# CSE 574: Introduction to Machine Learning (Fall 2018)
Instructor: Sargur N. Srihari

## Project 3: Image Classification with MNIST Dataset
November 19, 2018

Report By:
## Siddheswar Chandrasekhar

## Objective

The goal of this project is to implement machine learning methods to classify images in the MNIST dataset. We first implement an ensemble of four classifiers. Then the results of individual classifiers are to be combined to make a final decision.

**Plan of Action**

We train the following four classifiers using MNIST digit images:
1. Logistic Regression
2. Multilayer perceptron Neural Network
3. Random Forest Package
4. SVM (Support Vector Machine) Package

Based on the above implementations, we wish to answer the following questions:

1. We test the MNIST trained models on two different test sets: the test set from MNIST and a test set from the USPS dataset. Do our results support the "No Free Lunch" theorem?

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

3. Combine the results of the individual classifiers using a classifier combination method such as majority voting. Is the overall combined performance better than that of any individual classifier?

## Preparing the Datasets

The MNIST dataset is divided into three parts:
1.  Training dataset having 50,000 samples
2.  Validation dataset having 10,000 samples
3.  Testing dataset having 10,000 samples

We train our models using the MNIST training dataset (i.e. our models are trained over 50,000 samples). We then test our model over the MNIST testing dataset and the USPS dataset, both individually and combined. The USPS dataset has 20,000 samples (i.e. our models are tested over 30,000 samples)

## Logistic Regression (Linear Regression vs Logistic Regression [1] )

It's tempting to use the linear regression output as probabilities but it's a mistake because the output can be negative, and greater than 1 whereas probability cannot. As regression might actually produce probabilities that could be less than 0, or even bigger than 1, logistic regression was introduced.

- **Outcome**

    In linear regression, the outcome (dependent variable) is continuous. It can have any one of an infinite number of possible values.

    In logistic regression, the outcome (dependent variable) has only a limited number of possible values.

- **The dependent variable**

    Logistic regression is used when the response variable is categorical in nature whereas Linear regression is used when your response variable is continuous.

- **Equation**

    Linear regression gives an equation which is of the form $Y = mX + C$, means equation with degree 1.

    However, logistic regression gives an equation which is of the form $Y = e^X + e^{-X}$

- **Coefficient Interpretation**

    In linear regression, the coefficient interpretation of independent variables are quite straightforward (i.e. holding all other variables constant, with a unit increase in this variable, the dependent variable is expected to increase/decrease by xxx).

    However, in logistic regression, depends on the family (binomial, Poisson, etc.) and link (log, logit, inverse-log, etc.) you use, the interpretation is different.

- **Error minimization technique**

    Linear regression uses ordinary least squares method to minimise the errors and arrive at a best possible fit, while logistic regression uses maximum likelihood method to arrive at the solution.

    Linear regression is usually solved by minimizing the least squares error of the model to the data, therefore large errors are penalized quadratically.

    Logistic regression is just the opposite. Using the logistic loss function causes large errors to be penalized to an asymptotically constant.

# Logistic Regression Approach

1. Features $\rightarrow$ X

2. Target $\rightarrow$ y

3. Initialize weights (W) "normally" (i.e. using Normal Distribution)

4. $z = W^T X$

5.
$$\sigma = \frac{1}{1 + e^{-x}}$$

6. Gradient Descent
    $X^T (\hat{y} - y) \rightarrow$ / number of features $\rightarrow$ * learning rate $\rightarrow$ $- W \rightarrow$ new W

7. Repeat steps 4 thru 6 for N number of epochs

8. Test accuracy on unseen data

## Logistic Regression Observations

With learning rate as 0.01 and 500 epochs, we get the following results using Logistic Regression on our dataset:

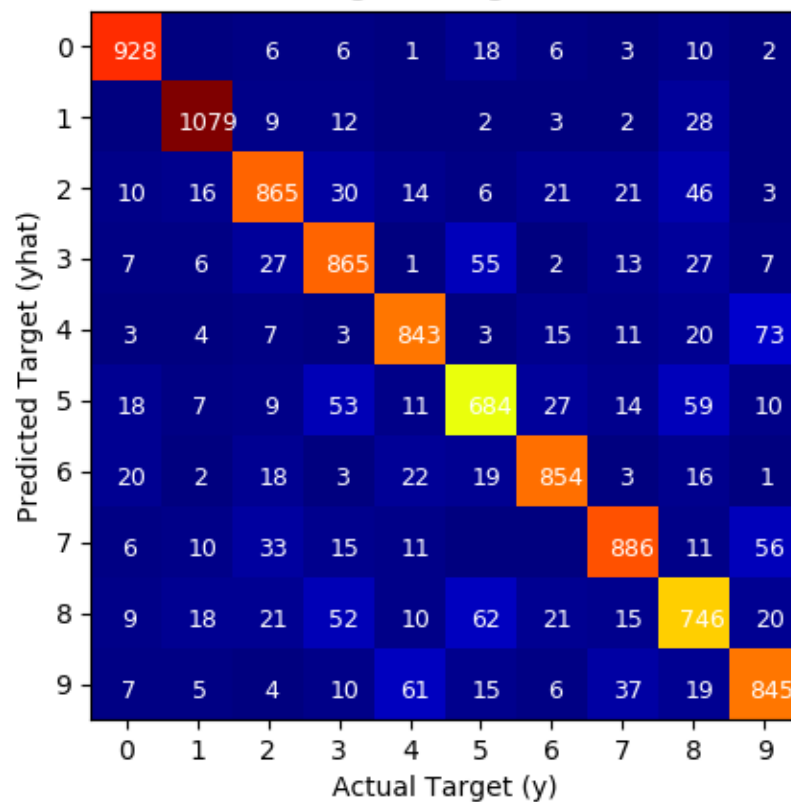| Testing on | Accuracy |
| --- | --- |
| MNIST dataset | 85.95 % |



Fig 1 – Logistic Regression on MNIST dataset

We see some pretty large inconsistencies, such as nine being predicted as four over seventy times. This is somewhat expected as Logistic Regression is not a good choice of model for multi-class regression.

| Testing on | Accuracy |
|---|---|
| USPS dataset | 25.12 % |



Fig 2 – Logistic Regression on USPS dataset

We clearly see poor accuracy when testing our model on the USPS dataset. By looking at the confusion matrix, we can deduce that several numbers in the USPS dataset are classified incorrectly by our model that has been trained over the MNIST dataset. An accuracy of 25% is proof of this.

