

Q1.1

$$\text{softmax}(x) = \text{softmax}(x + c)$$

$$\text{softmax}(x_i) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

$$\text{softmax}(x + c) = \frac{e^{x_i + c}}{\sum_j e^{x_j + c}} = \frac{e^{x_i} \cancel{e^c}}{\sum_j e^{x_j} \cancel{e^c}}$$

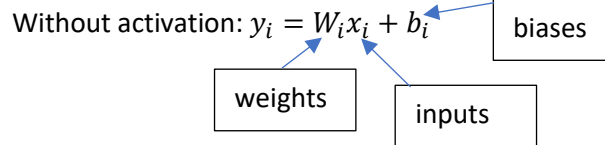
$$\text{softmax}(x + c) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

Q1.2

- Range (0,1] and sum over all the elements is equal to 1
- "...a probability distribution of the input x-vector"
- $1 - s = e^{x_i}$, calculates the exponential results of x_i and preserves the order
 - $2 - s = \sum s_i$, calculates the sum of s_i (vector normalization value)
 - $3 - \frac{1}{s}s_i$, Normalizes vector (between 0 and 1) and outputs the probability

Q1.3

Forward propagation of a multi-layer network, in solution below a 3-layer network will be used to show calculation and resulting linear regression problem equation:



$$i = 1: y_1 = W_1 x_1 + b_1 \text{ (1)}$$

$$i = 2: y_2 = W_2 y_1 + b_2 \text{ (2)}$$

$$i = 3: y_3 = W_3 y_2 + b_3 \text{ (3)}$$

Substitute equations (1) and (2) into equation 3 to get:

$$y_3 = W_3(W_2(W_1 x_1 + b_1) + b_3)$$

$$y_3 = W_3 W_2 W_1 x_1 + W_3 W_2 b_1 + W_3 b_2 + b_3$$

This is the same as solving a linear regression problem, as it is just a linear combination of the weights, biases, and inputs ($y = Wx + b$).

Q1.4

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\begin{aligned}\frac{d}{dx} \sigma(x) &= \frac{d}{dx} \left(\frac{1}{1 + e^{-x}} \right) = - \frac{1(e^{-x} * -1)}{(1 + e^{-x})^2} \\ &= \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1}{\underbrace{1 + e^{-x}}_{\boxed{\sigma(x)}}} * \frac{e^{-x}}{(1 + e^{-x})}\end{aligned}$$

$$= \frac{1}{1 + e^{-x}} * \left(1 - \frac{1}{1 + e^x} \right) = \sigma(x)(1 - \sigma(x))$$

Q1.5

Find $\frac{\partial J}{\partial W}, \frac{\partial J}{\partial x}, \frac{\partial J}{\partial b}$:

$$y = Wx + b$$

Scalar:

$$\text{Using: } y_i = \sum_{j=1}^d (x_j W_{ij} + b_i)$$

$$\frac{\partial J}{\partial W} = \sum_j \frac{\partial J}{\partial y} * \frac{\partial y}{\partial W} = \delta x^T \quad (1)$$

$$\frac{\partial J}{\partial x} = \sum_j \frac{\partial J}{\partial y} * \frac{\partial y}{\partial x} = W^T \delta \quad (2)$$

$$\frac{\partial J}{\partial b} = \sum_j \frac{\partial J}{\partial y} * \frac{\partial y}{\partial b} = \delta \quad (3)$$

Matrix:

$$\begin{bmatrix} \delta_1 x_1 & \delta_1 x_2 & \dots & \delta_1 x_d \\ \delta_2 x_1 & \delta_2 x_2 & \dots & \delta_2 x_d \\ \vdots & \vdots & \ddots & \vdots \\ \delta_k x_1 & \delta_k x_2 & \dots & \delta_k x_d \end{bmatrix} \quad (\text{Matrix for equation 1, } \in R^{k \times d})$$

$$\begin{bmatrix} W_1 \delta_1 \\ W_2 \delta_2 \\ \vdots \\ W_d \delta_d \end{bmatrix} \quad (\text{Matrix for equation 2, } \in R^{d \times 1})$$

$$\begin{bmatrix} \delta_1 \\ \delta_2 \\ \vdots \\ \delta_k \end{bmatrix} \quad (\text{Matrix for equation 3, } \in R^{k \times 1})$$

Q1.6.1

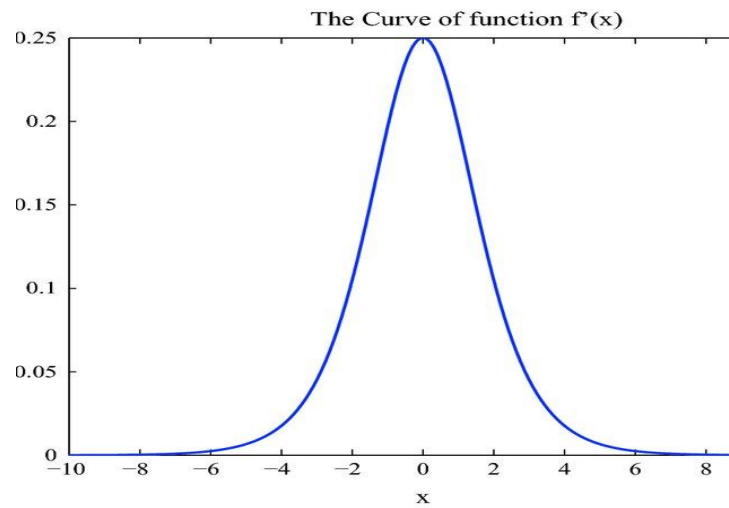


Figure 1: [Link 1](#)

In the plot shown above, the derivative of sigmoid (σ') is shown and the values range from $(0, 0.25]$ whereas the values for sigmoid will range from $(0, 1]$. After applying the activation function (this is done more than once) to multiple layers the values in range of 0.25 will decrease quickly. This might cause the “vanishing gradient” problem, as these values will not affect the training implementation.

Q1.6.2

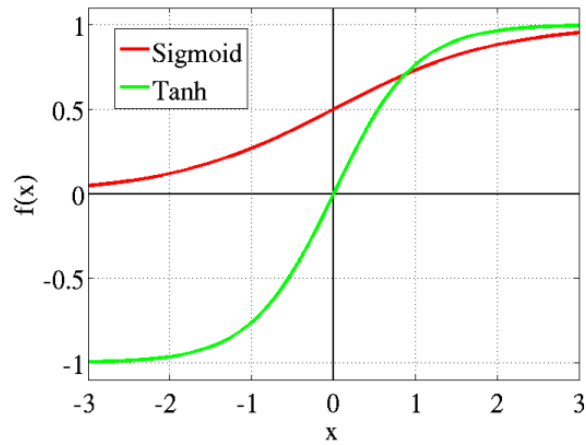


Figure 2: [Link 2](#)

As seen in the figure above the output range of sigmoid is $(0,1]$ and tanh is $(-1,1)$; tanh is better because the output will have an average value near zero. Also, tanh will correspond to x “sign” (i.e. when x -negative, tanh-negative and x -positive, tanh-positive).

