

# CS 615 Assignment2

Himanshu Gupta

May 2020

## 1 Theory

### 1.1

Given this hyperbolic tangent function:

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

In order to calculate the function:

$$\frac{\partial \tanh(x\theta)}{\partial \theta_j}$$

Let the  $z = x\theta$  so using chain rule the function becomes  $\frac{\partial \tanh(z)}{\partial \theta_j}$

$$\frac{\partial \tanh(z)}{\partial \theta_j} = \frac{\partial \tanh(z)}{\partial z_j} \frac{\partial (z_j)}{\partial \theta_j}$$

$$\begin{aligned} \frac{\partial \tanh(z)}{\partial z_j} &= \frac{((e^{z_j} + e^{-z_j})(e^{z_j} + e^{-z_j}) - (e^{z_j} - e^{-z_j})(e^{z_j} - e^{-z_j}))}{(e^{z_j} + e^{-z_j})^2} = \frac{((e^{z_j} + e^{-z_j})^2 - (e^{z_j} - e^{-z_j})^2)}{(e^{z_j} + e^{-z_j})^2} = 1 - \frac{(e^{z_j} - e^{-z_j})^2}{(e^{z_j} + e^{-z_j})^2} \\ &= 1 - (\tanh(z_j))^2 \end{aligned}$$

$$\frac{\partial (z_j)}{\partial \theta_j} = \frac{\partial (x\theta)}{\partial \theta_j} = x_j$$

Using the above two results:

$$\frac{\partial \tanh(z)}{\partial \theta_j} = (1 - (\tanh(z_j))^2)(x_j)$$

### 1.2

#### 1.2.1

$$\frac{\partial J}{\partial \theta_{j,k}} = \frac{\partial J}{\partial g(\text{net}_{o_k})} \cdot \frac{\partial g(\text{net}_{o_k})}{\partial \text{net}_{o_k}} \cdot \frac{\partial \text{net}_{o_k}}{\partial \theta_{j,k}}$$

In this case, our J (objective function is cross entropy), the output layer has softmax activation function, and the hidden layer has the linear activation function.

The Cross Entropy objective function is given below:

$$J = \sum_{k=1}^K -y_k \log \hat{y}_k$$

Since, we are using one-hot encoded distribution so only one  $y$  is 1. Let  $y_a = 1$  then

$$J = -\log \hat{y}_a$$

$$\text{so } \frac{\partial J}{\partial g(\text{net}_{o_k})} = \frac{-1}{\hat{y}_a} = -\frac{\sum_{k=1}^K e^{\text{net}_{o_k}}}{e^{\text{net}_{o_a}}}$$

$$g(\text{net}_{o_k}) = \frac{e^{\text{net}_{o_a}}}{\sum_{k=1}^K e^{\text{net}_{o_k}}}$$

$$\text{if } j \neq a, -\frac{g(\text{net}_{o_k})}{\text{net}_{o_k}} = -\frac{(e^{\text{net}_{o_a}})(e^{\text{net}_{o_k}})}{(\sum_{k=1}^K e^{\text{net}_{o_k}})^2}$$

$$\text{if } j = a, \frac{g(\text{net}_{o_k})}{\text{net}_{o_k}} = \frac{(e^{\text{net}_{o_a}})(\sum_{k=1}^K e^{\text{net}_{o_k}}) - (e^{\text{net}_{o_a}})(e^{\text{net}_{o_a}})}{(\sum_{k=1}^K e^{\text{net}_{o_k}})^2}$$

$$\frac{\partial \text{net}_{o_k}}{\partial \theta_{j,k}} = x_j \text{ for } j = a, \text{ otherwise zero}$$

Therefore, using the results of the above equations:

Case 1:  $j \neq a$

$$\frac{\partial J}{\partial \theta_{j,k}} = \left( \frac{\sum_{k=1}^K e^{\text{net}_{o_k}}}{e^{\text{net}_{o_a}}} \right) \left( \frac{(e^{\text{net}_{o_a}})(e^{\text{net}_{o_k}})}{(\sum_{k=1}^K e^{\text{net}_{o_k}})^2} \right) (x_j) = (x_j)(\hat{y}_k)$$