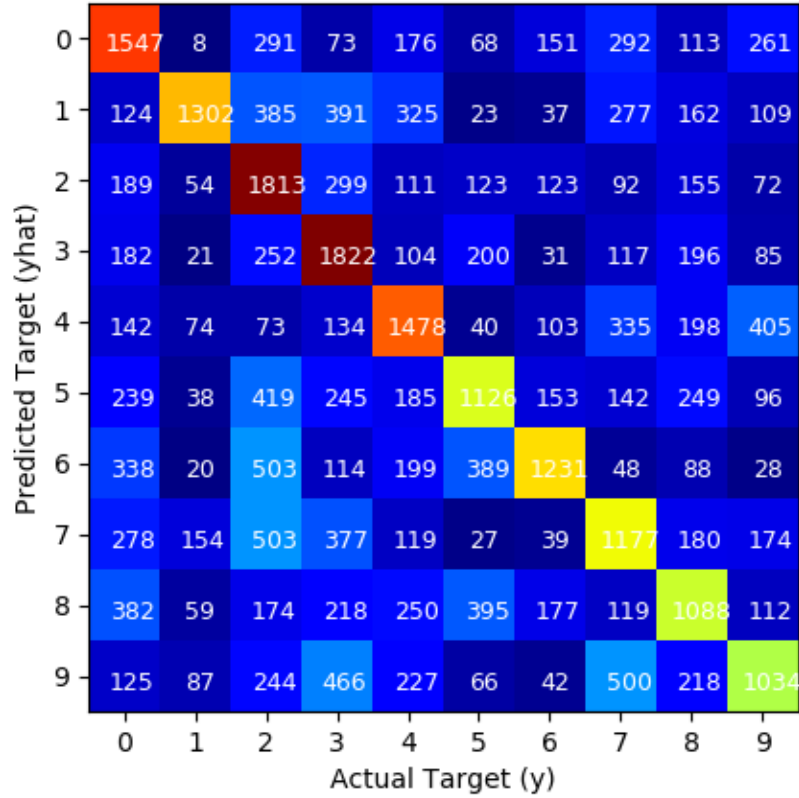| Testing on | Accuracy |
|---|---|
| MNIST + USPS datasets | 45.39 % |



Fig 3 – Logistic Regression on Combined datasets

As expected, the accuracy results on combined dataset is better than that of testing on USPS dataset alone, but worse than testing on MNIST dataset alone. We also notice that the accuracy is more inclined towards that of the accuracy when testing on USPS dataset alone. This is simply because the USPS dataset has twice the number of samples than that of the MNIST dataset. Hence the average is weighed more towards the USPS dataset testing accuracy.

## Neural Network

A neural network has a simple architecture that consists of an input layer, a black box which is a set of hidden layers, and an output layer. All the neurons in the network, connect with every neuron in its adjacent layer forming a mesh-like architecture.
The figure below is a good depiction of a simple neural network with three input neurons, one hidden layer of two neurons, and a single output neuron.
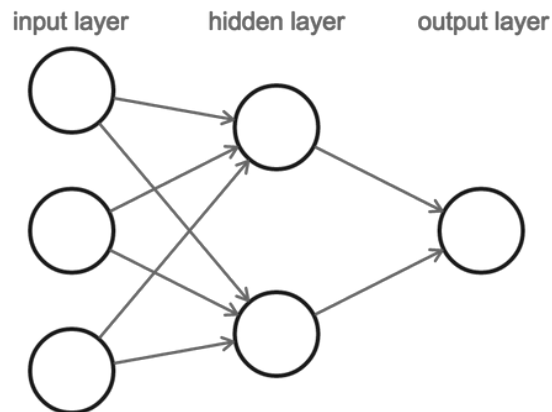


Fig 4 – Neural Network Architecture

We use ReLU (Rectified Linear Unit) as our Activation function for neural network which has the form:
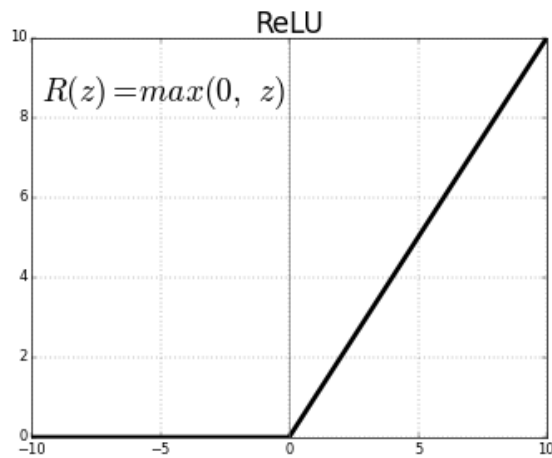


Fig 5 – ReLU Activation Function

# Neural Network Observations

With 500 epochs, 500 batch size and a decaying learning rate, we get the following results performing Neural Network on the MNIST dataset:

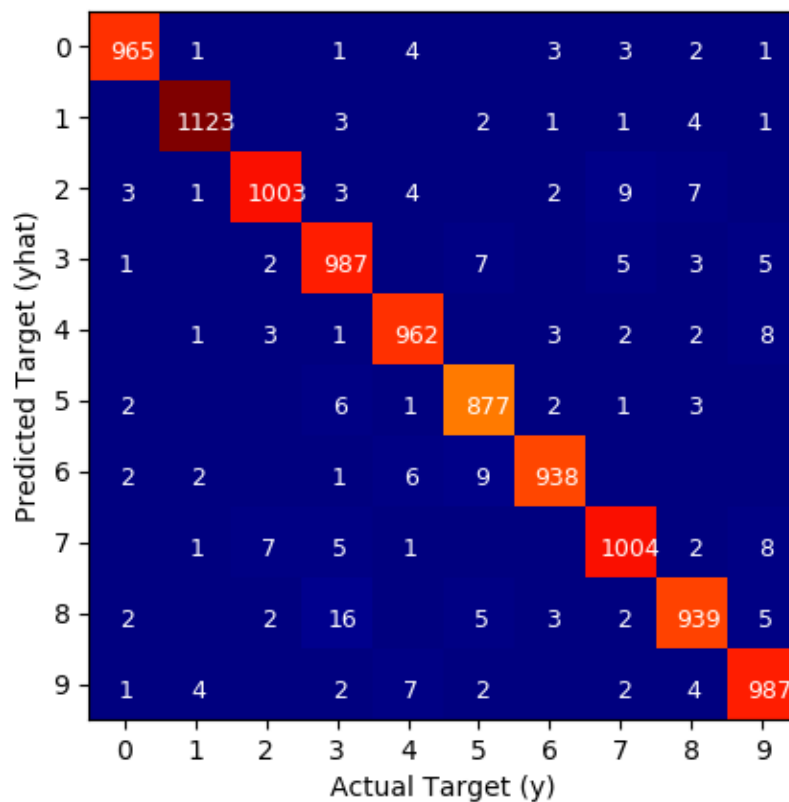| Activation Function | Testing on | Accuracy |
|---|---|---|
| ReLU | MNIST dataset | 97.85 % |



Fig 6 – Neural Network on MNIST dataset

The results are pretty impressive, with the only major inconsistencies being in four being wrongly predicted as nine, and vice versa.
We can hence deduce that the digits four and nine, look somewhat similar when handwritten.

