

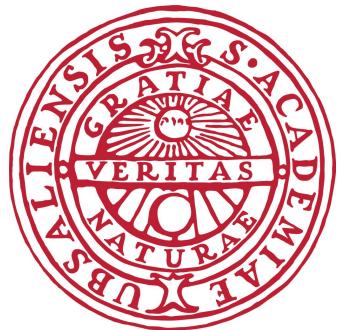
UPPSALA UNIVERSITY

DEEP LEARNING

Hand-in assignment (2)

Tong You

May 10, 2021



Contents

1	Introduction	2
2	PyTorch and MNIST	2
2.1	Fully-connected networks	2
2.2	Convolutional neural networks	5
3	PyTorch and semantic segmentation	14
4	Appendix	20
4.1	PyTorch and MNIST	20
4.2	PyTorch and semantic segmentation	30

1 Introduction

In the previous assignment we had to code a softmax classifier, mini-batched softmax classifier, and a neural network from scratch for the MNIST dataset. For the first part of this assignment we will return to the MNIST dataset, but this time we will be using PyTorch instead. This will also give us CUDA support which will drastically speed-up training. All results will be done on an RTX 2080ti. I will recreate the 2-layer and 4-layer neural networks from assignment 1 in PyTorch instead using hyper parameters that are similar so a one-to-one comparison can be made. We will then explore a convolutional neural network (CNN) for MNIST and try to gain as high accuracy as possible. In the second part of the assignment we will be looking at semantic segmentation on biomedical images, more specifically the Warwick biomedical dataset. For this we will adapt our CNN for MNIST to a fully-convolutional network (FCN). All code is attached to the Appendix.

2 PyTorch and MNIST

2.1 Fully-connected networks

As mentioned in the previous section, in this section we will be exploring MNIST using PyTorch as our framework. To start with we will re-implement the 2-layer and 4-layer fully-connected networks first. This will allow us to establish a baseline against to compare the CNN and various architectural changes, regularization, and different optimization methods. In the previous assignment we used a mini-batch size of 1500, so we will use this batch-size here as well. We also used quite a high learning rate 0.8 in the previous assignment. This was all with plain SGD and since the 4-layer network design (100 hidden units in all hidden layers) was not optimal, we might expect similar instability in the PyTorch implementation. We will be using ReLU as an activation function for both of these networks. The loss and accuracy plots are shown in Figure 1.

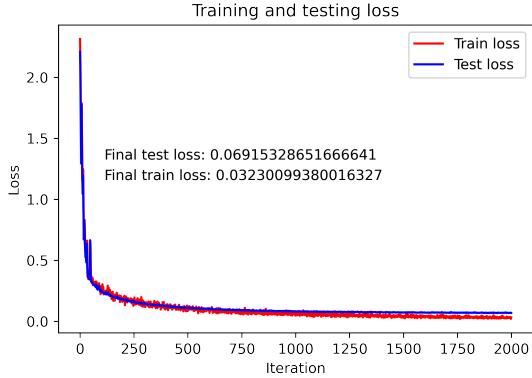
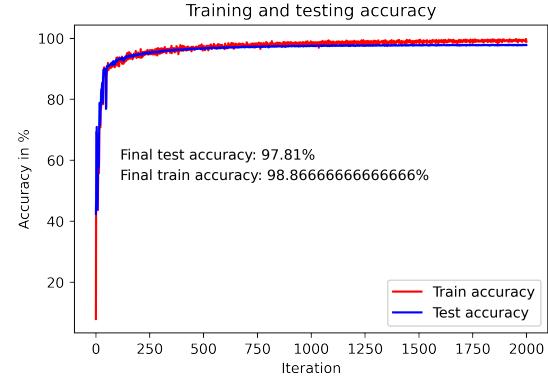
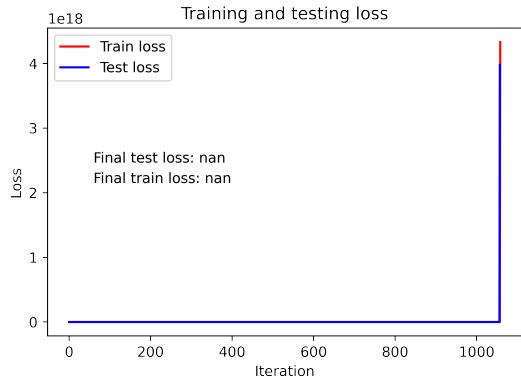
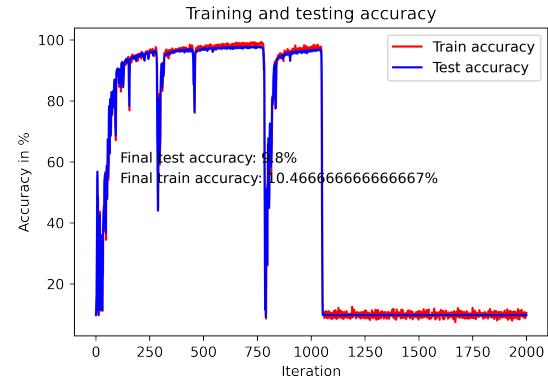
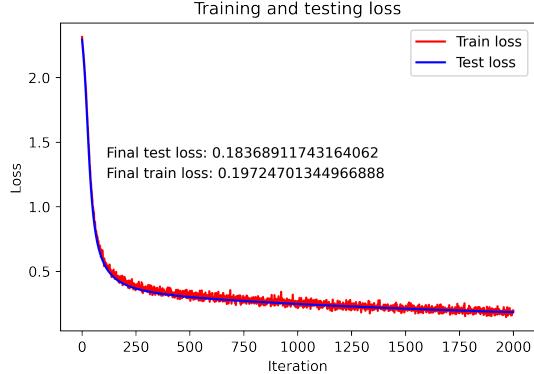
(a) $\gamma = 0.8$ for 2-layer network(b) $\gamma = 0.8$ for 2-layer network(c) $\gamma = 0.8$ for 4-layer network(d) $\gamma = 0.8$ for 4-layer network

Figure 1: Accuracy and loss plots for 2-layer and 4-layer networks. Note that for these figures we don't plot the accuracy and loss plots every k-th iteration, we plot the value every iteration.

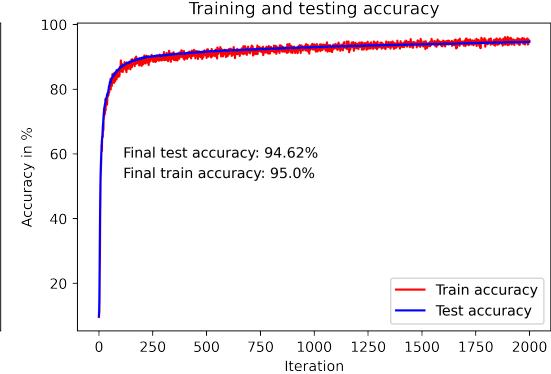
There are several interesting things to observe. We note first of all that the 2-layer network gains a final testing accuracy of 97.81% and the loss and accuracy plots converge in a stable way. For the 4-layer network however, we can see that the network is not stable. In fact the training and testing loss diverge near the end. The large increase in training/testing loss after 1000 iterations corresponds to a corresponding decrease in accuracy to 9.8%. However, it is quite apparent that for the earlier iterations, the accuracy does

reach a value of about 98%. Therefore to investigate whether a lower learning rate would fix such a bug we train both the 2-layer and 4-layer networks with a learning rate of 0.1. The results are shown in Figure 2

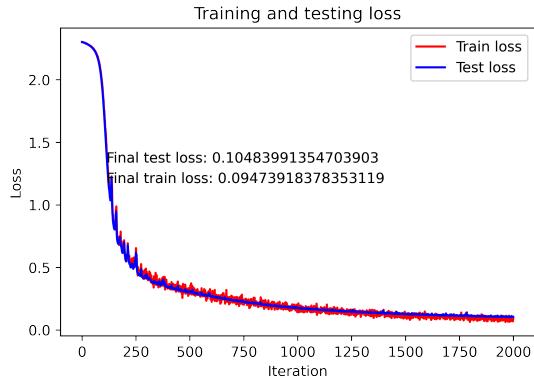
(a) $\gamma = 0.1$ for 2-layer network



(b) $\gamma = 0.1$ for 2-layer network



(c) $\gamma = 0.1$ for 4-layer network



(d) $\gamma = 0.1$ for 4-layer network

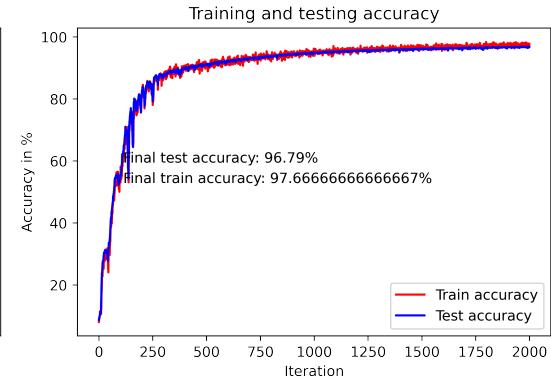


Figure 2: Accuracy and loss plots for 2-layer and 4-layer networks. Note that for these figures we don't plot the accuracy and loss plots every k -th iteration, we plot the value every iteration. The only difference between these plots and the ones in the previous figure are that we change the learning rate to 0.1.

We can see that the 2-layer network again displays stable convergence of the

training/testing loss and accuracies. We do note that the loss and accuracy plots are more stable compared to the 2-layer network with a learning rate of 0.8. The final testing accuracy is also slightly lower at 94.62%. The biggest difference is that the 4-layer network displays much more stable convergence of the loss and accuracies. There is also no huge blow-up anymore. In fact, the 4-layer network converges to a larger final testing accuracy of 96.79% than the 2-layer network. Since the number of hidden units and therefore the width of the networks are the same, we can see that by adding more layers we do get final higher testing accuracy when keeping all other parameters the same. Next we will look how high accuracy we can get by swapping the fully-connected architecture for a convolutional one.