Similarly, Case 2:  $j = a$

$$\frac{\partial J}{\partial \theta_{j,k}} = -\left( \frac{\sum_{k=1}^K e^{\text{net}_{o_k}}}{e^{\text{net}_{o_a}}} \right) \left( \frac{(e^{\text{net}_{o_a}})(\sum_{k=1}^K e^{\text{net}_{o_k}}) - (e^{\text{net}_{o_a}})(e^{\text{net}_{o_a}})}{(\sum_{k=1}^K e^{\text{net}_{o_k}})^2} \right) (x_j) = (x_j)(\hat{y}_k - 1)$$

$$\text{Hence, } \frac{\partial J}{\partial \theta_{j,k}} = x_j(\hat{y}_k - y_k)$$

### 1.2.2

$$\frac{\partial J}{\partial \beta_{i,j}} = \frac{\partial J}{\partial g(\text{net}_o)} \cdot \frac{\partial g(\text{net}_o)}{\partial \text{net}_o} \cdot \frac{\partial \text{net}_o}{\partial g(\text{net}_{h_j})} \cdot \frac{\partial g(\text{net}_{h_j})}{\partial \text{net}_{h_j}} \cdot \frac{\partial \text{net}_{h_j}}{\partial \beta_{i,j}}$$

The first two terms are already calculated in the previous question.

$$\frac{\partial \text{net}_o}{\partial g(\text{net}_{h_j})} = \frac{\partial \text{net}_o}{\partial h_j} = \theta_j$$

As activation function of hidden layer is linear ( $g(\text{net}_{h_j}) = \text{net}_{h_j}$ ).

$$\frac{\partial g(\text{net}_{h_j})}{\partial \text{net}_{h_j}} = 1$$

$$\frac{\partial \text{net}_{h_j}}{\partial \beta_{i,j}} = x_i$$

Combining all the results from the above equations, we get

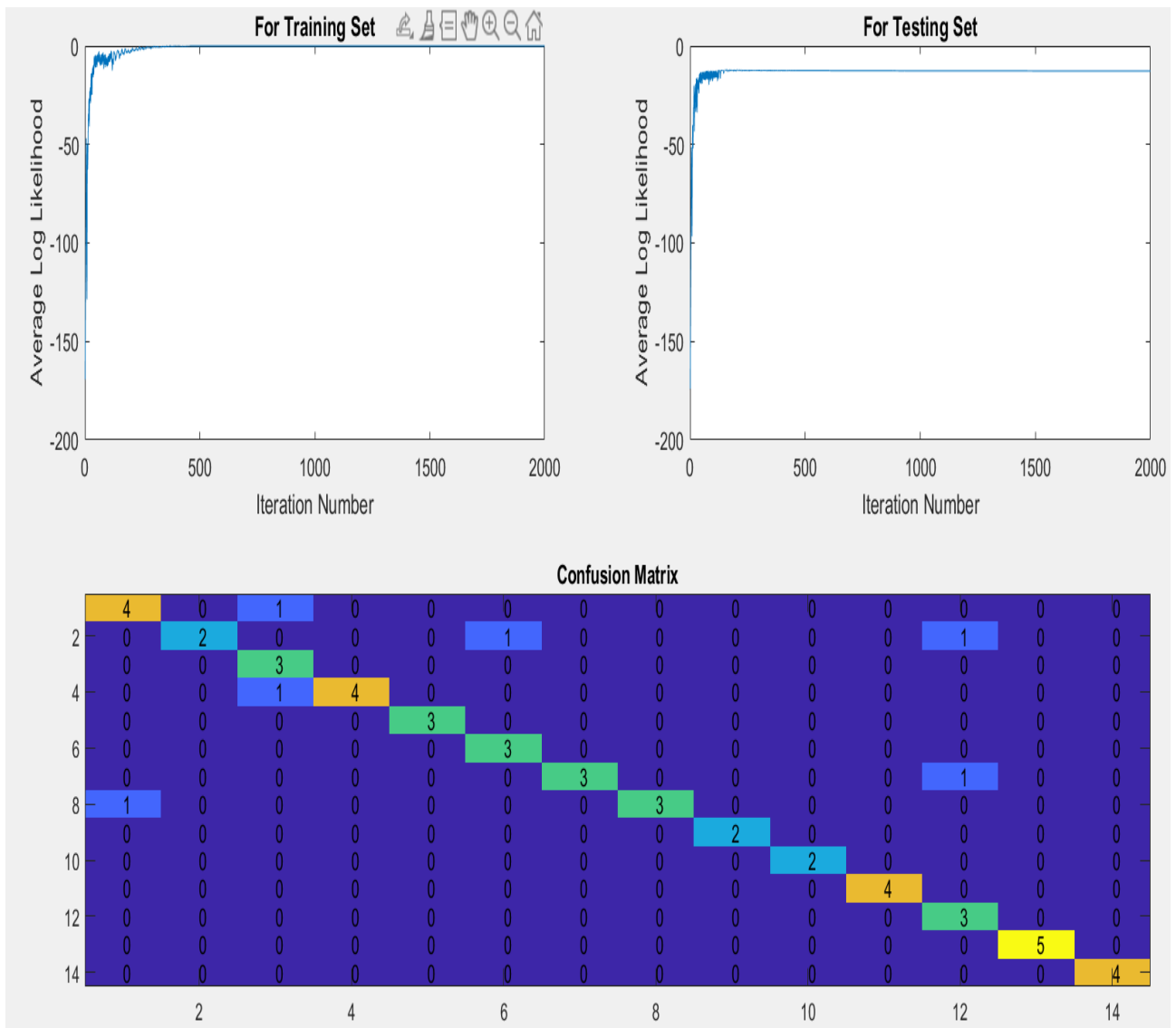
$$\text{Case 1: } j \neq a, \frac{\partial J}{\partial \beta_{i,j}} = \left( \frac{\sum_{k=1}^K e^{\text{net}_{o_k}}}{e^{\text{net}_{o_a}}} \right) \left( \frac{(e^{\text{net}_{o_a}})(e^{\text{net}_{o_k}})}{(\sum_{k=1}^K e^{\text{net}_{o_k}})^2} \right) (\theta_j)(x_i) = (x_i)(\hat{y}_k)(\theta_j)$$

