

## Table of Contents

<u>Title</u>	<u>Page no.</u>
Introduction	2
Methodology	2
Results	4
Conclusion and References	18

## ❏ INTRODUCTION

### ➔ **What is the Deep Neural Network?**

A Deep Neural Network consists of one or more hidden layers it is also known as multi layers perceptron which is the most trivial form of neural network. Deep Neural Networks are preferred to work with unstructured data such as images, audio and video. They can be also used to handle structured ( or tabular) data but most of the time this category of data is dealt with machine learning algorithms. The reason why Deep Neural Networks are so popular is that they can work with huge amounts of data

### ➔ **MNIST Dataset**

The dataset consists of 60000 handwritten images each of which is  $28 * 28$  pixels in size. All images are in grayscale format. We will use 10000 images to train and 1000 images to test your model.

## ❏ METHODOLOGY

### **1. Normalization**

The original data set is in a grayscale format in which each pixel values ranges between [0-255]. I normalized the dataset to bring all the feature values in the range [0-1].

### **2. Reshaping the data**

The original dataset is in the form (60000, 28, 28). I converted all the (28, 28) images to a vector of  $28*28=784$  pixels to train a neural network.

### **3. Define structure of neural network**

The network consists of an input layer, a hidden layer and an output layer. The input layers consist of 784 nodes for 784 features in each training example. For hidden layers, I will experiment with multiple numbers of nodes. And, finally the output layer consists of 10 nodes. Each output node represents a probability.

## 4. Initialize Parameters

Let us use  $n[0], n[1], n[2]$  to denote the number of nodes in each layer of the neural network. In my case  $n[0] = 784$ , and  $n[2] = 10$ . For the hidden layers I initialized the weights and bias matrices as  $W^{[1]}$  and  $b^{[1]}$ . The size of  $W^{[1]}$  is  $(n[1], n[0])$  and size of  $b^{[1]}$  is  $(n[1], 1)$ . For the output layer I will define parameters  $W^{[2]}$  and  $b^{[2]}$  of size  $(n[2], n[1])$  and  $(n[2], 1)$  respectively.

## 5. Training

In the forward propagation part I will first calculate output for the hidden layer. The equations will be

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = g(Z^{[1]})$$

For output layers I'll calculate the final output using result of the previous layers as

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g(Z^{[2]})$$

Here  $g$  is the activation function, in this case it will be ReLU activation function. In Backpropagation I will calculate the gradient of final output with respect to the initialized parameters in previous layers.

$A^{[2]}$  is the final output. The loss is defined as

$$\frac{dL}{dZ^{[2]}} = A^{[2]} - Y$$

Using the above formula we can calculate the gradients in the output and hidden layers. Gradient for the output layer