Q1.6.3

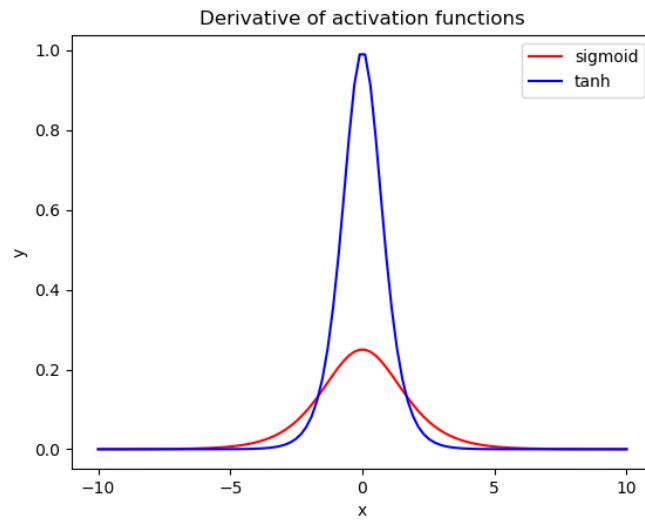


Figure 3: [Link 3](#)

The derivative of tanh will range from (0,1) so the gradient will have less of a vanishing gradient as it decreases slower because of the larger value (this is not seen in the derivative of the sigmoid value which ranges from (0,0.25) and vanishes quickly).

Q1.6.4

$$\sigma(x) = \frac{1}{1+e^{-x}}; \tanh(x) = \frac{1-e^{-2x}}{1+e^{-2x}}$$

$$\frac{1-e^{-2x}}{1+e^{-2x}} = \frac{2-(1+e^{-2x})}{1+e^{-2x}} = 2\sigma(x) - 1$$

$$2\sigma(x) - 1 = \tanh(x) = \frac{1-e^{-2x}}{1+e^{-2x}}$$

$$\tanh(x) = 2\sigma(x) - 1$$

Q2.1.1

Initializing a network with zeros will cause all nodes to have the same weights (and biases) so all the outputs will be the same for every input. This means that no meaningful information will be gathered. During backward propagation, the gradients will be gathered. During backward propagation, the gradients will be the same value for each node in the layer and the model won't learn. Instead, it will learn only one function versus learning many.

Q2.1.2

```
6 ##### Q 2.1 #####
7 # initialize b to 0 vector
8 # b should be a 1D array, not a 2D array with a singleton dimension
9 # we will do  $XW + b$ .
10 # X be [Examples, Dimensions]
11 def initialize_weights(in_size,out_size,params,name=''):
12     W, b = None, None
13     b = np.sqrt(6/(in_size+out_size))
14     W = np.random.uniform(-b,b,(in_size,out_size))
15
16     b = np.zeros((out_size))
17     params['W' + name] = W
18     params['b' + name] = b
19
20 ##### Q 2.2.1 #####
21 # x is a matrix
22 # a sigmoid activation function
23
24 def sigmoid(x):
25     res = None
26     res = 1/(1+np.exp(-x))
27
28     return res
29
30 ##### Q 2.2.1 #####
31 def forward(X,params,name='',activation=sigmoid):
32     """
33     Do a forward pass
34
35     Keyword arguments:
36     X -- input vector [Examples x D]
37     params -- a dictionary containing parameters
38     name -- name of the layer
39     activation -- the activation function (default is sigmoid)
40     """
41     pre_act, post_act = None, None
42
43     # get the layer parameters
44     W = params['W' + name]
45     b = params['b' + name]
46
47     #Activation..
48     pre_act = np.dot(X,W)
49     pre_act = pre_act+b
50
51     #Every value in pre_act
52     post_act = activation(pre_act)
53
54
55     # store the pre-activation and post-activation values
56     # these will be important in backprop
57     params['cache_' + name] = (X, pre_act, post_act)
58
59     return post_act
```

Q2.1.3

In order to prevent symmetry during training and eliminate issues of initializing with zeros. This also helps to avoid having large values in forward propagation, so the network will converge faster, and exploding during backward propagations can be avoided. Scaling the initialization by layer size, will allow for similar values to be computed in each layer (variance in the weight gradients will be the same). Again, this helps get values that aren't very large and helps avoid the vanishing gradient problem.

Q2.2.1

```
24 def sigmoid(x):
25     res = None
26     res = 1/(1+np.exp(-x))
27
28     return res
29
30 ##### Q 2.2.1 #####
31 def forward(X,params,name='',activation=sigmoid):
32     """
33     Do a forward pass
34
35     Keyword arguments:
36     X -- input vector [Examples x D]
37     params -- a dictionary containing parameters
38     name -- name of the layer
39     activation -- the activation function (default is sigmoid)
40     """
41     pre_act, post_act = None, None
42
43     # get the layer parameters
44     W = params['W' + name]
45     b = params['b' + name]
46
47     #Activation..
48     pre_act = np.dot(X,W)
49     pre_act = pre_act+b
50
51     #Every value in pre_act
52     post_act = activation(pre_act)
53
54
55     # store the pre-activation and post-activation values
56     # these will be important in backprop
57     params['cache_' + name] = (X, pre_act, post_act)
58
59     return post_act
60
```

Q2.2.2

```
61 ##### Q 2.2.2 #####
62 # x is [examples,classes]
63 # softmax should be done for each row
64 def softmax(x):
65     res = None
66
67     #get max x [examples,classes]
68     x_max = np.max(x,axis=1,keepdims=True) #Keepdims will return a 2D array, when not using this parameter cannot perform subtraction of matrix
69     shifted_x = x-x_max
70     x_exp = np.exp(shifted_x)
71
72     x_sum = np.exp(shifted_x).sum(axis=1, keepdims=True)
73     res = x_exp/x_sum
74
75     return res
```

Q2.2.3

```
77 ##### Q 2.2.3 #####
78 # compute total loss and accuracy
79 # y is size [examples,classes]
80 # probs is size [examples,classes]
81
82 def compute_loss_and_acc(y, probs):
83     loss, acc = None, None
84     acc = []
85     loss = 0
86
87
88     for i in range (0,y.shape[0]):
89         #Compute loss
90         loss = loss-np.sum(y[i]*np.log(probs[i]))
91         probability = np.argmax(probs[i])
92         y_2= np.argmax(y[i])
93
94         #Find where probs and y are equal
95         if probability == y_2:
96             acc.append(True)
97         else:
98             acc.append(False)
99     acc = np.mean(acc)
100
101     return loss, acc
102
```

