

Contents

1	MNIST Classification Using CNN	1
1.1	Model 1 - 1 Convolutional Layer	1
1.1.1	Detailed Architecture	1
1.1.2	Plot of Randomly Selected Test Images	2
1.2	Model 2 - 2 Convolutional Layer	2
1.2.1	Detailed Architecture	2
1.2.2	Plot of Randomly Selected Test Images	3
1.3	Model 3 - 2 Convolutional Layer + 1 hidden fully connected layer	4
1.3.1	Detailed Architecture	4
1.3.2	Plot of Randomly Selected Test Images	5
1.4	Comparison of the three models	5
1.4.1	Training Loss Plot	6
1.4.2	Validation Loss Plot	6
1.4.3	Test Accuracy	7
2	Generating Adversarial Examples	7
2.1	Learning Curves	7
2.1.1	Plot of Training Loss for all 10 classes	7
2.1.2	Plot of Validation Loss for all 10 classes	8
2.1.3	Plot of Test Accuracy for all 10 classes	8
2.2	Plot of All Noise Patterns	9
2.3	Examples Showing Misclassification Due to Addition of Noise	9
2.3.1	Misclassification with Class 0 Noise	9
2.3.2	Misclassification with Class 1 Noise	10
2.3.3	Misclassification with Class 2 Noise	10
2.3.4	Misclassification with Class 3 Noise	11
2.3.5	Misclassification with Class 4 Noise	11
2.3.6	Misclassification with Class 5 Noise	12
2.3.7	Misclassification with Class 6 Noise	12
2.3.8	Misclassification with Class 7 Noise	13
2.3.9	Misclassification with Class 8 Noise	13
2.3.10	Misclassification with Class 9 Noise	14
3	Visualizing the CNN	14
3.1	Plot of x_{init} for 10 Output Neurons	14
3.2	Plot of x_{init} for 10 Feature Maps after 2^{nd} Max Pooling	14

1 MNIST Classification Using CNN

In this section, we consider three architectures of the CNN which are presented below:

- Model 1 - 1 Convolutional Layer
- Model 2 - 2 Convolutional Layer
- Model 3 - 2 Convolutional Layer + 1 hidden fully connected layer

The fully description of the architecture is given in each corresponding section.

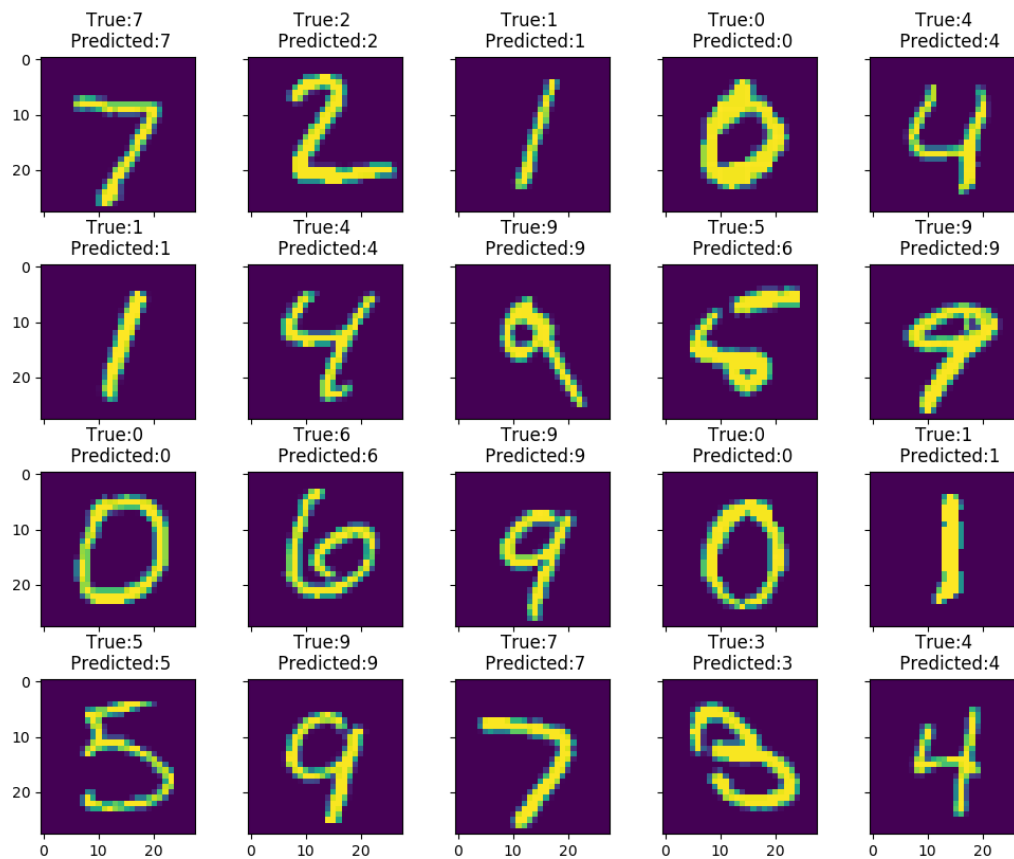
1.1 Model 1 - 1 Convolutional Layer

1.1.1 Detailed Architecture

The 1 Convolutional Layer CNN has the following architecture:

