

Data Mining Assignment 3

Zhengyuan Zhu
CSE5334 - Data Mining
UNIVERSITY OF TEXAS AT ARLINGTON

May 3, 2020

Problem 1: Logistic Regression

Solution for 1.1

For implement logistic regression, we need to use 'numpy' to make up a toy dataset which generates 1000 training instances in two sets of random data points(500 in each) from multi-variate normal distribution with

$$\mu_1 = [1, 0], \mu_2 = [0, 1.5], \Sigma_1 = \begin{bmatrix} 1 & 0.75 \\ 0.75 & 1 \end{bmatrix}, \Sigma_2 = \begin{bmatrix} 1 & 0.75 \\ 0.75 & 1 \end{bmatrix}$$

and label them with 0 and 1. For the convenience of later steps, we also need to stack two distributions data vertically.

```
1 # build train set
2 mul1, sigma1 = [1, 0], [[1, 0.75], [0.75, 1]]
3 mul2, sigma2 = [0, 1], [[1, 0.75], [0.75, 1]]
4 train_size = 500
5 test_size = 250
6
7 train1 = np.random.multivariate_normal(mean=mul1, cov=sigma1, size=
                                         train_size)
8 train1_label = np.zeros((train_size, 1))
9
10 train2 = np.random.multivariate_normal(mean=mul2, cov=sigma2, size=
                                         train_size)
11 train2_label = np.ones((train_size, 1))
12 X_train = np.vstack([train1, train2])
13 y_train = np.vstack([train1_label, train2_label])
14 print("Train set samples: \n", X_train[:5], X_train[-5:])
15 print("Train set labels: \n", y_train[:5], y_train[-5:])
16
17 test1 = np.random.multivariate_normal(mean=mul1, cov=sigma1, size=
                                         test_size)
18 test1_label = np.zeros((test_size, 1))
19 test2 = np.random.multivariate_normal(mean=mul2, cov=sigma2, size=
                                         test_size)
20 test2_label = np.ones((test_size, 1))
21 X_test = np.vstack([test1, test2])
22 y_test = np.vstack([test1_label, test2_label])
```

Then, we implement the logistic regression through **pytorch**. The implementation of model is fairly simple that only contains one fully connected layer and takes two number as input and output a single number. Then using 'Sigmoid' as activation function and binary cross entropy as loss function.

We also need to define a class for loading the data. The codes are shown below:

```

1 class LogisticReg(torch.nn.Module):
2     def __init__(self):
3         super(LogisticReg, self).__init__()
4         self.fc = torch.nn.Linear(2, 1)
5     def forward(self, x):
6         x = self.fc(x)
7         return F.sigmoid(x)
8
9 class NormDataset(torch.utils.data.Dataset):
10     def __init__(self, x, y):
11         self.len = x.shape[0]
12         self.device = 'cuda' if torch.cuda.is_available() else 'cpu'
13         self.x_data = torch.as_tensor(x, device=self.device, dtype=torch.
14                                         float)
15         self.y_data = torch.as_tensor(y, device=self.device, dtype=torch.
16                                         float)
17     def __getitem__(self, index):
18         return self.x_data[index], self.y_data[index]
19     def __len__(self):
20         return self.len
21
22 loss_func = torch.nn.BCELoss()

```

During the training process, firstly we create a object for normal distribution dataset. And use pytorch dataloader to transform numpy object of train set and test set to torch tensor object. Then use a for loop to iterate 4 different learning rate $\eta = \{1, 0.1, 0.01, 0.001\}$.

