1}{m} \cdot W^T \cdot \text{Probabilities}$$

```

def gradient_test_byX(function, gradient):
    # Gradient Check Setup
    batchSize = 2
    EPSILON = 0.1
    outputnum = 2

    X , C =pickRandomBach(data_training_X,data_training_C,batchSize)

    W1 = np.random.randn(2, 2)

    d = np.random.randn(*X.shape)

    def F(X):
        Z = W1.dot(X)
        return function(Z , C)

    def g_F(X):
        return gradient( X , W1 , C)

    F0 = F(X)
    g0 = g_F(X)
    y0 = np.zeros(8)
    y1 = np.zeros(8)

    print("K\terror order 1\t\terror order 2")

    for k in range(1, 9):
        epsk = EPSILON * (0.5 ** k)
        Fk = F(X + epsk * d)
        F1 = F0 + epsk * np.sum(g0 * d)
        y0[k-1] = abs(Fk - F0)
        y1[k-1] = abs(Fk - F1)
        print(f"\t{k}\t{y0[k-1]}\t{y1[k-1]}")

    # Plotting
    plt.semilogy(range(1, 9), y0, label="Zero order approx")
    plt.semilogy(range(1, 9), y1, label="First order approx")
    plt.legend()
    plt.title("Gradient Test for Softmax Regression By X")
    plt.xlabel("K")
    plt.ylabel("Error")
    plt.show()

gradient_test_byX(compute_loss_for_testing,softmax_regression_gradient_by_X)

```

```

def softmax_regression_gradient_by_X(X, W, C):
    m = X.shape[1]
    A = np.dot(W, X)
    prob = softmax(A)
    P = prob - C
    grad_X = np.dot(W.T, P) / m
    return grad_X

```

k	error order 1	error order 2
1	0.010403109720281245	0.00018141415968974432
2	0.00515630416646836	4.153868167240887e-05
3	0.002566968418748805	1.274528762350522e-05
4	0.0012805278873428971	2.815942268974414e-06
5	0.0006395596179034624	7.03645366528732e-07
6	0.0003196803533473944	1.758670788982733e-07
7	0.000159759353174605	4.396838223521583e-08
8	7.98678543879376e-05	0.098887170473318e-08

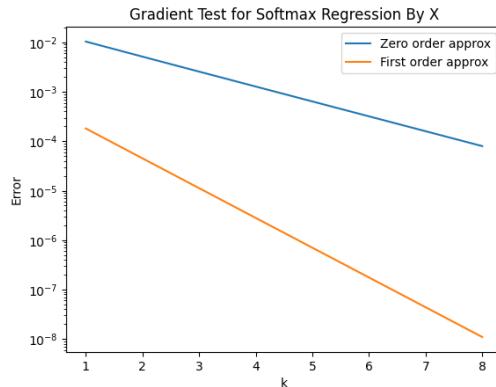


Figure 3: SM gradient test by X code

Figure 4: SM gradient by X code + results

2 Part 1.2 :

2.1 SGD on Least Squares

In this section, we used Least Squares in order to check our SGD with Momentum function we wrote the code for LS and LS gradient and checked with Gradient test for correctness.

Gradient Test for Least Squares

k	error order 1	error order 2
1	0.20001922201421962	0.004400367965265772
2	0.09890951901579381	0.001100091913168872
3	0.049179736510066796	0.0002758229978278895
4	0.02452112505575593	6.875574945652829e-05
5	0.012243367315423193	1.718893736243895e-05
6	0.006117386423371129	4.297234341699152e-06
7	0.003857618903180412	1.074308582027434e-06
8	0.0015285408744934612	2.6857714630068585e-07

```

def gradient_test_LS(function, gradient):
    # Gradient Check Setup
    batchSize = 2
    EPSILON = 0.1
    outputnum = 2

    X , C =pickRandomBach(data_training_X,data_training_C,batchSize)
    X = np.vstack([X, np.ones((1, X.shape[1]))])

    W = np.random.randn(outputnum, X.shape[0])

    d = np.random.randn(*W.shape)

    def F(W):
        return function(X, W, C)

    def g_F(W):
        return gradient(X, W, C)

    F0 = F(W)
    g0 = g_F(W)
    y0 = np.zeros(8)
    y1 = np.zeros(8)

    print("k\t|error order 1|\t|error order 2|")

    for k in range(1, 9):
        epsk = EPSILON * (0.5 ** k)
        Fk = F(W + epsk * d)
        F1 = F0 + epsk * np.sum(g0 * d)
        y0[k-1] = abs(Fk - F0)
        y1[k-1] = abs(Fk - F1)
        print(f'{k}\t{y0[k-1]}\t{y1[k-1]}')

    # Plotting
    plt.semilogy(range(1, 9), y0, label="Zero order approx")
    plt.semilogy(range(1, 9), y1, label="First order approx")
    plt.legend()
    plt.title("Gradient Test for Least Squares")
    plt.xlabel("k")
    plt.ylabel("Error")
    plt.show()

gradient_test_LS(least_squares_loss, least_squares_gradient)

```

Figure 5: LS gradient test code

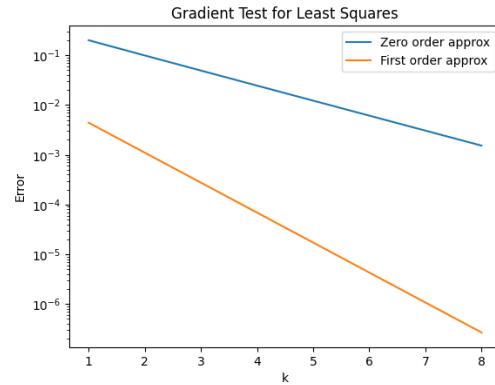


Figure 6: LS code + Results

2.2 Stochastic gradient descent with Momentum

learning_rate : The learning rate determines the step size in the weight updates during optimization.
momentum : The momentum controls the contribution of the previous velocity to the current update.
epochs : Represents how often the model will iterate over the entire training dataset.

Output:

network : The final network containing the weights after the optimization process.

We choose to use gradient descent with Momentum to try to avoid local minimum.

Code implementation :

$$\begin{aligned} w &:= w - \eta \nabla Q_i(w) + \alpha \Delta w \\ \Delta w &:= \alpha \Delta w - \eta \nabla Q_i(w) \end{aligned}$$

```
def SGD_with_momentum(network, gradients, velocity, learning_rate, momentum):
    for key in network.keys():
        velocity[key] = momentum * velocity[key] - learning_rate * gradients[key]
        network[key] += velocity[key]
```

Figure 7: SGD optimizer code

```

def testSGD(lossFunction , gradientLossFunction):
    batchSize = 2
    EPSILON = 0.1
    outputnum = 2
    network = {}

    X , C =pickRandomBach(data_training_X,data_training_C,batchSize)
    X = np.vstack([X, np.ones((1, X.shape[1]))])
    W1 = np.random.randn(outputnum, X.shape[0])
    network = {'W1' : W1}
    velocity = {'W1' :np.zeros_like(W1)}

    success_percentages_list = []
    success_percentages = 0
    losses = []

    # Set hyperparameters
    learning_rate = 0.05
    momentum = 0.8
    epochs = 80
    batch_Size = 20

    for i in range(epochs):
        W1 = network['W1']
        X , C =pickRandomBach(data_training_X,data_training_C,batchSize)
        X = np.vstack([X, np.ones((1, X.shape[1]))])
        gradients = {'W1': least_squares_gradient(X, W1, C)}

        SGD_with_momentum(network, gradients, velocity, learning_rate, momentum)

        losses.append(least_squares_loss(X, W1, C))

        success_percentages_list.append(calculate_success_percentage(W1.dot(X), C))

    # Plot the loss over epochs
    plt.plot(range(epochs), losses)
    plt.title("Losses over Epochs (SGD with Momentum) with LS")
    plt.xlabel("Epoch")
    plt.ylabel("losses")
    plt.show()

    # Plot the success percentages over epochs
    plt.plot(range(epochs), success_percentages_list)
    plt.title("success percentages over Epochs (SGD with Momentum) with LS")
    plt.xlabel("Epoch")
    plt.ylabel("success")
    plt.show()

```

Figure 8: SGD LS test code

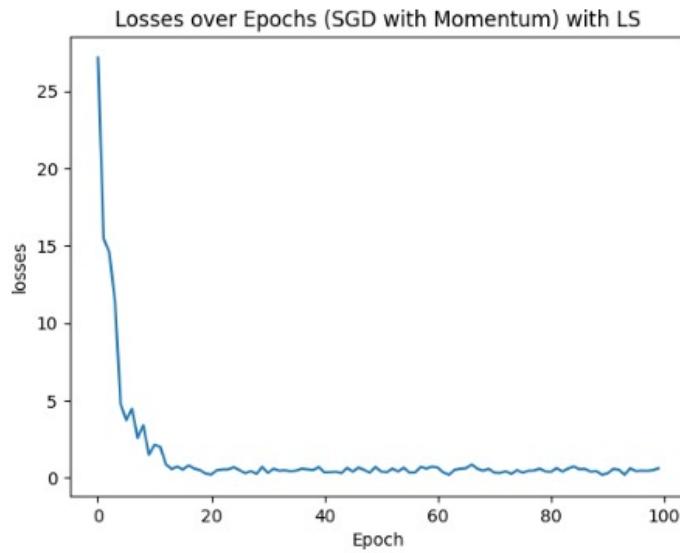


Figure 9: SGD LS test

3 Part 1.3:

3.1 Stochastic gradient descent with Momentum on SoftMax

On the left side is the graph for the Training Data .
And on the right side the graph for the Validation Data .