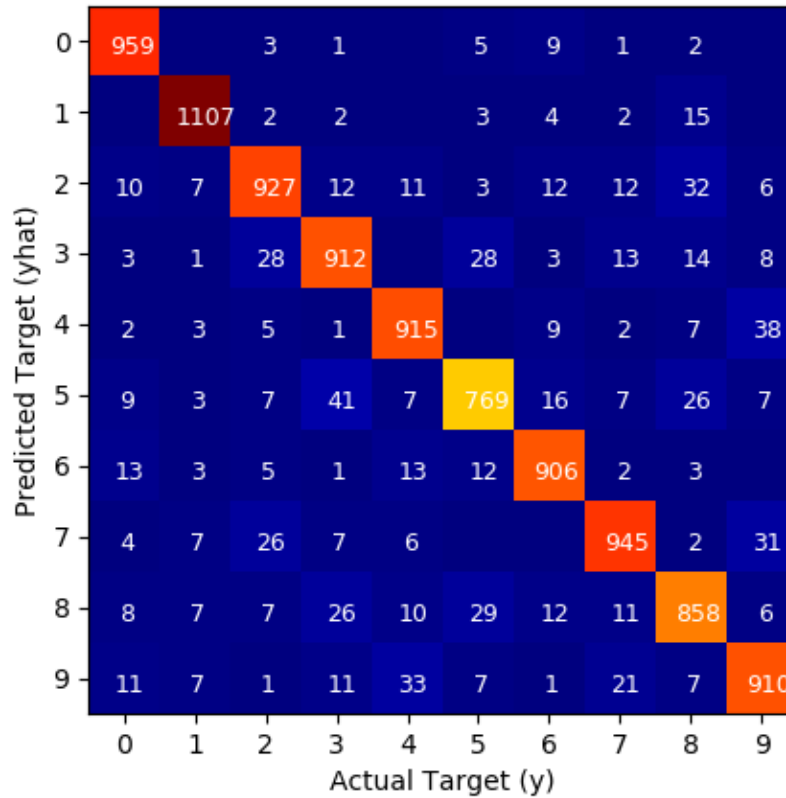| Activation Function | Testing on | Accuracy |
|---|---|---|
| Sigmoid | MNIST dataset | 92.08 % |



Fig 7 – Neural Network on MNIST dataset using Sigmoid

We see a drop in the accuracy. We observe similar results by using other activation functions such as tanh. We hence decide to stick with ReLU.

| Activation Function | Testing on | Accuracy |
|---|---|---|
| ReLU | USPS dataset | 46.89 % |

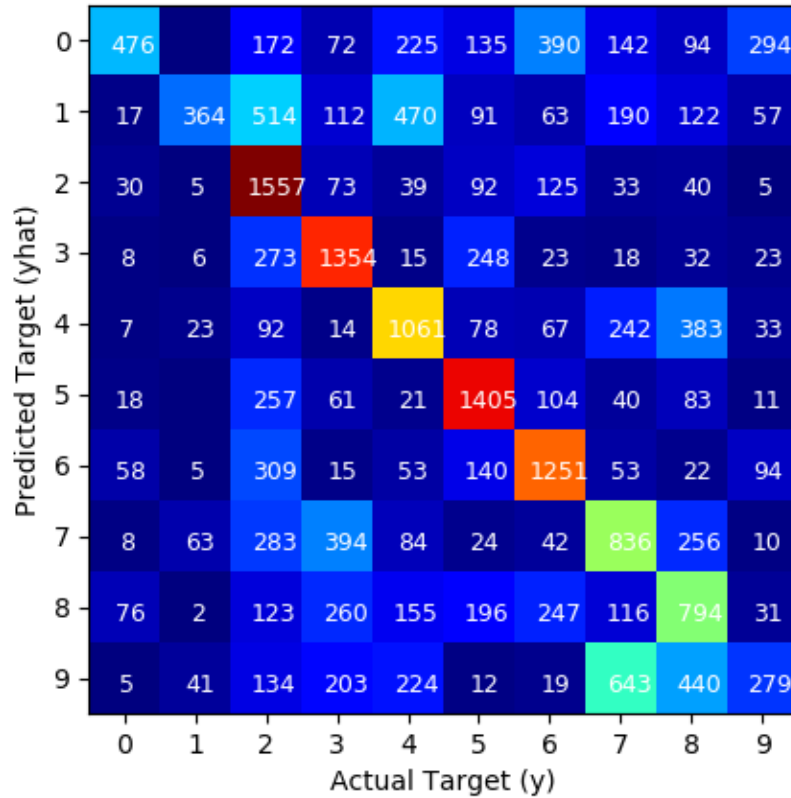Confusion Matrix for Neural Network on USPS dataset



Fig 8 – Neural Network on USPS dataset

Like in Logistic Regression, we clearly see poor accuracy when testing our model on the USPS dataset. By looking at the confusion matrix, we can deduce that several numbers in the USPS dataset are classified incorrectly by our model that has been trained over the MNIST dataset. Some notable inaccuracies are 9 being predicted as 0, seven being predicted as one or nine, three being predicted as seven or nine etc. As a matter of fact, we notice that seven is being incorrectly predicted as one almost as many times as it is being correctly predicted as seven. Moreover, nine is being incorrectly predicted as zero, almost three times that of it being correctly predicted as nine.

