

## Submission Assignment #1

Coordinator: Jakub Tomczak

Name: Bartłomiej Boczek, Netid: bbk205

## 1 Question answers

**Question 1** Work out the local derivatives of both, in scalar terms. Show the derivation. Assume that the target class is given as an integer value.

In this section I will work out the derivation of the **cross entropy loss** function. The loss is defined as follows:

$$loss = \sum_i l_i$$

$$l_i = \begin{cases} -\log y_c & \text{if } c = i \\ 0 & \text{otherwise} \end{cases}$$

Work out the derivative:

$$\frac{\partial l}{\partial y_i} = \begin{cases} -\frac{1}{y_c} & \text{if } c = i \\ 0 & \text{otherwise} \end{cases}$$

In the following section derivation of the **softmax function** will be presented. The softmax function is defined as:

$$y_i = \frac{e^{o_i}}{\sum_{j=0}^{n_{class}} e^{o_j}}$$

Because of the sum in the softmax formula,  $y_i$  depends on all  $o_j$ , not just  $o_i$ , it is helpful to work out and separately cases where  $i = j$  and  $i \neq j$ .

$$\frac{\partial y_i}{\partial o_j} = \frac{\partial \frac{e^{o_i}}{\sum_{j=0}^{n_{class}} e^{o_j}}}{\partial o_j}$$

*Quotient rule* is needed to calculate the derivative of softmax, because it is defined as a fraction.

- $i = j$

$$\frac{\partial y_i}{\partial o_j} = \frac{\partial \frac{e^{o_i}}{\sum_{j=0}^{n_{class}} e^{o_j}}}{\partial o_j} = \frac{e^{o_i} * \sum_{j=0}^{n_{class}} e^{o_j} - e^{o_j} * e^{o_i}}{[\sum_{j=0}^{n_{class}} e^{o_j}]^2} = \frac{e^{o_i}}{\sum_{j=0}^{n_{class}} e^{o_j}} \frac{(\sum_{j=0}^{n_{class}} e^{o_j} - e^{o_j})}{\sum_{j=0}^{n_{class}} e^{o_j}} = y_i(1 - y_i)$$

- $i \neq j$

In this case the derivative  $\frac{\partial e^{o_i}}{\partial o_j} = 0$ , thus the equation looks like this:

$$\frac{\partial y_i}{\partial o_j} = \frac{\partial \frac{e^{o_i}}{\sum_{j=0}^{n_{class}} e^{o_j}}}{\partial o_j} = \frac{0 * \sum_{j=0}^{n_{class}} e^{o_j} - e^{o_j} * e^{o_i}}{[\sum_{j=0}^{n_{class}} e^{o_j}]^2} = \frac{-e^{o_j} * e^{o_i}}{[\sum_{j=0}^{n_{class}} e^{o_j}]^2} = -y_j y_i$$

Finally the equation for the local derivative of the softmax function looks as follows:

$$\frac{\partial y_i}{\partial o_j} = \begin{cases} y_i(1 - y_i) & \text{if } i = j \\ -y_j y_i & \text{if } i \neq j \end{cases}$$

**Question 2\*2** Work out the derivative  $\frac{\partial l}{\partial o_i}$ . Why is this not strictly necessary for a neural2network, if we already have the two derivatives we worked out above?

**Question 3** Implement the network drawn in the image below, including the weights. Perform one forward pass, up to the loss on the target value, and one backward pass. Show the relevant code in your report. Report the derivatives on all weights (including biases). Do not use anything more than plain python and the math package.