For the batch training part, we use stochastic gradient descent(SGD) as optimizer and calculate L1 normalization of gradients as a threshold to early stop the training process. If L1 normalization of model's weight and bias has not change in 10 epochs, then the train process will stop. The codes and experiment results are shown below:

```

1 train_set, test_set = NormDataset(X_train, y_train), NormDataset(X_test,
2                             y_test)
3 train_loader = torch.utils.data.DataLoader(dataset=train_set, batch_size=
4                                     32, shuffle=False)
5 X_test_tsr, y_test_tsr = Variable(torch.from_numpy(X_test).float(),
6                                     requires_grad=False), Variable(torch.
7                                     from_numpy(y_test).float(),
8                                     requires_grad=False)
9
10 writer = SummaryWriter()
11 for lr in learning_rates:
12     model = LogisticReg()
13     optimizer = torch.optim.SGD(model.parameters(), lr=lr)
14     prev_norm, norms, cnt = torch.tensor(0), torch.tensor(0), 0
15     print("Parameters before training:")
16     for name, param in model.named_parameters():
17         if param.requires_grad:

```

Learning rate	1	0.1	0.01	0.001
Accuracy	0.864	0.914	0.912	0.912
Gradient Norm	0.0505	0.1222	0.1340	0.1354
Loss	0.0355	0.0923	0.1027	0.1039
Iterations	51	344	2440	151771

Table 1: Experiment results of batch training

```

13         print(name, param.data)
14     for epoch in tqdm(range(num_epochs)):
15         early_stop = False
16         for i, data in enumerate(train_loader):
17             X_train_tsr, y_train_tsr = data
18             y_pred = model(X_train_tsr)
19             loss = loss_func(y_pred, y_train_tsr)
20             optimizer.zero_grad()
21             loss.backward()
22             optimizer.step()
23             norms = torch.norm(model.fc.weight.grad)+torch.norm(model.fc.bias.grad)
24
25             if prev_norm.data==norms.data and cnt<10:
26                 cnt += 1
27             if cnt==10:
28                 print('Early stopping at {} epoch when norms={}'.format(epoch,
29                                     norms.data))
30                 break
31             writer.add_scalar('Loss/lr='+str(lr), loss, epoch)
32             writer.add_scalar('GradNorm/lr='+str(lr), norms, epoch)
33             prev_norm = norms
34             test_pred = model.forward(X_test_tsr).data.numpy()
35             test_pred = np.where(test_pred>0.5, 1., 0.)
36             acc = accuracy_score(test_pred, y_test_tsr.data.numpy())
37             print("\nParameters after training:")
38             for name, param in model.named_parameters():
39                 if param.requires_grad:
40                     print(name, param.data)
41             print('\nWhen lr={}, the accuracy is {}'.format(lr, acc))
42             print('-----'*10)

```

According to Table 1, it shows that the model achieves the best performance(accuracy: 91.4%) when we set learning rate to 0.1. And the figures of changes of norm of gradient and losses are shown in the Figure 1. It will take more iterations to converge when set a small learning rate.

