

EE6132 Assignment No. 1: Report

-Akash Reddy A, EE17B001

September 5, 2019

Overview

The aim of this assignment has been to implement a Multilayer Feedforward Neural Network (MLP) with Backpropagation (BP) learning. Additionally, the effects of different activations, noise, regularisation, and comparison with other Machine Learning algorithms have been experimented with.

1 MLP Architecture and Training (baseline model)

Note: Average Error Rate has been computed as the mean of the error rates for each digit, where the error rate for each digit i has been evaluated as:

$$error_i = 1 - \frac{TP_i + TN_i}{TP_i + FP_i + FN_i + TN_i}$$

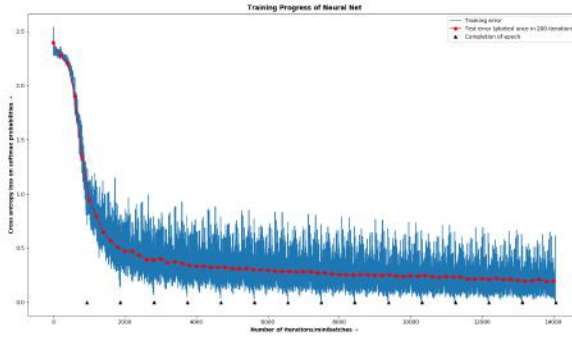
where for example for the digit 1, TP_1 = True Positives (all 1s predicted as 1), FP_1 = False Positives (all non-1s predicted as 1), FN_1 = False Negatives (all 1s predicted as non-1), TN_1 = True Negatives (all non-1s predicted as non-1).

Standard Deviation of the Error Rate is the standard deviation of the 10 error rates evaluated for each digit.

Training Loss, Test Loss Plots and Metrics for Sigmoidal Activation:

In each of the plots, each iteration is considered to be one run of gradient descent over a single minibatch of 64 samples. 15 epochs of a 60000 strong set would therefore lead to roughly 14000 iterations.

First, we explore different learning rates for the sigmoidal activated baseline MLP. The metrics for each algorithm are reported in the cropped terminal window screenshot to the right of the plot.

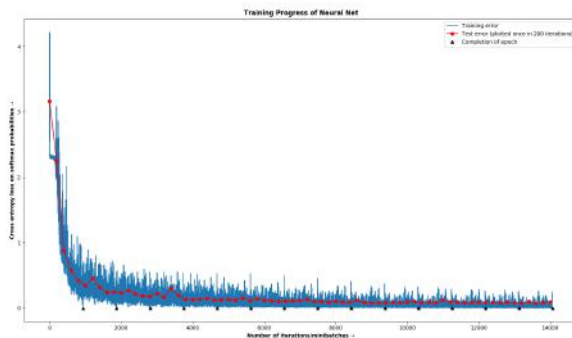


```
Total # of epochs: 15
Epoch #1 done...
Accuracy on test set after epoch #1 = 86.41%
Epoch #2 done...
Accuracy on test set after epoch #2 = 87.48%
Epoch #3 done...
Accuracy on test set after epoch #3 = 86.49%
Epoch #4 done...
Accuracy on test set after epoch #4 = 88.27000000000001%
Epoch #5 done...
Accuracy on test set after epoch #5 = 89.2%
Epoch #6 done...
Accuracy on test set after epoch #6 = 90.07%
Epoch #7 done...
Accuracy on test set after epoch #7 = 90.83%
Epoch #8 done...
Accuracy on test set after epoch #8 = 91.45%
Epoch #9 done...
Accuracy on test set after epoch #9 = 91.86%
Epoch #10 done...
Accuracy on test set after epoch #10 = 92.42%
Epoch #11 done...
Accuracy on test set after epoch #11 = 92.82000000000001%
Epoch #12 done...
Accuracy on test set after epoch #12 = 93.11%
Epoch #13 done...
Accuracy on test set after epoch #13 = 93.45%
Epoch #14 done...
Accuracy on test set after epoch #14 = 93.76%
Epoch #15 done...
Accuracy on test set after epoch #15 = 94.01%
Confusion Matrix:
[[ 960.  0.  1.  1.  0.  8.  6.  1.  2.  1.]
 [ 0. 1111.  2.  2.  1.  1.  3.  1. 13.  1.]
 [14.  0. 943.  5.  9.  1. 14.  8. 29.  1.]
 [ 2.  1. 13. 929.  0. 26.  1.  9. 21.  8.]
 [ 1.  2.  4.  0. 925.  0. 10.  1.  4. 35.]
 [ 7.  2.  0. 13.  3. 874. 13.  8. 22.  8.]
 [ 9.  3.  1.  0.  8. 16. 915.  6.  6.  0.]
 [ 2. 15. 15.  4.  5.  1.  0. 933.  1. 52.]
 [ 5.  6.  1. 12.  9. 10.  5.  3. 912.  5.]
 [10.  0.  1.  5. 23.  5.  0.  4.  0. 949.]]
...
Average error rate = 1.197999999999999%
Standard deviation of error rate = 0.3349865689485234%

Final accuracy on test set = 94.01%
Final precision on test set = 94.08855846997193%
Final recall on test set = 93.97312665766175%
Final F1-score on test set = 0.9395735579156144
```