Q2.3

```
103 ##### Q 2.3 #####
104 # we give this to you
105 # because you proved it
106 # it's a function of post_act
107 def sigmoid_deriv(post_act):
108     res = post_act*(1.0-post_act)
109     return res
110
111 def backwards(delta,params,name='',activation_deriv=sigmoid_deriv):
112     """
113     Do a backwards pass
114
115     Keyword arguments:
116     delta -- errors to backprop
117     params -- a dictionary containing parameters
118     name -- name of the layer
119     activation_deriv -- the derivative of the activation_func
120     """
121     grad_X, grad_W, grad_b = None, None, None
122     # everything you may need for this layer
123     W = params['W' + name]
124     b = params['b' + name]
125     X, pre_act, post_act = params['cache_' + name]
126     # det = np.matmul(delta,activation_deriv(post_act))
127     det = delta*activation_deriv(post_act)
128     grad_W = np.dot(np.transpose(X),det)
129     grad_b = np.sum(det, axis=0)
130     grad_X = np.dot(det,np.transpose(W))
131
132     # store the gradients
133     params['grad_W' + name] = grad_W
134     params['grad_b' + name] = grad_b
135     return grad_X
```


Q2.4

```
137 ##### Q 2.4 #####
138 # split x and y into random batches
139 # return a list of [(batch1_x, batch1_y)...]
140 def get_random_batches(x, y, batch_size):
141     batches = []
142     # Takes x and y as inputs and splits it into random batches
143     batches_num = x.shape[0] / batch_size
144
145     # Get random batches
146     random_batches_x_y = np.random.choice(x.shape[0], y.shape[0], False)
147
148     i = 0
149     while i < batches_num:
150         if len(random_batches_x_y) > batch_size:
151             indexes = random_batches_x_y[0:batch_size]
152             random_batches_x_y = random_batches_x_y[batch_size:]
153         else:
154             indexes = random_batches_x_y
155
156         batches.append((x[indexes], y[indexes]))
157         i += 1
158
159     return batches
160
```

```
85 # WRITE A TRAINING LOOP HERE
86 max_iters = 500
87 learning_rate = 1e-3
88 # with default settings, you should get loss < 35 and accuracy > 75%
89 for itr in range(max_iters):
90     total_loss = 0
91     avg_acc = 0
92     for xb, yb in batches:
93         # forward
94         yl = forward(xb, params, 'layer1', sigmoid)
95         probas = forward(yl, params, 'output', softmax)
96         # loss
97         loss, acc = compute_loss_and_acc(yb, probas)
98
99         # be sure to add loss and accuracy to epoch totals
100         total_loss += loss
101         # avg_acc += 1
102         avg_acc += acc
103         # backward
104         delta = probas - yb
105         delta_2 = backwards(delta, params, 'output', linear_deriv)
106         backwards(delta_2, params, 'layer1', sigmoid_deriv)
107         # apply gradient
108
109         # W Layer
110         params['Wlayer1'] = params['Wlayer1'] - learning_rate * params['grad_Wlayer1']
111         params['Woutput'] = params['Woutput'] - learning_rate * params['grad_Woutput']
112
113         # B layer
114         params['boutput'] = params['boutput'] - learning_rate * params['grad_boutput']
115         params['bplayer1'] = params['bplayer1'] - learning_rate * params['grad_bplayer1']
116
117         # gradients should be summed over batch samples
118     avg_acc = avg_acc / len(batches)
```

Q2.5

```

147 eps = 1e-6
148 N = v.shape[0]
149
150 for k,v in params.items():    # for each value inside the parameter
151     if '_' in k:
152         continue
153     grad_ = params['grad_'+k]
154     if 'b' in k: #Bias
155         for i in range(0,N):
156             v_val = v[i]
157             v[i] = v[i] - eps
158             #run the network
159             h_ = forward(x,params,'layer1')
160             prob = forward(h_,params,'output',softmax)
161
162             #Don't reassign loss,acc; get loss
163             loss_new,acc_new = compute_loss_and_acc(y,prob)
164             v[i] = v_val
165
166             #add epsilon
167             v[i] = v[i] + eps
168             h_2 = forward(x,params,'layer1')
169             prob_2 = forward(h_2,params,'output',softmax)
170
171             loss3,acc3 = compute_loss_and_acc(y,prob_2)
172             # compute derivative with central diffs
173             grad_[i] = (loss3-loss_new)/(2*eps)
174             # restore the original parameter value
175             v[i] = v[i]-eps
176     elif 'W' in k:
177         for i in range(v.shape[0]):
178             for j in range(v.shape[1]):
179                 v_val = v[i,j]
180                 v[i][j] = v[i][j]-eps
181
182                 #run the network
183                 h_ = forward(x,params,'layer1')

```

```

184                 prob = forward(h_,params,'output',softmax)
185                 loss_new, acc_new = compute_loss_and_acc(y, prob)
186
187                 v[i][j] = v_val
188
189                 #add epsilon
190                 v[i][j] += eps
191                 h_2 = forward(x,params,'layer1')
192                 prob_2 = forward(h_2,params,'output',softmax)
193
194                 loss3,acc3 = compute_loss_and_acc(y,prob_2)
195
196                 # compute derivative with central diffs
197                 grad_[i] = (loss3-loss_new)/(2*eps)
198                 # restore the original parameter value
199                 v[i][j] = v_val
200
201

```

Q3.1

Batch size = 32; Learning rate = $3e-3$

```
Validation accuracy: 0.7747222222222222  
Test accuracy: 0.7916666666666666
```