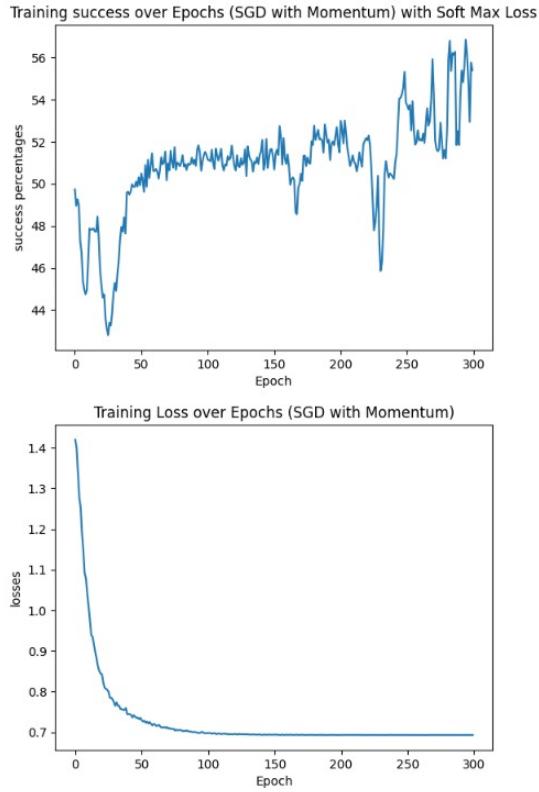


Figure 10: Loss and Success-Percentage - Training Data

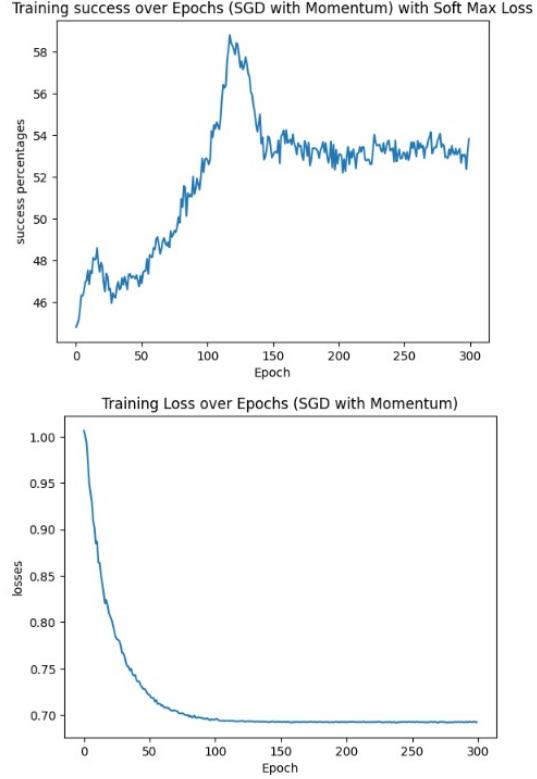


Figure 11: Loss and Success-Percentage - Validation Data

3.2 Batch size

In our testing, we observed that increasing the batch size resulted in smoother and less noisy trends in our results, with a gradual upward or downward movement. Conversely, using a smaller batch size led to more pronounced changes with each epoch. Despite both scenarios resulting in lower losses, the success percentage remained consistently between 45 percent to 60 percent. This outcome was anticipated, considering we haven't yet implemented additional layers in our model.

examples :

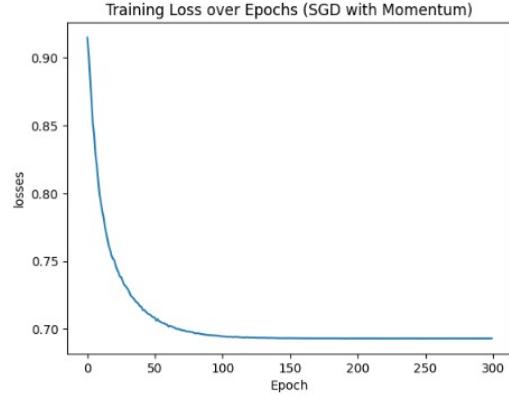
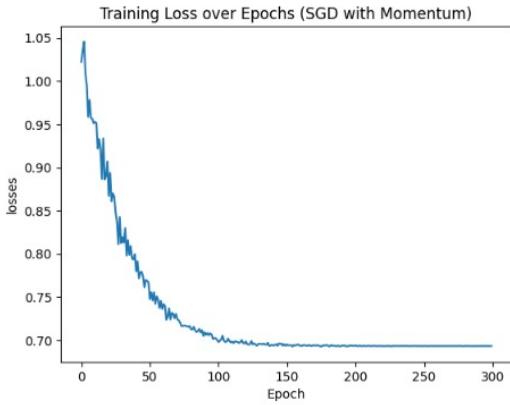
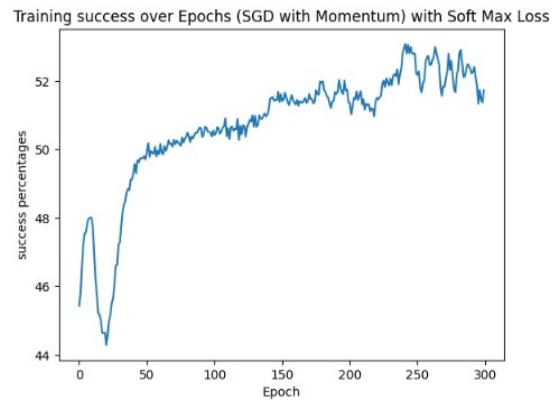
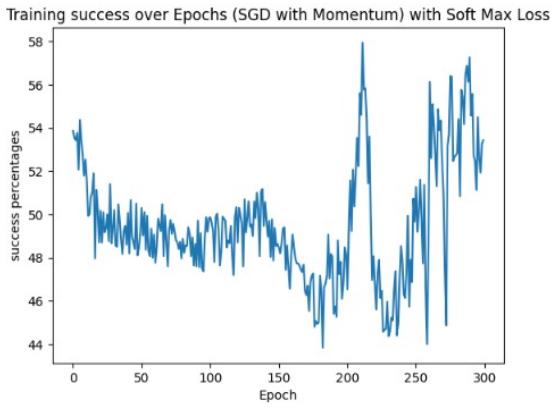


Figure 12: Lower batch size

Figure 13: Higher batch size

3.3 Learning rate

Furthermore, Increasing the learning rate caused a substantial increase or decrease in the success percentage, with each epoch producing notable changes. The losses also decreased rapidly. On the other hand, decreasing the learning rate led to more subtle shifts in the success percentage graph. In the loss graph, we observed a slower decline.

examples :

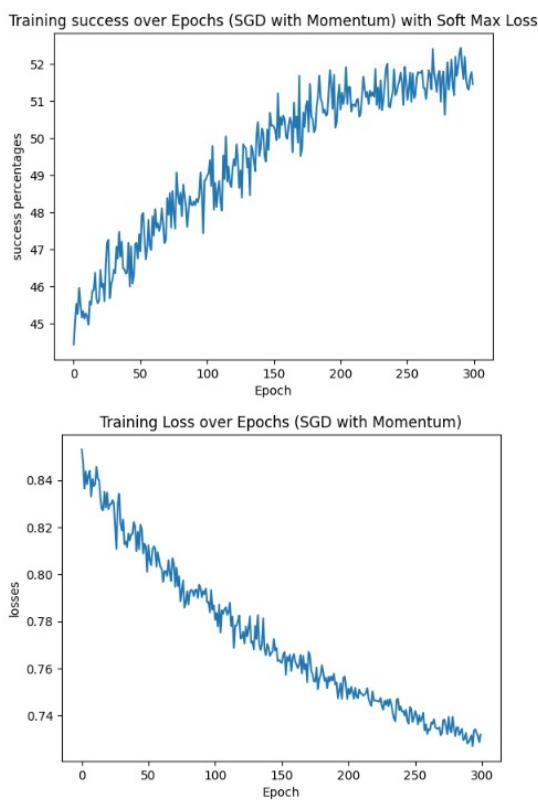


Figure 14: Lower learning rate

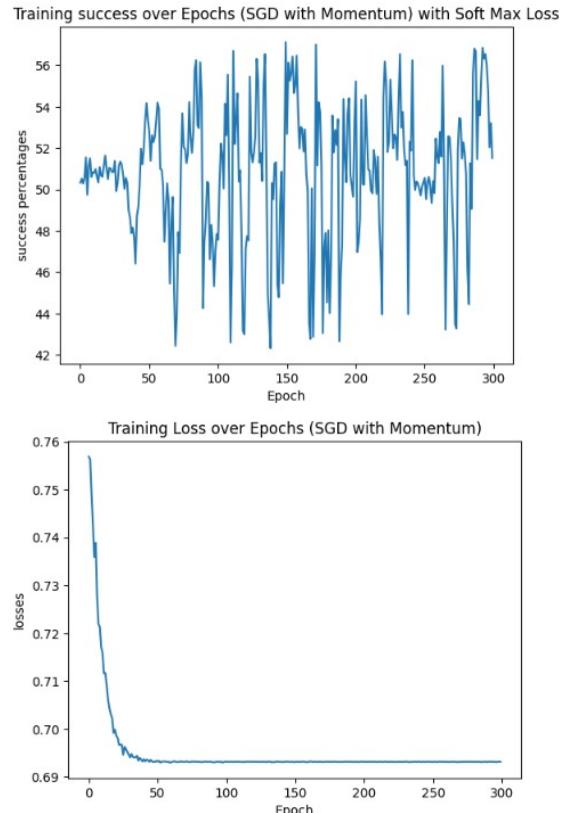


Figure 15: Higher learning rate

4 Part 2.1 :

4.1 activation functions

we implemented two activation functions, hyperbolic tangent (tanh) and Rectified Linear Unit (ReLU), along with their derivatives.

4.2 One step forward for the standard neural network

The `single_step_forward_class` function computes a forward step in a standard neural network using a selected activation function.

1. **Select the activation function:** The activation function is chosen based on the input parameter.

$$\text{activation_func} = \begin{cases} \text{relu} & \text{if activation == 'relu'} \\ \text{tanh} & \text{otherwise} \end{cases}$$

2. **Compute and return the activated output:** The output is computed as follows:

$$\text{activation}(\mathbf{W} \cdot \mathbf{X})$$

```
def single_step_forward_class(X, W, activation='relu'):
    activation_func = relu if activation == 'relu' else tanh

    Z = W.dot(X)
    A = activation_func(Z)
    return A
```

Figure 16: Classic single step forward code

4.3 Verification of Derivatives via Jacobian Test for standard network

This section demonstrates that the computed derivatives with respect to both \mathbf{X} and \mathbf{W} in the neural network model satisfy the Jacobian test.