$$\frac{dL}{dZ^{[2]}} = \frac{dL}{dA^{[2]}} \times \frac{dA^{[2]}}{dZ^{[2]}}$$

$$\frac{dL}{db^{[2]}} = \frac{dL}{dZ^{[2]}} \times \frac{dZ^{[2]}}{db^{[2]}}$$

Gradients for the hidden layer

$$j\,dZ^{[1]0)}$$



$$= W^{[2]T} \times \frac{dL}{dZ^{[2]}} - \frac{dL}{dZ^{[2]}}$$

$$\frac{dL}{dW^{[1]}} \times \frac{dL}{dZ^{[1]}} \times X^{T[1]}$$

$$\frac{dL}{db^{[1]}} \times \sum_j \frac{dL}{dZ^{[1](j)}} \quad [1](j)$$

Lastly using the above equation you can update the parameters as follows

$$W^{[1]} = W^{[1]} - \alpha \times \frac{dL}{dW^{[1]}}$$

$$b^{[1]} = b^{[1]} - \alpha \times \frac{dL}{db^{[1]}}$$

$$W^{[2]} = W^{[2]} - \alpha \times \frac{dL}{dW^{[2]}}$$

$$b^{[2]} = b^{[2]} - \alpha \times \frac{dL}{db^{[2]}}$$

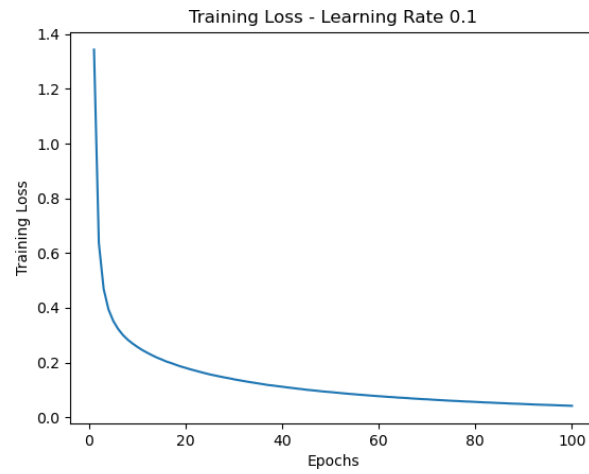
We will run the above code using multiple learning rates  $\alpha$  and test with various nodes in the hidden layers.

## ❑ RESULTS:

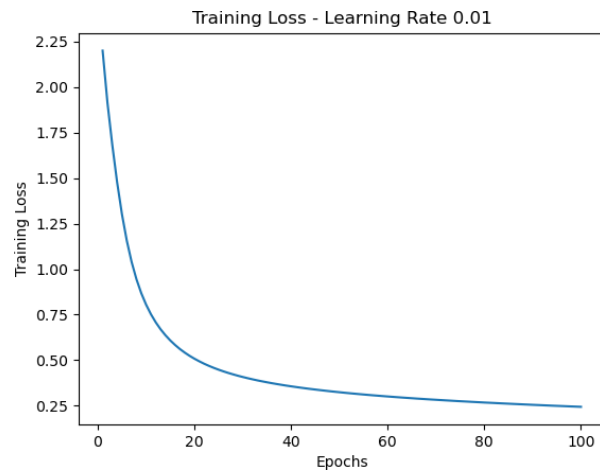
The overall accuracy of models run on binary images is lesser than the model's accuracy with grayscale images.

## ----- TASK 1 -----

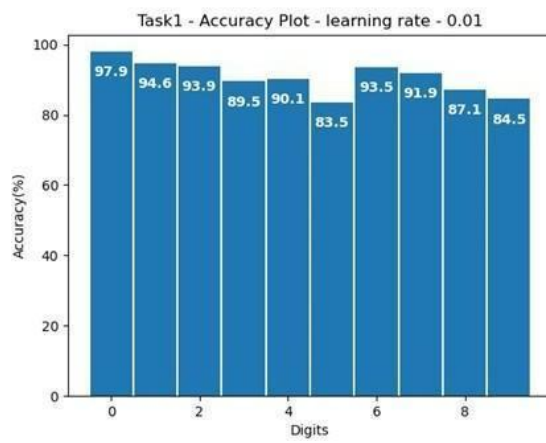
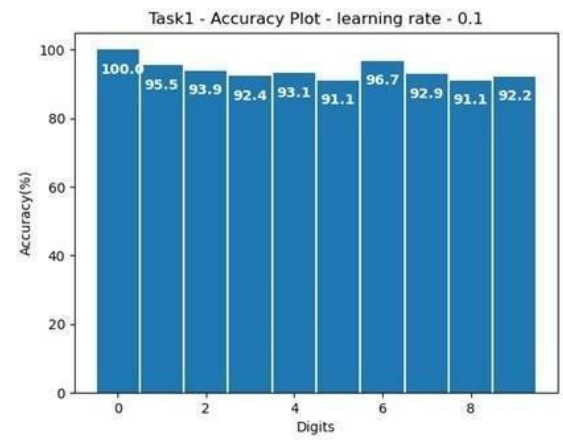
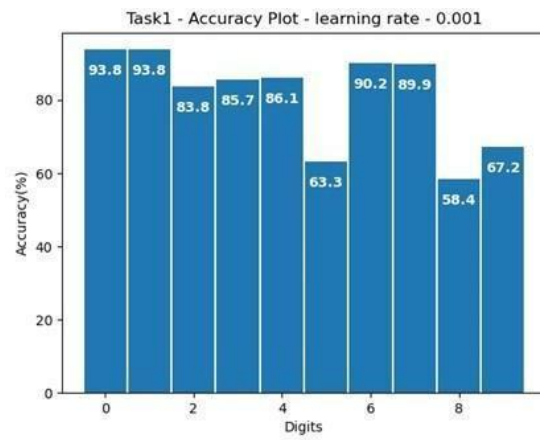
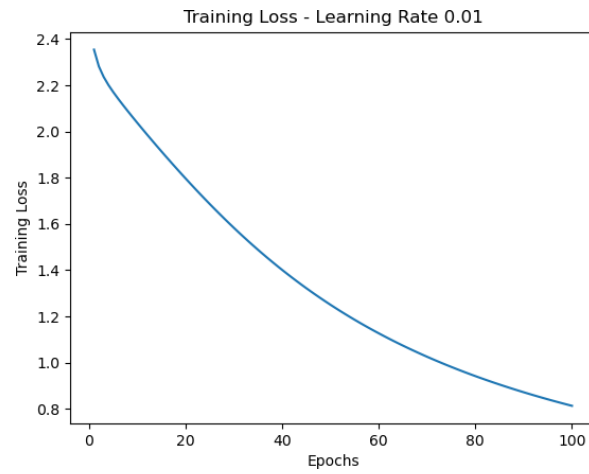
In task 1, the learning rate is 0.1, and the training is done for 100 epochs. The accuracy of this model is 93.3%.