Figure 1: Learning Rate $\alpha=0.1$

For $\alpha = 0.1$, we have a decent convergence of the losses, leading to a final test accuracy of $\approx 94\%$.



```
Total # of epochs: 15
Epoch #1 done...
Accuracy on test set after epoch #1 = 84.25%
Epoch #2 done...
Accuracy on test set after epoch #2 = 89.75%
Epoch #3 done...
Accuracy on test set after epoch #3 = 92.92%
Epoch #4 done...
Accuracy on test set after epoch #4 = 94.37%
Epoch #5 done...
Accuracy on test set after epoch #5 = 95.37%
Epoch #6 done...
Accuracy on test set after epoch #6 = 96.08%
Epoch #7 done...
Accuracy on test set after epoch #7 = 96.73%
Epoch #8 done...
Accuracy on test set after epoch #8 = 97.1%
Epoch #9 done...
Accuracy on test set after epoch #9 = 97.22%
Epoch #10 done...
Accuracy on test set after epoch #10 = 97.31%
Epoch #11 done...
Accuracy on test set after epoch #11 = 97.38%
Epoch #12 done...
Accuracy on test set after epoch #12 = 97.39999999999999%
Epoch #13 done...
Accuracy on test set after epoch #13 = 97.54%
Epoch #14 done...
Accuracy on test set after epoch #14 = 97.53%
Epoch #15 done...
Accuracy on test set after epoch #15 = 97.58%
Confusion Matrix:
[[ 973.  0.  1.  0.  1.  1.  0.  1.  3.  0.]
 [ 0. 1130.  1.  1.  0.  1.  0.  0.  1.  1.]
 [ 7.  3. 1066.  3.  4.  0.  3.  4.  2.  0.]
 [ 0.  0.  0. 989.  0.  2.  0.  5.  4.  4.]
 [ 1.  0.  1.  0. 993.  0.  4.  2.  0. 11.]
 [ 4.  1.  0. 10.  2. 871.  1.  0.  0.  3.]
 [ 3.  3.  2.  0. 14. 20. 911.  0.  3.  0.]
 [ 1.  7.  8.  1.  2.  0.  0. 997.  1. 11.]
 [ 5.  1.  3.  7.  0.  5.  2.  3. 938.  4.]
 [ 4.  2.  0.  5. 16.  1.  0.  1.  0. 980.]]
...
Average error rate = 0.484%
Standard deviation of error rate = 0.12026637102698357%

Final accuracy on test set = 97.58%
Final precision on test set = 97.57375804263823%
Final recall on test set = 97.54656186253188%
Final F1-score on test set = 0.975512428851021
```

Figure 2: Learning Rate $\alpha=1$

For $\alpha = 1$, we have almost the best and fastest convergence of the losses, leading to a high final test accuracy of $\approx 97\%$.

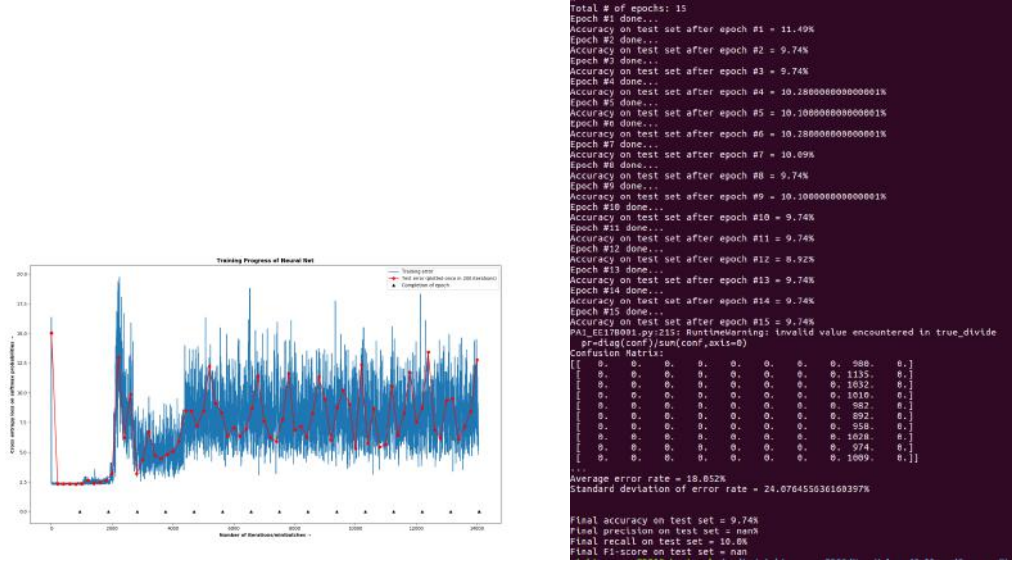


Figure 3: Learning Rate $\alpha=7$

For the high $\alpha = 7$, the convergence of the losses is much worse. This is because the learning rate is too large and it cannot take small enough jumps to reach the minimum during gradient descent. Hence, it tends to overshoot to remain in the higher values of the cost function.