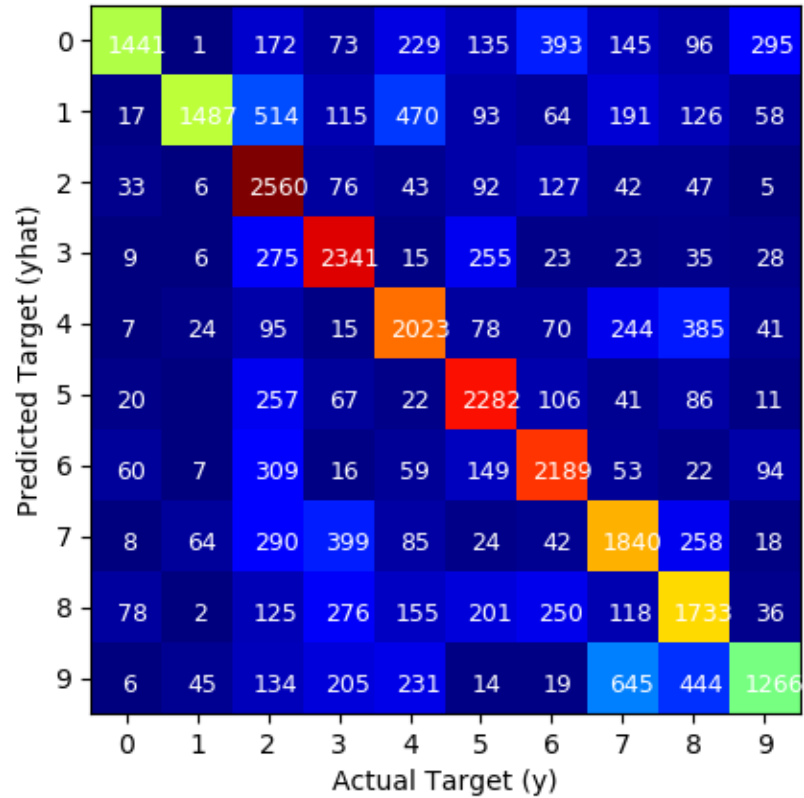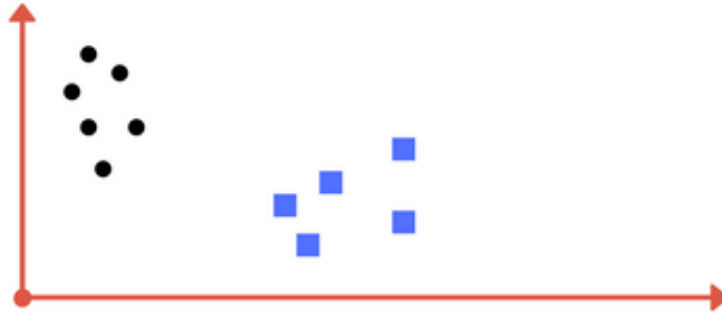| Activation Function | Testing on | Accuracy |
|---|---|---|
| ReLU | MNIST + USPS datasets | 63.88 % |



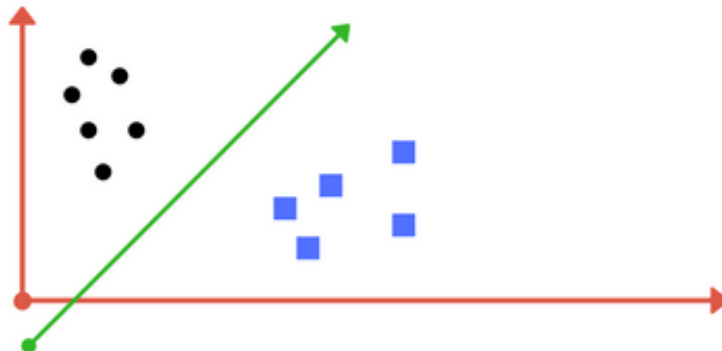Fig 9 – Neural Network on Combined datasets

As expected, the accuracy results on combined dataset is better than that of testing on USPS dataset alone, but worse than testing on MNIST dataset alone. Like in Logistic Regression, we also notice that the accuracy is more inclined towards that of the accuracy when testing on USPS dataset alone.

## Support Vector Machine (SVM)

*"A Support Vector Machine (SVM) is a discriminative classifier formally defined by a separating hyperplane. In other words, given labeled training data (supervised learning), the algorithm outputs an optimal hyperplane which categorizes new examples. In two dimentional space this hyperplane is a line dividing a plane in two parts where in each class lay in either side."*



We need to separate the two classes shown in the image above. An approach to this is to draw a line between them as shown in the image below:



Any point that is left of line falls into black circle class and on right falls into blue square class. "Separation of classes", that's what SVM does. It finds out a line/hyper-plane (in multidimensional space that separate outs classes)

**Tuning Parameters of SVM**

1. Kernel:

   The learning of the hyperplane in linear SVM is done by transforming the problem using some linear algebra. This is where the kernel plays role.

For linear kernel the equation for prediction for a new input using the dot product between the input (x) and each support vector (xi) is calculated as follows:

f(x) = B(0) + sum(ai * (x,xi))

This is an equation that involves calculating the inner products of a new input vector (x) with all support vectors in training data. The coefficients B0 and ai (for each input) must be estimated from the training data by the learning algorithm.

The polynomial kernel can be written as
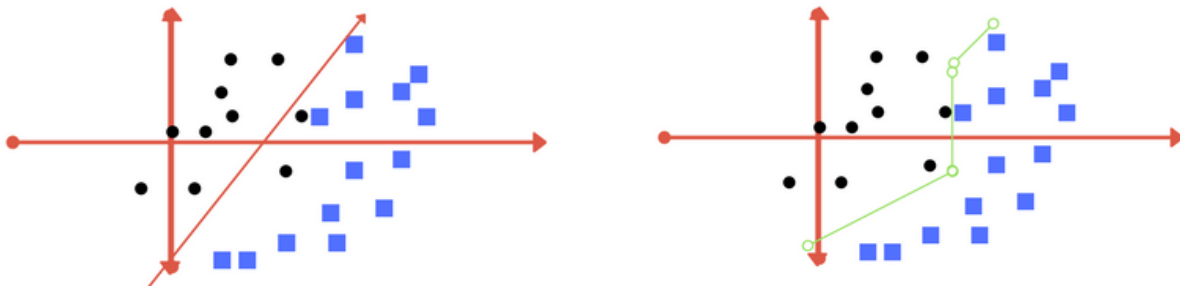
K(x,xi) = 1 + sum(x * xi)^d

and exponential as

K(x,xi) = exp(-gamma * sum((x—xi²)).