The second learning rate is 0.01. Decreasing the value of the learning rate decreased accuracy. The final accuracy is 90.1%.



The third value of the learning parameter is 0.001. The model does not converge in 100 epochs, as shown in the graph below. The accuracy decreased to 80.8%.





The percentage accuracy for each model with each learning rate are shown above. On the same test set each time the model gave different accuracy for different learning rates. For learning rate 0.1 each digit for correctly predicted more than 90%. For learning rates 0.01 the accuracy slightly decreased but seven out of ten digits are classified correctly more than 90% of the time

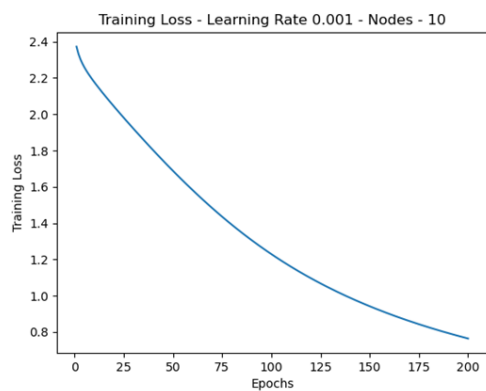
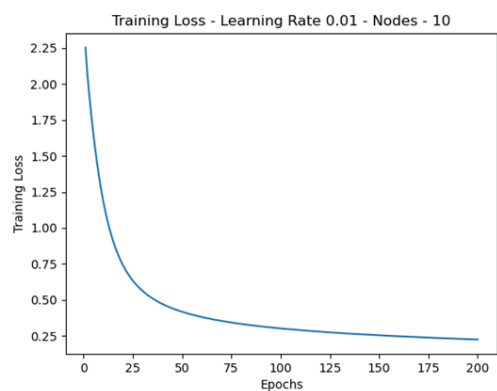
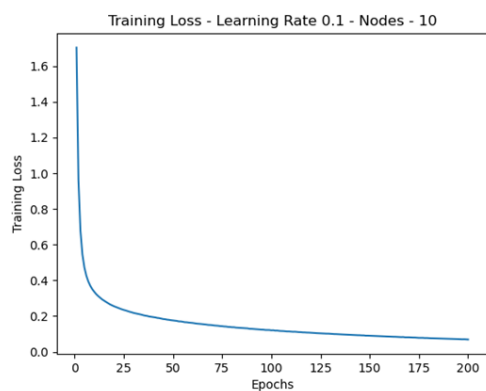
while remaining were classified correctly at least 80% of the time. When the learning rate is 0.001 the accuracy falls below 70% as shown in the graph in case of digits 5, 8 and 9. For the remaining digits the model accurately predicted more than 90% of the time.