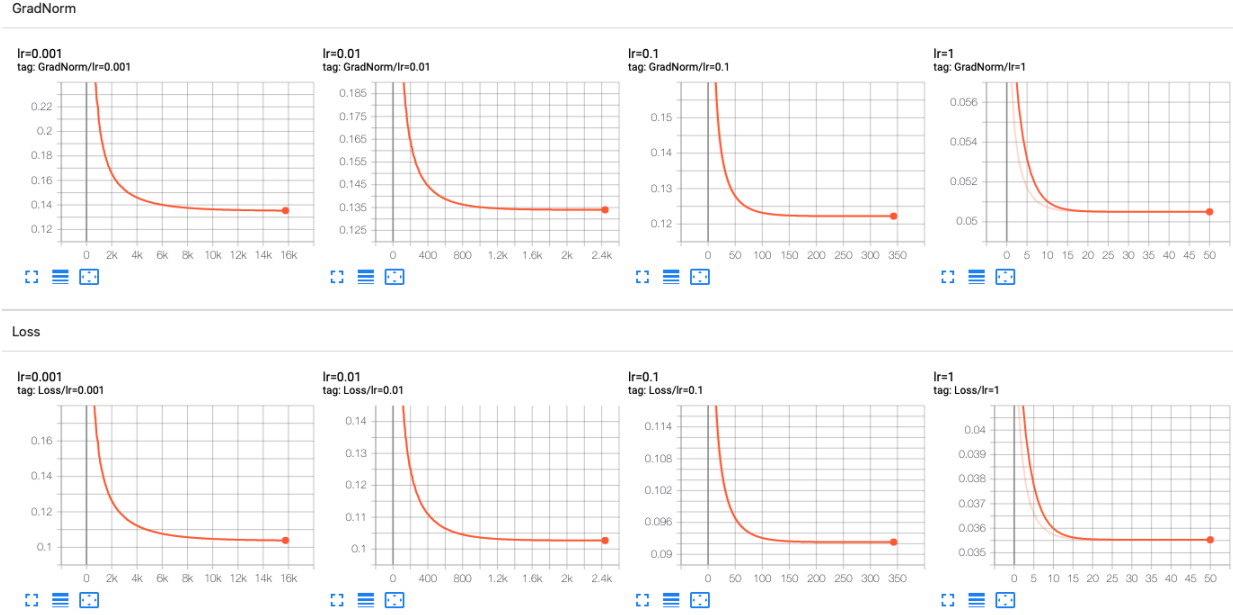


Figure 1: Figures show L1 normalization of gradients and losses w.r.t iteration.

We also visualise the scatter plot of the testing data and the trained decision boundary in Figure 2

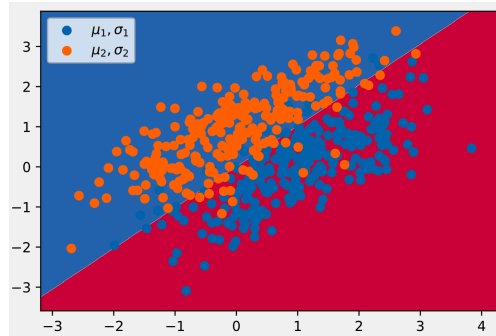


Figure 2: Scatter plot of the testing data and the trained decision boundary of batch training.

Then we do some visualizations of the probability density function. The figure below shows that the different densities with different bandwidth. The red line represents the ground truth distribution and the histogram represents the samples distribution. It shows that the distribution curve becomes more flat when bandwidths get larger and larger.

Solution for 1.2

For the problem 1.2, the methodology essentially is the same as problem 1.1. The only different is that we set batch size from 32 to 1 to make it a online training. The printed log records are shown below:

Learning rate	1	0.1	0.01	0.001
Accuracy	0.552	0.77	0.922	0.911
Gradient Norm	0.0003	0.0748	0.4951	0.6689
Loss	0.0001	0.0446	0.3404	0.4939
Iterations	15	27	103	726

Table 2: Experiment results of online training

```

1 Parameters before training:
2 fc.weight tensor([[ -0.7055,  0.0255]])
3 fc.bias tensor([ -0.2273])
4 Early stopping at 15 epoch when norms=0.003028314560651779
5 Parameters after training:
6 fc.weight tensor([[ -2.8413,  3.3302]])
7 fc.bias tensor([5.6395])
8 When lr=1, the accuracy is 0.552
9 -----
10 Parameters before training:
11 fc.weight tensor([[ -0.0246, -0.4560]])
12 fc.bias tensor([ -0.1986])
13 Early stopping at 27 epoch when norms=0.0748012363910675
14 Parameters after training:
15 fc.weight tensor([[ -2.6917,  3.3682]])
16 fc.bias tensor([2.3081])
17 When lr=0.1, the accuracy is 0.77
18 -----
19 Parameters before training:
20 fc.weight tensor([[ 0.2296, -0.2112]])
21 fc.bias tensor([ -0.6161])
22 Early stopping at 103 epoch when norms=0.49514836072921753
23 Parameters after training:
24 fc.weight tensor([[ -3.1573,  3.4086]])
25 fc.bias tensor([0.2957])
26 When lr=0.01, the accuracy is 0.922
27 -----
28 Parameters before training:
29 fc.weight tensor([[0.4855, 0.6238]])
30 fc.bias tensor([0.4780])
31 Early stopping at 726 epoch when norms=0.6689199209213257
32 Parameters after training:
33 fc.weight tensor([[ -3.3658,  3.3036]])
34 fc.bias tensor([ -0.0130])
35 When lr=0.001, the accuracy is 0.91
36 -----

```

The experiment results in Table 2 shows the same conclusion that smaller learning rate will lead to more iterations to converge. And the best performance(Accuracy 92.2%) appears when setting the learning rate to 0.01 during online training. We should notice that although the loss of *learningrate* = 1 is pretty low(0.0001), the accuracy is the worst among all the experiments.

