

# HOMEWORK 4

MADHUSHREE NIJAGAL  
9084490524

**Instructions:** Use this latex file as a template to develop your homework. Submit your homework on time as a single zip (containing the pdf and the code) file to Canvas. Please, check Piazza for updates about the homework. It is recommended that you use Python for your solutions. You are allowed to use the libraries numpy, matplotlib, pandas and pytorch.

## 1 Questions (50 pts)

1. (10 Points) Suppose the world generates a single observation  $x \sim \text{multinomial}(\theta)$ , where the parameter vector  $\theta = (\theta_1, \dots, \theta_k)$  with  $\theta_i \geq 0$  and  $\sum_i^k \theta_i = 1$ . Note  $x \in \{1, \dots, k\}$ . You know  $\theta$  and want to predict  $x$ . Call your prediction  $\hat{x}$ . What is your expected 0-1 loss:

$$\mathbb{E}[1\{\hat{x} \neq x\}]$$

using the following two prediction strategies respectively? Prove your answer.

- (a) (5 pts) Strategy 1:  $\hat{x} \in \text{argmax}_x \theta_x$ , the outcome with the highest probability.

As given in the problem statement, we know the value of  $\theta$  and from the strategy we can imply that the value of  $\theta$  is maximum. When  $\theta$  is maximum, we have  $P(\hat{x} = x)$  and the loss is 0 and when we have values of  $\theta$  less than the maximum value we have  $P(\hat{x} \neq x)$  and the loss is 1. That is when  $x$  is misclassified ( $\hat{x} \neq x$ ), the loss is 1 and 0 when correctly classified ( $\hat{x} = x$ ). Therefore we can write the expected loss function

$$\mathbb{E}[1\{\hat{x} \neq x\}]$$

as:

$$\begin{aligned} &= \sum x * P(x) = 1 * P(\hat{x} \neq x) + 0 * P(\hat{x} = x) \\ &= 1 * P(\hat{x} \neq x) = 1 - P(\hat{x} = x) = 1 - \theta_{max} \end{aligned}$$

- (b) (5 pts) Strategy 2: You mimic the world by generating a prediction  $\hat{x} \sim \text{multinomial}(\theta)$ . (Hint: your randomness and the world's randomness are independent)

As mentioned above, the loss function can be written as a function of all the values that are misclassified, that is when loss is 1 and values that are correctly classified, that is when loss is 0. Here the values of  $\hat{x}$  and  $x$  are independent. Therefore the expected loss function can be written as

$$\begin{aligned} &\mathbb{E}[1\{\hat{x} \neq x\}] \\ &= \sum x * P(x) = 1 * P(\hat{x} \neq x) + 0 * P(\hat{x} = x) \\ &= 1 - [P(\hat{x} = 1, x = 1) + P(\hat{x} = 2, x = 2) + \dots + P(\hat{x} = k, x = k)] \\ &= 1 - [\theta_1^2 + \theta_2^2 + \dots + \theta_k^2] \\ &= 1 - \sum_{i=1}^k \theta_i^2 \end{aligned}$$

2. (10 points) Like in the previous question, the world generates a single observation  $x \sim \text{multinomial}(\theta)$ . Let  $c_{ij} \geq 0$  denote the loss you incur, if  $x = i$  but you predict  $\hat{x} = j$ , for  $i, j \in \{1, \dots, k\}$ .  $c_{ii} = 0$  for all  $i$ . This is a way to generalize different costs on false positives vs false negatives from binary classification to multi-class classification. You want to minimize your expected loss:

$$\mathbb{E}[c_{x\hat{x}}]$$

Derive your optimal prediction  $\hat{x}$ .

We are given that the  $x = i$  and  $\hat{x} = j$  and associated cost is  $c_{ij} \geq 0$  and if  $x = \hat{x} = i$ , the associated cost is

$c_{ii} = 0$ . We can represent the same information in the following matrix as:

$$\begin{bmatrix} 0 & c_{12} & c_{13} \\ c_{21} & 0 & c_{23} \\ c_{31} & c_{32} & 0 \end{bmatrix}$$

Therefore the loss function can be represented as  $\mathbb{E}[c_{x\hat{x}}] = \sum_{j=1}^k \sum_{i=1}^k c_{ij} \theta_i \theta_j$ .