- Input
- Conv Layer
 - Number of Filters = 32
 - Filter Size - 3×3
 - Stride = 1
 - Zero Padding of 1
- Max Pool Layer
 - 2×2 Max Pooling
 - Stride = 2
- Fully Connected (10 Outputs)
- Softmax Classifier

1.1.2 Plot of Randomly Selected Test Images



1.2 Model 2 - 2 Convolutional Layer

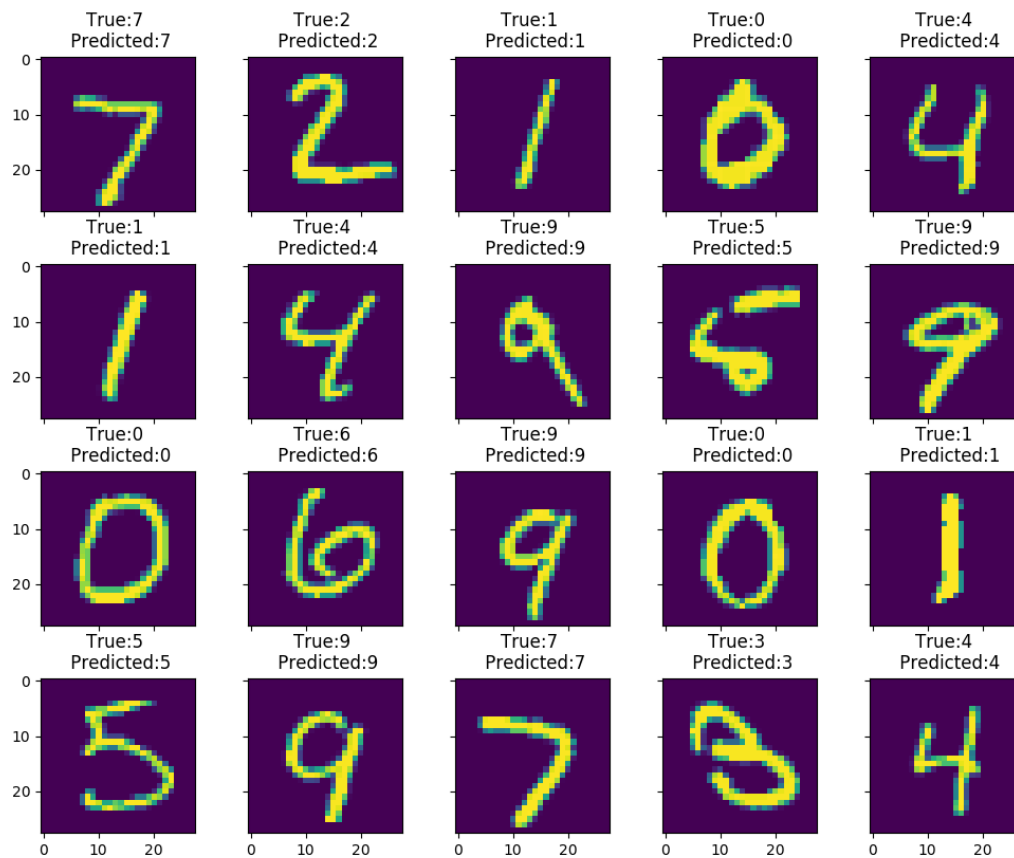
1.2.1 Detailed Architecture

The 2 Convolutional Layer CNN has the following architecture:

- Input
- Conv Layer 1
 - Number of Filters = 32

- Filter Size - 3×3
- Stride = 1
- Zero Padding of 1
- Max Pool Layer 1
 - 2×2 Max Pooling
 - Stride = 2
- Conv Layer 2
 - Number of Filters = 32
 - Filter Size - 3×3
 - Stride = 1
 - Zero Padding of 1
- Max Pool Layer 2
 - 2×2 Max Pooling
 - Stride = 2
- Fully Connected (10 Outputs)
- Softmax Classifier