Derivative with respect to W^T

$$\frac{\partial y}{\partial \mathbf{W}}^T \mathbf{v} = (\sigma'(\mathbf{WX})\mathbf{v})\mathbf{X}^T$$

X : Matrix representing the outputs of the previous layer.

W : Weight matrix associated with the current layer, used to transform the input matrix X before applying the activation function.

u : matrix containing the outputs from the forward pass of the current layer.

activation ('relu' by default): Specifies the type of activation function to use.

Returns: Computes and returns the Jacobian transpose matrix time u , This matrix is used in gradient computations for optimizing the weights during the backpropagation process.

Code implementation :

$$\begin{aligned} Z &= W \cdot X \\ A &= \text{activation}(Z) \\ \text{return : } &= (\mathbf{A} \circ \mathbf{V}) \mathbf{X}^T \end{aligned}$$

Figure 17: Jacobian transposed test code for \mathbf{W}'

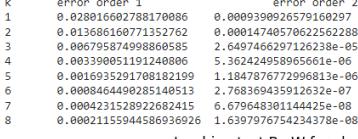


Figure 18: Jacobian by W function + test

Derivative with respect to X^T

$$\frac{\partial y}{\partial \mathbf{X}}^T \mathbf{v} = \mathbf{W}^T(\sigma'(\mathbf{W}\mathbf{X})\mathbf{v})$$

X : Matrix representing the outputs of the previous layer.

W : Weight matrix associated with the current layer, used to transform the input matrix X before applying the activation function.

u : matrix containing the outputs from the forward pass of the current layer.

activation ('relu' by default): Specifies the type of activation function to use.

Returns: Computes and returns the Jacobian transpose matrix time u.

Code implementation :

$$Z = W \cdot X$$

$$A = \text{activation}(Z)$$

return : $\mathbf{W}^T(\mathbf{A} \circ \mathbf{V})$

```
def jacobian_test_by_X_for_classic():
```

```

# Jacobian Check Setup
num_classes = 3
input_size = 2
num_features = 5
EPSILON = 0.1
W = np.random.randn(num_classes, num_features + 1)
X = np.random.randn(num_features, input_size)
X = np.vstack([X, np.ones((1, X.shape[1]))])
d = np.random.randn(num_features + 1, input_size)
u = np.random.randn(num_classes, input_size)

def F(X):
    return single_step_forward_class(X, W, tanh)

def g_F(X):
    return np.dot(u.flatten(), F0.flatten())

F0 = F(X)
g0 = g_F(X)
grad_g0 = JacobianTByX(X, W, u, tanh)

y0 = np.zeros(8)
y1 = np.zeros(8)

print("k\terror order 1 \t\t\t error order 2")
for k in range(1, 9):
    epsk = EPSILON * (0.5 ** k)
    fk = F(X + (epsk * d))
    gk = np.dot(fk.flatten(), u.flatten())
    g1 = g0 + (epsk * (np.dot(d.flatten(), grad_g0.flatten())))
    y0[k - 1] = abs(gk - g0)
    y1[k - 1] = abs(gk - g1)
    print(f"\t{k}\t{y0[k-1]}\t\t{y1[k-1]}")

# Plotting
plt.semilogy(range(1, 9), y0, label="Zero order approx")
plt.semilogy(range(1, 9), y1, label="First order approx")
plt.legend()
plt.title("Jacobian test By X for classic")
plt.xlabel("k")
plt.ylabel("Error")
plt.show()

def JacobianTByX(X, W, u, activation='relu'):

    activationD_func = relu_derivative if activation == 'relu' else tanh_derivative

    Z = W.dot(X)
    A = activationD_func(Z)

    return np.dot(W.T, np.multiply(A,u))

```

k	error order 1	error order 2
1	0.0011903876359053434	0.0007912456134713253
2	0.00616183848913525	0.0001857221423487987
3	0.0031287955309935445	4.4984962137673e-05
4	0.001575205457798713	1.1069700795681924e-05
5	0.00079060995047810561	2.745618501776015e-06
6	0.000396083886653649667	6.836951048638795e-07
7	0.0001981906947143841	1.705861063536778e-07
8	9.913803602878968e-05	4.260438157821511e-08

Jacobian test By X for classic

Figure 19: Jacobian transposed test code for X'

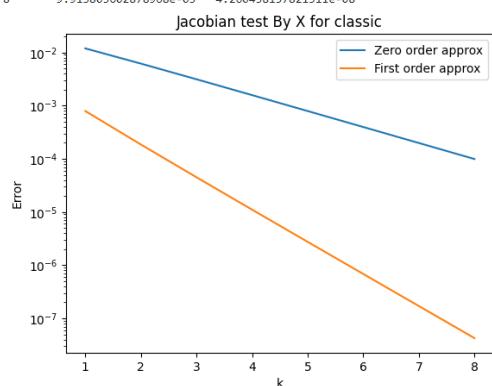


Figure 20: Jacobian by X function + test

5 Part 2.2 :

5.1 One step forward for the residual neural network

The `single_step_forward_res` function computes a forward step in a residual block using a selected activation function, and using the following formula :

$$\text{step_forward} = \mathbf{X} + \mathbf{W}_2 (\sigma(\mathbf{W}_1 \mathbf{X}))$$

Code implementation :

$$Z_{\text{res1}} = \mathbf{W}_1 \cdot \mathbf{X}$$

$$A_{\text{res1}} = \text{activation}(Z_{\text{res1}})$$

$$Z_{\text{res2}} = \mathbf{W}_2 \cdot A_{\text{res1}}$$

return: $Z_{\text{res2}} + \mathbf{A}$

```
def single_step_forward_res(X, W1, W2, activation='relu'):
    activation_func = relu if activation == 'relu' else tanh

    #step 1 :
    Z_res1 = W1.dot(X)
    A_res1 = activation_func(Z_res1)

    #step 2 :
    Z_res2 = W2.dot(A_res1)

    A = X[:-1, :]
    Z_res2 += A

    return Z_res2
```

Figure 21: Resnet single step forward code

5.2 Verification of Derivatives via Jacobian Test for residual network

This section verifies that the derivatives with respect to \mathbf{X} , \mathbf{W}_1 , and \mathbf{W}_2 in the ResNet model successfully pass the Jacobian test.

Derivative with respect to \mathbf{W}_1

$$\frac{\partial y}{\partial \mathbf{W}_1}^T \mathbf{v} = (\sigma'(\mathbf{W}_1 \mathbf{X}) \circ (\mathbf{W}_2^T \mathbf{u})) \mathbf{X}^T$$

Code implementation :

$$A_{\text{res1}} = \text{activation derivative}(\mathbf{W}_1 \mathbf{X})$$

$$Z_{\mathbf{u}} = (\mathbf{W}_2^T \mathbf{u})$$

$$\text{return: } (A_{\text{res1}} \circ Z_{\mathbf{u}}) \mathbf{X}^T$$

```
def jacobian_test_res_w1():
    # Jacobian Check Setup
    num_classes = 3
    input_size = 2
    num_features = num_classes
    EPSILON = 0.1
    W1 = np.random.randn(num_classes, num_features + 1)
    W2 = np.random.randn(num_classes, num_features )
    X = np.random.randn(num_features, input_size)
    X = np.vstack([X, np.ones((1, X.shape[1]))])
    d = np.random.randn(num_classes, num_features + 1)
    u = np.random.randn(num_classes, input_size)

    def F(W1):
        return single_step_forward_res(X, W1, W2, tanh)

    def g_F(W1):
        return np.dot(u.flatten(), F0.flatten())

    F0 = F(W1)
    g0 = g_F(W1)
    grad_g0 = JacobianT_W1_Res(X, W1, W2, u, tanh)

    y0 = np.zeros(8)
    y1 = np.zeros(8)

    print("k\terror order 1 \t\t\t error order 2")
    for k in range(1, 9):
        epsk = EPSILON * (0.5 ** k)
        fk = F(W1 + (epsk * d))
        gk = np.dot(fk.flatten(), u.flatten())
        g1 = g0 + (epsk * (np.dot(d.flatten(), grad_g0.flatten())))
        y0[k - 1] = abs(gk - g0)
        y1[k - 1] = abs(gk - g1)
        print(f'{k}\t{y0[k-1]}\t{y1[k-1]}')

    # Plotting
    plt.semilogy(range(1, 9), y0, label="Zero order approx")
    plt.semilogy(range(1, 9), y1, label="First order approx")
    plt.legend()
    plt.title("Jacobian test By W1 for ResNet")
    plt.xlabel("k")
    plt.ylabel("Error")
    plt.show()

def JacobianT_W1_Res(X, W1, W2, u, activation='relu'):
    activationD_func = relu_derivative if activation == 'relu' else tanh_derivative
    Z_res1 = W1.dot(X)
    A_res1 = activationD_func(Z_res1)

    Z_u = W2.T.dot(u)
    Z_res2 = np.multiply(A_res1, Z_u)

    return Z_res2.dot(X.T)
```

k	error order 1	error order 2
1	0.11309063807108277	0.002403180039270758
2	0.05595436527765063	0.0006106362617446237
3	0.027825704396100548	0.00015383988814754446
4	0.01387453673855905	3.869447658254884e-05
5	0.006927635115197006	9.668988208311191e-06
6	0.003461482532044744	2.414665584052e-06
7	0.001730096675813364	6.05144066412322e-07
8	0.0008648970864602035	1.5132058628353207e-07

jacobian test By W1 for ResNet

Figure 22: Jacobian transposed test code for \mathbf{W}_1

Figure 23: Jacobian by \mathbf{W}_1 function + test

Derivative with respect to W_2

$$\frac{\partial y}{\partial \mathbf{W}_2}^T \mathbf{v} = \mathbf{u} \left(\sigma(\mathbf{W}_1 \mathbf{X})^T \right)$$

Code implementation :

$$A = \text{activation}(\mathbf{W}_1 \mathbf{X})$$

return: (\mathbf{uA}^T)

Figure 24: Jacobian transposed test code for W2'

Test results : In the test, the differing results can be attributed to the variable $W_2 \cdot Z$ not undergoing a linear process.