2.2 Convolutional neural networks

We will now explore the MNIST dataset in terms of CNNs. We will start by implementing the following network:

```

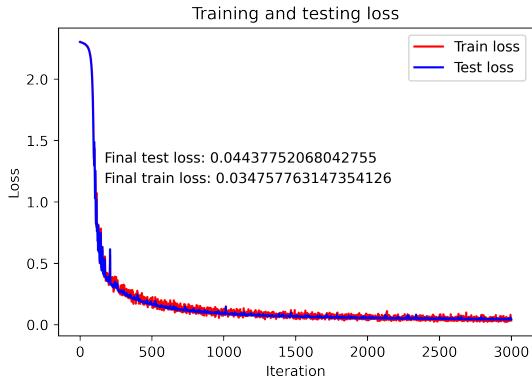
1 conv_mod = nn.Sequential(nn.Conv2d(in_channels = 1,
2     ↳ out_channels = 8, kernel_size = (3,3), stride = 1, padding
2     ↳ = 1), nn.ReLU(), nn.MaxPool2d(kernel_size = 2, stride =
2     ↳ 2),
2
2     ↳ nn.Conv2d(in_channels = 8,
2     ↳     ↳ out_channels = 16, kernel_size =
2     ↳     ↳ (3,3), stride = 1, padding = 1),
2     ↳     ↳ nn.ReLU(),
2     ↳     ↳ nn.MaxPool2d(kernel_size = 2,
2     ↳     ↳ stride = 2),
2
3     ↳ nn.Conv2d(in_channels = 16,
3     ↳     ↳ out_channels = 32, kernel_size =
3     ↳     ↳ (3,3), stride = 1, padding = 1),
3     ↳     ↳ nn.ReLU(), nn.Flatten(),
3     ↳     ↳ nn.Linear(32 * 7 * 7, M))

```

We will first train this network using plain SGD and look at the final accuracy. In order to ensure a good trade-off between convergence and waiting too long we will be using a [batch-size of 1000](#). We know that a small batch-size means more iterations per epoch, since we need more passes to sample the entire training set. The opposite holds for larger batch-sizes. It is also well-known

that smaller batch-sizes lead to noisier gradient updates, which also needs to be compensated for by changing the learning rate for example. We end up choosing a learning rate of 0.05. The results are shown in Figure 3.

(a) $\gamma = 0.05$ using SGD



(b) $\gamma = 0.05$ using SGD

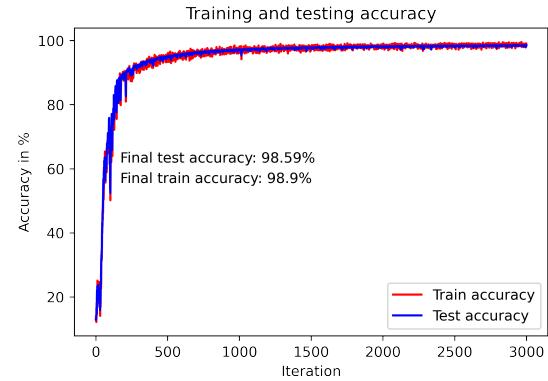


Figure 3: Accuracy and loss plots for CNN. The network is trained with plain SGD.

We can see that the CNN gains a final accuracy of 98.59%. This is higher than the final accuracies reached by the 2-layer and 4-layer networks. Since we take into account spatial relations by not flattening the images into vectors, for a network with optimal settings we should expect to gain higher accuracies. This was all done using plain SGD. What will happen if we choose a different optimizer, while keeping all the other parameters the same. We will investigate this by using [ADAM](#) (Adaptive Moment Estimation). For this we will be using the default parameters suggested by PyTorch, and the code looks like:

```
1 adam_opt = optim.Adam(mod.parameters(), lr = 0.001 , betas =
  ↵ (0.9, 0.999), eps = 1e-08, weight_decay = 0.0, amsgrad =
  ↵ False)
```

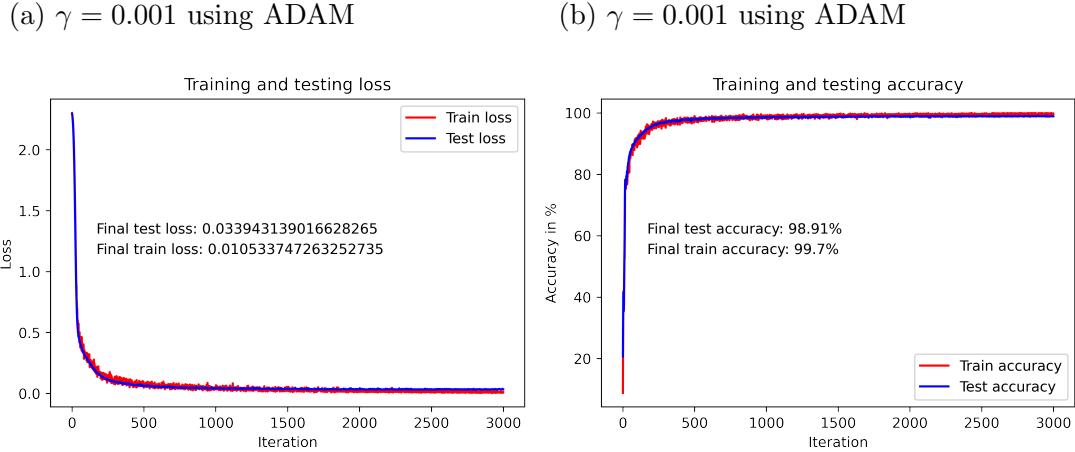


Figure 4: Accuracy and loss plots for CNN. The network is trained with ADAM with default values as suggested by PyTorch.

As we can see, there is a big difference between the CNN trained with plain SGD and the CNN trained with ADAM. For the same batch-size, we can see that training/testing loss and accuracy plots converge in a much more stable way. The final testing accuracy is slightly higher than the CNN optimized with plain SGD: the CNN trained with ADAM gets about 98.91%.

Next we will look at what happens when we swap the order of the activation and pooling layers. Since the pooling layer reduces the dimensionality after application, we will perform a ReLU function on a smaller representation. This should mean that the ReLU should perform faster. Since we don't change anything else, we don't expect a large change in accuracy (except for differences in weight initialisation). The results are shown in Figure 5. Note that we train with SGD since we want to compare to the CNN trained with SGD.

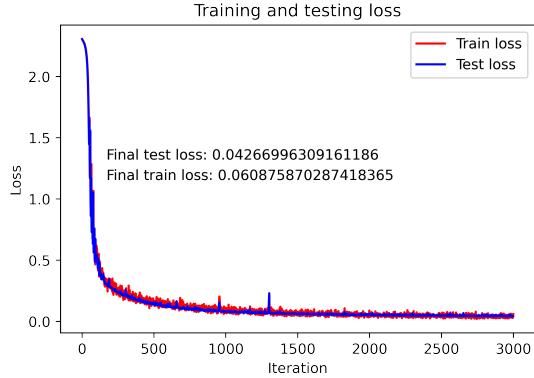
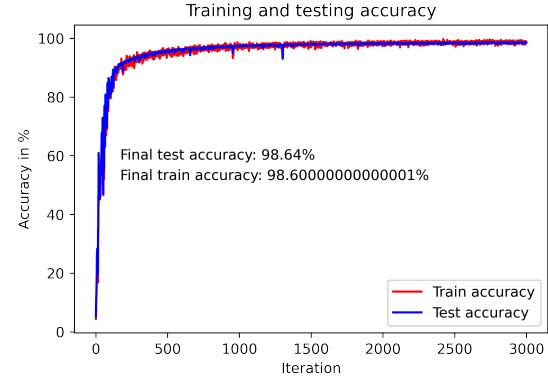
(a) $\gamma = 0.05$ using SGD(b) $\gamma = 0.05$ using SGD

Figure 5: Accuracy and loss plots for CNN. The network is trained with SGD and swapped ReLU and pooling layers.

We can see that compared to Figure 3 the training/testing loss and accuracy plots converge slightly more unstable. However, this could just be due to a difference in initialisation. The final testing accuracy however is quite similar at 98.64%, which is only a 0.05% difference. We will next see whether a swapping of activation and pooling layers is more apparent for a hyperbolic tangent activation function. This is shown in Figure 6.

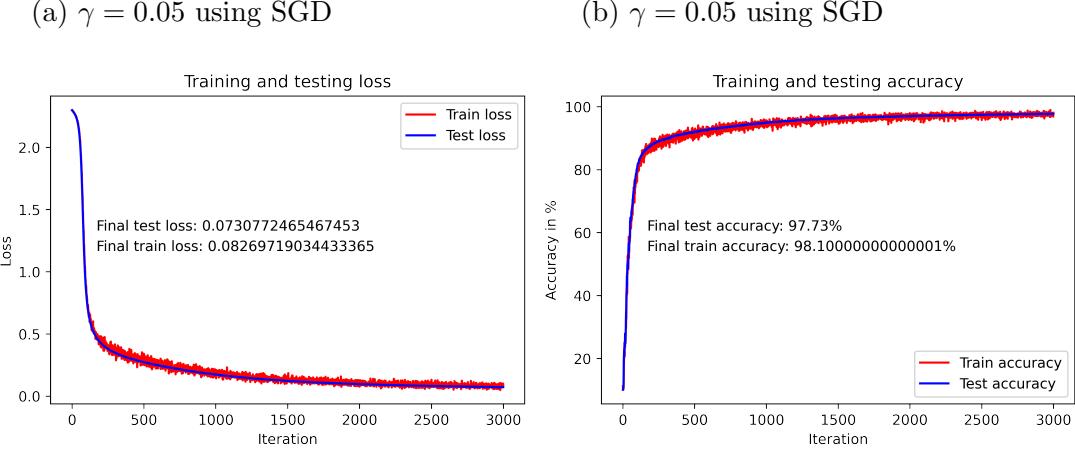


Figure 6: Accuracy and loss plots for CNN. The network is trained with SGD and swapped hyperbolic tangent and pooling layers.