1.2.2 Plot of Randomly Selected Test Images



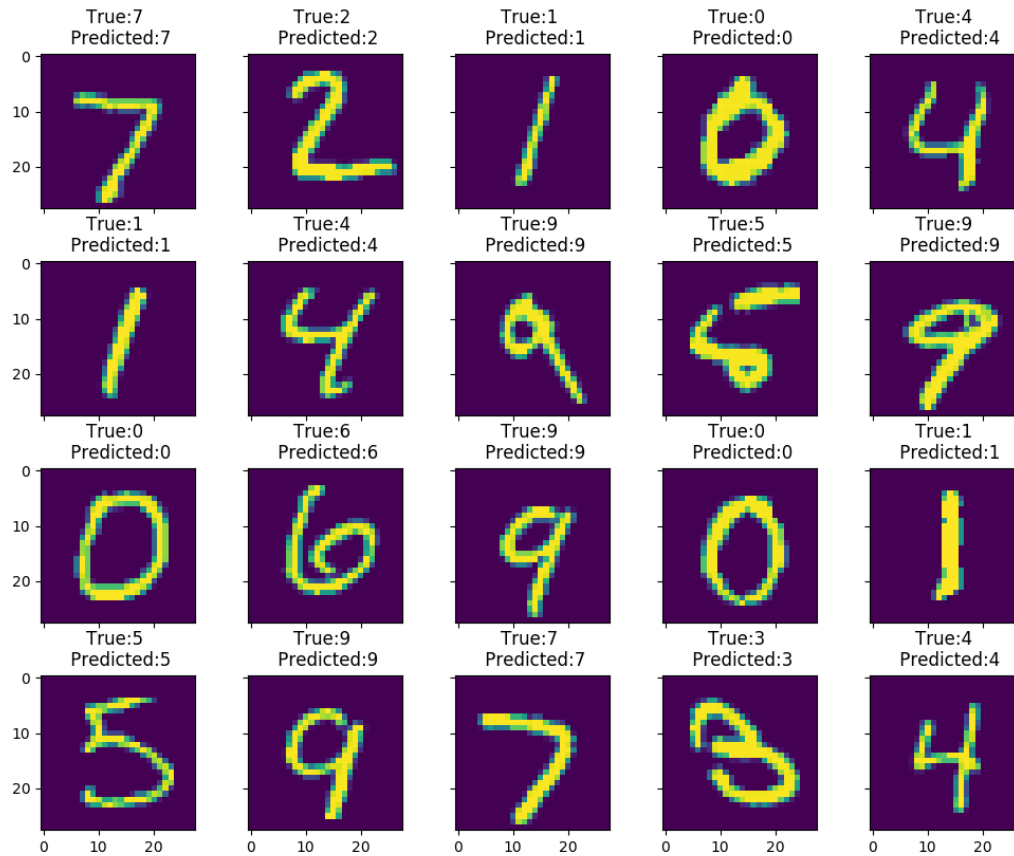
1.3 Model 3 - 2 Convolutional Layer + 1 hidden fully connected layer

1.3.1 Detailed Architecture

The 2 Convolutional Layer + 1 hidden fully connected layer CNN has the following architecture:

- Input
- Conv Layer 1
 - Number of Filters = 32
 - Filter Size - 3×3
 - Stride = 1
 - Zero Padding of 1
- Max Pool Layer 1
 - 2×2 Max Pooling
 - Stride = 2
- Conv Layer 2
 - Number of Filters = 32
 - Filter Size - 3×3
 - Stride = 1
 - Zero Padding of 1
- Max Pool Layer 2
 - 2×2 Max Pooling
 - Stride = 2
- Fully Connected (500 Outputs)
- Fully Connected (10 Outputs)
- Softmax Classifier

1.3.2 Plot of Randomly Selected Test Images

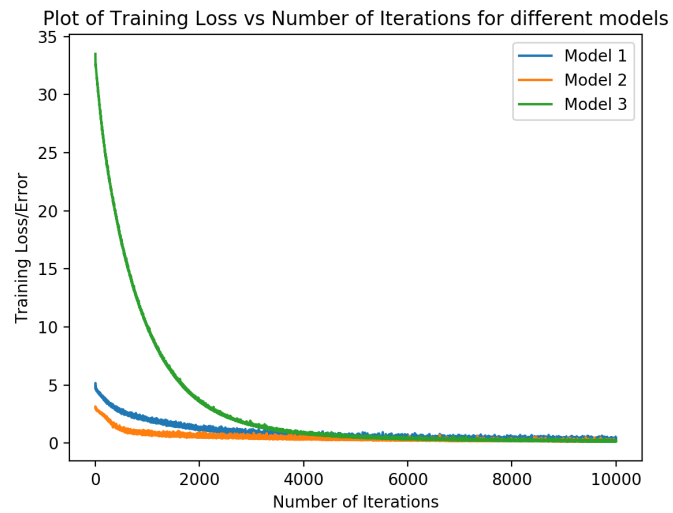


1.4 Comparison of the three models

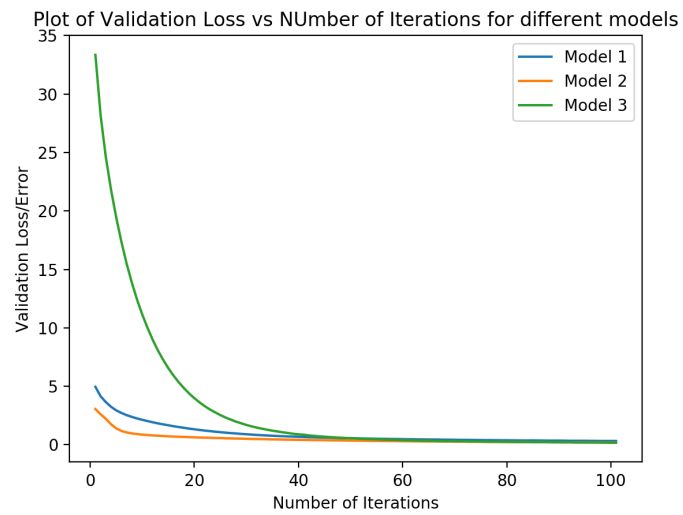
The three models are compared on the following three metrics :