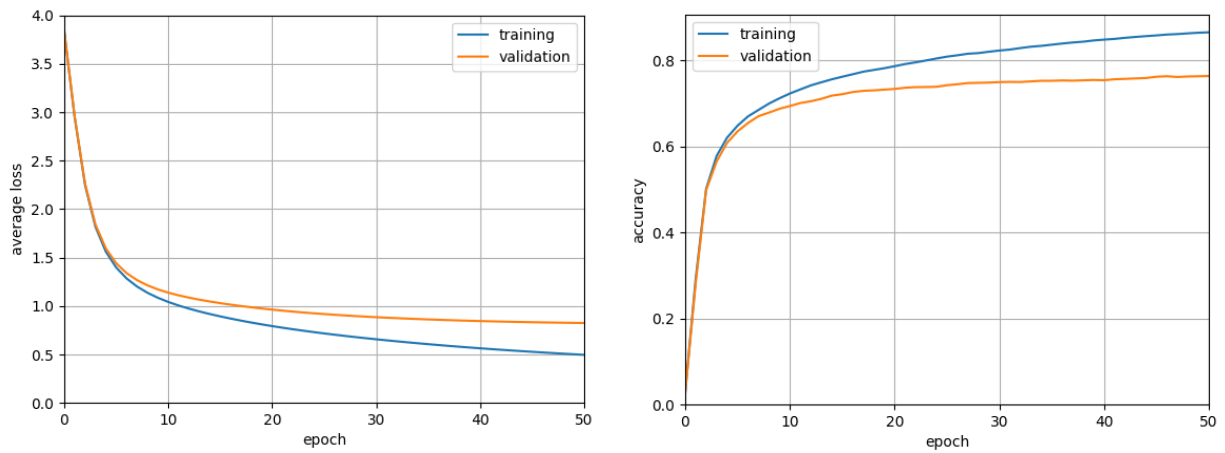


Figure 4: Learning rate of $3e-3$

Q3.2

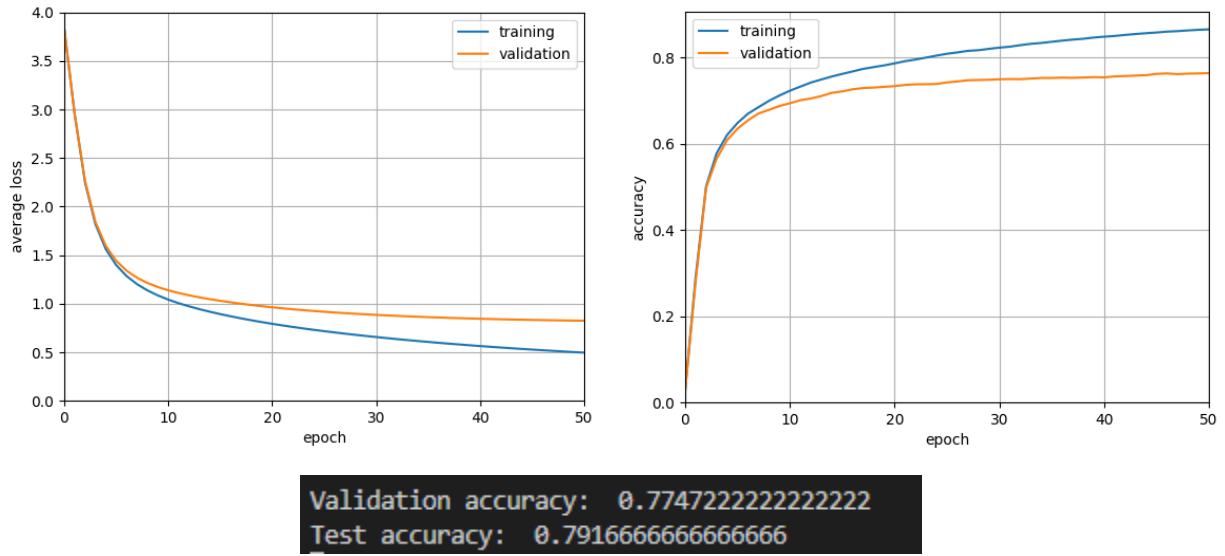


Figure 5: Best Learning rate of $3e-3$

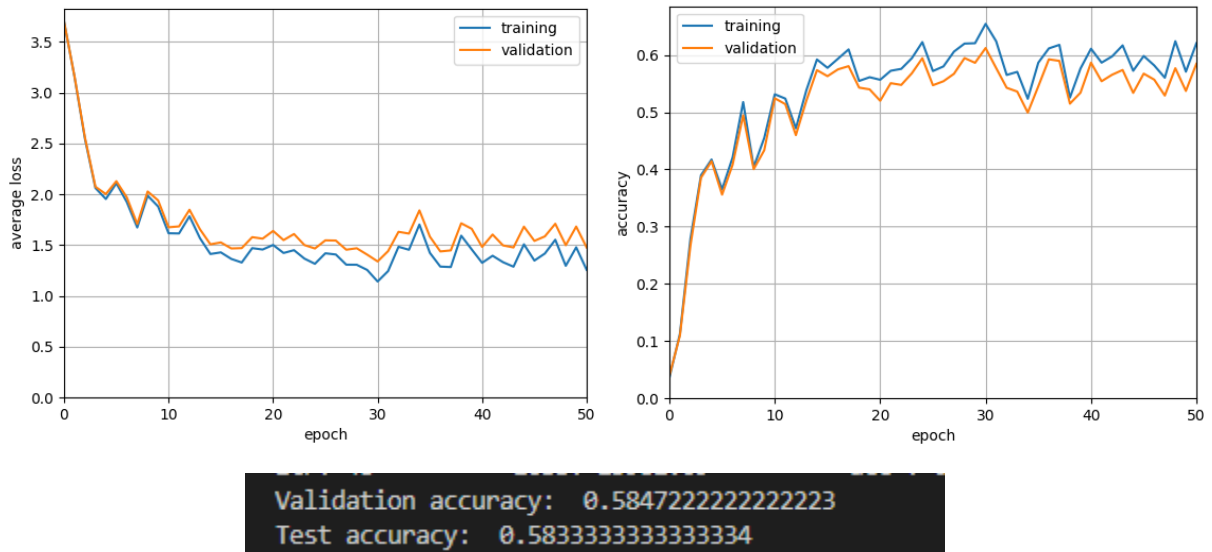


Figure 6: Learning rate of $3e-3 \times 10$

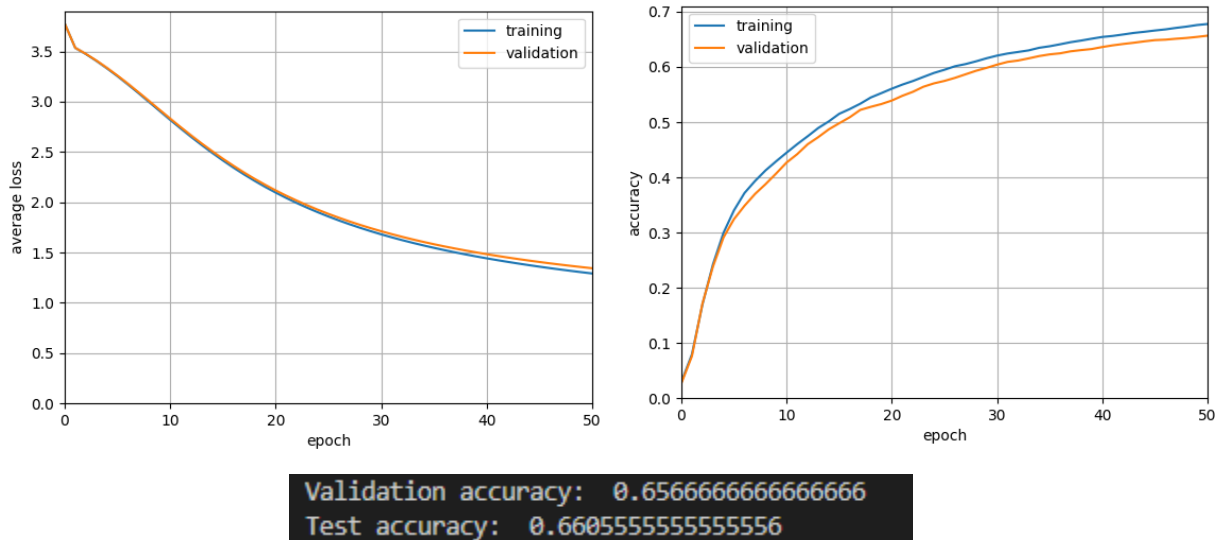


Figure 7: Learning rate of $3e-3 \times 1/10$

The **best learning rate was the tuned learning rate of $3e-3$** with a validation accuracy of 77.47% and test accuracy of 79.17%. When the learning rate is larger (see figure 6) there was more oscillations in the loss and the accuracy as