2. Regularization:

The Regularization parameter (often termed as C parameter in python's sklearn library) tells the SVM optimization how much you want to avoid misclassifying each training example.
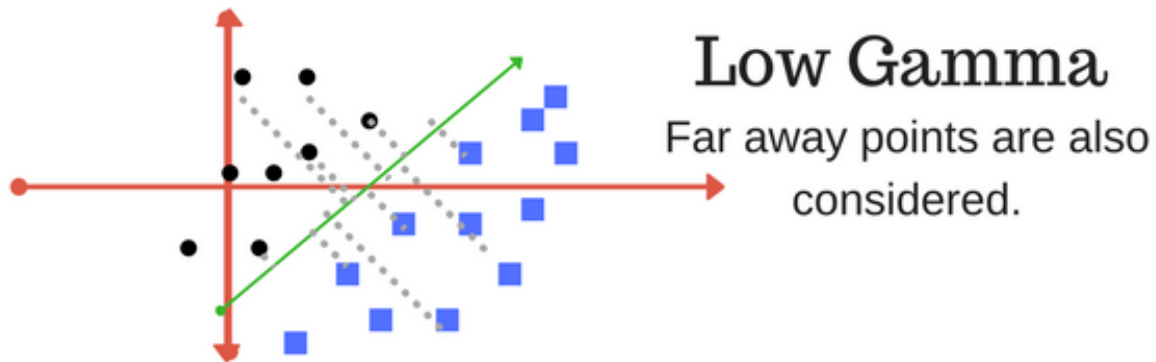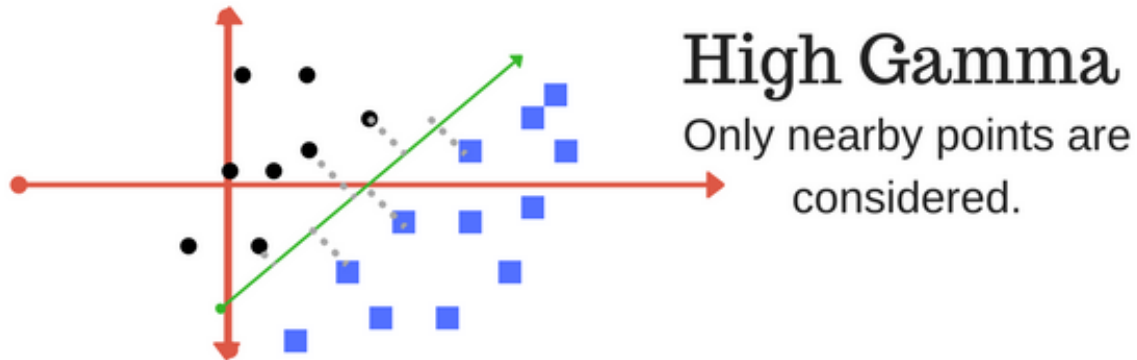
For large values of C, the optimization will choose a smaller-margin hyperplane if that hyperplane does a better job of getting all the training points classified correctly. Conversely, a very small value of C will cause the optimizer to look for a larger-margin separating hyperplane, even if that hyperplane misclassifies more points.

The images below are example of two different regularization parameter. Left one has some misclassification due to lower regularization value. Higher value leads to results like right one.
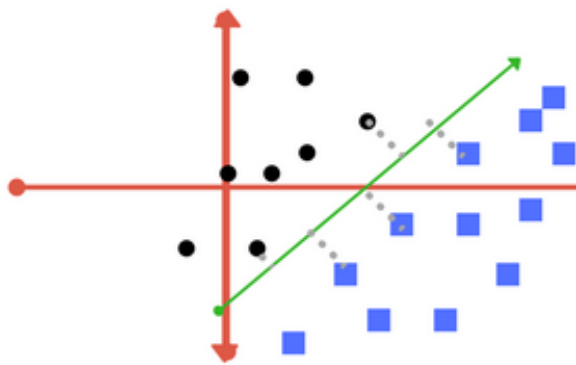
3. Gamma

The gamma parameter defines how far the influence of a single training example reaches, with low values meaning 'far' and high values meaning 'close'. In other words, with low gamma, points far away from plausible seperation line are considered in calculation for the seperation line. Where as high gamma means the points close to plausible line are considered in calculation.



**High Gamma**
Only nearby points are considered.



**Low Gamma**
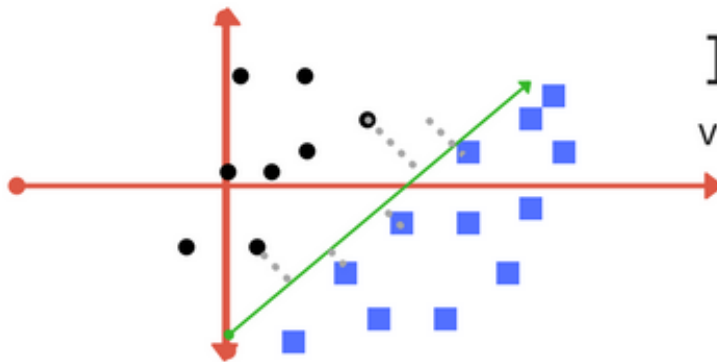Far away points are also considered.

4. Margin

A margin is a separation of line to the closest class points.

A good margin is one where this separation is larger for both the classes. Images below gives to visual example of good and bad margin. A good margin allows the points to be in their respective classes without crossing to other class.

# Good margin

equidistant as as far as
possible for both side.

# Bad margin

very close to blue class.

## SVM Observations

With *C* as 1.0 and using RBF (Radial Basisi Function) as the kernel, we get the following results performing SVM on our datasets:

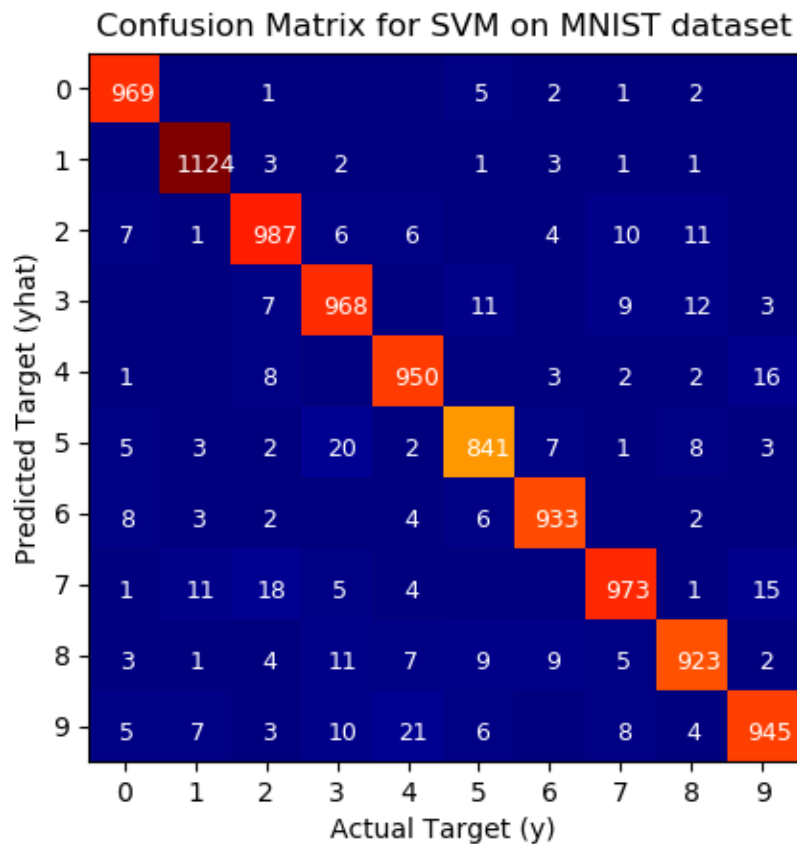| Testing on | Accuracy |
|---|---|
| MNIST dataset | 96.13 % |