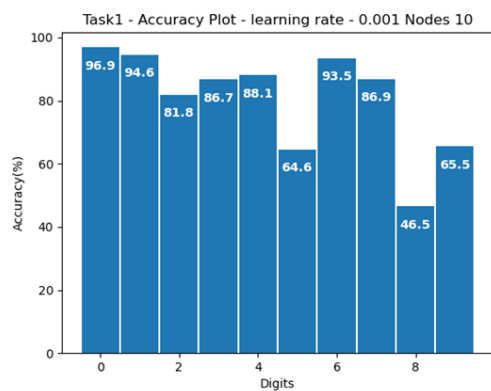
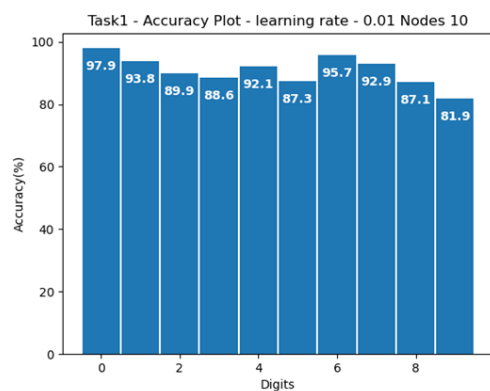
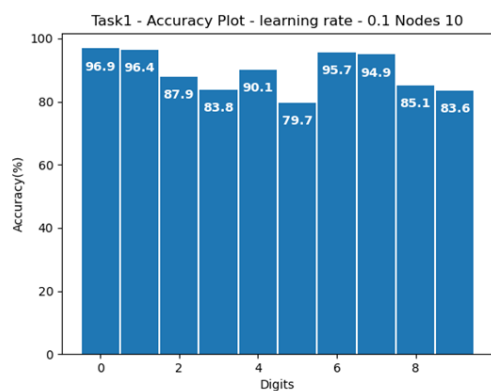
## ----- TASK 2 -----

For the task we will use 0.1, 0.01 and 0.01 as learning parameters with different combinations of nodes in the hidden layer.

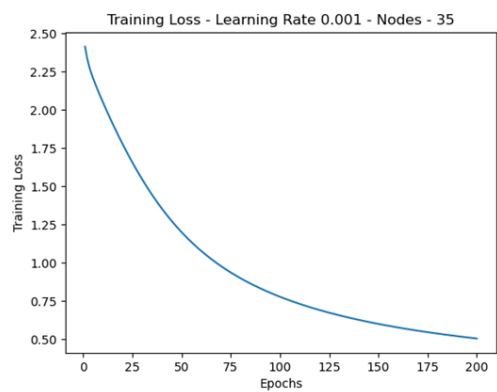
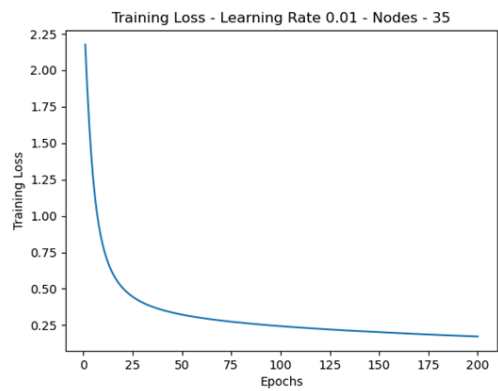
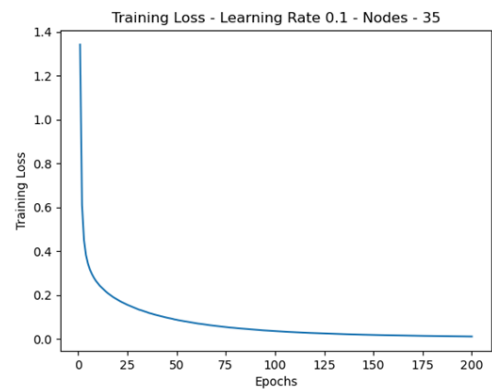
**Nodes = 10**