Figure 25: Jacobian by W2 function + test

Derivative with respect to X

$$\frac{\partial y}{\partial \mathbf{X}}^T \mathbf{v} = \mathbf{u} + \mathbf{W}_1^T (\sigma'(\mathbf{W}_1 \mathbf{X}) \circ (\mathbf{W}_2^T \mathbf{u}))$$

Code implementation :

$$A_{\text{res1}} = \text{activation derivative}(\mathbf{W}_1 \mathbf{X})$$

$$Z_{\mathbf{u}} = \mathbf{W}_2^T \mathbf{u}$$

$$\text{Jac} = \mathbf{W}_1^T (A_{\text{res1}} \circ Z_{\mathbf{u}})$$

return: $\text{Jac} + \mathbf{u}$

The line $\text{Jac} = \text{Jac}[:-1,:]$ removes the last row from the Jacobian matrix Jac . This row corresponds to the gradients related to the bias term, which is always set to one and does not change during training.

The figure shows a log-linear plot of Error versus k. The x-axis (k) ranges from 1 to 8, and the y-axis (Error) ranges from 10^{-1} to 10^{-6} . Two data series are plotted: 'Zero order approx' (blue line) and 'First order approx' (orange line). The 'First order approx' line shows a steeper negative slope, indicating faster convergence than the 'Zero order approx' line.

k	Zero order approx	First order approx
1	0.13425010692488937	0.00810507824945117
2	0.06507505468657748	0.002025403488583837
3	0.03203329582793191	0.0004970386590723574
4	0.01589109035307812	0.00012377476887803773
5	0.007914945288419517	3.0880996204629696e-05
6	0.003949744417679875	7.71228560543827e-06
7	0.00197294314266059	1.927668212337346e-06
8	0.000985989677696919	4.816411708581114e-07

Figure 26: Jacobian transposed test code for X'

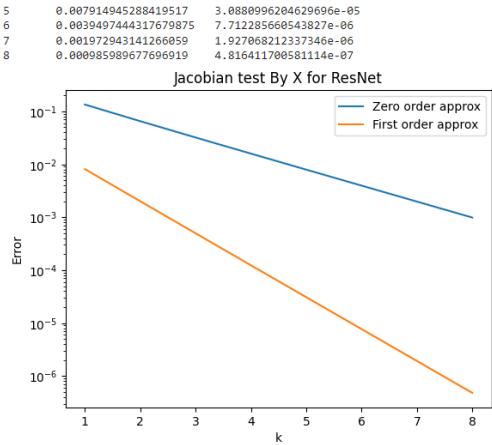


Figure 27: Jacobian by X function + test

6 Part 2.3 :

6.1 Initialize network

The Initialize network function configures a dictionary network with weight matrices for a neural network. It creates standard weights sized by each layer and its previous layer, and for residual layers, it generates two weights ensuring size compatibility between the intermediate and main layers.

6.2 Forward pass

The forward pass function conducts the data propagation through a neural network, applying different processing steps for standard and residual layers, and stores the intermediate outputs.

network : The weight matrices for each layer initialized by initialize network.

X : The initial data matrix fed into the network.

activation : Specifies the activation function (relu or tanh) used.

res places : A list identifying which layers have a residual structure for alternative processing steps.

Outputs:

A : The output from the last layer after applying the softmax function, representing the network's final prediction.

cache: A dictionary storing outputs from each layer, used later in the back propagation process.

```
def forward_pass(network, X, activation='relu', res_places=[]):
    activation_func = relu if activation == 'relu' else tanh
    cache = {'A0': X}
    A = X
    L = len(network) - (len(res_places))
    for i in range(1, L):
        # bias
        A_aug = np.vstack([A, np.ones((1, A.shape[1]))])
        # resNet
        if i-1 in res_places:
            W1 = network['W' + str(i) + '_res1']
            W2 = network['W' + str(i) + '_res2']
            A = single_step_forward_res(A_aug, W1, W2, activation)

        # classic
        else:
            W = network['W' + str(i)]
            A = single_step_forward_class(A_aug, W, activation)

        # caching
        cache['A' + str(i)] = A

    #last layer
    A_aug = np.vstack([A, np.ones((1, A.shape[1]))])
    A = softmax(network['W' + str(L)].dot(A_aug))
    cache['A' + str(L)] = A

    return A, cache
```

Figure 28: Forward pass

6.3 Back propagation

The back propagation function calculates gradients for updating the weights based on the error between the predicted outputs and actual outputs. It utilizes derivatives defined and tested above, operating on both standard and residual layers.

gradients : A dictionary to store the computed gradients for each weight matrix in the network.

Y : The actual output matrix used for training, against which the predicted outputs are compared.

return : The updated dictionary containing gradients for all weight matrices

```
def backpropagation(gradients,network, cache, Y, res_places = [], activation='relu'):  
    activation_func = relu if activation == 'relu' else tanh  
  
    L = len(network) - (len(res_places))  
    A = cache['A' + str(L - 1)]  
    A = np.vstack([A, np.ones((1, A.shape[1]))])  
    W_prime = softmax_regression_gradient_by_W(A ,network['W'+str(L)], Y)  
    gradients['W'+str(L)] = W_prime  
    dA = softmax_regression_gradient_by_X(A ,network['W'+str(L)], Y)  
    dA = dA[:-1,:]  
  
    for i in reversed(range(1, L)):  
        A = cache['A'+str(i-1)]  
        A_aug = np.vstack([A, np.ones((1, A.shape[1]))])  
  
        #resBack  
        if i - 1 in res_places:  
            W1 = network['W'+str(i)+'_res1']  
            W2 = network['W'+str(i)+'_res2']  
  
            W1_prime = JacobianT_W1_Res(A_aug, W1, W2, dA, activation)  
            gradients['W'+str(i)+'_res1'] = W1_prime  
  
            W2_prime = JacobianT_W2_Res(A_aug, W1, W2, dA, activation)  
            gradients['W'+str(i)+'_res2'] = W2_prime  
  
            dA = JacobianT_X_Res(A_aug, W1, W2, dA, activation)  
  
        #clasicBack  
        else:  
            W = network['W'+str(i)]  
            W_prime = JacobianTByW(A_aug, W, dA, activation)  
            gradients['W'+str(i)] = W_prime  
            dA = JacobianTByX(A_aug, W, dA, activation)  
            dA = dA[:-1,:]  
  
    return gradients
```

Figure 29: Back propagation

6.4 Gradient test for the complete network

The graph below illustrates the results of our gradient testing across all layers of the neural network, following both the forward pass and back propagation processes.

```

def Gradient_test_for_complete_network():
    # Gradient test Setup
    layer_sizes = [2,4,4,4,4,4,2]
    res_places = [1,2,3,4,5]
    batchSize = 3
    gradients = {}
    EPSILON = 0.1
    network = initialize_network(layer_sizes, res_places)
    X, C = pickRandomBatch(data_training_X, data_training_C, batchSize)

    # Generate a random vector d for each layer
    d = {layer: np.random.randn(network[layer].shape[0], network[layer].shape[1]) for layer in network}

    def F(network):
        A, cache = forward_pass(network, X, 'tanh', res_places)
        return compute_loss(A,C), cache

    def g_F(network, cache):
        return backpropagation(gradients, network, cache, C, res_places, 'tanh')

    F0, cache = F(network)
    g0 = g_F(network, cache)
    y0 = np.zeros(8)
    y1 = np.zeros(8)

    for k in range(1, 9):
        epsk = EPSILON * (0.5 ** k)
        # Perturb weights by epsk * d
        W_original = {layer: np.copy(network[layer]) for layer in network}
        for layer in network:
            network[layer] += epsk * d[layer]

        Fk, _ = F(network)
        temp = 0
        for key in g0:
            temp += np.sum(g0[key] * (epsk * d[key]))
        F1 = F0 + temp
        y0[k-1] = abs(Fk - F0)
        y1[k-1] = abs(Fk - F1)
        # Reset weights
        for layer in network:
            network[layer] = W_original[layer]

        print(f"(k)\t{y0[k-1]}\t{y1[k-1]}")

    # Plotting
    plt.semilogy(range(1, 9), y0, label="Zero order approx")
    plt.semilogy(range(1, 9), y1, label="First order approx")
    plt.legend()
    plt.title("Gradient Test for complete network")
    plt.xlabel("k")
    plt.ylabel("Error")
    plt.show()

```

Figure 30: Jacobian transposed test code for complete NN

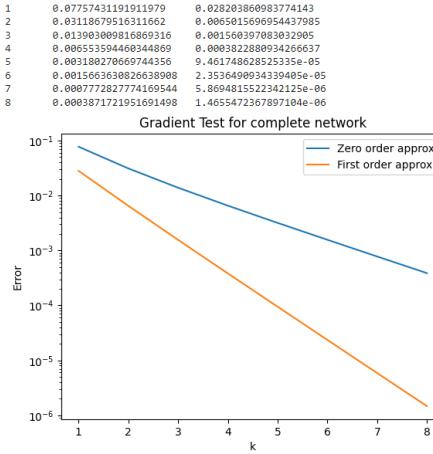


Figure 31: Test results of Jacobian transposed for complete NN

7 Part 2.4 :

7.1 train network

The train network function manages the training of a neural network through multiple iterations, using specified hyper parameters and settings. It performs forward and backward passes, applies stochastic gradient descent with momentum, and calculates losses and success rates for each epoch. The function outputs the trained network along with the success rates and loss values per iteration. These metrics are collected to evaluate and visualize the network's performance.

7.2 data validation test

The data validation test function evaluates the performance of the trained neural network on the validation dataset. The function takes the trained network and processes the validation data through the network. returns the calculated loss and success rate.