- Training Loss Plot
- Validation Loss Plot
- Test Accuracy

1.4.1 Training Loss Plot

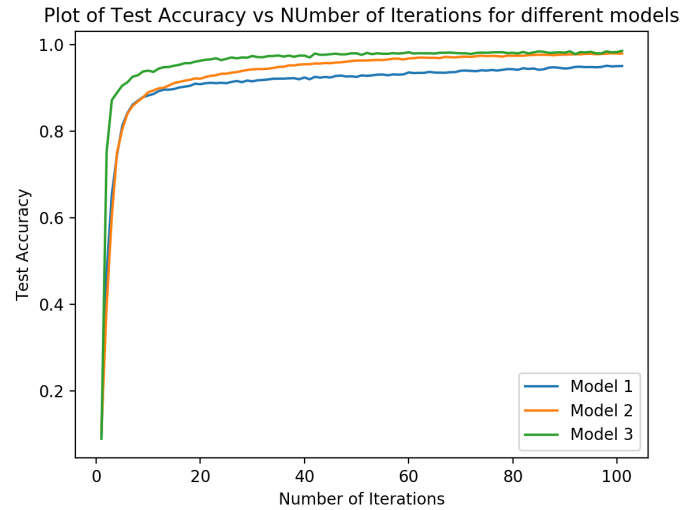


1.4.2 Validation Loss Plot



Note that since the validation loss plot is coming down with the increasing number of iterations, we can conclude that none of the models are over-fitting and regularization seems to be working. Also, note that even though it may not be very apparent from the graph, but Model 3 gives the least validation loss/error after 10000 iterations of training.

1.4.3 Test Accuracy



Model	Test Accuracy
Model 1 - 1 Convolutional Layer	95.06%
Model 2 - 2 Convolutional Layer	97.96%
Model 3 - 2 Convolutional Layer + 1 Hidden Layer	98.5%

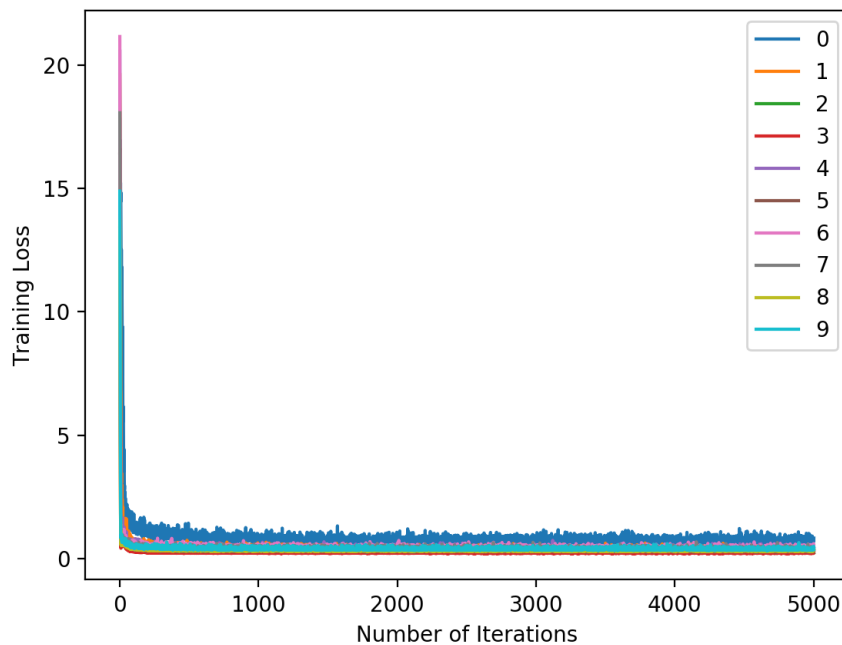
We can clearly see that Model 3 outperforms the other models in terms of **Test Accuracy** and hence is the best model learnt. We make use of this learnt model in our future experimentation.

