

Computer Assignment 3: Comparing SVM, BPNN and CNN performance for Pattern Recognition on MNIST Dataset

1 Introduction

The goal of this computer assignment is to design a Support Vector Machine (SVM) classifier and compare its performance for classifying MNIST digit database samples against the performance of two previously trained neural networks, one using a back propagation neural network (BPNN) and one using convolutional neural network (CNN). The MNIST database of handwritten digits, containing 60,000 training and 10,000 testing samples, is used for comparing the performance of all classifiers. Python and OpenCV were used as the programming language. The OpenCV package contains a very complete implementation of SVMs, allowing for different kernels and learning parameters among other things.

2 Theory

The theory of Back Propagation Neural Nets as well as Convolutional Neural Nets is omitted since it was discussed in the previous report.

2.1 Support Vector Machines

Support Vector Machines are learning machines that use the structure of the data to form decision boundaries instead of an empirical risk minimization technique.

Given a training data set $\{\mathbf{x}_p, d_p\}_{p=1}^P$ where $\mathbf{x}_p \in \mathcal{R}^N$ is the p th input pattern, $d_p = \pm 1$ is the desired label for two-class problems. If classes are linearly separable, then $N - D$ hyperplanes (DF) $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ exist such that

$$\mathbf{x}_p \in C_1 \implies g(\mathbf{x}_p) \geq 0 \text{ then } d_p = +1 \quad (1)$$

$$\mathbf{x}_p \in C_2 \implies g(\mathbf{x}_p) \leq 0 \text{ then } d_p = -1 \quad (2)$$

The goal is to find \mathbf{w}_{opt} and b_{opt} such that if $d_p = 1$ then $\mathbf{w}_{opt}^T \mathbf{x}_p + b_{opt} \geq 1$ while if $d_p = -1$ then $\mathbf{w}_{opt}^T \mathbf{x}_p + b_{opt} \leq -1$. Samples that meet the equality of previous formulas are called *support vectors*. SVMs aim to maximize the distance between each *canonical hyperplane*, the hyperplanes that satisfies the equality conditions.

The cost function used to find the weight vectors is $J(\mathbf{w}, b, \mathbf{x}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \sum_{p=1}^P \alpha_p [(\mathbf{w}^T \mathbf{x}_p + b)d_p - 1]$, which has $\mathbf{w}_{opt} = \sum_{p=1}^P \alpha_p d_p \mathbf{x}_p$.

The SVM often introduces the *kernel trick* to implicitly map the data to a higher dimensional space. The idea is that if the data is not linearly separable in the original dimensional space, than it might be linearly separable in a higher dimensional space. Several different kernel mappings are available, including Gaussian and polynomial as examined in this assignment.

Modifications to SVMs have been discovered such that the decision boundaries determined by an SVM can be extended beyond the 2 class case. The OpenCV implementation used in this assignment allows for any number of classes, and hence a 10 class SVM was used.

3 Network Configuration

Several network configurations were tested throughout the training phase for the SVM classifier. In particular, two different kernel types, polynomial and radial basis function (RBF), were used and are ultimately compared. The network configurations described below are representative of the classifiers that provided the best overall prediction accuracy on the validation set. No discussion of other kernel parameter values will be discussed as they provided worse classification accuracy and are therefore not worth mentioning.

The first kernel type used is a Gaussian kernel, which falls into the RBF family. The OpenCV implementation allows for multi-class SVMs with relatively little effort. A series of RBF kernel SVMs were trained and cross validation was used to identify the optimal values for C and γ which for this problem are 2 and 0.0014, respectively.

The second kernel type used is a polynomial kernel. Again, cross validation was used to identify a polynomial degree of 3 that provided the best classification accuracy.

4 Data Processing and Training

The data training set from MNIST contains 28×28 pixel images covering digits 0 – 9 with a class representation label corresponding to the digit, i.e. a 4 has a class label of ‘4’. Some random examples of the MNIST data set can be seen in Figure 1. The 60,000 training images and 10,000 testing images were combined into a total set of 70,000 unique samples. From the combined data set a training set, testing set, and validation set were formed randomly with 15%, 42.5% and 42.5% of the total number of samples, respectively. Each set of images contains an equal amount of samples from each class. Each training image was vectorized into a 784×1 vector for input into the SVM.

