# PR: SVM

## November 17, 2022



**S.M. Erfan Moosavi M.**

- ***Parsing Data ([0_parser.ipynb](0_parser.ipynb))***

For starter, I wrote a simple parser to parse input data.

This parser reads binary files which was downloaded from the provided link ([http://yann.lecun.com/exdb/mnist/](http://yann.lecun.com/exdb/mnist/)) and then convert them to CSV datasets. In this way we can use NumPy built-in functions to read data.

So, let's talk about parser. When you open up the parser notebook, first you should set up some parameters, these parameters are the path to binary datasets and save folder for CSV dataset.

There are 4 cells related to data parsing and 2 cells for checking parsed test and train images.

You can run all the cells altogether or you can run them one by one. In testing cells, you can pass the index to image and see it's visualization along with its label.

This code could have been integrated with other codes but I preferred to make the code modular so I would have more control over each module.

In this module, we are going to preprocess the CSV data that was generated by previous module in two steps.

In the first step, we're going to normalize the data. For this task we're going to use min-max normalization. This normalization applies the following formula on each sample:

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$ from Wikipedia

The normalization is usually done when the dataset consists of features with different ranges, but for our case, all our features are pixel values between 0 to 255. But empirically I found out that normalized features tend to be faster when it comes to model training and it doesn't affect the accuracy of the models so much. So, I decided to do the normalization.

I used built-in MinMaxScalar of sklearn library.

The second preprocessing step is to run PCA on dataset. Before this step, I concatenated the test and train data so they could be feature extracted together. Basically, if we do PCA over test and train separately we usually end up with different features, even if the number of features be the same. So, we can do two things, first is to use the extracted matrix from PCA operations to extract new features from new samples, the second thing is to concatenate test and train data and do one PCA over them, I did the latter and I kept a Boolean mask vector which keeps track of test and train samples.

We use PCA with two feature counts. Keeping 30 and 100 features.

So, we end up with 3 different set of data. First is the data that is not fed to PCA and have all 784 features, but it got scaled. Second one is the set of data that is scaled and fed to PCA30 and the third one is the

one that is scaled and fed to PCA100. We shall call them Full set, PCA30 and PCA100.

Then we stored these sets along with output set (let's call it y, and this output is a concatenation of test and train output but with the same order of inputs). And we also saved the mask vector.

3

Introduction of questions sheet asks us to train three different models corresponds to three different sets. In this section and next two sections we are talking about evaluations that are done over every one of these sets.

The first test is done over full dataset. We loaded the preprocessed datasets and then using the mask vector, separated the test and train sets. After that, using LinearSVC class of sklearn, we trained a multiclass SVM classifier using One VS Rest technique. To solve the SVM optimization there are two options, first is to solve the dual problem and the second option is to solve the primal problem. We choose to solve the primal because empirically the dual option causes the algorithm to not converge. Also, it is not suggested by sklearn documentation since $n\_samples > n\_features$. Then, using the trained SVM classifier, we predicted labels for our train and test samples. Here are the results:

Model: Linear SVM (LinearSVM class from sklearn)

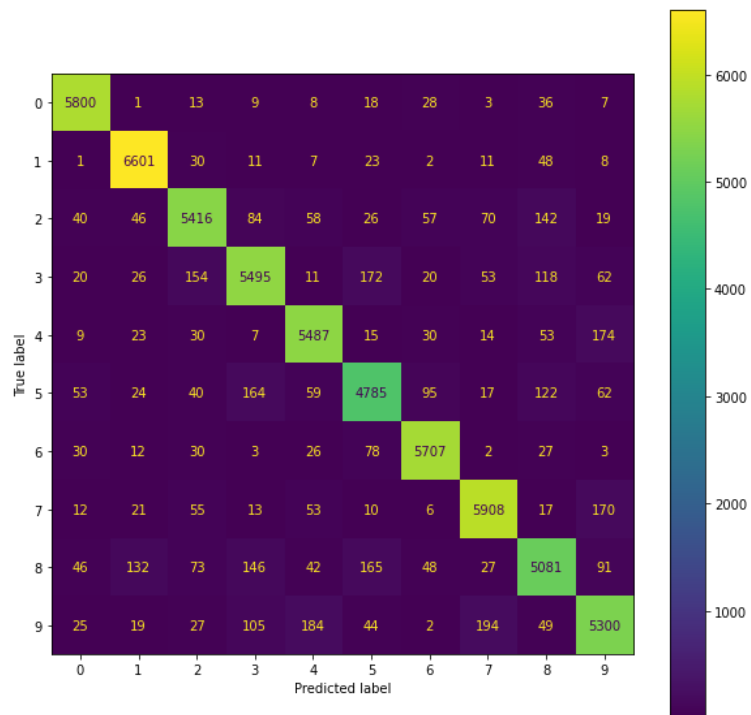Train Accuracy: 92.63%

Test Accuracy: 91.81%

Train Precision and Recall:

| CLASSES | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| **PRECISION (TRAIN)** | 0.96 | 0.95 | 0.92 | 0.91 | 0.92 | 0.89 | 0.95 | 0.93 | 0.89 | 0.89 |
| **RECALL (TRAIN)** | 0.97 | 0.97 | 0.90 | 0.89 | 0.93 | 0.88 | 0.96 | 0.94 | 0.86 | 0.89 |

Test Precision and Recall:

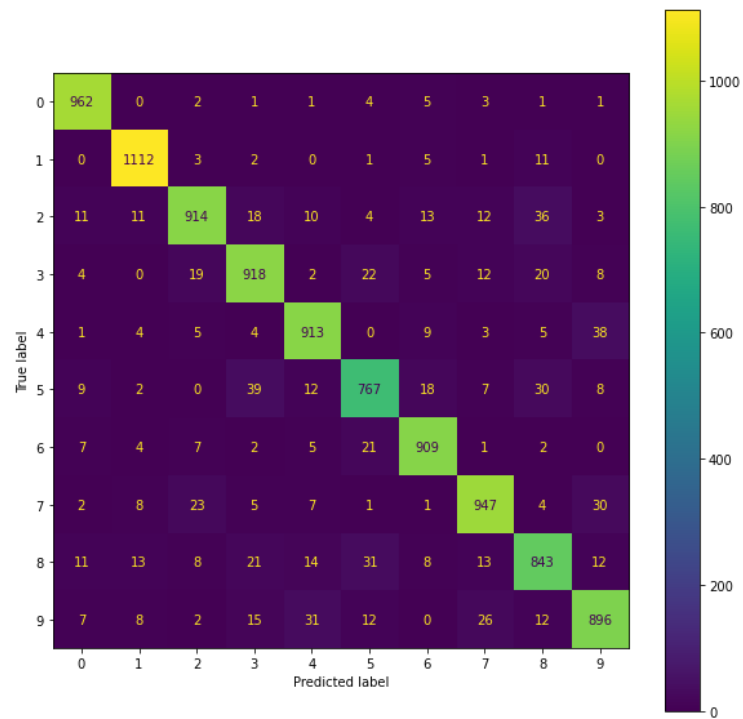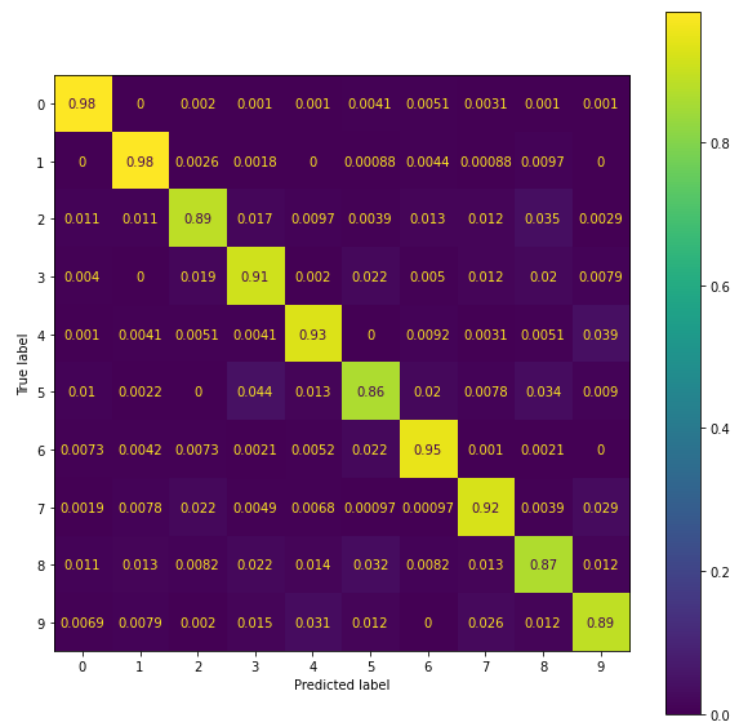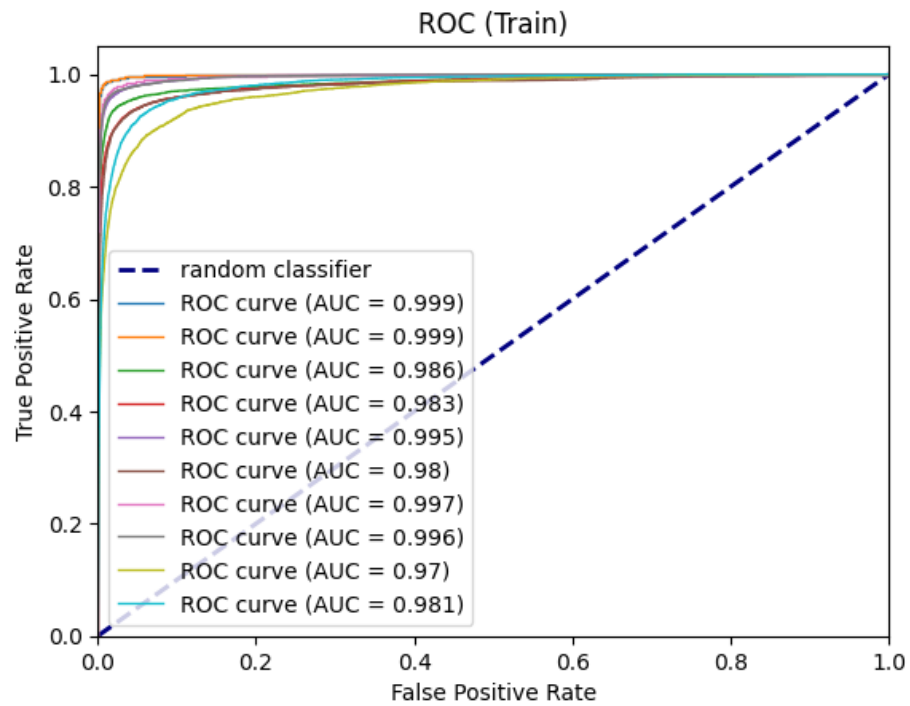| CLASSES | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| **PRECISION (TEST)** | 0.94 | 0.95 | 0.92 | 0.89 | 0.91 | 0.88 | 0.93 | 0.92 | 0.87 | 0.89 |
| **RECALL (TEST)** | 0.98 | 0.97 | 0.88 | 0.90 | 0.92 | 0.85 | 0.94 | 0.92 | 0.86 | 0.88 |

