# Project 3 Report

**Redwan Ibne Seraj Khan(UBITname: rikhan)**
Department of Computer Science and Engineering
University at Buffalo
Buffalo, NY-14228
*rikhan@buffalo.edu*

## Abstract

In this project different machine learning models like logistic regression, multilayer perceptron neural networks, support vector machine and random forest was used to solve a problem of determining the similarity between the samples of the MNIST dataset. After training the machine learning model on MNIST dataset the model was later tested on the USPS dataset. At first an ensemble of four classifiers was implemented and later the results of the individual classifiers were combined to make a final decision.
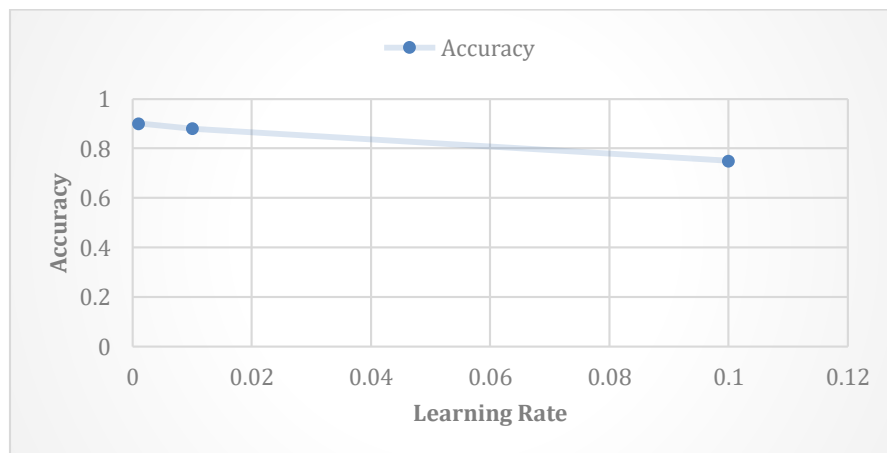
## 1    Logistic Regression

The softmax function was used in the logistic regression model. This was done as sigmoid function only accounted for binary classification but in MNIST and USPS dataset were have to classify the digits into 10 class. After calculations the loss function was obtained which was used to update the weights by taking the differential of the loss function. The model was tested with various hyperparamters. The best model that was built had iterations of 1000 and learning rate of 0.01. The output from running the model in python notebook is given below.

```
prediction accuracy of MNIST set:0.902
prediction accuracy of USPS set:0.10000500025
```
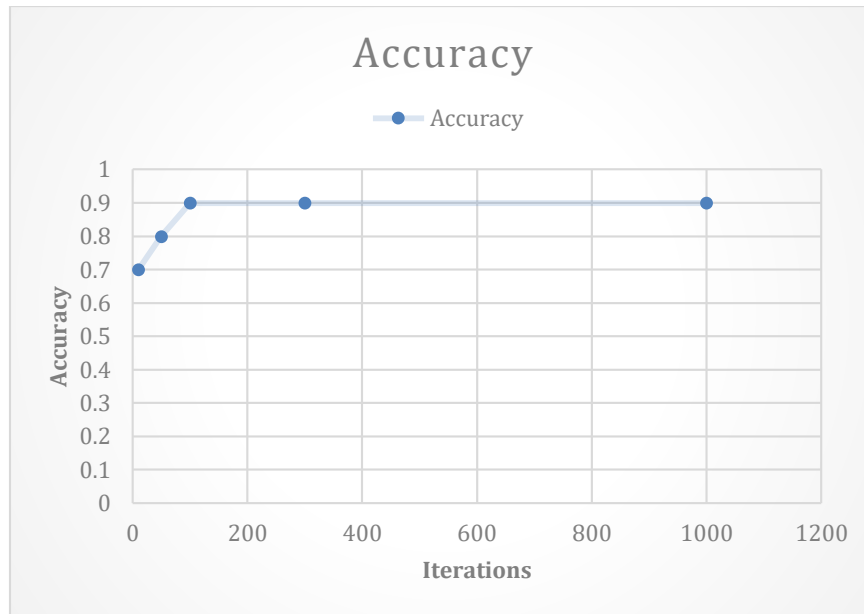
### 1.1    Learning rate

