

# Deep Learning for Image Analysis

## DL4IA – Report for Assignment 2

Student Linus Falk

April 12, 2023

### 1 Introduction

Second assignment in the course Deep learning for image analysis

### 2 Mathematical exercises

Given the linear model and multinomial cross-entropy loss:

$$z_i = \mathbf{W}\mathbf{x}_i + \mathbf{b}, \quad i = 1, 2, \dots, n \quad (1)$$

$$L_i = \ln \left( \sum_{l=1}^M e^{z_{il}} \right) - \sum_{m=1}^M \tilde{y}_{im} z_{im}, \quad J = \frac{1}{n} \sum_{i=1}^n L_i \quad (2)$$

Where  $\tilde{y}_{im}$  is the hot encoding of the true label  $y_i$  for data point  $i$

**Exercise 1.1** Derive the expressions for  $\frac{\partial J}{\partial b_m}$  and  $\frac{\partial J}{\partial w_{mj}}$  in terms of:  $\frac{\partial J}{\partial z_{im}}$ ,  $\frac{\partial z_{im}}{\partial b_m}$ ,  $\frac{\partial z_{im}}{\partial w_{mj}}$ :

$$\begin{aligned} \frac{\partial J}{\partial b_m} &= \frac{\partial J}{\partial z_{im}} \frac{\partial z_{im}}{\partial b_m} \\ \frac{\partial J}{\partial w_{mj}} &= \frac{\partial J}{\partial z_{im}} \frac{\partial z_{im}}{\partial w_{mj}} \end{aligned} \quad (3)$$

We start with  $\frac{\partial J}{\partial z_{im}}$

$$\begin{aligned} \frac{\partial J}{\partial z_{im}} &= \frac{\partial}{\partial z_{im}} \ln \left( \sum_{l=1}^M e^{z_{il}} \right) - \tilde{y}_{im} z_{im} = \frac{\partial}{\partial z_{im}} \ln \left( \sum_{l=1}^M e^{z_{il}} \right) - \frac{\partial}{\partial z_{im}} \tilde{y}_{im} z_{im} = \\ &= \frac{e^{z_{im}}}{\sum_{l=1}^M e^{z_{il}}} - \tilde{y}_{im} = \hat{y}_{im} - \tilde{y}_{im} \end{aligned} \quad (4)$$

Where we recognize the fraction as the softmax activation function. We then continue with  $\frac{\partial z_{im}}{\partial b_m}$  and  $\frac{\partial z_{im}}{\partial w_{mj}}$

$$\begin{aligned}\frac{\partial z_{im}}{\partial b_m} &= \frac{\partial}{\partial b_m} \mathbf{W}_m \mathbf{x}_{im} + \mathbf{b}_m = 1 \\ \frac{\partial z_{im}}{\partial w_{mj}} &= \frac{\partial}{\partial b_m} \mathbf{W}_m \mathbf{x}_{im} + \mathbf{b}_m = \mathbf{x}_{im}\end{aligned}\tag{5}$$

### 3 Code exercises

**Exercise 1.2** Using the **numpy** library in **Python**, a feed forward neural network was implemented for solving a classification problem. Following functions were implemented:

1. The *initialize\_parameters* function is used to initialize the weights and offsets by method X. The weight are initialized randomly with mean 0 and variance 0.01 to break the symmetry in the network, making sure that all neurons doesn't learn the same function.

```

1 class NeuralNetwork:
2     def __init__(self, features, learningRate, X_train, Y_train, X_test, Y_test):
3         self.features = features
4         self.learningRate = learningRate
5         self.layers = []
6         self.layersIO = []
7         self.layerGradients = []
8         self.training_history = []
9         self.test_history = []
10        self.accuracy = []
11        self.test_accuracy = []
12        self.iterations = []
13
14        def create_layer(self, nodes, activation):
15            '''
16            param arg 1: Number of nodes in layer
17            param arg 2: Activation functions for layer
18
19            return : Initializes a layer with weights and bias, IO (input/output) list for
20                    keeping track and gradients for backward propagation
21            '''
22            np.random.seed(42)
23            weights = np.random.rand(nodes[1], nodes[0]) * 0.01
24            bias = np.random.rand(nodes[1], 1) * np.sqrt(1 / (nodes[0]))
25            self.layers.append([weights, bias, activation])
26            self.layersIO.append([None, None, None]) # [input, Z, A]
27            self.layerGradients.append([None, None, None, None]) #[dZ, dW, dB, dA]
```