2 Generating Adversarial Examples

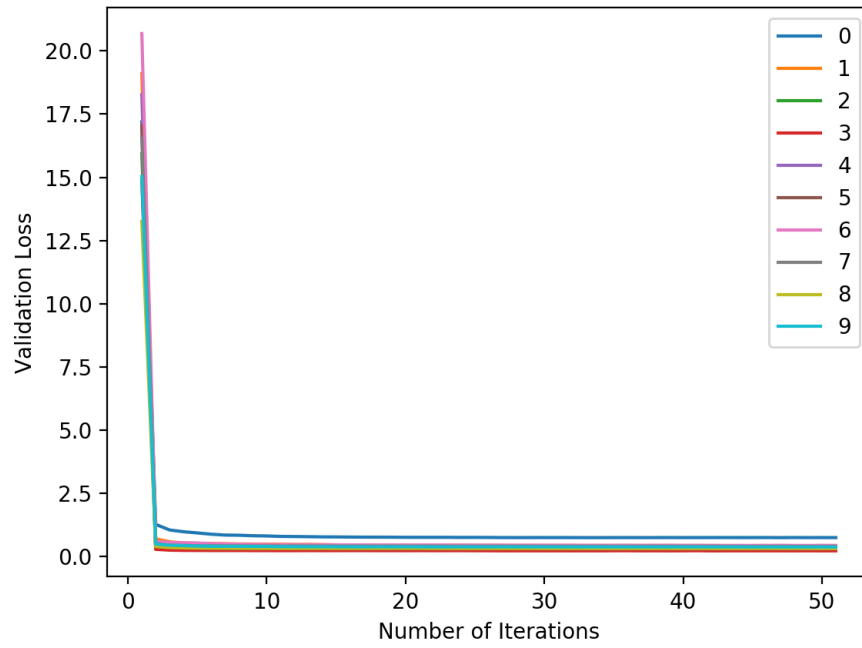
2.1 Learning Curves

Presented below are the learning curves for generating the noise masks for each of the classes from 0 to 9.

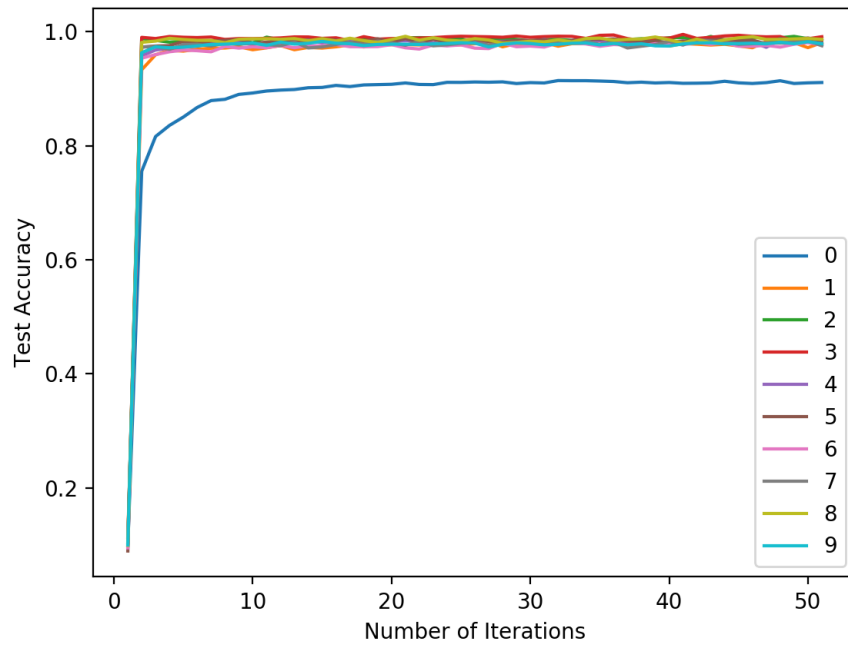
2.1.1 Plot of Training Loss for all 10 classes



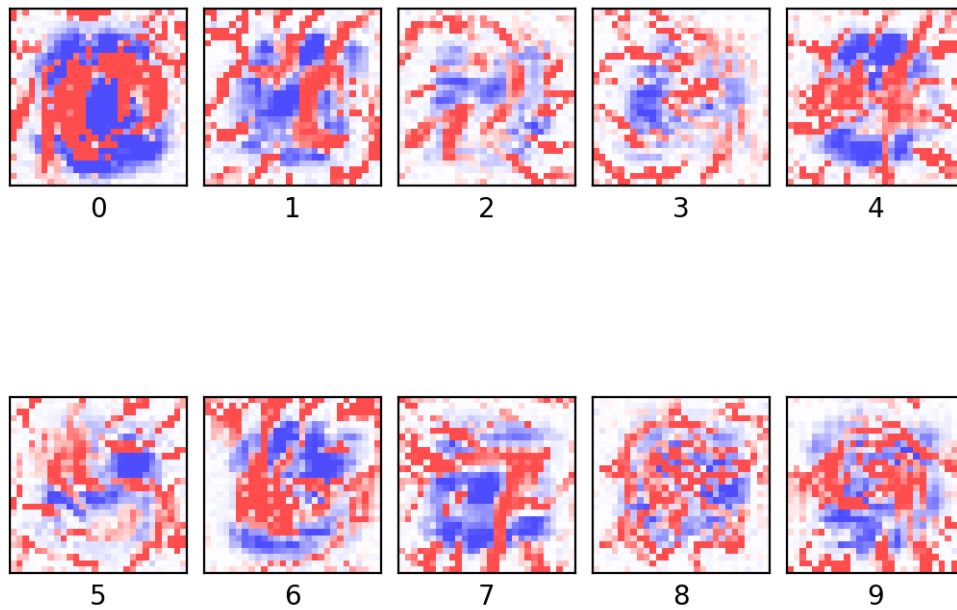
2.1.2 Plot of Validation Loss for all 10 classes



