

UPPSALA UNIVERSITY

DEEP LEARNING

---

# Hand-in assignment (1)

---

Tong You

April 19, 2021



# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Softmax regression</b>	<b>3</b>
2.1	Results . . . . .	7
<b>3</b>	<b>Softmax regression with mini-batch SGD</b>	<b>11</b>
3.1	Results . . . . .	11
<b>4</b>	<b>Full Neural Network</b>	<b>15</b>
4.1	Results . . . . .	16
4.2	Discussion . . . . .	17
<b>5</b>	<b>Appendix</b>	<b>18</b>
5.1	Loading function . . . . .	18
5.2	Softmax regression . . . . .	19
5.3	Softmax regression with mini-batching . . . . .	24
5.4	Full Neural Network . . . . .	32

# 1 Introduction

In this document we will be highlighting the different steps that go into building a classifier for hand-written digits from scratch. We will be focusing on the MNIST dataset. This is a well-known data set that has 60000 training samples and 10000 testing samples. The samples are digits from 0 until 9. There are therefore 10 classes. We note that the images (28 x 28 pixels) have been flattened to a vector with 784 entries.



Figure 1: 16 examples of each class from the MNIST data set. Image retrieved from <https://commons.wikimedia.org/wiki/File:MnistExamples.png>

Compared to the previous assignment, this is not a binary classification task anymore and we need a different approach to the problem. To this end, we will be implementing a feedforward neural network to solve a multiclass classification problem. We will be dividing the main problem into three parts. First, a softmax output needs to be added. Second, the regular gradient descent method will be replaced with mini-batch training (SGD). And finally, the full neural network will be implemented with multiple layers. Important intermediate results will be shown in their respective sections. To make development understandable, the Python code will be implemented in 3 separate scripts. Development will be done in a [Github repository](#). All

Python code will also be attached in the Appendix.

## 2 Softmax regression

In order to do multiclass classification on the MNIST dataset we will need to update the linear regression model with a softmax output. Before we do so, we will need to get our dimensions in order. We will be using row-column notation and our first dimensions will be reserved for the samples. Our training data set has dimensions 60000 x 784 and our test data set has dimensions 10000 x 784. Our training labels have dimensions 60000 x 10 and test data set has dimensions 10000 x 10, where we are using one-hot encoding: the  $m$ -th column entry is "1" if it corresponds to a certain digit. Therefore we have the following indices:

$$i = 1, 2, \dots, n \text{ (number of samples)} \quad (1)$$

$$j = 1, 2, \dots, p \text{ (number of pixels)} \quad (2)$$

$$m = 1, 2, \dots, M \text{ (number of classes)} \quad (3)$$

Note that the number of samples can refer to the training and testing data. In the Python code, the distinction is made clearly. Our weight matrix will have dimensions  $M \times p$ , and the bias term will have dimensions  $1 \times M$ . Therefore, our linear model for the log odds will be:

$$z_{im} = \sum_{j=1}^p w_{mj} x_{ij} + b_m \quad (4)$$

Since we don't want to loop over array indices we will vectorize the equation as follows:

$$\mathbf{z}_i = \mathbf{W} \mathbf{x}_i + \mathbf{b}_m \quad (5)$$

In the Python implementation we will be transposing the weight matrix in order to get the correct dimensions for the subsequent matrix multiplication (see the [Appendix](#) for implementations per task). We will be using gradient descent to update our weights and biases. Before we can implement the Python code we need to derive the update equations for gradient descent. Compared to the pre-course assignment, we will be using the cross-entropy as our loss function. The loss:

$$L_i = - \sum_{m=1}^M \tilde{y}_{im} \ln(p_{im}) \quad (6)$$

In the previous equation,  $\tilde{y}_{im}$  are the one-hot encoded class labels. This one-hot encoding has already been done during the loading in the provided [Loading function](#). The one-hot encoding is defined as:

$$\tilde{y}_{im} = \begin{cases} 1, & \text{if } y_i = m. \\ 0, & \text{if } y_i \neq m. \end{cases} \quad (7)$$

The  $p_{im}$  are the probabilities that the  $i$ -th data point belongs to the  $m$ -th class (or digit). It is also called the softmax function. It is written as:

$$p_{im} = \frac{e^{z_{im}}}{\sum_{l=1}^M e^{z_{il}}} \quad (8)$$

The cost can be determined as:

$$J = \frac{1}{n} \sum_{i=1}^n L_i \quad (9)$$

Or, with the loss function filled in:

$$J = -\frac{1}{n} \sum_{i=1}^n \sum_{m=1}^M \tilde{y}_{im} \ln(p_{im}) \quad (10)$$

Before we extend these equations and derive the update rules for the weights and biases, we need to remark two things. The naive implementation of the softmax function will lead to numerical overflow/underflow if we have large/small  $z_{im}$  respectively. To counteract numerical overflow, we will normalize the  $z_{im}$  by subtracting the maximum and to counteract numerical underflow, we will be calculating the softmax function and the loss in one go.

We want the derivatives for the cost function with respect to the  $b_m$  and  $w_{mj}$ :

$$\frac{dJ}{db_m} = \sum_{i=1}^n \frac{dJ}{dz_{im}} \frac{dz_{im}}{db_m} \quad (11)$$

$$\frac{dJ}{dw_{mj}} = \sum_{i=1}^n \frac{dJ}{dz_{im}} \frac{dz_{im}}{dw_{mj}} \quad (12)$$

The loss function with the softmax function filled in:

$$L_i = \sum_{m=1}^M \left( \tilde{y}_{im} \ln \left( \sum_{l=1}^M e^{z_{il}} \right) - \tilde{y}_{im} z_{im} \right) \quad (13)$$

We also note that  $\frac{dJ}{dz_{im}}$  can be written as:

$$\frac{dJ}{dz_{im}} = \frac{1}{n} \frac{dL_i}{dz_{im}} = \frac{1}{n} \frac{dL_i}{dp_{im}} \frac{dp_{im}}{dz_{im}} \quad (14)$$

Before going any further, note that we know what the dimensionality of the derivatives will be simply by inspecting the indices underneath the symbols. This will also be helpful when implementing the softmax gradient descent itself. We will do the simple derivatives first before deriving the more complicated ones.

From equation 4, we can see that:

$$\frac{dz_{im}}{db_m} = 1 \quad (15)$$

Again from equation 4, we also see that:

$$\frac{dz_{im}}{dw_{mj}} = x_{ij} \quad (16)$$

Next, using equation 6 and the derivative of the natural logarithm:

$$\frac{dL_i}{dp_{im}} = -\frac{\tilde{y}_{im}}{p_{im}} \quad (17)$$

The most complicated derivative is the derivative of the softmax function with respect to the  $z_{im}$ . To reduce the amount of equations, we will not list all steps of the derivation. Note however that we will be using equation 8 and the [quotient rule](#) for derivatives. After some algebra we get:

$$\frac{dp_{im}}{dz_{im}} = \frac{\sum_{l=1}^M e^{z_{il}} e^{z_{im}} - e^{z_{im}} e^{z_{im}}}{\left( \sum_{l=1}^M e^{z_{il}} \right)^2} \quad (18)$$

This is simplified further into:

$$\frac{dp_{im}}{dz_{im}} = \frac{e^{z_{im}}}{\sum_{l=1}^M e^{z_{il}}} \frac{\sum_{l=1}^M e^{z_{il}} - e^{z_{im}}}{\sum_{l=1}^M e^{z_{il}}} \quad (19)$$

And further into:

$$\frac{dp_{im}}{dz_{im}} = p_{im} \left( 1 - \frac{e^{z_{im}}}{\sum_{l=1}^M e^{z_{il}}} \right) \quad (20)$$

Finally, we get:

$$\frac{dp_{im}}{dz_{im}} = p_{im} (1 - p_{im}) \quad (21)$$

Now combining equations 14, 17, and 21 the derivative for the cost function becomes:

$$\frac{dJ}{dz_{im}} = -\frac{1}{n} \frac{\tilde{y}_{im}}{p_{im}} p_{im} (1 - p_{im}) = \frac{\tilde{y}_{im} p_{im} - \tilde{y}_{im}}{n} \quad (22)$$

Using equations 11, 15, 16, and 22 we get:

$$\frac{dJ}{db_m} = \sum_{i=1}^n \frac{\tilde{y}_{im} p_{im} - \tilde{y}_{im}}{n} \quad (23)$$

And using equations 12, 15, 16, and 22:

$$\frac{dJ}{dw_{mj}} = \sum_{i=1}^n \frac{(\tilde{y}_{im} p_{im} - \tilde{y}_{im}) x_{ij}}{n} \quad (24)$$

Now we are ready to give the update equations for the weights and biases. Using  $\leftarrow$  to denote a variable update and  $\gamma$  the learning rate, we get the gradient descent update rule for the biases  $b_m$ :

$$b_m \leftarrow b_m - \gamma \frac{dJ}{db_m} \quad (25)$$

The update rule for weight  $w_{mj}$

$$w_{mj} \leftarrow w_{mj} - \gamma \frac{dJ}{dw_{mj}} \quad (26)$$

With the derivatives fully written out:

$$b_m \leftarrow b_m - \frac{\gamma}{n} \sum_{i=1}^n (\tilde{y}_{im} p_{im} - \tilde{y}_{im}) \quad (27)$$

And for  $w_{mj}$

$$w_{mj} \leftarrow w_{mj} - \frac{\gamma}{n} \sum_{i=1}^n (\tilde{y}_{im} p_{im} - \tilde{y}_{im}) x_{ij} \quad (28)$$

In the next section we will show the results for the softmax classifier on MNIST.