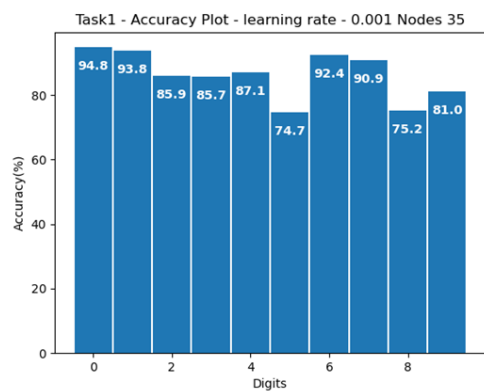
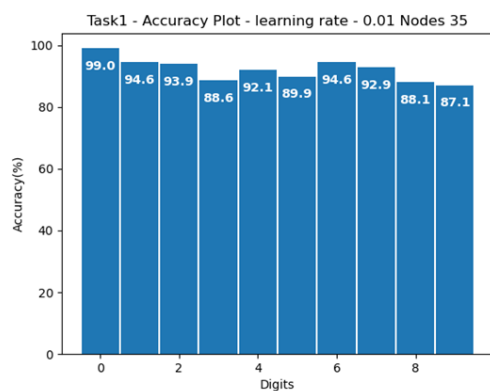
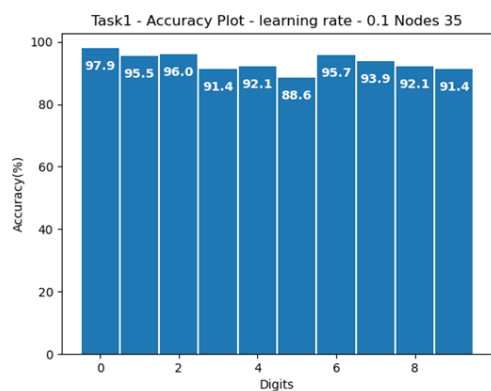
The accuracy plots are as follows:



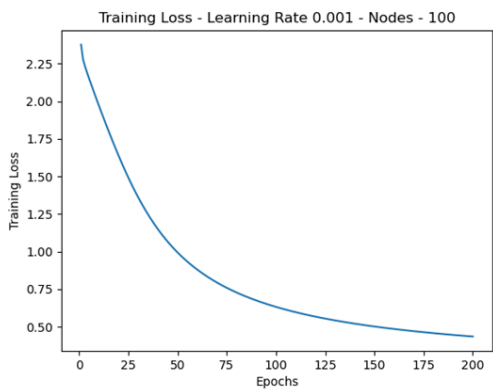
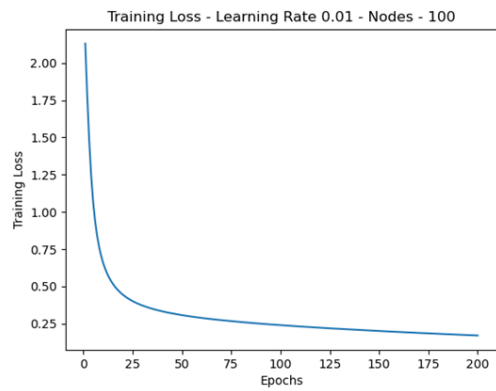
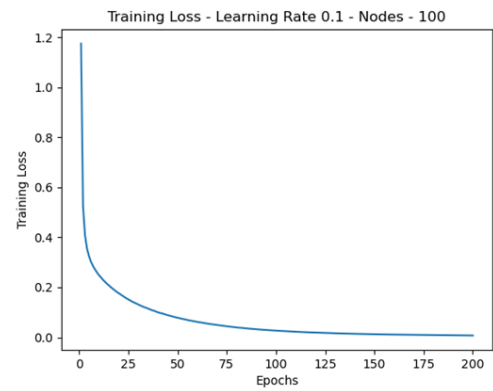