the system adjusts, the accuracy was lower compared to the best learning rate and there was also more loss. Then with the decreased learning rate (see figure 7) the accuracy decreased but the system graphs overall are smooth (smoothest graphs out of all learning rates).

Q3.3

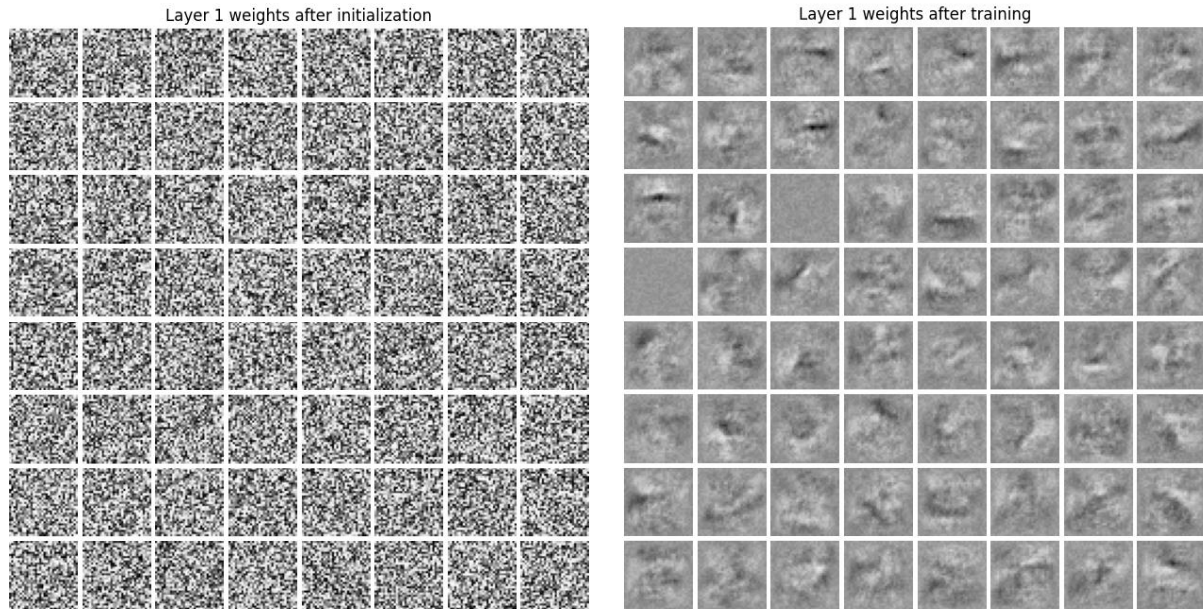


Figure 8: Network of weight after initialization (left) and after training (right)

The weights after initialization show random distribution, while the first layer after training there is more clarity to the weights (can start to see what we are looking for).