2 Activation Function

ReLU

When it comes to the ReLU, we find that the previous near-optimal learning rate $\alpha = 1$ is too large. Upon playing with the learning rate, we get a good convergence for $\alpha = 0.1$ with a final test accuracy of $\approx 98\%$.

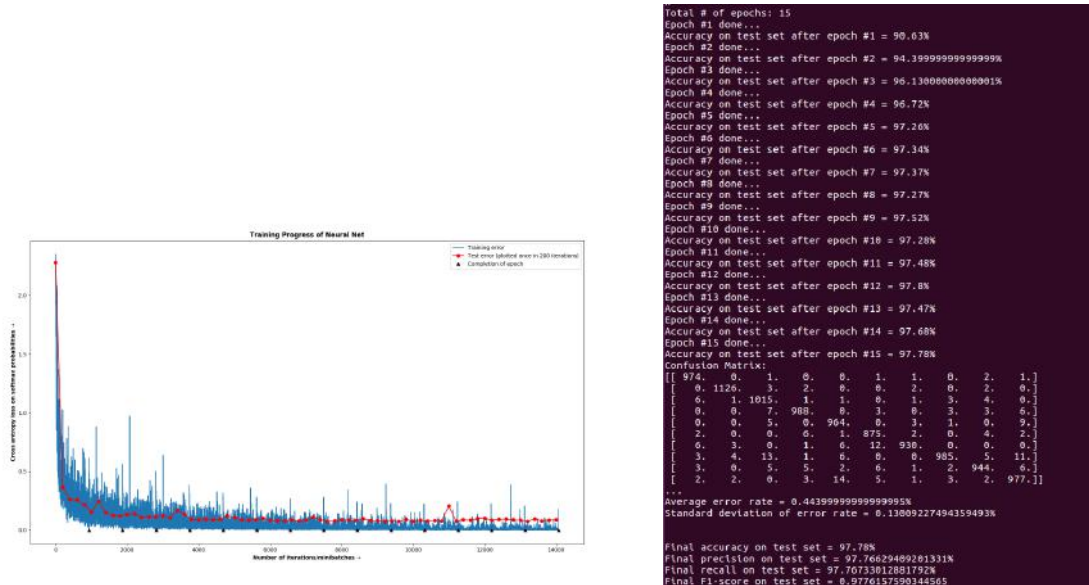


Figure 4: Learning Rate $\alpha=0.1$

tanh

When it comes to tanh, we find that the previous near-optimal learning rate $\alpha = 1$ is slightly large, as the final test accuracy gets stuck at roughly 89% when $\alpha = 1$ is used. Upon playing with the learning rate, we get a good convergence for $\alpha = 0.1$ with a final test accuracy of $\approx 97\%$.

We also observe a slower convergence in tanh than in ReLU. This can be explained by the fact that the ReLU neurons which enter into the negative values for z are made equal to the fixed 0. They change no more after this. Therefore, the neurons in the ReLU algorithm become inactive quicker, leading to faster convergence.