2.1.3 Plot of Test Accuracy for all 10 classes

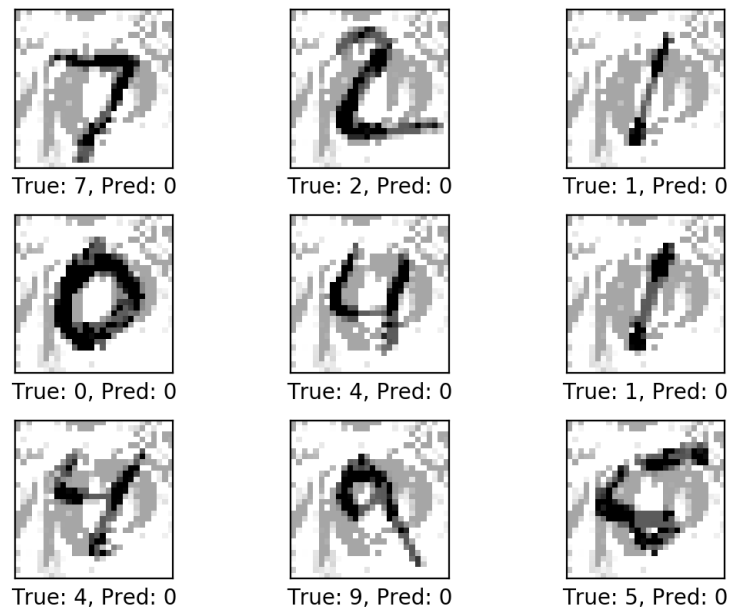


2.2 Plot of All Noise Patterns

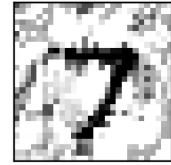


2.3 Examples Showing Misclassification Due to Addition of Noise

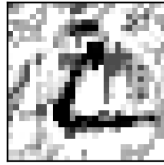
2.3.1 Misclassification with Class 0 Noise



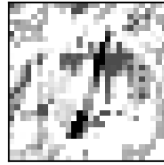
2.3.2 Misclassification with Class 1 Noise



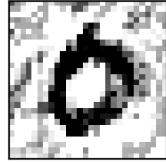
True: 7, Pred: 1



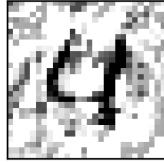
True: 2, Pred: 1



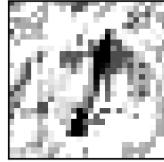
True: 1, Pred: 1



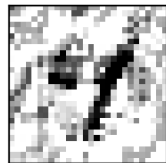
True: 0, Pred: 1



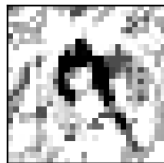
True: 4, Pred: 1



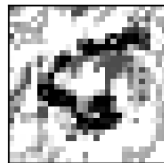
True: 1, Pred: 1



True: 4, Pred: 1

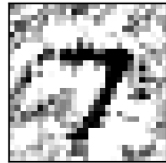


True: 9, Pred: 1

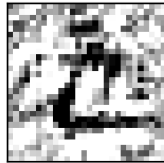


True: 5, Pred: 1

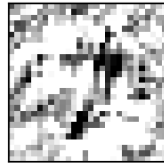
2.3.3 Misclassification with Class 2 Noise



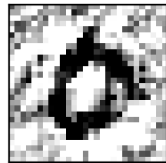
True: 7, Pred: 2



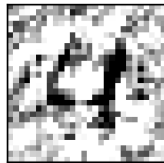
True: 2, Pred: 2



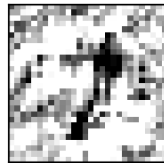
True: 1, Pred: 2



True: 0, Pred: 2



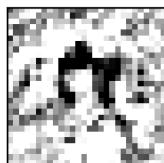
True: 4, Pred: 2



True: 1, Pred: 2



True: 4, Pred: 2

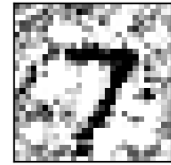


True: 9, Pred: 2

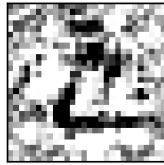


True: 5, Pred: 2

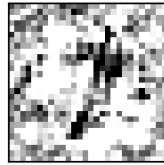
2.3.4 Misclassification with Class 3 Noise



