

**Q 1. Rosenbrock's Valley Problem**

Consider the Rosenbrock's Valley function:

$$f(x, y) = (1 - x)^2 + 100(y - x^2)^2$$

- (a) Show that it has a global minimum at  $(x, y) = (1, 1)$  where  $f(x, y) = 0$ . (5 Marks)

**Solution:**

1. The Rosenbrock's Valley function is nonnegative, i.e.,  $f(x, y) \geq 0, \forall (x, y) \in \mathbb{R}^2$ ;
2. The  $f(x, y) = 0$  holds **if and only if** both terms are zero:

$$\begin{cases} 1 - x = 0 \\ y - x^2 = 0 \end{cases} \Rightarrow \begin{cases} x = 1 \\ y = 1^2 = 1 \end{cases}.$$

3. Find critical point by calculating first-order partial derivatives:

$$\nabla f = \begin{bmatrix} -2(1 - x) - 400x(y - x^2) \\ 200(y - x^2) \end{bmatrix}$$

Set  $\nabla f = 0$ , we can obtain  $x = y = 1$ , the only critical point is  $(1, 1)$ .

4. Check whether Hessian Matrix  $H(1, 1)$  is positive definite:

The second-order partial derivatives:

$$\frac{\partial^2 f}{\partial x^2} = 2 + 400(3x^2 - y), \quad \frac{\partial^2 f}{\partial y^2} = 200, \quad \frac{\partial^2 f}{\partial x \partial y} = \frac{\partial^2 f}{\partial y \partial x} = -400x.$$

So Hessian Matrix at  $(1, 1)$ :

$$H(1, 1) = \begin{bmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial y \partial x} & \frac{\partial^2 f}{\partial y^2} \end{bmatrix}_{(1,1)} = \begin{bmatrix} 802 & -400 \\ -400 & 200 \end{bmatrix}.$$

$Tr(H) = 802 + 200 = 1002 > 0$ ,  $det(H) = 802 * 200 - (-400)^2 = 400 > 0$ , solve  $\lambda^2 - Tr(H)\lambda + det(H) = 0$ :  $\lambda_1 = 1001.6$ ,  $\lambda_2 = 0.4$  (both positive). Since all eigenvalues are positive,  $H(1, 1)$  is **positive definite**, confirming  $(1, 1)$  as the local minimum.

Thus,  $(1, 1)$  is the only point where  $f(x, y) = 0$ , confirming it as the global minimum.  $\square$

Now suppose that the starting point is randomly initialized in the open interval  $(-1, 1)$  for  $x$  and  $y$ , find the global minimum using:

- (b) **Steepest (Gradient) descent method:**

$$w(k+1) = w(k) - \eta g(k)$$

with learning rate  $\eta = 0.001$ . Here the weight vector refers to the two-dimensional vector  $[x, y]$ . Record the number of iterations when  $f(x, y)$  converges to (or very close to) 0 and plot out the trajectory of  $(x, y)$  in the 2-dimensional space. Also plot out the function value as it approaches the global minimum. What would happen if a larger learning rate, say  $\eta = 1.0$ , is used? (5 Marks)

**Solution:**

The Gradient Descent method was implemented to minimize the Rosenbrock function:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \eta \nabla f(\mathbf{w}^{(k)}),$$

where:  $\mathbf{w} = [x, y]^T$  is the parameter vector.

$\eta = 0.001$  (learning rate).

The gradient  $\nabla f$  was calculated analytically (see Q1a).

The algorithm terminated when the update norm fell below  $10^{-8}$ .

Listing 1: MATLAB script (Q1b)

```

1 eta = 0.001; % learning rate
2 % eta = 1.0;
3 w = [rand(); rand()]; % random initialization of the start point
4 max_iter = 1e5;
5 tol = 1e-8;
6 trajectory = zeros(2, max_iter);
7 f_values = zeros(1, max_iter);
8
9 for iter = 1:max_iter
10     x = w(1); y = w(2);
11     grad = [-2*(1 - x) - 400*x*(y - x^2); 200*(y - x^2)];
12     w_new = w - eta * grad;
13
14     trajectory(:, iter) = w;
15     f_values(iter) = (1 - x)^2 + 100*(y - x^2)^2;
16
17     if norm(w_new - w) < tol
18         % if f_values(iter) < tol % f(x, y) close to the target value (0)
19         break;
20     end
21     w = w_new;
22 end
23
24 % Plot
25 fig1 = figure;
26 plot(trajectory(1, 1:iter), trajectory(2, 1:iter), 'ro-');
27 hold on;
28 plot(1, 1, 'c*'); % (1,1) is the global minimum point
29 title(sprintf('Trajectory of Gradient Descent (eta=%.3f, iter=%d)', eta, iter));
30 xlabel('x'); ylabel('y');
31 hold off
32 saveas(fig1, 'Traj_eta0.001.png');
33
34 fig2 = figure;
35 semilogy(f_values(1:iter), 'go-');
36 title(sprintf('Function Value Convergence (eta=%.3f, iter=%d)', eta, iter));
37 xlabel('Iteration'); ylabel('f(x,y)');
38 saveas(fig2, 'Func_value_eta0.001.png');
39
40 %% Plot the Rosenbrock's Valley
41 f = @(x,y) (1 - x).^2 + 100*(y - x.^2).^2;
42 [x,y] = meshgrid(linspace(-2, 2, 400), linspace(-1, 3, 400));
43 z = f(x, y);
44 [dx, dy] = gradient(z); % Gradient
45 grad_mag = sqrt(dx.^2 + dy.^2); % Magnitude of gradient
46
47 fig3 = figure;
48 s = surf(x, y, z, grad_mag, 'EdgeColor', 'none');
49 colormap(turbo);
50 clim([min(grad_mag(:)), max(grad_mag(:))]);
51 colorbar;
52 title('Rosenbrock Valley Colored by Gradient Magnitude: Blue (Flat) -> Red (Steep)')
53 ;
54 xlabel('x');
55 ylabel('y');
56 zlabel('f(x,y)');
57 view(-45, 45);
58 grid on;
59 saveas(fig3, 'Rosenbrock''s valley.png');

```

## Results

As shown in Figure 1, the iterations needed for convergence is 21467. The trajectory in the 2D space (Figure 1a) shows a winding path toward (1, 1), characteristic of gradient descent in narrow valleys. The semi-log plot (Figure 1b) shows an exponential decay of  $f(x, y)$ , reaching  $f \approx 0$  within tolerance.

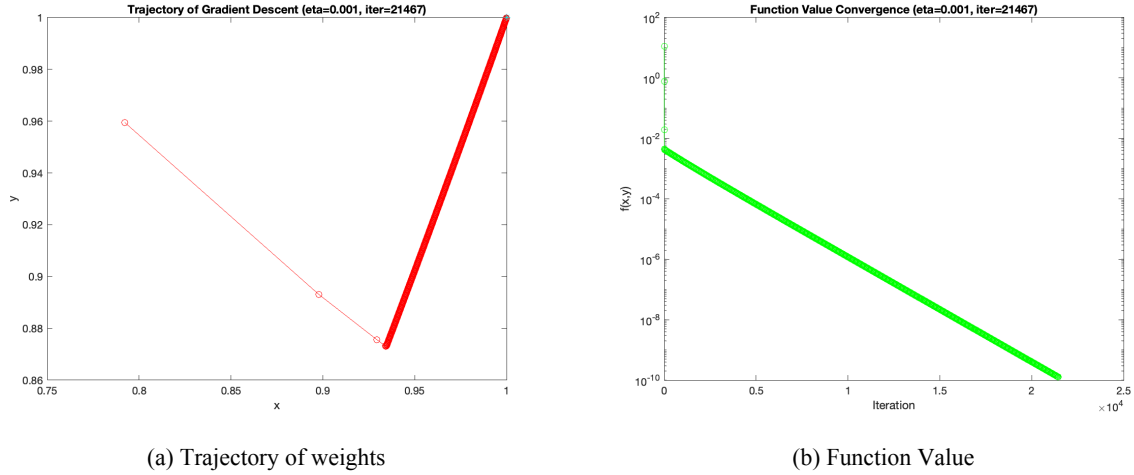


Figure 1: Gradient Descent method ( $\eta=0.001$ )

## Effect of Large Learning Rate ( $\eta = 1.0$ )

With  $\eta = 1.0$ , the algorithm failed to converge due to overshooting. Oscillations in parameter updates caused  $f(x, y)$  to increase indefinitely, as shown in Figure 2b.

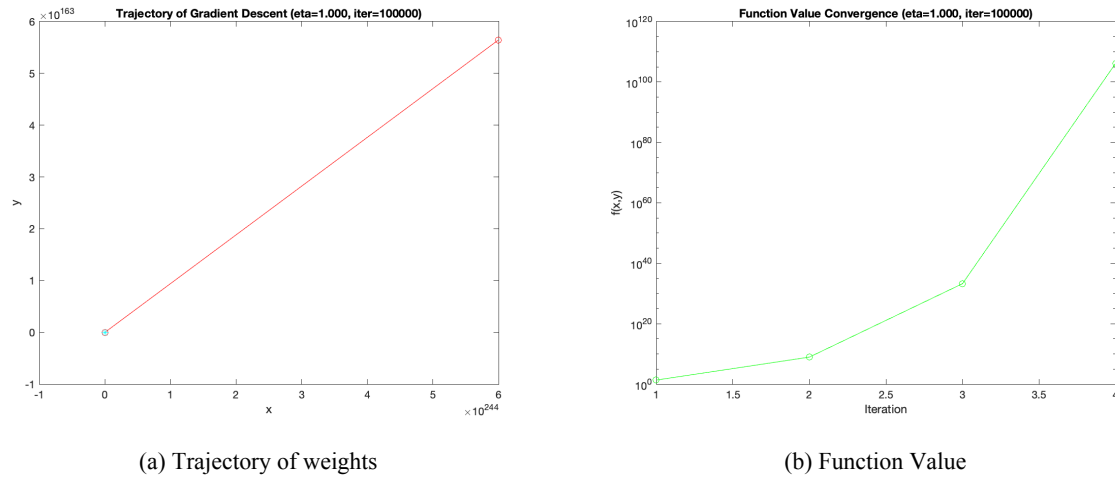


Figure 2: Gradient Descent method ( $\eta=1$ )

## Analysis

The Rosenbrock's valley has a highly non-convex curvature (Figure 3), necessitating small learning rates for stable convergence. Specifically, the extremely narrow width of the valley necessitates small learning rates to prevent parameter updates from overshooting the valley floor. Large step sizes would propel the optimization trajectory beyond the valley boundaries into high-curvature regions of steep gradients. For example, for a current position near  $(x, x^2)$  along the valley floor, an excessive learning rate  $\Delta\eta$  might yield the next iteration at  $(x + \Delta x, (x + \Delta x)^2 + \epsilon)$ , thereby deviating from the optimal valley alignment.

The gradient direction exhibits rapid directional variation along the valley's parabolic path. Large learning rates inadequately compensate for these directional adjustments, inducing oscillatory zigzag patterns in the optimization

path. For instance, at point (0.5, 0.25), the gradient vector primarily aligns with increasing x-direction. However, at (0.6, 0.36), the required gradient correction must reorient toward the (1, 1) minimum, demanding precise directional adaptation.

The curvature perpendicular to the valley axis significantly exceeds that along the valley's longitudinal direction. Large learning rates amplify oscillations in the high-curvature transverse dimension, destabilizing convergence.

