

GR 5241 Final Project Milestone II Report

Contents

<i>Milestone I: Machine Learning</i>	<i>1</i>
<i>Milestone II: Deep Learning.....</i>	<i>2</i>
<i>3. Train with ANN.....</i>	<i>2</i>
<i>(a) Cross-Entropy Error Plot</i>	<i>3</i>
<i>(b) Misclassification Error Plot.....</i>	<i>4</i>
<i>(c) Best Result Visualization</i>	<i>5</i>
<i>(d) Different Hyper-parameters.....</i>	<i>6</i>
<i>4. Train with CNN.....</i>	<i>6</i>
<i>(a) Cross-Entropy Error Plot</i>	<i>7</i>
<i>(b) Misclassification Error Plot.....</i>	<i>8</i>
<i>(c) Best Result Visualization</i>	<i>9</i>
<i>(d) Different Hyper-parameters.....</i>	<i>9</i>
<i>5. Favorite Deep Learning Architecture</i>	<i>12</i>
<i>More about Deep Learning.....</i>	<i>15</i>

Side Note:

As mentioned in Ed that this report will be focusing on Milestone II

Milestone II: Deep Learning

In this part, I had implemented two different artificial neural network structures (ANN and CNN) with the MNIST training and test sets.

A brief introduction about MNIST dataset: The MNIST database of handwritten digits is one of the most used datasets for training various image processing systems and machine learning algorithms. MNIST has a training set of 60,000 examples, and a test set of 10,000 examples.

3. Train with ANN

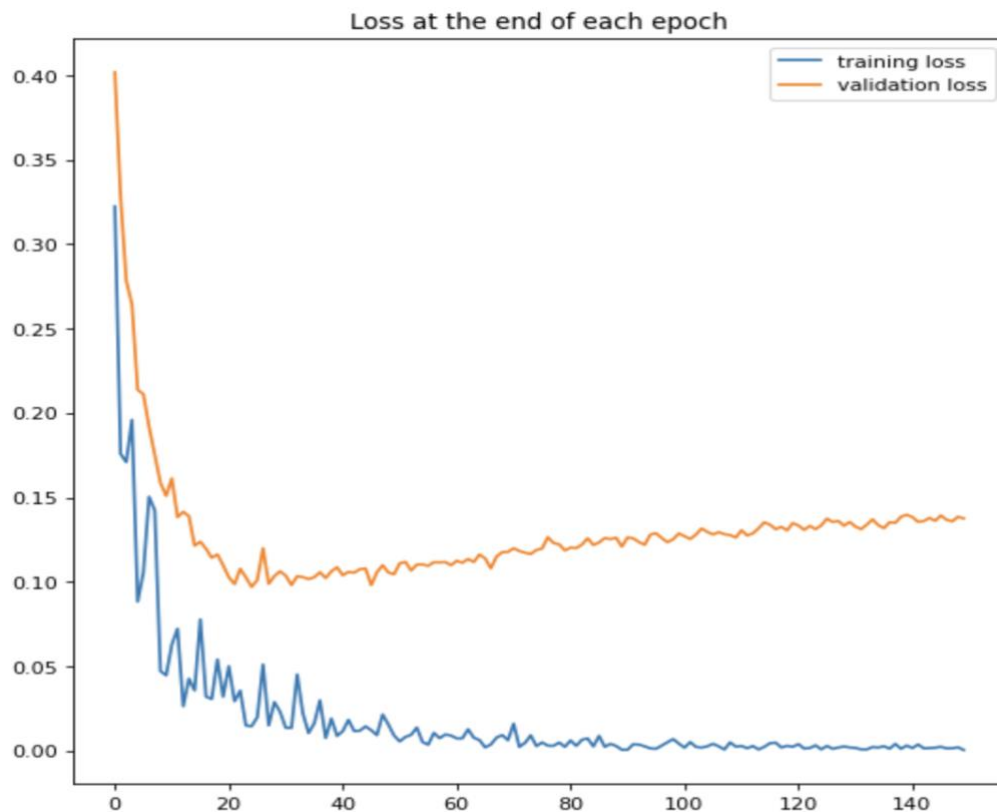
In this part I have constructed my model as required: a single layer neural network with 100 hidden units that is with architecture (input features: 784, hidden layer units: 100 and output features: 10). Below is the screenshot of my model and this model remains valid until **Train with CNN**.

```
1 class Model(nn.Module):
2
3     def __init__(self, in_features = 784, h1 = 100, out_features = 10):
4         # one single layer neural network
5         super().__init__()
6         self.fc1 = nn.Linear(in_features, h1)
7         self.out = nn.Linear(h1, out_features)
8
9     def forward(self, x):
10
11         x = F.relu(self.fc1(x))
12         x = self.out(x)
13
14         return F.log_softmax(x, dim=1)
```

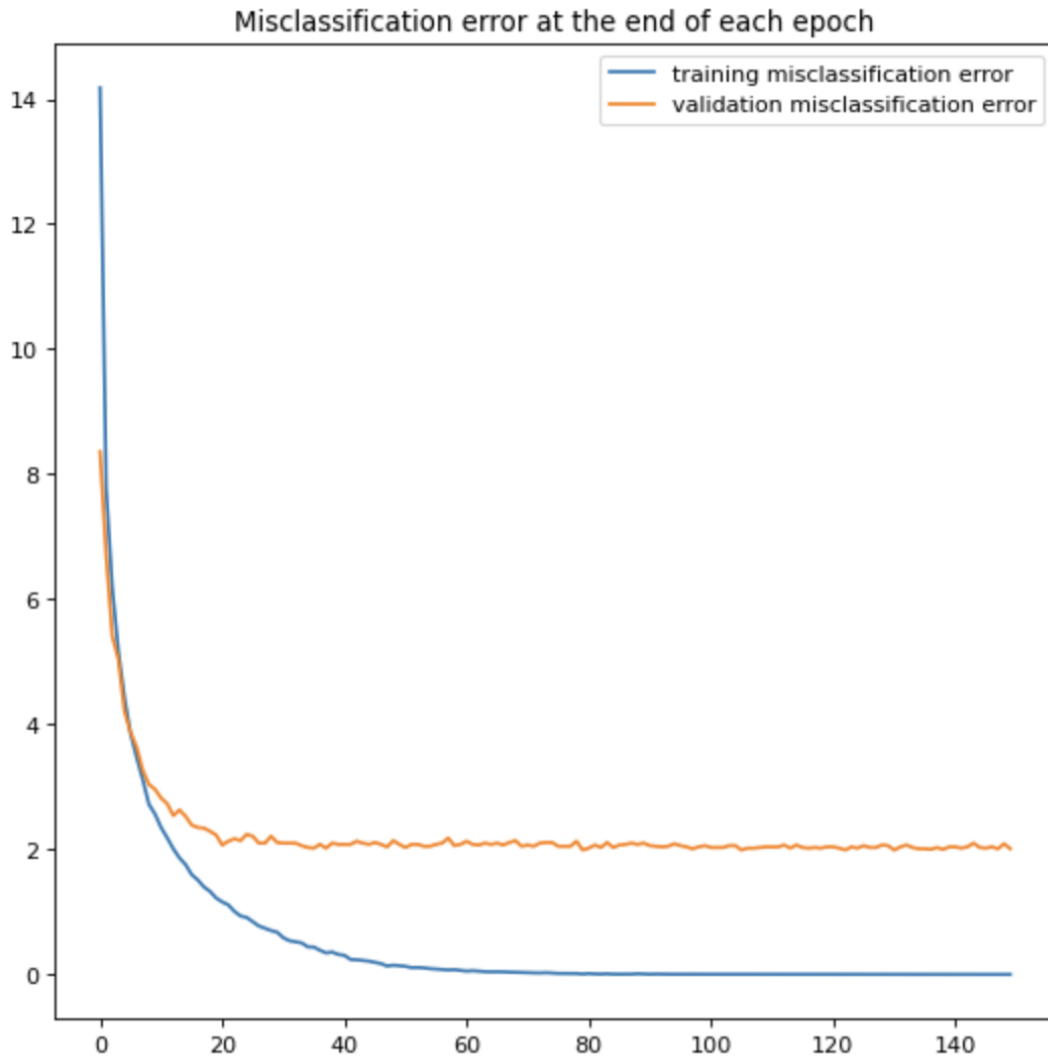
Firstly, for the default setup, I choose SGD as my optimizer with learning rate = 0.1 as our baseline model. Then I have trained this model 6 times with 150 epochs each time, and based on the results for the validation curve, the misclassification error converges during the first