Train Confusion Matrix:


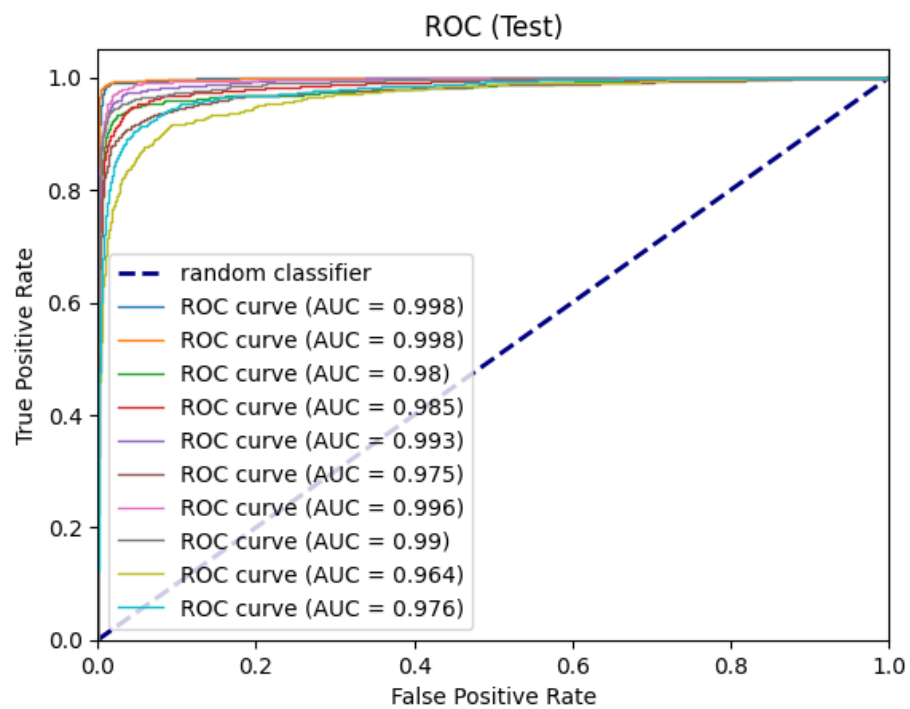
Normalized Train Confusion Matrix:

Test Confusion Matrix:



Normalized Test Confusion Matrix:

Train ROC Plot & AUC Values:



Train ROC Plot & AUC Values:

The second test is done over PCA30 dataset. We loaded the preprocessed datasets (Scaled and PCA30) and then using the mask vector, separated the test and train sets. After that, using LinearSVC class of sklearn, we trained a multiclass SVM classifier using One VS Rest technique. To solve the SVM optimization there are two options, first is to solve the dual problem and the second option is to solve the primal problem. We choose to solve the primal because empirically the dual option causes the algorithm to not converge. Also, it is not suggested by sklearn documentation since n_samples > n_features. Then, using the trained SVM classifier, we predicted labels for our train and test samples. Here are the results:

Model: Linear SVM (LinearSVM class from sklearn)

Train Accuracy: 87.66%

Test Accuracy: 88.38% (seems the test set was easier to classify)
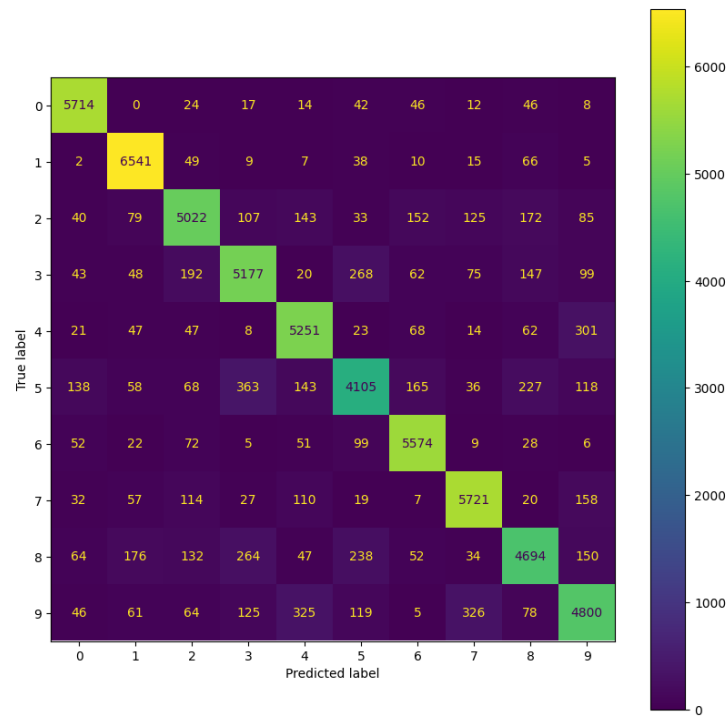
Train Precision and Recall:

| CLASSES | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| PRECISION (TRAIN) | 0.92 | 0.92 | 0.86 | 0.84 | 0.85 | 0.82 | 0.90 | 0.89 | 0.84 | 0.83 |
| RECALL (TRAIN) | 0.96 | 0.97 | 0.84 | 0.84 | 0.89 | 0.75 | 0.94 | 0.91 | 0.80 | 0.80 |

Test Precision and Recall:

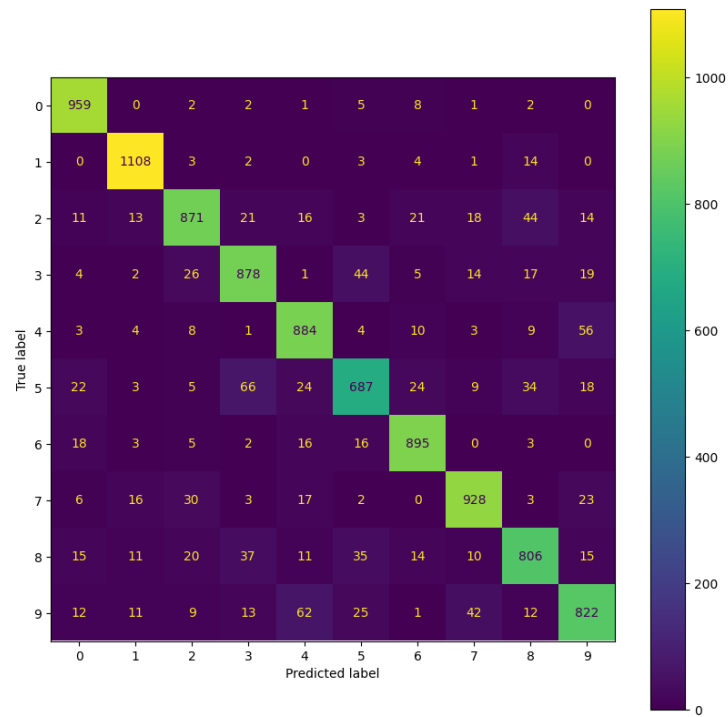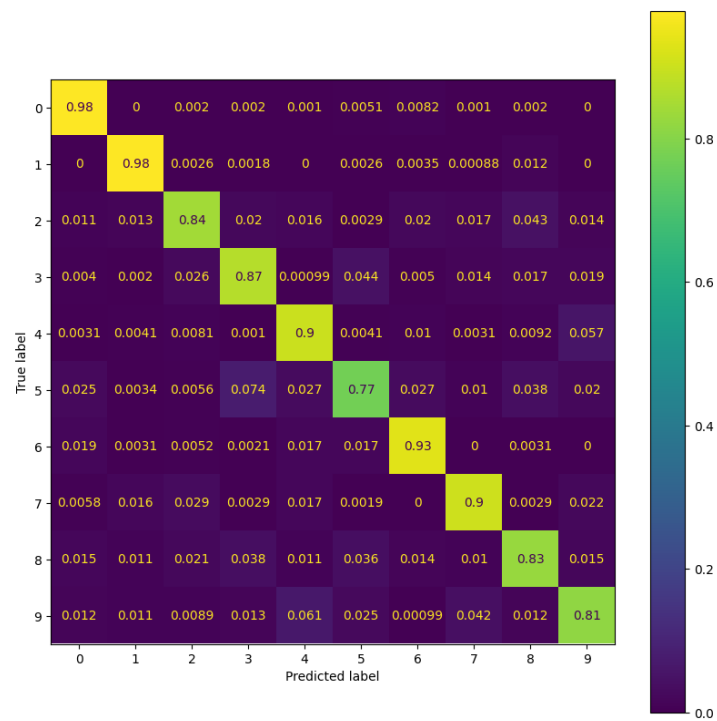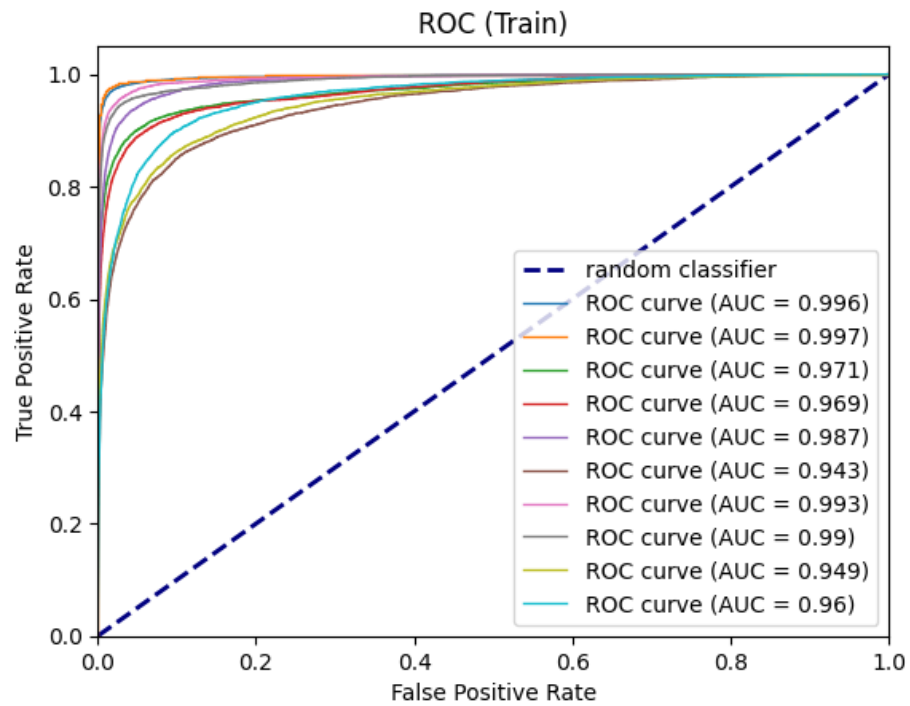| CLASSES | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| PRECISION (TEST) | 0.91 | 0.94 | 0.88 | 0.85 | 0.85 | 0.83 | 0.91 | 0.90 | 0.85 | 0.85 |
| RECALL (TEST) | 0.97 | 0.97 | 0.84 | 0.86 | 0.90 | 0.77 | 0.93 | 0.90 | 0.82 | 0.81 |