The  $\theta_i$  and  $\theta_j$  are the probabilities of  $x$  and  $\hat{x}$  respectively. To minimize the expected loss, we need to minimize the function by finding the smallest value of  $c_{ij}$  for every  $i$  value and vary  $j$ , the function is represented as:

$$\hat{x} \in \operatorname{argmin}_j (\sum_{i=1}^k \theta_i c_{ij})$$

3. (30 Points) The Perceptron Convergence Theorem shows that the Perceptron algorithm will not make too many mistakes as long as every example is “far” from the separating hyperplane of the target halfspace. In this problem, you will explore a *variant* of the Perceptron algorithm and show that it performs well (given a little help in the form of a good initial hypothesis) as long as every example is “far” (in terms of angle) from the separating hyperplane of the *current hypothesis*.

Consider the following variant of Perceptron:

- Start with an initial hypothesis vector  $\mathbf{w} = \mathbf{w}^{\text{init}}$ .
- Given example  $\mathbf{x} \in \mathbb{R}^n$ , predict according to the linear threshold function  $\mathbf{w} \cdot \mathbf{x} \geq 0$ .
- Given the true label of  $\mathbf{x}$ , update the hypothesis vector  $\mathbf{w}$  as follows:
  - If the prediction is correct, leave  $\mathbf{w}$  unchanged.
  - If the prediction is incorrect, set  $\mathbf{w} \leftarrow \mathbf{w} - (\mathbf{w} \cdot \mathbf{x})\mathbf{x}$ .

So the update step differs from that of Perceptron shown in class in that  $(\mathbf{w} \cdot \mathbf{x})\mathbf{x}$  (rather than  $\mathbf{x}$ ) is added or subtracted to  $\mathbf{w}$ . (Note that if  $\|\mathbf{x}\|_2 = 1$ , then this update causes vector  $\mathbf{w}$  to become orthogonal to  $\mathbf{x}$ , i.e., we add or subtract the multiple of  $\mathbf{x}$  that shrinks  $\mathbf{w}$  as much as possible.)

Suppose that we run this algorithm on a sequence of examples that are labeled according to some linear threshold function  $\mathbf{v} \cdot \mathbf{x} \geq 0$  for which  $\|\mathbf{v}\|_2 = 1$ . Suppose moreover that

- Each example vector  $\mathbf{x}$  has  $\|\mathbf{x}\|_2 = 1$ ;
- The initial hypothesis vector  $\mathbf{w}^{\text{init}}$  satisfies  $\|\mathbf{w}^{\text{init}}\|_2 = 1$  and  $\mathbf{w}^{\text{init}} \cdot \mathbf{v} \geq \gamma$  for some fixed  $\gamma > 0$ ;
- Each example vector  $\mathbf{x}$  satisfies  $\frac{|\mathbf{w} \cdot \mathbf{x}|}{\|\mathbf{w}\|_2} \geq \delta$ , where  $\mathbf{w}$  is the current hypothesis vector when  $\mathbf{x}$  is received. (Note that for a unit vector  $\mathbf{x}$ , this quantity  $\frac{|\mathbf{w} \cdot \mathbf{x}|}{\|\mathbf{w}\|_2}$  is the cosine of the angle between vectors  $\mathbf{w}$  and  $\mathbf{x}$ .)

Show that under these assumptions, the algorithm described above will make at most  $\frac{2}{\delta^2} \ln(1/\gamma)$  many mistakes.