7.3 Optimal Network Configurations: Evaluating Different Lengths

We conducted experiments with three different neural network architectures, varying the number of layers to study the impact on training performance and epoch requirements. The details of the hyperparameters used in each configuration are as follows:

Tests on "Swiss Roll"

Network Configurations:

- **Batch Size:** 64
- **Learning Rate:** 0.09
- **Momentum:** 0.2
- **Activation Function:** ReLU
- **Epochs Tested:** 25, 50, 100, 200

Each network configuration also employed a specific architecture setup:

- Each network started with a standard layer that scaled from an input size of 2 units to 8 units.
- The middle layers consisted of residual network (ResNet) blocks with 8 units.
- The final layer was a standard layer that reduced the size back to 2 units.

First Test: 4 layers

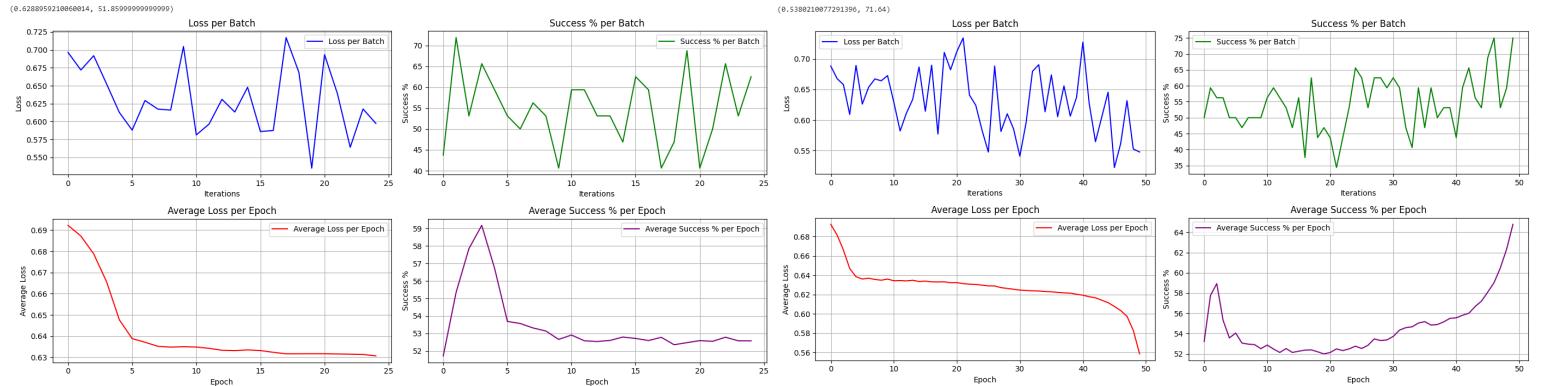


Figure 32: Results for 25 epochs - 51.8% success

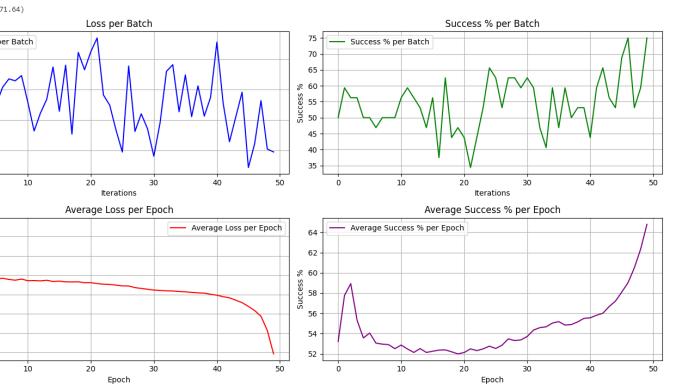


Figure 33: Results for 50 epochs - 71.6% success

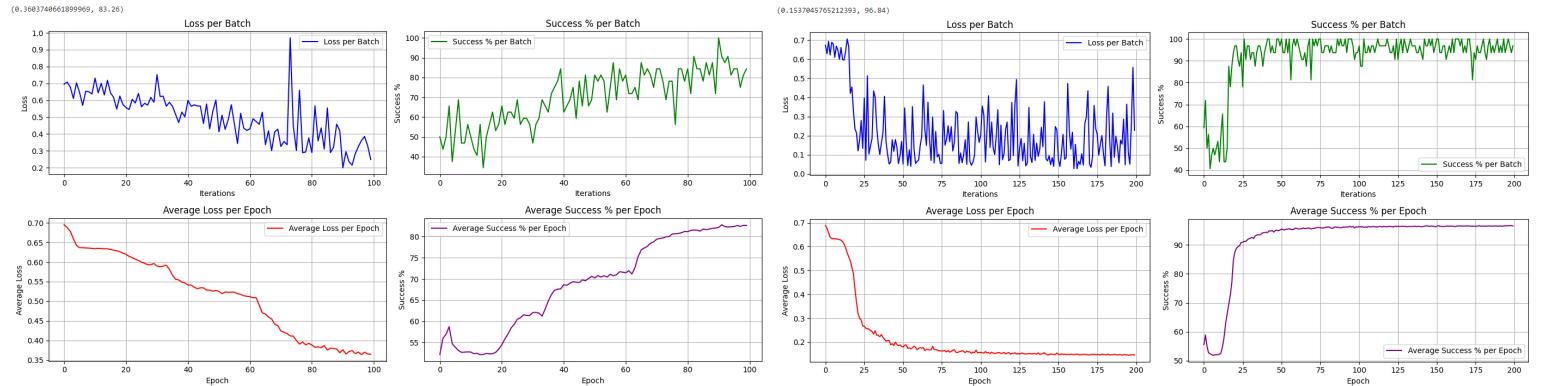


Figure 34: Results for 100 epochs - 83.2% success

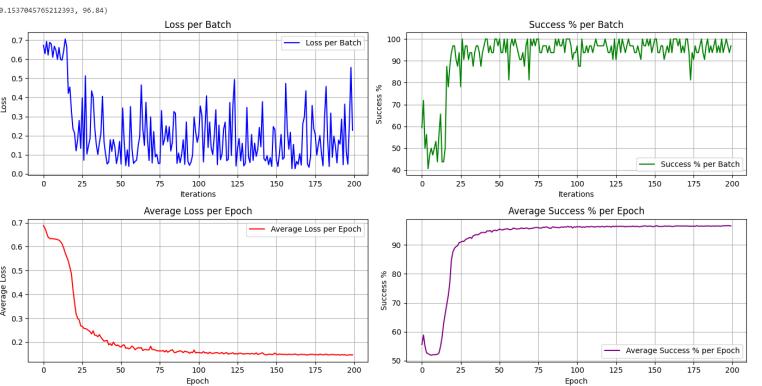


Figure 35: Results for 200 epochs - 96.8% success

Second Test: 7 layers

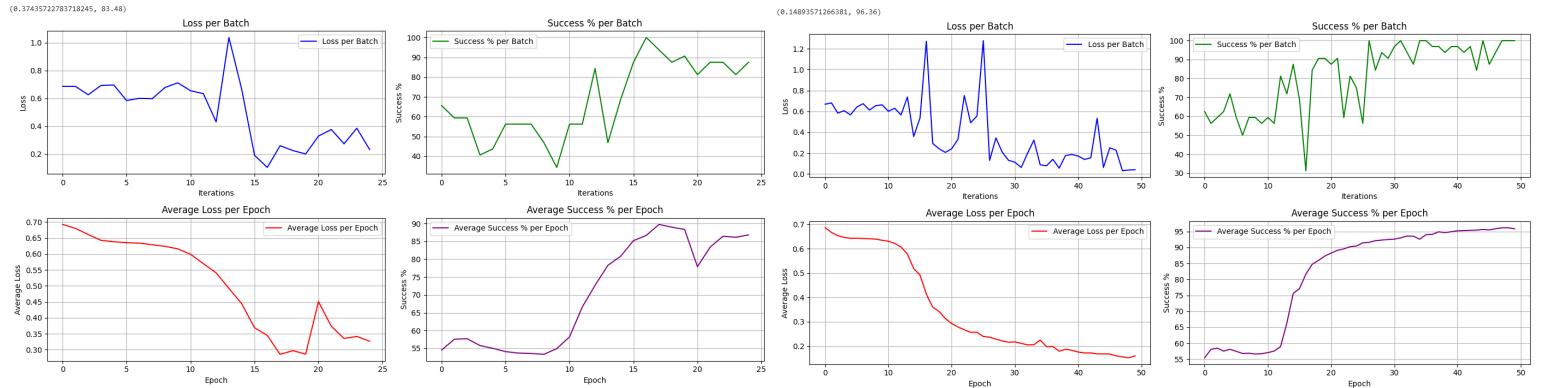


Figure 36: Results for 25 epochs - 83.4% success

Figure 37: Results for 50 epochs - 96.3% success

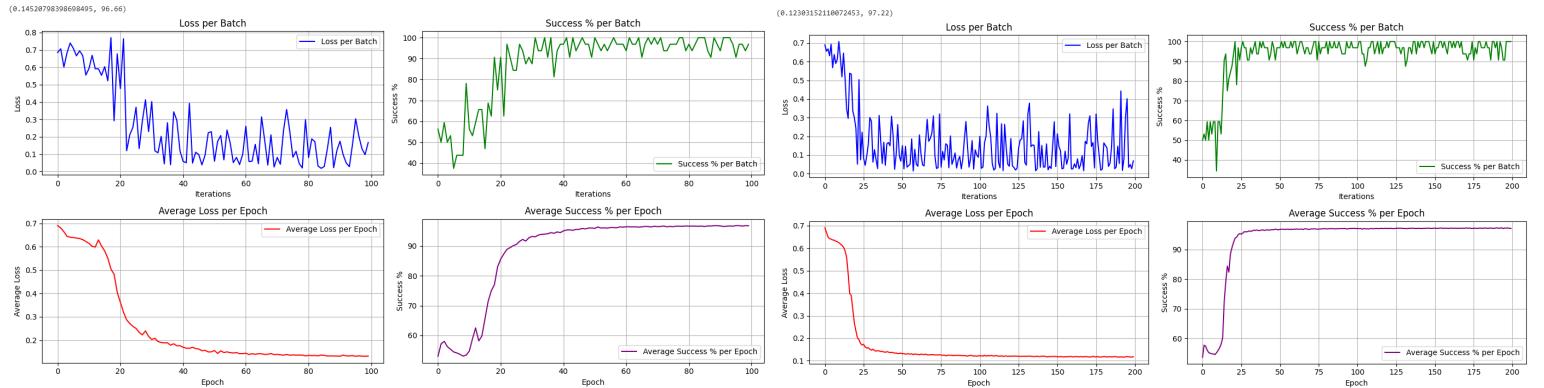


Figure 38: Results for 100 epochs - 96.6% success

Figure 39: Results for 200 epochs - 97.2% success

Third Test: 12 layers

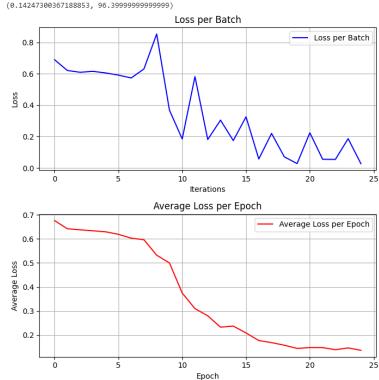


Figure 40: Results for 25 epochs - 96.3% success