As the learning rate was decreased the accuracy increased but the time to run the model increased.
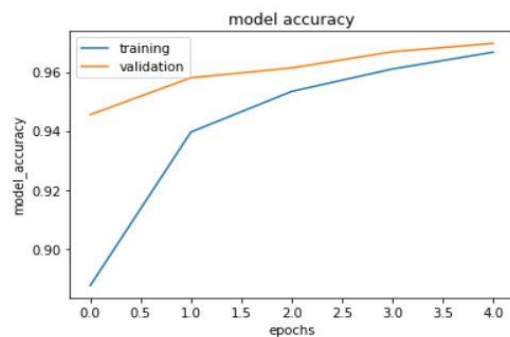
35  ## 1.2    Number of Iterations
36
37  As the number of iterations increased the accuracy increased for some time but then it became
38  concentrated. The time also increased with the increase in iterations.
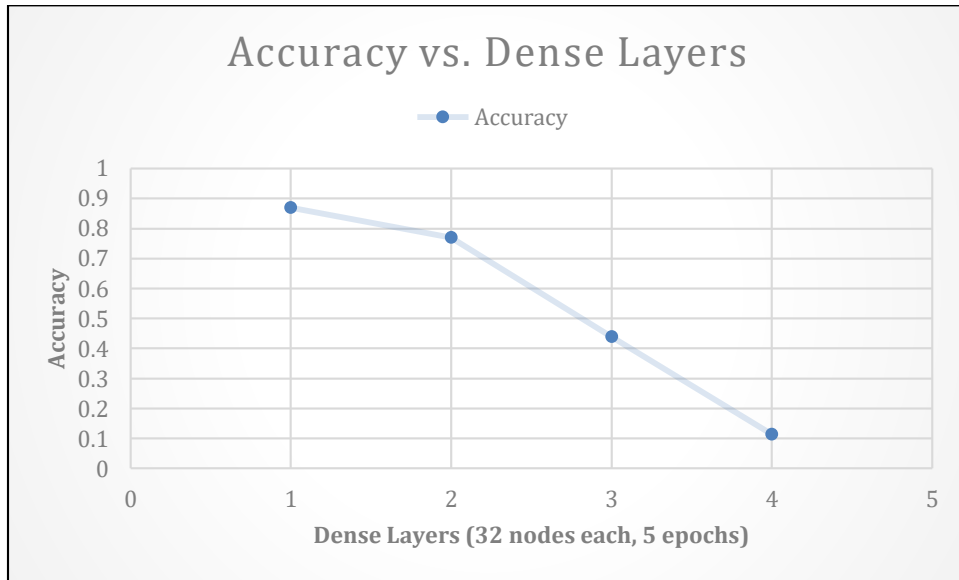39



40
41
42  # 2    Neural Networks
43
44  A sequential model was used for neural networks. Four dense layers having a combination of
45  different activation function and nodes was used. The first layer had 1024 nodes and relu
46  activation function. The second layer had 512 nodes and the sigmoid activation function. The third
47  layer had 128 nodes with relu activation function. The last layer had 10 nodes and used softmax
48  activation function.



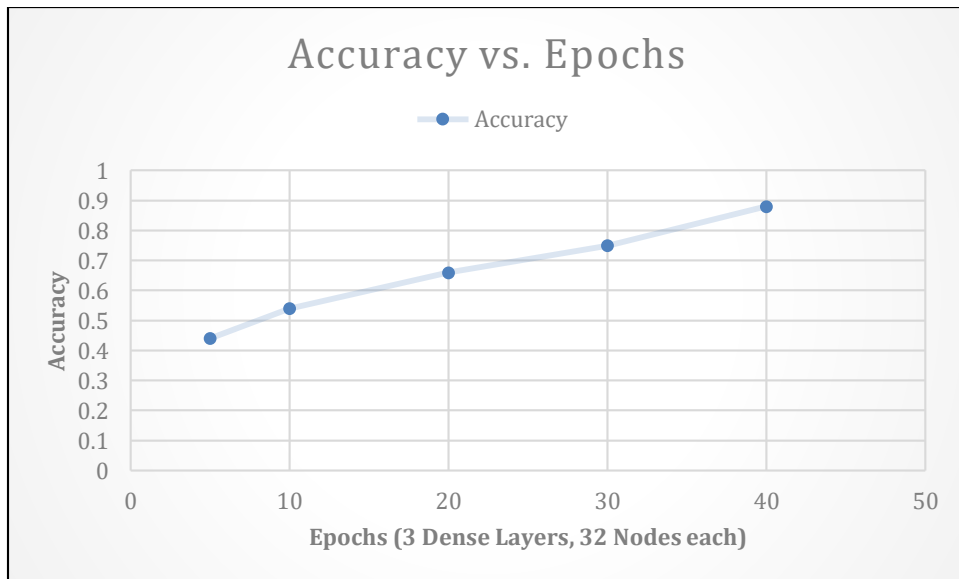Test loss: 0.1108
Test accuracy: 0.9657

49
50  The model was tested with different hyperparameters. The results have been shown below.
51
52
53
54
55
56
57
58