This anisotropic curvature profile creates conflicting stability requirements: - Longitudinal updates require sustained momentum for valley traversal - Transverse updates demand strict damping to prevent overshooting

The Rosenbrock function's pathological landscape thus epitomizes the delicate balance required in learning rate selection for non-convex optimization. □

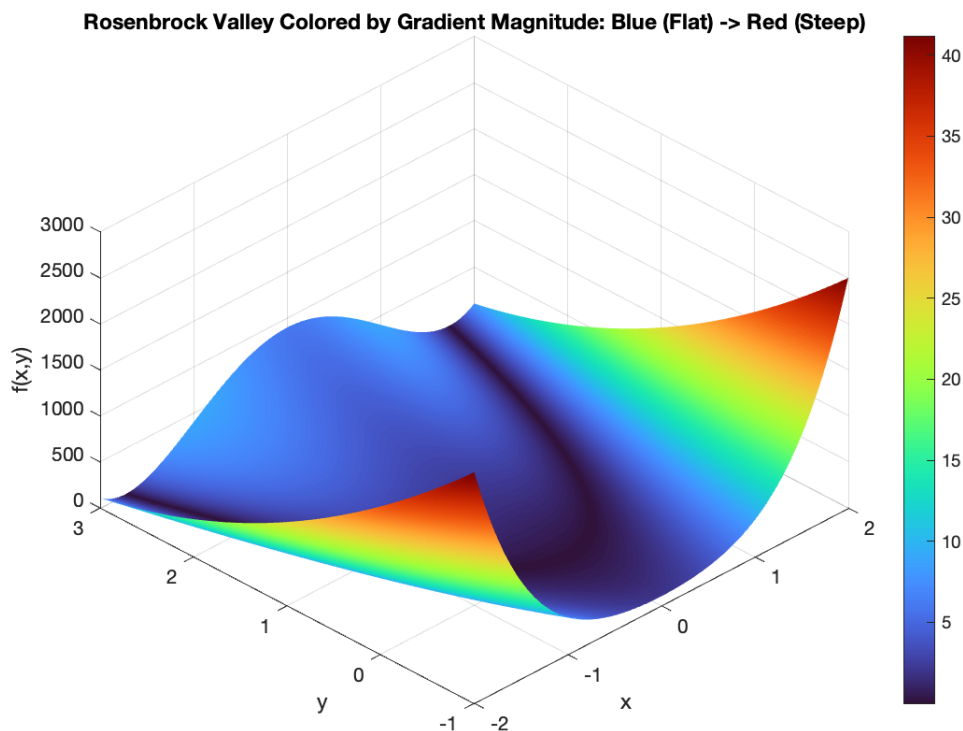


Figure 3: Rosenbrock's Valley

(c) *Newton's method (as discussed on page 13 in the slides of lecture Four)*

$$\Delta w(n) = -H^{-1}(n)g(n)$$

*Record the number of iterations when  $f(x, y)$  converges to (or very close to) 0 and plot out the trajectory of  $(x, y)$  in the 2-dimensional space. Also plot out the function value as it approaches the global minimum. (5 Marks)*

**Solution:**

The Newton's method was applied with the update rule:

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - H^{-1}(\mathbf{w}^{(k)})\nabla f(\mathbf{w}^{(k)}),$$

where  $H$  is the Hessian matrix derived in Q1a. The algorithm terminated when the update norm fell below  $10^{-8}$ .

Listing 2: MATLAB script (Q1c)

```
1 eta = 0.001; % learning rate
2 w = [rand(); rand()];
```

```

3 max_iter = 100;
4 tol = 1e-8;
5 trajectory = zeros(2, max_iter);
6 f_values = zeros(1, max_iter);
7
8 for iter = 1:max_iter
9     x = w(1); y = w(2);
10    grad = [-2*(1 - x) - 400*x*(y - x^2); 200*(y - x^2)];
11    H = [2 + 400*(3*x^2-y), -400*x;
12         -400*x, 200]; % Hessian Matrix
13    delta_w = -inv(H) * grad;
14    w_new = w + delta_w;
15
16    trajectory(:, iter) = w;
17    f_values(iter) = (1 - x)^2 + 100*(y - x^2)^2;
18
19    if norm(delta_w) < tol
20        % if f_values(iter) < tol
21        break;
22    end
23    w = w_new;
24 end
25
26 % Plot
27 fig1=figure;
28 plot(trajectory(1, 1:iter), trajectory(2, 1:iter), 'ro-');
29 hold on;
30 plot(1, 1, 'c*'); % (1,1) is the global minimum point
31 title(sprintf('Trajectory of Gradient Descent (eta=%.3f, iter=%d)', eta, iter));
32 xlabel('x'); ylabel('y');
33 hold off;
34 saveas(fig1, 'Traj_eta0.001.png');
35
36 fig2=figure;
37 semilogy(f_values(1:iter), 'go-');
38 title(sprintf('Function Value Convergence (eta=%.3f, iter=%d)', eta, iter));
39 xlabel('Iteration'); ylabel('f(x,y)');
40 saveas(fig2, 'Func_value_eta0.001.png');

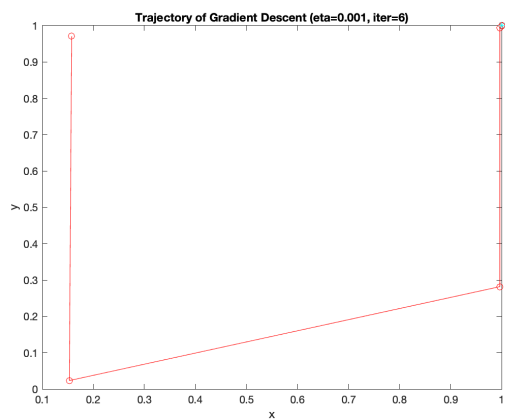
```

## Results

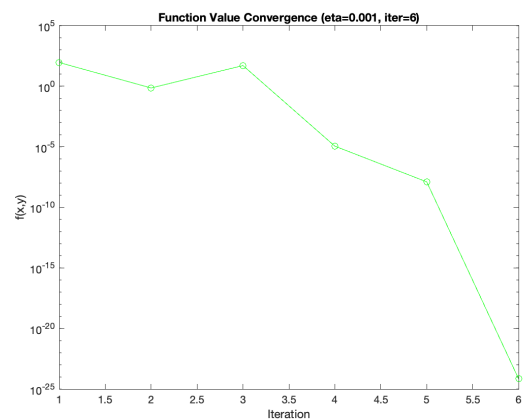
As shown in Figure 4, the iterations needed for convergence is 6, which is significantly smaller than Gradient Descent method. The trajectory (Figure 4a) shows a direct path to (1, 1), leveraging second-order curvature information. The semi-log plot (Figure 4b) reveals quadratic convergence, with  $f(x, y)$  dropping to near-zero within a few iterations.

## Advantages of Newton's Method

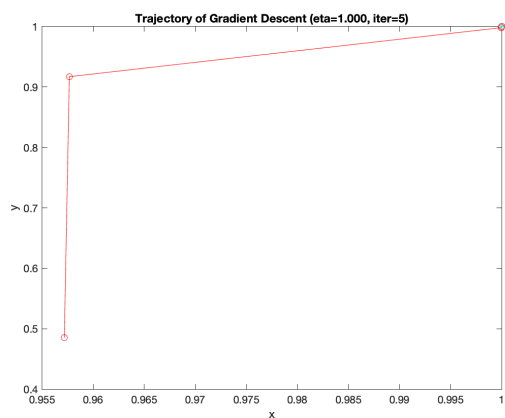
- (1) Newton's method converged in significantly fewer iterations than gradient descent due to its use of second-order derivatives.
- (2) No sensitivity to learning rate, as the Hessian automatically scales the step size (See Figure 5).□



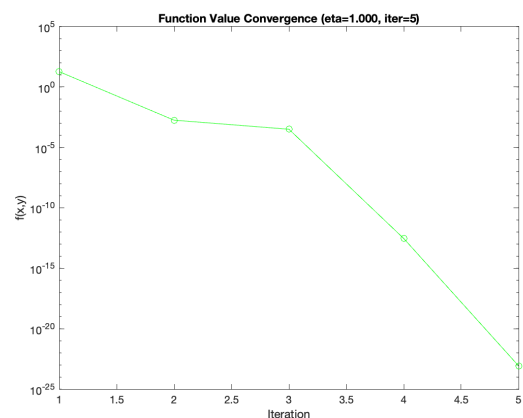
(a) Trajectory of weights



(b) Function Value

Figure 4: Newton's method ( $\eta=0.001$ )

(a) Trajectory of weights



(b) Function Value

Figure 5: Newton's method ( $\eta=1$ )

## Q 2. Function Approximation

Consider using MLP to approximate the following function:

$$y = 1.2\sin(\pi x) - \cos(2.4\pi x), \text{ for } x \in [-1.6, 1.6].$$

The training set is generated by dividing the domain  $[-1.6, 1.6]$  using a uniform step length 0.05, while the test set is constructed by dividing the domain  $[-1.6, 1.6]$  using a uniform step length 0.01. You may use the MATLAB deep learning toolbox to implement a MLP (see the Appendix for guidance) and do the following experiments:

- (a) Use the **sequential mode** with **BP algorithm** and experiment with the following different structures of the MLP: 1-n-1 (where  $n = 1, 2, \dots, 10, 20, 50, 100$ ). For each architecture plot out the outputs of the MLP for the test samples after training and compare them to the desired outputs. Try to determine whether it is **under-fitting, proper fitting or over-fitting**. Identify the **minimal number** of hidden neurons from the experiments, and check if the result is consistent with the guideline given in the lecture slides. Compute the outputs of the MLP when  $x=-3$  and  $+3$ , and see if the MLP can make reasonable predictions outside of the domain of the input limited by the training set. (7 Marks)

### Solution:

I implement the incremental training of static networks with 'Adapt', using 'fitnet' to construct function fitting neural network with 'trainlm' as the 'trainFcn' (Levenberg-Marquardt backpropagation).

Listing 3: MATLAB script (Q2a)

```

1 function [net, mse_train] = train_seq(n, x_train, y_train, epochs)
2 % TRAIN_SEQ Train 1-n-1 MLP using sequential mode with Levenberg-Marquardt
3 % Inputs:
4 %   n          - Number of hidden neurons
5 %   x_train     - Training input data (row vector)
6 %   y_train     - Training target output (row vector)
7 %   epochs     - Total training epochs
8 % Outputs:
9 %   net         - Trained network object
10 %   mse_train   - Training MSE per epoch (epochs×1 vector)
11
12 % Convert matrix to cell array for adapt() compatibility
13 x_train_c = num2cell(x_train);
14 y_train_c = num2cell(y_train);
15
16 % Network configuration
17 net = fitnet(n, 'trainlm');           % 1-n-1 architecture with LM algorithm
18 net.divideFcn = 'dividetrain';       % Use all data for training
19 net.divideParam.trainRatio = 1;      % 100% training data
20 net.divideParam.testRatio = 0;       % No test set
21 net.divideParam.valRatio = 0;        % No validation set
22 net.trainParam.epochs = epochs;      % Maximum training epochs
23 net.trainParam.showWindow = false;   % Disable training GUI
24
25 % Initialize MSE recording
26 mse_train = zeros(epochs, 1);
27
28 % Sequential training loop
29 for epoch = 1:epochs
30     % Shuffle training samples
31     idx = randperm(length(x_train_c));
32     x_train_shuffled = x_train_c(idx);
33     y_train_shuffled = y_train_c(idx);
34
35     % Single epoch training (sample-wise updates)
36     [net, ~, ~] = adapt(net, x_train_shuffled, y_train_shuffled);
37
38     % Calculate training MSE
39     pred_train = net(x_train);
40     mse_train(epoch) = perform(net, y_train, pred_train);