So we are given that for every true value of  $\mathbf{x}$ , if the prediction is incorrect we have  $\mathbf{w}$  updated by  $\mathbf{w} - (\mathbf{w} \cdot \mathbf{x})\mathbf{x}$

The convergence upper bound proof:

Therefore for every mistake the update done can be represented as below:  $\mathbf{w}_{t+1} = \mathbf{w}_t - (\mathbf{w}_t \cdot \mathbf{x})\mathbf{x}$

Squaring both the sides we get:

$$\begin{aligned} \|\mathbf{w}_{t+1}\|^2 &= \|\mathbf{w}_t - (\mathbf{w}_t \cdot \mathbf{x})\mathbf{x}\|^2 \\ &= \|\mathbf{w}_t - (\mathbf{w}_t \cdot \mathbf{x})\mathbf{x}\|^T \|\mathbf{w}_t - (\mathbf{w}_t \cdot \mathbf{x})\mathbf{x}\| \end{aligned}$$

Expanding the above we get:

$$\begin{aligned} \|\mathbf{w}_{t+1}\|^2 &= \|\mathbf{w}_t\|^2 - 2\|\mathbf{w}_t(\mathbf{w}_t \cdot \mathbf{x})\mathbf{x}\|^2 + \|(\mathbf{w}_t \cdot \mathbf{x})\mathbf{x}\|^2 \\ \|\mathbf{w}_{t+1}\|^2 &= \|\mathbf{w}_t\|^2 - |\mathbf{w}_t \cdot \mathbf{x}|^2; \text{ since we have } \|\mathbf{x}\|_2 = 1 \\ \frac{\|\mathbf{w}_{t+1}\|^2}{\|\mathbf{w}_t\|^2} &= 1 - \frac{|\mathbf{w}_t \cdot \mathbf{x}|^2}{\|\mathbf{w}_t\|^2} \leq 1 - \delta^2 \end{aligned}$$

After  $m_t$  mistakes the update to  $\mathbf{w}$  changes to:

$$\|\mathbf{w}_{t+1}\|^2 \leq (1 - \delta^2)^{m_t} \rightarrow \text{Proof 1}$$

The convergence lower bound proof:

$$\mathbf{w}^{init} = 1$$

$$\mathbf{w}_{t+1} = \mathbf{w}^{init} - (\mathbf{w} \cdot \mathbf{x})\mathbf{x}$$

$$\mathbf{w}_{t+1} - \mathbf{w}^{init} = -(\mathbf{w} \cdot \mathbf{x})\mathbf{x}$$

After  $m_t$  mistakes we have:

$$\begin{aligned} \mathbf{w}_{t+1} - \mathbf{w}^{init} &= \sum_{i=0}^{m_t} -(\mathbf{w} \cdot \mathbf{x})\mathbf{x} \\ \mathbf{w}_{t+1} - \mathbf{w}^{init} &= -\sum_{i=0}^{m_t} (\mathbf{w} \cdot \mathbf{x})\mathbf{x} \end{aligned}$$

Taking the dot product of  $\mathbf{v}$  on both sides:

$$\begin{aligned} \mathbf{w}_{t+1} \cdot \mathbf{v} - \mathbf{w}^{init} \cdot \mathbf{v} &= -\sum_{i=0}^{m_t} (\mathbf{w} \cdot \mathbf{x})\mathbf{x} \cdot \mathbf{v} \\ \mathbf{w}_{t+1} \cdot \mathbf{v} &= -\sum_{i=0}^{m_t} (\mathbf{w} \cdot \mathbf{x})\mathbf{x} \cdot \mathbf{v} + \mathbf{w}^{init} \cdot \mathbf{v} \end{aligned}$$

We have  $\mathbf{w}^{init} \cdot \mathbf{v} \geq \gamma$  and  $\|\mathbf{v}\|_2 = 1$

$$\|\mathbf{w}_{t+1}\| \geq \gamma \rightarrow \text{Proof 2}$$

Combining both Proof 1 and Proof 2, we have

$$(1 - \delta^2)^{m_t} \geq \gamma^2$$

Taking  $\ln$  on both sides:

$$m_t \ln(1 - \delta^2) \geq 2 \ln \gamma$$

We know that  $-\ln(1 - x) \geq x$  and replacing it the above equation we get:

$$-m_t \delta^2 \geq 2 \ln(\gamma)$$

$$-m_t \geq 2 \ln(\gamma) / \delta^2$$

Simplifying, we get  $m_t \leq 2 / \delta^2 * \ln(1/\gamma)$

## 2 Programming (60 pts)

In this exercise, you will derive, implement back-propagation for a simple neural network, and compare your output with some standard library's output. Consider the following 3-layer neural network.