The dual optimization problem was solved to find the optimal α s and thus the corresponding weights. In both the polynomial and RBF kernel cases, the training data set was able to get classified with 100% accuracy. The selected parameters for each kernel that were the parameters that provided the highest classification accuracy.

5 Results and Discussion

The SVM was able to essentially match the performance of the best CNN while using the polynomial kernel of degree 3. This can be seen in the ROC curve and confusion matrix in Figure 2a and 2b. The classification accuracy of the SVM with the polynomial kernel is 99.63%. The accuracy of the SVM with the RBF kernel was essentially identical to the accuracy with the polynomial kernel. The best classification accuracy using the RBF kernel was only 99.41% using $C = 2$ and $\gamma = 0.0014$. The ROC curve and confusion matrix for the RBF kernel can be seen in Figure 2c and 2d.

The confusion matrix for the BPNN can be seen in Figure 3b, with the overall accuracy of the classifier being 96.3%.

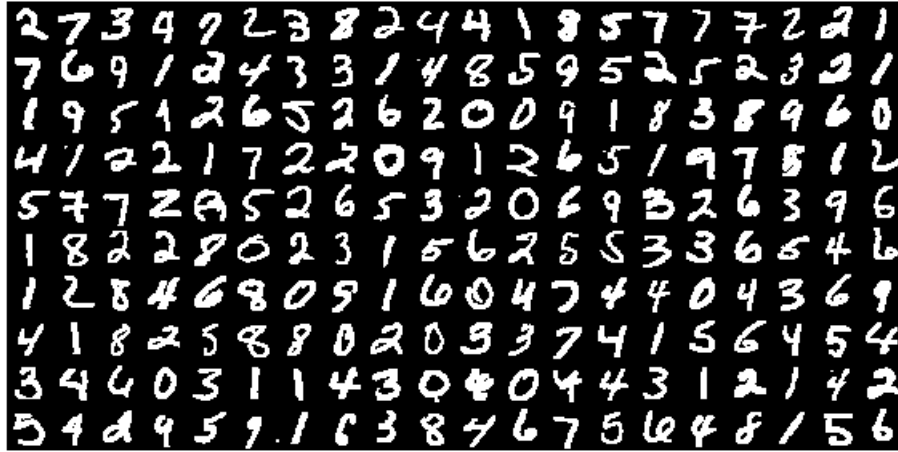
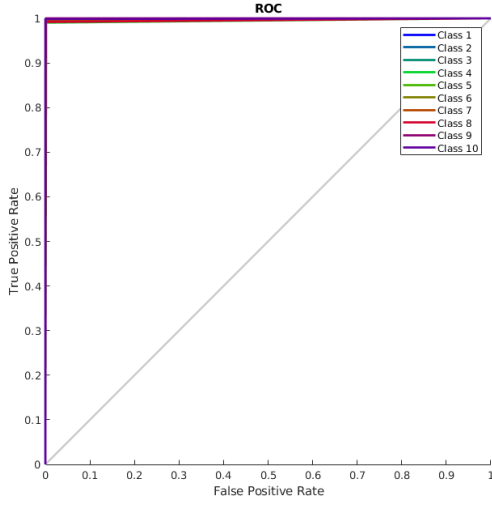


Figure 1: Samples from the MNIST data set

It should be noted that due to limitations in scripts used for generating the confusion matrix that all of the class labels have been incremented by 1, i.e. class 0 is represented by 1 in the figure. During the training phase there were times when the accuracy was fairly poor, around 60%, due to the uniform random sampling of the entire dataset not being truly random resulting in the training set missing representatives from one or more classes. This was readily apparent when looking at the ROC curves for the outputs and seeing one or more classes that had a concave curve. The ROC curves for the best BPNN can be seen in Figure 3a.

The confusion matrix for the CNNs can be seen in Figure 4b and 4d for the two and three layer networks. The two layer was able to classify correctly 98.6% of the time while the three layer was able to get another 1% performance boost on top of that. This is likely due to an additional amount of location invariance from the third pair of convolutional and pooling layers that extracted lower level features. The ROC curves for the best CNNs can be seen in Figure 4b and 4d.