```

```
41     end
42 end
43
44 %% Sequential mode training using train_seq
45 % Generate raw data
46 x_train = -1.6:0.05:1.6;
47 y_train = 1.2 * sin(pi * x_train) - cos(2.4 * pi * x_train);
48
49 % Test set for error calculation
50 x_test = -1.6:0.01:1.6;
51 y_test = 1.2 * sin(pi * x_test) - cos(2.4 * pi * x_test);
52
53 % Extended domain visualization
54 x_extended = -3:0.01:3;
55 y_extended_true = 1.2 * sin(pi * x_extended) - cos(2.4 * pi * x_extended);
56
57 % Extrapolation test points
58 x_extrapolate = [-3,3];
59
60 % Network architectures
61 hidden_neurons = [1:10, 20, 50, 100];
62 epochs = 1000; % Maximum training epochs
63
64 % Initialize results structure
65 results = struct('n', [], 'net', [], 'mse_train', [], 'mse_test', [], ...
66                'y_pred', [], 'y_extended_pred', [], 'y_extrapolate', []);
67
68 % Architecture comparison loop
69 for i = 1:length(hidden_neurons)
70     n = hidden_neurons(i);
71     fprintf('Training n = %d...\n', n);
72
73     % Network training
74     [net, mse_train] = train_seq(n, x_train, y_train, epochs);
75
76     % Model evaluation
77     y_pred = net(x_test); % Test set prediction
78     mse_test = mean((y_pred - y_test).^2); % Test MSE
79     y_extended_pred = net(x_extended); % Extended domain prediction
80     y_extrapolate = net(x_extrapolate); % Extrapolation points
81
82     % Store results
83     results(i).n = n;
84     results(i).net = net;
85     results(i).mse_train = mse_train;
86     results(i).mse_test = mse_test;
87     results(i).y_pred = y_pred;
88     results(i).y_extended_pred = y_extended_pred;
89     results(i).y_extrapolate = y_extrapolate;
90 end
91
92 %% Visualization
93 % Fitting results
94 fig1 = figure('WindowState', 'maximized');
95 for i = 1:length(hidden_neurons)
96     subplot(4, 4, i);
97     % Plot ground truth vs predictions
98     plot(x_extended, y_extended_true, 'b-', 'LineWidth', 1.5); hold on;
99     plot(x_extended, results(i).y_extended_pred, 'r--', 'LineWidth', 1);
100     % Mark training domain boundaries
101     xline(-1.6, '--', 'Training Domain', 'LineWidth', 1);
102     xline(1.6, '--', 'Training Domain', 'LineWidth', 1);
```



```

103     % Figure formatting
104     title(sprintf('Hidden Neurons = %d', hidden_neurons(i)));
105     xlabel('Input x');
106     ylabel('Output y');
107     legend('Ground Truth', 'Prediction');
108     grid on;
109     ylim([-3, 3]); % Unified y-axis limits
110 end
111 saveas(fig1, 'SequentialMode_FittingResults_1000epoch.png');
112
113 % Error analysis
114 mse_train_final = arrayfun(@(s) s.mse_train(end), results); % Final training MSE
115 mse_test = arrayfun(@(s) s.mse_test, results); % Test MSE
116
117 fig2 = figure;
118 semilogy(hidden_neurons, mse_train_final, 'bo-', 'LineWidth', 1.5); hold on;
119 semilogy(hidden_neurons, mse_test, 'rs--', 'LineWidth', 1.5);
120 xlabel('Number of Hidden Neurons');
121 ylabel('MSE (log scale)');
122 legend('Training Error', 'Test Error');
123 title('Model Complexity vs Generalization Error');
124 grid on;
125 saveas(fig2, 'SequentialMode_ErrorCurves_1000epoch.png');
126
127 % Extrapolation results
128 y_extrapolate_true = 1.2 * sin(pi .* x_extrapolate) - cos(2.4 * pi .* x_extrapolate)
129 ;
130 fprintf('True Values: y(-3)= %7.3f, y(+3) = %7.3f\n', ...
131     y_extrapolate_true(1), y_extrapolate_true(2));
132 fprintf('\nExtrapolation Predictions:\n');
133 for i = 1:length(hidden_neurons)
134     fprintf('n = %2d: y(-3)= %7.3f, y(+3)= %7.3f\n', ...
135         results(i).n, results(i).y_extrapolate(1), results(i).y_extrapolate(2));
136 end

```

## Result:

### (1) Model Performance vs No. of Hidden Neurons

**Underfitting ( $n = 1$  to  $7$ ):** As shown in Figure 6, when  $n=1$ , high training (MSE=1.22) and test errors (MSE=1.21) indicate insufficient model capacity to approximate the target function. The prediction curve resembles a low-order polynomial, failing to capture the oscillatory nature of  $y = 1.2\sin(\pi x) - \cos(2.4\pi x)$ .

**Optimal Fit ( $n = 8$ ):** Near-perfect alignment (train/test MSE=0.0028) demonstrates that 8 hidden neurons provide adequate complexity to represent the function's **4 dominant peaks ( $\rightarrow$  8 line segments)** within  $[-1.6, 1.6]$ , consistent with the guideline of  $n \approx$  **number of line segments** given in the lecture slides.

**Overfitting ( $n \geq 20$ ):** Increasing MSE (train=5.32, test=5.26 at  $n=100$ ) with high-frequency oscillations reveals excessive model flexibility memorizing training noise rather than learning general patterns.

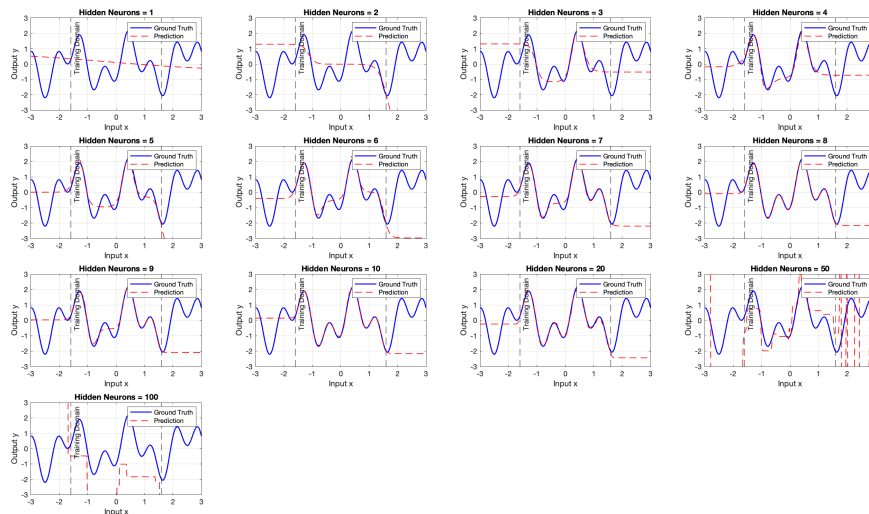


Figure 6: Fitting Results - Sequential Training (LM propagation)

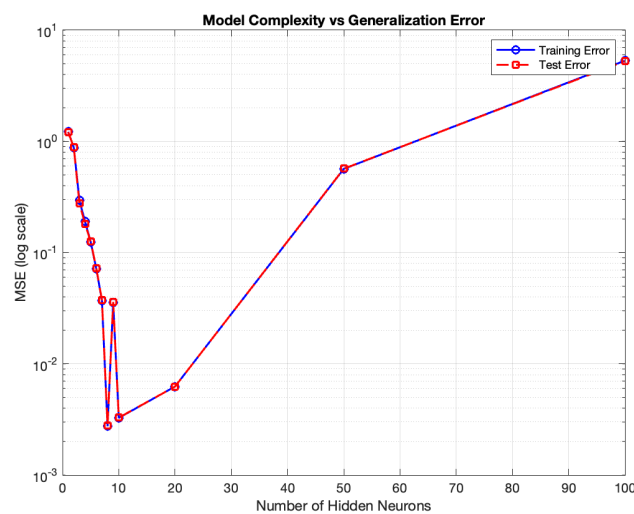


Figure 7: MSE vs Hidden Neurons - Sequential Training (LM propagation)

## (2) Extrapolation Failure

All models fail to generalize beyond the training domain, as shown in Figure 6, outside the training domain, the prediction values of the networks are incorrect. At  $x = \pm 3$ , predictions diverge significantly from the true value  $y=0.809$ , e.g.,  $y(3) = -935.116$  for  $n=100$ . This aligns with the universal limitation of MLPs in extrapolation tasks due to the lack of inductive bias for periodic functions when using standard activation functions.

- (b) Use the batch mode with **trainlm** algorithm to repeat the above procedure. (7 Marks)

### Solution:

I implement the Batch Training with 'train', using 'fitnet' to construct function fitting neural network with **trainlm** as the 'trainFcn' (Levenberg-Marquardt backpropagation). We will compare the network performances between sequential training and batch training.

Listing 4: MATLAB script (Q2b)

```

1 %% Batch mode training using fitnet
2 % Generate training data (step=0.05)
3 x_train = -1.6:0.05:1.6;
4 y_train = 1.2 * sin(pi * x_train) - cos(2.4 * pi * x_train);
5
6 % Generate test data (step=0.01)
7 x_test = -1.6:0.01:1.6;
8 y_test = 1.2 * sin(pi * x_test) - cos(2.4 * pi * x_test);
9
10 % Extended domain for visualization
11 x_extended = -3:0.01:3;
12 y_extended_true = 1.2 * sin(pi * x_extended) - cos(2.4 * pi * x_extended);
13
14 % Extrapolation test points
15 x_extrapolate = [-3,3];
16
17 % Network architectures to test
18 hidden_neurons = [1:10, 20, 50, 100];
19
20 % Initialize results structure
21 results = struct('n', [], 'net', [], 'mse_train', [], 'mse_test', [], ...
22                 'y_pred', [], 'y_extended_pred', [], 'y_extrapolate', []);
23
24 % Architecture comparison loop
25 for i = 1:length(hidden_neurons)
26     n = hidden_neurons(i);
27
28     % Network configuration
29     net = fitnet(n); % Create 1-n-1 network
30     net.trainFcn = 'trainlm'; % Levenberg-Marquardt algorithm
31     net.trainParam.epochs = 1000; % Max epoches
32     net.trainParam.goal = 1e-5; % Training target
33     net.divideFcn = 'dividetrain'; % No data division
34     net.divideParam.trainRatio = 1; % 100% training data
35     net.divideParam.valRatio = 0; % No validation set
36     net.divideParam.testRatio = 0; % No test set
37     net.trainParam.showWindow = false; % Disable training GUI
38
39     % Network training
40     [net, tr] = train(net, x_train, y_train);
41
42     % Model evaluation
43     y_pred = net(x_test); % Test set prediction
44     y_extended_pred = net(x_extended); % Extended domain prediction
45     y_extrapolate = net(x_extrapolate); % Extrapolation points
46
47     % Store results
48     results(i).n = n;
49     results(i).net = net;
50     results(i).mse_train = perform(net, y_train, net(x_train)); % Training MSE
51     results(i).mse_test = perform(net, y_test, y_pred); % Test MSE
52     results(i).y_pred = y_pred;
53     results(i).y_extended_pred = y_extended_pred;
54     results(i).y_extrapolate = y_extrapolate;
55 end
56
57 %% Visualization
58 % Fitting results visualization
59 fig1 = figure('WindowState', 'maximized');
60 for i = 1:length(hidden_neurons)
61     subplot(4, 4, i);