$$\text{Similarly, Case 2: } j = a, \frac{\partial J}{\partial \beta_{i,j}} = -\left( \frac{\sum_{k=1}^K e^{\text{net}_{o_k}}}{e^{\text{net}_{o_a}}} \right) \left( \frac{(e^{\text{net}_{o_a}})(\sum_{k=1}^K e^{\text{net}_{o_k}}) - (e^{\text{net}_{o_a}})(e^{\text{net}_{o_a}})}{(\sum_{k=1}^K e^{\text{net}_{o_k}})^2} \right) (x_i)(\theta_j)$$

$$= (x_i)(\hat{y}_k - 1)(\theta_j)$$

$$\text{Hence, } \frac{\partial J}{\partial \beta_{i,j}} = (x_i)(\hat{y}_k - y_k)(\theta_j)$$

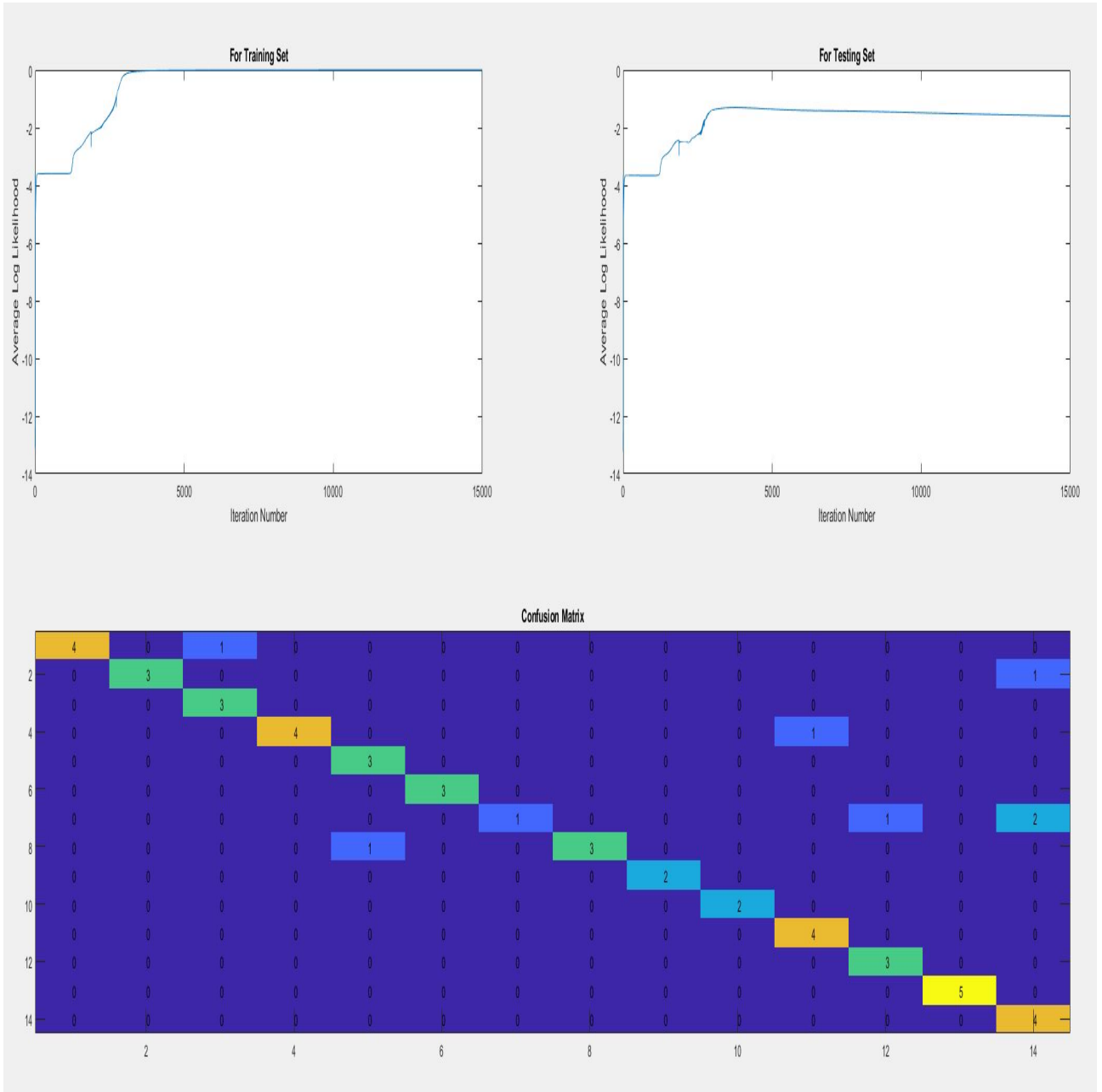
## 2 Shallow ANN



The initial values of  $\theta$ s are randomly chosen between  $[0,1]$ , then made even smaller by multiplying with 0.001. Accuracy = 88.24%, Termination Criteria is max iterations = 2000, Learning Rate ( $\eta$ ) = 0.8, and L2 Regularization term ( $\alpha$ ) = 0.00001.