## 2.1    Number of Dense Layers and Activation Functions
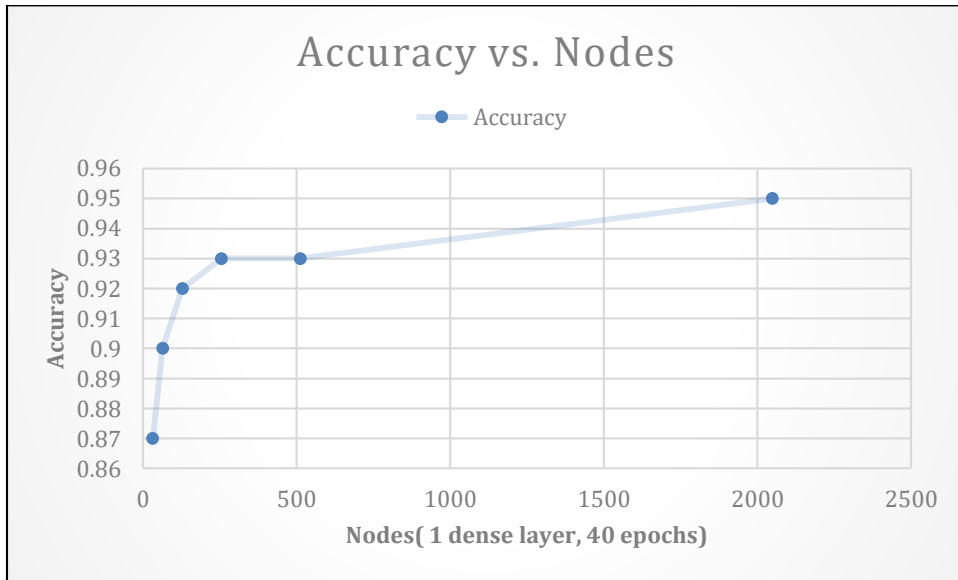


Accuracy vs. Dense Layers

The number of dense layers was gradually increased keeping the number of nodes constant. Sigmoid activation function was used in the dense layers. It was seen that the performance deteriorated when increasing the number of layers keeping the number of nodes constant.

## 2.2    Number of Epochs



Accuracy vs. Epochs

It was seen that the model performed better if the number of epochs was increased. The number of nodes was kept constant at 32 Nodes and 3 dense layers were used at various combinations of activation functions like relu and sigmoid.

## 2.3    Number of Nodes



The number of dense layers was kept constant at one and the number of epochs was also kept constant at 40. It was seen that as the number of nodes were increased there was a remarkable increase in performance.


# 3    Random Forest

The random forest classifier was run on the MNIST dataset with different hyperparameters. It was seen the best performance came from not restricting the max_leaf_nodes and using a high number of n_estimator like 1000. The n_jobs parameter was set to -1. The results that were obtained are given below.

```
Random forest accuracy on MNIST data:  0.9708
random forest accuracy on USPS data:  0.112455622781
```


## 3.1    Number of Trees in Forest (n_estimators)

It was seen that as the number of trees increased the accuracy of the model also increased.

Accuracy vs. n_estimator

x-axis: n_estimator( n_jobs =-1)
y-axis: Accuracy

## 3.2    Number of Max_Leaf_Nodes and n_jobs

It was seen that if max_leaf_nodes are decreased and n_jobs is set at value 1, then the accuracy decreases. It is reasonable as the trees cannot flourish to their maximum potential.