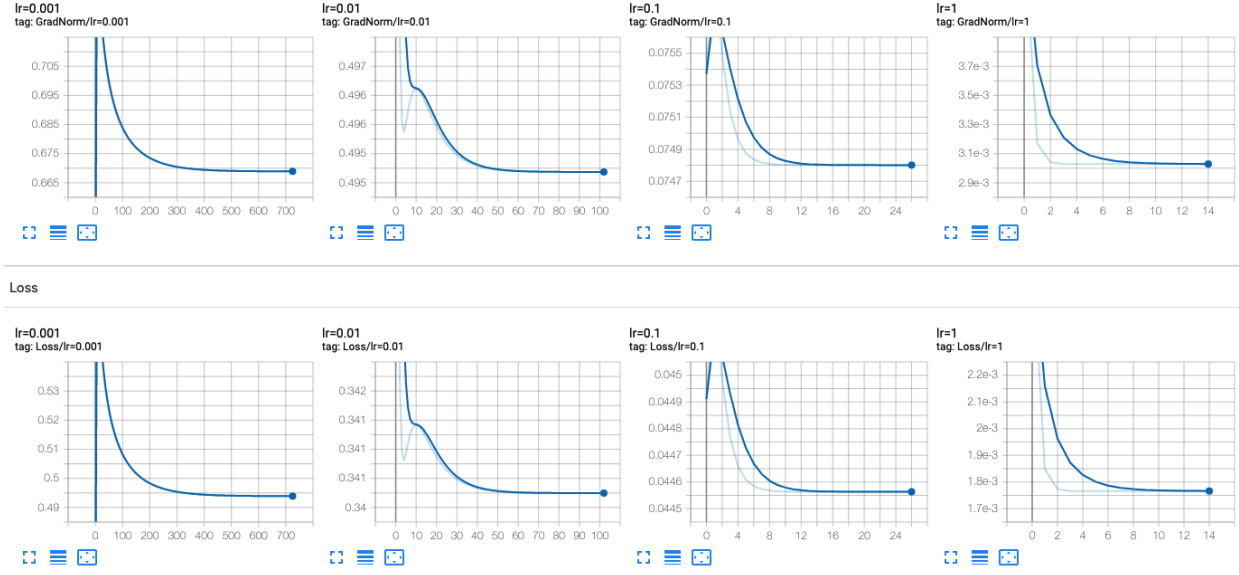


Figure 3: Figures show L1 normalization of gradients and losses w.r.t iteration.

We also visualise the scatter plot of the testing data and the trained decision boundary in Figure 3

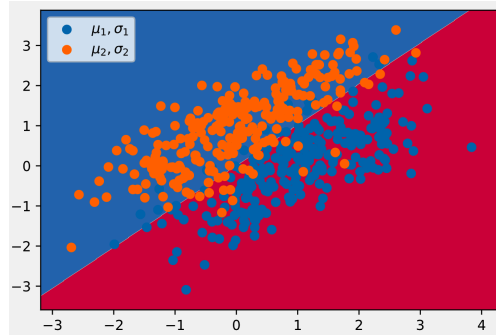


Figure 4: Scatter plot of the testing data and the trained decision boundary of online training.

Problem 2 Artificial Neural Network

Solution 2.1

Use *PyTorch* to design a Neural Network for Image Classification on *CIFAR-10 dataset*. First, we need to convert color images to gray scale and make sure that the pixel intensities are normalized to stay in $[0, 1]$. It is very convenient to implement transform converter via *torch.vision*, but do not forget there are three channels(RGB) in *CIFAR-10 datasets* and

only one channel in *fashion MNIST*. As a result, we will use two different transformers that take images as input and transform them into tensor then normalize those tensor:

```
1 cifar10_transform = transforms.Compose(
2     [transforms.ToTensor(),
3      transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))]
4 )
5 fashion_transform = transforms.Compose(
6     [transforms.ToTensor(),
7      transforms.Normalize([0.5], [0.5])]
8 )
```

We modify the sample codes to complete the first problem. For initiating the Net class of problem 2.1, the first fully connected layer should change is input size from *image_cifar10_width*image_cifar10_height* to *image_cifar10_width*image_cifar10_height*color_channels*. We don't need to modify too much sample codes to finish the first task. The result are shown in Figure 5 and Figure 6.