training at the 35th epochs with a peak test accuracy 98.03% and so is the loss error reaching the lowest point at the 35th epochs with a lowest value around 0.10 and starts to increase after that point. As the training process keeps going, we notice that for the validation curve the loss error keeps increasing until it converges to around 0.14 in the 5th training and the misclassification error tends to be stable during the rest trainings. Unsurprisingly, for the training curve, the loss error keeps having the decreasing trends and the misclassification error goes to 0 around the 80th epochs. Since here we are using the SGD optimizer, the training curve for loss error has fluctuations instead of having the smooth curve. Below are the loss error / misclassification error plots for the first training:

(a) Cross-Entropy Error Plot



(b) Misclassification Error Plot

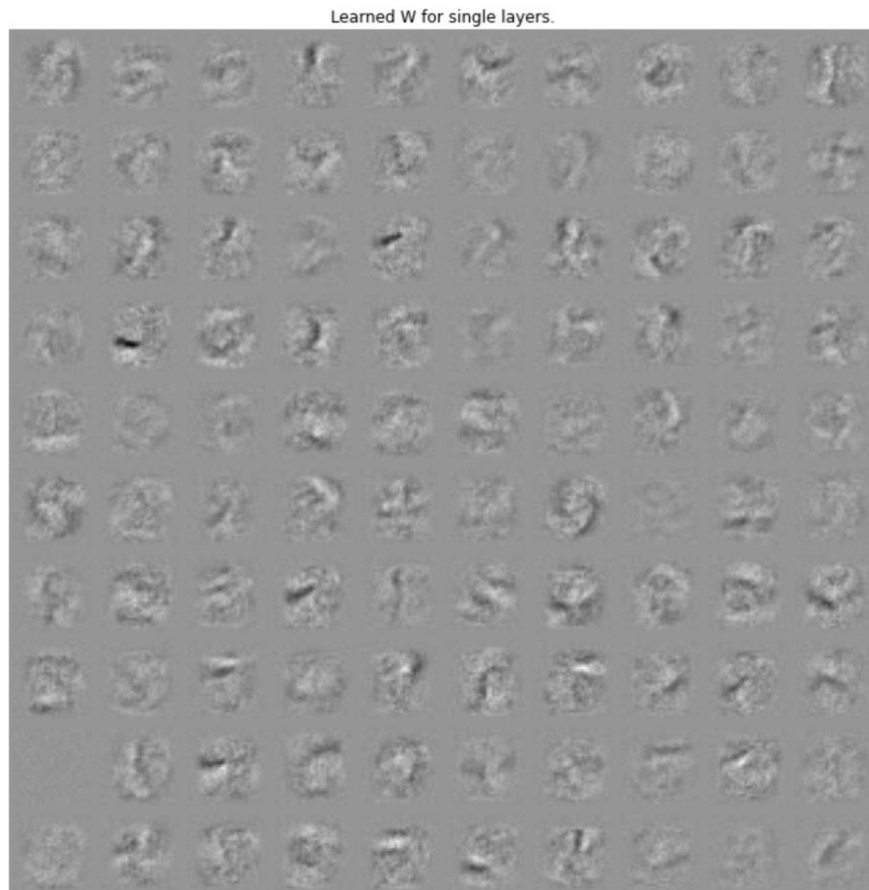


By comparing to those two different types of plots we can find out that for plots in (a), the validation loss starts to increase after a point. In my model at learning rate 0.1, it starts the increasing trend in the middle of first training at the 35th epoch and loss gradually getting more and more during the rest training.

For the plots in (b), we will find out that the misclassification error (in percentage) eventually converged to a specific percentage, in this model it's around 2%, and it does not change much as we keep repeating more trainings.

As a result, this model obtains its best prediction in the first training. Since the ANN model reaches its best performance at the 35th epochs with test accuracy of 98.03%, the below visualization of my best results of the learned weight as one hundred 28 by 28 images will be the result from the first training.

(c) Best Result Visualization



The ANN takes flattened tabular inputs and therefore loses the information stored in the spatial features. Based on the plot above, it is hard to recognize any patterns or structures by human eyes.

(d) Different Hyperparameters

Concerned with the computational cost, I did not train the rest model for more than 5 times as it is a result based on the findings from part a and b, the ANN model is able to converge during the first training.

In this section I have tried 6 different hyperparameters: learning rate of 0.01, 0.2 and 0.5; momentum of 0.0, 0.5 and 0.9. As a conclusion, we can find out that as the learning rate increases, both of our errors (loss, misclassification) will have a faster rate to converge. For example, at learning rate 0.01, it takes more than 150 epochs for our model to reach the convergence and yet for learning rate 0.2 and 0.5, our model takes about 10 epochs and 4 epochs to be able to converge. After the convergent point, loss error starts to increase, and misclassification error remains stable. For the momentum, we can find out it plays the same role as learning rate does, that is the increment of momentum brings a faster rate for our model to touch the converge point.