Accuracy vs. n_estimator

x-axis: n_estimator (n_jobs = 1, max_leaf_nodes = 100)
y-axis: Accuracy

# 4    Support Vector Machine

LinearSVC was used to calculate the accuracy at first setting every other parameter at the default value. The results obtained are given below.

```
133   Linear SVM accuracy for MNIST data:  0.8538
134   Linear SVM accuracy for USPS data:   0.123756187809
135
136
137   4.1      SVM with different kernels
138
139   Using the default value for kernel produced an overfitting model that could not predict the digits
140   of the test set accurately. When SVC with linear kernel was used it produced a good model. The
141   accuracy was seen to be around 86% for MNIST data. The reason for this is that empty spaces
142   between different classes are created with a Gaussian function when using radial basis functions
143   and in case of linear kernel, the model is less tunable and is basically a linear interpolation.
144   Moreover, the nonlinear function has a higher variance and linear kernel has a lower variance.
145
146   4.2      SVM with radial basis functions having different gamma values
147
148   When the gamma value was set to 1, then the model became hugely overfitted with a training test
149   accuracy of 100%. As it was not a general model it performed poorly on the test datasets of
150   MNIST and USPS. The model showed an accuracy of around 11% on the test set of data. A better
151   illustration of the model would be through a confusion matrix which is shown below.
152
```
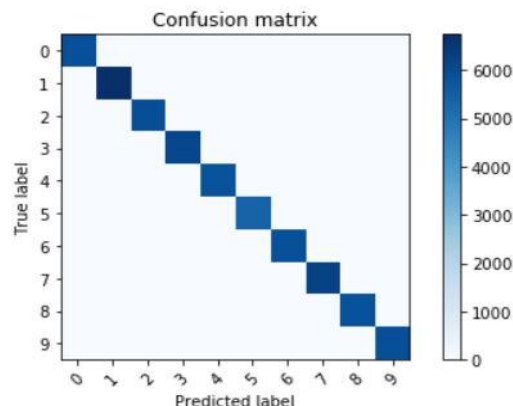
```
Confusion Matrix of MNIST Train data(Gamma = 1)

Confusion matrix, without normalization
[[5923    0    0    0    0    0    0    0    0    0]
 [   0 6742    0    0    0    0    0    0    0    0]
 [   0    0 5958    0    0    0    0    0    0    0]
 [   0    0    0 6131    0    0    0    0    0    0]
 [   0    0    0    0 5842    0    0    0    0    0]
 [   0    0    0    0    0 5421    0    0    0    0]
 [   0    0    0    0    0    0 5918    0    0    0]
 [   0    0    0    0    0    0    0 6265    0    0]
 [   0    0    0    0    0    0    0    0 5851    0]
 [   0    0    0    0    0    0    0    0    0 5949]]
```
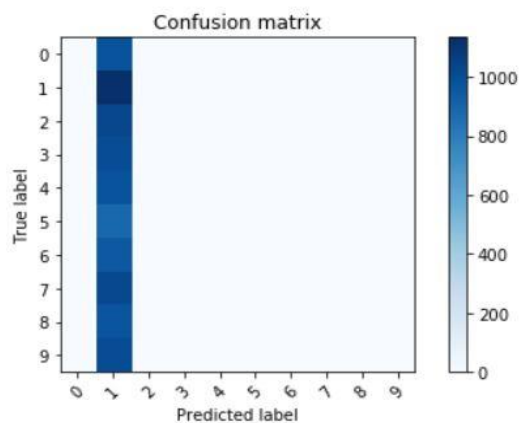


Confusion matrix

```
153
154
155
156   It can be seen that the model predicted everything correctly in the training set when gamma =1. It
157   shows the model is hugely overfitted. The prediction of the model can be seen in the later
158   confusion matrices which shows that the model almost always predicted every digit as 1 both in
159   the MNIST and the USPS dataset.
```

Confusion Matrix of MNIST Test data(Gamma = 1)

Confusion matrix, without normalization
```
[[   0  980    0    0    0    0    0    0    0    0]
 [   0 1135    0    0    0    0    0    0    0    0]
 [   0 1032    0    0    0    0    0    0    0    0]
 [   0 1010    0    0    0    0    0    0    0    0]
 [   0  982    0    0    0    0    0    0    0    0]
 [   0  892    0    0    0    0    0    0    0    0]
 [   0  958    0    0    0    0    0    0    0    0]
 [   0 1028    0    0    0    0    0    0    0    0]
 [   0  974    0    0    0    0    0    0    0    0]
 [   0 1009    0    0    0    0    0    0    0    0]]
```
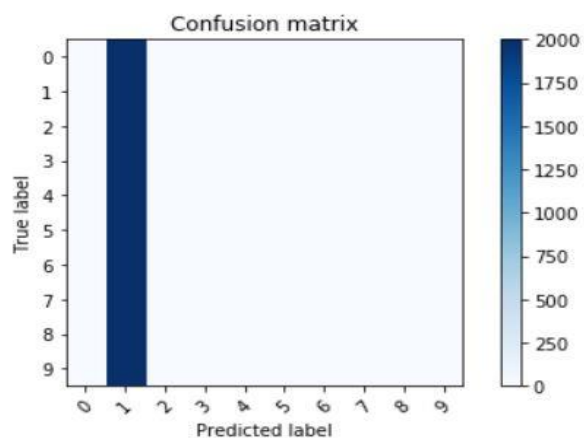


160
161
162

Confusion Matrix of USPS Test data(Gamma = 1)

Confusion matrix, without normalization
```
[[   0 2000    0    0    0    0    0    0    0    0]
 [   0 2000    0    0    0    0    0    0    0    0]
 [   0 1999    0    0    0    0    0    0    0    0]
 [   0 2000    0    0    0    0    0    0    0    0]
 [   0 2000    0    0    0    0    0    0    0    0]
 [   0 2000    0    0    0    0    0    0    0    0]
 [   0 2000    0    0    0    0    0    0    0    0]
 [   0 2000    0    0    0    0    0    0    0    0]
 [   0 2000    0    0    0    0    0    0    0    0]
 [   0 2000    0    0    0    0    0    0    0    0]]
```