## 2.1 Results

Previously, we have derived the gradient descent update equations for the weights and biases. This is very important, but how we choose to implement the equations in Python is also important. The full code to generate the plots can be found in [Softmax regression](#). We will highlight the most important part and discuss some brief implementation details before showing the results. In the main optimization loop, the weights and biases are updated as follows:

```
1      z_im = xtrain @ w_mj.T + b_m
2
3      y_im = ytrain
4
5      z_im_norm = z_im - np.max(z_im, axis = 1, keepdims =
6      ↪ True)
7
8      p_im = np.exp(z_im_norm) / np.sum(np.exp(z_im_norm),
9      ↪ axis = 1,
10     keepdims = True)
11
12     dJdzim = (1/n_train) * (y_im * p_im - y_im)
13
14     dJdbm = np.sum(dJdzim, axis = 0)
15     dJdwmj = dJdzim.T @ xtrain
16
17     b_m = b_m - lr * dJdbm
18     w_mj = w_mj - lr * dJdwmj
```

The choice has been made to try to stay as close as possible to the mathematical equations we have used previously. The normalization of the  $z_{im}$  is explicit in the Python code. The numerical-stable version of the loss was also implemented:

```
1      L_i = np.sum(y_im * np.log(np.sum(np.exp(z_im_norm), axis
2      ↪ = 1,
3      keepdims = True))) - y_im * z_im_norm, axis = 1)
```

The testing loss is implemented in a similar way, with the only difference



being that we use testing data to evaluate (not training data!!!) our model. Before we show the results, we will briefly talk about parameter initialization. For the biases  $b_m$  it was sufficient to initialize all 10 values to zero. For the weights  $w_{mj}$  we will initialize the entire matrix with normally-distributed values with a mean of 0 and a standard deviation of 0.01.

An important value to tune is the learning rate  $\gamma$ . Too high a value means there will be no convergence, whereas a low value means that the learning will be slow. We will show results for  $\gamma = 0.01$ ,  $\gamma = 0.02$ , and  $\gamma = 0.03$ . These learning rates resulted in accuracies close to 90%. We ran the training for a maximum of 3000 iterations. In Figures 2 and 3 the weights and accuracy/cost plots are shown for the three different learning rates. From the weight matrices themselves no clear distinction can be made between the choice of the three different learning rates. It is only when we look at the cost/accuracy plots that a clear distinction can be made between the choice of the various learning rates. First of all, we note that for all the learning rates the training accuracy is always slightly below the testing accuracy. In the cost plots, this is reverse where the training cost is higher than the testing cost. This is consistent with each other. We also see that the testing accuracy does increase going from a learning rate of 0.01 to 0.03. However, from testing (results not shown) it was difficult to reach a testing accuracy of 90%. Finally, we note that the cost increases after arriving at a minimum value for both the training and testing data when learning rates of 0.02 and 0.03 are used. Note that we don't train on the testing data (can also be seen in the code). We only evaluate the testing data using the weights and bias that is being trained on the training data. Since we have more training than testing data it is also no unusual for the testing accuracy to be slightly higher than the training accuracy.

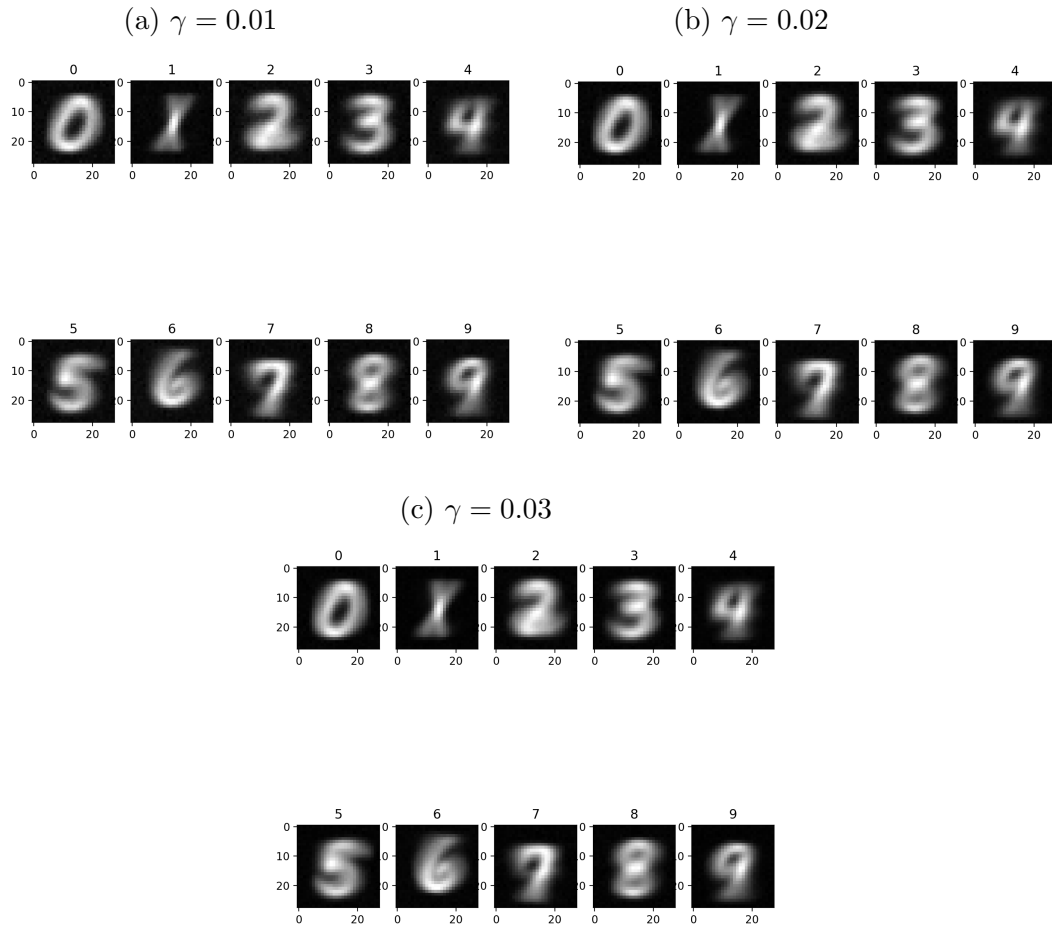
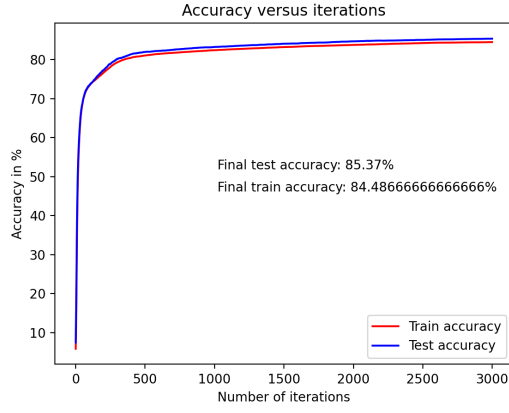
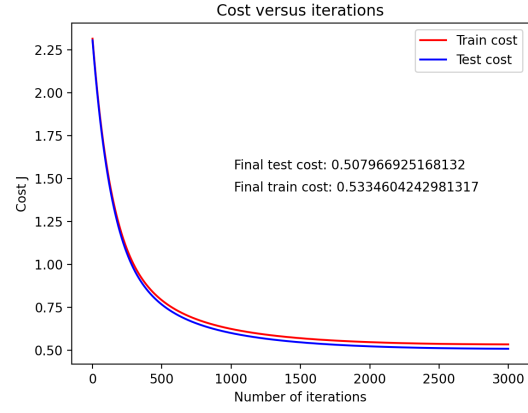


Figure 2: All 10 weight matrices for three different learning rates. The learning rates are indicated in the subcaptions. No post-processing was performed on the weights obtained. The weight matrices were only reshaped from 784 to 28 x 28 images.