Q3.4

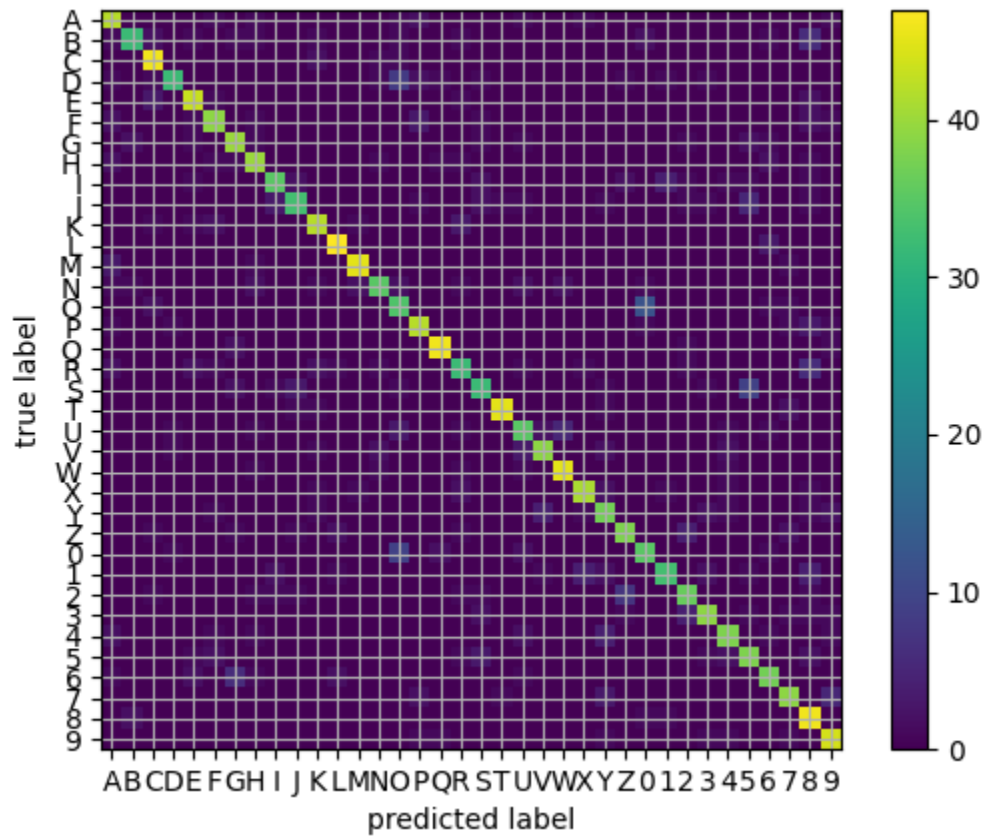


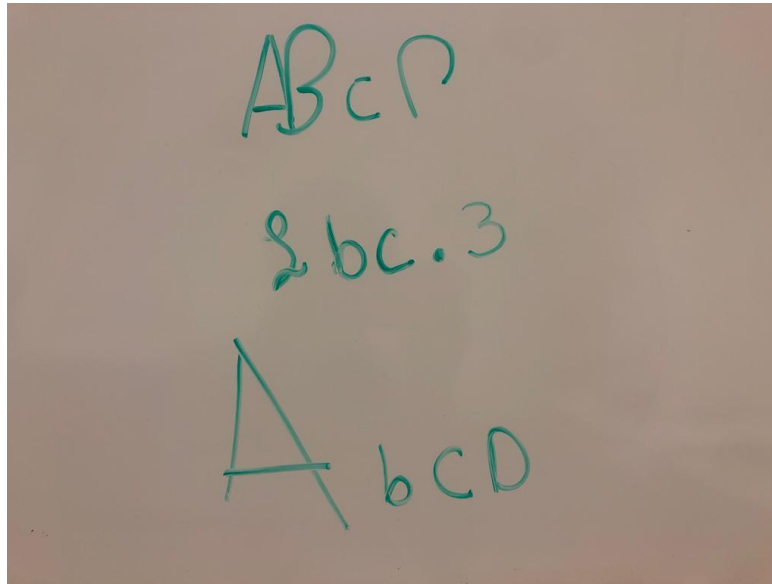
Figure 9: Confusion Matrix

The most commonly confused sets are (from looking at the image above): ['G',6], ['O',0], ['S',5],['I',1] and ['D',0]. This is because the shapes look very similar, and the trained network (done by the computer) will have a difficult time distinguishing between the two labels.

Q4.1

Assumptions made:

- Letters are of similar sizes and all text is either numbers or letters, that we have training data available for
- Letters are separate from each other, so no overlapping of characters and each character is fully connected



In the image above I have drawn out some examples of the overlapping characters and incomplete characters (first row). The second row of the image shows two symbols instead of characters, and the last row shows different sizes of characters.

Q4.2

```
def findLetters(image):
    #Draw rectangles on the image to show the bounding boxes
    # insert processing in here
    # one idea estimate noise -> denoise -> greyscale -> threshold -> morphology -> label -> skip small boxes
    # this can be 10 to 15 lines of code using skimage functions
    bboxes = []
    bw_img = None
    area_box = 0
    area_box_list = []

    #denoise image
    img_noise = skimage.restoration.denoise_bilateral(image, multichannel=True)
    # convert image to grayscale
    img_gray = skimage.color.rgb2gray(img_noise)
    # get a binary image first
    threshold = skimage.filters.threshold_otsu(img_gray)

    #Process the binary image
    mask = img_gray>threshold

    #Use morphology: # Thicken letters; skimage.morphology dilation step
    mask = skimage.morphology.erosion(mask)
    mask = skimage.morphology.erosion(mask)
    mask = skimage.morphology.dilation(mask)

    bw_img = mask
    # mask = skimage.morphology.remove_small_objects(mask, 50)
    # mask = skimage.morphology.remove_small_holes(mask, 50)

    # skimage to find boundaries of the letters; skimage.measure.regionprops
    labels = skimage.measure.label(bw_img, background=1, connectivity=2)
    bb = skimage.measure.regionprops(labels,bw_img)

    # check if boundary box is larger than average box size
    # skip the small boxes
    # determine threshold (100-200) for comparing the top edge of the bounding boxes to the next bbox
    # compare against the letter next to it
```

```
51
52 ✓   for box in bb:
53       area_box = area_box+ box.area
54       area_box_list.append(area_box)
55   avg_box_size = area_box/len(bb)
56 ✓   for i in bb:
57       if i.area >avg_box_size/4:
58           bboxes.append(i.bbox)
59           # print("here")
60   # bw_img = (~bw_img).astype(np.float)
61   return bboxes, bw_img
```