4. Train with CNN

In this part I have constructed my model as required: one 2-D convolutional layer, Relu activation function and Maxpooling with appropriate hyperparameters, below is the screenshot of my CNN model.

```

1 class ConvolutionalNetwork(nn.Module):
2     def __init__(self):
3         # one 2-D convolutional layers
4         super().__init__()
5         self.conv1 = nn.Conv2d(1, 16, 3, 1)
6         self.out = nn.Linear(13*13*16,10)
7
8     def forward(self, X):
9         # one 2-D convolutional layers -> Relu activation -> Maxpooling
10        X = F.relu(self.conv1(X))
11        X = F.max_pool2d(X, 2, 2)
12        X = X.view(-1,13*13*16)
13        X = self.out(X)
14        return F.log_softmax(X, dim=1)

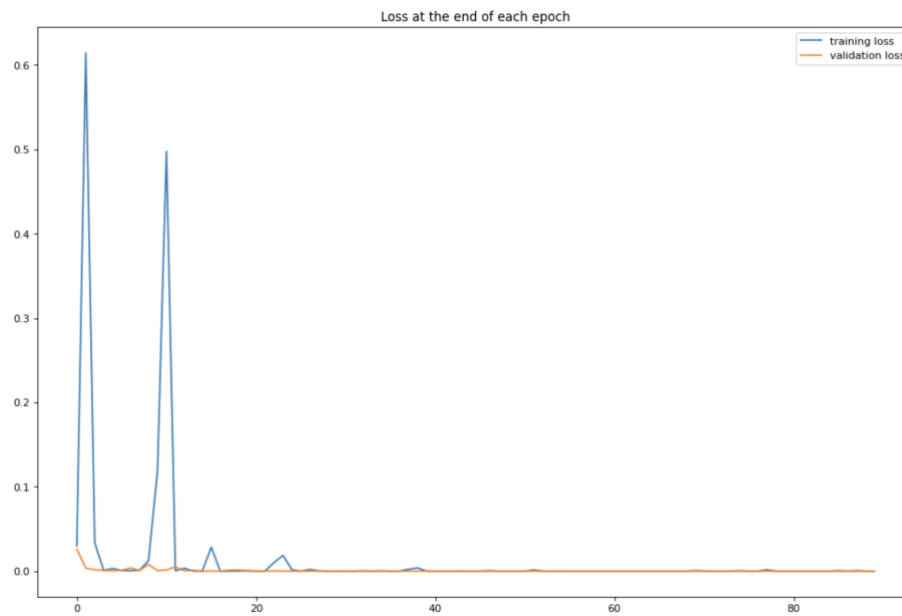
```

The CNN model performs better than my previous ANN model and yet it is more computationally expensive hence I only trained CNN model once and with 90 epochs per time. For the default setup, I choose SGD as my optimizer with learning rate = 0.10 as my baseline model.

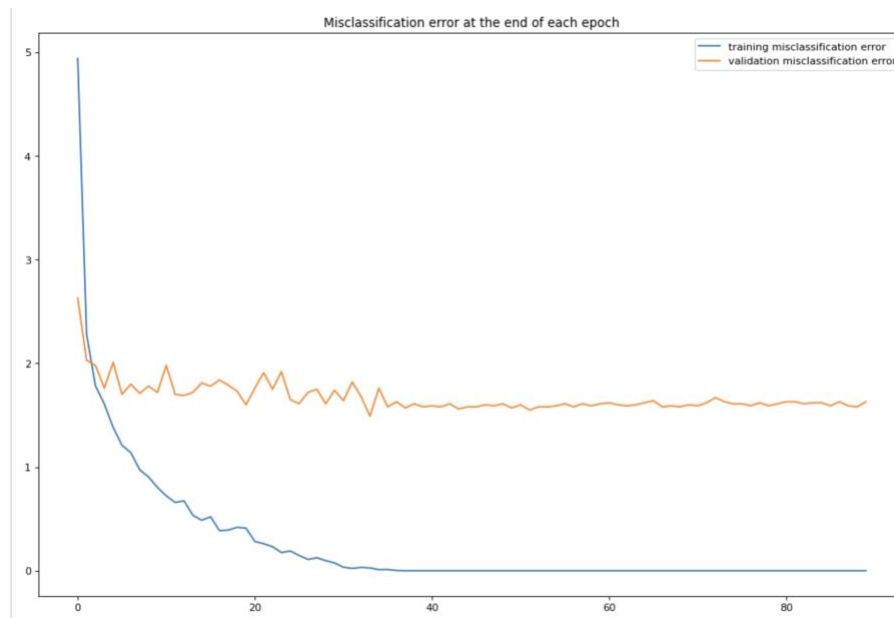
ANN uses weights and an activation function for the bulk of its method as it reconstructs how a brain's neural network works. Comparatively, there is no neuron or weights in CNN. CNN instead casts multiple layers on images and uses filtration to analyze image inputs.

By comparing to plots from ANN model, CNN model takes less epochs to converge and it has a better performance than ANN's first training. From the plots, CNN seems to converge in the very beginning of the training process that it for the training set it converges around the 20th epochs and for the testing set it converges around the 10th epochs. Also, as a comparison to the ANN model: the misclassification error for the first training of ANN is slight above 2% and yet the misclassification error for the first training of CNN is slight below 2%.

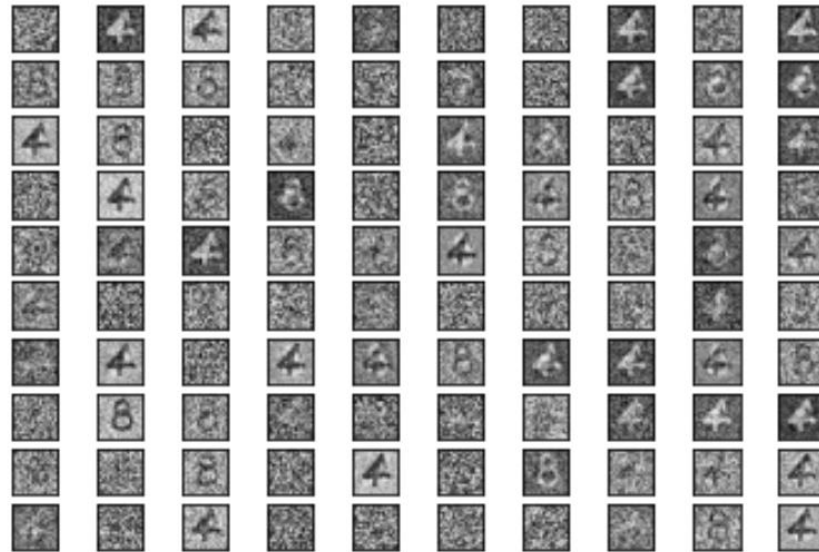
(a) Cross-Entropy Error Plot



(b) Misclassification Error Plot



(c) Best Result Visualization



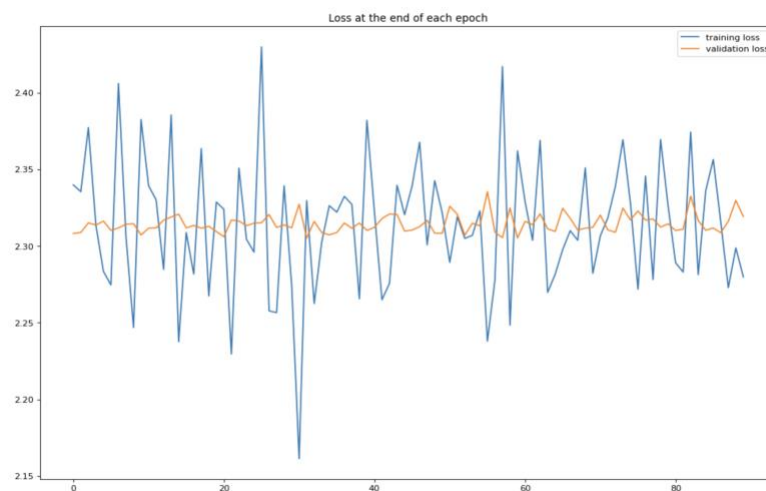
When using CNN, these spatial features are extracted from image input. This makes the CNN model ideal when thousands of features need to be extracted. Instead of having to measure each individual feature, and so CNN gathers these features on its own. Based on the plot above, we can find out that CNN exhibits learned features with some structures.