One of the biggest differences between the swapped activation and pooling layers is indeed more visible for hyperbolic tangent activation than ReLU activation. The first thing we note is that the training/testing loss and accuracy plots converge in a more stable way. The biggest difference is that the final testing accuracy is 97.73%. This is a difference of 0.86% which is an order-of-magnitude different. If we look at the PyTorch documentation for the [ReLU](#) and [hyperbolic tangent activation functions](#) and the equations:

$$ReLU(x) = \max(0, x) \quad (1)$$

and

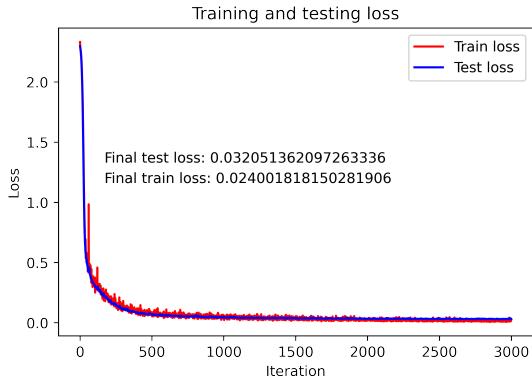
$$Tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (2)$$

we can see that the Tanh activation saturates at a much smaller value for positive inputs whereas the ReLU activation will increase proportionally for larger positive input values.

Finally, we will look at 3 different ways to try to improve the testing accuracy even more. Since ADAM performs much better than plain SGD, we want to see if there is an improvement when we add batch normalization to the CNN. We will be using a batch-size of 1000 and default ADAM values.

There has been some discussion as to whether we should place the batch normalization layer before or after the non-linearity, in this work we will choose the former. Results are shown in Figure 7.

(a) $\gamma = 0.001$ using ADAM



(b) $\gamma = 0.001$ using ADAM

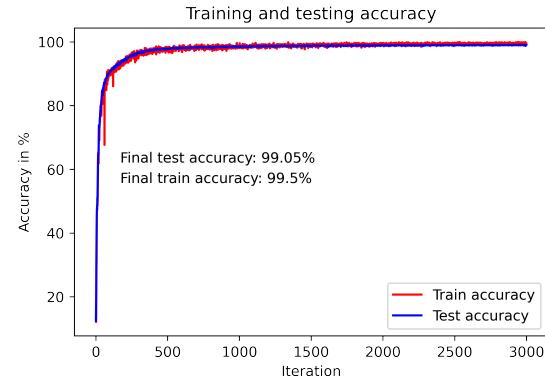


Figure 7: Accuracy and loss plots for CNN. The network is trained with SGD and ReLU with batch normalization before the ReLU.

The final testing accuracy for an ADAM-trained CNN is 99.05%. This is the highest final testing accuracy compared to all previous results. We will explore two more approaches to try to improve the testing accuracy. Next, we will look at adding dropout to the previous model. Following this [discussion](#) we will be placing the dropout layer after the ReLU activation and use a dropout probability of 0.5 following the [original dropout paper](#).

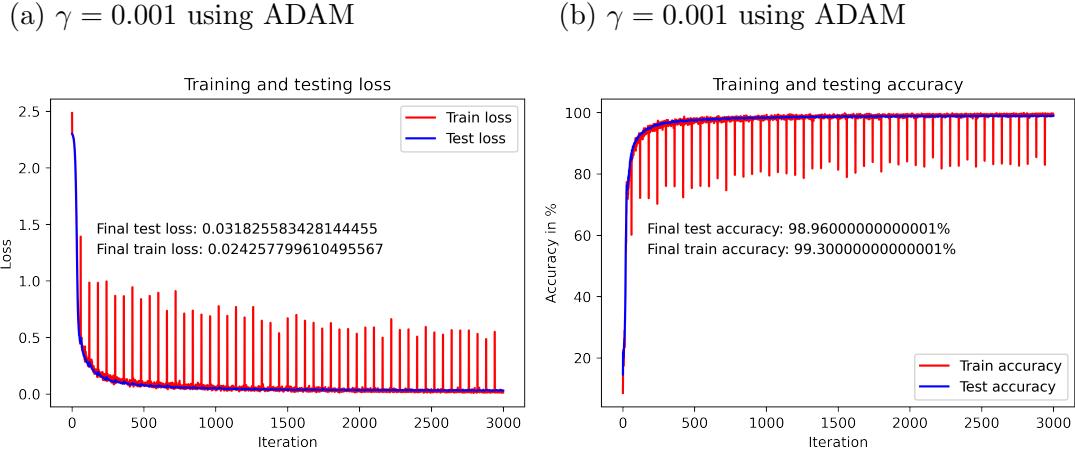


Figure 8: Accuracy and loss plots for CNN. Dropout layers are added after the ReLU activation.

One the biggest difference we can see is that the training loss and accuracy show a very oscillatory behaviour, converging to a final testing accuracy of 98.96%. To test a different type of activation function and see if we can gain back some accuracy, we will use the **SELU** activation function. We tested the SELU activation both with and without batch normalization and only display the best results. Whereas batch normalization does an explicit normalization to mean zero and unit variance, the SELU activation is supposed to "automatically converge towards zero mean and unit variance" (literal quote).

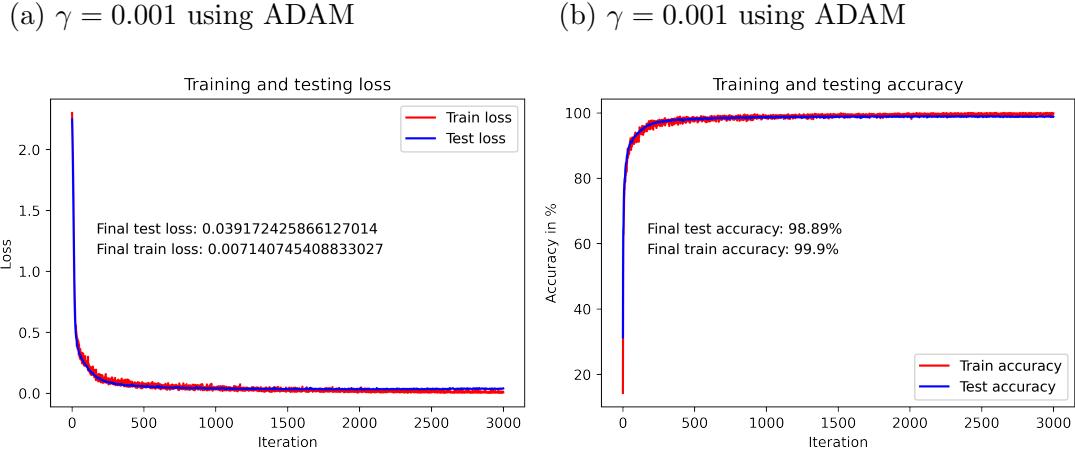


Figure 9: Accuracy and loss plots for CNN. The SELU activation function is used instead of ReLU.

The SELU CNN does not manage to give the best classification results at only a final testing accuracy of 98.89%. The next explorations would be to change the network architecture, but it is interesting to try to optimize the current architecture as much as possible using different optimizers and etc.

To finalize this section, we will be looking at the best-performing model and displaying a confusion matrix which denotes the number of correct classifications on the main diagonal and wrong classifications on the off-diagonals. A wrong classification would be for example an "1" being classified as a "7", or a "6" being classified as an "8" for example. We also show some wrong classifications with the predicted and correct class labels. Since the CNN trained with ADAM gave the best results, we will show the confusion matrix and some wrongly classified images next.

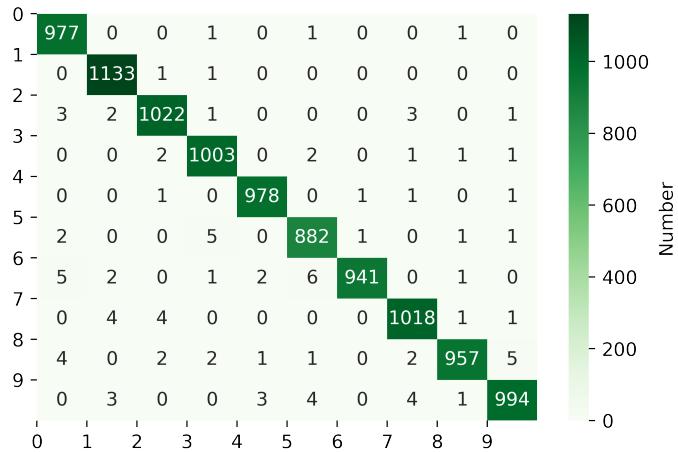


Figure 10: Confusion matrix of best network using ADAM.



Figure 11: Several wrong classifications using ADAM. From left to right and top to bottom. The correct class is 4 - predicted is 9. The correct class is 6 - predicted is 1. Correct class is 1 - predicted is 3. Correct class is 9 - predicted is 1.

Looking at the distribution of the predictions in the confusion matrix there is not a huge tendency of the network to predict some classes more wrong than others. The example images that were classified wrongly are indeed quite ambiguous so a different approach might yield better results. In the next part we will leave behind MNIST and look at biomedical image segmentation.