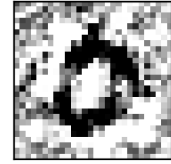
True: 7, Pred: 3



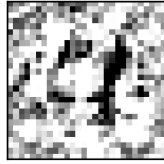
True: 2, Pred: 3



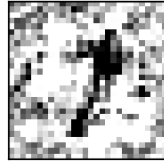
True: 1, Pred: 3



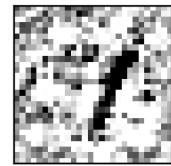
True: 0, Pred: 3



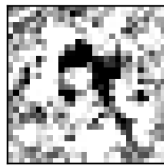
True: 4, Pred: 3



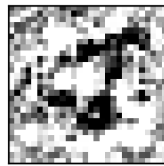
True: 1, Pred: 3



True: 4, Pred: 3

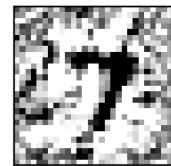


True: 9, Pred: 3

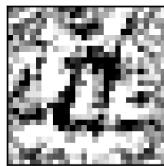


True: 5, Pred: 3

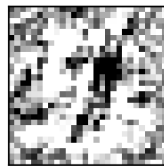
2.3.5 Misclassification with Class 4 Noise



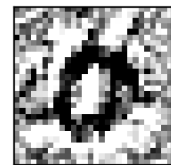
True: 7, Pred: 4



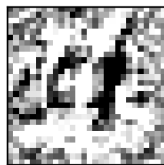
True: 2, Pred: 4



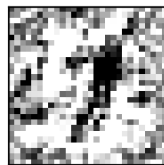
True: 1, Pred: 4



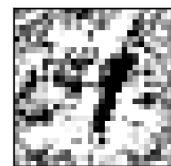
True: 0, Pred: 4



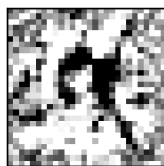
True: 4, Pred: 4



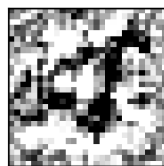
True: 1, Pred: 4



True: 4, Pred: 4



True: 9, Pred: 4

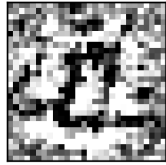


True: 5, Pred: 4

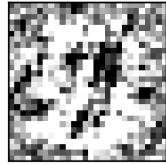
2.3.6 Misclassification with Class 5 Noise



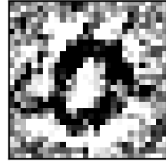
True: 7, Pred: 5



True: 2, Pred: 5



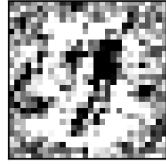
True: 1, Pred: 5



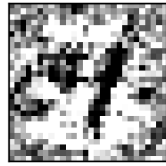
True: 0, Pred: 5



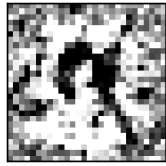
True: 4, Pred: 5



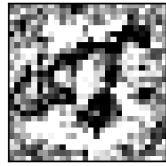
True: 1, Pred: 5



True: 4, Pred: 5

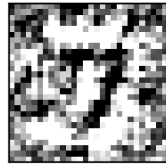


True: 9, Pred: 5

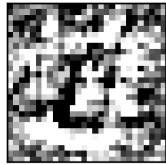


True: 5, Pred: 5

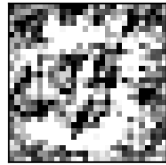
2.3.7 Misclassification with Class 6 Noise



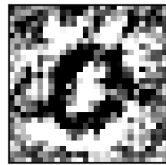
True: 7, Pred: 6



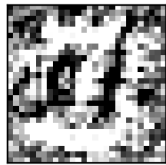
True: 2, Pred: 6



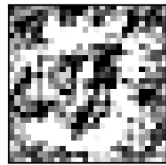
True: 1, Pred: 6



True: 0, Pred: 6



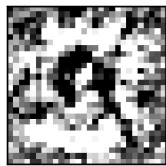
True: 4, Pred: 6



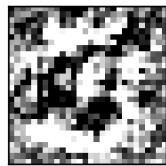
True: 1, Pred: 6



True: 4, Pred: 6

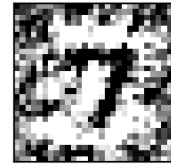


True: 9, Pred: 6



True: 5, Pred: 6

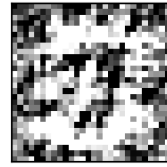
2.3.8 Misclassification with Class 7 Noise



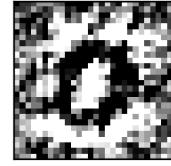
True: 7, Pred: 7



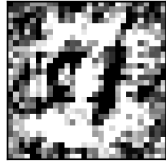
True: 2, Pred: 7



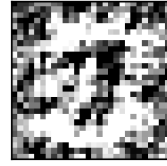
True: 1, Pred: 7



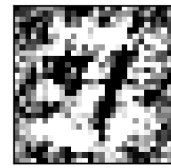
True: 0, Pred: 7



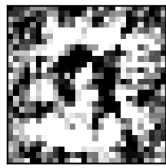
True: 4, Pred: 7



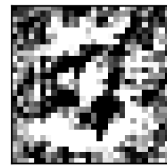
True: 1, Pred: 7



True: 4, Pred: 7

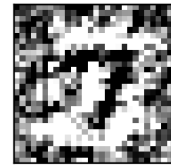


True: 9, Pred: 7

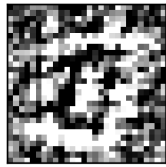


True: 5, Pred: 7

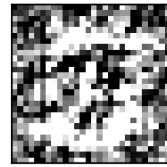
2.3.9 Misclassification with Class 8 Noise