(d) Different Hyper-parameters

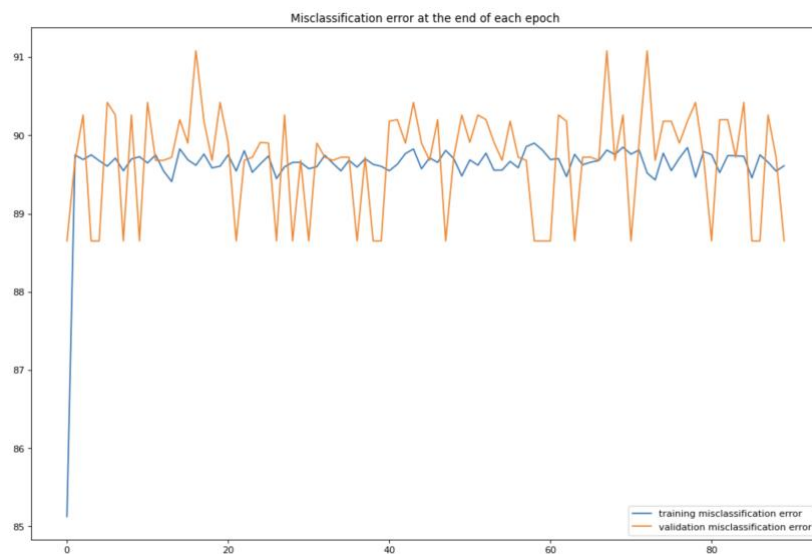
In this section I have also tried 6 different hyperparameters: learning rate of 0.01, 0.2 and 0.5; momentum of 0.0, 0.5 and 0.9. There are two differences between CNN and ANN towards parameters tuning. First, CNN is able to converge at learning rate of 0.01 while ANN could not, and CNN touches the convergent point around 30th epochs. On top of that, CNN shows the same trend as ANN does within a certain extent and yet its performance drops drastically when the learning rate or the momentum is bigger than a threshold. That is, we find out as the learning rate / momentum increases, both of our errors (loss, misclassification) will be unable to reach the

convergence and those two errors fluctuate around some values. For example, at learning rate 0.5, the validation loss error is bouncing around 2.3 and training loss error is bouncing between 2.15 and 2.45. And it shows any sign to converge even till the last epochs. Besides, the misclassification error is high for both training and validation curve that training curve is bouncing around 90% and validation is bouncing between 88% to 91%. Also, it shows no sign to converge. Below are the error plots at learning rate of 0.5 for CNN model:

Loss error (learning rate 0.5)

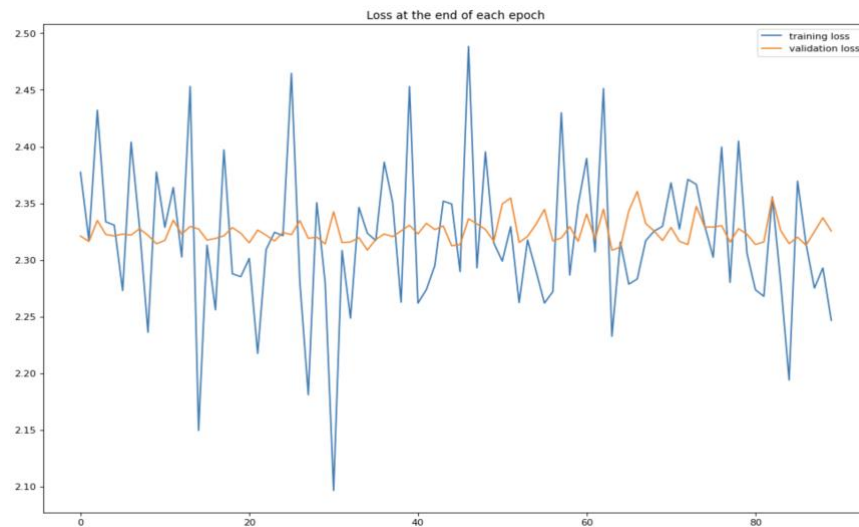


Misclassification error (learning rate 0.5)

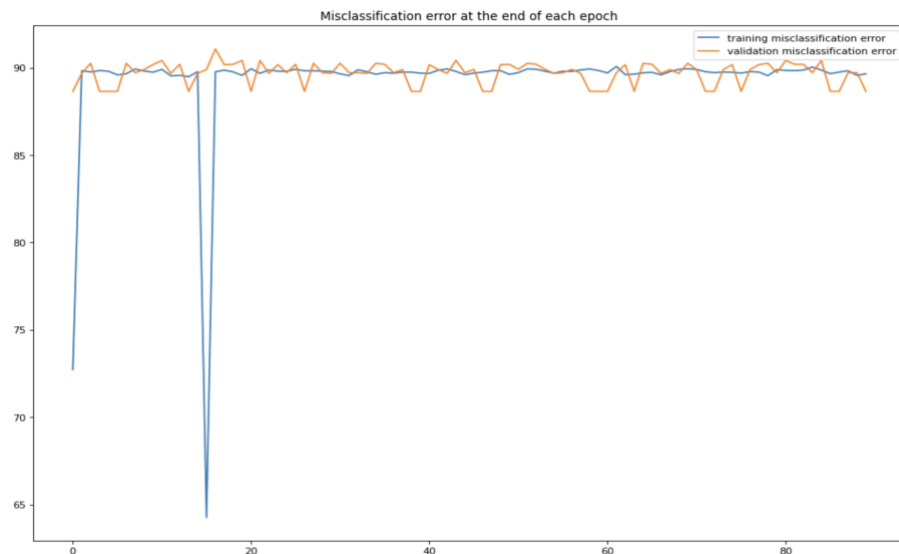


The same phenomenon shows in the momentum too that for the momentum of 0.9 the validation loss error is bouncing around 2.33 and training loss error is bouncing between 2.15 and 2.50. And it shows any sign to converge even till the last epochs. Besides, the misclassification error is high for both training and validation curve that both curves are bouncing around 90%. Below are the error plots at momentum of 0.9 for CNN model:

Loss error (momentum: 0.9)



Misclassification error (momentum: 0.9)



5. Favorite Deep Learning Architecture

The main goal in this step is to beat the performance of SVM with Gaussian Kernel, that is to have a test error rate lower than 1.4%. Here, I construct two CNN models with different output channels towards convolutional layers and each of them have two convolutional layers, two pulling layers and two fully connected layers.