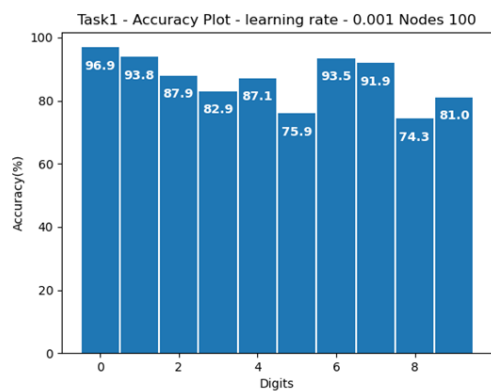
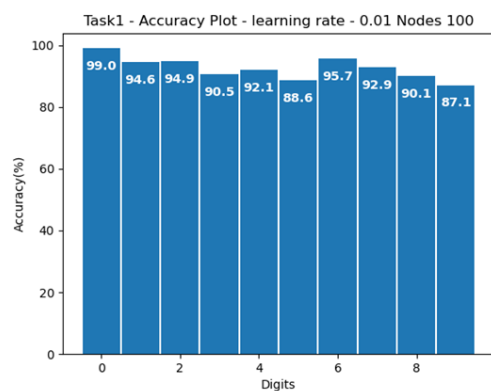
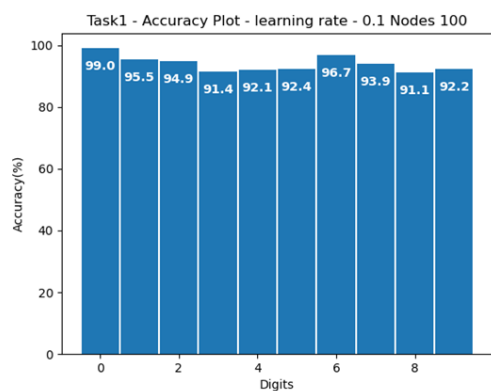
Nodes = 35



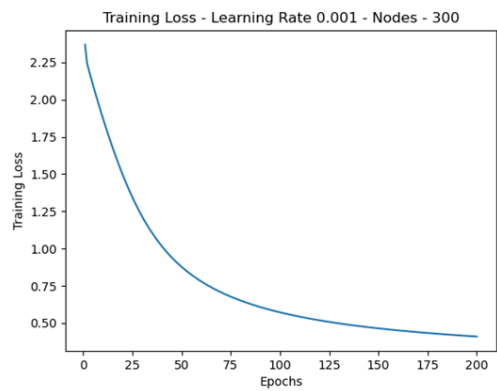
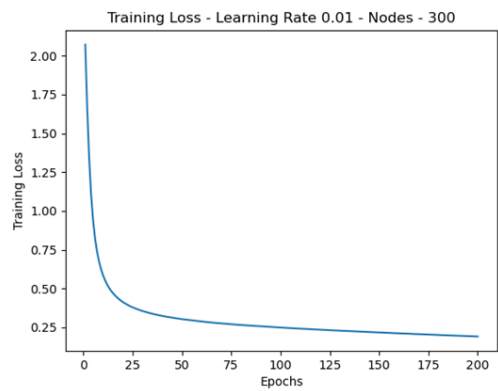
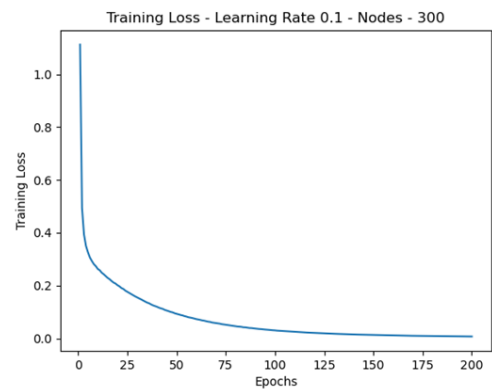