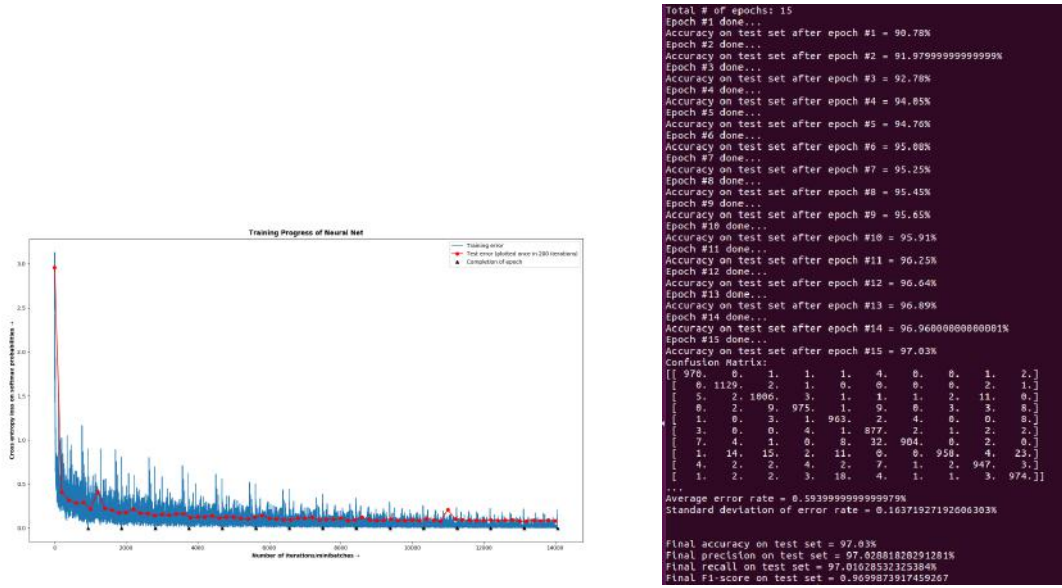


Figure 5: Learning Rate $\alpha=0.1$

Inactive Neurons

When the gradients of neurons don't change much with iterations, this means that they are influencing the change in the neural network parameters lesser and lesser. Therefore, the percentage of inactive neurons is an indicator of convergence of the algorithm to a cost minimum.

It can be clearly seen that in all three cases, the overall percentage of inactive neurons increases as a general trend with iterations. This is because, as the algorithm converges, more and more neurons become inactive.

The percentage of inactive neurons for the sigmoidal activation for its best learning rate, and also the ReLU for its best learning rate are more or less equal. This is because the rate of convergence (if one observes the training/test loss plots) is similar for both these activations. On the other hand, while the percentage of inactive neurons is also increasing for tanh, it doesn't increase as much due to slower convergence.

The percentage of inactive neurons has been plotted once every 200 iterations for better visibility.