Fig 10 – SVM on MNIST dataset

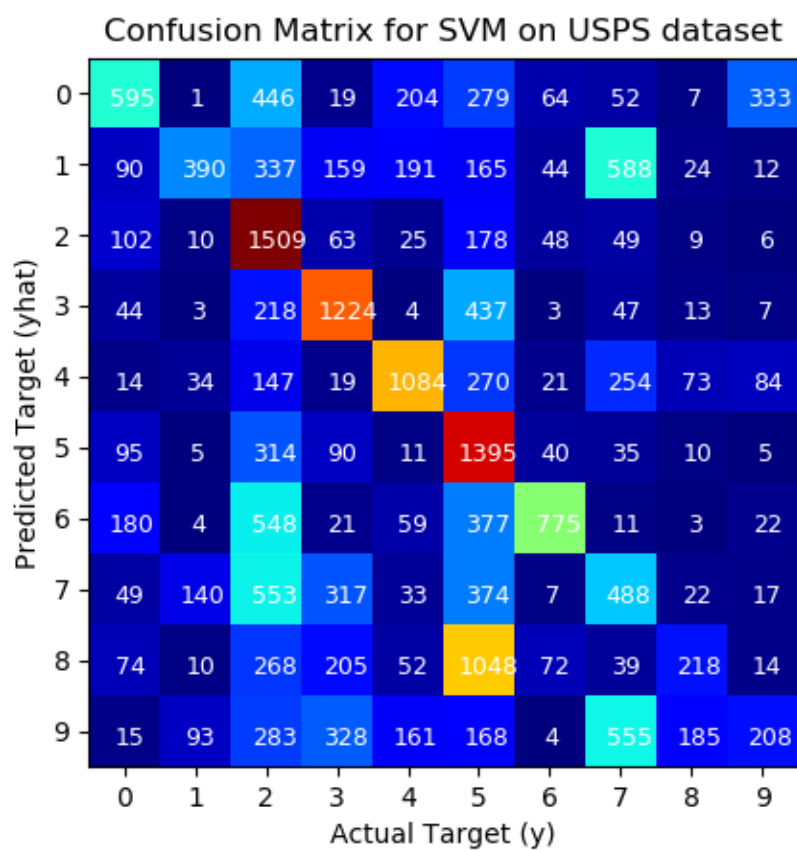| Testing on | Accuracy |
|---|---|
| USPS dataset | 39.43 % |



Fig 11 – SVM on USPS dataset

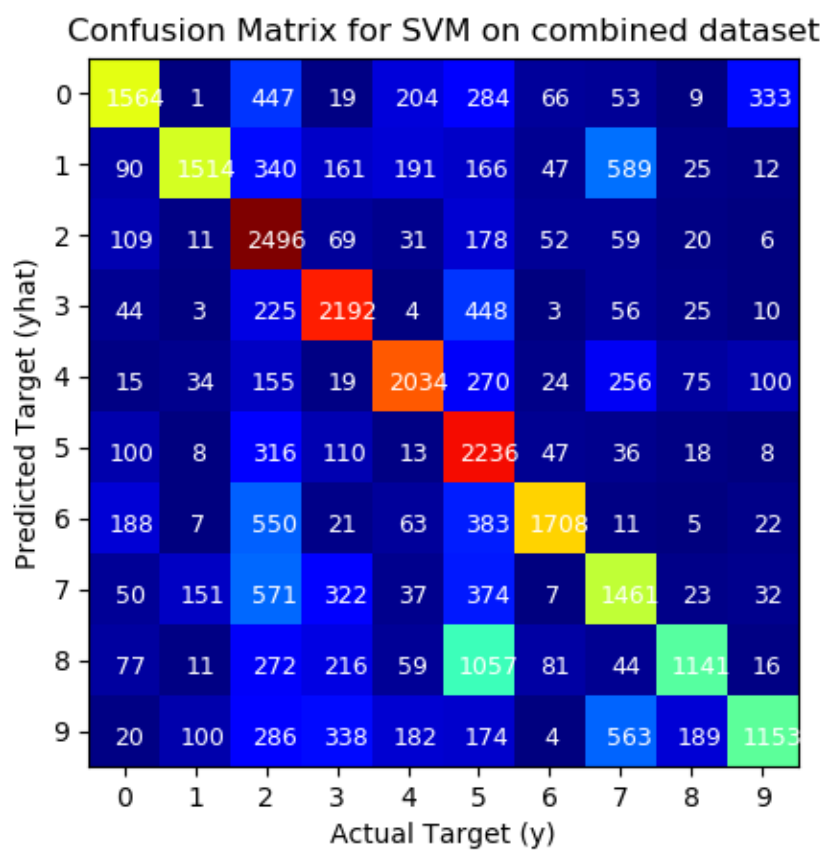| Testing on | Accuracy |
|---|---|
| MNIST + USPS datasets | 58.33 % |



Fig 12 – SVM on Combined datasets

# Random Forest

Random Forest is a supervised learning algorithm. Like the name suggests, the algorithm creates a forest and makes it somehow random. The "forest" it builds, is an ensemble of Decision Trees, most of the time trained with the "bagging" method. The general idea of the bagging method is that a combination of learning models increases the overall result.

To say it in simple words: Random forest builds multiple decision trees and merges them together to get a more accurate and stable prediction.

One big advantage of random forest is, that it can be used for both classification and regression problems, which form the majority of current machine learning systems. I will talk about random forest in classification, since classification is sometimes considered the building block of machine learning. Below you can see how a random forest would look like with two trees:



Fig 13 – Random Forest

Random Forest has nearly the same hyperparameters as a decision tree or a bagging classifier. Fortunately, you don't have to combine a decision tree with a bagging classifier and can just easily use the classifier-class of Random Forest. [3]

Random Forest adds additional randomness to the model, while growing the trees. Instead of searching for the most important feature while splitting a node, it searches for the best feature

among a random subset of features. This results in a wide diversity that generally results in a better model.

Therefore, in Random Forest, only a random subset of the features is taken into consideration by the algorithm for splitting a node. [3]

**Random Forest Observations**

With *n_estimators* as 100, *max_depth* as 20 and *random_state* as 0, we get the following results performing Random Forest on the MNIST dataset.