3 PyTorch and semantic segmentation

In this section we will explore semantic segmentation of the Warwick biomedical dataset. The dataset contains 80 training images and 60 test images of 128 by 128 pixels. The images provided have the blue channel set to a constant value of zero. In the code attached in the Appendix, this channel is not removed. Our task is a binary one: predict 1 where there is a gland in the testing data and 0 where there is no gland (or background). The provided test labels are binary masks. The training and testing images and segmentation masks are all normalized by dividing by 255. There is an additional step done before we can actually start implementing the training loop. In Pytorch the mini-batches need to have the number of input channels in the second dimensions. Therefore we need to permute the normalized images as follows:

```
1 xtra_torch = torch.from_numpy(xtra).permute(0, 3, 1,
    ↵ 2).float()
2 ytra_torch = torch.from_numpy(ytra)
3 xtes_torch = torch.from_numpy(xtes).permute(0, 3, 1,
    ↵ 2).float()
4 ytes_torch = torch.from_numpy(ytes)
```

To start with we will start from the CNN architecture for the previous MNIST classification task and adapt it into a fully-connected network. This will be done by replacing the fully-connected layer with a 1x1 convolution with the appropriate number (in our case 2) of output channels. After each pooling layers, our spatial resolution is reduced by 2 (using the same pooling layers as the MNIST CNN). To get an one-to-one correspondence between our CNN output and the segmentation masks we need to upsample the information encoded in the first part of the network. To do this, we will use transposed convolutions. In order to avoid [checkerboard artifacts](#) we will use transposed

convolutional layers such that the filter size is an integer multiple of the stride. In what follows, we will use a kernel size of 4x4 and a stride of 2 for each transposed convolution layer. To evaluate the segmentation results on the testing data, we will use the Sørensen–Dice coefficient (DSC for short):

$$DSC(A, B) = \frac{2|A \cap B|}{|A| + |B|} \quad (3)$$

In the above equation A is the segmentation output of the network and B is the provided ground truth segmentation masks. Before we show the final segmentation results and some segmentation outputs we will discuss briefly about the initial network architectures and design choices that were tried. As mentioned earlier, the initial network was an one-to-one copy of the MNIST CNN provided earlier where the fully-connected layer was removed and the rest was kept the same. Using the [binary cross entropy loss with logits](#) as our loss function we initially trained with a batch-size of 10 and an ADAM optimizer using the default parameters suggested by PyTorch:

```
1 adam_opt = optim.Adam(model.parameters(), lr = 0.001, betas =
  ↪ (0.9, 0.999), eps = 1e-08, weight_decay = 0.0, amsgrad =
  ↪ False)
```

We also trained for 50 epochs, and after it was observed that the training was quite fast we extended this to 300 epochs to ensure convergence is achieved (important for later). However, the loss didn't decrease and the DSC coefficient stayed constant with spikes at some iterations. Therefore, a different network architecture was used after experimentation with several different candidate networks. This is the final network:

```
1 segnetbn = nn.Sequential(nn.Conv2d(in_channels = 3,
  ↪ out_channels = 8, kernel_size = (3,3), stride = 1, padding
  ↪ = 1), nn.BatchNorm2d(8), nn.ReLU(),
  ↪ nn.MaxPool2d(kernel_size = 2, stride = 2),
  ↪ nn.Conv2d(in_channels = 8,
  ↪ out_channels = 16, kernel_size =
  ↪ (3,3), stride = 1, padding = 1),
  ↪ nn.BatchNorm2d(16), nn.ReLU(),
  ↪ nn.MaxPool2d(kernel_size = 2,
  ↪ stride = 2),
```

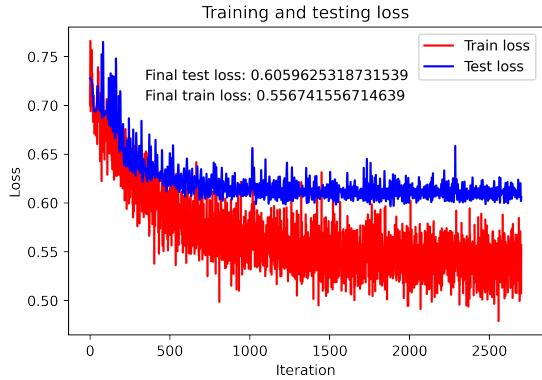
```

3      nn.Conv2d(in_channels = 16,
4          ←  out_channels = 32, kernel_size =
4          ←  (3,3), stride = 1, padding = 1),
4          ←  nn.BatchNorm2d(32), nn.ReLU(),
4          ←  nn.MaxPool2d(kernel_size = 2,
4          ←  stride = 2),
5      nn.Conv2d(in_channels = 32,
6          ←  out_channels = 64, kernel_size =
6          ←  (3,3), stride = 1, padding = 1),
6          ←  nn.BatchNorm2d(64), nn.ReLU(),
7      nn.ConvTranspose2d(in_channels = 64,
8          ←  out_channels = 32, kernel_size =
8          ←  (4,4), stride = 2, padding = 1),
9      nn.ConvTranspose2d(in_channels = 32,
10     ←  out_channels = 16, kernel_size =
10     ←  (4,4), stride = 2, padding = 1),
11     nn.ConvTranspose2d(in_channels = 16,
12     ←  out_channels = 8, kernel_size =
12     ←  (4,4), stride = 2, padding = 1),
13     nn.Conv2d(in_channels = 8,
14     ←  out_channels = 2, kernel_size =
14     ←  (1,1), stride = 1, padding = 0),
15     ←  nn.Sigmoid())

```

Importantly, batch-normalization layers were added before each ReLU non-linearity and an additional convolutional layer was added to increase the number of channels to 64. This network was trained with ADAM for 300 epochs with default parameters (as mentioned previously). The mini-batch size was 10 and evaluation was done on the entire test data set. A final curiosity about the network is that despite the fact that the used loss function actually combines a sigmoid layer and a binary cross entropy loss function the network still outputs negative predictions. We reduced the loss by adding a sigmoid function to the final layer. Since we know that the sigmoid outputs values between 0 and 1, and our segmentation masks are also between 0 and 1 the network will most likely converge in a more stable way. Results are shown in Figure 12.

(a) Binary cross entropy loss



(b) Dice coefficient

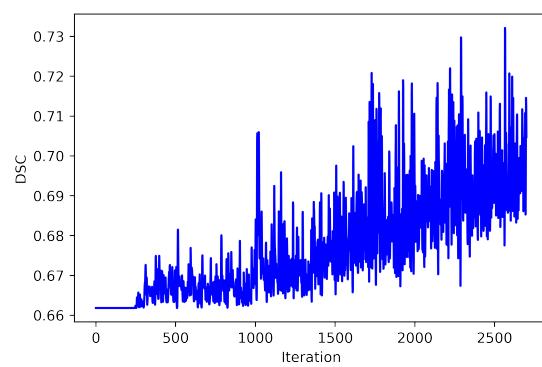


Figure 12: Training and testing loss of the segmentation network on the left. The testing loss plateaus at a larger value than the training loss. There is no indication of overfitting. On the right, the dice coefficient evaluated on the entire testing data is shown. The final dice coefficient is about 0.7047.

The final training and testing losses are indicated in the above caption, as well as the final dice coefficient. We can see that the dice coefficient does increase over time, however the increase is slow. Next we will show one example of a prediction on the training data with the correct mask.

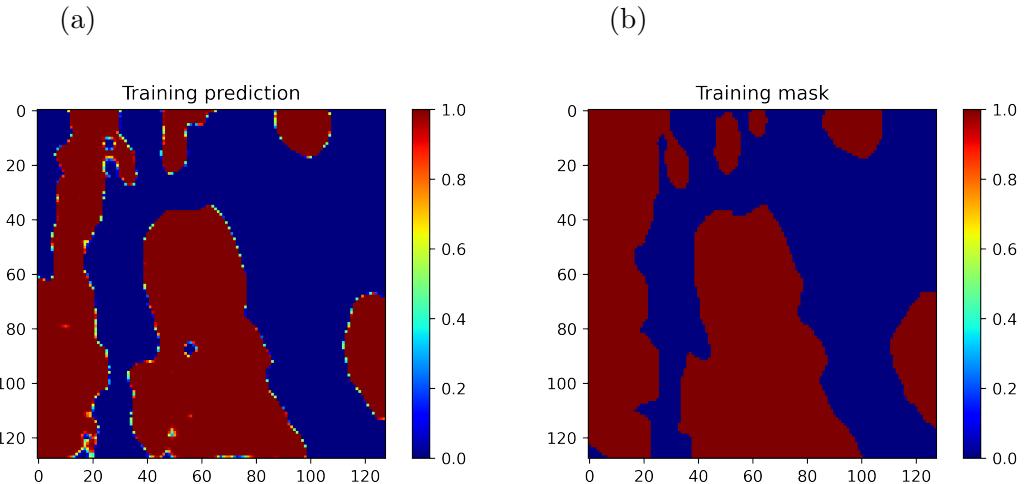


Figure 13: An example from the first of the 10 final mini-batches with the network prediction on the left and the ground truth mask on the right. The color map was chosen such that the differences between prediction and ground truth mask can be observed clearly.

We finally look at two examples out of the sixty testing images with network prediction and ground truth masks as a comparison. Since the testing data was not shuffled the two chosen test images correspond directly to the provided data set.