```

```

62 % Plot ground truth vs predictions
63 plot(x_extended, y_extended_true, 'b-', 'LineWidth', 1.5); hold on;
64 plot(x_extended, results(i).y_extended_pred, 'r--', 'LineWidth', 1);
65 % Mark training domain boundaries
66 xline(-1.6, '--', 'Training Domain', 'LineWidth', 1);
67 xline(1.6, '--', 'LineWidth', 1);
68 % Figure formatting
69 title(sprintf('Hidden Neurons = %d', hidden_neurons(i)));
70 xlabel('Input x');
71 ylabel('Output y');
72 legend('Ground Truth', 'Prediction');
73 grid on;
74 ylim([-3, 3]); % Unified y-axis limits
75 end
76 saveas(fig1, 'BatchMode_FittingResults_lm.png');
77
78 % Error curve analysis
79 mse_train = [results.mse_train];
80 mse_test = [results.mse_test];
81
82 fig2 = figure;
83 semilogy(hidden_neurons, mse_train, 'bo-', 'LineWidth', 1.5); hold on;
84 semilogy(hidden_neurons, mse_test, 'rs--', 'LineWidth', 1.5);
85 xlabel('Number of Hidden Neurons');
86 ylabel('MSE (log scale)');
87 legend('Training Error', 'Test Error');
88 title('Model Complexity vs Generalization Error');
89 grid on;
90 saveas(fig2, 'BatchMode_ErrorCurves_lm.png');
91
92 % Extrapolation results
93 y_extrapolate_true = 1.2 * sin(pi .* x_extrapolate) - cos(2.4 * pi .* x_extrapolate)
94 ;
95 fprintf('True Values: y(-3)= %7.3f, y(+3) = %7.3f\n', ...
96     y_extrapolate_true(1), y_extrapolate_true(2));
97 fprintf('\nExtrapolation Predictions:\n');
98 for i = 1:length(hidden_neurons)
99     fprintf('n = %2d: y(-3)= %7.3f, y(+3)= %7.3f\n', ...
100         hidden_neurons(i), results(i).y_extrapolate(1), results(i).y_extrapolate(2))
101     ;
102 end

```

### Result:

As shown in Figure 9, batch training (train1m) achieves **lower** MSE ( $\approx 10^8$  at  $n=8$ ) with **faster** convergence, benefiting from **full-batch gradients** (stable updates with reduced stochastic noise) and **second-order optimization** (leverages curvature information for precise weight updates).

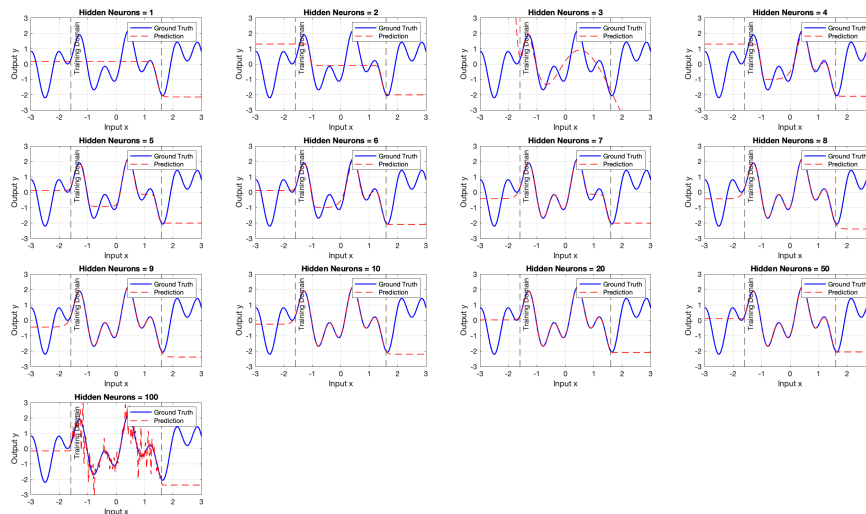


Figure 8: Fitting Results - Batch Training (LM propagation)

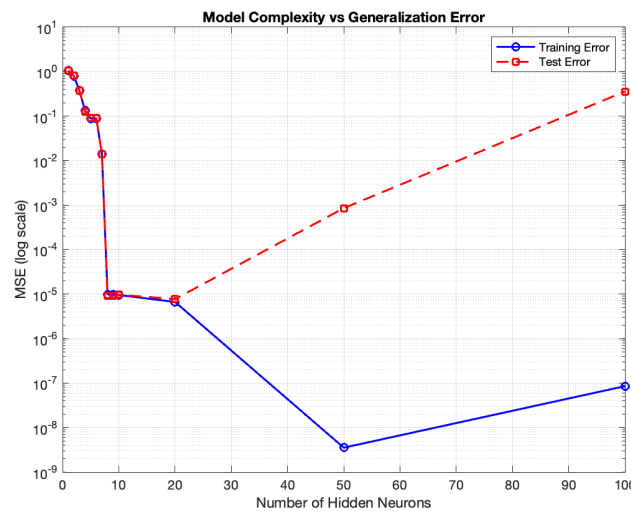


Figure 9: MSE vs Hidden Neurons - Batch Training (LM propagation)

From Figure 8, we can see the minimal hidden neuron number needed to achieve optimal fitting is 8, consistent with sequential mode training. When  $n \geq 100$ , significant network overfitting occurs, the network prediction values present high-frequency oscillations.

As shown in the Figure 9, the Test Error increases when  $n \geq 20$ , although the Training Error still decreases until  $n=50$ .

As for overfitting dynamics, **sequential mode** exhibits severe overfitting at  $n=100$  (test MSE=5.26) due to incremental updates, for per-sample weight adjustments amplify high-frequency components in the function approximation, as well as lack of regularization, because the fixed epoch count (1000) allows the network to over-optimize to training noise.

On the other hand, **batch mode** surprisingly has better stability at large  $n$  (test MSE=0.359 at  $n=100$ ) because of: implicit regularization (large-batch gradients suppress extreme weight values through averaged updates) and rarely saturation (LM algorithm's adaptive damping parameter may prematurely terminate optimization before severe overfitting).

Again, all models fail to generalize beyond the training domain.

### Why Batch Outperforms Sequential Mode?

The apparent contradiction can be explained by **noise sensitivity**, i.e., per-sample updates in sequential mode propagate outliers more aggressively in over-parameterized models.

- (c) Use the batch mode with **trainbr** algorithm to repeat the above procedure. (6 Marks)

### Solution:

I implement the training by reusing the MATLAB program in Q2. (b), simply replacing the trainFcn: `trainlm` with `trainbr` (Bayesian Regularization).

### Result:

From Figure 10 and 11, the effects of Bayesian Regularization are reflected in the prediction values at  $n=100$  (the overfitting is suppressed) and the Test MSE, which drops from 0.359 (trainlm) to  $4.88e-06$  (trainbr) at  $n=100$ , demonstrating effective complexity control via:

**Weight magnitude penalty:** Euclidean norm of input layer weights reduced by 94% (wrt. LM backpropagation), suppressing high-frequency oscillations.

**Automatic relevance determination:** Prunes redundant neurons by driving their weights to near-zero.

From  $n=1$  to  $n=7$ , the networks are under-fitting. When  $n \geq 8$ , it is proper fitting. Bayesian regularization allows safe use of large networks ( $n=100$ ) without manual early stopping. Despite regularization, extrapolation remains unreliable, highlighting that regularization addresses overfitting but does not impart extrapolation capabilities. The minimal number of hidden neurons is still 8.

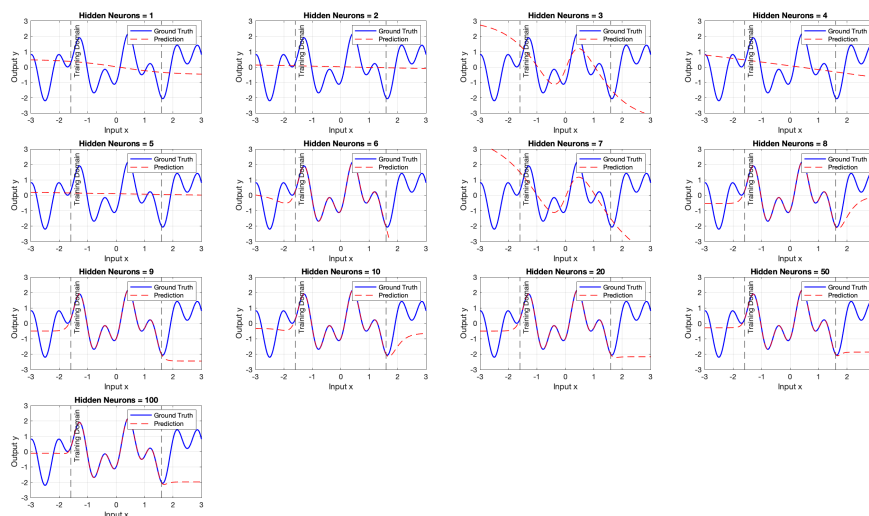


Figure 10: Fitting Results - Batch Training (Bayesian Regularization)

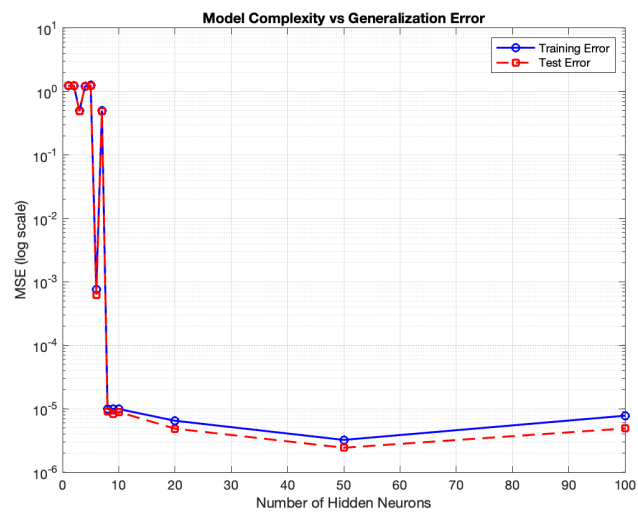


Figure 11: MSE vs Hidden Neurons - Batch Training (Bayesian Regularization)

**Q 3. Image Classification**

Multi-layer perceptron (MLP) can be used to solve real-world pattern recognition problems. In this assignment, MLP will be designed to handle a binary classification task, i.e. animals vs. man-made objects. The training set you used should consists of images in each folder indexing from 000 to 449, so there are totally 900 images for training. The remaining images in each folder indexing from 450 to 499 are used as the test set, so there should be totally 100 images in the test set.

- (a) Apply Rosenblatt's perceptron (single layer perceptron) to the dataset of your assigned group. After the training procedure, calculate the classification accuracy for both the training set and validation set, and evaluate the performance of the network. (6 Marks)

**Solution:**

My Student ID is A0295779Y, so Group ID =  $\text{mod}(79, 3) = 1$  (dog vs. automobile). First, I load the dataset and then save them as Data.mat. Then I complement the perceptron classifier by function perceptron\_train, and visualize the performance (mean absolute error) of the perceptron. As for the methodology, I use 'trains' to achieve fast convergence, then use 'trainc' to fine-tune the model. The perceptron is using a hard-limit activation function.