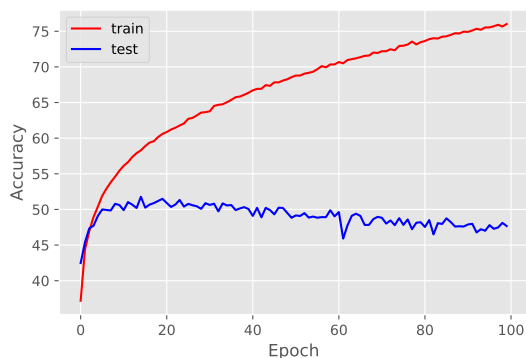


Figure 5: Accuracy of vanilla neural network

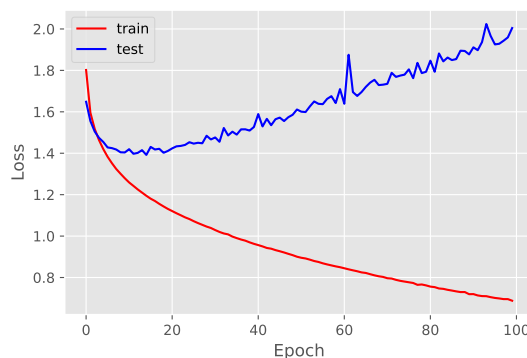


Figure 6: Loss of vanilla neural network

Solution 2.2

We can simply add a hidden layer through 'torch.nn.Linear()' function to '__init__' in the Net class. The results in Figure 7 and Figure 8 show that the one more hidden layer model outperforms the vanilla one. The shape of figures pretty much look the same, and the accuracy increased by 5% and although loss value is higher at the early stage.

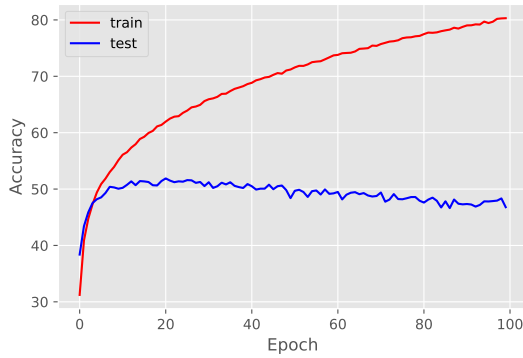


Figure 7: Accuracy of neural network has one hidden layer

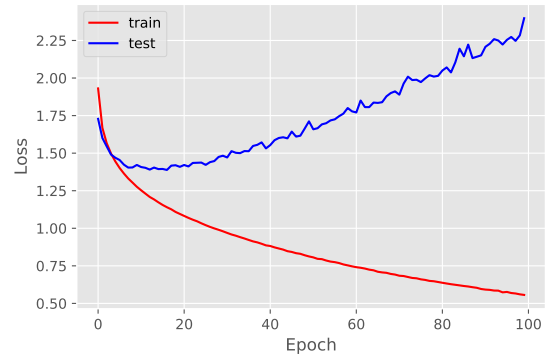


Figure 8: Losses of neural network has one hidden layer

Solution 2.3

After add dropout at the rate of 0.2 for the hidden layer. The idea under the hood is simple. It deactivates the neurons randomly at each training step instead of training the data on the original network, train the data on the network with dropped out nodes. In the next iteration of the training step, the hidden neurons which are deactivated by dropout changes because of its probabilistic behavior. In this way, by applying dropout i.e...deactivating certain individual nodes at random during training we can simulate an ensemble of neural network with different architectures.