It is clear that the CNNs are able to out perform the BPNN when trained and validated on the same data, while both kernel type SVMs can match the performance of the CNNs. One major advantage of the SVM classifiers over both the BPNN and CNN is that it is able to find the optimal decision boundary *much* faster during the training phase. This is due, in part, to an easier optimization problem. Only one weight vector and one bias value are required to be stored in order to make decisions for a two class SVM. The maximum number of elements in the weight vector is N , the number of pixels in the training image. This means that a total of $N + 1$ elements must be stored per class, for a total of $10 \times (N + 1) = 7850$ elements requiring storage. This is more than the 1189 elements of the smaller CNN, and more than the 5797 elements of the larger CNN.

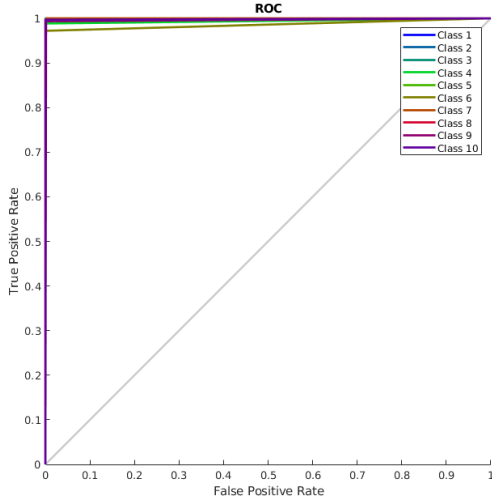


(a) ROC curve for polynomial kernel

Confusion Matrix

Output Class \ Target Class	1	2	3	4	5	6	7	8	9	10	Accuracy
1	617 10.2%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100%
2	0 0.0%	597 9.8%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	4 0.1%	0 0.0%	0 0.0%	0 0.0%	99.3%
3	0 0.0%	0 0.0%	606 10.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	99.8%
4	0 0.0%	0 0.0%	0 0.0%	610 10.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100%
5	0 0.0%	0 0.0%	0 0.0%	0 0.0%	608 10.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100%
6	5 0.1%	0 0.0%	0 0.0%	6 0.1%	0 0.0%	597 9.8%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	98.2%
7	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	593 9.8%	0 0.0%	0 0.0%	0 0.0%	99.8%
8	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	597 9.8%	0 0.0%	0 0.0%	99.8%
9	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	3 0.0%	0 0.0%	605 10.0%	1 0.0%	99.3%
10	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	610 10.1%	100%
Overall	99.2%	99.8%	100%	99.0%	99.8%	100%	99.5%	99.2%	100%	99.8%	99.6%

(b) Confusion matrix for polynomial kernel



(c) ROC curve for RBF kernel

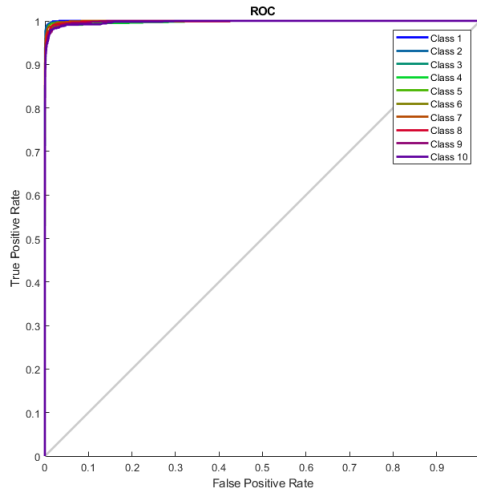
Confusion Matrix

Output Class \ Target Class	1	2	3	4	5	6	7	8	9	10	Accuracy
1	644 10.6%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	99.8%
2	0 0.0%	595 9.8%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100%
3	1 0.0%	0 0.0%	607 10.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	99.7%
4	0 0.0%	0 0.0%	1 0.0%	595 9.8%	0 0.0%	4 0.1%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	99.0%
5	0 0.0%	0 0.0%	0 0.0%	0 0.0%	572 9.4%	5 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	99.0%
6	0 0.0%	0 0.0%	0 0.0%	2 0.0%	0 0.0%	591 9.7%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	99.5%
7	0 0.0%	0 0.0%	0 0.0%	3 0.0%	1 0.0%	2 0.0%	607 10.0%	0 0.0%	2 0.0%	1 0.0%	98.5%
8	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	614 10.1%	0 0.0%	0 0.0%	100%
9	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	596 9.8%	1 0.0%	98.7%
10	2 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	628 10.3%	99.5%
Overall	99.5%	100%	99.7%	98.8%	99.8%	97.2%	100%	99.7%	99.3%	99.7%	99.4%