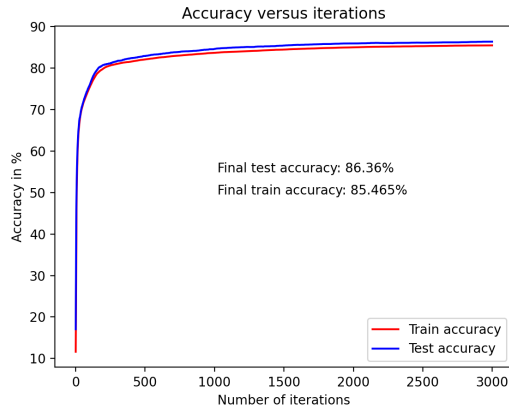
(a)  $\gamma = 0.01$



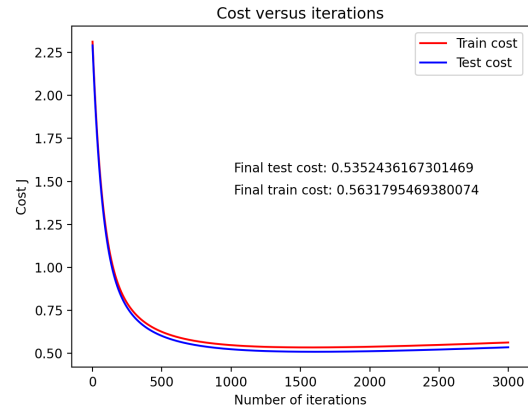
(b)  $\gamma = 0.01$



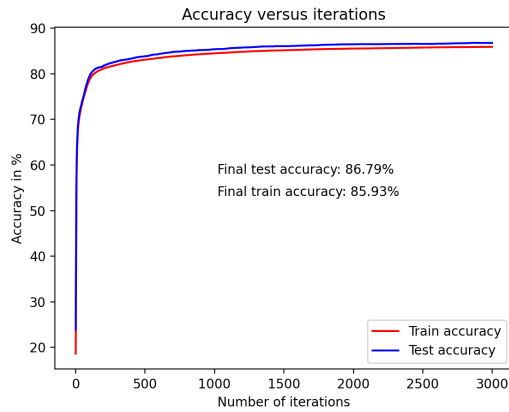
(c)  $\gamma = 0.02$



(d)  $\gamma = 0.02$



(e)  $\gamma = 0.03$



(f)  $\gamma = 0.03$

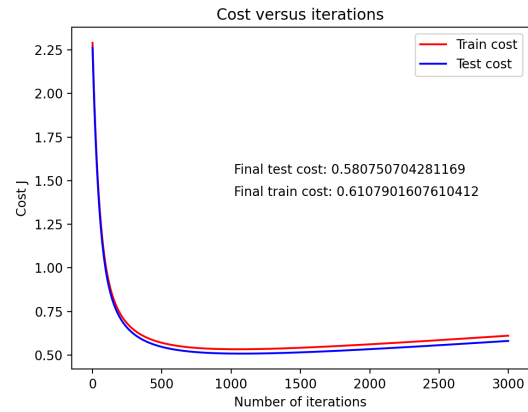


Figure 3: Accuracy and cost for the corresponding weights in Figure 2.

### 3 Softmax regression with mini-batch SGD

There are no major differences compared to the softmax-only implementation. We will use the same update equations for the weights and biases. There will be minor changes to the equations and Python code used previously. The first one will be in the gradient calculations. Instead of calculating gradients over the entire training set, we will calculate gradients over a smaller batch (or subset) of the training data. This means that the gradient updates will be noisier, but our training should be significantly faster! We need to make sure that the gradients are properly scaled with the batch size. This means that equations 27 and 28 need to be updated to:

$$b_m \leftarrow b_m - \frac{\gamma}{n_{batch}} \sum_{i \in batch} (\tilde{y}_{im} p_{im} - \tilde{y}_{im}) \quad (29)$$

$$w_{mj} \leftarrow w_{mj} - \frac{\gamma}{n_{batch}} \sum_{i \in batch} (\tilde{y}_{im} p_{im} - \tilde{y}_{im}) x_{ij} \quad (30)$$

The final thing to keep in mind is that we need to decrease the learning rate when using mini-batch SGD. The learning rate needs to decrease over time so that convergence is guaranteed. We will be using a linear schedule:

$$\gamma_{it} = (1 - \frac{it}{\tau})\gamma_0 + \frac{it}{\tau}\gamma_\tau \quad (31)$$

Here  $\gamma_0$  is the initial learning rate that we start with,  $\gamma_\tau$  is the learning rate that we use after  $\tau$  iterations ( $it$  is iteration counter). The values for the learning rates and important parts of the Python implementation will be given below.

#### 3.1 Results

We will highlight major differences between the softmax-only and mini-batch softmax implementations before showing the results. One difference is that we have two loops instead of one loop. In the current implementation we have a while-loop over the epochs, while the optimization loop is a for-loop. Each epoch, we will loop  $\frac{n_{train}}{n_{batch}}$  times over the training data (an iteration is a gradient update using a mini-batch). The most important code is the mini-batching code. We shuffle our training data each epoch randomly. We also need to shuffle our training labels. The shuffling is done as follows:

```

1      ### Shuffling indices
2      ind = np.arange(n_train)
3      np.random.shuffle(ind) # shuffle indices
4
5      ### Shuffling training data and labels
6      xt = xtrain[ind, :]
7      yt = ytrain[ind, :]
8
9      ytrue_shuff = np.argmax(yt, axis = 1)

```

The next important implementation detail is the mini-batching done each iteration. This is done as follows:

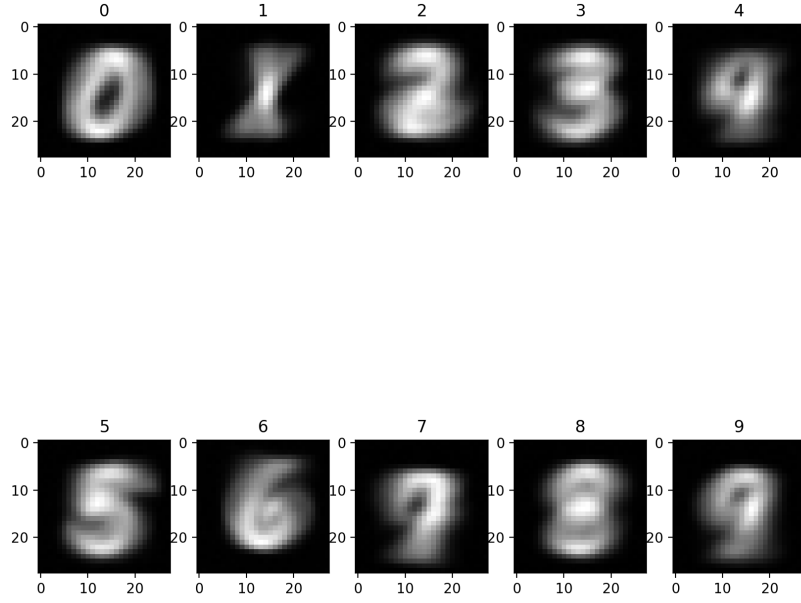
```

1      mini_batch = np.random.randint(0, n_train, size = nb)
      ↪ # batch indices

```

The rest of the Python implementation is similar, where we need to be careful to only calculate the  $z_{im}$  and etc. for the current mini-batch. For more details, please refer to the full [Softmax regression with mini-batching](#) code. Since a learning rate of 0.01 was found to be desirable we will be using this value for  $\gamma_0$ . For  $\gamma_\tau$  we will use a rule-of-thumb to set it to 1% of  $\gamma_0$ . The value of  $\tau$  is slightly more difficult to set. The choice has been made to keep the learning rate constant for the final 5 iterations. For the mini-batch size of 1500 that was settled on, each epoch lasts 40 iterations. We will train for 300 epochs and calculate/plot the training and testing cost/accuracy every 5 iterations. In order to compare differences, we will look at a smaller batch-size of 150. For a batch-size equal to the training data size we expect similar results to the softmax-only implementation.

(a)  $\gamma = 0.01$  and batch-size of 150



(b)  $\gamma = 0.01$  and batch-size of 1500

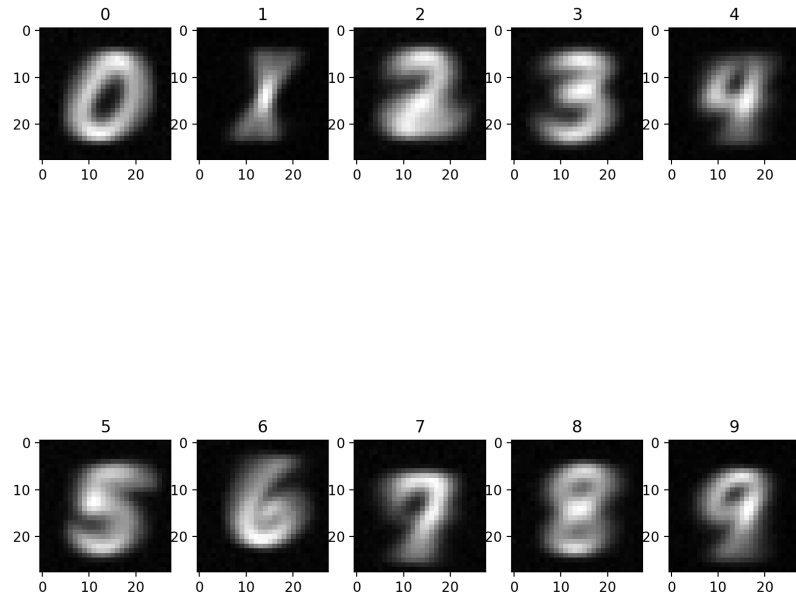
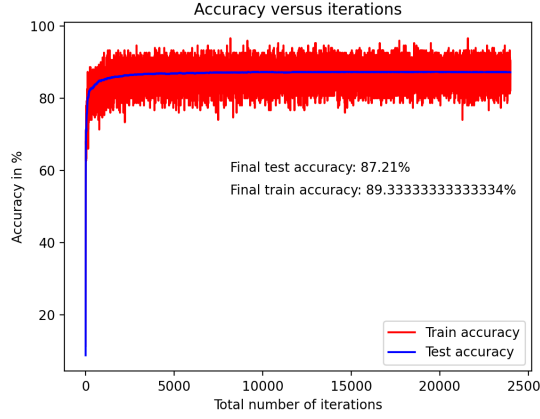
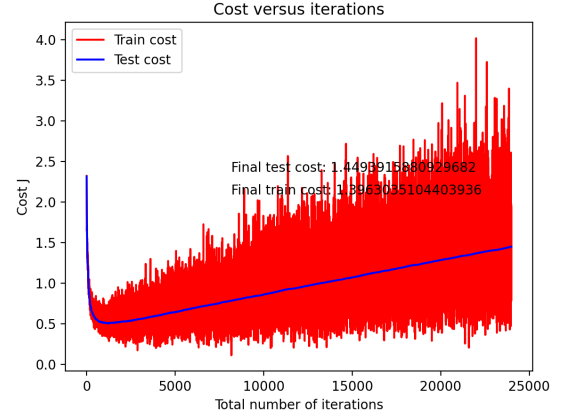


Figure 4: All 10 weight matrices for two different learning rates. The learning rate and batch-size are indicated in the subcaptions. No post-processing was performed on the weights obtained. The weight matrices were only reshaped from 784 to 28 x 28 images.