The result shows that although it solves the over-fitting to some extent. But the performance of model is effected by dropout. It seems like the drop rate should be adjust to a reasonable value to make it better performance.

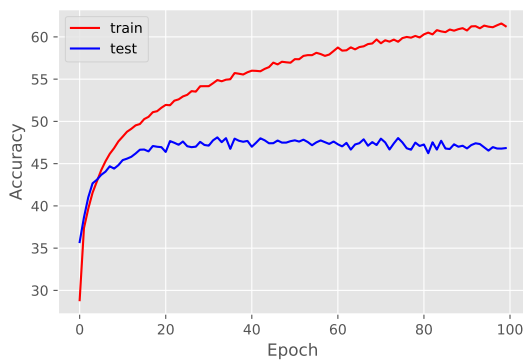


Figure 9: Accuracy of vanilla neural network of sub-task 3

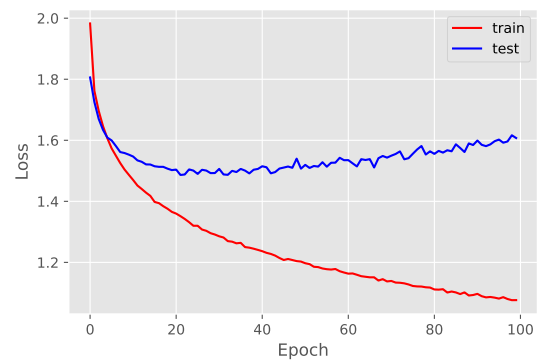


Figure 10: Loss of vanilla neural network of sub-task 3

Solution 2.4

In this problem, we freely change the architecture of ANN. So that we add two more hidden layer and *BatchNormalization* to the *Net* class. But unfortunately the performance have not improved. The final modified *Net* class are shown below:

```
1 class Net(torch.nn.Module):
2     def __init__(self, one_more_hid=False, freestyle=False, vanilla=False,
3                 logistic=False):
4         super(Net, self).__init__()
5         self.vanilla = vanilla
6         self.logistic = logistic
7         self.dropout = dropout
8         self.fc1 = torch.nn.Linear(image_cifar10_width*
9                                   image_cifar10_height*
10                                  color_channels, hidden_nodes)
11
12         self.fc1_fashion = torch.nn.Linear(image_fashion_mnist_width*
13                                             image_fashion_mnist_height,
14                                             hidden_nodes)
15
16         self.fc2 = None
17         self.fc3 = None
18         self.bn1 = None
19         self.bn2 = None
20         self.bn3 = None
21         self.non_hidden = torch.nn.Linear(image_fashion_mnist_width*
22                                           image_fashion_mnist_height,
23                                           class_number_fashion_mnist)
24
25         if one_more_hid:
26             self.fc2 = torch.nn.Linear(hidden_nodes, hidden_nodes)
27
28         if freestyle:
29             self.bn1 = torch.nn.BatchNorm1d(num_features=64)
30             self.fc2 = torch.nn.Linear(hidden_nodes, 128)
31             self.bn2 = torch.nn.BatchNorm1d(num_features=128)
32             self.fc3 = torch.nn.Linear(128, 64)
33             self.bn3 = torch.nn.BatchNorm1d(num_features=64)
34
35         if logistic:
36             self.log_fc = torch.nn.Linear(image_fashion_mnist_width*
37                                           image_fashion_mnist_height,
38                                           1)
39
40         self.out_cifar10 = torch.nn.Linear(hidden_nodes,
41                                             class_number_cifar10)
42
43         self.out_fashion_mnist = torch.nn.Linear(hidden_nodes,
44                                                  class_number_fashion_mnist)
45
46     def forward(self, x):
47         if dataset_name == "cifar10":
48             x = x.view(-1, image_cifar10_width*image_cifar10_height*
49                       color_channels)
50
51             x = self.fc1(x)
52             if self.bn1:
53                 x = self.bn1(x)
54             x = torch.nn.functional.relu(x)
55             x = torch.nn.Dropout(self.dropout)(x)
56             if self.fc2:
```

```

37         x = self.fc2(x)
38         if self.bn2:
39             x = self.bn2(x)
40         x = torch.nn.functional.relu(x)
41         x = torch.nn.Dropout(self.dropout)(x)
42     if self.fc3:
43         x = self.fc3(x)
44         if self.bn3:
45             x = self.bn3(x)
46         x = torch.nn.functional.relu(x)
47         x = torch.nn.Dropout(self.dropout)(x)
48     return self.out_cifar10(x)