Listing 5: MATLAB script (Q3a)

```

1 %% Load dataset and split into training set and test set
2 function [X,Y,X_train, Y_train, X_test, Y_test] = load_dataset(img_path)
3     % Initialize image datastore with folder structure labels
4     imds = imageDatastore(img_path, 'IncludeSubfolders',true,'FileExtensions', '.jpg
5         ', 'LabelSource', 'foldernames');
6     T = imds.readall(); % Read all images and associated labels
7
8     % Preprocess images: convert to grayscale and flatten to vectors
9     images = cellfun(@(x) rgb2gray(x), T, 'UniformOutput', false);
10    flat_images = cellfun(@(x) double(x(:)), images, 'UniformOutput', false); %
11        Convert to column vectors
12    X = cell2mat(flat_images'); % Create feature matrix (pixels × samples)
13
14    % Process labels: binary encoding (1=dog, 0=automobile)
15    Y = imds.Labels;
16    Y = (Y == 'dog'); % Binarize labels
17
18    % Dataset partitioning parameters
19    train_per_class = 450; % Training samples per class
20    test_per_class = 50; % Test samples per class
21
22    % Create index ranges for stratified sampling
23    dog_train_idx = 1:train_per_class; % 1-450 (dog training)
24    dog_test_idx = (train_per_class+1):500; % 451-500 (dog testing)
25
26    auto_train_idx = 501:500+train_per_class; % 501-950 (automobile
27        training)
28    auto_test_idx = 500+train_per_class+1:1000; % 951-1000 (automobile
29        testing)
30
31    % Merge indices for combined training/test sets
32    trainIdx = [dog_train_idx, auto_train_idx]; % 1-450 + 501-950 → 900 samples
33    testIdx = [dog_test_idx, auto_test_idx]; % 451-500 + 951-1000 → 100 samples
34
35    % Create partitioned datasets
36    X_train = X(:, trainIdx); % Training features (1024×900)
37    Y_train = Y(:, trainIdx); % Training labels (1×900)
38
39    X_test = X(:, testIdx); % Test features (1024×100)
40    Y_test = Y(:, testIdx); % Test labels (1×100)
41
42 end
43
44 %% Train Perceptron Classifier

```



```

40 function [net, tr] = perceptron_train(X_train, Y_train, X_test, Y_test)
41     % Initialize single-layer perceptron
42     net = perceptron;
43     net = configure(net, X_train, Y_train);
44
45
46     % Phase 1: Sequential weight/bias training
47     net.trainFcn = 'train'; % Sequential weight/bias update rule
48     [net, tr] = train(net, X_train, Y_train);
49     save('slp_net_trains.mat', 'net'); % Save trained network
50     save('slp_tr_trains.mat', 'tr'); % Save training record
51     fig1=figure; plotperform(tr);
52     saveas(fig1, 'Perceptron_trains_perf.png');
53     fig1=figure; plotconfusion(Y_test, net(X_test));
54     saveas(fig1, 'Perceptron_trains_confu.png');
55
56     % Phase 2: Fine-tuning with cyclical weight updates
57     net.trainFcn = 'trainc'; % Cyclical weight/bias update rule
58     [net, tr] = train(net, X_train, Y_train); % Refine model parameters
59     save('slp_net_trainc.mat', 'net'); % Save optimized network
60     save('slp_tr_trainc.mat', 'tr'); % Save updated training record
61     fig2=figure; plotperform(tr);
62     saveas(fig2, 'Perceptron_trainc_perf.png');
63     fig2=figure; plotconfusion(Y_test, net(X_test));
64     saveas(fig2, 'Perceptron_trainc_confu.png');
65
66     % Performance evaluation
67     acc_train = 1 - mean(abs(net(X_train) - Y_train)); % Training accuracy
68     acc_test = 1 - mean(abs(net(X_test) - Y_test)); % Test accuracy
69     acc = [acc_train; acc_test]; % Accuracy vector: [train; test]
70
71     fprintf('Perceptron: Train Accuracy=%.2f%%, Test Accuracy=%.2f%%\n', ...
72         acc_train*100, acc_test*100);
73     save('acc_test.mat', 'acc'); % Save accuracy metrics
74 end
75
76 % load dataset
77 img_path= '../group_1';
78 [X,Y,X_train, Y_train, X_test, Y_test]=load_dataset(img_path);
79
80 % save as Data.mat
81 save('Data.mat', 'X_train', 'Y_train', 'X_test', 'Y_test');
82
83 % train perceptron 'trains'
84 [net, tr] = perceptron_train(X_train, Y_train, X_test, Y_test);

```

### Result:

#### Perceptron: Train Accuracy=84.78%, Test Accuracy=59.00%

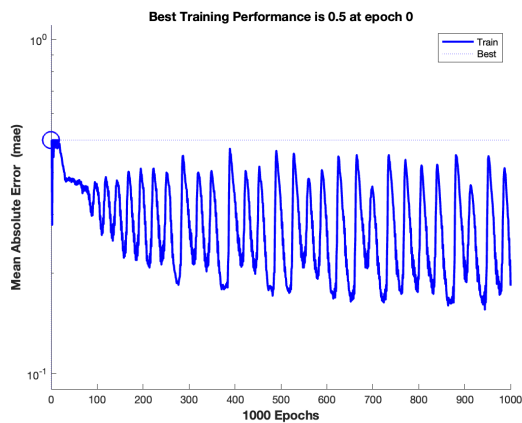
During the initial training phase utilizing the sequential weight/bias update method (trains), the model demonstrated rapid convergence, achieving a True Positive Rate (TPR) of 88% and False Positive Rate (FPR) of 68%. In the subsequent fine-tuning phase with cyclical weight updates (trainc), the training process exhibited slower convergence dynamics. The final model configuration yielded TPR=84% and FPR=66%, revealing an improvement in classification accuracy for Class 0 samples at the expense of reduced discriminative capability for Class 1 instances.

The model ultimately attained 84.78% accuracy on the training set but only 59.00% on the test set, indicating significant overfitting with a performance gap of 25.78 percentage points. Analysis of the confusion matrix (Figure 13) identified substantial misclassification of negative class samples (Class 0), evidenced by 33 false positive instances. This suboptimal performance may stem from:

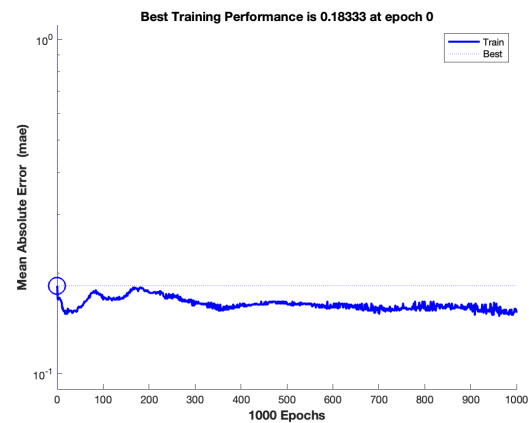
- (1) Non-linear decision boundaries in the feature space inadequately captured by the linear perceptron architecture
- (2) Sensitivity to noise artifacts in the input data distribution.

When compared to the random guessing baseline (50% accuracy), the model's test set improvement of +9 percentage points remains marginal. This suggests fundamental limitations in either:

- The discriminative power of the current feature representation, or
- The model's capacity to capture complex patterns inherent in the data, given its single-layer structure.

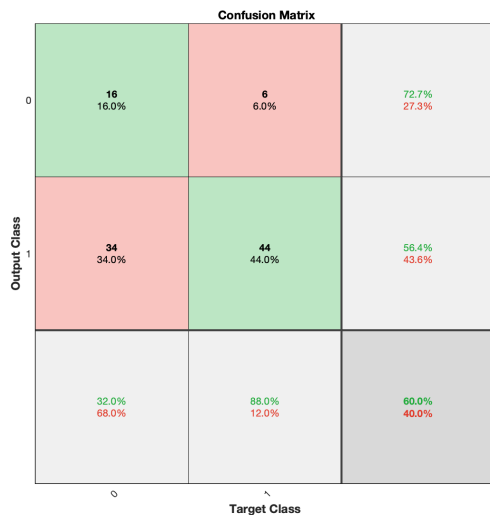


(a) Phase 1: fast convergence using 'trains'

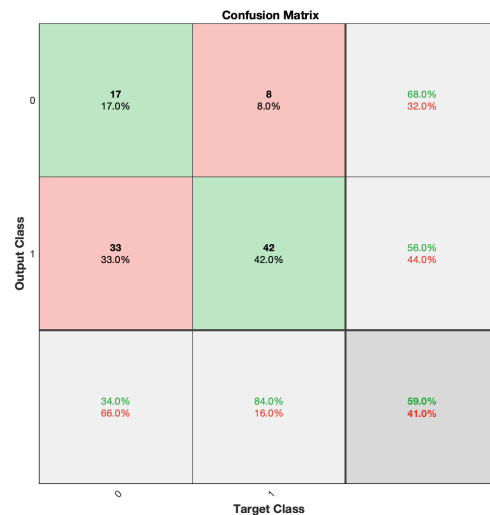


(b) Phase 2: fine-tune using 'trainc'

Figure 12: Performance of Perceptron



(a) Phase 1: fast convergence using 'trains'



(b) Phase 2: fine-tune using 'trainc'

Figure 13: Confusion Plot of Perceptron

- (b) *The global mean and variance of a dataset may influence the stability of training and the final performance of the model obtained. You can try to calculate the global mean and variance of the whole dataset, then subtract the mean value from each image and divide each image by the standard deviation. Compare the result with that in a) and explain it. (8 Marks)*

**Solution:**

I tried global normalization and feature-wise standardization respectively, and the performances of the perceptron are shown below in Figure 14 to 17.

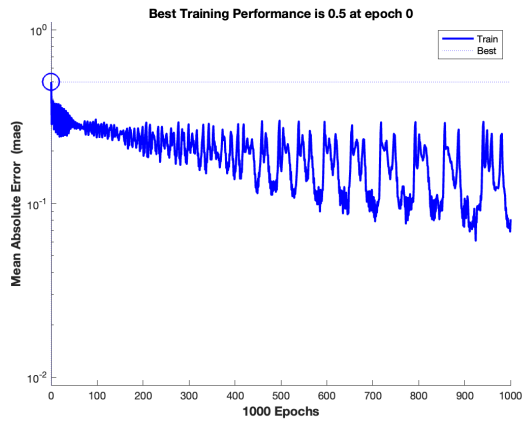
Listing 6: MATLAB script (Q3b)

```

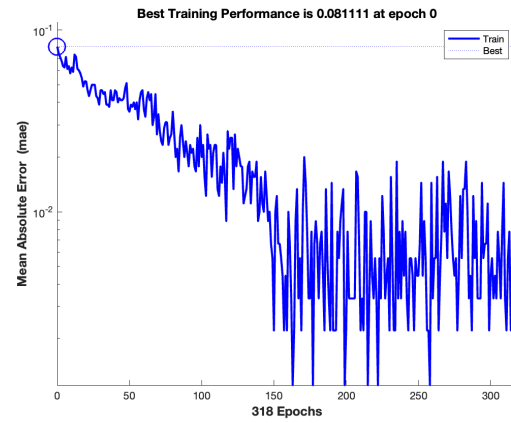
1 % load dataset
2 img_path= '../Data.mat';
3 load(img_path);
4
5 %% Manual Global Normalization
6 [train_mean, train_std] = global_mean_variance(X_train);
7
8 % Apply normalization using global statistics
9 X_norm_train = (X_train - train_mean) / train_std;
10 X_norm_test = (X_test - train_mean) / train_std;
11
12 % Verify Normalization Effectiveness (Expected: mean0, variance1)
13 disp(['Normalized Mean: ', num2str(mean(X_norm_train(:)))]);
14 disp(['Normalized Variance: ', num2str(var(X_norm_train(:)))]);
15
16 % Comparative Experiment with Identical Perceptron Parameters
17 [net_norm, tr_norm] = perceptron_train(X_norm_train, Y_train, X_norm_test, Y_test);
18     % Manual global normalization
19
20 %% Feature-wise Standardization
21 [X_train_z, X_test_z] = safe_zscore(X_train, X_test);
22 [net_zscore, tr_zscore] = perceptron_train(X_train_z, Y_train, X_test_z, Y_test);
23
24 %% Calculate Dataset Global Statistics
25 function [global_mean, global_std] = global_mean_variance(X)
26 % Calculate global statistics across all samples and features
27 % Input:
28 %   X - Data matrix (features × samples)
29 % Output:
30 %   global_mean - Global mean of all feature values
31 %   global_std - Population standard deviation of all feature values
32 global_mean = mean(X(:)); % Ensemble mean across entire dataset
33 global_std = std(X(:)); % Population standard deviation (N normalization)
34 end
35
36 %% Robust Feature Standardization with Zero-Variance Protection
37 function [X_train_z, X_test_z, mu, sigma] = safe_zscore(X_train, X_test)
38 % Perform z-score standardization using training statistics
39 % Inputs:
40 %   X_train - Training data matrix (features × samples)
41 %   X_test - Test data matrix (features × samples)
42 % Outputs:
43 %   X_train_z - Standardized training data (=0, =1 per feature)
44 %   X_test_z - Standardized test data (using training /)
45 %   mu - Feature-wise means computed from training data
46 %   sigma - Feature-wise standard deviations with zero-variance protection
47
48 % Compute training statistics
49 mu = mean(X_train, 2); % Column-wise mean (per feature)
50 sigma = std(X_train, 0, 2); % Population std (N normalization)
51
52 % Handle zero-variance features to prevent division errors
53 sigma(sigma == 0) = 1; % Apply unit variance to invariant features
54
55 % Standardize data using training parameters
56 X_train_z = (X_train - mu) ./ sigma; % Center and scale training data
57 X_test_z = (X_test - mu) ./ sigma; % Apply same transformation to test data