Below there is a code that implements forward calculation of forward and backward pass from the image. Weights of layers are stored in *W1* and *W2* lists, biases in *b1* and *b2* lists. *X* is the input lists and *Y* is the output list. All has been implemented using basic math library, in scalar terms. At the end, resulting gradients are printed. In the results *V=W2*, *W=W1*, *b=b1* and *c=b2* to keep consistency with the assignment text.

```

1 import math
2 from typing import List
3
4 X = [1.0, -1.0]
5 Y = [1, 0]
6 b1 = [0.0, 0.0, 0.0]
7 W1 = [[1.0, 1.0, 1.0], [-1.0, -1.0, -1.0]]
8 b2 = [0.0, 0.0]
9 W2 = [[1.0, 1.0], [-1.0, -1.0], [-1.0, -1.0]]
10
11 n_input=2
12 n_hidden = 3
13 n_out = 2
14
15 def sigmoid_scalar(Z: List[float]) -> List[float]:
16     return [1 / (1 + math.exp(-x)) for x in Z]
17
18 def softmax_scalar(Z: List[float]) -> List[float]:
19     return [math.exp(oi) / sum([math.exp(oj) for oj in Z]) for oi in Z]
20
21 # Forward pass
22 # Linear transformation with hidden layer
23 Z1 = [None] * n_hidden
24 for j in range(n_hidden):
25     Z1[j] = sum(X[i] * W1[i][j] for i in range(n_input)) + b1[j]
26
27 # Activations for hidden layer
28 A1 = sigmoid_scalar(Z1)
29
30 # Linear transformation with output layer
31 Z2 = [None] * n_out
32 for j in range(n_out):
33     Z2[j] = sum(A1[i] * W2[i][j] for i in range(n_hidden)) + b2[j]
34
35 # Activations for output layer
36 A2 = softmax_scalar(Z2)
37
38 # Calculate loss
39 trg_idx = Y.index(1)
40 yc = A2[trg_idx]
41 loss = -math.log(yc)
42
43 # Backward pass
44 # Gradient of activated output layer
45 dZ2 = [A2[i] - Y[i] for i in range(n_out)]
46
47 # Gradient of second layer weights
48 dW2 = [[None]*n_out for _ in range(n_hidden)]

```

```

49 for j in range(n_out):
50     for i in range(n_hidden):
51         dW2[i][j] = dZ2[j]*A1[i]
52
53 # Gradient of second layer bias
54 db2 = dZ2
55
56 # Gradient of first activated output
57 dZ1 = [None]*n_hidden
58 for j in range(n_hidden):
59     dZ1[j] = A1[j] * (1 - A1[j]) * sum(dZ2[i]*W2[j][i] for i in range(n_out))
60
61 # Gradient of first layer weights
62 dW1 = [[None]*n_hidden for _ in range(n_input)]
63 for j in range(n_hidden):
64     for i in range(n_input):
65         dW1[i][j] = dZ1[j]*X[i]
66
67 # Gradient of first layer bias
68 db1 = dZ1
69
70 print(f'dW = {dW1}')
71 print(f'db = {db1}')
72 print(f'dV = {dW2}')
73 print(f'dc = {db2}')

```

The resulting derivatives are as follows:

```

1 dW = [[0.0, 0.0, 0.0], [-0.0, -0.0, -0.0]]
2 db = [0.0, 0.0, 0.0]
3 dV = [[-0.44039853898894116, 0.44039853898894116], [-0.44039853898894116,
4         0.44039853898894116], [-0.44039853898894116, 0.44039853898894116]]
5 dc = [-0.5, 0.5]

```

**Question 4** *Implement a training loop for your network and show that the training loss drops as training progresses.*

## 1.1 Dataset

For training synthetic dataset provided in the assignment is used. It contains 2 classes, which distribution is presented in the Figure 1

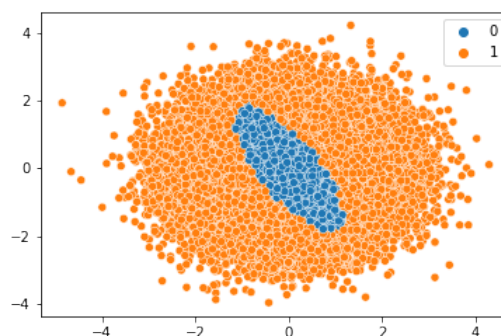


Figure 1: Synthetic data distribution

## 1.2 Implementation

To implement a training loop additional MLP python class has been added. It stores weights and biases, but also caches activations, input and output used in forward method, so they can be used in backward pass and parameter updates. The interface with the class is very simple. ‘forward’ method returns a loss of