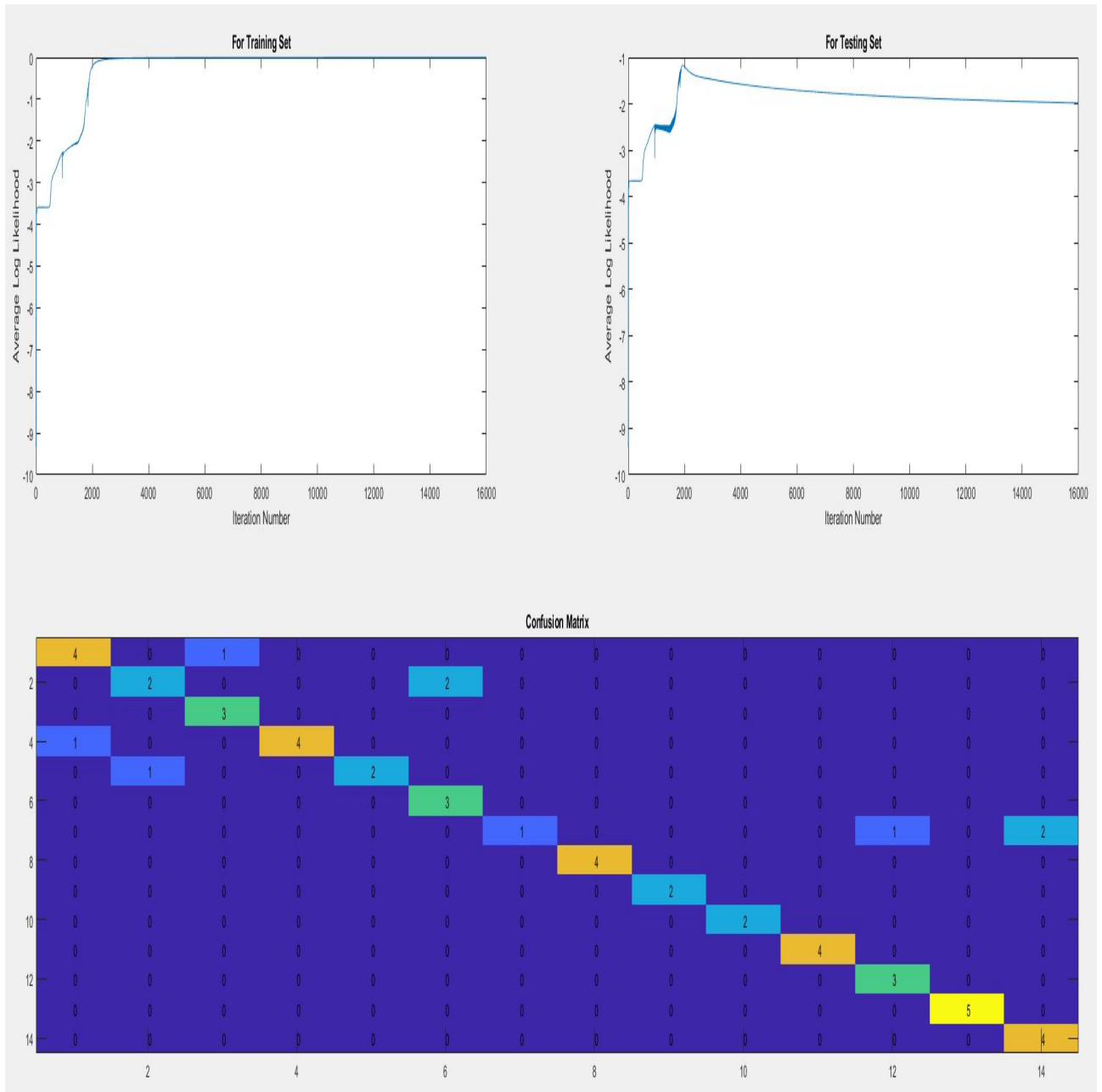
### 3 Multi-layers ANN

#### 3.1 Network Configuration 1



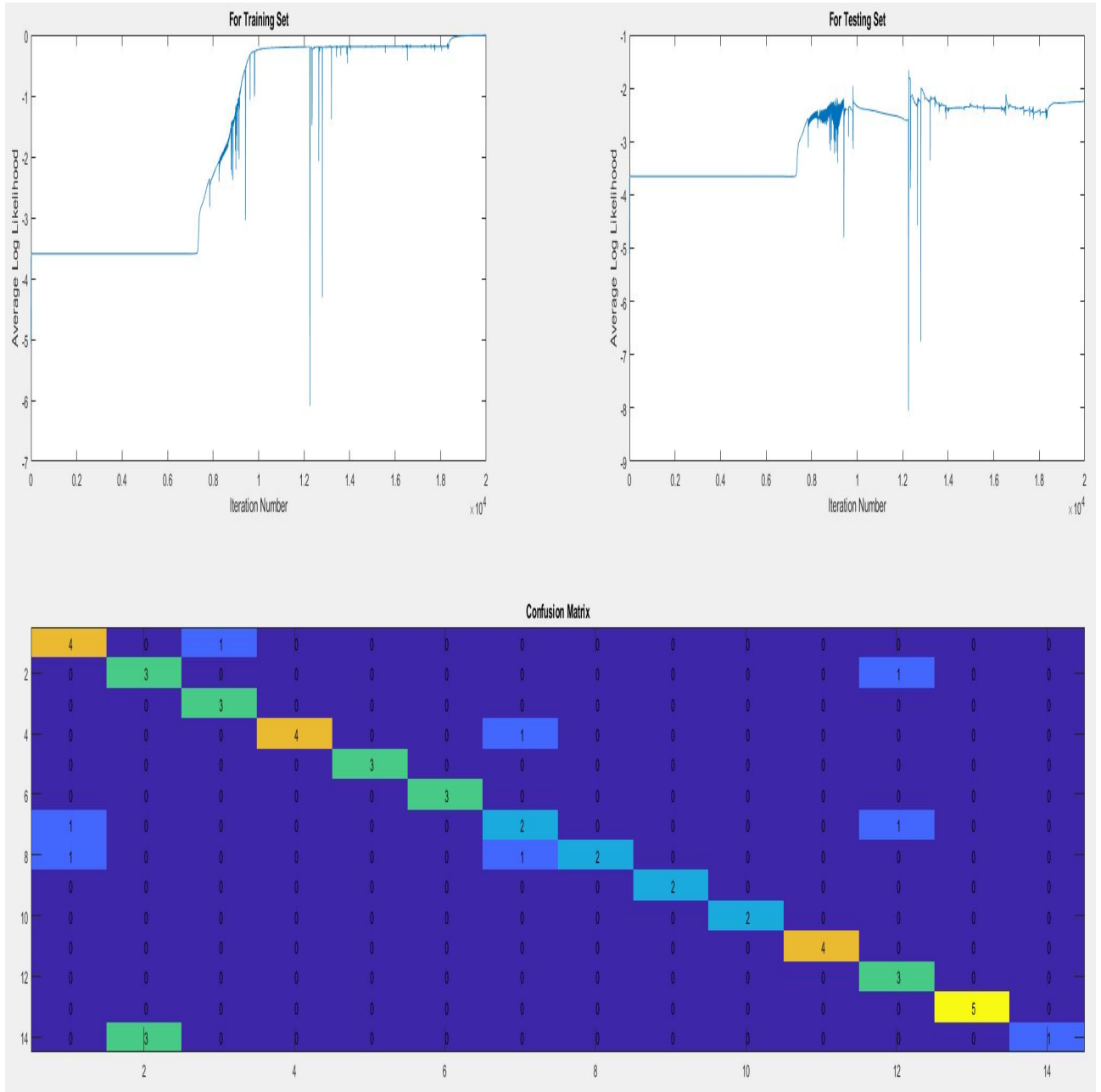
The initial values of  $\theta$ s are randomly chosen between  $[0,1]$ , then made even smaller by multiplying with 0.001. Hidden Layers nodes are  $[800, 400, 200]$ , Learning rate ( $\eta$ ) = 0.55, Termination Criteria is max iterations = 15000, Accuracy = 86%, and L2 regularization term ( $\alpha$ ) = 0.00001.