Sigmoidal Activation

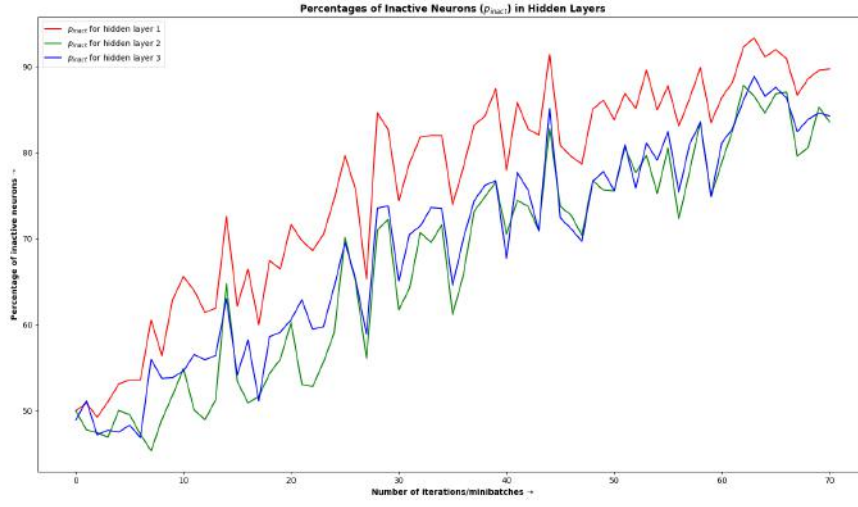


Figure 6: Learning Rate $\alpha=1$

ReLU Activation

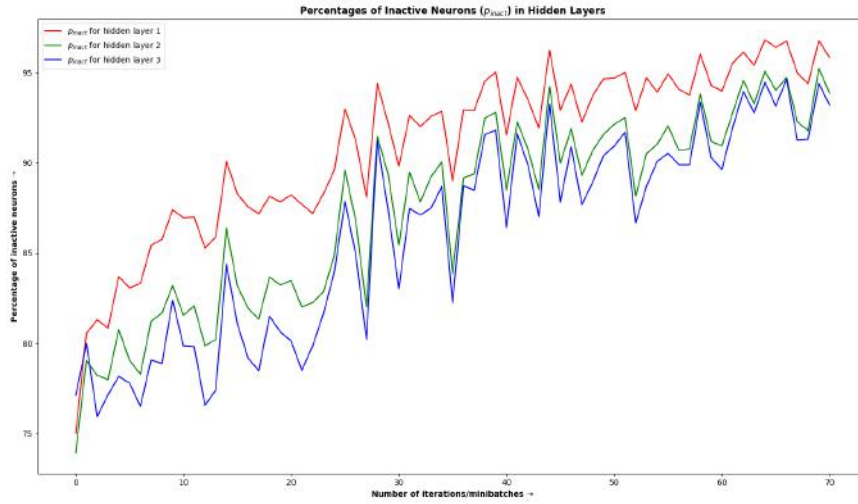


Figure 7: Learning Rate $\alpha=0.1$

tanh Activation

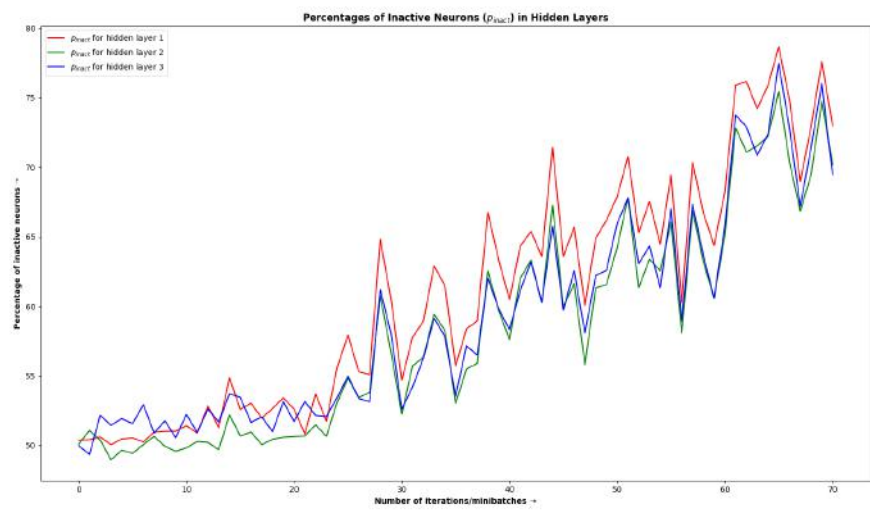


Figure 8: Learning Rate $\alpha=0.1$

3 Regularisation

1

The addition of Gaussian noise to the hidden layers in the training of the neural network is a mild regularisation method. This is because, by adding noise, we make each of our hidden layer neuron values a stochastic model, with a probability of a Gaussian distribution centred around the original value. This makes the hidden layer values more "general", accounting for deviations from the training data. This is what creates the regularising effect.

However, in our particular case, the regularisation isn't much visible. This is because the training set is already quite large and varied, and hence not much overfitting occurs during training.

It is seen that adding noise in the forward pass is better- in the sense that it converges faster and to a slightly better test set accuracy- than adding noise in the backward pass. This is probably because, adding noise in the forward pass is like adding noise to the input data itself, because the noise features in the calculation of the a_i 's during forward pass.

```
Total # of epochs: 15
Epoch #1 done...
Accuracy on test set after epoch #1 = 65.46%
Epoch #2 done...
Accuracy on test set after epoch #2 = 83.8%
Epoch #3 done...
Accuracy on test set after epoch #3 = 86.8%
Epoch #4 done...
Accuracy on test set after epoch #4 = 88.53%
Epoch #5 done...
Accuracy on test set after epoch #5 = 89.41%
Epoch #6 done...
Accuracy on test set after epoch #6 = 90.27%
Epoch #7 done...
Accuracy on test set after epoch #7 = 91.25%
Epoch #8 done...
Accuracy on test set after epoch #8 = 91.53999999999999%
Epoch #9 done...
Accuracy on test set after epoch #9 = 92.11%
Epoch #10 done...
Accuracy on test set after epoch #10 = 92.46%
Epoch #11 done...
Accuracy on test set after epoch #11 = 92.86%
Epoch #12 done...
Accuracy on test set after epoch #12 = 93.11%
Epoch #13 done...
Accuracy on test set after epoch #13 = 93.53%
Epoch #14 done...
Accuracy on test set after epoch #14 = 93.76%
Epoch #15 done...
Accuracy on test set after epoch #15 = 94.06%
Confusion Matrix:
[[1922, 0, 2, 2, 0, 10, 16, 2, 4, 2],
 [ 0, 2218, 4, 4, 2, 2, 8, 2, 28, 2],
 [ 24, 22, 1876, 14, 16, 2, 28, 16, 64, 2],
 [ 2, 4, 22, 1888, 0, 28, 2, 16, 42, 16],
 [ 4, 2, 6, 0, 1852, 0, 20, 2, 10, 68],
 [ 16, 4, 0, 44, 8, 1622, 18, 0, 54, 18],
 [ 20, 6, 4, 0, 16, 26, 1830, 0, 14, 9],
 [ 4, 30, 30, 10, 10, 2, 0, 1844, 4, 122],
 [ 10, 10, 0, 26, 16, 22, 12, 4, 1840, 8],
 [ 20, 10, 2, 8, 30, 10, 0, 4, 14, 1920]]
...
Average error rate = 1.1880000000000002%
Standard deviation of error rate = 0.33388921804234014%

Final accuracy on test set = 94.06%
Final precision on test set = 94.11656168690851%
Final recall on test set = 94.01430868516578%
Final F1-score on test set = 0.9401514018846825

Total # of epochs: 15
Epoch #1 done...
Accuracy on test set after epoch #1 = 60.27%
Epoch #2 done...
Accuracy on test set after epoch #2 = 82.50999999999999%
Epoch #3 done...
Accuracy on test set after epoch #3 = 86.04%
Epoch #4 done...
Accuracy on test set after epoch #4 = 88.02%
Epoch #5 done...
Accuracy on test set after epoch #5 = 88.99000000000001%
Epoch #6 done...
Accuracy on test set after epoch #6 = 89.8%
Epoch #7 done...
Accuracy on test set after epoch #7 = 90.53999999999999%
Epoch #8 done...
Accuracy on test set after epoch #8 = 91.3%
Epoch #9 done...
Accuracy on test set after epoch #9 = 91.82000000000001%
Epoch #10 done...
Accuracy on test set after epoch #10 = 92.12%
Epoch #11 done...
Accuracy on test set after epoch #11 = 92.44%
Epoch #12 done...
Accuracy on test set after epoch #12 = 92.66%
Epoch #13 done...
Accuracy on test set after epoch #13 = 93.2%
Epoch #14 done...
Accuracy on test set after epoch #14 = 93.30000000000001%
Epoch #15 done...
Accuracy on test set after epoch #15 = 93.67%
Confusion Matrix:
[[1922, 0, 2, 4, 0, 14, 10, 2, 4, 2],
 [ 0, 2220, 4, 4, 2, 4, 8, 2, 24, 2],
 [ 28, 22, 1860, 16, 14, 4, 26, 20, 70, 4],
 [ 2, 4, 26, 1856, 0, 42, 2, 16, 58, 14],
 [ 4, 4, 6, 0, 1842, 0, 20, 2, 12, 74],
 [ 16, 4, 2, 32, 10, 1622, 24, 0, 58, 16],
 [ 10, 0, 2, 0, 10, 38, 1820, 0, 12, 0],
 [ 0, 28, 32, 8, 10, 2, 0, 1856, 4, 110],
 [ 10, 12, 0, 20, 18, 36, 10, 4, 1828, 10],
 [ 22, 14, 2, 6, 38, 8, 0, 8, 18, 1902]]
...
Average error rate = 1.2659999999999993%
Standard deviation of error rate = 0.38024202818731173%

Final accuracy on test set = 93.67%
Final precision on test set = 93.76109247168298%
Final recall on test set = 93.6230176038257%
Final F1-score on test set = 0.936146289705809
```