forward pass. ‘backward’ calculates gradients and updates neural network weights. In my implementation I use one-hot encoding - i.e. the output vector is binary and has 1 at the position of true class. Weights are initialized with normal distribution and it uses stochastic gradient descent - loss is calculated over one instance at a time.

```

1 import math
2 from random import gauss
3 from datasets import load_synth
4 from typing import List
5
6 def sigmoid_scalar(Z: List[float]) -> List[float]:
7     return [1 / (1 + math.exp(-x)) for x in Z]
8
9 def softmax_scalar(Z: List[float]) -> List[float]:
10    return [math.exp(oi) / sum([math.exp(oj) for oj in Z]) for oi in Z]
11
12 class MLP:
13     def __init__(self, n_input, n_hidden, n_out, lr) -> None:
14         self.n_input, self.n_hidden, self.n_out = n_input, n_hidden, n_out
15         self.lr = lr
16
17         self.b1 = [0.0 for _ in range(n_hidden)]
18         self.W1 = [[gauss(0,1) for _ in range(n_hidden)] for _ in range(n_input)]
19         self.b2 = [0.0 for _ in range(n_out)]
20         self.W2 = [[gauss(0,1) for _ in range(n_out)] for _ in range(n_hidden)]
21
22         self.A1 = None
23         self.A2 = None
24         self.X = None
25         self.Y = None
26
27     def forward(self, X, Y):
28         Z1 = [None] * self.n_hidden
29         for j in range(self.n_hidden):
30             Z1[j] = sum(X[i] * self.W1[i][j] for i in range(self.n_input)) + self.b1[j]
31
32         A1 = sigmoid_scalar(Z1)
33
34         # Linear transformation with output layer
35         Z2 = [None] * self.n_out
36         for j in range(self.n_out):
37             Z2[j] = sum(A1[i] * self.W2[i][j] for i in range(self.n_hidden)) + self.b2[j]
38
39         # Activations for output layer
40         A2 = softmax_scalar(Z2)
41
42         # Calculate loss
43         trg_idx = Y.index(1)
44         yc = A2[trg_idx]
45         loss = -math.log(yc)
46
47         # Cache for backward
48         self.X = X
49         self.Y = Y
50         self.A1 = A1
51         self.A2 = A2
52         return loss
53
54     def backward_and_update(self):
55         # Backward pass
56         # Gradient of activated output layer
57         dZ2 = [self.A2[i] - self.Y[i] for i in range(self.n_out)]
58
59         # Gradient of second layer weights
60         dW2 = [[None]*self.n_out for _ in range(self.n_hidden)]
61         for j in range(self.n_out):
62             for i in range(self.n_hidden):
63                 dW2[i][j] = dZ2[j]*self.A1[i]
64
65         # Gradient of second layer bias
66         db2 = dZ2
67
68         # Gradient of first activated output

```

```

69     dZ1 = [None]*self.n_hidden
70     for j in range(self.n_hidden):
71         dZ1[j] = self.A1[j] * (1 - self.A1[j]) * sum(dZ2[i]*self.W2[j][i] for i in
range(self.n_out))
72
73     # Gradient of first layer weights
74     dW1 = [[None]*self.n_hidden for _ in range(self.n_input)]
75     for j in range(self.n_hidden):
76         for i in range(self.n_input):
77             dW1[i][j] = dZ1[j]*self.X[i]
78
79     # Gradient of first layer bias
80     db1 = dZ1
81
82     for j in range(self.n_hidden):
83         for i in range(self.n_input):
84             self.W1[i][j] -= self.lr * dW1[i][j]
85             self.b1[j] -= self.lr * db1[j]
86
87     for j in range(self.n_out):
88         for i in range(self.n_hidden):
89             self.W2[i][j] -= self.lr * dW2[i][j]
90             self.b2[j] -= self.lr * db2[j]
91
92
93 n_input=2
94 n_hidden = 3
95 n_out = 2
96 epochs = 5
97 lr = 0.001
98
99 (xtrain, ytrain), (xval, yval), num_cls = load_synth()
100
101 n_train = len(xtrain)
102 one_hots = [[0]* n_input for _ in range(n_train)]
103 for i, yi in enumerate(ytrain):
104     one_hots[i][yi] = 1
105
106 mlp = MLP(n_input, n_hidden, n_out, lr)
107
108 losses = list()
109 for epoch in range(epochs):
110     epoch_losses = list()
111     for x, y in zip(xtrain, one_hots):
112         loss = mlp.forward(x,y)
113         epoch_losses.append(loss)
114         mlp.backward_and_update()
115     avg_loss = sum(epoch_losses)/len(epoch_losses)
116     print(f'Epoch: {epoch}, loss: {avg_loss}')
117     losses.append(avg_loss)

```

As can be observe in Figure 2 the training loss drops from epoch to epoch.