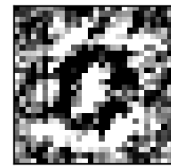
True: 7, Pred: 8



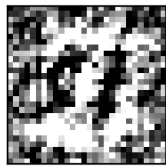
True: 2, Pred: 8



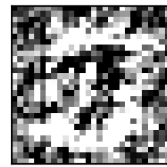
True: 1, Pred: 8



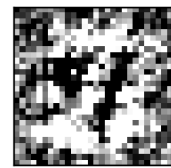
True: 0, Pred: 8



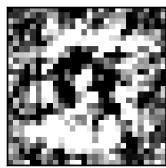
True: 4, Pred: 8



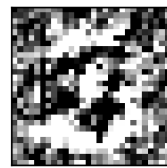
True: 1, Pred: 8



True: 4, Pred: 8

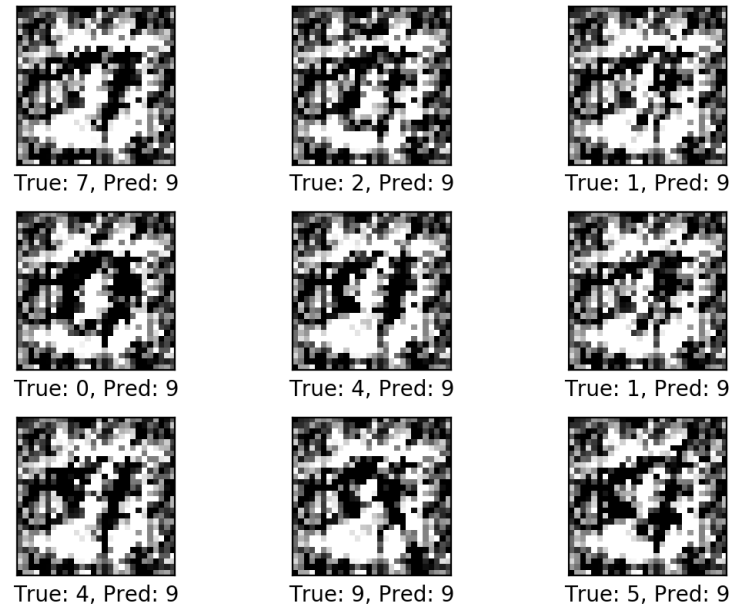


True: 9, Pred: 8



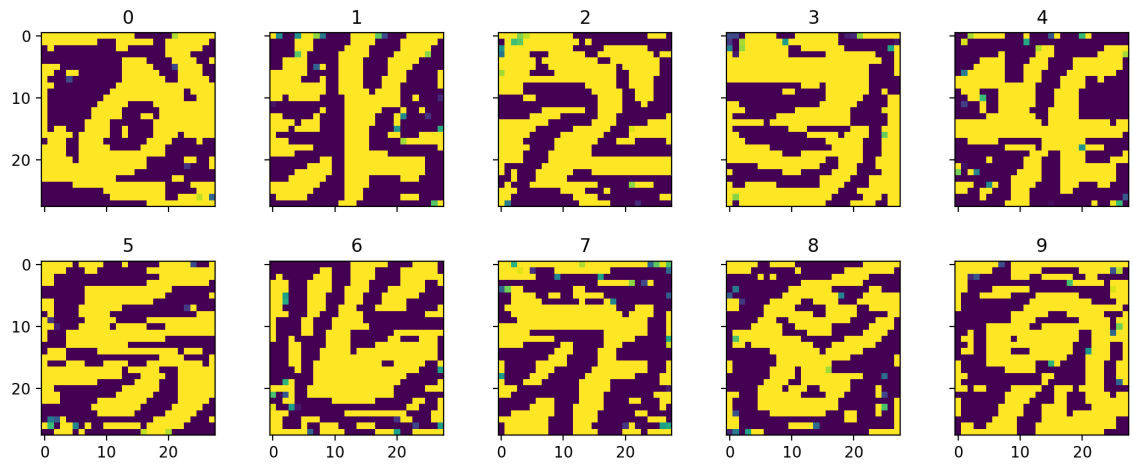
True: 5, Pred: 8

2.3.10 Misclassification with Class 9 Noise



3 Visualizing the CNN

3.1 Plot of x_{init} for 10 Output Neurons



3.2 Plot of x_{init} for 10 Feature Maps after 2^{nd} Max Pooling