$$\mathbf{y} = \mathbf{f}(\mathbf{x}) = \mathbf{g}(\mathbf{W}_2 \sigma(\mathbf{W}_1 \mathbf{x}))$$

Suppose  $\mathbf{x} \in \mathbb{R}^d$ ,  $\mathbf{W}_1 \in \mathbb{R}^{d_1 \times d}$  and  $\mathbf{W}_2 \in \mathbb{R}^{k \times d_1}$  i.e.,  $f: \mathbb{R}^d \mapsto \mathbb{R}^k$ . Let  $\sigma(\mathbf{z}) = [\sigma(z_1), \dots, \sigma(z_n)]$  for any  $\mathbf{z} \in \mathbb{R}^n$  where  $\sigma(t) = \frac{1}{1 + \exp(-t)}$  is the sigmoid (logistic) activation function and  $\mathbf{g}(\mathbf{z})_i = \frac{\exp(\mathbf{z}_i)}{\sum_{i=1}^k \exp(\mathbf{z}_i)}$  is the softmax function. Suppose that the true pair is  $(\mathbf{x}, \mathbf{y})$  where  $\mathbf{y} \in \{0, 1\}^k$  with exactly one of the entries equal to 1 and you are working with the cross-entropy loss function given below,

$$L(\mathbf{x}, \mathbf{y}) = - \sum_{i=1}^k \mathbf{y} \log(\hat{\mathbf{y}}) .$$

1. Derive backpropagation updates for the above neural network. (10 pts)

We are given cross entropy loss function is given by  $L(\mathbf{x}, \mathbf{y}) = - \sum_{i=1}^k \mathbf{y} \log(\hat{\mathbf{y}}) .$

Some of the notations can be defined as:

$$\begin{aligned} \mathbf{z}^{(1)} &= \mathbf{W}^{(1)} \mathbf{x} \\ \mathbf{a}^{(1)} &= \sigma(\mathbf{z}^{(1)}) \\ \mathbf{z}^{(2)} &= \mathbf{W}^{(2)} \mathbf{a}^{(1)} \\ \mathbf{a}^{(2)} &= \hat{\mathbf{y}} = \text{softmax}(\mathbf{z}^{(2)}) \end{aligned}$$

To derive the back-propagation update, we need to calculate the partial derivative of the loss function with respect to the  $W^{(2)}$

By chain rule we can express the partial derivative as:

$$\frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial W^{(2)}}$$

$$\frac{\partial z^{(2)}}{\partial W^{(2)}} = a^{(1)}$$

The softmax function is given by:

$$\text{softmax}(z^{(i)}) = \frac{\exp(z_i)}{\sum_{i=1}^k \exp(z_i)}$$

The derivative of the softmax is expressed as below: (The notes in the lecture were referred to arrive at the derivative):

$$\frac{\partial \text{softmax}(z^{(i)})}{\partial z^{(i)}} = \frac{(\sum e^{(z_i)})(e^{(z_i)}) - (e^{(z_i)})^2}{(\sum e^{(z_i)})^2} = s(z_i)(1 - s(z_i))$$

$$\frac{\partial \text{softmax}(z^{(i)})}{\partial z^{(j)}} = \frac{0 - (e^{(z_i)})(e^{(z_j)})}{(\sum e^{(z_i)})^2} = -s(z_i)s(z_j)$$