Figure 9: Metrics after adding Forward and Backward Noise respectively- $\alpha=0.1$ and a Sigmoidal Activation (the confidence matrices have been doubled by mistake)

Data Augmentation: Data augmentation has been done by adding noise to the training set and augmenting the original training set with the noisy one. This increases 60000 training examples to 120000 of them, and hence takes nearly twice the time to learn for 15 epochs. This is one of the main disadvantages of data augmentation.

We use sigmoidal activation and learning rate $\alpha = 0.1$ to observe the effects of regularisation. Hence, the beauty of this can be appreciated by comparing with Figure 1- the metrics of the *unregularised*, baseline neural network with sigmoidal activation and learning rate $\alpha = 0.1$.

We can very clearly observe a jump in the accuracy during the first few epochs, as compared to the unregularised case. The final test accuracy is also significantly better. Therefore, data augmentation is very effective at regularising the neural network for better test set accuracy.

```
Total # of epochs: 15
Epoch #1 done...
Accuracy on test set after epoch #1 = 84.27%
Epoch #2 done...
Accuracy on test set after epoch #2 = 88.98%
Epoch #3 done...
Accuracy on test set after epoch #3 = 90.8%
Epoch #4 done...
Accuracy on test set after epoch #4 = 91.7%
Epoch #5 done...
Accuracy on test set after epoch #5 = 92.46%
Epoch #6 done...
Accuracy on test set after epoch #6 = 93.15%
Epoch #7 done...
Accuracy on test set after epoch #7 = 93.85%
Epoch #8 done...
Accuracy on test set after epoch #8 = 94.47%
Epoch #9 done...
Accuracy on test set after epoch #9 = 95.04%
Epoch #10 done...
Accuracy on test set after epoch #10 = 95.47%
Epoch #11 done...
Accuracy on test set after epoch #11 = 95.83%
Epoch #12 done...
Accuracy on test set after epoch #12 = 96.00999999999999%
Epoch #13 done...
Accuracy on test set after epoch #13 = 96.19%
Epoch #14 done...
Accuracy on test set after epoch #14 = 96.41%
Epoch #15 done...
Accuracy on test set after epoch #15 = 96.5%
Confusion Matrix:
[[ 964.  0.  0.  2.  0.  4.  2.  2.  2.  4.]
 [ 0. 1121.  2.  2.  0.  1.  2.  2.  5.  0.]
 [ 8.  2. 994.  9.  4.  0.  2.  5.  7.  1.]
 [ 2.  0.  4. 977.  0.  6.  0.  7.  10.  4.]
 [ 1.  0.  4.  0. 939.  1.  7.  2.  4. 24.]
 [ 0.  2.  0.  4.  2. 858.  0.  0.  9.  5.]
 [ 0.  3.  1.  0.  8. 14. 915.  0.  9.  0.]
 [ 1. 13. 12.  5.  6.  0.  0. 971.  2. 18.]
 [ 4.  1.  1.  2.  4.  0.  4.  3. 947.  2.]
 [ 0.  3.  0.  9. 14.  2.  0.  1. 10. 964.]]
...
Average error rate = 0.7000000000000006%
Standard deviation of error rate = 0.17216271373325748%

Final accuracy on test set = 96.5%
Final precision on test set = 96.4864097035871%
Final recall on test set = 96.47286915594118%
Final F1-score on test set = 0.9647124713461718
```