```

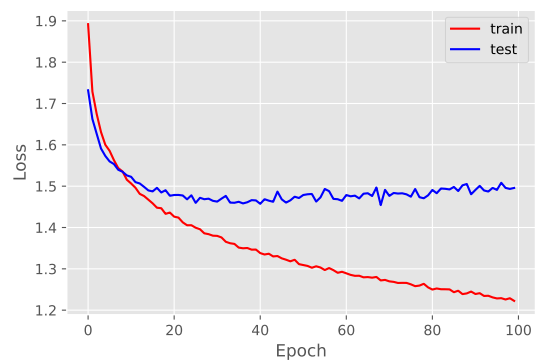
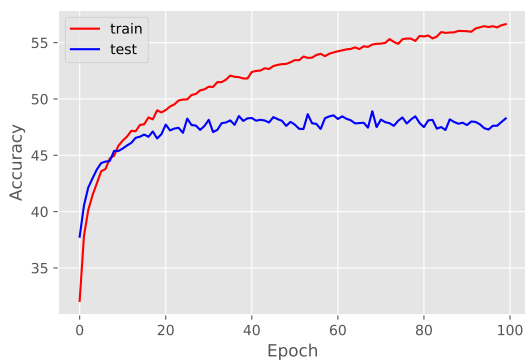


Figure 11: Accuracy of freestyle neural net- Figure 12: Loss of freestyle neural network work

Solution 2.5

When we dive into *fashion MNIST* dataset. There are nothing much to change for the codes. The results shown in Figure 13 and 14.

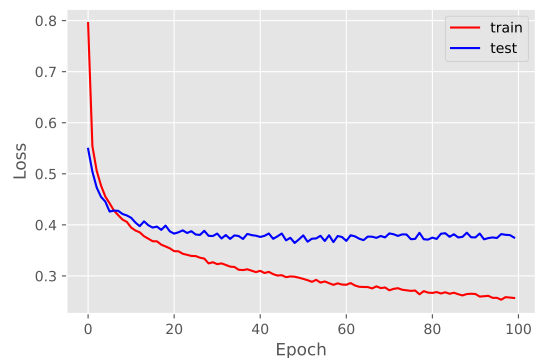
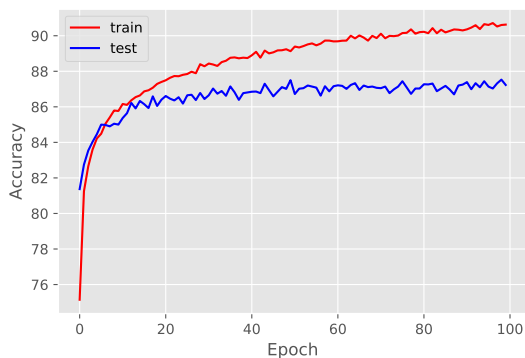


Figure 13: Accuracy of neural network for fashion MNIST Figure 14: Losses of vanilla neural network for fashion MNIST

Solution 2.6

If we remove all the hidden layer and train the model. The results shown in Figure 13 and 14.