163

## 5    Answer to Questions

**Q.1    We test the MNIST trained models on two different test sets: the test set from MNIST and a test set from the USPS data set. Do your results support the "No Free Lunch" theorem?**
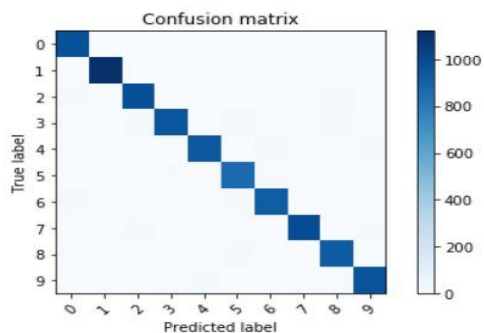
The **No Free Lunch** theorem says that there is no superior black-box optimization strategy. That is, we cannot device a model in such a way that it can give optimal results in all test data sets. In our case, the results support the theorem. This is because even though MNIST and USPS datasets have the same numerical digits (0-9), after training our model on MNIST dataset, when we tested the model against the USPS dataset the accuracy in prediction was really poor. Whereas we could obtain around 90% accuracy in all of our models with MNIST dataset, the accuracy in prediction for the USPS dataset was only around 20-30 %. These results go on to show that training a machine learning model in one set of data cannot give the same performance in another set of data.

**Q.2    Observe the confusion matrix of each classifier and describe the relative strengths/weaknesses of each classifier. Which classifier has the overall best performance?**

While passing the training and testing sets in the model first to_categorical of keras.utils was used to convert class vectors to binary class matrices. After training of the model the predictions were converted back to digits from binary class matrices using a decode function so that we can create a confusion matrix.

**Neural Network Confusion Matrix (MNIST dataset)**

```
Confusion matrix, without normalization
[[ 972    0    1    2    0    1    2    1    1    0]
 [   0 1120    2    1    0    1    3    1    7    0]
 [  11    4  987    4    1    3    4    5   12    1]
 [   2    0    9  954    0   23    0    6   11    5]
 [   1    0    4    0  947    0    9    1    4   16]
 [   4    1    0    8    2  863    6    1    4    3]
 [  12    3    0    0    5    9  926    0    3    0]
 [   1    7   10    2    1    0    0  995    2   10]
 [   6    2    1    4    6   10    5    4  930    6]
 [   3    4    1    5   15    3    1    5    6  966]]
```



Confusion matrix

## Neural Network Confusion Matrix (USPS dataset)

```
Confusion matrix, without normalization
[[ 125  769  591    0   12  400    0   40    0   63]
 [   0 1737  131    0    0   62    0   70    0    0]
 [   0 1000  879    0    0  119    0    1    0    0]
 [   0 1251  172   46    0  525    0    6    0    0]
 [   2 1518   97    0   98  189    0   91    0    5]
 [   0  546   91    0    1 1359    0    3    0    0]
 [   4 1322  276    0    3  390    1    4    0    0]
 [   0 1645  117    0    0  117    0  121    0    0]
 [   1 1300   90    0    0  598    0   11    0    0]
 [   0 1555  146    2    1  129    0  157    0   10]]
```



Confusion matrix

From the confusion matrix in my neural network model it can be seen that the model could almost accurately predict the digits in case of MNIST dataset, but when trying to predict the digits in USPS dataset the model had a tendency to predict all of the digits as the digit one and digit 5 in some cases. The neural network model was really successful in predicting the digit 1 in the MNIST dataset.

## Random Forest Confusion Matrix (MNIST Dataset)
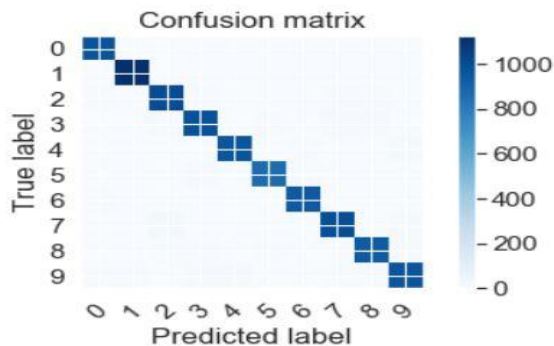
```
Confusion matrix, without normalization
[[ 970    0    0    0    0    2    2    1    4    1]
 [   0 1124    2    3    0    2    2    0    1    1]
 [   6    0  999    6    3    0    4    8    6    0]
 [   0    0    9  976    0    5    0    9    9    2]
 [   1    0    1    0  958    0    5    0    2   15]
 [   3    0    0   11    3  862    6    2    4    1]
 [   7    3    0    0    2    3  940    0    3    0]
 [   1    2   18    1    1    0    0  992    1   12]
 [   4    0    6    7    2    5    3    4  934    9]
 [   5    5    2    9   10    3    1    5    5  964]]
```