2. Activation functions *relu* and *sigmoid*

```

1     def sigmoid(self, Z):
2         A = 1 / (1 + np.exp(-Z))
3         return A
4
5     def relu(self, Z):
6         return np.maximum(Z, 0)
```

3. Forward propagation with *linear\_forward* computing the linear model of the node, *activation\_forward* for computing the activation after the linear computations and *model\_forward* that combine these and compute the forward pass for the whole model.

```

1     def linear_forward(self, layer, input):
2         '''
3         param arg 1: The layer that will be calculated in forward pass
4         param arg 2: The input to the layer that will be calculated in forward pass
5
6         return: Returns the calculated linear forward pass to the layers IO list
7         '''
8         Z = np.dot(self.layers[layer][0], input) + self.layers[layer][1]
9         self.layersIO[layer][0] = input #input
10        self.layersIO[layer][1] = Z #dotproduct
```

```

11
12 def activation_forward(self, layer):
13     '''
14     param arg 1: the layer to calculate the activation of
15
16     return: Activation of the calculated linear forward pass of the layer
17     '''
18     if self.layers[layer][2] == 'relu':
19         A = self.relu(self.layersIO[layer][1]) #activation of dotproduct
20         self.layersIO[layer][2] = A #save activated dotproduct array
21
22     elif self.layers[layer][2] == 'sigmoid':
23         A = self.sigmoid(self.layersIO[layer][1])
24         self.layersIO[layer][2] = A
25
26     elif self.layers[layer][2] == 'softmax':
27         A = self.softmax(self.layersIO[layer][1])
28         self.layersIO[layer][2] = A
29
30 def model_forward(self, minibatch_X):
31     '''
32     param arg 1: batch for forward pass of the whole model
33
34     return: the updated IO for all layers
35     '''
36     self.linear_forward(0, minibatch_X)
37     self.activation_forward(0)
38     for i in range(1, len(self.layers), 1):
39         self.linear_forward(i, self.layersIO[i-1][2])
40         self.activation_forward(i)
41
42

```

#### 4. A softmax activation function for classification *softmax* and a cross entropy loss function *cross\_entropy*

```

1 def softmax(self, Z):
2     Z_max = np.max(Z, axis=0)
3     Z_shifted = Z - Z_max
4     A = np.exp(Z_shifted) / np.sum(np.exp(Z_shifted), axis=0)
5     return A
6
7 def cross_entropy(self, X, Y):
8     self.model_forward(X)
9     batch_size = Y.shape[1]
10    cost = -np.sum(Y * np.log(self.layersIO[-1][2] + 1e-10)) / batch_size
11    return cost

```

#### 5. Backward propagation functions *linear\_backward* for the linear model, *sigmoid\_backward* and *relu\_backward* for the gradients of the activation functions and a *model\_backward* that computes the backward propagation for the whole model.

```

1 def linear_backward(self, layer, Y):
2     '''
3     param arg 1: the layer that the linear backward pass will be calculated on
4     param arg 2: the true labels for the pass
5
6     return: Calculates the gradients for the layer
7     '''
8     batch_size = self.layersIO[layer][0].shape[1]
9
10    if self.layers[layer][2] == 'softmax':
11        dZ = -Y + self.layersIO[layer][2]
12        dW = (1/batch_size) * np.dot(dZ, self.layersIO[layer-1][2].T)
13        dB = (1/batch_size) * np.reshape(np.sum(dZ,1), (dZ.shape[0],1))
14
15        self.layerGradients[layer][0:3] = dZ, dW, dB
16
17    else:
18        dA = self.layerGradients[layer][3]
19        dZ = np.dot(self.layers[layer+1][0].T, self.layerGradients[layer+1][0]) * dA
20        dW = (1/batch_size) * np.dot(dZ, self.layersIO[layer][0].T)
21        dB = (1/batch_size) * np.reshape(np.sum(dZ,1), (dZ.shape[0],1))
22