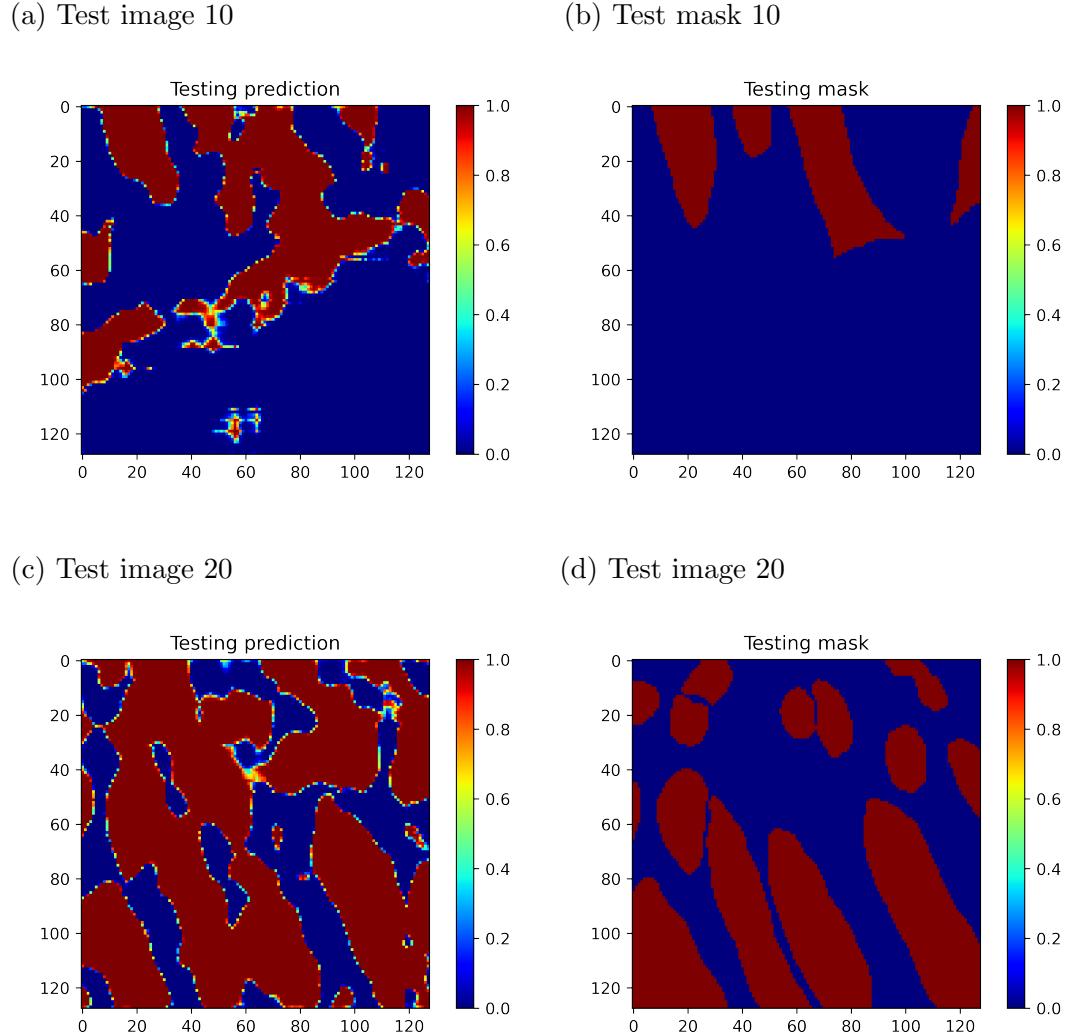


Figure 14: Network predictions and ground truth masks of the testing data. All network predictions are shown on the left and all masks on the right. The number of the image is indicated above the plot.

In the above figure we can see that the network predictions for both image 10 and image 20 are far off from the ground truth mask. In image 10 we can see in fact that the network predicted glands in locations where is not even a gland in the ground truth. This is less apparent for image 20 since the

ground truth is not as sparse as image 10. We suspect this has something to do with the fact that the convolutional kernels have too small kernel sizes. So to aggregate more information about the actual image it might be beneficial to try to increase the kernel size at the same time as we increase the channel depth. However, it is promising that at locations where there is actually a gland the network managed to predict a gland as well. It might be more harmful to miss glands than to wrongfully predict a gland, however it is obviously better to not predict a gland where there is none. The examples chosen represent a dense and sparse image with varying shapes. Test images that contained many similar objects seemed to be predicted correctly more often since the network can re-use the information it has encoded for one object. This is an advantage of the convolutional approach.

4 Appendix

All models are implemented in Jupyter notebooks to allow for easy experimentation. We will attach the code in this report as a single file. The code will also be separately attached as two separate Jupyter notebooks for the MNIST and semantic segmentation task.

4.1 PyTorch and MNIST

```
1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torch.nn.functional as F
5
6 from torch.utils.data import TensorDataset
7 from torch.utils.data import DataLoader
8
9 import numpy as np
10 import seaborn as sns
11
12 from matplotlib import pyplot as plt
13 from sklearn.metrics import confusion_matrix
14
15 from load_mnist import load_mnist
```

```

16
17 # Loading training and testing data
18 xtrain, ytrain, xtest, ytest = load_mnist()
19
20 # Reset function
21 def reset_model(model):
22     for layer in model.children():
23         if hasattr(layer, 'reset_parameters'):
24             layer.reset_parameters()
25
26 # Get dimensions of training and testing data
27 M = ytrain.shape[1]
28 p = xtrain.shape[1]
29
30 ntrain = xtrain.shape[0]
31 ntest = xtest.shape[0]
32
33 # Converting training and testing data to torch tensors
34 xtra_torch = torch.tensor(xtrain).float()
35 ytra_torch = torch.tensor(ytrain).float()
36 xtes_torch = torch.tensor(xtest).float()
37 ytes_torch = torch.tensor(ytest).float()
38
39 # Converting our flat vectors to images for the CNNs
40 xtra_conv = torch.tensor(xtrain).reshape(ntrain, 28,
41     ↳ 28).float().unsqueeze(1)
42 ytra_conv = torch.tensor(ytrain).float()
43 xtes_conv = torch.tensor(xtest).reshape(ntest, 28,
44     ↳ 28).float().unsqueeze(1)
45 ytes_conv = torch.tensor(ytest).float()
46
47 # Check if CUDA is available
48 device = torch.device("cuda") if torch.cuda.is_available()
49     ↳ else torch.device("cpu")
50 print("Device:", torch.cuda.get_device_name(device))
51
52 # FCN
53 num_batch = 1000

```

```

51 | training_data = TensorDataset(xtra_torch.to(device),
52 |                                ↵ ytra_torch.to(device))
53 |
54 | train_dat_fcn = DataLoader(training_data, shuffle = True,
55 |                                ↵ batch_size = num_batch)
56 |
57 | # CNN
58 | train_dat_im = TensorDataset(xtra_conv.to(device),
59 |                                ↵ ytra_conv.to(device))
60 | train_dat_cnn = DataLoader(train_dat_im, shuffle = True,
61 |                                ↵ batch_size = num_batch)
62 |
63 | xte_im = xtes_conv.to(device)
64 | yte_im = ytes_conv.to(device)
65 |
66 | # 2-layer network
67 | model_2L = nn.Sequential(nn.Linear(p, 100), nn.ReLU(),
68 |                                ↵ nn.Linear(100, M))
69 |
70 | # 4-layer network
71 | model_4L = nn.Sequential(nn.Linear(p, 256), nn.ReLU(),
72 |                                ↵ nn.Linear(256, 128), nn.ReLU(), nn.Linear(128, 64),
73 |                                ↵ nn.ReLU(), nn.Linear(64, M))
74 |
75 | # Convolutional network
76 | conv_mod = nn.Sequential(nn.Conv2d(in_channels = 1,
77 |                                ↵ out_channels = 8, kernel_size = (3,3), stride = 1, padding
78 |                                ↵ = 1), nn.ReLU(), nn.MaxPool2d(kernel_size = 2, stride =
79 |                                ↵ 2),
80 |                                ↵ nn.Conv2d(in_channels = 8,
81 |                                ↵ out_channels = 16, kernel_size =
82 |                                ↵ (3,3), stride = 1, padding = 1),
83 |                                ↵ nn.ReLU(),
84 |                                ↵ nn.MaxPool2d(kernel_size = 2,
85 |                                ↵ stride = 2),

```

```

73     nn.Conv2d(in_channels = 16,
    ↵   out_channels = 32, kernel_size =
    ↵   (3,3), stride = 1, padding = 1),
    ↵   nn.ReLU(), nn.Flatten(),
    ↵   nn.Linear(32 * 7 * 7, M))

74 conv_swap = nn.Sequential(nn.Conv2d(in_channels = 1,
    ↵   out_channels = 8, kernel_size = (3,3), stride = 1, padding
    ↵   = 1), nn.MaxPool2d(kernel_size = 2, stride = 2),
    ↵   nn.ReLU(),
75           nn.Conv2d(in_channels = 8,
    ↵   out_channels = 16, kernel_size =
    ↵   (3,3), stride = 1, padding = 1),
    ↵   nn.MaxPool2d(kernel_size = 2,
    ↵   stride = 2), nn.ReLU(),
76           nn.Conv2d(in_channels = 16,
    ↵   out_channels = 32, kernel_size =
    ↵   (3,3), stride = 1, padding = 1),
    ↵   nn.ReLU(), nn.Flatten(),
    ↵   nn.Linear(32 * 7 * 7, M))

77 conv_swap_htan = nn.Sequential(nn.Conv2d(in_channels = 1,
    ↵   out_channels = 8, kernel_size = (3,3), stride = 1, padding
    ↵   = 1), nn.MaxPool2d(kernel_size = 2, stride = 2),
    ↵   nn.Tanh(),
78           nn.Conv2d(in_channels = 8,
    ↵   out_channels = 16, kernel_size =
    ↵   (3,3), stride = 1, padding = 1),
    ↵   nn.MaxPool2d(kernel_size = 2,
    ↵   stride = 2), nn.Tanh(),
79           nn.Conv2d(in_channels = 16,
    ↵   out_channels = 32, kernel_size =
    ↵   (3,3), stride = 1, padding = 1),
    ↵   nn.Tanh(), nn.Flatten(),
    ↵   nn.Linear(32 * 7 * 7, M))

80
81
82