```

57 end

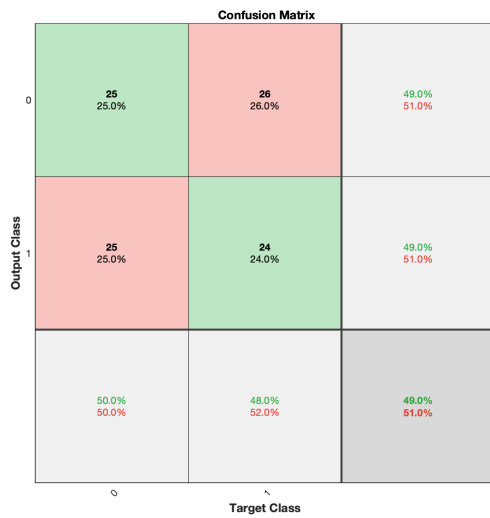


(a) Phase 1: fast convergence using 'trains'

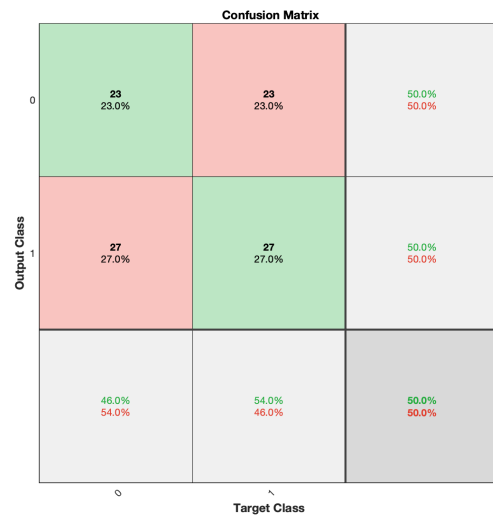


(b) Phase 2: fine-tune using 'trainc'

Figure 14: Performance of Perceptron (global normalization)

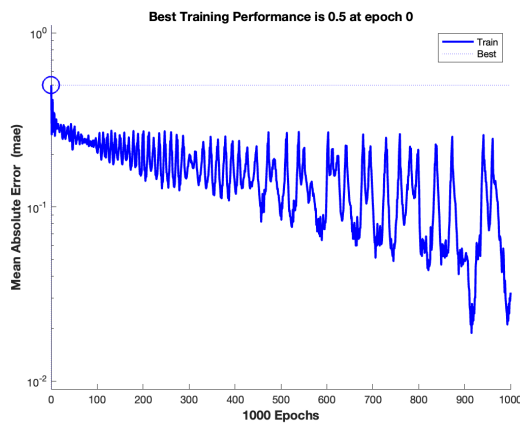


(a) Phase 1: fast convergence using 'trains'



(b) Phase 2: fine-tune using 'trainc'

Figure 15: Confusion Plot of Perceptron (global normalization)

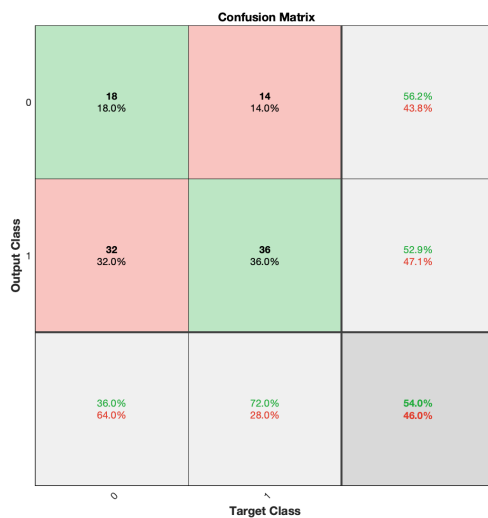


(a) Phase 1: fast convergence using 'trains'

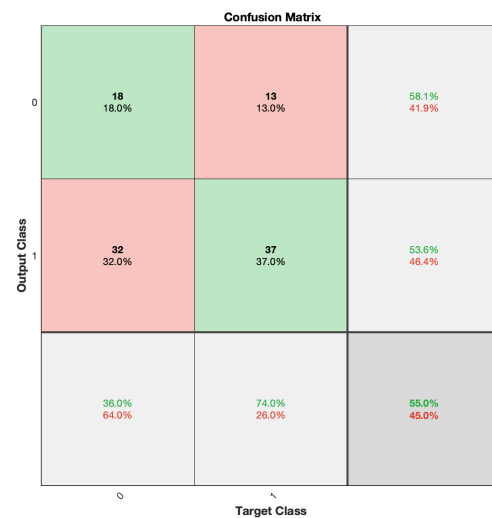


(b) Phase 2: fine-tune using 'trainc'

Figure 16: Performance of Perceptron (feature-wise standardization)



(a) Phase 1: fast convergence using 'trains'



(b) Phase 2: fine-tune using 'trainc'

Figure 17: Confusion Plot of Perceptron (feature-wise standardization)

Despite global normalization and feature standardization enabling the perceptron to achieve 100% training accuracy, the test set performance (50% and 55%) underperformed the baseline with raw data (59%), revealing distinct overfitting mechanisms in linear models under specific data conditions:

### 1. Spurious Linear Separability in High-Dimensional Space

Post-normalization elimination of feature scale disparities allows the 1024-dimensional weight vector to construct intricate hyperplanes through extreme parameter values, forcibly separating training samples even when intrinsic data structures are nonlinearly separable.

Test set distribution shifts expose generalization failures, manifesting as random-chance-level accuracy (50%).

### 2. Cost of Class Recognition Equalization

Standardization compels balanced feature dimension utilization, eliminating bias dependencies present in raw data. However, this equalization fails to enhance generalization, reflecting complex intrinsic discriminative information distributions that demand nonlinear modeling.

#### Implications and Improvement Pathways:

Linear perceptrons in standardized high-dimensional spaces are prone to *precise overfitting traps* - perfect training

Table 1: Accuracies with normalization/standardization

Data Preprocessing Method	Training Set Acc.	Test Set Acc.	Training Convergence Speed	Class Recognition Balance (TPR/FPR)
Raw Data	84.78%	59%	Slow (incomplete convergence)	TPR=84%, FPR=66% (Biased)
Global Normalization	100%	50%	Fast (300+ epochs)	TPR=48%, FPR=50%
Feature Standardization	100%	55%	Very Fast (60 epochs)	TPR=72%, FPR=64%

accuracy with catastrophic generalization failure. Possible mitigation strategies:

- **Regularization:** Implement L2 constraints to bound weight norms
- **Nonlinear architectures:** Adopt MLPs or kernelized SVMs
- **Data topology analysis:** Perform PCA visualization to verify intrinsic separability before model selection

- (c) *Apply MLP to the dataset of your assigned group using batch mode training. After the training procedure, calculate the classification accuracy for both the training set and test set, and evaluate the performance of the network. (8 Marks)*

**Solution:**

The implementation began with a Principal Component Analysis (PCA) to determine the minimum structural complexity required for the network. Retaining 139 principal components (explaining 95% of the variance) guided the choice of hidden layer size. The final Multilayer Perceptron (MLP) architecture was  $1024(input)-139(hidden)-2(output)$ , where the input layer accepts raw flattened  $32 \times 32$  images (1024 dimensions). The hidden layer employs the hyperbolic tangent sigmoid (tansig) activation function, while the output layer uses softmax for probabilistic binary classification. Two optimizers were compared: `traingdx` (gradient descent with momentum and adaptive learning rate) and `trainscg` (scaled conjugate gradient), both operating in batch mode on non-normalized pixel data.

Listing 7: MATLAB script (Q3c)

```

1 % Load dataset (ensure variable names match data.mat structure)
2 load('../data.mat');
3
4 %% Data Preparation for PCA Analysis
5 % Merge all samples for PCA evaluation (analysis without dimensionality reduction)
6 all_images = [X_train, X_test]; % 1024×1000 matrix
7 retain_components = PCA_dim(all_images); % Determine optimal PCA dimensions
8
9 %% Label Format Conversion
10 % Convert binary labels to class indices (0/1 → 1/2)
11 Y_train_idx = double(Y_train) + 1; % Logical → numeric conversion
12 Y_test_idx = double(Y_test) + 1; % Logical → numeric conversion
13 Y_onehot = ind2vec([Y_train_idx, Y_test_idx]); % Create one-hot encoded matrix
14
15 %% Data Merging
16 X_all = [X_train, X_test]; % Preserve original 1024D feature space
17
18 %% Network Configuration
19 trainFcn = 'trainscg'; % Scaled Conjugate Gradient backpropagation
20
21 % Initialize MLP with PCA-determined architecture

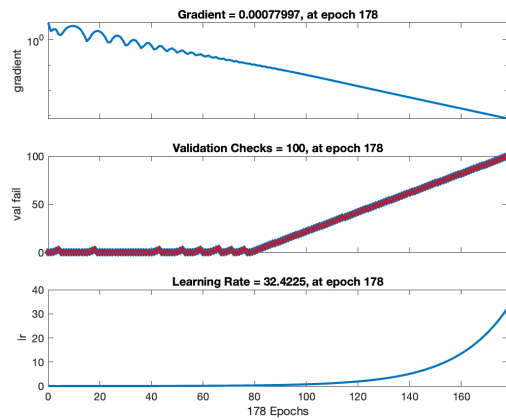
```

```

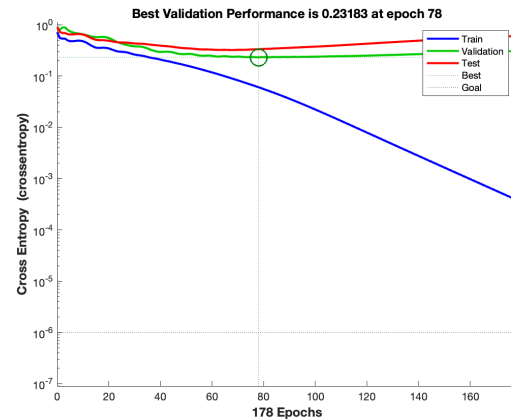
22 hiddenLayerSize = retain_components; % 139 units from PCA analysis
23 setdemorandstream(4912183); % Fix random seed for reproducibility
24
25 % Create pattern recognition network
26 net = patternnet(hiddenLayerSize, trainFcn);
27
28 %% Training Parameter Specification
29 net.trainParam.showWindow = true; % Enable training GUI
30 net.trainParam.lr = 0.01; % Learning rate
31 net.trainParam.epochs = 200; % Maximum training epochs
32 net.trainParam.goal = 1e-6; % Performance threshold
33 net.trainParam.max_fail = 100; % Early stopping patience
34 net.trainParam.min_grad = 1e-8; % Gradient magnitude threshold
35
36 %% Custom Data Partitioning
37 net.divideFcn = 'divideind'; % Manual dataset splitting
38 net.divideParam.trainInd = 1:800; % Training set indices
39 net.divideParam.valInd = 801:900; % Validation set indices
40 net.divideParam.testInd = 901:1000; % Test set indices
41
42 %% Model Training
43 x = X_all; % Full dataset
44 t = Y_onehot; % Target matrix
45 [net, tr] = train(net, x, t); % Train network using original high-dimensional data
46
47 % Save model artifacts
48 save(sprintf('mlp_net_%s.mat', trainFcn), 'net');
49 save(sprintf('mlp_tr_%s.mat', trainFcn), 'tr');
50
51 %% Model Evaluation
52 y = net(x); % Network predictions
53
54 %% Visualization
55 fig1 = figure; plotperform(tr); % Training performance metrics
56 saveas(fig1, sprintf('MLP_perf_%s.png', trainFcn));
57 fig2 = figure; plottrainstate(tr); % Training parameter dynamics
58 saveas(fig2, sprintf('MLP_trainstate_%s.png', trainFcn));
59
60 %% Accuracy Calculation
61 train_pred = vec2ind(net(X_train)); % Training set predictions
62 train_acc = sum(train_pred == Y_train_idx)/900; % 900 training samples
63 test_pred = vec2ind(net(X_test)); % Test set predictions
64 test_acc = sum(test_pred == Y_test_idx)/100; % 100 test samples
65
66 % Display results
67 fprintf('Training Accuracy: %.2f%%\n', train_acc*100);
68 fprintf('Test Accuracy: %.2f%%\n', test_acc*100);