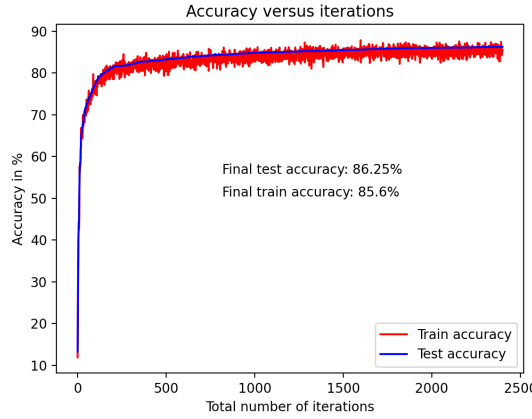
(a)  $\gamma = 0.01$  and batch-size of 150



(b)  $\gamma = 0.01$  and batch-size of 150



(c)  $\gamma = 0.01$  and batch-size of 1500



(d)  $\gamma = 0.01$  and batch-size of 1500

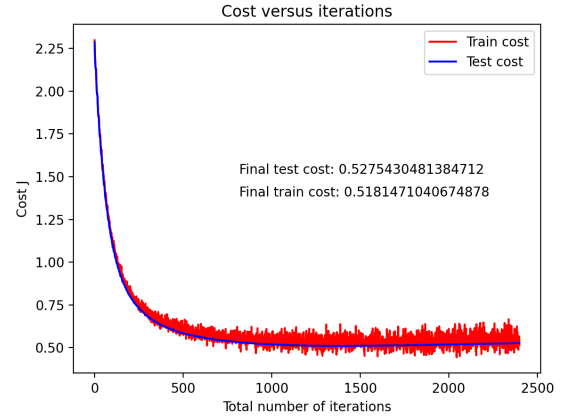


Figure 5: Accuracy and cost for the weights in Figure 4.

Compared to the softmax-only implementation we can see with the same parameters except batch-size there is a large difference between using a batch-size of 150 and 1500. In Figure 4 we can see that for most of the 10 classes, the weights are blurrier for a batch-size of 150 versus one of 1500. However,

for the digit "4" the differences are not as clear for a batch size of 150 and 1500. There is a much clearer difference between the two batch sizes when we look at Figure 5. In this we can see that both the training cost and accuracy fluctuates much more for a batch size of 150 compared to 1500. This is to be expected since we know that for a smaller batch-size the gradient updates are much noisier. This was confirmed by looking at the evolution of the gradient of the bias versus iterations. However, even though the training accuracy fluctuates a lot it does converge to a final value of about 89.3%. The final test accuracy for a batch-size of 150 was 87.21%. Compare this to the final training accuracy of 85.6% and final test accuracy of 86.25% for a batch-size of 1500. Although the test accuracy is higher than the softmax-only result, it is still not able to reach the 92% target. Maybe doing a parameter search using a function like [GridSearchCV](#) might result in optimal values for the mini-batch softmax implementation. Next, we will look at the full neural network implementation.

## 4 Full Neural Network

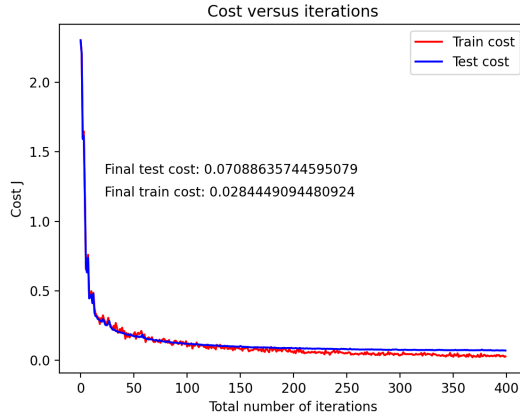
In this section we will implement an L-layer neural network (note that the previous two networks were one-layer networks). The process is similar to the softmax implementations. One of the major differences is that the mathematics is more involved. We need to obtain equations for the forward and backward propagation. During the forward propagation we just propagate the input sequentially through all layers and calculate the cost/loss. During the backward propagation we evaluate all the gradients from the last layer back to the input with respect to the hidden units. We will be using the results from the [Supervised Machine Learning](#) book in our Python implementation. We will follow backpropagation Algorithm 6.1 and use the derived equations in our implementation. This algorithm only gives us the cost function and the gradient(s) of the cost function. We will update the weights and biases using gradient descent (stochastic since we are using mini-batch) like previously.



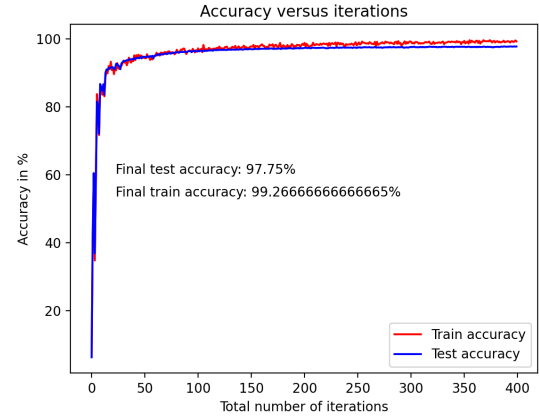
## 4.1 Results

Similar to the previous two sections, the full Python code can be found in the [Full Neural Network](#) Appendix section. The main optimization loop looks quite similar to the code for the previous two sections. The current implementation of the neural network only has code implemented for 2-layer and 4-layer neural networks. We will therefore make comparisons between the 2-layer and 4-layer results. We want to start by stating that an optimal network architecture for this task seems to be a 2-layer network. Also note that we will be showing results only for ReLU activations. We will also accumulate the cost/accuracy values every 5 iterations and use a batch-size of 1500 for both the 2-layer and the 4-layer network. The number of hidden units was also kept the same at 100 and 100,100, and 100 for the 2-layer and 4-layer network respectively. The 2-layer network seemed stable with various hidden units. The 4-layer network was much more difficult to configure, but the results that will be shown are the "optimal" ones. Before showing the results we also want to note that the learning rate was kept at a constant value for both the 2-layer and 4-layer network. Using a decaying learning rate we saw that both networks exhibited oscillations in the cost/accuracy. They did manage to converge to a "reasonable" final value, but these oscillations are clearly not desirable. The learning rate was set to a constant value for both networks. We ran both networks for 50 epochs. Note that for the various settings this means we have plots over 400 iterations. Both networks had their weight matrices initialized according to a normal distribution with mean 0 and standard deviation 0.01. The biases were all initialized to 0.

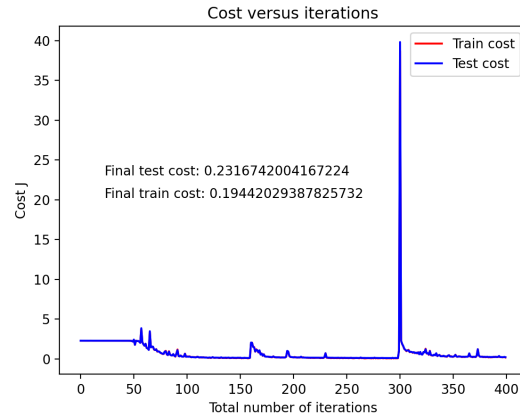
(a)  $\gamma = 0.8$  for 2-layer network



(b)  $\gamma = 0.8$  for 2-layer network



(c)  $\gamma = 0.8$  for 4-layer network



(d)  $\gamma = 0.8$  for 4-layer network

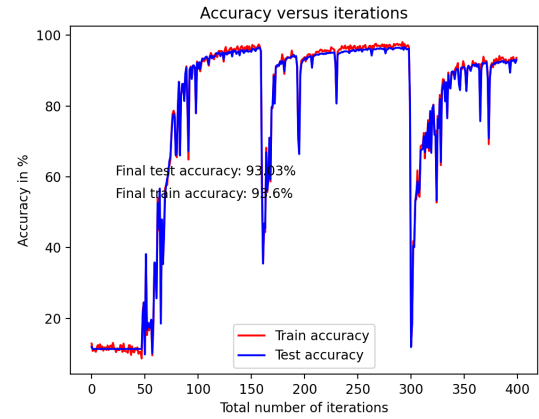


Figure 6: Accuracy and cost plots for 2-layer and 4-layer networks.

## 4.2 Discussion

One of the big differences between the 2-layer and 4-layer network is that for the same number of iterations we can see that the cost and accuracy for the 2-layer network converge nicely. The final testing accuracy for the 2-layer network is 97.75%, which is close to the target of 98%. The final

test accuracy for the 4-layer network is around 93%. However, note that the 4-layer network exhibits a much less stable cost/accuracy evolution. An increase at around 300 iterations occurs in both the training and testing cost, with a corresponding decrease in the training and testing accuracy. Despite this instability, the 4-layer network does manage to converge to a reasonable final accuracy. It is in fact much more than the softmax-only/mini-batch softmax implementations. The 4-layer network is also much more sensitive to initialization since running the 4-layer network with the same parameters resulted in wildly-varying final testing accuracies. Again, a parameter-search for the 4-layer network might result in a much more well-behaved network for this specific data set. A final thing of note is that the 2-layer network also exhibits slight fluctuation in the testing accuracy and testing cost. For the softmax-only and mini-batch softmax implementations this effect is not visible.