```

```

83 conv_mod_bn = nn.Sequential(nn.Conv2d(in_channels = 1,
84     ↵  out_channels = 8, kernel_size = (3,3), stride = 1, padding
85     ↵  = 1), nn.BatchNorm2d(8), nn.ReLU(),
86     ↵  nn.MaxPool2d(kernel_size = 2, stride = 2),
87             nn.Conv2d(in_channels = 8,
88                 ↵  out_channels = 16, kernel_size =
89                 ↵  (3,3), stride = 1, padding = 1),
90                 ↵  nn.BatchNorm2d(16), nn.ReLU(),
91                 ↵  nn.MaxPool2d(kernel_size = 2,
92                 ↵  stride = 2),
93             nn.Conv2d(in_channels = 16,
94                 ↵  out_channels = 32, kernel_size =
95                 ↵  (3,3), stride = 1, padding = 1),
96                 ↵  nn.BatchNorm2d(32), nn.ReLU(),
97                 ↵  nn.Flatten(), nn.Linear(32 * 7 *
98                 ↵  7, M))
99
100 conv_mod_bn_dropout = nn.Sequential(nn.Conv2d(in_channels = 1,
101     ↵  out_channels = 8, kernel_size = (3,3), stride = 1, padding
102     ↵  = 1), nn.BatchNorm2d(8), nn.ReLU(), nn.Dropout2d(0.5),
103     ↵  nn.MaxPool2d(kernel_size = 2, stride = 2),
104             nn.Conv2d(in_channels = 8,
105                 ↵  out_channels = 16, kernel_size =
106                 ↵  (3,3), stride = 1, padding = 1),
107                 ↵  nn.BatchNorm2d(16), nn.ReLU(),
108                 ↵  nn.Dropout2d(0.5),
109                 ↵  nn.MaxPool2d(kernel_size = 2,
110                 ↵  stride = 2),
111             nn.Conv2d(in_channels = 16,
112                 ↵  out_channels = 32, kernel_size =
113                 ↵  (3,3), stride = 1, padding = 1),
114                 ↵  nn.BatchNorm2d(32), nn.ReLU(),
115                 ↵  nn.Dropout2d(0.5), nn.Flatten(),
116                 ↵  nn.Linear(32 * 7 * 7, M))

```

```

91 conv_mod_bn_selu = nn.Sequential(nn.Conv2d(in_channels = 1,
92     ↳ out_channels = 8, kernel_size = (3,3), stride = 1, padding
93     ↳ = 1), nn.SELU(), nn.MaxPool2d(kernel_size = 2, stride =
94     ↳ 2),
95             nn.Conv2d(in_channels = 8,
96                 ↳ out_channels = 16, kernel_size =
97                 ↳ (3,3), stride = 1, padding = 1),
98                 ↳ nn.SELU(),
99                 ↳ nn.MaxPool2d(kernel_size = 2,
100                    ↳ stride = 2),
101                nn.Conv2d(in_channels = 16,
102                    ↳ out_channels = 32, kernel_size =
103                    ↳ (3,3), stride = 1, padding = 1),
104                    ↳ nn.SELU(), nn.Flatten(),
105                    ↳ nn.Linear(32 * 7 * 7, M))
106
107
108
109
110
111
112
113
114
115
# Selecting model
mod = conv_mod

# Resetting model when we switch architectures
reset_model(mod)

# Moving model to GPU
mod.to(device)

# Check dimensions through network
check_dims = False
if check_dims == True:
    X = torch.rand(size = (10, 1, 28, 28)).to(device)
    for layer in mod:
        X = layer(X)
        print(layer.__class__.__name__, 'output shape:\t',
              X.shape)

# Loss function
loss = F.cross_entropy
it = 0

```

```

116 loss_train , acc_train = [] , []
117 loss_test, acc_test = [], []
118 l_rate = 0.05
119
120 # Vanilla SGD optimizer
121 sgd_opt = optim.SGD(mod.parameters(), lr = l_rate)
122
123 # Adam optimizer
124 adam_opt = optim.Adam(mod.parameters(), lr = 0.001 , betas =
125   ↪ (0.9, 0.999), eps = 1e-08, weight_decay = 0, amsgrad =
126   ↪ False)
127
128 # Exercise 1.5 - choosing three different optimizers
129 # Adadelta optimizer
130 adad_opt = optim.Adadelta(mod.parameters(), lr = 1.0, rho =
131   ↪ 0.9, eps = 1e-06, weight_decay = 0)
132
133 # Adagrad optimizer
134 adag_opt = optim.Adagrad(mod.parameters(), lr = 0.01, lr_decay
135   ↪ = 0, weight_decay = 0, initial_accumulator_value = 0, eps
136   ↪ = 1e-10)
137
138 # Adamax optimizer
139 adamax_opt = optim.Adamax(mod.parameters(), lr = 0.002, betas
140   ↪ = (0.9, 0.999), eps = 1e-08, weight_decay = 0)
141
142 # ASGD optimizer
143 asgd_opt = optim.ASGD(mod.parameters(), lr = 0.01, lambd =
144   ↪ 0.0001, alpha = 0.75, t0 = 1000000.0, weight_decay = 0)
145
146 # RMSprop optimizer
147 rmsp_opt = optim.RMSprop(mod.parameters(), lr = 0.01, alpha =
148   ↪ 0.99, eps = 1e-08, weight_decay = 0, momentum = 0,
149   ↪ centered = False)
150
151 # Selecting optimizer
152 optr = adam_opt
153 full_n = True

```

```

145
146 if full_n == False:
147     training_dat = train_dat_fcn
148     xt = xte
149     yt = yte
150 else:
151     training_dat = train_dat_cnn
152     xt = xte_im
153     yt = yte_im
154
155 # Optimization loop
156 num_epochs = 50
157
158 y_predictions = []
159
160 for epoch in range(num_epochs):
161     mod.train()
162     for xbatch, ybatch in training_dat:
163         optr.zero_grad()
164         prediction = mod(xbatch)
165         pred = torch.argmax(prediction, dim = 1)
166
167         true_lab = torch.argmax(ybatch, dim = 1)
168
169         ce_loss = loss(prediction, true_lab)
170
171         acc_train.append(100 * (1/num_batch) * torch.sum(pred
172             == true_lab, dim = 0).item())
173
174         loss_train.append(ce_loss.item())
175
176         ce_loss.backward()
177         optr.step()
178
179         mod.eval()
180         with torch.no_grad():
181             ytrue_lab = torch.argmax(yt, dim = 1)
182             prediction_test = mod(xt)

```

```

182     pred_test = torch.argmax(prediction_test, dim = 1)
183
184     ce_test = loss(prediction_test, ytrue_lab)
185     loss_test.append(ce_test.item())
186     acc_test.append(100 * (1/ntest) *
187         → torch.sum(pred_test == ytrue_lab, dim =
188         → 0).item())
189
190     # Saving final predictions
191     if epoch == (num_epochs - 1):
192         y_predictions = mod(xt)
193
194     it += 1
195     print("Epoch %s/%s" % (epoch + 1, num_epochs))
196
197     plt.figure(1)
198     loss_tra, = plt.plot(loss_train, 'r')
199     loss_tes, = plt.plot(loss_test, 'b')
200     plt.title("Training and testing loss")
201     plt.xlabel("Iteration")
202     plt.ylabel("Loss")
203     plt.legend([loss_tra, loss_tes], ['Train loss', 'Test loss'])
204     plt.annotate("Final train loss: %s" % (loss_train[-1]),
205         → ,xycoords = 'figure fraction', xy = (0.2,0.5))
206     plt.annotate("Final test loss: %s" % (loss_test[-1]), xycoords
207         → = 'figure fraction', xy = (0.2,0.55))
208     print("Final training loss: %.2f" % loss_train[-1])
209     print("Final testing loss: %.2f" % loss_test[-1])
210     plt.savefig("loss_plots", dpi = 500)
211
212     plt.figure(2)
213     acc_tra, = plt.plot(acc_train, 'r')
214     acc_tes, = plt.plot(acc_test, 'b')
215     plt.title("Training and testing accuracy")
216     plt.xlabel("Iteration")
217     plt.ylabel("Accuracy in %")
218     plt.legend([acc_tra, acc_tes], ['Train accuracy', 'Test
219         → accuracy'])
```

```

215 plt.annotate("Final train accuracy: %s%%" % (acc_train[-1])
216     ,xycoords = 'figure fraction', xy = (0.2,0.5))
217 plt.annotate("Final test accuracy: %s%%" % (acc_test[-1]),
218     xycoords = 'figure fraction', xy = (0.2,0.55))
219 print()
220 print("Final training accuracy: %s%%." % acc_train[-1])
221 print("Final testing accuracy: %s%%." % acc_test[-1])
222 plt.savefig("acc_plots", dpi = 500)
223
224
225 # Retrieve data and labels from the GPU
226 y_test = yte_im.data.cpu()
227
228
229 # Remove singleton dimension
230 y_test = np.squeeze(y_test, axis = 1)
231 ytrue = torch.argmax(y_test, dim = 1)
232
233
234 ypreds = y_predictions.data.cpu()
235 ypredslab = torch.argmax(ypreds, dim = 1)
236
237 cf_mat = confusion_matrix(ytrue, ypredslab, normalize = None)
238
239 mnist_classes = np.arange(10)
240 mnist_lab = ['0', '1', '2', '3', '4', '5', '6', '7', '8', '9']
241
242 if True:
243     plt.figure(3)
244     cmat = sns.heatmap(cf_mat, cmap = 'Greens', annot = True,
245         fmt = "d", cbar_kws = {'label':'Number'}) # use .3f
246         for floating-point data
247     cmat.set_xticklabels(mnist_lab, rotation = 'horizontal')
248     cmat.set_yticklabels(mnist_lab, rotation = 'horizontal')
249     plt.xticks(mnist_classes, mnist_lab)
250     plt.yticks(mnist_classes, mnist_lab)
251     plt.savefig("confusion_mat", dpi = 500)
252
253
254 # Showing 3 examples of wrong classifications
255 wrong_examples = (ytrue != ypredslab)
256
257