The first one:

```
1  ## two 2-D convolutional layers,two fully connected hidden layers
2
3  class ConvolutionalNetwork(nn.Module):
4      def __init__(self):
5          super().__init__()
6          self.conv1 = nn.Conv2d(1, 16, 3, 1)
7          self.conv2 = nn.Conv2d(16, 25, 3, 1)
8          self.fc1 = nn.Linear(5*5*25, 200)
9          self.fc2 = nn.Linear(200, 100)
10         self.out = nn.Linear(100,10)
11
12     def forward(self, X):
13         X = F.relu(self.conv1(X))
14         X = F.max_pool2d(X, 2, 2)
15         X = F.relu(self.conv2(X))
16         X = F.max_pool2d(X, 2, 2)
17         X = X.view(-1, 5*5*25)
18         X = F.relu(self.fc1(X))
19         X = F.relu(self.fc2(X))
20         X = self.out(X)
21         return F.log_softmax(X, dim=1)
```

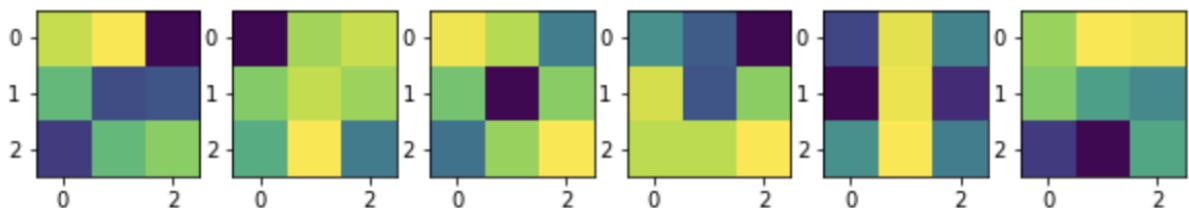
The second one:

```
1 class ConvolutionalNetwork(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.conv1 = nn.Conv2d(1, 6, 3, 1)
5         self.conv2 = nn.Conv2d(6, 16, 3, 1)
6         self.fc1 = nn.Linear(5*5*16, 120)
7         self.fc2 = nn.Linear(120, 84)
8         self.fc3 = nn.Linear(84, 10)
9
10    def forward(self, X):
11        X = F.relu(self.conv1(X))
12        X = F.max_pool2d(X, 2, 2)
13        X = F.relu(self.conv2(X))
14        X = F.max_pool2d(X, 2, 2)
15        X = X.view(-1, 5*5*16)
16        X = F.relu(self.fc1(X))
17        X = F.relu(self.fc2(X))
18        X = self.fc3(X)
19        return F.log_softmax(X, dim=1)
```

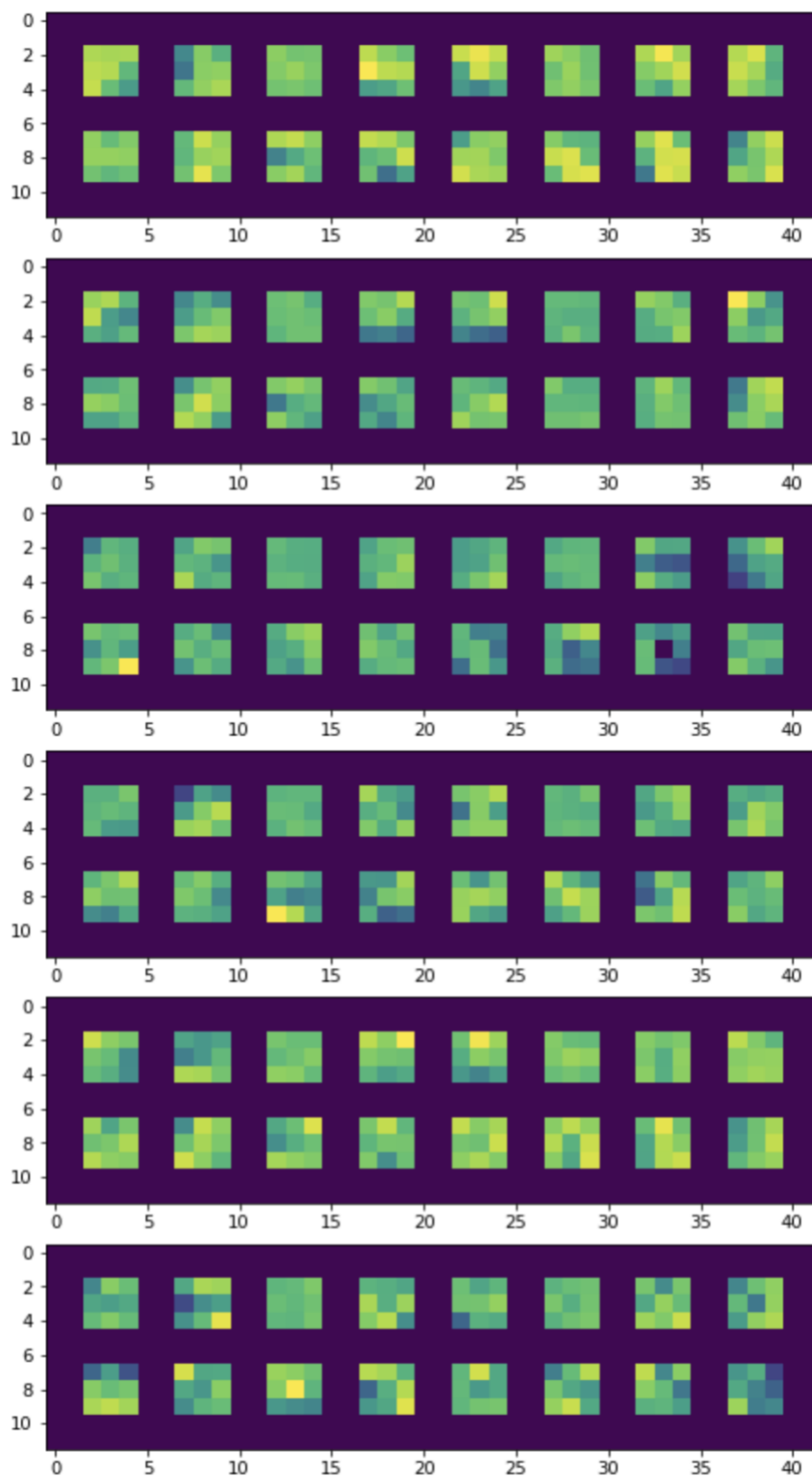
Turns out the second model has a better performance than the first model and it only takes about 3 epochs to converge as its loss error is practically 0 and the final test accuracy is 1.33% which is smaller than 1.4%.

Since for the new construction I have used multilayer Convolutional Network, below I will attach how the filter looks like in the CNN model.