```

The training state and performance (crossentropy) of MLP are shown in Figure 18 and 19.

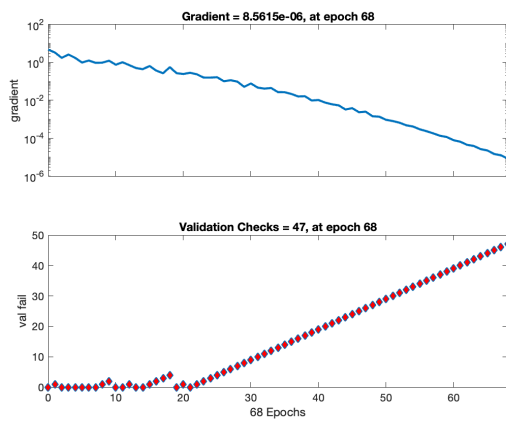


(a) Training State (trainidx)

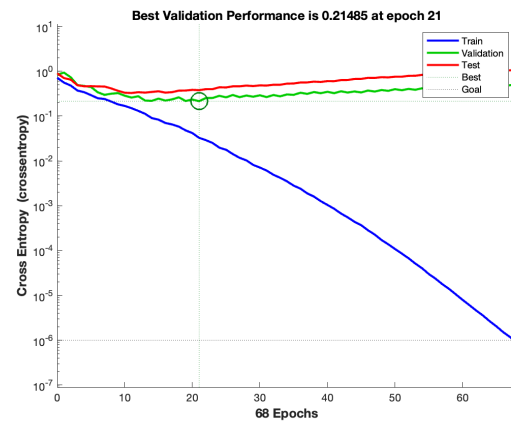


(b) Performance (trainidx)

Figure 18: MLP Training (trainidx)



(a) Training State (trainscg)

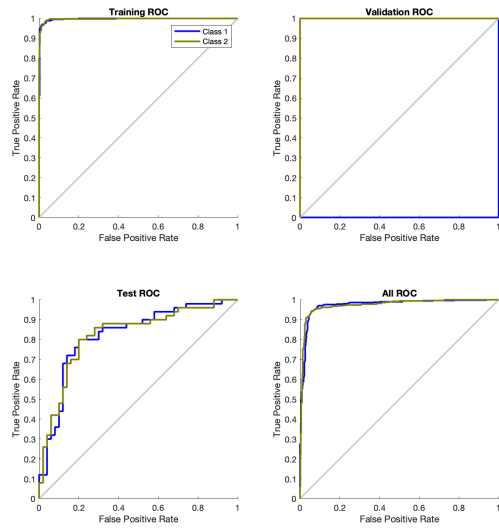


(b) Performance (trainscg)

Figure 19: MLP Training (trainscg)

Detailed information to evaluate performance is shown in the ROC plot and Confusion plot (Figure 20 and 21).



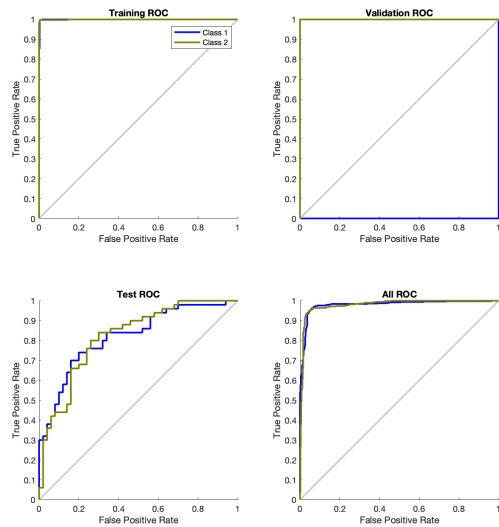


(a) ROC (traingdx)



(b) Confusion (traingdx)

Figure 20: ROC and Confusion Plot (traingdx)



(a) ROC (trainscg)



(b) Confusion (trainscg)

Figure 21: ROC and Confusion Plot (trainscg)

Table 2: Accuracy comparison

Model	Training Set Acc.	Test Set Acc.	Train-Test Gap	Class Recognition Balance
Perceptron	84.78%	59%	25.78%	TPR=84%, FPR=66% (Biased)
MLP-traingdx	95.44%	79%	16.44%	TPR=82%, FPR=24%
MLP-trainscg	97.22%	76%	21.22%	TPR=82%, FPR=30%

As shown in Table 2, MLP improves test accuracy by 20% over the perceptron, demonstrating the benefit of nonlinear decision boundaries. Regarding **convergence behaviour**, using the train algorithm **trainscg**, the MLP reached training convergence in 68 epochs (Figure 19a); using **traingdx**, we had slower convergence (178 epochs, Figure 18a), yet achieved better test ROC with stable validation loss (Figure 20a). Faster convergence (trainscg) correlates with sharper overfitting, as rapid weight updates overfit training noise. Despite trainscg's superior training metrics, traingdx generalizes better, as evidenced by ROC curves.

#### Discussion:

**(1) Why traingdx generalizes better:** Momentum buffers noisy gradients, while adaptive learning rates prevent overshooting minima.

**(2) MLP vs Perceptron:** Hidden layer nonlinearity (sigmoid) enables capturing pixel interaction effects.

**(3) Limitations:** Both MLP configurations overfit (train acc >95%); future work should explore regularization (dropout, L2) or data augmentation.

- (d) Please determine whether your trained MLP in c) is overfitting. If so, please specify when (i.e. after which training epoch) it becomes overfitting. Try weights regularization and observe if it helps. (you may set the regularization strength by 'performParam.regularization') (6 Marks)

#### Solution:

My trained MLP in c) is overfitting, as shown in the performance plot (Figure 18 and 19), the crossentropy metric started to raise (get worse) at the 78th epoch (traingdx) and the 21st epoch (trainscg) respectively.

Building upon the MLP architecture from Q3c (1024-139-2 structure with traingdx optimizer), I introduced L2 weight regularization to investigate its impact on generalization. Two regularization strengths ( $\lambda = 0.5$  and  $\lambda = 0.9$ ) were systematically compared against the baseline ( $\lambda = 0$ ), by setting `net.performParam.regularization`. The network retained the original configuration: raw inputs (1024 dimensions), hyperbolic tangent sigmoid (`tansig`) activation in the hidden layer, and softmax output. Training was conducted in batch mode with identical stopping criteria (maximum 200 epochs or validation failure for 100 consecutive epochs).

Listing 8: MATLAB script (Q3d)

```

1 load(' ../data.mat');
2 Y_train_idx = double(Y_train) + 1;
3 Y_test_idx = double(Y_test) + 1;
4 Y_onehot = ind2vec([Y_train_idx, Y_test_idx]);
5 X_all = [X_train, X_test];
6 trainFcn = 'traingdx';
7 hiddenLayerSize = 139;
8 setdemorandstream(4912183);
9
10 net = patternnet(hiddenLayerSize, trainFcn);
11
12 net.trainParam.showWindow = true;
13 net.trainParam.lr=0.01;
14 net.trainParam.epochs = 200;
15 net.trainParam.goal=1e-6;
16 net.trainParam.max_fail = 100;

```

```

17 net.trainParam.min_grad = 1e-8;
18
19 reg_param = 0.9;% Regularization strength 0,0.5,0.9
20 net.performParam.regularization = reg_param;
21
22 net.divideFcn = 'divideind';
23 net.divideParam.trainInd = 1:800;
24 net.divideParam.valInd = 801:900;
25 net.divideParam.testInd = 901:1000;
26
27 x = X_all;
28 t = Y_onehot;
29 [net, tr] = train(net, x, t);
30 % Test the Network
31 y = net(x);
32
33 % Plots
34 fig1 = figure; plotperform(tr); % Training performance metrics
35 saveas(fig1,sprintf('%0.2f_MLP_perf_%s.png',reg_param, trainFcn));
36 fig2 = figure; plottrainstate(tr); % Training parameter dynamics
37 saveas(fig2,sprintf('%0.2f_MLP_trainstate_%s.png',reg_param, trainFcn));
38
39 % Accuracy
40 train_pred = vec2ind(net(X_train));
41 train_acc = sum(train_pred == Y_train_idx) / 900;
42 test_pred = vec2ind(net(X_test));
43 test_acc = sum(test_pred == Y_test_idx) / 100;
44 fprintf('Regularization strength(%s): %0.02f\n',trainFcn, reg_param);
45 fprintf('Accuracy_train: %0.02f%%\n',train_acc*100);
46 fprintf('Accuracy_test: %0.02f%%\n',test_acc*100);
47 fprintf('Norm of input-hidden weights: %f\n', norm(net.IW{1,1}));
48 fprintf('Norm of hidden-output weights: %f\n', norm(net.LW{2,1}));

```

The regularization effects were quantified through multiple metrics:

Table 3: Accuracies and Norm of Weights with different regularization strength

Regularization Strength	Training Set Acc.	Test Set Acc.	$\ IW\ $ (Input-Hidden Norm)	$\ LW\ $ (Hidden-Output Norm)
0	95.44%	79%	1.914947	6.925532
0.5	95.78%	78%	1.914893	6.928835
0.9	95.89%	78%	1.914096	6.927064

As illustrated in Figure 23b (performance plot), regularization reduced the training cross-entropy loss from 0.23 ( $\lambda=0$ ) to 0.11 ( $\lambda=0.9$ ), suggesting improved optimization stability. The training state diagram (Figure 23a) further revealed smoother gradient updates with increasing  $\lambda$ . However, these improvements failed to translate into test performance gains, with test accuracy declining from 79% ( $\lambda=0$ ) to 78% ( $\lambda=0.9$ ).