Q4.3

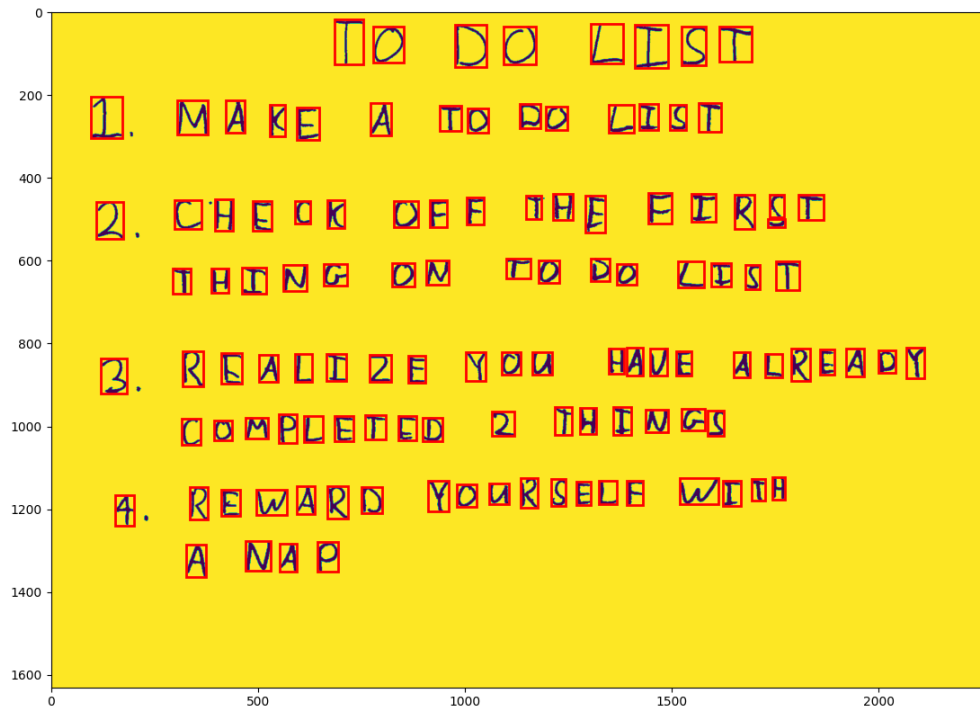


Figure 10: Accuracy of 1st image

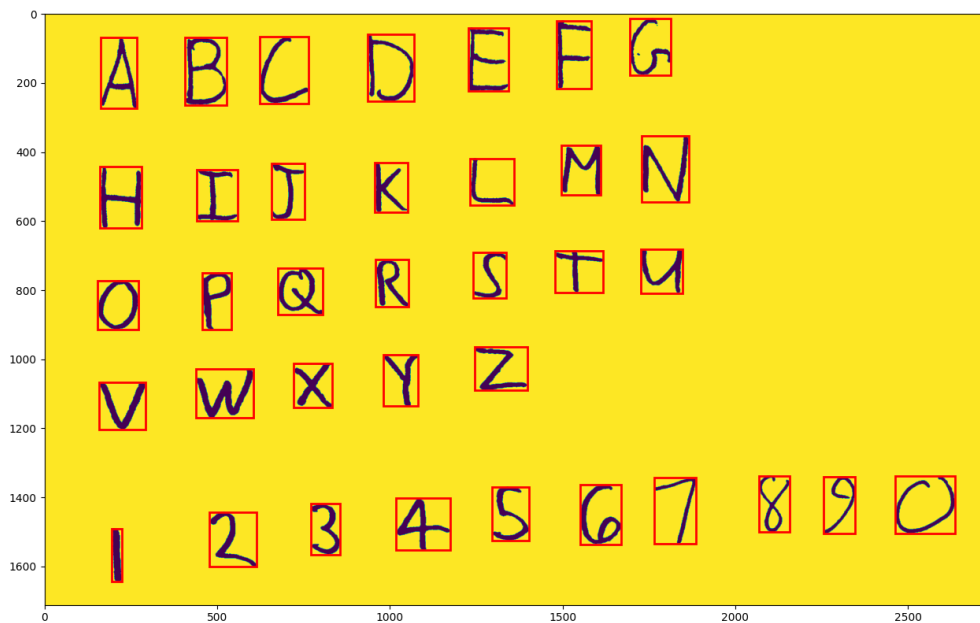


Figure 11: Accuracy of 2nd image



Figure 12: Accuracy of 3rd image

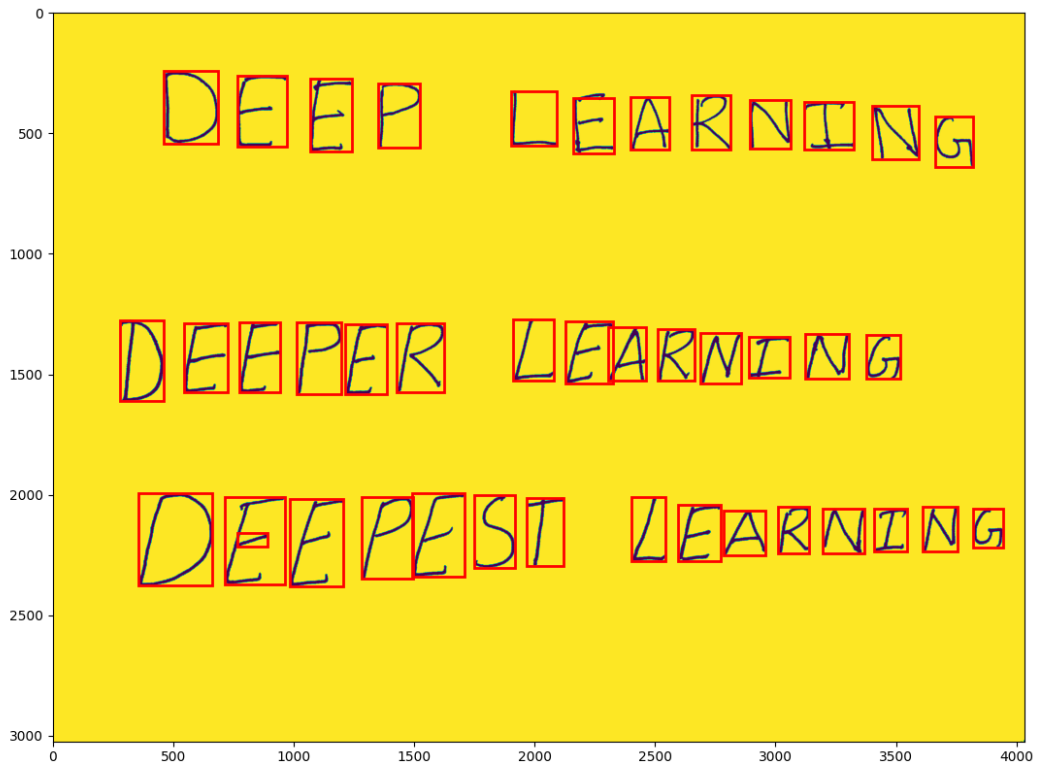


Figure 13: Accuracy of 4th image

Q4.4

01_List:

TO DO LIST

I MAXE A TD DOLIST

2 CH8CK OFE 7HE FIR9CT

7HING ON TO DO LIST

3 R2ALIZE YOU HAVE ALREAU7

COMPLET2D ZYHINGS

4 RE4ARD YOURSELF WITH

A NAP

02_letters:

A B C U E F G

H I J K L M U

Q P Q R J T U

V W X Y Z

8 Z 3 F S G 7 8 7 O

03_haiku:

HAIKUS ARE EMASY

BUT SQMETIMES THEY DON'T MAKG SENG

M

REFRIGERAT9OR

04_deep:

D5EP LC2RMING

DFFYFX L8AKNIXG

DFECPFST LFARNING

This was the best accuracy/classification of the images that I could get, even after performing dilation to make the characters thicker such that they would match (or look similar) to the thicker characters that the network was trained on in 3.1. I suspect there might not be enough dilation but there was not a noticeable difference in the characters even after multiple "rounds" of dilation.

Q5.1.1

```
28 # Q5.1 & Q5.2
29 # initialize layers here
30 initialize_weights(1024, hidden_size, params, 'layer1') #ReLU
31 initialize_weights(hidden_size, hidden_size, params, 'hidden1') #ReLU
32 initialize_weights(hidden_size, hidden_size, params, 'hidden2') #ReLU
33 initialize_weights(hidden_size, 1024, params, 'output') #Sigmoid
34
35 # should look like your previous training loops
36 losses = []
37 for itr in range(max_iters):
38     total_loss = 0
39     for xb, _ in batches:
40
41         #forward
42         h = forward(xb, params, 'layer1', relu)
43         h_1 = forward(h, params, 'hidden', relu)
44         h_2 = forward(h_1, params, 'hidden2', relu)
45         output_ = forward(h_2, params, 'output', sigmoid)
46
47         # loss
48         p_X = np.square(output_ - xb)
49         loss = p_X.sum()
50
51         # be sure to add loss and accuracy to epoch totals
52         total_loss += loss
53
54         # backward
55         delta = output_ - xb
56         delta_2 = backwards(2*delta, params, 'output', sigmoid_deriv)
57         delta_3 = backwards(delta_2, params, 'hidden2', relu_deriv)
58         delta_4 = backwards(delta_3, params, 'hidden', relu_deriv)
59         backwards(delta_4, params, 'layer1', relu_deriv)
60
61
62
63
64
65 #print(params.keys())
66 params_w_list = ['wlayer1', 'wlayer2', 'wlayer3', 'woutput', 'whidden1', 'whidden2']
67 params_b_list = ['blayer1', 'blayer2', 'blayer3', 'boutput', 'bhidden1', 'bhidden2']
68
69 # apply gradient (5.1.1)
70 #W Layer
71 for i in params_w_list:
72     params[i] -= learning_rate * params['grad_' + i]
73 #B layer
74 for j in params_b_list:
75     params[j] -= learning_rate * params['grad_' + j]
76
77
78
79 losses.append(total_loss / train_x.shape[0])
80 if itr % 2 == 0:
81     print("itr: {:02d} \t loss: {:.2f}".format(itr, total_loss))
82 if itr % lr_rate == lr_rate - 1:
83     learning_rate *= 0.9
```