First filters:



Second filters:



More about Deep Learning

In this part, we need to train our customized dataset through at least two of my favorite Deep Learning Algorithms, here is my choice: ANN and CNN with optim.Adam as my optimizer. The difference between this part and previous part is that I am using Adam as my optimizer and using accuracy for my plots instead of using misclassification error (in percentage). Therefore, I did not test how momentum will affect my model as Adam don't take momentum as one of its parameters.

As a brief introduction about our customized dataset, there are 3 txt files and have been split into 3 parts: train set contains 20,000 lines and validation set and test set are both containing 5000 lines in the same format. Each line contains 1569 coordinates, with the first 784 real-valued numbers correspond to the 784-pixel values for the first digit, and the next 784 real valued numbers correspond to the pixel values for the second digit.

A brief description about my workflow:

First of all, I define two customized dataset functions in order to turn the txt file into a similar dataset structure like MNIST dataset. As ANN and CNN takes different structured data inputs, that is ANN deals with flattened tabular data and CNN deals with data which keeps its spatial features, thus my two customized dataset functions are defined differently.

On top of that I build my ANN model and CNN model to train the dataset respectively and test my model with different parameters (learning rate of 0.2 and 0.5). Since here we are using the validation set instead of the test set like we did before, the validation curve below then stands for the test results for the validation set and test accuracy is calculated based on the test set.

Customized dataset function ANN

```
1 class CustomerDataset(Dataset):
2     def __init__(self, file_path):
3         file_out = pd.read_table(file_path, sep=',', names=range(1569))
4         X = file_out.iloc[:, :1568].values
5         y = file_out.iloc[:, 1568].values
6         # Standardize y
7         sc = StandardScaler()
8         X = sc.fit_transform(X)
9         # convert to tensor
10        self.X = torch.tensor(X, dtype=torch.float32)
11        self.y = torch.tensor(y)
12
13    def __len__(self):
14        return len(self.y)
15
16    def __getitem__(self, idx):
17        return self.X[idx], self.y[idx]
18
19
```

```
1 train_q5 = CustomerDataset(train_path)
2 val_q5 = CustomerDataset(val_path)
3 test_q5 = CustomerDataset(test_path)
```

```
1 train_loader = DataLoader(train_q5, batch_size=100, shuffle=True)
2 val_loader = DataLoader(val_q5, batch_size=100, shuffle=True)
3 test_loader = DataLoader(test_q5, batch_size=500, shuffle=False)
```

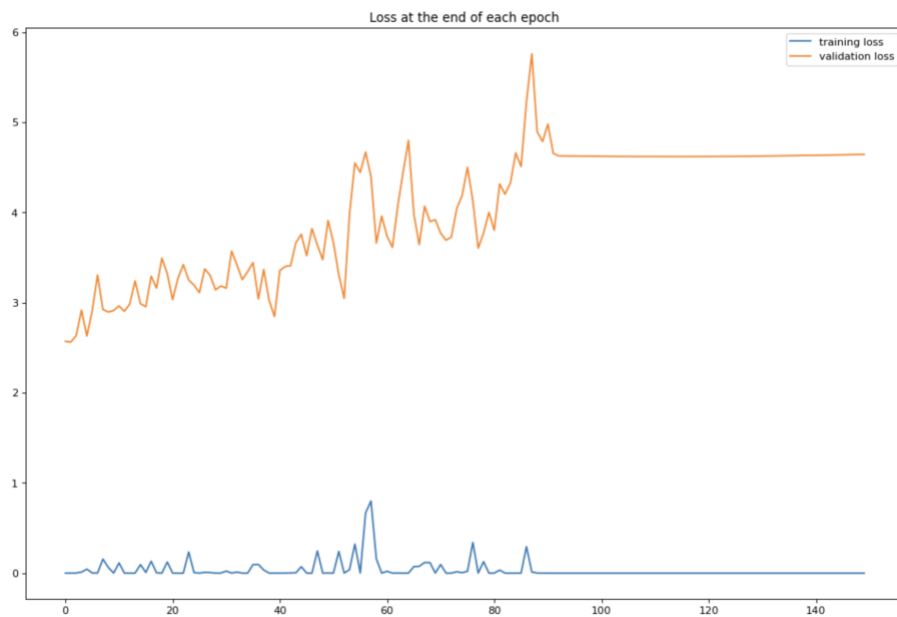
Starting with the ANN model, in this architecture, I have used 3 hidden layers with 784, 120, 84 units respectively and set up learning rate of 0.001 as my baseline model with 150 epochs per training.

Model structure

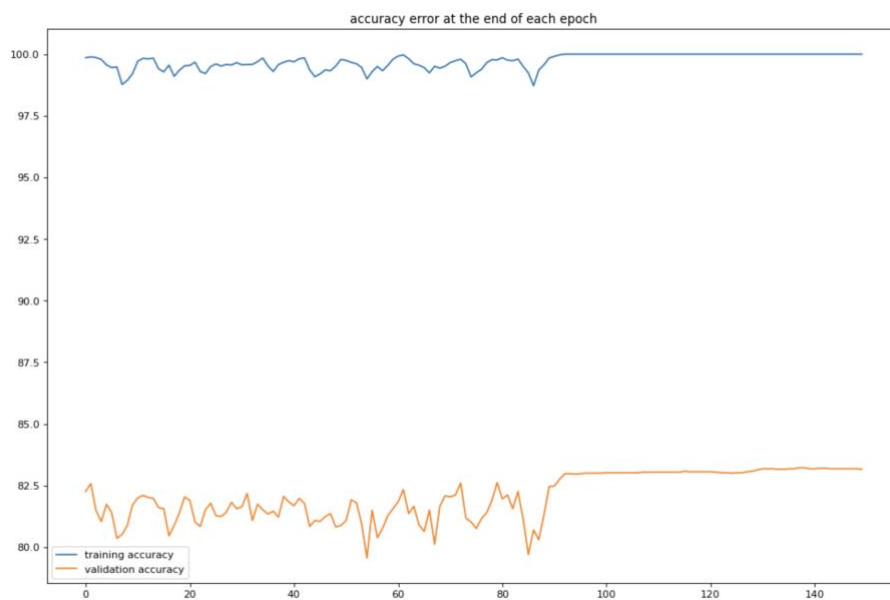
```
1 class MultilayerPerceptron(nn.Module):
2     def __init__(self, in_sz=1568, out_sz=19, layers=[784, 120, 84]):
3         super().__init__()
4         self.fc1 = nn.Linear(in_sz, layers[0])
5         self.fc2 = nn.Linear(layers[0], layers[1])
6         self.fc3 = nn.Linear(layers[1], layers[2])
7         self.fc4 = nn.Linear(layers[2], out_sz)
8
9     def forward(self, X):
10        X = F.relu(self.fc1(X))
11        X = F.relu(self.fc2(X))
12        X = F.relu(self.fc3(X))
13        X = self.fc4(X)
14        return F.log_softmax(X, dim=1)
```


At learning of 0.001, the model fails to converge at the first training, and it converges during the second training around the 85th epochs with highest test accuracy 83.16%

Loss Error (Second training)