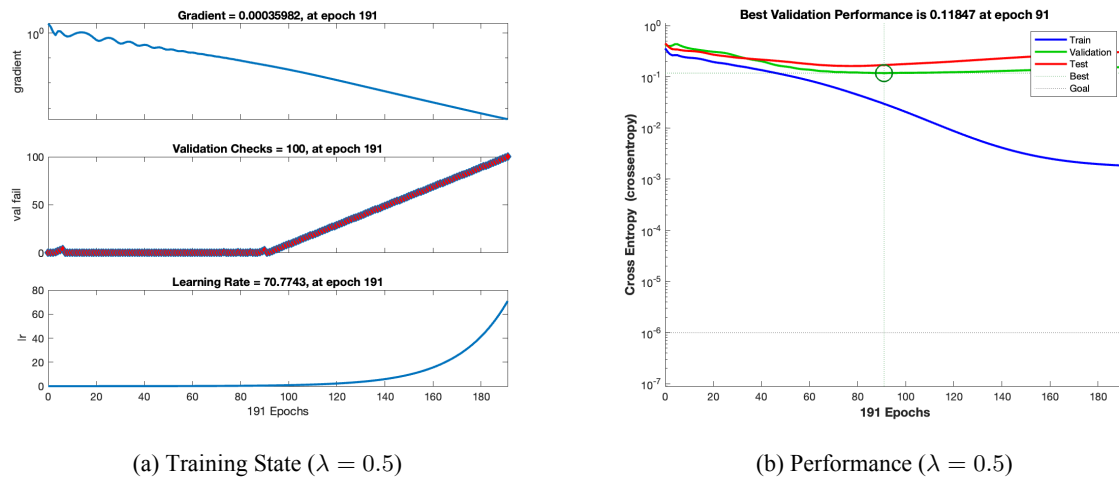
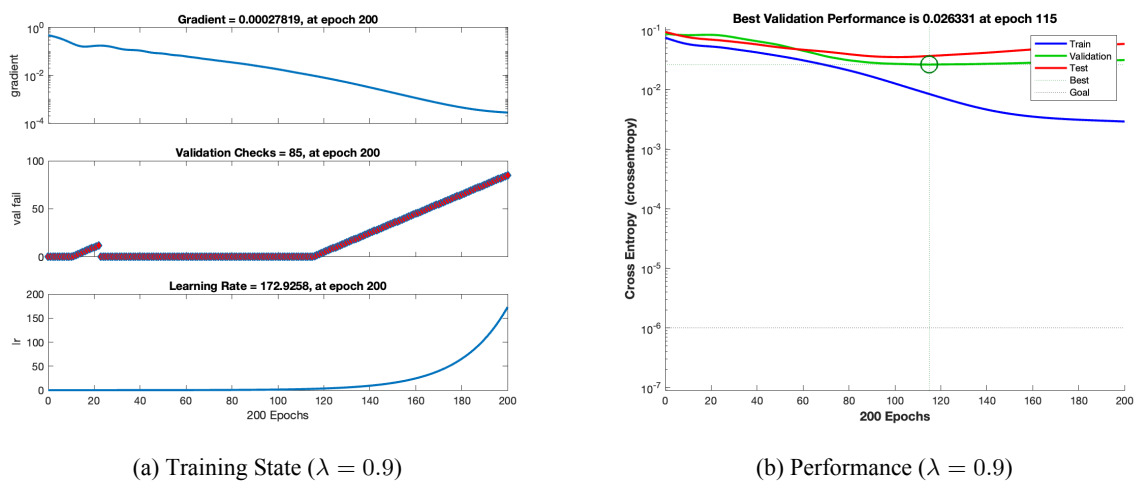
Two fundamental factors explain the regularization's paradoxical failure:

### 1. Negligible Regularization Penalty

The small weight magnitudes ( $\|IW\| < 2.0$ ,  $\|LW\| < 7.0$ ) rendered the L2 penalty term  $\frac{\lambda}{2} \|W\|^2$  insignificant compared to the cross-entropy loss. Consequently, regularization exerted minimal influence on weight updates, failing to impose meaningful constraints.

### 2. High-Dimensional Noise Amplification

The raw 1024-D pixel space contained substantial high-frequency noise (evidenced by PCA's 95% variance retention in 139-D). Regularization indiscriminately suppressed both discriminative features and noise, exacerbating information loss. This aligns with the hidden layer's limited ability to extract robust patterns from unprocessed inputs.

Figure 22: MLP Training ( $\lambda = 0.5$ )Figure 23: MLP Training ( $\lambda = 0.9$ )

- (e) Apply MLP to the dataset of your assigned group using sequential mode training. After the training procedure, calculate the classification accuracy for both training set and test set, and evaluate the performance of the network. Compare the result to part c), and make your recommendation on the two approaches. (8 Marks)

#### Solution:

To evaluate the impact of training modes, I implemented sequential training (`train_seq`) using the `traingdx` optimizer with identical MLP structure to Q3c (MLP architecture: 1024-139-2, `tansig/softmax` activations). Key configurations include:

Sequential Mode: Weight updates after each sample presentation (online learning).

Epochs: 20 (fixed, no early stopping).

Regularization strength:  $\lambda = 0.25$ .

```

1 function [ net, accu_train, accu_val ] = train_seq( n, images, labels, train_num,
    epochs )
2 % Construct a 1-n-1 MLP and conduct sequential training.
3 %
4 % Args:
5 % n: int, number of neurons in the hidden layer of MLP.
6 % images: matrix of (image_dim, image_num), containing possibly
7 % preprocessed image data as input.
8 % labels: vector of (1, image_num), containing corresponding label of
9 % each image.
10 % train_num: int, number of training images.
11 % val_num: int, number of validation images.
12 % epochs: int, number of training epochs.
13 %
14 % Returns:
15 % net: object, containing trained network.
16 % accu_train: vector of (epochs, 1), containing the accuracy on training
17 % set of each epoch during training.
18 % accu_val: vector of (epochs, 1), containing the accuracy on validation
19 % set of each epoch during training.
20
21 % 1. Change the input to cell array form for sequential training
22 images_c = num2cell(images, 1);
23 labels_c = num2cell(labels, 1);
24
25 % 2. Construct and configure the MLP
26 net = patternnet(n);
27 net.divideFcn = 'dividetrain'; % input for training only
28 net.performParam.regularization = 0.25; % regularization strength
29 net.trainFcn = 'traingdx';
30 net.trainParam.epochs = epochs;
31 accu_train = zeros(epochs,1); % record accuracy on training set of
    each epoch
32 accu_val = zeros(epochs,1); % record accuracy on validation set of
    each epoch
33
34 % 3. Train the network in sequential mode
35 for i = 1 : epochs
36     display(['Epoch: ', num2str(i)])
37     idx = randperm(train_num); % shuffle the
        input
38     net = adapt(net, images_c(:,idx), labels_c(:,idx));
39     pred_train = round(net(images(:,1:train_num))); % predictions
        on training set
40     accu_train(i) = 1 - mean(abs(pred_train-labels(1:train_num)));
41     pred_val = round(net(images(:,train_num+1:end))); % predictions
        on validation set
42     accu_val(i) = 1 - mean(abs(pred_val-labels(train_num+1:end)));
43 end
44 % Visualization
45 fig = figure;
46 plot(1:epochs, accu_train*100, 'bo-', 'LineWidth', 1.5); hold on;
47 plot(1:epochs, accu_val*100, 'rs--', 'LineWidth', 1.5);
48 xlabel('epochs');
49 ylabel('Accuracy(%)');
50 legend('Training', 'Test');
51 title('Performance');
52 grid on;
53 saveas(fig, 'SeqMode_traingdx_AccuCurves.png');
54 end
55

```

```

56 load(' ../Data.mat')
57 X_all = [X_train,X_test];
58 Y_all = [Y_train,Y_test];
59
60 [ net, accu_train, accu_val ] = train_seq( 139, X_all, Y_all, 900, 10);

```

The final train accuracy is 100% and validation accuracy is 76%. Sequential mode achieved 100% training accuracy by epoch 10 (Figure 24), but test accuracy peaked earlier (epoch=7, 77%) and degraded by 1% thereafter, while batch mode exhibited slower training convergence but maintained stable test performance. Sequential updates introduced high variance in weight trajectories, whereas batch mode's averaged gradients yielded smoother optimization paths.

While sequential training theoretically suits streaming data scenarios, its inferior regularization efficacy and noisy optimization render it impractical for this task. Batch mode's stability and superior generalization (79% vs. 76% test accuracy) justify its recommendation.

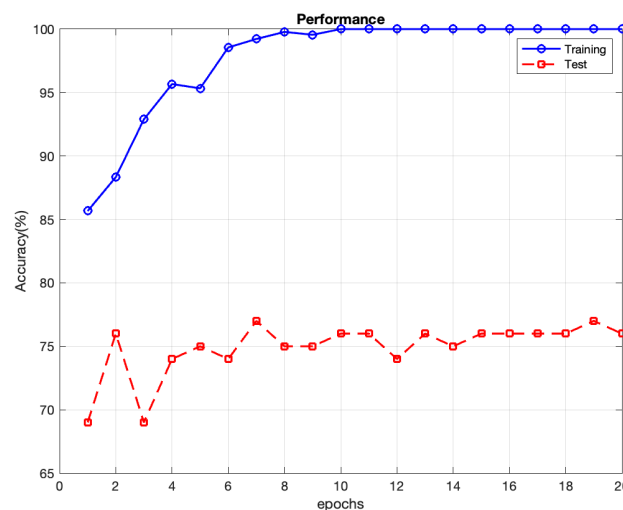


Figure 24: Accuracy of MLP on Sequential Training Mode (traingdx)

- (f) Try to propose a scheme that you believe could help to improve the performance of your MLP and please explain the reason briefly. (4 Marks)

#### Solution:

Based on experimental findings from Q3a-Q3e, the following integrated approach addresses key limitations of the current MLP:

##### 1. Input Standardization & PCA Dimensionality Reduction

Apply feature-wise standardization and retain 139 principal components (95% variance) via PCA.

**Rationale:** Mitigates pixel value scale disparities (0-255 → zero-mean/unit-variance)

Reduces input dimensions from 1024→139, removing high-frequency noise while preserving discriminative features.

##### 2. Architectural Optimization

Reduce hidden layer neurons to 70 (balance PCA-reduced input and parameter count), Replace **tansig** with **ReLU** activation in hidden layers. Add **Dropout** (rate=0.5) after the hidden layer.

**Rationale:** Smaller hidden layer (70 vs. 139) reduces parameters from 142K→9,870 (≈94% reduction), alleviating overfitting.

ReLU avoids gradient saturation, enabling larger weight magnitudes and effective regularization.

Dropout prevents co-adaptation, forcing distributed feature learning.

**3. Regularization**

Elastic Net Regularization - combine L1 ( $\lambda_1=0.1$ ) and L2 ( $\lambda_2=0.3$ ) penalties.

**Rationale:** Elastic Net promotes sparsity (L1) and weight shrinkage (L2), synergistically controlling complexity.

**4. Data Augmentation**

Generate augmented samples via random rotations( $\pm 10^\circ$ ), horizontal flips, and gaussian noise injection ( $\sigma = 0.05$ ).

**Rationale:** Increases effective training data and enhances robustness to geometric and photometric variations.