**Question 5** *Implement a neural network for the MNIST data. Use two linear layers as before, with a hidden layer size of 300, a sigmoid activation, and a softmax activation over the output layer, which has size 10.*

### 1.3 Dataset

For this task MNIST dataset has been used. Dataset contains 70000 samples with handwritten digits. In the experiments a network for digits classification will be trained. 60000 are in the training/validation sets and 10000 are used for testing. Data is randomly shuffled before training, both in batched and SGD cases.

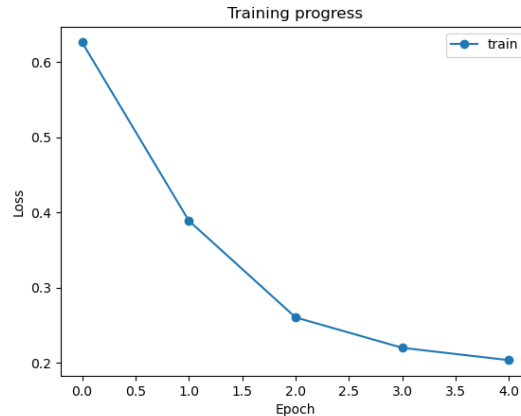


Figure 2: Training progress of scalar version of Neural Network

## 1.4 Implementation

In this section I will present the solution to transform scalar version from previous paragraph into vectorized version that works using *numpy* library. In *numpy* matrix operations are well optimized. What is more mathematical functions can be applied to whole matrices, without the need to do it element-wise and using for loops.

Similarly to previous implementations, in *forward* pass inputs, outputs and activations are caches, so they can be used in *backward* pass. What is more, minibatch gradient descent is implemented by accumulation of gradients from multiple iterations. When update of parameters is performed, average gradients from the minibatch are used. Training loop looks the same as in the previous paragraph.

Training progress can be seen in Figure 3. As we can see training loss is going down from epoch to epoch. Training is stable, because mini-batch of 8 was used.

```

1 import numpy as np
2
3 def sigmoid(x: np.ndarray) -> np.ndarray:
4     return 1 / (1 + np.exp(-x))
5
6 def softmax(x: np.ndarray) -> np.ndarray:
7     exp = np.exp(x)
8     return exp / exp.sum()
9
10 class MLP:
11     def __init__(self, n_input: int, n_hidden: int, n_out: int, lr: float, minibatch: int) -> None:
12         self.n_input, self.n_hidden, self.n_out = n_input, n_hidden, n_out
13         self.lr = lr
14
15         self.b1 = np.zeros(n_hidden)
16         self.W1 = np.random.normal(size=(n_input, n_hidden))
17         self.b2 = np.zeros(n_out)
18         self.W2 = np.random.normal(size=(n_hidden, n_out))
19
20         self.minibatch = minibatch
21         self.batch_counter = 0
22
23         # minibatch gradient accumulators
24         self.dW1 = 0
25         self.dW2 = 0
26         self.db1 = 0
27         self.db2 = 0
28
29         # cache
30         self.A1 = None
31         self.A2 = None
32         self.X = None
33         self.Y = None
34

```

```

35     def forward(self, X: np.ndarray, Y: np.ndarray) -> float:
36         Z1 = X.dot(self.W1) + self.b1
37         self.A1 = sigmoid(Z1)
38         Z2 = self.A1.dot(self.W2) + self.b2
39         self.A2 = softmax(Z2)
40
41         loss = np.sum(-Y * np.log(self.A2))
42
43         self.X = X
44         self.Y = Y
45         return loss
46
47     def backward_and_update(self):
48         dZ2 = self.A2 - self.Y
49         dW2 = self.A1[:, np.newaxis].dot(dZ2[np.newaxis, :])
50         db2 = dZ2
51
52         dZ1 = self.A1 * (1 - self.A1) * self.W2.dot(dZ2)
53         dW1 = self.X[:, np.newaxis].dot(dZ1[np.newaxis, :])
54         db1 = dZ1
55
56         self.dW1 += dW1
57         self.dW2 += dW2
58         self.db1 += db1
59         self.db2 += db2
60         self.batch_counter += 1
61         if self.batch_counter % self.minibatch == 0:
62             # Perform update
63             self.W1 -= self.lr * self.dW1 / self.minibatch
64             self.W2 -= self.lr * self.dW2 / self.minibatch
65             self.b1 -= self.lr * self.db1 / self.minibatch
66             self.b2 -= self.lr * self.db2 / self.minibatch
67
68             self.zero_grad()
69
70     def zero_grad(self):
71         self.dW1 = 0
72         self.dW2 = 0
73         self.db1 = 0
74         self.db2 = 0