Figure 10: Learning Rate $\alpha=0.1$ and a Sigmoidal Activation

L2 Regularisation: The implementation of L2 regularisation is done as follows:

- The baseline loss function is now modified to include an L2 weight term:

$$J'(W; x, y) = J(W; x, y) + \frac{\lambda}{2m} \sum ||W^2||$$

- The gradient now turns out to be, for each weight:

$$\frac{\partial J'(W; x, y)}{\partial W} = \frac{\partial J(W; x, y)}{\partial W} + \frac{\lambda}{m} W$$

The new values for loss and gradient step are implemented in the code, in order to implement L2 regularisation. The plot is shown below. We observe a faster reduction of test loss

during the initial iterations. Also, during the whole training process, the test loss is at the lower end of the training loss. This is due to good generalisation, as opposed to the plot for the baseline model.

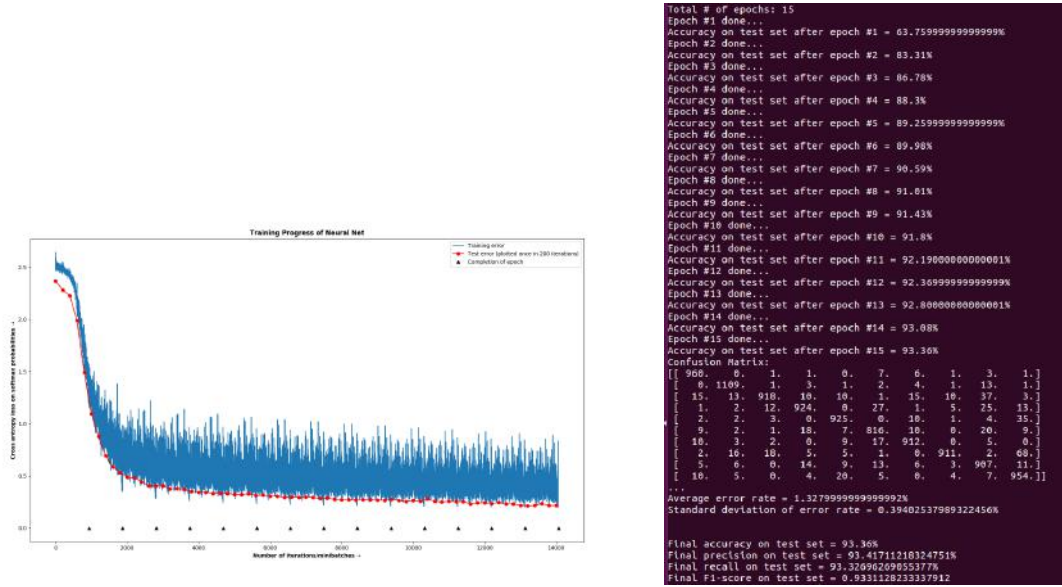


Figure 11: $\alpha=0.1$ and a Sigmoidal Activation

4 Hand-Crafted Features

I have used `scikit-learn` and `scikit-image` to implement both, my feature extractor-the Histogram of Oriented Gradients (HOG) and the K Nearest Neighbors and SVM algorithms.

HOG works by replacing each cell (of a certain number of pixels) by a single vector representing the histogram of the total gradient magnitudes of all pixels in that cell. The bins of the histogram are the bins of angle that a gradient can point towards, obtained by the division of the total angle. This makes HOG successfully and efficiently contain any edge-related information while compressing the amount of information there is to hold.

In our case, for example, an HOG when implemented on each 28×28 image with 9 orientations and a 7×7 cell gives 16 pixels, with each pixel containing a histogram (or 9 vector values). Therefore, each input has been reduced from 784 intensity values to 144 histogram values.