```

```

249 wrong_preds = np.squeeze(xtes_conv[wrong_examples])
250 wrong_labs = ypredslab[wrong_examples]
251 actual_labs = torch.argmax(ytes_conv[wrong_examples], dim = 1)
252
253 ind = np.random.choice(a = wrong_labs.shape[0], size = 4,
254    ↳ replace = True)
255
256 fig, ax = plt.subplots(2 , 2)
257 for a in fig.axes:
258     a.set_xticks([])
259     a.set_yticks([])
260
261 ax[0,0].imshow(wrong_preds[ind[0], :, :], cmap = 'gray')
262 ax[0,1].imshow(wrong_preds[ind[1], :, :], cmap = 'gray')
263 ax[1,0].imshow(wrong_preds[ind[2], :, :], cmap = 'gray')
264 ax[1,1].imshow(wrong_preds[ind[3], :, :], cmap = 'gray')
265 plt.savefig("wrong_preds", dpi = 500)
266
267 for i in ind:
268     print("Correct label: %d" % actual_labs[i].item())
269     print("Predicted label: %d" % wrong_labs[i].item())

```

4.2 PyTorch and semantic segmentation

```

1 import torch
2 import torch.nn as nn
3 import torch.optim as optim
4 import torch.nn.functional as F
5
6 import seaborn as sn
7 import numpy as np
8 from sklearn.metrics import confusion_matrix
9
10 from torch.utils.data import TensorDataset
11 from torch.utils.data import DataLoader
12
13 from load_warwick import load_warwick

```

```

14
15 from matplotlib import pyplot as plt
16
17 # Loading data
18 # Training data: 85 images of 128x128x3
19 # Training labels: 85 images of 128x128
20 # Testing data: 60 images of 128x128x3
21 # Testing labels: 60 images of 128x128
22 xtra, ytra, xtes, ytes = load_warwick()
23
24 # Resetting model function
25 # Credits:
26 #   ↪ https://discuss.pytorch.org/t/reset-model-weights/19180/4
27 def reset_model(model):
28     for layer in model.children():
29         if hasattr(layer, 'reset_parameters'):
30             layer.reset_parameters()
31
32 # Sorenson-Dice coefficient
33 def dsc(A, B):
34
35     a = A.bool()
36     b = B.bool()
37
38     intersect = torch.logical_and(a,b)
39
40     dice_coeff = (2.0 * intersect.sum()) / (a.sum() + b.sum())
41
42     return dice_coeff
43
44 # Preprocessing of train/test data by dividing with maximum
45 #   intensity
46 xtra = xtra / 255
47 ytra = ytra / 255
48
49 xtes = xtes / 255
50 ytes = ytes / 255

```

```

50 # Get some dimensions
51 pixels = 128
52 ntrain = xtra.shape[0]
53 ntest = xtes.shape[0]
54
55 # Convert to Torch tensors
56 xtra_torch = torch.from_numpy(xtra).permute(0, 3, 1,
57                                         ↵ 2).float()
57 ytra_torch = torch.from_numpy(ytra)
58 xtes_torch = torch.from_numpy(xtes).permute(0, 3, 1,
59                                         ↵ 2).float()
60 ytes_torch = torch.from_numpy(ytes)
61
62 # Check if CUDA is available
63 device = torch.device("cuda") if torch.cuda.is_available()
64                                         ↵ else torch.device("cpu")
63 print("Device:", torch.cuda.get_device_name(device))
64
65 # Creating data loader for training data
66 num_batch = 10
67 training_set = TensorDataset(xtra_torch.to(device),
68                             ↵ ytra_torch.to(device))
68 training_loader = DataLoader(training_set, shuffle = True,
69                             ↵ batch_size = num_batch)
69
70 xtest = xtes_torch.to(device)
71 ytest = ytes_torch.to(device)
72
73 # Convolutional network
74 segnet = nn.Sequential(nn.Conv2d(in_channels = 3, out_channels
75                         ↵ = 8, kernel_size = (3,3), stride = 1, padding = 1),
75                         ↵ nn.ReLU(), nn.MaxPool2d(kernel_size = 2, stride = 2),
75                         ↵ nn.Conv2d(in_channels = 8,
75                         ↵ out_channels = 16, kernel_size =
75                         ↵ (3,3), stride = 1, padding = 1),
75                         ↵ nn.ReLU(),
75                         ↵ nn.MaxPool2d(kernel_size = 2,
75                         ↵ stride = 2),

```

```

76     nn.Conv2d(in_channels = 16,
77             ← out_channels = 32, kernel_size =
78             ← (3,3), stride = 1, padding = 1),
79             ← nn.ReLU(),
80             ← nn.MaxPool2d(kernel_size = 2,
81             ← stride = 2),
82             nn.Conv2d(in_channels = 32,
83             ← out_channels = 64, kernel_size =
84             ← (3,3), stride = 1, padding = 1),
85             ← nn.ReLU(),
86             nn.ConvTranspose2d(in_channels = 64,
87             ← out_channels = 32, kernel_size =
88             ← (4,4), stride = 2, padding = 1),
89             nn.ConvTranspose2d(in_channels = 32,
90             ← out_channels = 16, kernel_size =
91             ← (4,4), stride = 2, padding = 1),
92             nn.ConvTranspose2d(in_channels = 16,
93             ← out_channels = 8, kernel_size =
94             ← (4,4), stride = 2, padding = 1),
95             nn.Conv2d(in_channels = 8,
96             ← out_channels = 2, kernel_size =
97             ← (1,1), stride = 1, padding = 0))

98 segnetbn = nn.Sequential(nn.Conv2d(in_channels = 3,
99             ← out_channels = 8, kernel_size = (3,3), stride = 1, padding
100             ← = 1), nn.BatchNorm2d(8), nn.ReLU(),
101             ← nn.MaxPool2d(kernel_size = 2, stride = 2),
102             nn.Conv2d(in_channels = 8,
103             ← out_channels = 16, kernel_size =
104             ← (3,3), stride = 1, padding = 1),
105             ← nn.BatchNorm2d(16), nn.ReLU(),
106             ← nn.MaxPool2d(kernel_size = 2,
107             ← stride = 2),

```

```

85     nn.Conv2d(in_channels = 16,
86             ← out_channels = 32, kernel_size =
87             ← (3,3), stride = 1, padding = 1),
88             ← nn.BatchNorm2d(32), nn.ReLU(),
89             ← nn.MaxPool2d(kernel_size = 2,
90             ← stride = 2),
91             nn.Conv2d(in_channels = 32,
92                     ← out_channels = 64, kernel_size =
93                     ← (3,3), stride = 1, padding = 1),
94                     ← nn.BatchNorm2d(64), nn.ReLU(),
95                     nn.ConvTranspose2d(in_channels = 64,
96                             ← out_channels = 32, kernel_size =
97                             ← (4,4), stride = 2, padding = 1),
98                             nn.ConvTranspose2d(in_channels = 32,
99                                 ← out_channels = 16, kernel_size =
100                                 ← (4,4), stride = 2, padding = 1),
101                                 nn.ConvTranspose2d(in_channels = 16,
102                                     ← out_channels = 8, kernel_size =
103                                     ← (4,4), stride = 2, padding = 1),
104                                     nn.Conv2d(in_channels = 8,
105                                         ← out_channels = 2, kernel_size =
106                                         ← (1,1), stride = 1, padding = 0),
107                                         ← nn.Sigmoid())
108
109 m_seg_selu = nn.Sequential(nn.Conv2d(in_channels = 3,
110                             ← out_channels = 8, kernel_size = (3,3), stride = 1, padding
111                             ← = 1), nn.SELU(), nn.MaxPool2d(kernel_size = 2, stride =
112                             ← 2),
113                             nn.Conv2d(in_channels = 8,
114                                 ← out_channels = 16, kernel_size =
115                                 ← (3,3), stride = 1, padding = 1),
116                                 nn.SELU(),
117                                 nn.MaxPool2d(kernel_size = 2,
118                                 ← stride = 2),
119

```

```

94     nn.Conv2d(in_channels = 16,
95             ← out_channels = 32, kernel_size =
96             ← (3,3), stride = 1, padding = 1),
97             ← nn.SELU(),
98             ← nn.MaxPool2d(kernel_size = 2,
99             ← stride = 2),
100            nn.Conv2d(in_channels = 32,
101            ← out_channels = 64, kernel_size =
102            ← (3,3), stride = 1, padding = 1),
103            ← nn.SELU(),
104            nn.ConvTranspose2d(in_channels = 64,
105            ← out_channels = 32, kernel_size =
106            ← (4,4), stride = 2, padding = 1),
107            nn.ConvTranspose2d(in_channels = 32,
108            ← out_channels = 16, kernel_size =
109            ← (4,4), stride = 2, padding = 1),
110            nn.ConvTranspose2d(in_channels = 16,
111            ← out_channels = 8, kernel_size =
112            ← (4,4), stride = 2, padding = 1),
113            nn.Conv2d(in_channels = 8,
114            ← out_channels = 2, kernel_size =
115            ← (1,1), stride = 1, padding = 0))

116 m_seg_k_big = nn.Sequential(nn.Conv2d(in_channels = 3,
117             ← out_channels = 8, kernel_size = (5,5), stride = 1, padding
118             ← = 2), nn.BatchNorm2d(8), nn.ReLU(),
119             ← nn.MaxPool2d(kernel_size = 2, stride = 2),
120             nn.Conv2d(in_channels = 8,
121             ← out_channels = 16, kernel_size =
122             ← (5,5), stride = 1, padding = 2),
123             ← nn.BatchNorm2d(16), nn.ReLU(),
124             ← nn.MaxPool2d(kernel_size = 2,
125             ← stride = 2),