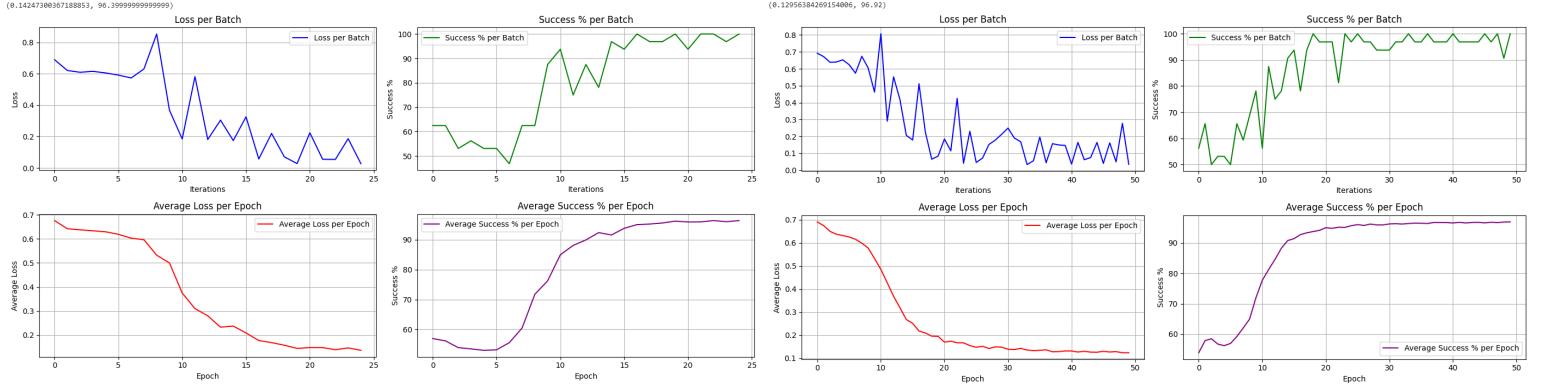


Figure 41: Results for 50 epochs - 96.9% success

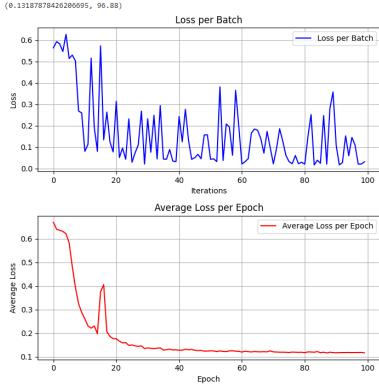


Figure 42: Results for 100 epochs - 96.8% success

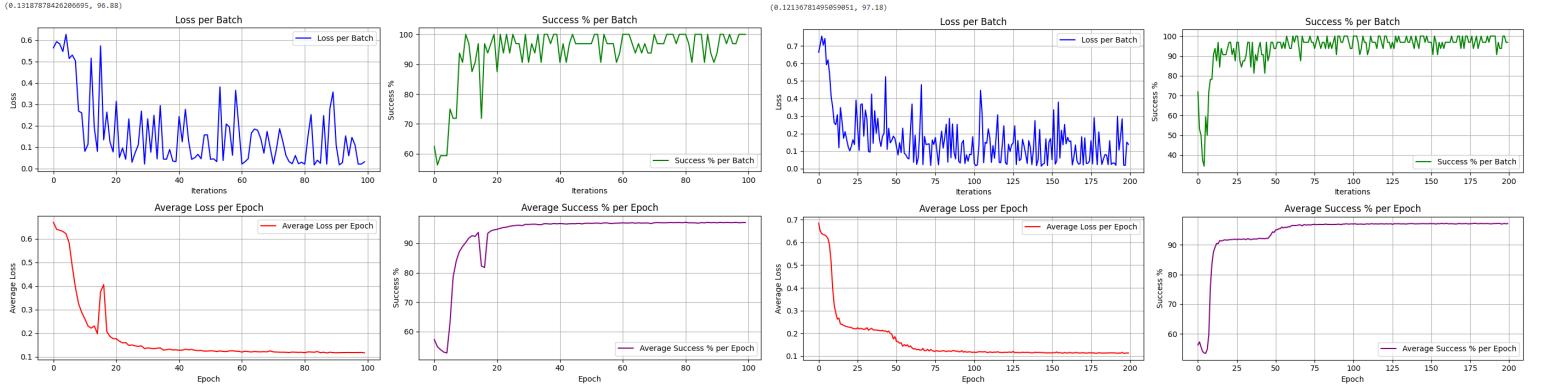


Figure 43: Results for 200 epochs - 97.18% success

In our experiments with varying sizes, we observed that the smaller NN required approximately 200 epochs to achieve satisfactory results. Conversely, the mid-sized layer NN reached comparable outcomes within just 50 epochs. The largest NN configuration completed training in merely 25 epochs, achieving similar results, however, additional training for another 175 epochs only improved performance marginally by approximately 0.4% to 1%.

Conclusion : While larger neural networks are more efficient in terms of epochs needed, they are costlier in terms of time and space. The medium-sized network presents a good compromise, balancing performance, time, and resource demands

The following section presents additional experiments where we tested various layer dimensions and batch sizes.

Tests on "GMM" - layer sizes

Network Configurations:

- **Number of layers:** 5
- **Batch Size:** 64
- **Learning Rate:** 0.09
- **Momentum:** 0.2
- **Activation Function:** ReLU
- **Epochs Tested:** 50, 100

Each network configuration also employed a specific architecture setup:

- Each network started with a standard layer that scaled from an input size of 5 units to 4/8/16 units.
- The middle layers consisted of residual network (ResNet) blocks with 4/8/16 units.
- The final layer was a standard layer that reduced the size back to 5 units.