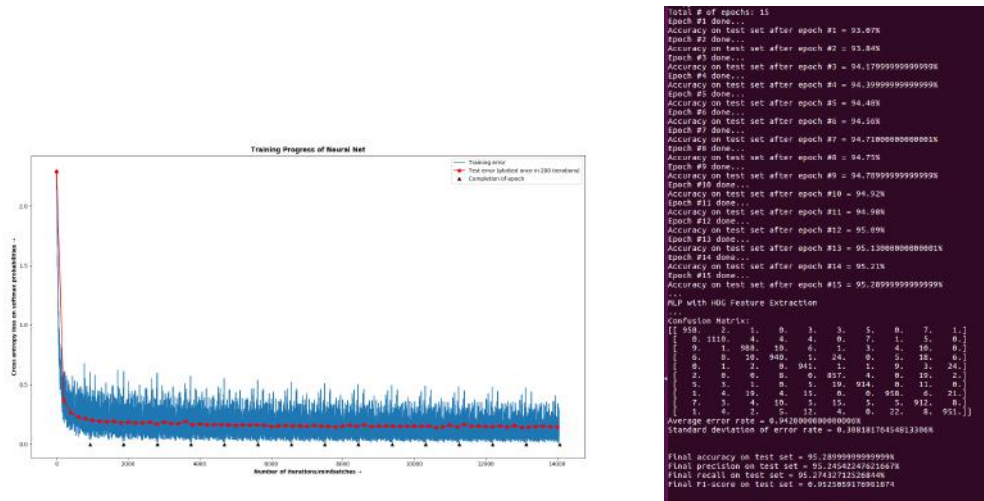


Figure 12: $\alpha=0.1$ and a tanh Activation

We observe that the K Nearest Neighbors algorithm has less accuracy, and the SVM has more accuracy. This can be explained by considering how the KNN algorithm only considers Euclidean distance as a metric. This can lead to a noisy 8 and a clean 8 being far off from each other and not being clustered together, for instance. A noisy 8 could even end up being closer to a clean 3. Hence, this leads to misclassification, higher error rate and lower accuracy.

```

...
K Nearest Neighbours with HOG Feature Extraction
...
Confusion Matrix:
[[ 962.    1.    1.    0.    3.    7.    1.    3.    1.]
 [   1. 1115.    2.    3.    3.    0.    6.    2.    0.]
 [  28.    2.  946.   34.    0.    0.    1.    6.   15.]
 [  16.    0.    5.  943.    0.   12.    0.   10.   17.]
 [    1.    6.    4.    0.  866.    0.    7.    3.    0.]
 [  11.    0.    0.   47.    2.  776.   26.    1.   24.]
 [   10.    6.    0.    0.    5.    4.  925.    0.    8.]
 [    1.    9.    9.    5.    8.    0.    0.  937.    1.]
 [   22.    4.    3.   23.   10.   23.    9.    5.  860.]
 [   12.    6.    1.   13.    6.    5.    0.   33.    8.]
Average error rate = 1.4899999999999999%
Standard deviation of error rate = 0.5695963483029007%

Final accuracy on test set = 92.55%
Final precision on test set = 92.70052942155239%
Final recall on test set = 92.42908395398152%
Final F1-score on test set = 0.9247543917137978

```

Figure 13: The Metrics of the KNN Algorithm

The SVM, however, does a better job as it finds the best margin hyperplane between supervised data sets with different labels. Besides, the linear SVM has less chance of overfitting, and hence does a better job of generalising to the test data.

```

...
SVM with HOG Feature Extraction
...
Confusion Matrix:
[[ 959.    1.    3.    1.    2.    2.    4.    1.    6.    1.]
 [   0. 1115.    3.    2.    3.    0.    6.    2.    4.    0.]
 [  11.    3.  966.   17.    6.    0.    4.    9.   16.    0.]
 [   2.    0.    8.  957.    4.   16.    0.    7.   14.    2.]
 [   2.    0.    7.    1.  919.    1.    5.    5.    6.   36.]
 [   2.    2.    5.   20.    1.  831.    6.    1.   23.    1.]
 [   9.    5.    0.    1.    8.    7.  918.    0.    8.    2.]
 [   0.    6.   22.    4.    9.    0.    0.  945.    5.   37.]
 [   9.    3.   10.   17.    9.   14.    4.    7.  890.   11.]
 [   5.    6.    2.   11.    7.    5.    1.   25.    7.  940.]]
Average error rate = 1.1200000000000001%
Standard deviation of error rate = 0.39864771415373673%

Final accuracy on test set = 94.39999999999999%
Final precision on test set = 94.38109196547127%
Final recall on test set = 94.34861318927302%
Final F1-score on test set = 0.9435873201063922

```

Figure 14: The Metrics of the SVM Algorithm