```

```

23     self.layerGradients[layer][0:3] = dZ, dW, dB
24
25     def relu_backward(self, Z):
26         return Z > 0
27
28
29     def sigmoid_backward(self, A):
30         return A * (1 - A)
31
32     def activation_backward(self, layer, Y):
33         '''
34         param arg 1: the layer to calculate the activation gradient
35         param arg 2: the true labels
36
37         return: updates the activation gradient
38         '''
39         if self.layers[layer][2] == 'relu':
40             dA = self.relu_backward(self.layersIO[layer][2])
41             self.layerGradients[layer][3] = dA
42
43         elif self.layers[layer][2] == 'sigmoid':
44             dA = self.sigmoid_backward(self.layersIO[layer][2])
45             self.layerGradients[layer][3] = dA
46
47     def model_backward(self, minibatch_Y):
48         '''
49         param arg 1: the true labels for the batch in backward pass
50
51         return: updated gradients for all layers
52         '''
53         for i in range(len(self.layers)-1, -1, -1):
54             self.activation_backward(i, minibatch_Y)
55             self.linear_backward(i, minibatch_Y)
56

```

6. A function *update\_parameters* for taking a step that updates the weight and biases

```

1     def update_parameters(self):
2         '''
3         return: takes a step using the calculated gradients
4         '''
5         for i in range(len(self.layers)):
6             self.layers[i][0] -= self.learningRate * self.layerGradients[i][1]
7             self.layers[i][1] -= self.learningRate * self.layerGradients[i][2]

```

7. A function *predict* for predicting, running the network forward checking the result against the labels

```

1     def predict(self, data, labels):
2         '''
3         param arg 1: data to make classification on
4         param arg 2: the true labels of the data
5
6         return: percentage
7         '''
8         self.model_forward(data.T)
9         percent = self.compare(self.layersIO[-1][2], labels.T)
10        return percent
11
12
13    def compare(self, arr1, arr2):
14        # get maximum values along the "label-array"
15        max_indices_arr1 = np.argmax(arr1, axis=0)
16        max_indices_arr2 = np.argmax(arr2, axis=0)
17
18        # compare and see if the index of max matches hot encoded label
19        matches = np.sum(max_indices_arr1 == max_indices_arr2)
20
21        # percentage of correct matches
22        percentage = (matches / arr1.shape[1]) * 100
23
24        return percentage

```

## 8. Functions to create minibatches:

```

1  def create_mini_batches(data, labels, num_batches):
2
3      # Calculate the batch size
4      batch_size = data.shape[0] // num_batches
5
6      mini_batches = []
7
8      for i in range(num_batches):
9          start_index = i * batch_size
10         end_index = (i + 1) * batch_size
11
12         data_batch = data[start_index:end_index]
13         labels_batch = labels[start_index:end_index]
14
15         mini_batches.append((data_batch, labels_batch))
16
17     # Take the remaining and make a batch
18     if data.shape[0] % batch_size != 0:
19         start_index = num_batches * batch_size
20         data_batch = data[start_index:]
21         labels_batch = labels[start_index:]
22
23         mini_batches.append((data_batch, labels_batch))
24
25     return mini_batches

```

## 9. A function *train\_model* to train the whole model with the mini batches

```

1  def train_model(self, mini_batches, epochs):
2      iter = 0
3      k = 400
4      for y in range(epochs):
5          print(str(y+1) + ' out of ' + str(epochs) + ' epochs')
6          for x in range(len(mini_batches)):
7              self.model.forward(mini_batches[x][0].T)
8              self.model.backward(mini_batches[x][1].T)
9              self.update_parameters()
10             if x % k == 0:
11                 self.training_history.append(self.cross_entropy(mini_batches[x]
12 ] [0].T, mini_batches[x][1].T))
13                 self.test_history.append(self.cross_entropy(X_test.T, Y_test.T))
14                 self.test_accuracy.append(nn.predict(mini_batches[x][0],
15 mini_batches[x][1]))
16                 self.accuracy.append(nn.predict(X_train, Y_train))
17                 iter += k
18                 self.iterations.append(iter)
19
20         print(nn.predict(X_test, Y_test)) #print test accuracy during training