## 5 Appendix

### 5.1 Loading function

```
1 import numpy as np
2 import imageio
3 import glob
4
5 # Small modification: added absolute path of Test and Train
6 → folders
7
8 def load_mnist():
9     # Loads the MNIST dataset from png images
10
11     NUM_LABELS = 10
12     # create list of image objects
13     test_images = []
14     test_labels = []
15
16     for label in range(NUM_LABELS):
```

```

16     for image_path in glob.glob("/Users/tongyou/Desktop/
    ↪ Deep-learning-course/Assignment-1/MNIST/Test/" +
    ↪ str(label) + "/*.png"):
17         image = imageio.imread(image_path)
18         test_images.append(image)
19         letter = [0 for _ in range(0, NUM_LABELS)]
20         letter[label] = 1
21         test_labels.append(letter)
22
23     # create list of image objects
24     train_images = []
25     train_labels = []
26
27     for label in range(NUM_LABELS):
28         for image_path in glob.glob("/Users/tongyou/Desktop/
    ↪ Deep-learning-course/Assignment-1/MNIST/Train/" +
    ↪ str(label) + "/*.png"):
29             image = imageio.imread(image_path)
30             train_images.append(image)
31             letter = [0 for _ in range(0, NUM_LABELS)]
32             letter[label] = 1
33             train_labels.append(letter)
34
35     X_train= np.array(train_images).reshape(-1,784)/255.0
36     Y_train= np.array(train_labels)
37     X_test= np.array(test_images).reshape(-1,784)/255.0
38     Y_test= np.array(test_labels)
39
40     return X_train, Y_train, X_test, Y_test

```

## 5.2 Softmax regression

```

1 import numpy as np
2 from load_mnist import load_mnist
3 from matplotlib import pyplot as plt
4
5 x_train, y_train, x_test, y_test = load_mnist()

```

```

6
7 M = 10
8 n_train = x_train.shape[0] # number of training examples -
  ↳ 60000
9 p = x_train.shape[1] # number of input pixels - 784
  ↳ (flattened 28x28 image)
10
11 n_test = x_test.shape[0]
12
13 def softmax_gd(xtrain, ytrain, xtest, ytest, n_train, n_test,
  ↳ lr, maxit):
14
15     w_mj = np.random.normal(scale = 0.02, size = (M, p)) #
  ↳ weight matrix
16     b_m = np.zeros(shape = (1, M)) # offset vector
17     z_im = np.zeros(shape = (n_train, M)) # model in (n x M)
18     dJdbm = np.zeros(shape = (1, M))
19     dJdwmj = np.zeros(shape = (M, p))
20
21     J = np.zeros(shape = (maxit, 1)) # cost vector
22     acc_train = np.zeros(shape = (maxit, 1)) # classification
  ↳ accuracy
23
24     J_test = np.zeros(shape = (maxit, 1)) # cost vector
25     acc_test = np.zeros(shape = (maxit, 1)) # classification
  ↳ accuracy
26
27     it = 0
28
29     while it != maxit:
30
31         z_im = xtrain @ w_mj.T + b_m
32
33         y_im = ytrain
34
35         z_im_norm = z_im - np.max(z_im, axis = 1, keepdims =
  ↳ True)

```

```

36
37 p_im = np.exp(z_im_norm) / np.sum(np.exp(z_im_norm),
    ↪ axis = 1, keepdims = True)
38
39 dJdzim = (1/n_train) * (y_im * p_im - y_im)
40
41 dJdbm = np.sum(dJdzim, axis = 0)
42 dJdwmj = dJdzim.T @ xtrain
43
44 b_m = b_m - lr * dJdbm
45 w_mj = w_mj - lr * dJdwmj
46
47 # Cost and accuracy for training mini-batch
48 L_i = np.sum(y_im * np.log(np.sum(np.exp(z_im_norm),
    ↪ axis = 1, keepdims = True))) - y_im * z_im_norm,
    ↪ axis = 1)
49 J[it] = (1/n_train) * np.sum(L_i)
50
51 ypred_train = np.argmax(p_im, axis = 1)
52 ytrue_train = np.argmax(y_train, axis = 1)
53 acc_train[it] = np.array([1 for i in range(0,n_train)
    ↪ if ypred_train[i] == ytrue_train[i]]).sum()
54
55 # Cost and accuracy for testing data
56 z_test = xtest @ w_mj.T + b_m
57 z_test_norm = z_test - np.max(z_test, axis = 1,
    ↪ keepdims = True)
58 L_i_test = np.sum(ytest *
    ↪ np.log(np.sum(np.exp(z_test_norm), axis = 1,
    ↪ keepdims = True))) - ytest * z_test_norm, axis = 1)
59 J_test[it] = (1/n_test) * np.sum(L_i_test)
60
61 p_test = np.exp(z_test_norm) /
    ↪ np.sum(np.exp(z_test_norm), axis = 1, keepdims =
    ↪ True)
62
63 ypred_test = np.argmax(p_test, axis = 1)

```

```

64     ytrue_test = np.argmax(y_test, axis = 1)
65     acc_test[it] = np.array([1 for i in range(0,n_test) if
    ↪ ypred_test[i] == ytrue_test[i]]).sum()
66
67     it += 1
68     print("Iteration: (%s/%s)" % (it, maxit))
69
70     return J, acc_train * (1/n_train) * 100, J_test, acc_test
    ↪ * (1/n_test) * 100, it, w_mj, b_m , z_im, p_im
71
72 lr = 0.01
73 iters = 5000
74
75 J, acc_train, J_test, acc_test, it, wmj, bm, zim, pim =
    ↪ softmax_gd(x_train, y_train, x_test, y_test, n_train,
    ↪ n_test, lr, iters)
76
77 plt.figure(1)
78 plt_J, = plt.plot(J, 'r')
79 plt_J_test, = plt.plot(J_test, 'b')
80 plt.legend([plt_J, plt_J_test], ['Train cost', 'Test cost'])
81 plt.annotate("Final train cost: %s" % (float(J[-1])), xycoords
    ↪ = 'figure fraction',xy = (0.4,0.5))
82 plt.annotate("Final test cost: %s" % (float(J_test[-1])),
    ↪ xycoords = 'figure fraction',xy = (0.4,0.55))
83 plt.xlabel('Number of iterations')
84 plt.ylabel('Cost J')
85 plt.title('Cost versus iterations')
86 print("Final train cost: %s" % float(J[-1]))
87 print("Final test cost: %s" % float(J_test[-1]))
88
89 plt.figure(2)
90 plt_acc_train, = plt.plot(acc_train, 'r')
91 plt_acc_test, = plt.plot(acc_test, 'b')
92 plt.legend([plt_acc_train, plt_acc_test], ['Train accuracy',
    ↪ 'Test accuracy'])

```

```

93 plt.annotate("Final train accuracy: %s%" %
    ↳ (float(acc_train[-1])), xycoords = 'figure fraction',xy =
    ↳ (0.4,0.5))
94 plt.annotate("Final test accuracy: %s%" %
    ↳ (float(acc_test[-1])), xycoords = 'figure fraction',xy =
    ↳ (0.4,0.55))
95 plt.xlabel('Number of iterations')
96 plt.ylabel('Accuracy in %')
97 plt.title('Accuracy versus iterations')
98 print("Final train accuracy: %s%" % float(acc_train[-1]))
99 print("Final test accuracy: %s%" % float(acc_test[-1]))
100
101 figw, axw = plt.subplots(2, 5, figsize = (8,8))
102 axw[0,0].imshow(wmj[0, :].reshape(28,28), cmap = 'gray')
103 axw[0,0].set_title('0')
104 axw[0,1].imshow(wmj[1, :].reshape(28,28), cmap = 'gray')
105 axw[0,1].set_title('1')
106 axw[0,2].imshow(wmj[2, :].reshape(28,28), cmap = 'gray')
107 axw[0,2].set_title('2')
108 axw[0,3].imshow(wmj[3, :].reshape(28,28), cmap = 'gray')
109 axw[0,3].set_title('3')
110 axw[0,4].imshow(wmj[4, :].reshape(28,28), cmap = 'gray')
111 axw[0,4].set_title('4')
112 axw[1,0].imshow(wmj[5, :].reshape(28,28), cmap = 'gray')
113 axw[1,0].set_title('5')
114 axw[1,1].imshow(wmj[6, :].reshape(28,28), cmap = 'gray')
115 axw[1,1].set_title('6')
116 axw[1,2].imshow(wmj[7, :].reshape(28,28), cmap = 'gray')
117 axw[1,2].set_title('7')
118 axw[1,3].imshow(wmj[8, :].reshape(28,28), cmap = 'gray')
119 axw[1,3].set_title('8')
120 axw[1,4].imshow(wmj[9, :].reshape(28,28), cmap = 'gray')
121 axw[1,4].set_title('9')
122
123 plt.show()