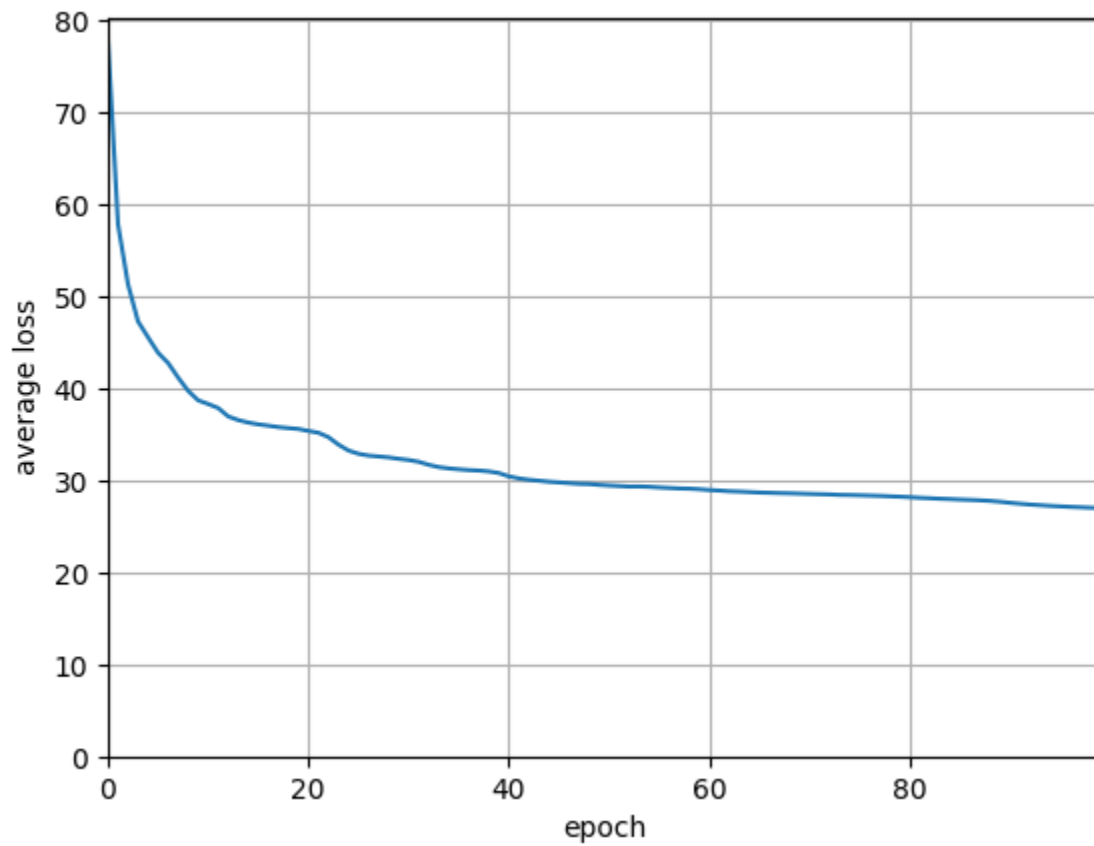
Q5.1.2

```

35 # should look like your previous training loops
36 losses = []
37 for itr in range(max_iters):
38     total_loss = 0
39     for xb, _ in batches:
40
41         #forward
42         h = forward(xb, params, 'layer1', relu)
43         h_1 = forward(h, params, 'hidden', relu)
44         h_2 = forward(h_1, params, 'hidden2', relu)
45         output_ = forward(h_2, params, 'output', sigmoid)
46
47         # loss
48         p_X = np.square(output_-xb)
49         loss = p_X.sum()
50
51         # be sure to add loss and accuracy to epoch totals
52         total_loss += loss
53
54         # backward
55         delta = output_ - xb
56         delta_2 = backwards(2*delta, params, 'output', sigmoid_deriv)
57         delta_3 = backwards(delta_2, params, 'hidden2', relu_deriv)
58         delta_4 = backwards(delta_3, params, 'hidden', relu_deriv)
59         backwards(delta_4, params, 'layer1', relu_deriv)
60
61         params_w_list = ['Wlayer1', 'Wlayer2', 'Wlayer3', 'Woutput']
62         params_b_list = ['blayer1', 'blayer2', 'blayer3', 'boutput']
63
64     # apply gradient (5.1.1)
65     #W Layer
66     # for i in params_w_list:
67     #     params[i] += learning_rate*params['grad_'+ i]
68     # #B layer
69
70 #print(params.keys())
71 params_w_list = ['Wlayer1', 'Wlayer2', 'Wlayer3', 'Woutput', 'Whidden1', 'Whidden2']
72 params_b_list = ['blayer1', 'blayer2', 'blayer3', 'boutput', 'bhidden1', 'bhidden2']
73
74 # apply gradient (5.1.1)
75 # #W Layer
76 # for i in params_w_list:
77 #     params[i] -= learning_rate*params['grad_'+ i]
78 # #B layer
79 # for j in params_b_list:
80 #     params[j] -= learning_rate*params['grad_'+ j]
81
82 #Momentum
83 #just use 'm_'+name variables
84 #W layer
85 for i in params_w_list:
86     params['m_'+i] = 0.9*params['m_'+i] - learning_rate*params['grad_'+i]
87     params[i] += params['m_'+ i]
88 #B layer
89 for j in params_b_list:
90     params['m_'+j] = 0.9*params['m_'+j] - learning_rate*params['grad_'+j]
91     params[j] += params['m_'+ j]
92
93 losses.append(total_loss/train_x.shape[0])

```

Q5.2



For the first 0-20 epochs the average loss shows a sharp decrease before it starts to slow down (near the 20th epoch), then decrease in loss remains steady (i.e., the rate of change begins to plateau and lull).

Q5.3.1



In the reconstructed images, the characters are less visible and more blurred and while there are matches between some of the characters this is not the case for all of them. The network is essentially taking the most notable features of the original character and attempting to reconstruct this in the reconstructed image.

Q5.3.2

15.57999246829877

The PSNR value achieved was around 15.58%.

Q6

- For this homework assignment I opted to do the extra credit question, as I did not have time to finish implementing question 6. I will be implementing after the deadline, for my own learning purposes as it is important for me to know but I unfortunately do not have the time needed to complete this question.