```

```

103     nn.Conv2d(in_channels = 16,
104             ← out_channels = 32, kernel_size =
105             ← (5,5), stride = 1, padding = 2),
106             ← nn.BatchNorm2d(32), nn.ReLU(),
107             ← nn.MaxPool2d(kernel_size = 2,
108             ← stride = 2),
109             nn.Conv2d(in_channels = 32,
110                     ← out_channels = 64, kernel_size =
111                     ← (5,5), stride = 1, padding = 2),
112                     ← nn.BatchNorm2d(64), nn.ReLU(),
113                     ← nn.MaxPool2d(kernel_size = 2,
114                     ← stride = 2),
115                     nn.Conv2d(in_channels = 64,
116                             ← out_channels = 128, kernel_size =
117                             ← (5,5), stride = 1, padding = 2),
118                             ← nn.BatchNorm2d(128), nn.ReLU(),
119                             nn.ConvTranspose2d(in_channels = 128,
120                                     ← out_channels = 64, kernel_size =
121                                     ← (4,4), stride = 2, padding = 1),
122                                     nn.ConvTranspose2d(in_channels = 64,
123                                         ← out_channels = 32, kernel_size =
124                                         ← (4,4), stride = 2, padding = 1),
125                                         nn.ConvTranspose2d(in_channels = 32,
126                                             ← out_channels = 16, kernel_size =
127                                             ← (4,4), stride = 2, padding = 1),
128                                             nn.ConvTranspose2d(in_channels = 16,
129                                                 ← out_channels = 8, kernel_size =
130                                                 ← (4,4), stride = 2, padding = 1),
131                                                 nn.Conv2d(in_channels = 8,
132                                         ← out_channels = 2, kernel_size =
133                                         ← (1,1), stride = 1, padding = 0))

m_seg_k_big_selu = nn.Sequential(nn.Conv2d(in_channels = 3,
    ← out_channels = 8, kernel_size = (5,5), stride = 1, padding
    ← = 2), nn.SELU(), nn.MaxPool2d(kernel_size = 2, stride =
    ← 2),

```

```

114     nn.Conv2d(in_channels = 8,
115             ← out_channels = 16, kernel_size =
116                 ← (5,5), stride = 1, padding = 2),
117                 ← nn.SELU(),
118                 ← nn.MaxPool2d(kernel_size = 2,
119                         ← stride = 2),
120                         ← nn.Conv2d(in_channels = 16,
121                             ← out_channels = 32, kernel_size =
122                                 ← (5,5), stride = 1, padding = 2),
123                                 ← nn.SELU(),
124                                 ← nn.MaxPool2d(kernel_size = 2,
125                                     ← stride = 2),
126                                     ← nn.Conv2d(in_channels = 32,
127                                         ← out_channels = 64, kernel_size =
128                                             ← (5,5), stride = 1, padding = 2),
129                                             ← nn.SELU(),
130                                             ← nn.MaxPool2d(kernel_size = 2,
131                                                 ← stride = 2),
132                                                 ← nn.Conv2d(in_channels = 64,
133                                                     ← out_channels = 128, kernel_size =
134                                                         ← (5,5), stride = 1, padding = 2),
135                                                         ← nn.SELU(),
136                                                         ← nn.ConvTranspose2d(in_channels = 128,
137                                                             ← out_channels = 64, kernel_size =
138                                                                 ← (4,4), stride = 2, padding = 1),
139                                                                 ← nn.ConvTranspose2d(in_channels = 64,
140                                                                     ← out_channels = 32, kernel_size =
141                                                                         ← (4,4), stride = 2, padding = 1),
142                                                                         ← nn.ConvTranspose2d(in_channels = 32,
143                                                                             ← out_channels = 16, kernel_size =
144                                                                                 ← (4,4), stride = 2, padding = 1),
145                                                                                 ← nn.ConvTranspose2d(in_channels = 16,
146                                                                 ← out_channels = 8, kernel_size =
147                                                                     ← (4,4), stride = 2, padding = 1),
148                                                                     ← nn.Conv2d(in_channels = 8,
149                         ← out_channels = 2, kernel_size =
150                             ← (1,1), stride = 1, padding = 0),
151                             ← nn.Sigmoid())

```

```
123  
124 # Selecting model  
125 model = segnetbn  
126  
127 # Resetting model when we switch architectures  
128 reset_model(model)  
129  
130 # Moving model to GPU  
131 model.to(device)  
132  
133 # Check dimensions through network  
134 # Credits:  
135     → https://d2l.ai/chapter\_convolutional-neural-networks/lenet.html  
136 check_dims = False  
137 if check_dims == True:  
138     X = torch.rand(size = (num_batch, 3, 128, 128)).to(device)  
139     for layer in model:  
140         X = layer(X)  
141         print(layer.__class__.__name__, 'output shape:\t',  
142             → X.shape)  
143  
144 # Loss function(s)  
145 bce = F.binary_cross_entropy  
146 bce_logits = F.binary_cross_entropy_with_logits  
147  
148 it = 0  
149 num_epochs = 300  
150 l_rate = 0.001 # 0.001 works good  
151  
152 sgd_opt = optim.SGD(model.parameters(), lr = l_rate,  
153     → weight_decay = 0, momentum = 0.0)  
154  
155 rmsp_opt = optim.RMSprop(model.parameters(), lr = l_rate,  
156     → alpha = 0.99, eps = 1e-08, weight_decay = 0, momentum =  
157     → 0.1, centered = False)
```

```

154 adam_opt = optim.Adam(model.parameters(), lr = l_rate, betas =
155     ↪ (0.9, 0.999), eps = 1e-08, weight_decay = 0.0, amsgrad =
156     ↪ False)
157
158 adad_opt = optim.Adadelta(model.parameters(), lr = 1.0, rho =
159     ↪ 0.9, eps = 1e-06, weight_decay = 0)
160
161 adag_opt = optim.Adagrad(model.parameters(), lr = 0.01,
162     ↪ lr_decay = 0, weight_decay = 0, initial_accumulator_value
163     ↪ = 0, eps = 1e-10)
164
165 adamax_opt = optim.Adamax(model.parameters(), lr = 0.002,
166     ↪ betas = (0.9, 0.999), eps = 1e-08, weight_decay = 0)
167
168 asgd_opt = optim.ASGD(model.parameters(), lr = 0.01, lambd =
169     ↪ 0.0001, alpha = 0.75, t0 = 1000000.0, weight_decay = 0)
170
171 #sched = torch.optim.lr_scheduler.ExponentialLR(optr, gamma,
172     ↪ last_epoch = -1, verbose = False)
173
174 optr = adam_opt
175 dscore = []
176 bcetrain = []
177 bcetest = []
178
179 for epoch in range(num_epochs):
180     model.train()
181     for xbatch, ybatch in training_loader:
182         optr.zero_grad()
183
184         prediction = model(xbatch)
185         pred = prediction.mean(1) # float32, ybatch float64
186         #pred = predictionamax(1)
187
188         bce_l = bce_logits(pred, ybatch)
189         bcetrain.append(bce_l.item())
190
191         bce_l.backward()

```

```

184     optr.step()
185     #sched.step()
186
187     model.eval()
188     with torch.no_grad():
189         testpreds = model(xtest)
190         testpred = testpreds.mean(1)
191         #estpred = testpredsamax(1)
192
193         bce_test = bce_logits(testpred, ytest)
194         bcetest.append(bce_test.item()) # ytest is
195         ↪ float64, testpred is float32
196
197         dice_score = dsc(testpred.double(), ytest)
198         dscore.append(dice_score.item())
199
200         if epoch == (num_epochs - 1):
201             te_pred = testpred
202             test_batch = ytest
203
204             tra_pred = pred
205             train_batch = ybatch
206
207             it += 1
208             print("Epoch %s/%s" % (epoch + 1, num_epochs))
209
210             plt.figure(1)
211             train, = plt.plot(bcetrain, 'r')
212             test, = plt.plot(bcetest, 'b')
213             plt.xlabel("Iteration")
214             plt.ylabel("Loss")
215             plt.title("Training and testing loss")
216             plt.legend([train, test], ['Train loss', 'Test loss'])
217             plt.annotate("Final train loss: %s" % (bcetrain[-1]), xycoords =
218             ↪ 'figure fraction', xy = (0.25,0.7))
219             plt.annotate("Final test loss: %s" % (bcetest[-1]), xycoords =
220             ↪ 'figure fraction', xy = (0.25,0.75))
221             print("Final training loss: %.2f" % bcetrain[-1])

```

```

219 print("Final testing loss: %s." % bcetest[-1])
220 plt.savefig("bce_loss_plots", dpi = 500)
221
222 num = 0
223 plt.figure(2)
224 plt.imshow(traj_pred[num, :, :].detach().cpu(), cmap = 'jet')
225 plt.title('Training prediction')
226 plt.colorbar()
227 plt.savefig("train_pred", dpi = 500)
228 plt.figure(3)
229 plt.imshow(train_batch[num, :, :].detach().cpu(), cmap =
   ↪ 'jet')
230 plt.title('Training mask')
231 plt.colorbar()
232 plt.savefig("train_mask", dpi = 500)
233 plt.figure(4)
234 plt.plot(dscore, 'b')
235 plt.xlabel("Iteration")
236 plt.ylabel("DSC")
237 plt.savefig("dice_plot", dpi = 500)
238 print("Final dice coefficient: %s." % dscore[-1])
239
240 ind = 10
241 plt.figure(5)
242 plt.imshow(te_pred[ind, :, :].detach().cpu(), cmap = 'jet')
243 plt.title('Testing prediction')
244 plt.colorbar()
245 plt.savefig("test_pred_1", dpi = 500)
246 plt.figure(6)
247 plt.imshow(test_batch[ind, :, :].detach().cpu(), cmap = 'jet')
248 plt.title('Testing mask')
249 plt.colorbar()
250 plt.savefig("test_mask_1", dpi = 500)
251
252 ind = 20
253 plt.figure(7)
254 plt.imshow(te_pred[ind, :, :].detach().cpu(), cmap = 'jet')
255 plt.title('Testing prediction')

```

```
256 plt.colorbar()
257 plt.savefig("test_pred_2", dpi = 500)
258 plt.figure(8)
259 plt.imshow(test_batch[ind, :, :].detach().cpu(), cmap = 'jet')
260 plt.title('Testing mask')
261 plt.colorbar()
262 plt.savefig("test_mask_2", dpi = 500)
```