**Nodes = 100**

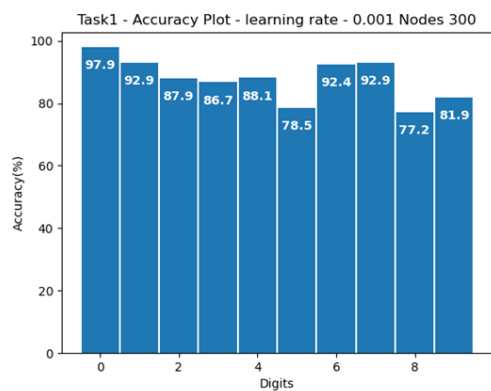
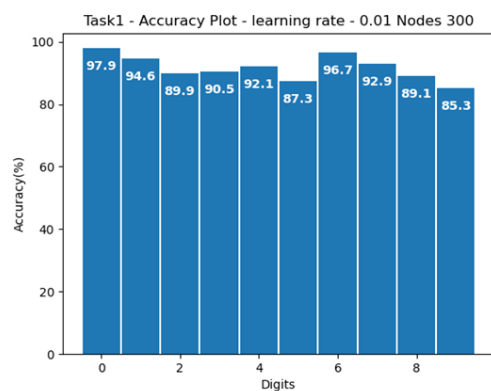
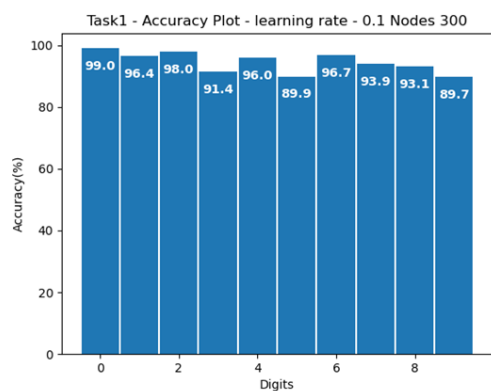




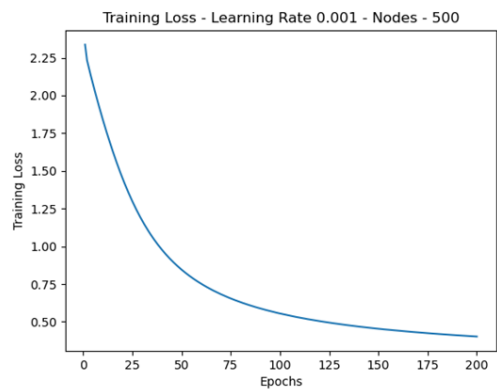
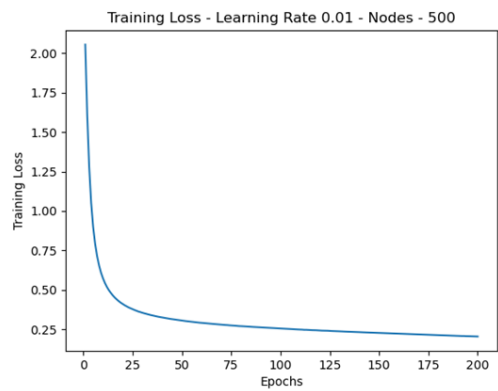
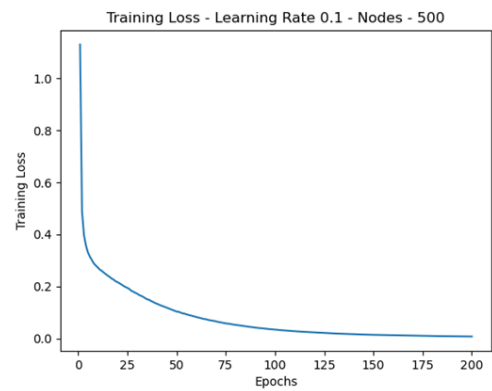
**Nodes = 300**