```

## 10. Main function, creating the network, train it and plot the training history:

```

1  # Create neural network
2  nn = NeuralNetwork(784, 1e-2, X_train, Y_train, X_test, Y_test)
3  nn.create_layer([784, 128], 'relu')
4  nn.create_layer([128, 64], 'relu')
5  nn.create_layer([64, 10], 'softmax')
6  nn.print_layer()
7
8  #Train the network and print test accuracy while training
9  epochs = 50
10 nn.train_model(mini_batches, epochs)
11 print(nn.predict(shuffled_testdata, shuffled_testlabels))
12
13
14 #Plot training cost and accuracy history after training
15 fig, (ax1, ax2) = plt.subplots(1, 2)
16 fig.set_figwidth(10)
17
18 ax1.set_title('Cost')
19 ax1.plot(nn.training_history, label='training cost')
20 ax1.plot(nn.test_history, label='test cost')
21 ax1.grid()

```

```

22 ax1.legend()
23
24 ax2.set_title('Accuracy')
25 ax2.set_xlabel('Iterations')
26 ax2.plot(nn.accuracy, label= 'accuracy')
27 ax2.plot(nn.test_accuracy, label= 'test accuracy')
28 ax2.grid()
29 ax2.legend()
30
31 plt.show()

```

### Exercise 1.3 & Exercise 1.4

Table 1: Sample table of results

Model	Activation	Accuracy (%)
1, [784, 10]	Softmax	92.06
2, [784, 128, 64, 10]	ReLU, ReLU, Softmax	97.27
3, [784, 128, 64, 10]	Sigmoid, Sigmoid, Softmax	89.64

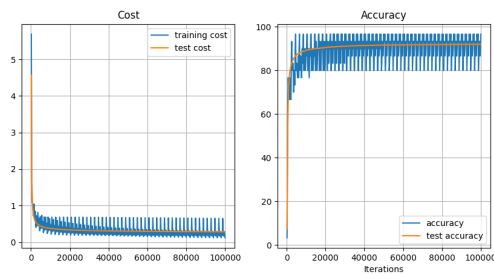


Figure 1: Training history, model 1

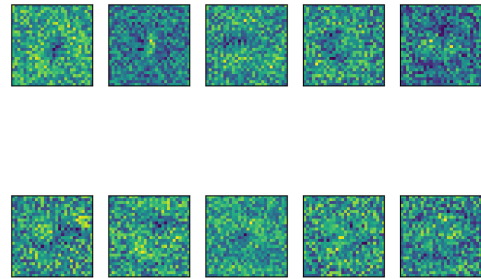


Figure 2: Reshaped weight matrix, model 1

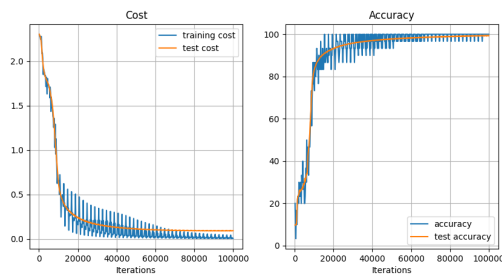


Figure 3: Training history, model 2

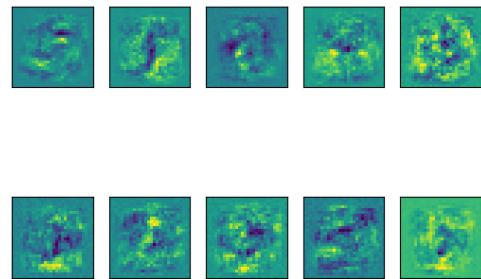


Figure 4: Reshaped weight matrix, model 2

We can see that the case of the one layer network, the network is not complex enough when the train accuracy never get close to 100% accuracy even though the batch size was really small, 30 images in each batch. In model with three layers we can see that the gap between the cost for training and test starts to

open up that can indicate over training. We also see that the network is complex enough as the training accuracy can hit 100% with this batch size (30 images). The model with the Sigmoid activation function was much harder to train and suffered a bit from the initialization not being optimal for this activation function. This lead to slow learning and poor accuracy.

Taking a look on the images of the weight we can see in the one layer model (if we squint with our eyes a bit) that they resemble the digits: 0, 1, 2 and 3. The rest of them are harder to see. This is not so surprising since we try to achieve a high softmax score for the correct digit, by taking the dot product with a similar digit will result in a mean value of the matrix. Quite similar to using the singular value decomposition method to classify the MNIST data set. In the case with the multilayer models it is harder to make out of how the model have tuned the weights since there are multiple layers after with each other.

## References

- [1] Asimov, Issac (1942). Runaround