(d) Confusion matrix for RBF kernel

Figure 2: ROC curve and Confusion matrix for SVMs

Some interesting things to note for the RBF kernel SVM is that the value of γ is incredibly small. I was originally using a much larger (> 1) value of γ which was causing the network to classify all samples in the test and validation sets as a single class while being able to correctly classify all samples in the test set correctly. Because of this poor performance, I decided to perform HOG feature extraction and train on those features instead of the raw image. Doing so provided a 100% accuracy on the test, validation and training sets. The results of the classification on the HOG features are not presented here though, as after inquiring with Professor Azimi it was determined that only networks operating on raw images



(a) ROC curve

Confusion Matrix

	0	1	2	3	4	5	6	7	8	9	10
0	968 9.7%	0 0.0%	6 0.1%	0 0.0%	2 0.0%	5 0.0%	6 0.1%	0 0.0%	8 0.1%	96.6%	3.4%
1	0 0.0%	1116 11.2%	2 0.0%	0 0.0%	3 0.0%	1 0.0%	1 0.0%	6 0.1%	3 0.0%	98.2%	1.8%
2	1 0.0%	1 0.0%	990 9.9%	9 0.1%	4 0.0%	2 0.0%	2 0.0%	11 0.1%	5 0.0%	96.4%	3.6%
3	0 0.0%	2 0.0%	10 0.1%	966 9.7%	2 0.0%	12 0.1%	1 0.0%	9 0.1%	10 0.1%	94.9%	5.1%
4	0 0.0%	0 0.0%	4 0.0%	1 0.0%	945 9.4%	2 0.0%	5 0.0%	2 0.0%	6 0.1%	96.0%	4.0%
5	4 0.0%	1 0.0%	0 0.0%	10 0.1%	0 0.0%	852 8.5%	5 0.0%	1 0.0%	3 0.0%	96.9%	3.1%
6	4 0.0%	6 0.1%	4 0.0%	2 0.0%	6 0.1%	8 0.1%	932 9.3%	0 0.0%	5 0.0%	96.4%	3.6%
7	1 0.0%	3 0.0%	9 0.1%	10 0.1%	1 0.0%	0 0.0%	1 0.0%	981 9.8%	4 0.0%	96.5%	3.5%
8	2 0.0%	6 0.1%	6 0.1%	8 0.1%	2 0.0%	8 0.1%	4 0.0%	1 0.0%	925 9.2%	96.0%	4.0%
9	0 0.0%	0 0.0%	1 0.0%	4 0.0%	17 0.2%	2 0.0%	1 0.0%	17 0.2%	7 0.1%	95.1%	4.9%
10	98.8% 1.2%	98.3% 1.7%	95.9% 4.1%	95.6% 4.4%	96.2% 3.8%	95.5% 4.5%	97.3% 2.7%	95.4% 4.6%	94.9% 5.1%	94.8% 5.2%	96.3% 3.7%
	0	1	2	3	4	5	6	7	8	9	10

Output Class

Target Class

(b) Confusion matrix

Figure 3: ROC curve and Confusion matrix for BPNN

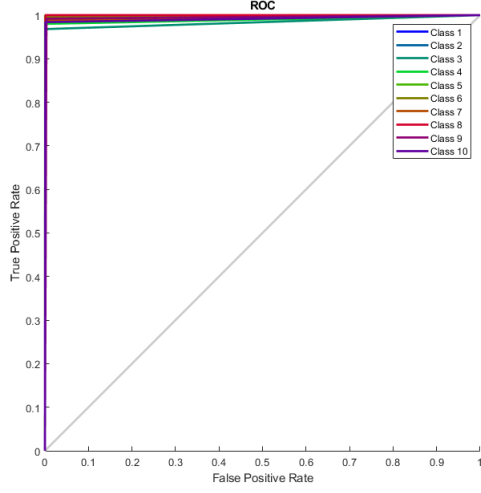
should be evaluated.