Confusion matrix

215 **Random Forest Confusion Matrix (USPS Dataset)**
216

```
Confusion matrix, without normalization
[[    0    53    0    0   18   83    0 1846    0    0]
 [    0   544    0    0    0   29    0 1427    0    0]
 [    0   248    1    0    1   71    0 1678    0    0]
 [    0    35    0    0    6  166    0 1793    0    0]
 [    0   146    0    0   15   44    0 1795    0    0]
 [    0   104    0    0    5  259    0 1632    0    0]
 [    0    98    0    0    1  210    0 1691    0    0]
 [    0   564    0    0    0   33    0 1403    0    0]
 [    0   136    0    0    6  550    0 1308    0    0]
 [    0   283    0    0    5   59    0 1653    0    0]]
```

Confusion matrix



217
218
219  From the confusion matrix in my random forest model it can be seen that the model could almost
220  accurately predict the digits in case of MNIST dataset, but when trying to predict the digits in
221  USPS dataset the model had a tendency to predict all of the digits as the digit seven. The random
222  forest model was really successful in predicting the digit 1 in the MNIST dataset.
223
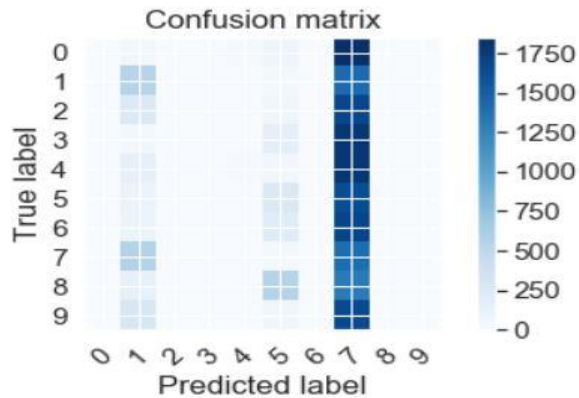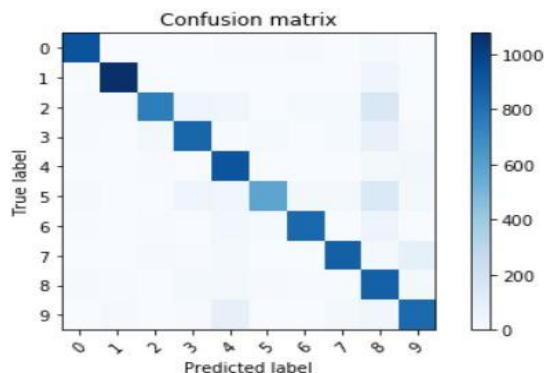224  **Support Vector Machine Confusion Matrix (MNIST Dataset)**
225

```
Confusion matrix, without normalization
[[ 928    0    3    3    5    8   11    4   16    2]
 [   0 1077    2    2    0    0    4    2   47    1]
 [   6    7  749   47   37    4   13   13  154    2]
 [   6    2   18  853    8   13    4   13   84    9]
 [   1    1    6    0  921    0    3    2   18   30]
 [  12    4    3   50   35  581   18   18  150   21]
 [   8    3    6    1   27   15  843    1   54    0]
 [   0    4    9    6   19    0    1  877   21   91]
 [  10    6    4   20   21   12    6    5  872   18]
 [   3   12    1    7   86    2    0   22   39  837]]
```

Confusion matrix



226

**Support Vector Machine Confusion Matrix (USPS Dataset)**

228
229

```
Confusion matrix, without normalization
[[  20   51    0    0    0 1414    0  515    0    0]
 [   0  404    0    0    0 1570    0   26    0    0]
 [   5  228    0    0    0 1730    0   36    0    0]
 [   5  167    0    0    0 1787    0   41    0    0]
 [   3  189    0    0    6 1534    0  268    0    0]
 [   3  144    0    0    0 1755    0   98    0    0]
 [   3   94    0    0    1 1820    3   79    0    0]
 [   3  340    0    0    0 1370    0  287    0    0]
 [   0   96    0    0    1 1694    0  209    0    0]
 [   4  240    0    0    2 1233    0  521    0    0]]
```