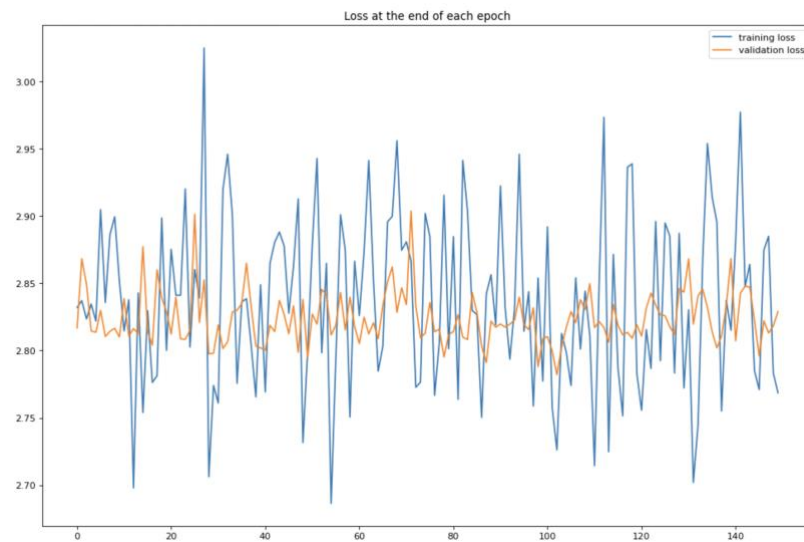


Accuracy (Second training)

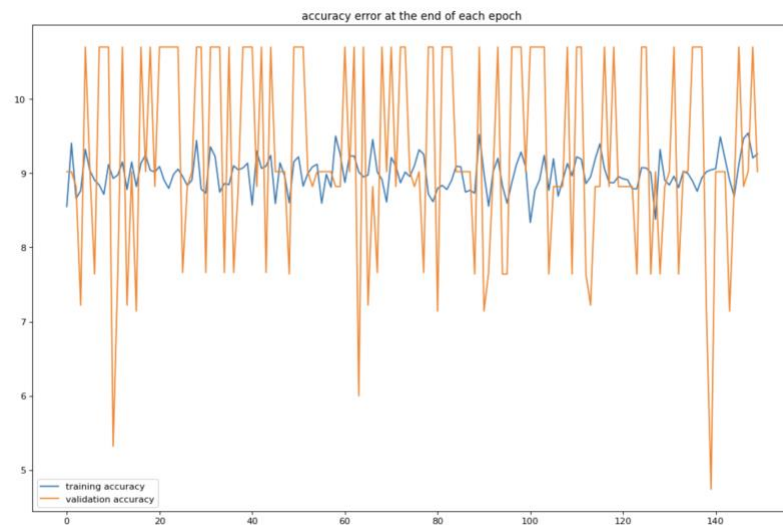


In this part, for both ANN and CNN big learning rate has a strong bad influence towards the training results and to certain extent, big learning rate will result in very bad results and leave the training performance looks unreasonable. For the ANN model we shall observe such phenomenon at learning rate of 0.5.

Loss Error (Learning rate 0.5)



Accuracy (Learning rate 0.5)



Then comes to the CNN model. The idea of creating the customized dataset function is to reshape the flattened data into a rectangular shape, that is restore the data into two 28 by 28 matrix and then concatenate those two matrices into one by the row.

Customized dataset function

```
1 class CustomerDataset(Dataset):
2     def __init__(self, file_path):
3         file_out = pd.read_table(file_path, sep=',', names=range(1569))
4         X = np.array(file_out.iloc[:, :1568].values).reshape(file_out.shape[0], 28, 56)
5         y = file_out.iloc[:, 1568].values
6         # convert to tensor
7         self.X = torch.tensor(X, dtype=torch.float32)
8         self.y = torch.tensor(y)
9
10    def __len__(self):
11        return len(self.y)
12
13    def __getitem__(self, idx):
14        return self.X[idx], self.y[idx]
15
```

For the CNN model I have chosen the same architecture as what I did in the previous part:

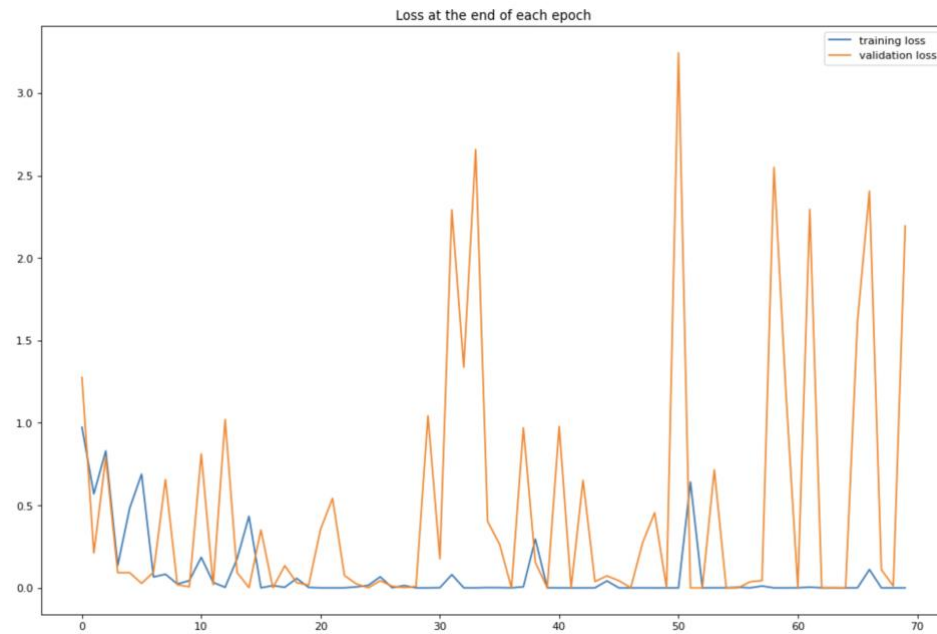
Favorite Deep Learning Architecture. Since the shape of data is a rectangular instead of a square and so the model has been modified slightly different. In here I also choose to set up learning rate of 0.001 as my baseline model with 70 epochs per training.