The partial derivative of the L with respect to  $z^{(2)}$

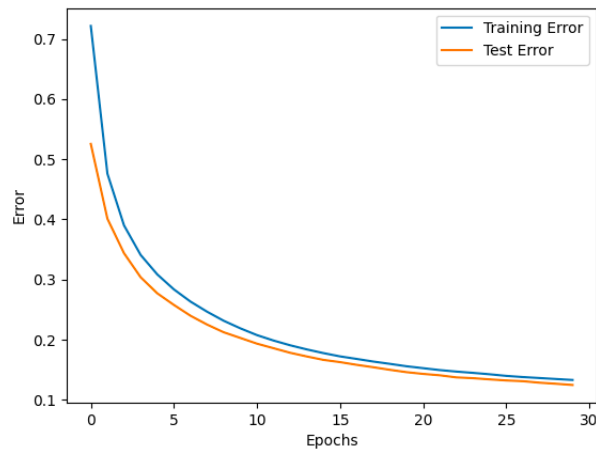
$$\begin{aligned} \frac{\partial L}{\partial z^{(2)}} &= - \sum \frac{y^{(i)}}{\hat{y}^i} \frac{\partial \hat{y}^{(i)}}{\partial z^{(2)}} \\ &= - \left[ \frac{y^{(i)}}{s(z_i)} s(z_i)(1 - s(z_i)) + \frac{\sum_{i \neq j} y^{(i)}}{s(z_i)} - s(z_i)s(z_j) \right] \\ &= -(y^{(i)}(1 - s(z_i)) + \sum_{i \neq j} -y^{(i)}s(z_j)) \\ &= -y^{(i)} + y^{(i)}s(z_i) + \sum_{i \neq j} -y^{(i)}s(z_j) \\ &= -y^{(i)} + \sum_{\forall i} -y^{(i)}s(z_i) \\ &= -y^{(i)} + s(z_i) \sum_{\forall i} -y^{(i)} \\ &= s(z_i) - y^{(i)} \\ &\Rightarrow \frac{\partial L}{\partial z^{(2)}} = a^{(2)} - y \\ &\Rightarrow \frac{\partial L}{\partial w^{(2)}} = (a^{(2)} - y)(a^{(1)}) \\ \frac{\partial L}{\partial W^{(1)}} &= \frac{\partial L}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial a^{(1)}} \frac{\partial a^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial W^{(1)}} \\ \frac{\partial a^{(1)}}{\partial z^{(1)}} &= \sigma(z^{(1)})(1 - \sigma(z^{(1)})) ; \frac{\partial z^{(2)}}{\partial a^{(1)}} = W^{(2)} ; \frac{\partial z^{(1)}}{\partial W^{(1)}} = x \\ &\Rightarrow \frac{\partial L}{\partial W^{(1)}} = (a^{(2)} - y)(a^{(1)})(\sigma(z^{(1)})(1 - \sigma(z^{(1)}))(x)) \end{aligned}$$

2. Implement it in `numpy` or `pytorch` using basic linear algebra operations. (e.g. You are not allowed to use `auto-grad`, built-in optimizer, model, etc. in this step. You can use library functions for data loading, processing, etc.). Evaluate your implementation on MNIST dataset, report test error, and learning curve. (25 pts)

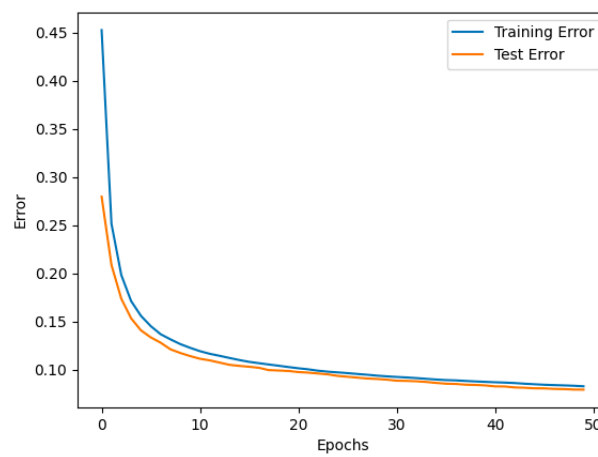
d	Learning Rate	Batch Size	Epoch	Test Error
300	0.01	128	30	0.19977666666666666
300	0.001	32	50	0.105526
300	0.005	64	100	0.107887

Table 1: Please find the respective graphs of each row below in order.