```



### 5.3 Softmax regression with mini-batching

```
1 import numpy as np
2 from load_mnist import load_mnist
3 from matplotlib import pyplot as plt
4
5 x_train, y_train, x_test, y_test = load_mnist()
6
7 M = 10
8 n_train = x_train.shape[0] # number of training examples -
   ↳ 60000
9 p = x_train.shape[1] # number of input pixels - 784
   ↳ (flattened 28x28 image)
10
11 n_test = x_test.shape[0] # number of testing examples -
   ↳ 10000
12 ytrue_test = np.argmax(y_test, axis = 1)
13
14 def softmax_gd_minibatch(xtrain, ytrain, xtest, ytest, ep, nb,
   ↳ lr_init, tau, n_train, n_test, k):
15
16     w_mj = np.random.normal(scale = 0.01, size = (M, p)) #
   ↳ weight matrix
17     b_m = np.zeros(shape = (1, M))
18     z_im = np.zeros(shape = (n_train, M))
19     dJdbm = np.zeros(shape = (1, M))
20     dJdwmj = np.zeros(shape = (M, p))
21
22     J_train = np.zeros(shape = (ep, 1))
23     acc_train = np.zeros(shape = (ep, 1))
24
25     J_test = np.zeros(shape = (ep, 1))
26     acc_test = np.zeros(shape = (ep, 1))
27
28     e_p = 0 # epoch counter
29     lr0 = lr_init # initial learning rate
30     lrt = 0.01 * lr0 # final learning rate
31     t_tau = tau
```

```

32
33     tot_it = 0
34     it_k = 0
35
36     Jtrainiter = np.zeros(shape = (((n_train//nb)//(k)) *
37     ↪ ep), 1))
37     Jtestniter = np.zeros(shape = (((n_train//nb)//(k)) *
38     ↪ ep), 1))
38
39     acctrainiter = np.zeros(shape = (((n_train//nb)//(k)) *
40     ↪ ep), 1))
40     acctestiter = np.zeros(shape = (((n_train//nb)//(k)) *
41     ↪ ep), 1))
41
42     while e_p != ep:
43
44         ### Shuffling indices
45         ind = np.arange(n_train)
46         np.random.shuffle(ind) # shuffle indices
47
48         ### Shuffling training data and labels
49         xt = xtrain[ind, :]
50         yt = ytrain[ind, :]
51
52         ytrue_shuff = np.argmax(yt, axis = 1)
53
54         it = 0
55
56         Jtrainaccum = []
57         acctrainaccum = []
58
59         Jtestaccum = []
60         acctestaccum = []
61
62         while it != n_train//nb:
63
64             lr = (1 - (it/t_tau)) * lr0 + (it/t_tau) * lrt
65

```

```

66     mini_batch = np.random.randint(0, n_train, size =
    ↪ nb) # batch indices
67
68     z_im[mini_batch, :] = xt[mini_batch, :] @ w_mj.T
    ↪ + b_m
69
70     y_im = yt[mini_batch, :]
71
72     z_im_norm = z_im[mini_batch, :] -
    ↪ np.max(z_im[mini_batch, :], axis = 1, keepdims
    ↪ = True)
73
74     p_im = np.exp(z_im_norm) /
    ↪ np.sum(np.exp(z_im_norm), axis = 1, keepdims =
    ↪ True)
75
76     dJdzim = (1/nb) * (y_im * p_im - y_im)
77     dJdbm = np.sum(dJdzim, axis = 0)
78     dJdwmj = dJdzim.T @ xt[mini_batch, :]
79
80     b_m = b_m - lr * dJdbm
81     w_mj = w_mj - lr * dJdwmj
82
83     # Calculate the cost and accuracy every k-th
    ↪ iteration for averaging per epoch
84     if it % k == 0:
85         #Cost and accuracy for training data
86         L_i = np.sum(y_im *
    ↪ np.log(np.sum(np.exp(z_im_norm), axis = 1,
    ↪ keepdims = True)) - y_im * z_im_norm, axis
    ↪ = 1)
87         ypred = np.argmax(p_im, axis = 1)
88
89         Jtrainaccum.append((1/nb) * np.sum(L_i))
90         acctrainaccum.append((1/nb) * np.sum(ypred ==
    ↪ ytrue_shuff[mini_batch]))
91

```

```

92     # Cost and accuracy for testing data
93     z_test = xtest @ w_mj.T + b_m
94     z_test_norm = z_test - np.max(z_test, axis =
    ↪ 1, keepdims = True)
95     p_test = np.exp(z_test_norm) /
    ↪ np.sum(np.exp(z_test_norm), axis = 1,
    ↪ keepdims = True)
96     L_i_test = np.sum(y_test *
    ↪ np.log(np.sum(np.exp(z_test_norm), axis =
    ↪ 1, keepdims = True)) - y_test *
    ↪ z_test_norm, axis = 1)
97     ypred_test = np.argmax(p_test, axis = 1)
98
99     Jtestaccum.append((1/n_test) *
    ↪ np.sum(L_i_test))
100     acctestaccum.append((1/n_test) *
    ↪ np.sum(ypred_test == ytrue_test))
101
102     # Calculate the cost and accuracy every k-th
    ↪ iteration and accumulate over all epochs
103     if tot_it % k == 0:
104         #Cost and accuracy for training data
105         L_i = np.sum(y_im *
    ↪ np.log(np.sum(np.exp(z_im_norm), axis = 1,
    ↪ keepdims = True)) - y_im * z_im_norm, axis
    ↪ = 1)
106         ypred = np.argmax(p_im, axis = 1)
107         if nb == n_train:
108             Jtrainiter = []
109             acctrainiter = []
110             Jtrainiter.append(((1/nb) * np.sum(L_i)))
111             acctrainiter.append(((1/nb) * np.sum(ypred
    ↪ == ytrue_shuff[mini_batch])))
112         else:
113             Jtrainiter[it_k] = ((1/nb) * np.sum(L_i))
114             acctrainiter[it_k] = ((1/nb) *
    ↪ np.sum(ypred ==
    ↪ ytrue_shuff[mini_batch]))

```

```

115         # Cost and accuracy for testing data
116         z_test = xtest @ w_mj.T + b_m
117         z_test_norm = z_test - np.max(z_test, axis =
↪ 1, keepdims = True)
118         p_test = np.exp(z_test_norm) /
↪ np.sum(np.exp(z_test_norm), axis = 1,
↪ keepdims = True)
119         L_i_test = np.sum(y_test *
↪ np.log(np.sum(np.exp(z_test_norm), axis =
↪ 1, keepdims = True)) - y_test *
↪ z_test_norm, axis = 1)
120         ypred_test = np.argmax(p_test, axis = 1)
121         if nb == n_train:
122             Jtestniter = []
123             acctestiter = []
124             Jtestniter.append(((1/n_test) *
↪ np.sum(L_i_test)))
125             acctestiter.append(((1/n_test) *
↪ np.sum(ypred_test == ytrue_test)))
126         else:
127             Jtestniter[it_k] = ((1/n_test) *
↪ np.sum(L_i_test))
128             acctestiter[it_k] = ((1/n_test) *
↪ np.sum(ypred_test == ytrue_test))
129         it_k += 1
130
131         Jtrainaccum_av = np.mean(Jtrainaccum)
132         acctrainaccum_av = np.mean(acctrainaccum)
133         Jtestaccum_av = np.mean(Jtestaccum)
134         acctestaccum_av = np.mean(acctestaccum)
135
136         it += 1
137         tot_it += 1
138
139         print("Epoch: (%s/%s), iteration: %s" % (e_p + 1,
↪ ep, it))
140

```

```

141         J_train[e_p] = Jtrainaccum_av
142         acc_train[e_p] = acctrainaccum_av
143         J_test[e_p] = Jtestaccum_av
144         acc_test[e_p] = acctestaccum_av
145
146         e_p += 1
147
148         return J_train, 100 * acc_train, J_test, 100 * acc_test,
149             ↪ Jtrainiter, Jtestniter, 100 * acctrainiter, 100 *
150             ↪ acctestiter, it, w_mj, b_m , z_im, p_im
151
152     n_batch = 1500 # batch size ---> 30 iterations per epoch
153     epochs = 300 # epochs
154     lr0 = 0.01 # initial learning rate
155     tau_it = (n_train//n_batch) - 5 # decay
156     k_plot = 5 # storing accuracy/cost values each k-th
157             ↪ iteration
158
159     Jtrain, acc_train, Jtest, acc_test, J_trainiter, J_testniter,
160     ↪ acc_trainiter, acc_testiter, it, wmj, bm, zim, pim =
161     ↪ softmax_gd_minibatch(x_train, y_train, x_test, y_test,
162     ↪ epochs, n_batch, lr0, tau_it, n_train, n_test, k_plot)
163
164     # Costs and accuracies averaged over k-th iteration per
165     ↪ epoch
166     plt.figure(1)
167     plt_J, = plt.plot(Jtrain, 'r')
168     plt_J_test, = plt.plot(Jtest, 'b')
169     plt.legend([plt_J, plt_J_test], ['Train cost', 'Test cost'])
170     plt.annotate("Final train cost: %s" % (float(Jtrain[-1])),
171             ↪ xycoords = 'figure fraction', xy = (0.4,0.5))
172     plt.annotate("Final test cost: %s" % (float(Jtest[-1])),
173             ↪ xycoords = 'figure fraction', xy = (0.4,0.55))
174     plt.xlabel('Number of epochs')
175     plt.ylabel('Cost J')
176     plt.title('Average cost versus epochs')
177     print("Final train cost: %s" % float(Jtrain[-1]))
178     print("Final test cost: %s" % float(Jtest[-1]))