Train Confusion Matrix:



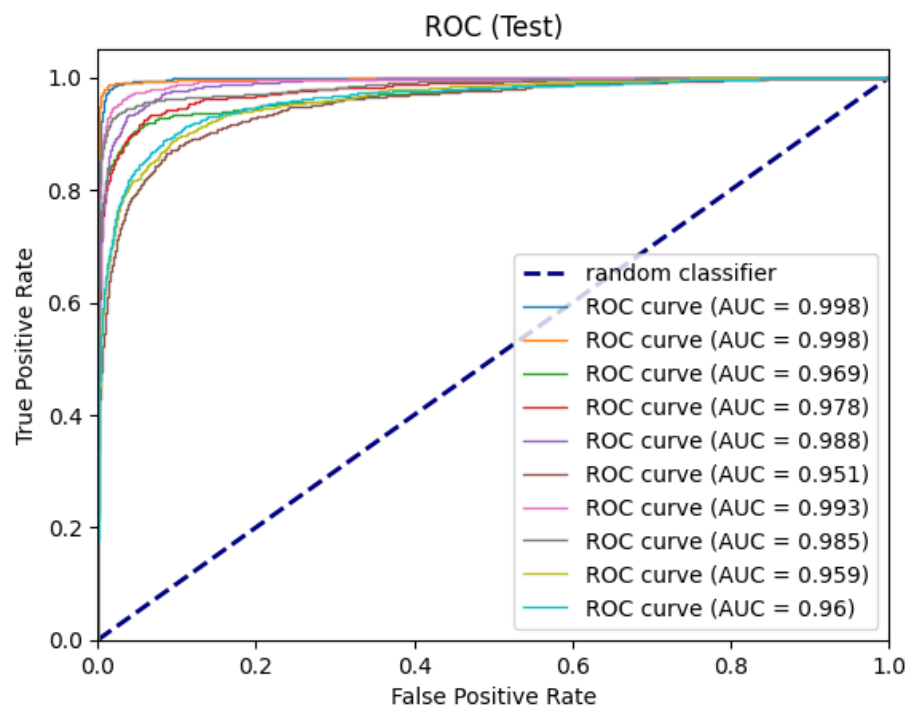Normalized Train Confusion Matrix:

Test Confusion Matrix:



Normalized Test Confusion Matrix:

Train ROC Plot & AUC Values:



Train ROC Plot & AUC Values:

The third test is done over PCA100 dataset. We loaded the preprocessed datasets (Scaled and PCA100) and then using the mask vector, separated the test and train sets. After that, using LinearSVC class of sklearn, we trained a multiclass SVM classifier using One VS Rest technique. To solve the SVM optimization there are two options, first is to solve the dual problem and the second option is to solve the primal problem. We choose to solve the primal because empirically the dual option causes the algorithm to not converge. Also, it is not suggested by sklearn documentation since $n\_samples > n\_features$. Then, using the trained SVM classifier, we predicted labels for our train and test samples. Here are the results:

Model: Linear SVM (LinearSVM class from sklearn)