```

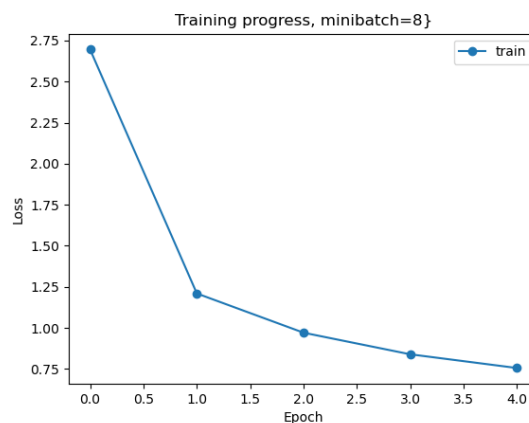


Figure 3: Training progress of vectorized version of Neural Network trained on MNIST with minibatch=8

**Question6** Work out the vectorized version of a batched forward and backward. That is, work out how you can feed your network a batch of images in a single 3-tensor, and still perform each multiplication by a weight matrix, the addition of a bias vector, computation of the loss, etc. in a single numpy call.

To implement batched forward and backward descent one needs to decide the dimension for batch. I decided to make the first dimension of the 3-d tensor, a *batch dimension*, thus the input vector has the size of *batch* x 784. The first dimension means batch in every tensor operation during forward and backward pass. Update is performed on average gradient of the batch.

Below crucial set of code is presented. Tensors can be multiplied with `np.matmul` for tensor operations. During calculation of  $dZ1$ , reduction of the last excessive dimension is needed with `squeeze(-1)` method.

```

1  def forward(self, X: np.ndarray, Y: np.ndarray) -> np.ndarray:
2      Z1 = np.dot(X, self.W1) + self.b1
3      self.A1 = sigmoid(Z1)
4      Z2 = np.dot(self.A1, self.W2) + self.b2
5      self.A2 = softmax(Z2)
6      loss = np.sum(-Y * np.log(self.A2), axis=0)
7      self.X = X
8      self.Y = Y
9      return loss
10
11  def backward_and_update(self):
12      dZ2 = self.A2 - self.Y
13      dW2 = np.matmul(self.A1[:, :, np.newaxis], dZ2[:, np.newaxis, :])
14      db2 = dZ2
15      dZ1 = self.A1 * (1 - self.A1) * np.matmul(self.W2[np.newaxis, :, :], dZ2[:, :, np.newaxis]).squeeze(-1)
16      dW1 = np.matmul(self.X[:, :, np.newaxis], dZ1[:, np.newaxis, :])
17      db1 = dZ1
18      self.W1 -= self.lr * dW1.mean(axis=0)
19      self.W2 -= self.lr * dW2.mean(axis=0)
20      self.b1 -= self.lr * db1.mean(axis=0)
21      self.b2 -= self.lr * db2.mean(axis=0)

```

**Question 7** Train the network on MNIST and plot the loss of each batch or instance against the timestep. This is called a learning curve or a loss curve. You can achieve this easily enough with a library like matplotlib in a jupyter notebook, or you can install a specialized tool like tensorboard. We'll leave that up to you.

**1** Compare the training loss per epoch to the validation loss per epoch. What does the difference tell you?

In this experiment dataset is split into 2 groups - validation and training. Model is trained only on the training part and that, after each epoch, forward pass is performed in the validation part to compute loss and accuracy of predictions. From the plot we can see that validation loss is lower than training loss after even 15 epochs. Train accuracy is lower than validation accuracy. It might suggest that validation data is not a good representation of the problem and might contains too easy examples, that are not challenging enough for the model.

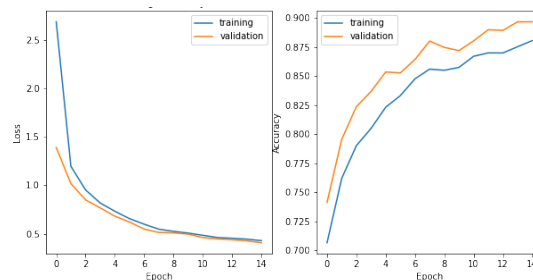


Figure 4: Training vs validation loss and accuracy

**2** Run the SGD method multiple times (at least 3) and plot an average and a standard deviation of the objective value in each iteration (e.g., see :here). What does this tell you?