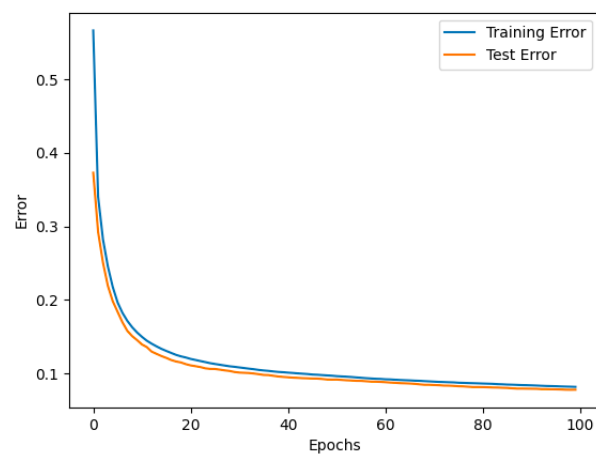
For  $d1 = 300$ :



$lr = 0.01$  batch-size = 128 epoch = 30



$lr = 0.001$  batch-size = 32 epoch = 50



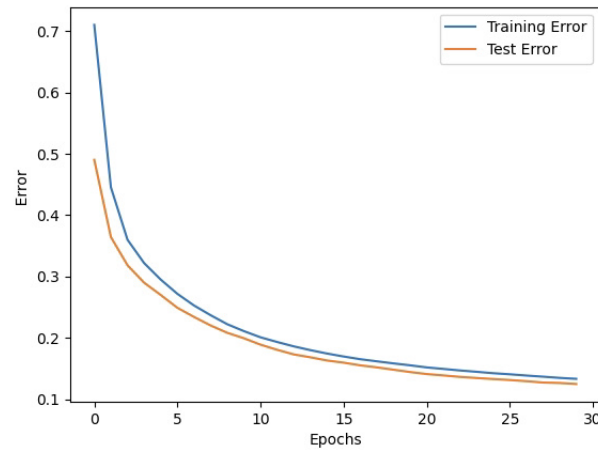
$lr = 0.005$  batch-size = 64 epoch = 100

- Implement the same network in `pytorch` (or any other framework). You can use all the features of the framework e.g. `auto-grad` etc. Evaluate it on MNIST dataset, report test error, and learning curve. (20 pts)

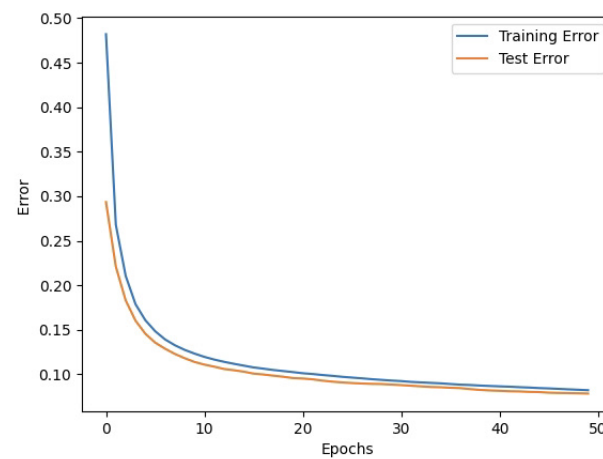
For  $d1 = 300$ :

d	Learning Rate	Batch Size	Epoch	Test Error
300	0.01	128	30	0.193653
300	0.001	32	50	0.10417
300	0.005	64	100	0.106197