Model structure

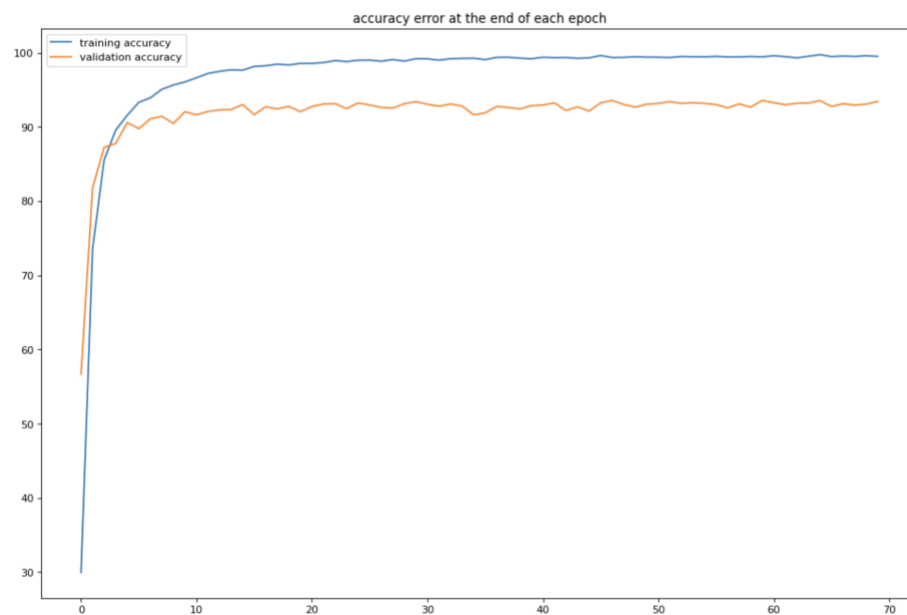
```
1 class ConvolutionalNetwork(nn.Module):
2     def __init__(self):
3         super().__init__()
4         self.conv1 = nn.Conv2d(1, 6, 3, 1)
5         self.conv2 = nn.Conv2d(6, 16, 3, 1)
6         self.fc1 = nn.Linear(5*12*16, 120)
7         self.fc2 = nn.Linear(120, 84)
8         self.fc3 = nn.Linear(84, 19)
9
10    def forward(self, X):
11        X = F.relu(self.conv1(X))
12        X = F.max_pool2d(X, 2, 2)
13        X = F.relu(self.conv2(X))
14        X = F.max_pool2d(X, 2, 2)
15        X = X.view(-1, 5*12*16)
16        X = F.relu(self.fc1(X))
17        X = F.relu(self.fc2(X))
18        X = self.fc3(X)
19        return F.log_softmax(X, dim=1)
```

At learning of 0.001, the CNN model is able to converge at the 46th epochs during the first training with highest test accuracy 93.96%.

Loss Error

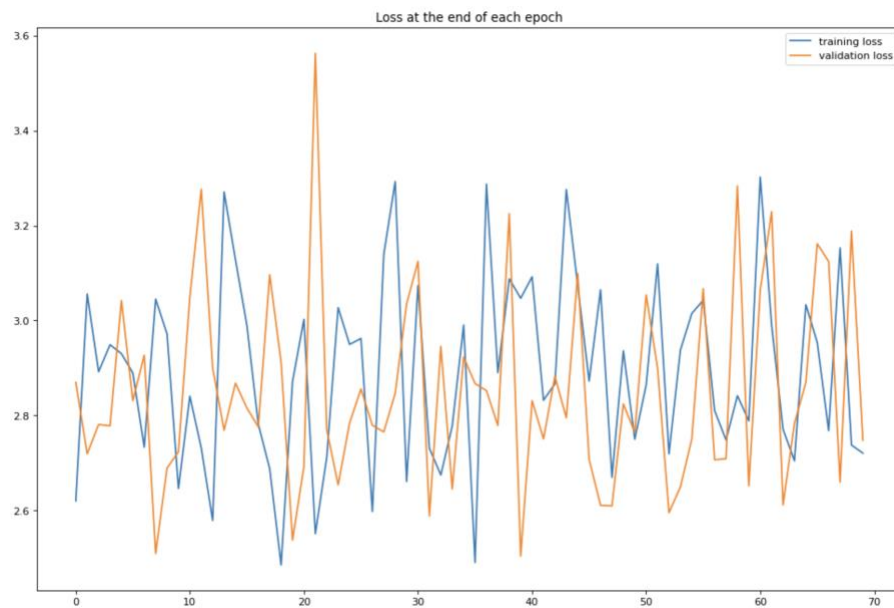


Accuracy

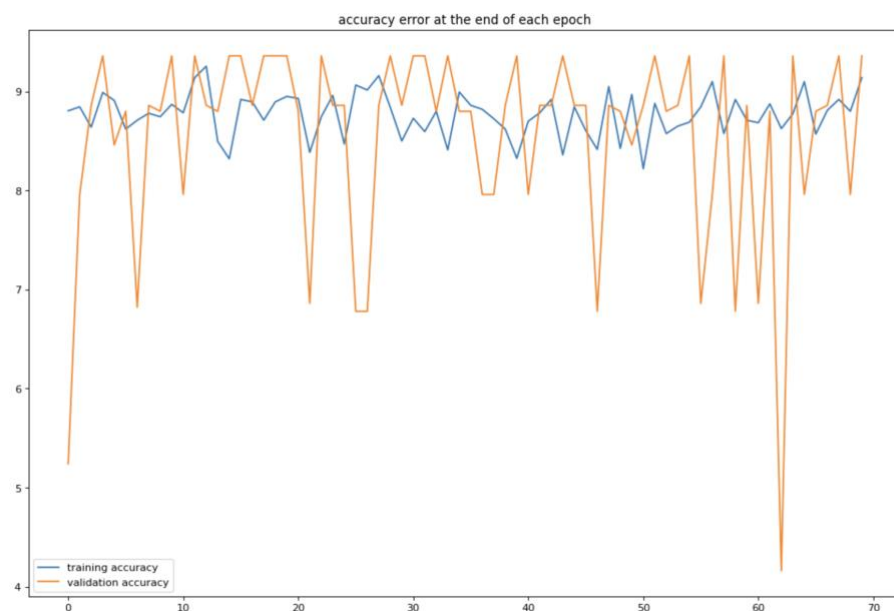


As mentioned before, big learning rate has even a stronger influence towards the CNN model. The result shows that for both learning rate of 0.2 and 0.5, the performance of CNN drops drastically. Below are the plots from learning rate of 0.5.

Loss Error (learning rate 0.5)



Accuracy (learning rate 0.5)



Hence the best performance is performed by CNN. Compared to the MNIST, the customized dataset has a poorer performance as the test accuracy is a lot less than what we got in MNIST dataset.

Reasons why for my model the test misclassification could not go below 1% based on variance perspective and bias perspective.

Handling variance

1. From the plot we can see that the training accuracy is practically 100% and yet the validation accuracy is not which means our model has some overfitting problems and so we might want to add some regulation terms in our convolutional neural network.
2. Theoretically speaking, some of the images in the training set might be randomly rotated in both directions by some degrees or zoomed by some percents or shifted in different directions.
3. For my model, I did not add any dropout layers. If I could add some dropout layers after the pooling layers and some fully connected layers, the performance will probably get better.
4. The model did not use dynamical learning rate. Dynamically reducing the learning rate helps the algorithm to converge faster and reach closer to the global minima, since in this case the CNN model is very sensitive to the learning rate, hence I would consider using dynamical learning rate for future models and trainings.

Handling bias

1. The model did not have enough layers. The number of hidden in the dense layers were also increased to accommodate the larger input due to the increase in the volume

of the convolutional layers since the customized dataset has more information contained than the MNIST dataset.

2. The number of filters is insufficient. By increasing the number of layers and filters, we will obtain a deeper and denser network which allows the model to learn more complex features of the handwritten digits from our customized dataset.