First Test: 4 units per hidden layer

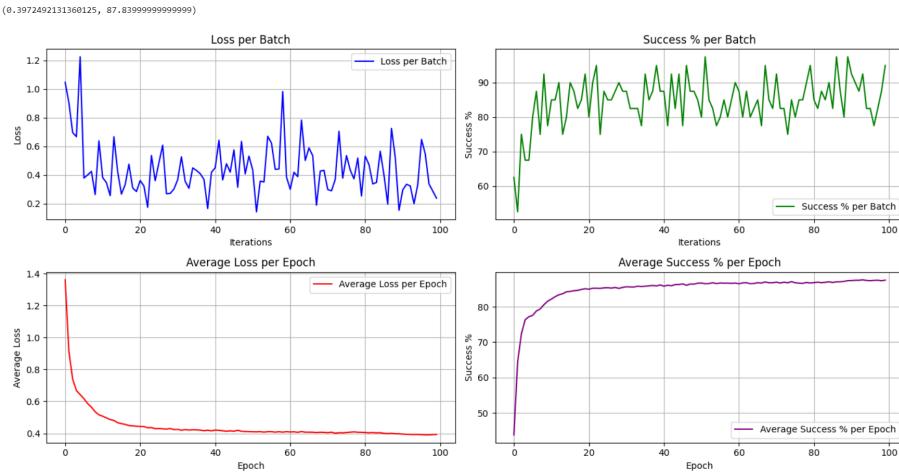


Figure 44: 4 units per hidden layer - 87.8% success

Second Test: 8 units per hidden layer

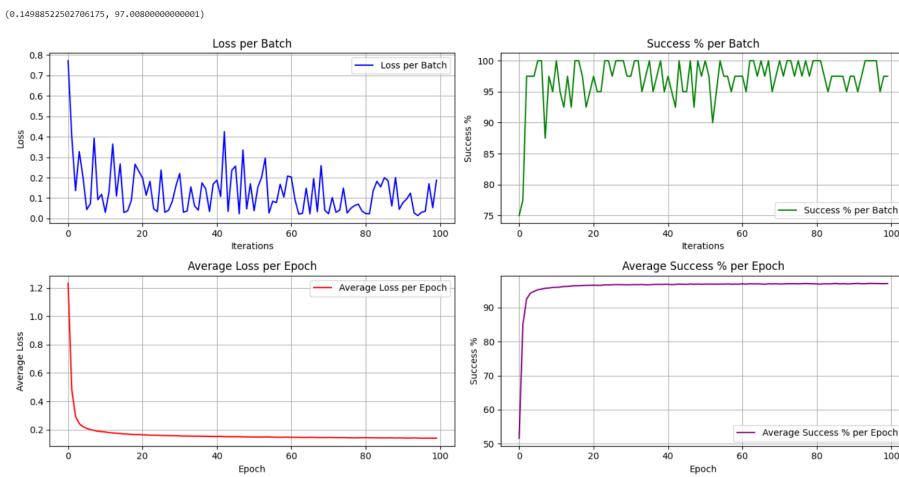


Figure 45: 8 units per hidden layer - 97% success

Third Test: 16 per hidden layer

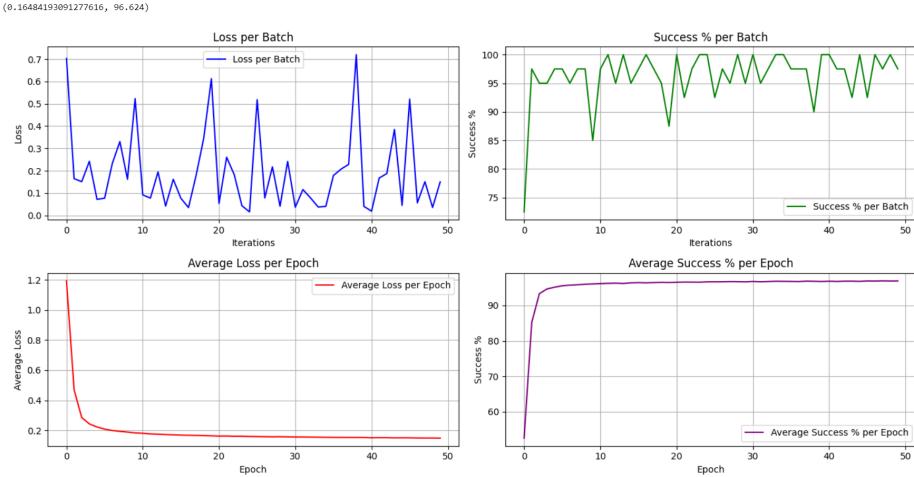


Figure 46: 16 units per hidden layer - 96.6% success

Conclusion : Our findings indicate that larger layers in neural networks yield higher performance outcomes; however, they also demand greater costs in terms of computational time and resource usage. Additionally, beyond a certain point, increasing the layer size does not result in improved results, suggesting a threshold beyond which further enlargement is inefficient.

Tests on "GMM" - batch sizes

Network Configurations:

- **Number of layers:** 5
- **Learning Rate:** 0.09
- **Momentum:** 0.2
- **Activation Function:** ReLU
- **Epochs:** 100

Each network configuration also employed a specific architecture setup:

- Each network started with a standard layer that scaled from an input size of 5 units to 16 units.
- The middle layers consisted of residual network (ResNet) blocks with 16 units.
- The final layer was a standard layer that reduced the size back to 5 units.

First Test: 250 Batch size

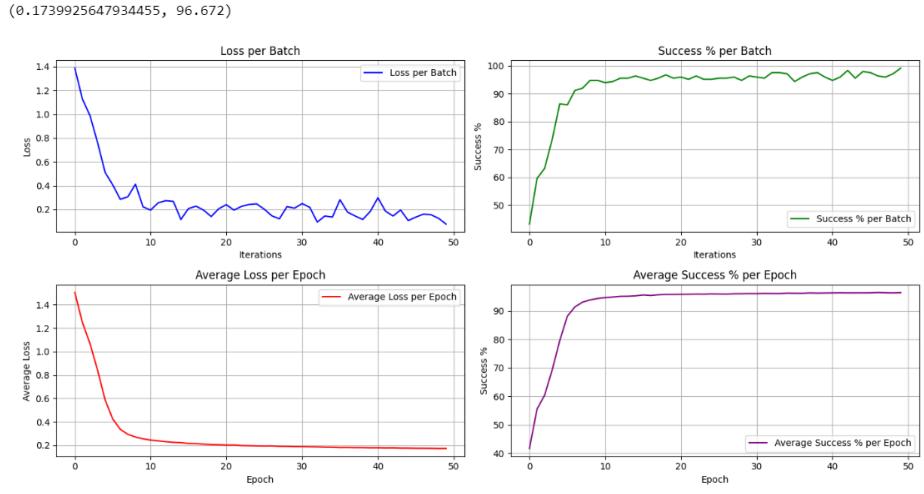


Figure 47: Batch size 250 - 96.6% success

Second Test: 1000 Batch size

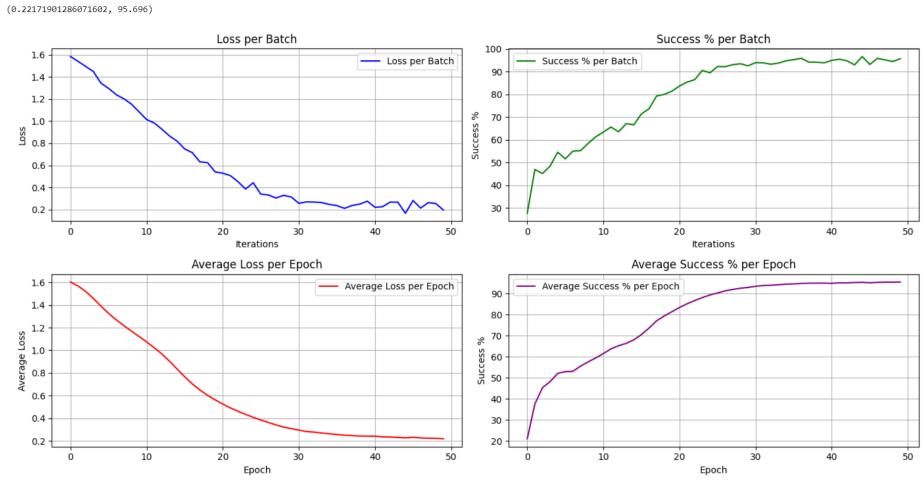


Figure 48: Batch size 1000 - 95.6% success

Third Test: 2500 Batch size

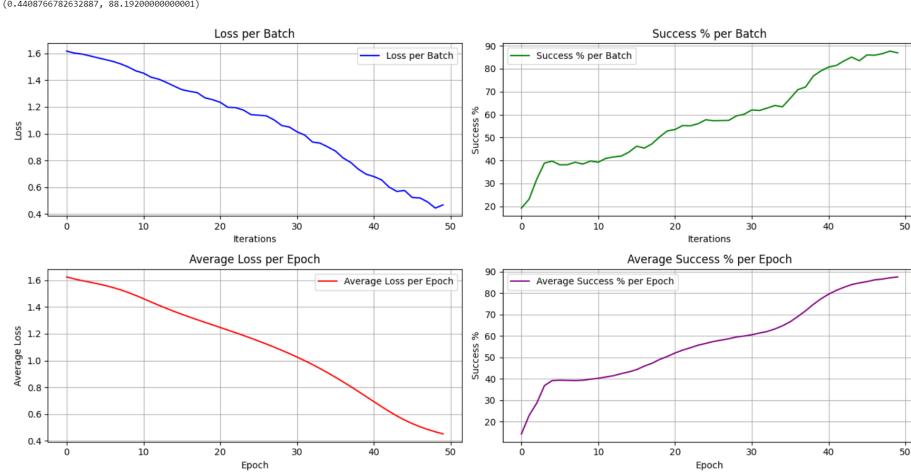


Figure 49: Batch size 2500 - 88.1% success

Fourth Test: 5000 Batch size

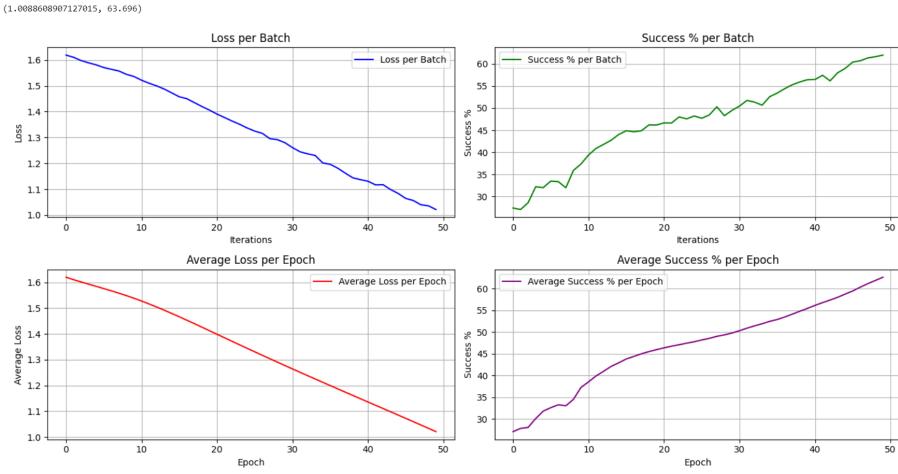


Figure 50: Batch size 5000 - 63.6% success

Conclusion : Our experiments reveal that smaller batch sizes generally lead to better results. Furthermore, excessively large batch sizes require significantly more epochs to train effectively. The graph above illustrates that networks with smaller batch sizes not only minimize loss more effectively but also achieve faster increases in success percentage. Based on these observations, we conclude that a much smaller batch size is preferable for efficient training and enhanced performance.

8 Part 2.5

8.1 Optimizing Neural Network Efficiency - "Swiss Roll"

Next, we focused on maximizing performance while keeping the network size below the threshold of $100C$ scalar parameters in terms of weight count. We explored various configurations, ranging from 'long' networks with many layers to more compact networks with up to four layers. The optimal configuration emerged as a balance between network length and weight density—neither the longest nor the heaviest in terms of weights.

Optimal Performance Achieved:

- **Layer Sizes:** [2, 6, 6, 6, 6, 2].
- **Batch Size:** 64 .
- **Number of Epochs:** 40 .
- **Learning Rate:** 0.09,
- **Momentum:** 0.2.
- **Activation Function:** 'ReLU'.