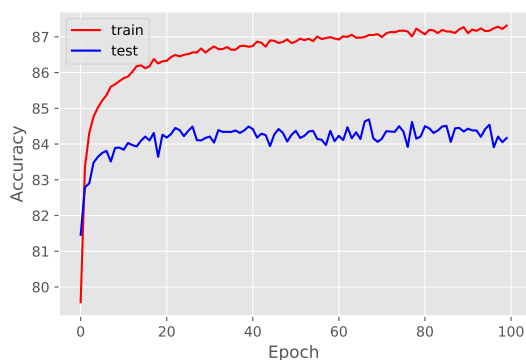


Figure 15: Accuracy of none layer neural network

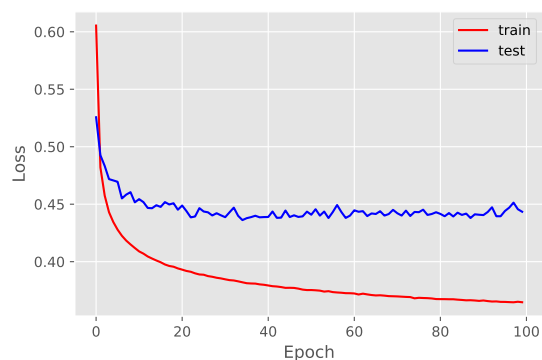


Figure 16: Loss of none layer neural network

The trained model contains the weight directly from input to output nodes. The visualization of the learned weights to the image dimension show that the weights have pretty much the same shape with train set images.

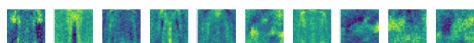


Figure 17: Network weights at epoch 0

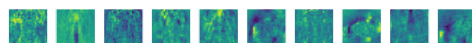


Figure 18: Network weights at epoch 25

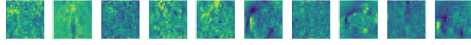


Figure 19: Network weights at epoch 75

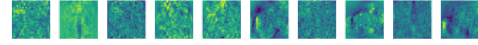


Figure 20: Network weights at epoch 95

Solution 2.7

For this problem, we select the first two classes and apply ANN model. The ROC curve along with the other results are shown below separately.

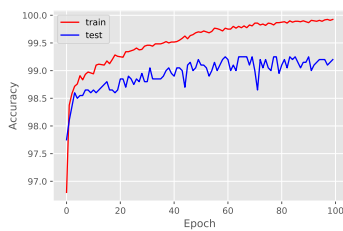


Figure 21: Accuracy of ANN

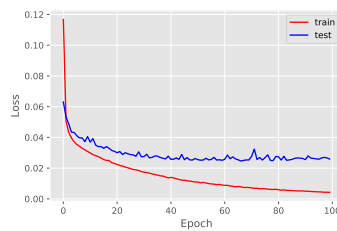


Figure 22: Loss of ANN

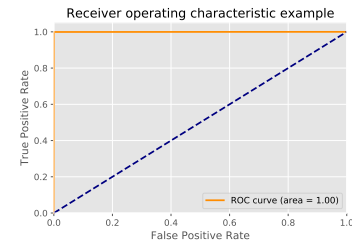


Figure 23: ROC curve of ANN

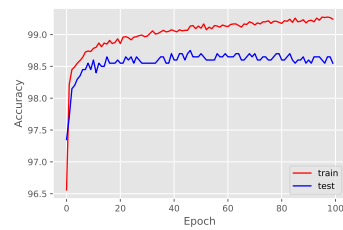


Figure 24: Accuracy of Logistic Regression

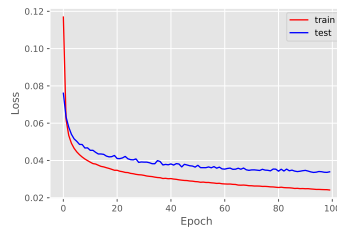


Figure 25: Loss of Logistic Regression

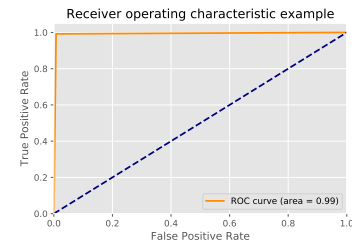


Figure 26: ROC curve of Logistic Regression

From those figures, we can easily obtain that the performance of both two models have good validity. The neural network method has better accuracy compare to logistic regression.

Reference

- <https://medium.com/biaslyai/pytorch-linear-and-logistic-regression-models-5c5f0da2c>
- https://scikit-learn.org/stable/auto_examples/model_selection/plot_roc.html#sphx-glr-auto-examples-model-selection-plot-roc-py
- https://scikit-learn.org/stable/modules/model_evaluation.html#roc-metrics