| Testing on | Accuracy |
| --- | --- |
| MNIST dataset | 96.74 % |



Fig 14 – Random Forest on MNIST dataset

| Testing on | Accuracy |
|---|---|
| USPS dataset | 38.89 % |



Fig 15 – Random Forest on USPS dataset

| Testing on | Accuracy |
|---|---|
| MNIST + USPS datasets | 58.18 % |

## Confusion Matrix for Random Forest on combined dataset

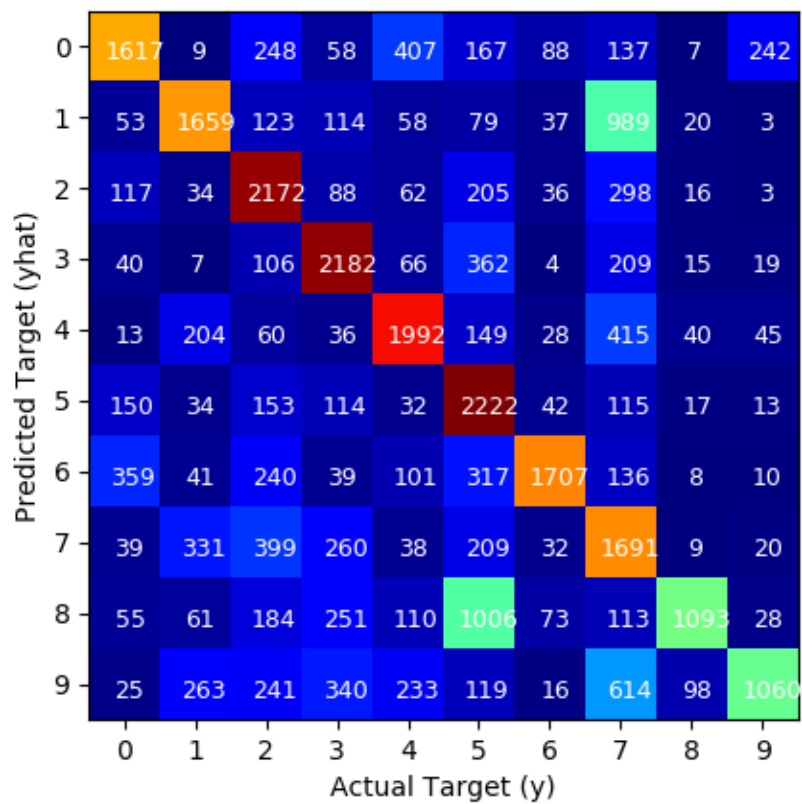| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1617 | 9 | 248 | 58 | 407 | 167 | 88 | 137 | 7 | 242 |
| 1 | 53 | 1659 | 123 | 114 | 58 | 79 | 37 | 989 | 20 | 3 |
| 2 | 117 | 34 | 2172 | 88 | 62 | 205 | 36 | 298 | 16 | 3 |
| 3 | 40 | 7 | 106 | 2182 | 66 | 362 | 4 | 209 | 15 | 19 |
| 4 | 13 | 204 | 60 | 36 | 1992 | 149 | 28 | 415 | 40 | 45 |
| 5 | 150 | 34 | 153 | 114 | 32 | 2222 | 42 | 115 | 17 | 13 |
| 6 | 359 | 41 | 240 | 39 | 101 | 317 | 1707 | 136 | 8 | 10 |
| 7 | 39 | 331 | 399 | 260 | 38 | 209 | 32 | 1691 | 9 | 20 |
| 8 | 55 | 61 | 184 | 251 | 110 | 1006 | 73 | 113 | 1093 | 28 |
| 9 | 25 | 263 | 241 | 340 | 233 | 119 | 16 | 614 | 98 | 1060 |

Predicted Target (yhat) / Actual Target (y)

Fig 16 – Random Forest on Combined datasets

# Ensemble of Classifiers

We combine our models using the concept of majority voting. In simple words, we find the 'mode' of the classifier predictions for each test sample.

# Ensemble Observations

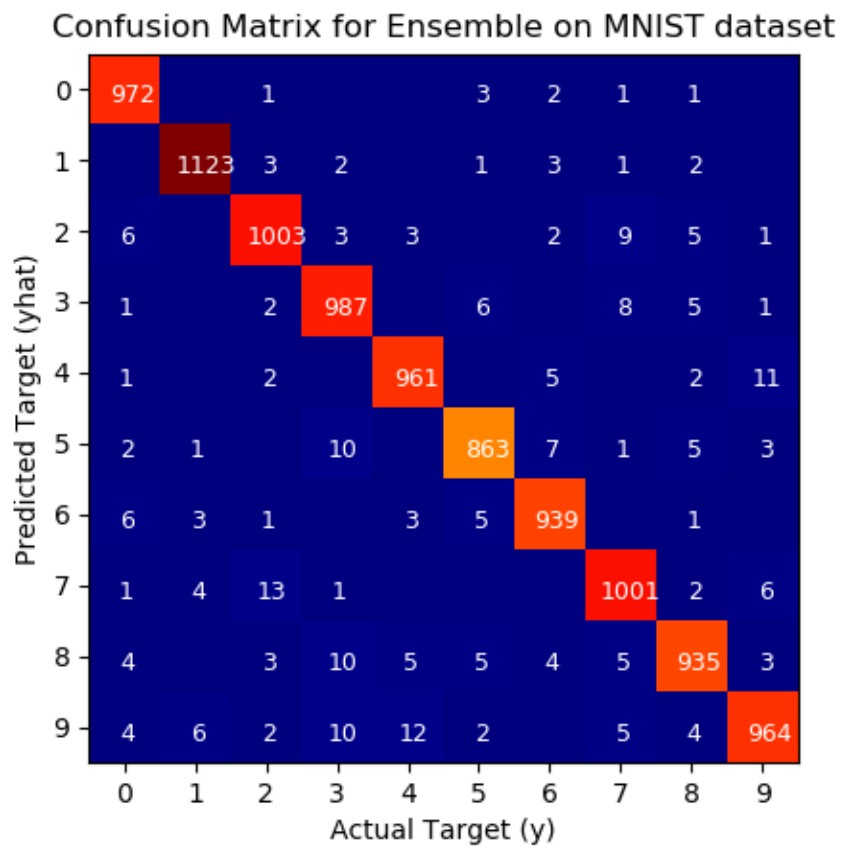| Testing on | Accuracy |
|---|---|
| MNIST dataset | 97.48 % |



Fig 17 – Ensemble on MNIST dataset

The accuracy is very similar to that of Neural Network, but slightly lower than that of Neural Network.

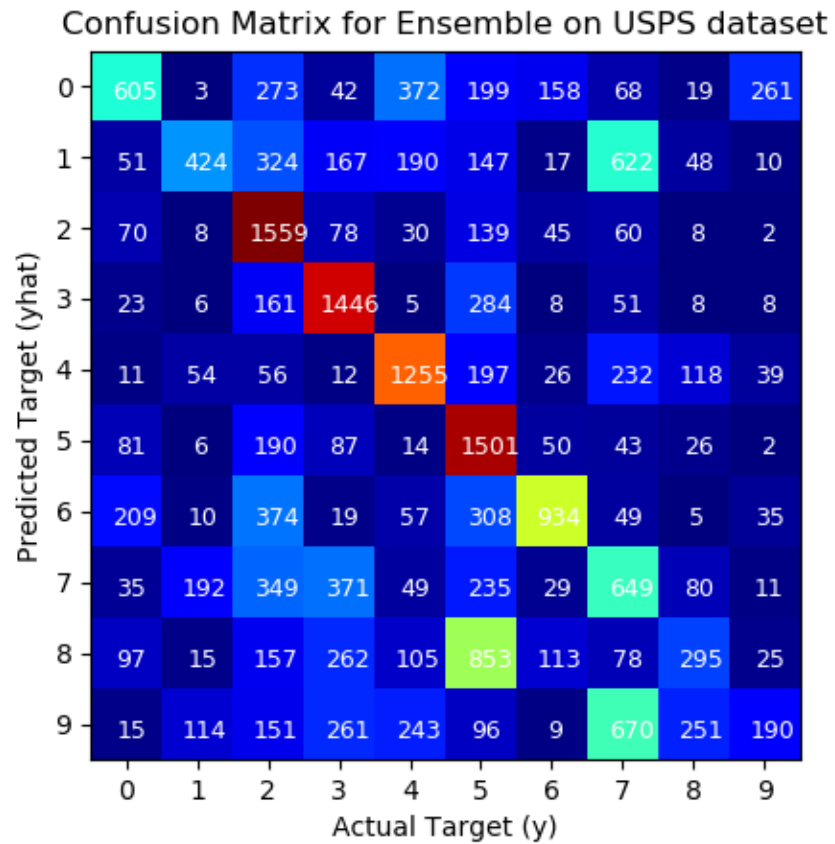| Testing on | Accuracy |
|---|---|
| USPS dataset | 44.29 % |