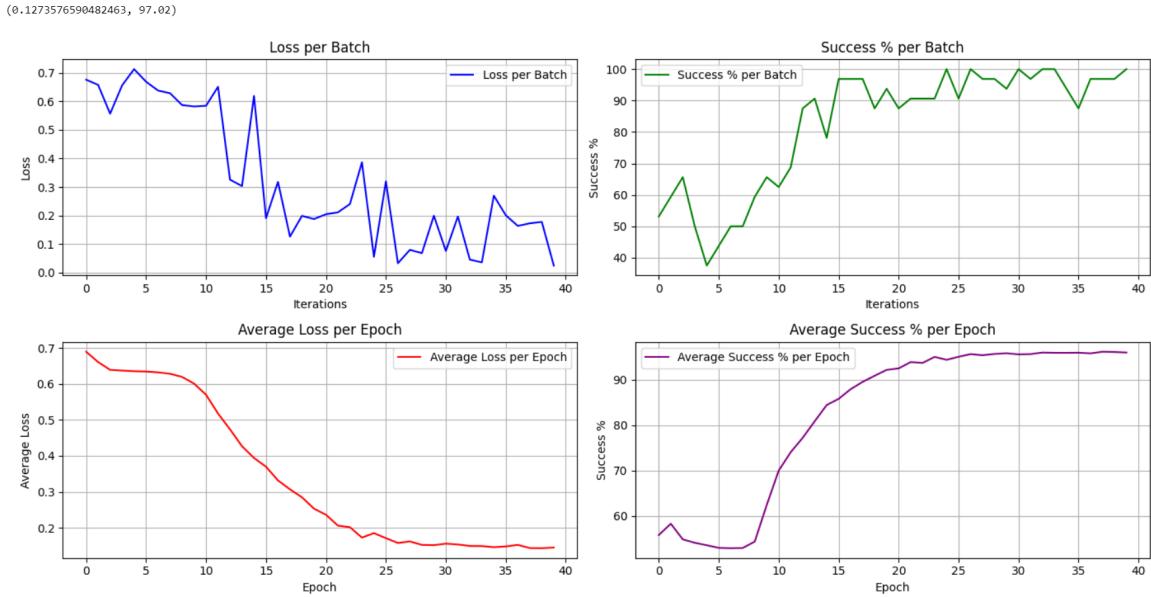


Figure 51: Peak performance achieved with a success rate of 97.0%

To calculate the total number of parameters in our neural network, we multiply the number of neurons in each layer by the number of neurons in the preceding layer to determine the number of weights between them. This calculation is performed for each consecutive layer pair, and the resulting values are then summed to obtain the total parameter count. Our objective was to identify the architecture that delivers the best performance using the fewest epochs, ensuring an efficient training process.

Another notable result is :

- **Layer Sizes:** [2, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 4, 2].
- **Batch Size:** 64 .
- **Number of Epochs:** 40 .
- **Learning Rate:** 0.09,
- **Momentum:** 0.2.
- **Activation Function:** 'ReLU'.

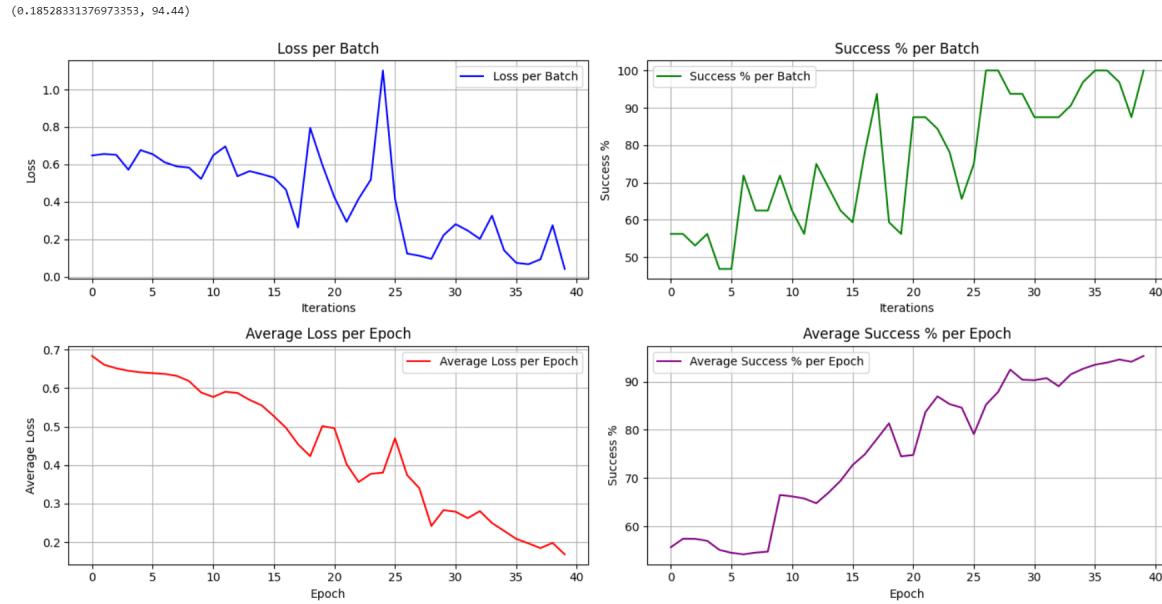


Figure 52: Peak performance achieved with a success rate of 94.4%

9 Part 2.5

9.1 Optimizing Neural Network Efficiency - "GMM"

We conducted the same experiments on the GMM dataset and achieved notable outcomes. Here are the details of the best results :

Optimal Performance Achieved:

- **Layer Sizes:** [5, 17, 17, 5].
- **Batch Size:** 64 .
- **Number of Epochs:** 10 .
- **Learning Rate:** 0.09,
- **Momentum:** 0.2.
- **Activation Function:** 'ReLU'.

(0.20173706958750973, 96.096)

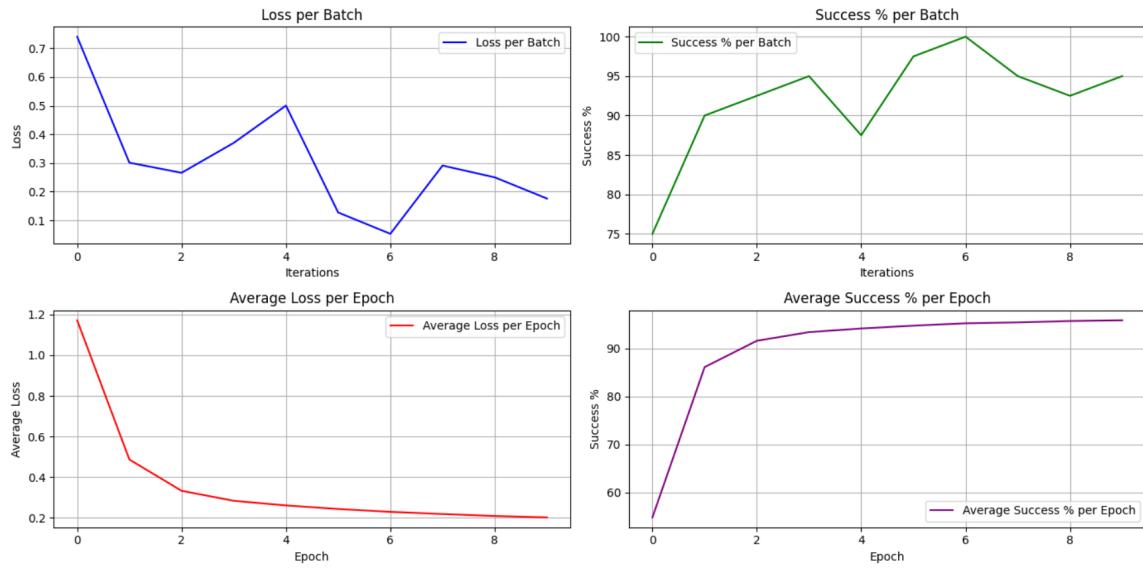


Figure 53: Peak performance achieved with a success rate of 96.0%

Our experiments demonstrated that using the initial settings allowed the network to reach optimal performance within just 10 epochs. Alternative configurations failed to achieve similar success rates and consistently required a significantly higher number of epochs to converge.

Another notable result is :

- **Layer Sizes:** [5, 13, 13, 13, 5].
- **Batch Size:** 64 .
- **Number of Epochs:** 20 .
- **Learning Rate:** 0.09,
- **Momentum:** 0.2.
- **Activation Function:** 'ReLU'.

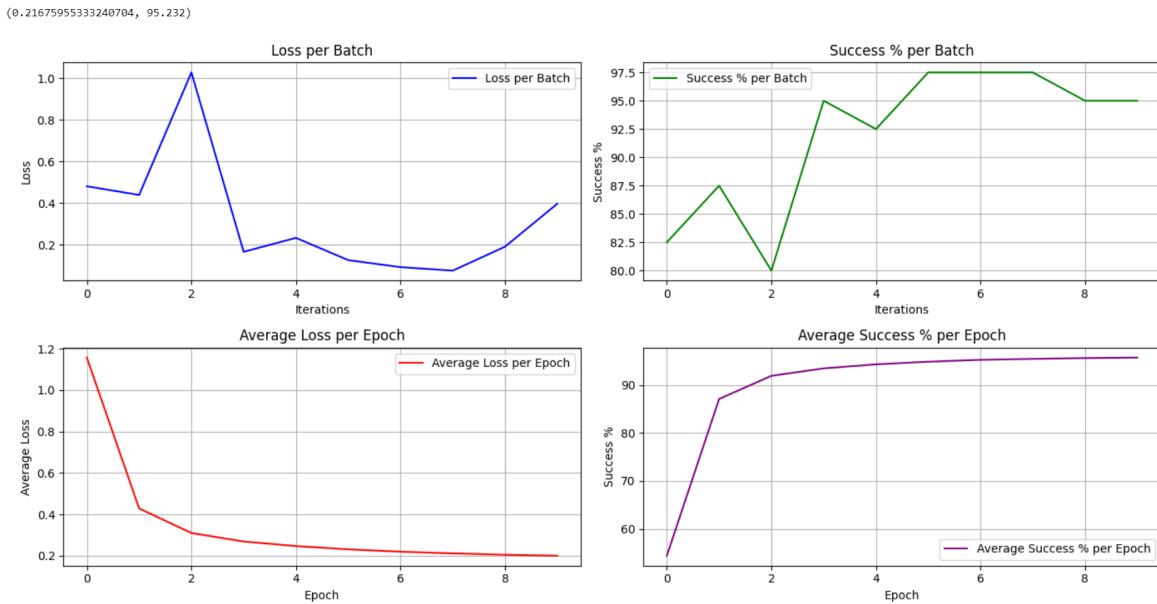


Figure 54: Peak performance achieved with a success rate of 95.2%

10 Code implementation here

[Our Code implementation here](#)