```

```

170
171 plt.figure(2)
172 plt_acc_train, = plt.plot(acc_train, 'r')
173 plt_acc_test, = plt.plot(acc_test, 'b')
174 plt.legend([plt_acc_train, plt_acc_test], ['Train accuracy',
    ↪ 'Test accuracy'])
175 plt.annotate("Final train accuracy: %s%%" %
    ↪ (float(acc_train[-1])), xycoords = 'figure fraction', xy =
    ↪ (0.4,0.5))
176 plt.annotate("Final test accuracy: %s%%" %
    ↪ (float(acc_test[-1])), xycoords = 'figure fraction', xy =
    ↪ (0.4,0.55))
177 plt.xlabel('Number of epochs')
178 plt.ylabel('Accuracy in %')
179 plt.title('Average accuracy versus epochs')
180 print("Final train accuracy: %s%%" % float(acc_train[-1]))
181 print("Final test accuracy: %s%%" % float(acc_test[-1]))
182
183 ### Costs and accuracies per k-th iteration
184 plt.figure(3)
185 plt_Jtrain_it, = plt.plot(J_trainiter, 'r')
186 plt_Jtest_it, = plt.plot(J_testniter, 'b')
187 plt.legend([plt_Jtrain_it, plt_Jtest_it], ['Train cost', 'Test
    ↪ cost'])
188 plt.annotate("Final train cost: %s" %
    ↪ (float(J_trainiter[-1])), xycoords = 'figure fraction', xy
    ↪ = (0.4,0.5))
189 plt.annotate("Final test cost: %s" % (float(J_testniter[-1])),
    ↪ xycoords = 'figure fraction', xy = (0.4,0.55))
190 plt.xlabel('Total number of iterations')
191 plt.ylabel('Cost J')
192 plt.title('Cost versus iterations')
193 print("Final train cost: %s" % float(J_trainiter[-1]))
194 print("Final test cost: %s" % float(J_testniter[-1]))
195
196 plt.figure(4)
197 plt_acc_train_it, = plt.plot(acc_trainiter, 'r')
198 plt_acc_test_it, = plt.plot(acc_testiter, 'b')

```

```

199 plt.legend([plt_acc_train_it, plt_acc_test_it], ['Train
    ↳ accuracy', 'Test accuracy'])
200 plt.annotate("Final train accuracy: %s%%" %
    ↳ (float(acc_trainiter[-1])), xycoords = 'figure fraction',
    ↳ xy = (0.4,0.5))
201 plt.annotate("Final test accuracy: %s%%" %
    ↳ (float(acc_testiter[-1])), xycoords = 'figure fraction',
    ↳ xy = (0.4,0.55))
202 plt.xlabel('Total number of iterations')
203 plt.ylabel('Accuracy in %')
204 plt.title('Accuracy versus iterations')
205 print("Final train accuracy: %s%%" % float(acc_trainiter[-1]))
206 print("Final test accuracy: %s%%" % float(acc_testiter[-1]))
207
208 figw, axw = plt.subplots(2, 5, figsize = (10,10))
209
210 axw[0,0].imshow(wmj[0, :].reshape(28,28), cmap = 'gray')
211 axw[0,0].set_title('0')
212
213 axw[0,1].imshow(wmj[1, :].reshape(28,28), cmap = 'gray')
214 axw[0,1].set_title('1')
215
216 axw[0,2].imshow(wmj[2, :].reshape(28,28), cmap = 'gray')
217 axw[0,2].set_title('2')
218
219 axw[0,3].imshow(wmj[3, :].reshape(28,28), cmap = 'gray')
220 axw[0,3].set_title('3')
221
222 axw[0,4].imshow(wmj[4, :].reshape(28,28), cmap = 'gray')
223 axw[0,4].set_title('4')
224
225 axw[1,0].imshow(wmj[5, :].reshape(28,28), cmap = 'gray')
226 axw[1,0].set_title('5')
227
228 axw[1,1].imshow(wmj[6, :].reshape(28,28), cmap = 'gray')
229 axw[1,1].set_title('6')
230
231 axw[1,2].imshow(wmj[7, :].reshape(28,28), cmap = 'gray')

```



```

232 axw[1,2].set_title('7')
233
234 axw[1,3].imshow(wmj[8, :].reshape(28,28), cmap = 'gray')
235 axw[1,3].set_title('8')
236
237 axw[1,4].imshow(wmj[9, :].reshape(28,28), cmap = 'gray')
238 axw[1,4].set_title('9')
239
240 plt.show()

```

## 5.4 Full Neural Network

```

1  import numpy as np
2  from load_mnist import load_mnist
3  from matplotlib import pyplot as plt
4
5  x_train, y_train, x_test, y_test = load_mnist()
6
7  def relu(x):
8
9      return np.maximum(0, x)
10
11 def sigmoid(x):
12
13     return (np.exp(x)) / (1 + np.exp(x))
14
15 def relu_deriv(x):
16
17     x[x <= 0] = 0
18     x[x > 0] = 1
19
20     return x
21
22 def sigmoid_deriv(x):
23
24     return sigmoid(x) * (1-sigmoid(x))
25

```

```

26 def softmax(x):
27
28     x_norm = x - np.max(x, axis = 1, keepdims = True)
29     p_x = np.exp(x_norm) / np.sum(np.exp(x_norm), axis = 1,
    ↪ keepdims = True)
30
31     return p_x
32
33 def init_params(M, p, n_hidden):
34
35     W1 = np.random.normal(scale = 0.01, size = (n_hidden[0],
    ↪ p))
36     W2 = np.random.normal(scale = 0.01, size = (n_hidden[1],
    ↪ n_hidden[0]))
37     W3 = np.random.normal(scale = 0.01, size = (n_hidden[2],
    ↪ n_hidden[1]))
38     W4 = np.random.normal(scale = 0.01, size = (M,
    ↪ n_hidden[2]))
39
40     b1 = np.zeros(shape = (n_hidden[0], 1))
41     b2 = np.zeros(shape = (n_hidden[1], 1))
42     b3 = np.zeros(shape = (n_hidden[2], 1))
43     b4 = np.zeros(shape = (M, 1))
44
45     return W1, b1, W2, b2, W3, b3, W4, b4
46
47 def calc_cost(nb, mini_batch, batchBool, y_L, z_L, sz):
48     if batchBool == True:
49         z_L_norm = z_L - np.max(z_L, axis = 1, keepdims =
    ↪ True)
50         loss = np.sum(y_L[mini_batch, :] *
    ↪ np.log(np.sum(np.exp(z_L_norm), axis = 1, keepdims
    ↪ = True)) - y_L[mini_batch, :] * z_L_norm, axis =
    ↪ 1, keepdims = True)
51         cost = (1/nb) * np.sum(loss, axis = 0, keepdims =
    ↪ True)
52
53         dz_L = - y_L[mini_batch, :] + sz

```

```

54         return cost, dz_L
55     else:
56         z_L_norm = z_L - np.max(z_L, axis = 1, keepdims =
57             ↪ True)
58         loss = np.sum(y_L * np.log(np.sum(np.exp(z_L_norm),
59             ↪ axis = 1, keepdims = True)) - y_L * z_L_norm, axis
60             ↪ = 1, keepdims = True)
61         cost = (1/nb) * np.sum(loss, axis = 0, keepdims =
62             ↪ True)
63
64         return cost
65
66 def forward(xt, mb, batchBool, w1, b1, w2, b2, w3, b3, w4,
67     ↪ b4):
68
69     if batchBool == True:
70         z_1 = xt[mb, :] @ w1.T + b1.T
71         q_1 = relu(z_1)
72         z_2 = q_1 @ w2.T + b2.T
73         q_2 = relu(z_2)
74         z_3 = q_2 @ w3.T + b3.T
75         q_3 = relu(z_3)
76         z = q_3 @ w4.T + b4.T
77         softmax_z = softmax(z)
78     else:
79         z_1 = xt @ w1.T + b1.T
80         q_1 = relu(z_1)
81         z_2 = q_1 @ w2.T + b2.T
82         q_2 = relu(z_2)
83         z_3 = q_2 @ w3.T + b3.T
84         q_3 = relu(z_3)
85         z = q_3 @ w4.T + b4.T
86         softmax_z = softmax(z)
87
88     return z_1, q_1, z_2, q_2, z_3, q_3, z, softmax_z
89
90 def backward(q1,q2,q3, z1, z2, z3, dzl, w2, w3, w4, xt, mb):

```

```

86
87     dq_3 = dz1 @ w4
88
89     dz_3 = np.multiply(dq_3, relu_deriv(z3))
90     dq_2 = dz_3 @ w3
91
92     dz_2 = np.multiply(dq_2, relu_deriv(z2))
93     dq_1 = dz_2 @ w2
94
95     dz_1 = np.multiply(dq_1, relu_deriv(z1))
96
97     dW_4 = (1/n_batch) * dz1.T @ q3
98     dW_3 = (1/n_batch) * dz_3.T @ q2
99     dW_2 = (1/n_batch) * dz_2.T @ q1
100    dW_1 = (1/n_batch) * dz_1.T @ xt[mb, :]
101
102    db_4 = (1/n_batch) * np.sum(dz1, axis = 0, keepdims =
103    ↪ True)
104    db_3 = (1/n_batch) * np.sum(dz_3, axis = 0, keepdims =
105    ↪ True)
106    db_2 = (1/n_batch) * np.sum(dz_2, axis = 0, keepdims =
107    ↪ True)
108    db_1 = (1/n_batch) * np.sum(dz_1, axis = 0, keepdims =
109    ↪ True)
110
111    return dW_1, db_1.T, dW_2, db_2.T, dW_3, db_3.T, dW_4,
112    ↪ db_4.T
113
114    def init_params_2(M, p, n_hidden):
115
116        W1 = np.random.normal(scale = 0.01, size = (n_hidden[0],
117        ↪ p))
118        W2 = np.random.normal(scale = 0.01, size = (M,
119        ↪ n_hidden[0]))
120
121        b1 = np.zeros(shape = (n_hidden[0], 1))
122        b2 = np.zeros(shape = (M, 1))