6 Conclusion

In this assignment we looked at the performance of SVM networks for use in pattern recognition with the MNIST digit database and compared the previously attained results from BPNN and CNN networks. The two different kernel type SVMs used in this assignment were able to provide comparable performance to the slightly better CNN while training in a much faster time and with much less model complexity.

7 References

Neural Networks and Learning Machines
ECE656 Class Notes

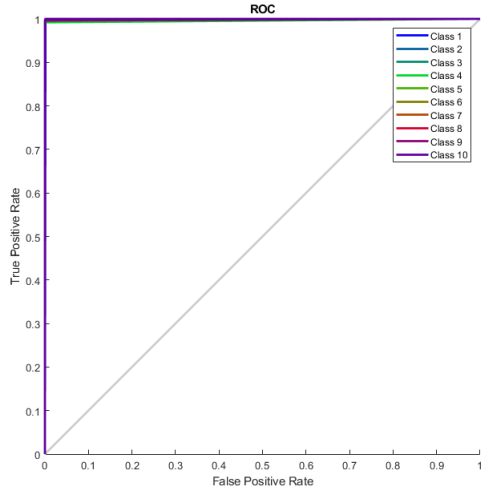


(a) ROC curve for 2 layer

Confusion Matrix

Output Class \ Target Class	0	1	2	3	4	5	6	7	8	9	Accuracy
0	250 10.0%	0 0.0%	3 0.1%	1 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	3 0.1%	3 0.1%	95.8% 4.2%
1	0 0.0%	250 10.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	2 0.1%	1 0.0%	0 0.0%	0 0.0%	98.8% 1.2%
2	0 0.0%	0 0.0%	246 9.8%	2 0.1%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	2 0.1%	0 0.0%	98.4% 1.6%
3	0 0.0%	0 0.0%	0 0.0%	240 9.6%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
4	0 0.0%	0 0.0%	0 0.0%	0 0.0%	249 10.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	99.6% 0.4%
5	0 0.0%	0 0.0%	0 0.0%	2 0.1%	0 0.0%	250 10.0%	0 0.0%	0 0.0%	1 0.0%	1 0.0%	98.4% 1.6%
6	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	247 9.9%	0 0.0%	0 0.0%	0 0.0%	99.6% 0.4%
7	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	248 9.9%	0 0.0%	4 0.2%	98.0% 2.0%
8	0 0.0%	0 0.0%	0 0.0%	4 0.2%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	244 9.8%	1 0.0%	98.0% 2.0%
9	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	241 9.6%	99.6% 0.4%
Accuracy	100%	100%	98.4%	96.0%	99.6%	100%	98.8%	99.2%	97.6%	96.4%	98.6% 1.4%

(b) Confusion matrix for 2 layer



(c) ROC curve for 3 layer

Confusion Matrix

Output Class \ Target Class	0	1	2	3	4	5	6	7	8	9	Accuracy
0	249 10.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
1	0 0.0%	249 10.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	99.6% 0.4%
2	0 0.0%	0 0.0%	249 10.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	99.6% 0.4%
3	0 0.0%	0 0.0%	0 0.0%	249 10.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	100% 0.0%
4	0 0.0%	0 0.0%	0 0.0%	0 0.0%	250 10.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	99.6% 0.4%
5	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	250 10.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	99.6% 0.4%
6	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	250 10.0%	0 0.0%	0 0.0%	0 0.0%	99.6% 0.4%
7	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	247 9.9%	0 0.0%	1 0.0%	99.2% 0.8%
8	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	250 10.0%	0 0.0%	100% 0.0%
9	0 0.0%	0 0.0%	1 0.0%	0 0.0%	0 0.0%	0 0.0%	0 0.0%	1 0.0%	0 0.0%	248 9.9%	99.2% 0.8%
Accuracy	99.6% 0.4%	99.6% 0.4%	99.6% 0.4%	99.6% 0.4%	100% 0.0%	100% 0.0%	100% 0.0%	98.8% 1.2%	100% 0.0%	99.2% 0.8%	99.6% 0.4%

(d) Confusion matrix for 3 layer

Figure 4: ROC curve and Confusion matrix for CNNs