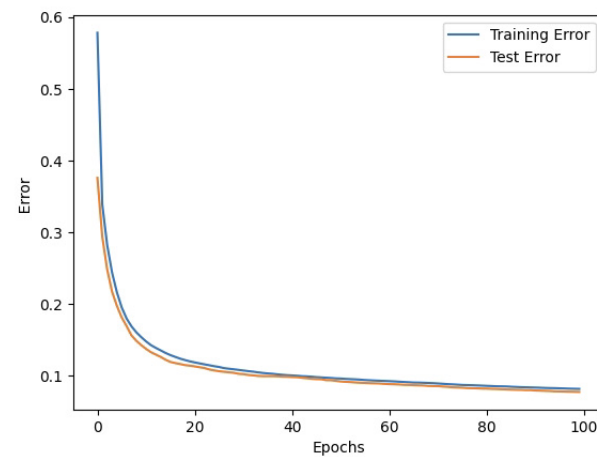
Table 2: Please find the respective graphs of each row below in order



lr = 0.01 batch-size = 128 epoch = 30



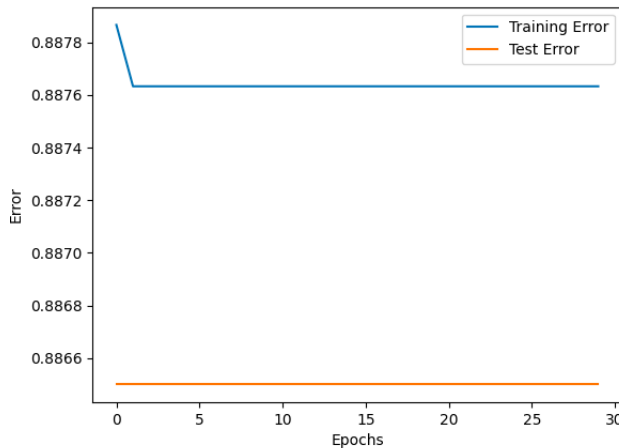
lr = 0.001 batch-size = 32 epoch = 50



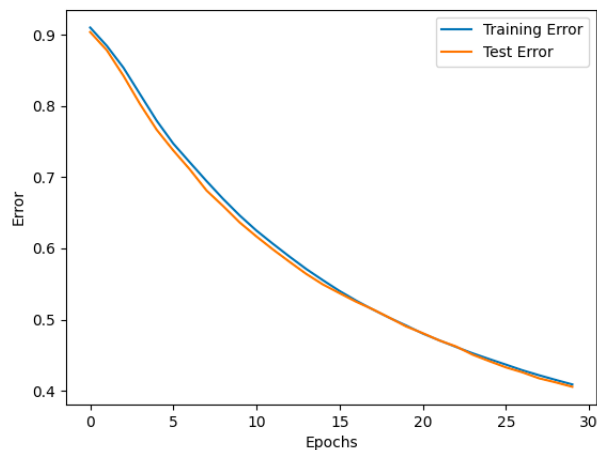
lr = 0.005 batch-size = 64 epoch = 100

4. Try different weight initializations a) all weights initialized to 0, and b) Initialize the weights randomly between -1 and 1. Report test error and learning curves for both. ( You can use either of the implementations) (5 pts)

a) Test error: 0.8865



b) Test error: 0.5833



You should play with different hyper-parameters like learning rate, batch size, etc. For Questions 2-3; you should report the answers for at least 3 different sets of hyperparameters. You should mention the values of those along with the results. Use  $d_1 = 300, d_2 = 200$ . For optimization use SGD (Stochastic gradient descent) without momentum, with some batch-size say 32, 64 etc. MNIST can be obtained from here (<https://pytorch.org/vision/stable/datasets.html>)

Check Piazza for a sample PyTorch project

To load the dataset:

```
import torch
import torch.nn.functional as F
from torch.utils.data import TensorDataset, DataLoader
import torchvision

mnist_data_train = torchvision.datasets.MNIST('path/to/download',
                                             train=True, download=True)
data_loader = torch.utils.data.DataLoader(mnist_data,
                                           batch_size=,
                                           shuffle=True)

mnist_data_test = torchvision.datasets.MNIST('path/to/download',
                                             train=False, download=True)
```