230
231 From the confusion matrix in my support vector machine model it can be seen that the model
232 could almost accurately predict the digits in case of MNIST dataset, but when trying to predict the
233 digits in USPS dataset the model had a tendency to predict all of the digits as the digit five. The
234 neural network model was really successful in predicting the digit 1 in the MNIST dataset
235 however in comparison to neural network and random it performed relatively worse in predicting
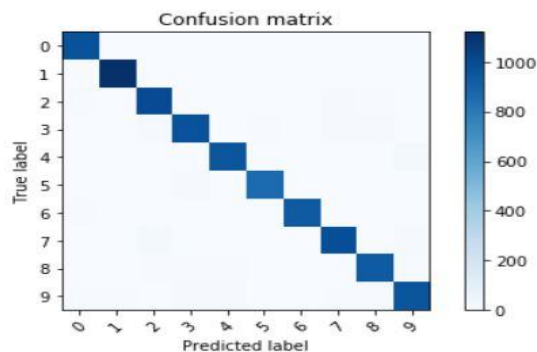236 all of the other digits.

237
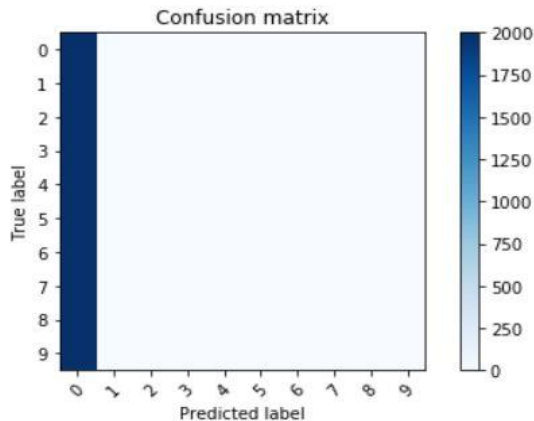238 **Logistic Regression Confusion Matrix (MNIST Dataset)**
239

```
Confusion matrix, without normalization
[[ 971    0    0    0    0    3    3    1    2    0]
 [   0 1120    3    3    0    2    3    1    3    0]
 [   6    0  999    3    3    1    4   10    6    0]
 [   0    0   10  971    0    8    0    9    9    3]
 [   1    0    2    0  954    0    4    1    2   18]
 [   3    0    0   13    3  862    4    1    4    2]
 [   6    3    1    0    3    4  939    0    2    0]
 [   1    3   23    2    0    0    0  986    3   10]
 [   4    0    7    5    5    4    3    4  934    8]
 [   7    5    2   12    9    1    1    5    6  961]]
```



240

241    **Logistic Regression Confusion Matrix (USPS Dataset )**
242

```
Confusion matrix, without normalization
[[2000    0    0    0    0    0    0    0    0    0]
 [2000    0    0    0    0    0    0    0    0    0]
 [1999    0    0    0    0    0    0    0    0    0]
 [1999    1    0    0    0    0    0    0    0    0]
 [2000    0    0    0    0    0    0    0    0    0]
 [2000    0    0    0    0    0    0    0    0    0]
 [2000    0    0    0    0    0    0    0    0    0]
 [2000    0    0    0    0    0    0    0    0    0]
 [2000    0    0    0    0    0    0    0    0    0]
 [2000    0    0    0    0    0    0    0    0    0]]
```



243
244
245    From the confusion matrix in my logistic regression model it can be seen that the model could
246    almost accurately predict the digits in case of MNIST dataset, but when trying to predict the digits
247    in USPS dataset the model had a tendency to predict all of the digits as the digit zero. Similar to
248    the other models that I have designed, logistic regression model was really successful in predicting
249    the digit 1 in the MNIST dataset.
250
251    From the analysis of the different models it can be seen that the best performance was achieved by
252    the neural network and random forest model. However, these models excelled in predicting
253    different digits accurately. The results are summarized as follows.
254
255    Digit 0: Neural Network
256    Digit 1: Random Forest
257    Digit 2: Random Forest
258    Digit 3: Random Forest
259    Digit 4: Random Forest
260    Digit 5: Neural Network
261    Digit 6: Random Forest
262    Digit 7: Neural Network
263    Digit 8: Random Forest
264    Digit 9: Neural Network
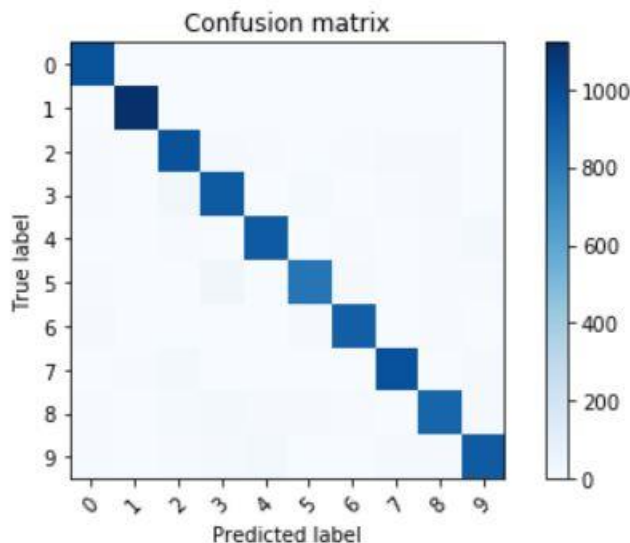265
266
267
268
269
270
271
272
273