### 3.2 Network Configuration 2



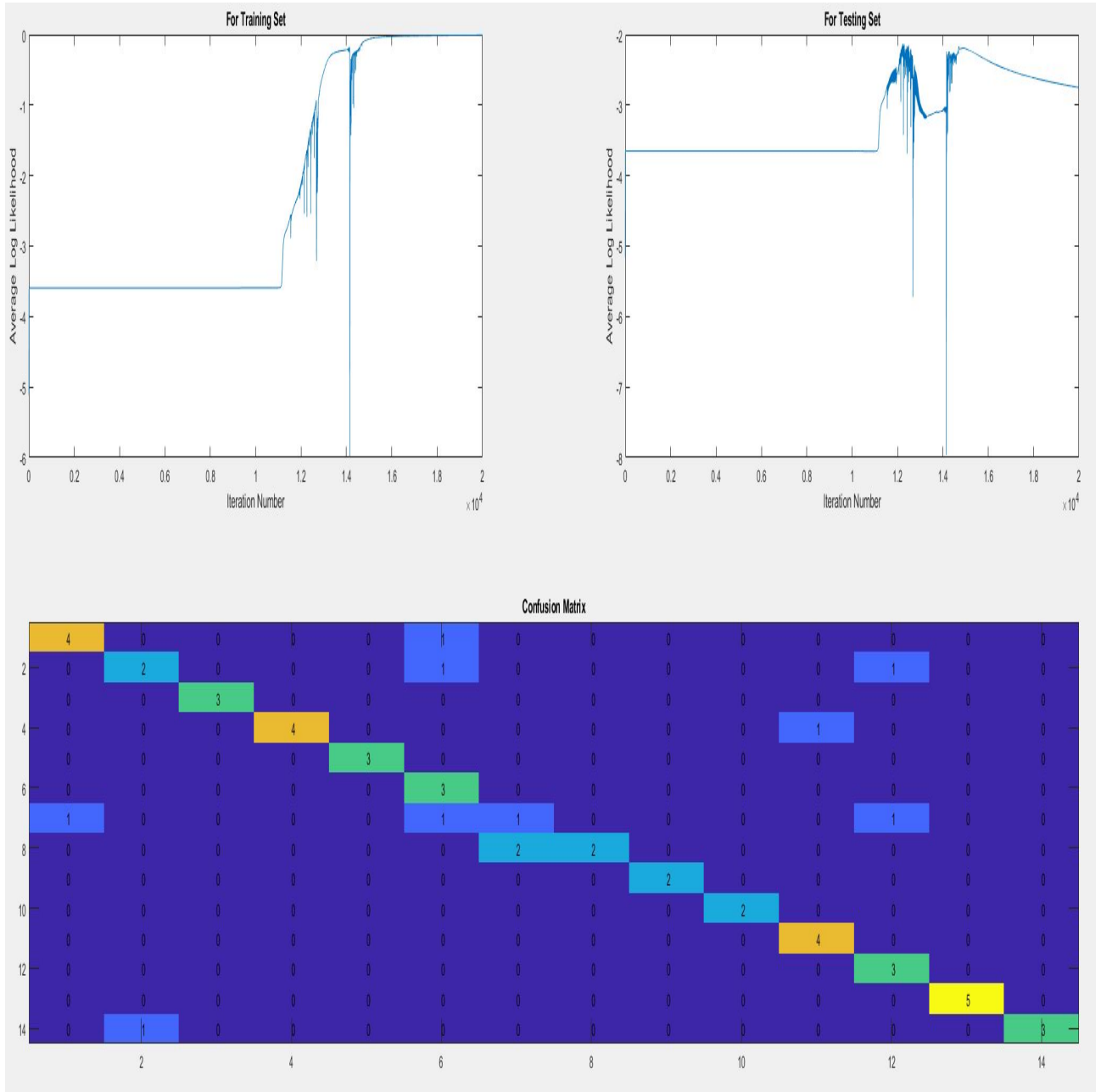
The initial values of  $\theta$ s are randomly chosen between  $[0,1]$ , then made even smaller by multiplying with 0.001. Hidden layers nodes are  $[1000,400,100]$ , Accuracy = 84%, Learning rate ( $\eta$ ) = 0.8, Termination Criteria is max iterations = 16000, and L2 regularization term ( $\alpha$ ) = 0.00001.

### 3.3 Network Configuration 3



The initial values of  $\theta_s$  are randomly chosen between  $[0,1]$ , then made even smaller by multiplying with 0.001. Hidden layers nodes are  $[1200 \ 400 \ 100 \ 50]$ , Accuracy = 80%, Learning rate ( $\eta$ ) = 0.91, Termination Criteria is max iteration = 20000 and L2 regularization term ( $\alpha$ ) = 0.00001.

### 3.4 Network Configuration 4



The initial values of  $\theta$ s are randomly chosen between  $[0,1]$ , then made even smaller by multiplying with 0.001. Hidden layers are  $[1000,200,100,50]$ , 80.39%, Learning Rate ( $\eta$ ) = 0.92, Termination Criteria is max iteration = 20000 and L2 regularization term ( $\alpha$ ) = 0.00001.