Fig 18 – Ensemble on USPS dataset

Like when tested on MNIST dataset, the accuracy of ensemble classifier when tested on USPS dataset is very similar to, but lower than that of Neural Network, which is the model with the highest accuracy on our test samples.

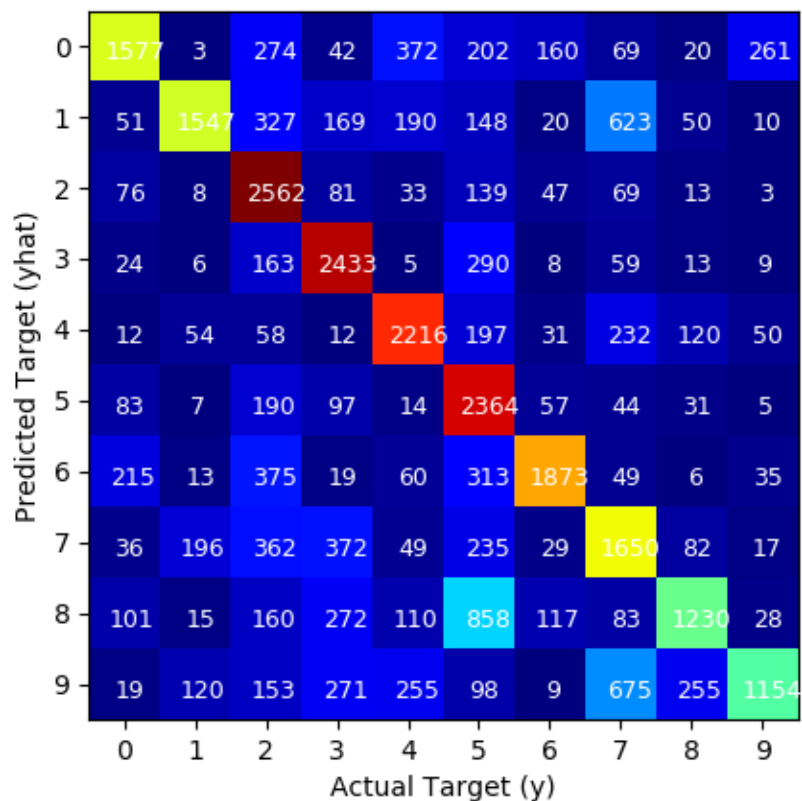| Testing on | Accuracy |
|---|---|
| MNIST + USPS datasets | 62.02 % |



Fig 19 – Ensemble on MNIST + USPS datasets

## Inferences

| Classifier / Test on | MNIST dataset | USPS dataset | MNIST + USPS dataset |
|---|---|---|---|
|  |  |  |  |
| **Logistic Regression** | 85.95 % | 25.12 % | 45.39 % |
| **Neural Network** | 97.85 % | 46.89 % | 63.88 % |
| **SVM** | 96.13 % | 39.43 % | 58.33 % |
| **Random Forest** | 96.74 % | 38.61 % | 57.99 % |
| **Ensemble** | 97.48 % | 44.29 % | 62.02 % |

Table 1 – Accuracy Table

- We clearly observe that Logistic Regression has the poorest performance compared to the other three models. We can hence deduce that Logistic Regression is not a good choice for a model when we are required to perform multi-class regression.

- We observe that our models perform well when tested with MNIST datasets, however their accuracies drop considerably when tested on USPS dataset. Hence, we can deduce that the MNIST and the USPS datasets are quite different from each other, and training a model over only one of them, will result in lower accuracies when tested over the other.

- When we test on a combination of MNIST and USPS models, as we would expected, the accuracy results is better than that of testing on USPS dataset alone, but worse than testing on MNIST dataset alone. We also notice that the accuracy is more inclined towards that of the accuracy when testing on USPS dataset alone. This is simply because the USPS dataset has twice the number of samples than that of the MNIST dataset. Hence the average is weighed more towards the USPS dataset testing accuracy.

- One of the major inconsistencies we can observe in the confusion matrices of our models is the considerably high count of four being wrongly predicted as nine, and vice versa.
Another similar inaccuracy is observed between classifying two and seven.
We can hence deduce that the digits four and nine, and two and seven, look somewhat similar when handwritten.

- We observe that Neural Network has the best performance, both, on each test dataset individually and overall. Hence, we can say that the results we obtained, do not support the "No Free Lunch" Theorem. [^]

- Our ensemble of classifiers performs better than Logistic Regression, SVM and Random Forest classifiers, however is slightly poorer than that of Neural Network's performance. This holds true for all the three test scenarios i.e. testing on MNIST samples, testing on USPS samples and testing on a combination of MNIST and USPS samples.

  One possible explanation to this could be that the other classifiers (i.e. Logistic Regression, SVM and Random Forest) tend to predict incorrectly more often than Neural Network (as can be seen by their individual accuracies) but moreover, also tend to make the same incorrect prediction. Hence, they have a stronger vote in deciding the prediction of the ensemble.

---

^ The "No Free Lunch" Theorem states that if one algorithm performs better than another for some class of problems, then it must perform worse in the remaining problems.

# References

1.  Sayali Sonawane
    "What is the difference between Linear Regression and Logistic Regression"
    *https://stackoverflow.com/questions/12146914/what-is-the-difference-between-linear-regression-and-logistic-regression*

2.  Savan Patel
    "Machine Learning 101 – Chapter 2"
    *https://medium.com/machine-learning-101/chapter-2-svm-support-vector-machine-theory-f0812effc72*

3.  Niklas Donges
    "The Random Forest Algorithm"
    *https://towardsdatascience.com/the-random-forest-algorithm-d457d499ffcd*