With set learning rate to 0.001 and number of epochs to 5, the same experiment but with differently initialized weights was repeated 3 times to determine how stable is the SGD algorithm. As we can see on the graph below, SGD behaves very stable.

**3** Run the SGD with different learning rates (e.g., 0.001, 0.003, 0.01, 0.03). Analyze how the learning rate value influences the final performance.

In this experiment training was run one time for different values of learning rate, namely [0.001, 0.003, 0.01, 0.03]. We can see in the Figure 6, that the best learning rate in terms training loss was 0.003 and worse was 0.03.



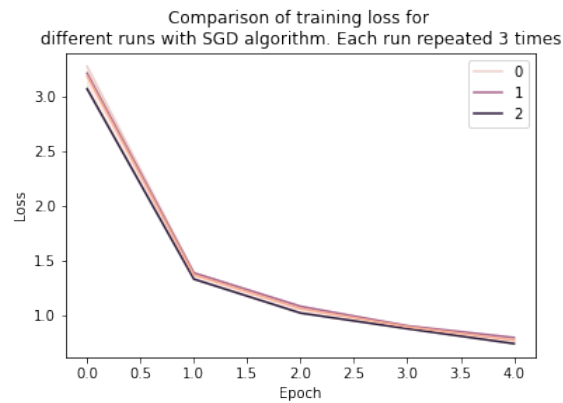


Figure 5: Training loss from 3 SGD runs

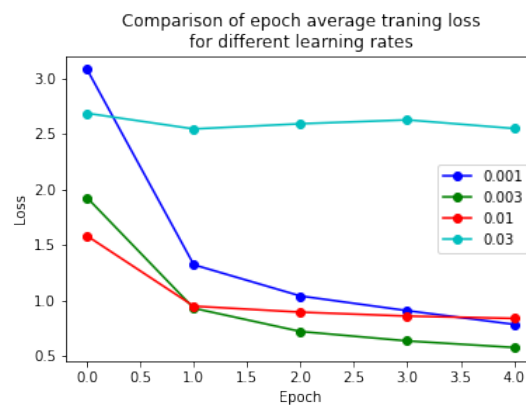


Figure 6: Training loss from 4 different learning rates

4 Based on these experiments, choose a final set of hyperparameters, load the full training data with the canonical test set, train your model with the chosen hyperparameters and report the accuracy you get. For the final experiment 0.003 learning rate was used. The model was trained with minibatch=8 for 10 epochs. As we can see after 10 epochs loss is not dropping significantly. After training model was evaluated on the testset. We can see from classification report (Figure 8), that the model performs best in classification of 1 and 6, and worst in detecting 5 and 8.

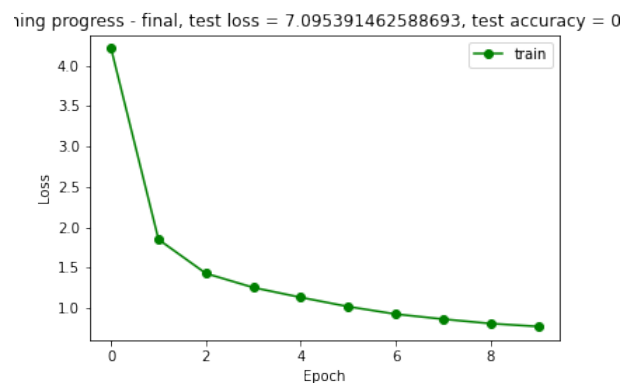


Figure 7: Training loss final experiment.

	precision	recall	f1-score	support
0	0.892604	0.898980	0.895780	980.0000
1	0.949513	0.944493	0.946996	1135.0000
2	0.823587	0.818798	0.821186	1032.0000
3	0.800801	0.792079	0.796416	1010.0000
4	0.879824	0.812627	0.844891	982.0000
5	0.773893	0.744395	0.758857	892.0000
6	0.884298	0.893528	0.888889	958.0000
7	0.840160	0.818093	0.828980	1028.0000
8	0.748538	0.788501	0.768000	974.0000
9	0.757052	0.824579	0.789374	1009.0000
accuracy	0.835700	0.835700	0.835700	0.8357
macro avg	0.835027	0.833607	0.833937	10000.0000
weighted avg	0.836928	0.835700	0.835939	10000.0000

Figure 8: Classification report on the testset

## References

## Appendix