274  **Q.3**  **Combine the results of the individual classifiers using a classifier combination**
275  **method such as majority voting. Is the overall combined performance better than**
276  **that of any individual classifier?**
277
278  The ensemble classifier model that I have used is the majority voting classifier. It was a
279  combination of the four models: Multilayer Perceptron Neural Network, Logistic Regression,
280  Random Forest and Support Vector Machine. Through the combination of the four models it was
281  seen that the combined performance showed a better performance than support vector machine
282  and logistic regression but performed a little worse than neural network and random forest. This is
283  with respect to the best model that I could make from neural network and random forest. If I had
284  used less n_estimators in random forest or less nodes and dense layers in multilayer perceptron
285  neural network the majority voting classifier would be at the top. However, it was seen that
286  majority voting performed really worse than the individual models when trying to predict the
287  USPS dataset. To understand the effect of the models on the classification accuracy I decided to
288  make a combined model using Logistic Regression, Random Forest and Neural Network. It was
289  seen that the accuracy of the digit prediction was around 92%. This showed the neural network
290  contributed a lot in determining a higher accuracy of the entire model.
291
292  The confusion matrices of the majority voting classifier for the MNIST and USPS datasets are
293  given below.
294

```
Confusion matrix, without normalization
[[ 968    1    1    0    0    3    4    1    2    0]
 [   0 1120    4    4    0    2    3    0    2    0]
 [   8    2  970    9    7    4    5   13   11    3]
 [   5    2   27  936    1   16    0    8   12    3]
 [   3    2    8    1  935    1    6    2    5   19]
 [   7    3    4   36    4  814    9    1    7    7]
 [  12    3    3    2    4   10  918    0    6    0]
 [   3    5   22    2    3    1    0  976    4   12]
 [   6    0   17   21    9    9    6    3  891   12]
 [   7    4    5   13   23    4    1    9   13  930]]
```
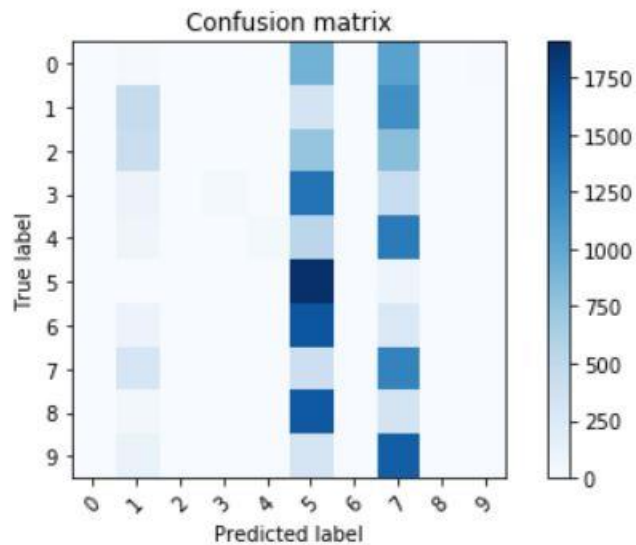


Confusion matrix

295

```
Confusion matrix, without normalization
[[    2   19    1    0    5  914    0 1051    0    8]
 [    0  466    0    0    0  333    0 1201    0    0]
 [    0  428    7    0    0  755    0  809    0    0]
 [    0  107    0   43    0 1408    0  442    0    0]
 [    0   77    0    0   43  532    0 1347    0    1]
 [    0    5    0    0    0 1905    0   90    0    0]
 [    0  106    2    0    1 1625    0  266    0    0]
 [    0  300    0    0    0  406    0 1294    0    0]
 [    0   58    0    0    0 1606    0  336    0    0]
 [    0  122    0    0    1  317    0 1558    0    2]]
```



Confusion matrix

296
297
298    From the confusion matrix it can be seen that the ensemble classifier predicted almost accurately
299    the MNIST dataset but in case of USPS dataset it predicted it to be 5 and 7 in most cases.
300
301    The accuracy of Majority Voting Classifier in MNIST Dataset:  95%
302    The accuracy of Majority Voting Classifier in USPS Dataset:  31%
303
304