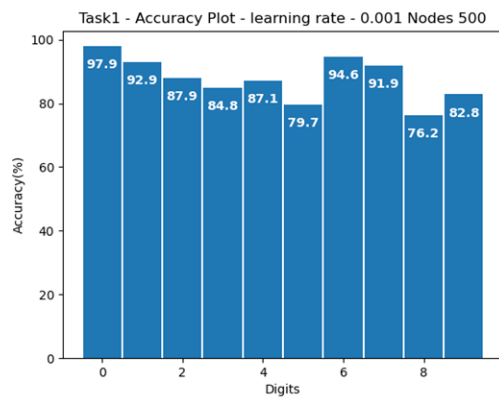
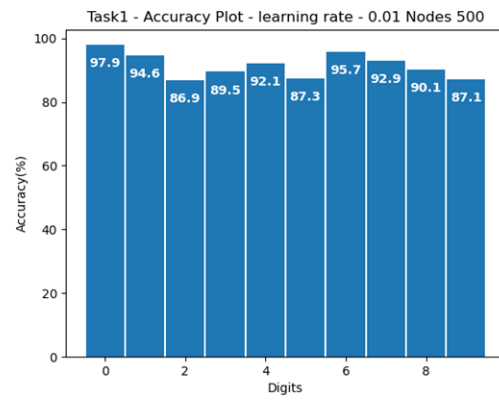
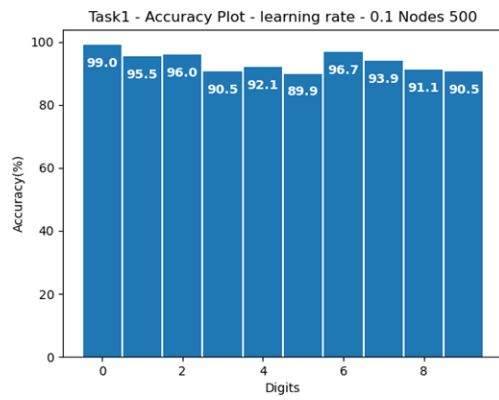






**Nodes = 500**





The accuracy line plots confirm the fact that lowering the value of the learning parameter also slows learning. From the plots you can find that the learning rate 0.01 is most suited for modelling as the training loss curve is smooth.

**Accuracy Table**

Accuracy(%)	Learning Rate		
Nodes	0.1	0.01	0.001
10	89.8	89.7	83.9
35	93.9	91.1	86
100	94	91.8	86.9
300	94.1	91.2	87.1
500	94.3	91.2	87.4

## ----- TASK 4 -----

### **Comparing SVM and Single Layer Perceptron Results**

In our SVM model we achieved maximum accuracy using RBF kernel taking C equals to 100 and gamma equals to 0.01 on grayscale images. The maximum accuracy in case of Single Layer Perceptron is 89.2 percent on grayscale images with learning rate as 0.5.

In case of binary images using threshold value as 150 the maximum accuracy we will get in SVM model is 88.6 percent whereas in case of perceptron the maximum accuracy is 86.4 percent. For binary images both models performed equally well.

Perceptron models irrespective of learning rate and on the same dataset take more time than SVM. This is due to the fact that perceptron requires more iterations to converge to minimum as compared to SVMs.

Your best MLP model is the one with 500 nodes and learning rates set to 0.1 which is better than both SVM and Perceptron models.



## ❑ CONCLUSION:

The multi-layer perceptron is a strong classifier in comparison with both Support Vector Machines and Perceptron models. The learning in case of MLP is highly affected by the learning rate and the number of nodes in the hidden units. In the Multilayer perceptron we have found that the accuracy increases with the number of nodes and it is directly proportional to the learning rate. Lower the learning rate lower is the model's accuracy.

## ❑ REFERENCES:

[https://scikit-learn.org/stable/user\\_guide.html](https://scikit-learn.org/stable/user_guide.html)  
<http://yann.lecun.com/exdb/mnist/>  
<https://seaborn.pydata.org/>  
<https://matplotlib.org/stable/contents.html>  
<https://www.tensorflow.org/guide>