```

```

116
117     return W1, b1, W2, b2
118
119 def backward_2(q1, z1, w1, w2, xt, mb, dz1):
120     dq_1 = dz1 @ w2
121     dz_1 = np.multiply(dq_1, relu_deriv(z1))
122
123     dW_2 = (1/n_batch) * dz1.T @ q1
124
125     db_2 = (1/n_batch) * np.sum(dz1, axis = 0, keepdims =
        ↪ True)
126
127     dW_1 = (1/n_batch) * dz_1.T @ xt[mb, :]
128
129     db_1 = (1/n_batch) * np.sum(dz_1, axis = 0, keepdims =
        ↪ True)
130
131     return dW_1, db_1.T, dW_2, db_2.T
132
133 def forward_2(xt, mb, batchBool, w1, b1, w2, b2):
134
135     if batchBool == True:
136         z_1 = xt[mb, :] @ w1.T + b1.T
137         q_1 = relu(z_1)
138         z = q_1 @ w2.T + b2.T
139         softmax_z = softmax(z)
140     else:
141         z_1 = xt @ w1.T + b1.T
142         q_1 = relu(z_1)
143         z = q_1 @ w2.T + b2.T
144         softmax_z = softmax(z)
145
146     return z_1, q_1, z, softmax_z
147
148 def neural_network(epochs, nb, M, p, k, xtrain, ytrain, xtest,
        ↪ ytest, ntrain, ntest):
149

```

```

150 ytrue_test = np.argmax(ytest, axis = 1) # labels for
    ↪ testing data
151
152 e_p = 0 # epoch counter
153 #lr0 = 0.8 # initial learning rate
154 #lrt = 0.01 * lr0 # final learning rate
155 #t_tau = 30 # iterations until learning rate is set to
    ↪ constant lrt value
156
157 tot_it = 0 # total iteration counter
158 it_k = 0 # k-th iteration counter
159
160 #n_hidden_4 = np.array([50, 50, 50]) # hidden units per
    ↪ layer ---> L - 1 hidden layers
161 #w1,b1,w2,b2,w3,b3,w4,b4 = init_params(M, p ,
    ↪ n_hidden_4)
162
163 n_hidden = np.array([100]) # hidden units per layer --->
    ↪ L - 1 hidden layers
164 w1,b1,w2,b2 = init_params_2(M, p , n_hidden)
165
166 acctrain = np.zeros(shape = (((n_train//nb)//k) *
    ↪ epochs,1))
167 costtrain = np.zeros(shape = (((n_train//nb)//k) *
    ↪ epochs,1))
168
169 acctest = np.zeros(shape = (((n_train//nb)//k) *
    ↪ epochs,1))
170 costtest = np.zeros(shape = (((n_train//nb)//k) *
    ↪ epochs,1))
171
172 while e_p != epochs:
173
174     ### Shuffling indices
175     ind = np.arange(n_train)
176     np.random.shuffle(ind)
177
178     ### Shuffling training data and labels

```

```

179 xt = xtrain[ind, :]
180 yt = ytrain[ind, :]
181
182 ytrue_train = np.argmax(yt, axis = 1) # labels for
    ↪ training data
183
184 it = 0 # iteration counter for an epoch
185
186 while it != n_train//nb:
187
188     #lr = (1 - (it/t_tau)) * lr0 + (it/t_tau) * lrt
189     lr = 0.8
190
191     mini_batch = np.random.randint(0, n_train, size =
    ↪ nb) # batch indices
192
193     #2-layer code
194     z1, q1, z, softz = forward_2(xt, mini_batch, True,
    ↪ w1,b1,w2,b2)
195     train_cost, dzL = calc_cost(nb, mini_batch, True,
    ↪ yt, z, softz)
196     dw1, db1, dw2, db2 = backward_2(q1, z1, w1, w2,
    ↪ xt, mini_batch, dzL)
197     _, _, z_test, softz_test = forward_2(xtest,
    ↪ mini_batch, False, w1,b1,w2,b2)
198
199     # 4-layer code
200     #z1, q1, z2, q2, z3, q3, z, softz = forward(xt,
    ↪ mini_batch, True, w1,b1,w2,b2,w3,b3,w4,b4)
201     #train_cost, dzL = calc_cost(nb, mini_batch,
    ↪ True, yt, z, softz)
202     #dw1, db1, dw2, db2, dw3, db3, dw4, db4 =
    ↪ backward(q1, q2, q3, z1, z2, z3, dzL, w2, w3,
    ↪ w4, xt, mini_batch)
203     #_, _, _, _, _, z_test, softz_test =
    ↪ forward(xtest, mini_batch, False, w1, b1, w2,
    ↪ b2, w3, b3, w4, b4)
204

```

```

205         test_cost = calc_cost(n_test, mini_batch, False,
    ↪ ytest, z_test, softz_test)
206
207         w1 = w1 - lr * dw1
208         w2 = w2 - lr * dw2
209         b1 = b1 - lr * db1
210         b2 = b2 - lr * db2
211
212         #w3 = w3 - lr * dw3
213         #w4 = w4 - lr * dw4
214         #b3 = b3 - lr * db3
215         #b4 = b4 - lr * db4
216
217         if tot_it % k == 0:
218             y_predtrain = np.argmax(softz, axis = 1)
219             acctrain[it_k] = 100 * ((1/nb) *
    ↪ np.sum(y_predtrain ==
    ↪ ytrue_train[mini_batch]))
220             costtrain[it_k] = train_cost
221
222             y_predtest = np.argmax(softz_test, axis = 1)
223             acctest[it_k] = 100 * ((1/n_test) *
    ↪ np.sum(y_predtest == ytrue_test))
224             costtest[it_k] = test_cost
225
226             it_k += 1
227
228             tot_it += 1
229             it += 1
230
231             print("Epoch: (%s/%s), iteration: %s" % (e_p + 1,
    ↪ epochs, it))
232
233             e_p += 1
234
235     #return w1,w2,w3,w4,b1,b2,b3,b4, costtrain, acctrain,
    ↪ costtest, acctest

```



```

236     return w1,w2, softz, costtrain, acctrain, costtest,
        ↪ acctest
237
238 M = 10 # number of classes/ digits
239 p = x_train.shape[1] # number of input pixels - 784
        ↪ (flattened 28x28 image)
240
241 n_train = x_train.shape[0] # number of training examples -
        ↪ 60000
242 n_test = x_test.shape[0] # number of testing examples -
        ↪ 10000
243
244 n_batch = 1500 # batch size
245 epochs = 50 # number of epochs
246
247 k_acc = 5
248
249 w1,w2,sz, costtrain, acctrain, costtest, acctest =
        ↪ neural_network(epochs, n_batch, M, p, k_acc, x_train,
        ↪ y_train, x_test, y_test, n_train, n_test)
250 #w1,w2,w3,w4,b1,b2,b3,b4, costtrain, acctrain, costtest,
        ↪ acctest = neural_network(epochs, n_batch, M, p, k_acc,
        ↪ x_train, y_train, x_test, y_test, n_train, n_test)
251
252 plt.figure(1)
253 plt_Jtrain_it, = plt.plot(costtrain, 'r')
254 plt_Jtest_it, = plt.plot(costtest, 'b')
255 plt.legend([plt_Jtrain_it, plt_Jtest_it], ['Train cost', 'Test
        ↪ cost'])
256 plt.annotate("Final train cost: %s" % (float(costtrain[-1])),
        ↪ xycoords = 'figure fraction', xy = (0.2,0.5))
257 plt.annotate("Final test cost: %s" % (float(costtest[-1])),
        ↪ xycoords = 'figure fraction', xy = (0.2,0.55))
258 plt.xlabel('Total number of iterations')
259 plt.ylabel('Cost J')
260 plt.title('Cost versus iterations')
261 print("Final train cost: %s" % float(costtrain[-1]))
262 print("Final test cost: %s" % float(costtest[-1]))

```

```

263
264 plt.figure(2)
265 plt_acc_train_it, = plt.plot(acctrain, 'r')
266 plt_acc_test_it, = plt.plot(acctest, 'b')
267 plt.legend([plt_acc_train_it, plt_acc_test_it], ['Train
    ↪ accuracy', 'Test accuracy'])
268 plt.annotate("Final train accuracy: %s%%" %
    ↪ (float(acctrain[-1])), xycoords = 'figure fraction', xy =
    ↪ (0.2,0.5))
269 plt.annotate("Final test accuracy: %s%%" %
    ↪ (float(acctest[-1])), xycoords = 'figure fraction', xy =
    ↪ (0.2,0.55))
270 plt.xlabel('Total number of iterations')
271 plt.ylabel('Accuracy in %')
272 plt.title('Accuracy versus iterations')
273 print("Final train accuracy: %s%%" % float(acctrain[-1]))
274 print("Final test accuracy: %s%%" % float(acctest[-1]))
275
276 plt.show()

```