Train Accuracy: 90.84%

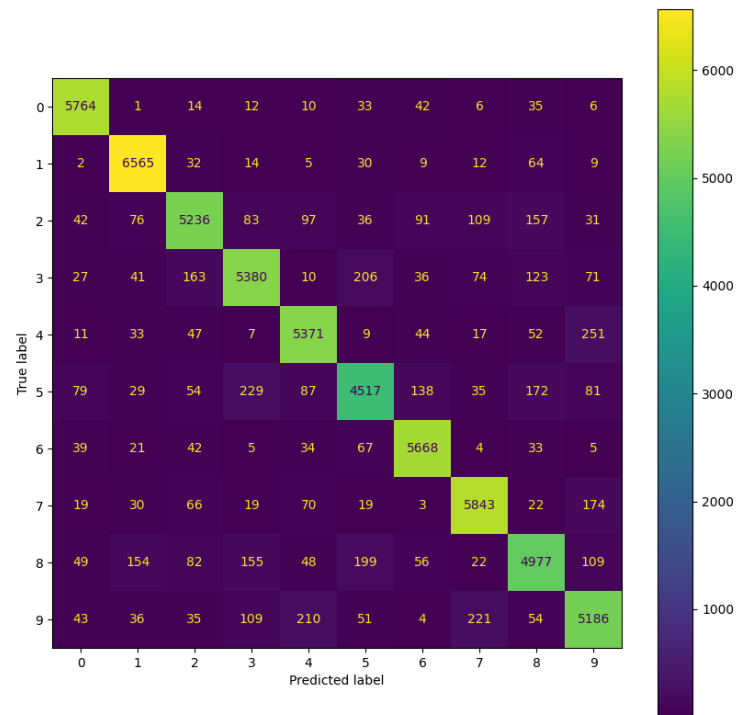Test Accuracy: 91.3% (again, it seems the test set was easier to classify)

Train Precision and Recall:

| CLASSES | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| **PRECISION (TRAIN)** | 0.94 | 0.93 | 0.90 | 0.89 | 0.90 | 0.87 | 0.93 | 0.92 | 0.87 | 0.87 |
| **RECALL (TRAIN)** | 0.97 | 0.97 | 0.87 | 0.87 | 0.91 | 0.83 | 0.95 | 0.93 | 0.85 | 0.87 |

Test Precision and Recall:

| CLASSES | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|---|
| **PRECISION (TEST)** | 0.94 | 0.95 | 0.92 | 0.90 | 0.90 | 0.87 | 0.92 | 0.91 | 0.87 | 0.88 |
| **RECALL (TEST)** | 0.97 | 0.97 | 0.88 | 0.90 | 0.91 | 0.84 | 0.95 | 0.92 | 0.85 | 0.87 |

Train Confusion Matrix:



Normalized Train Confusion Matrix:

Test Confusion Matrix:



Normalized Test Confusion Matrix:

Train ROC Plot & AUC Values:



ROC (Train)

True Positive Rate / False Positive Rate
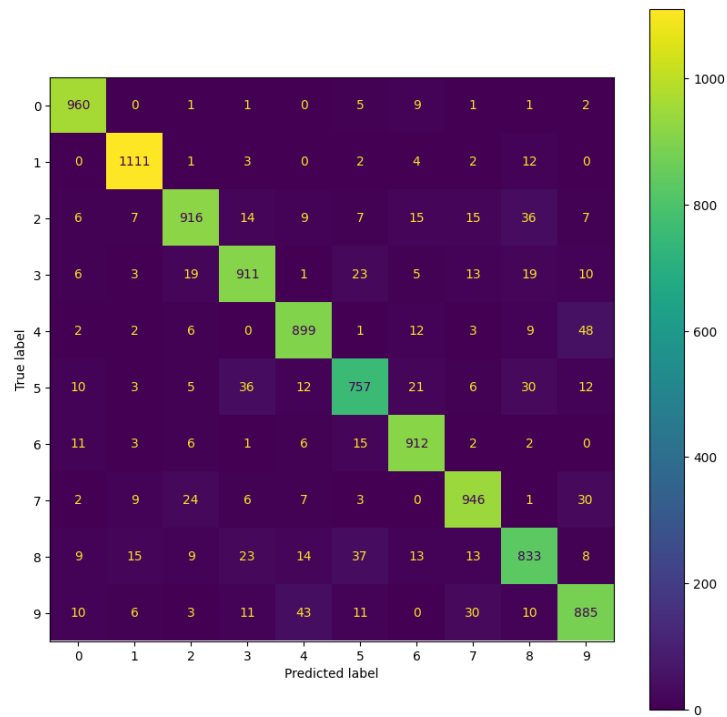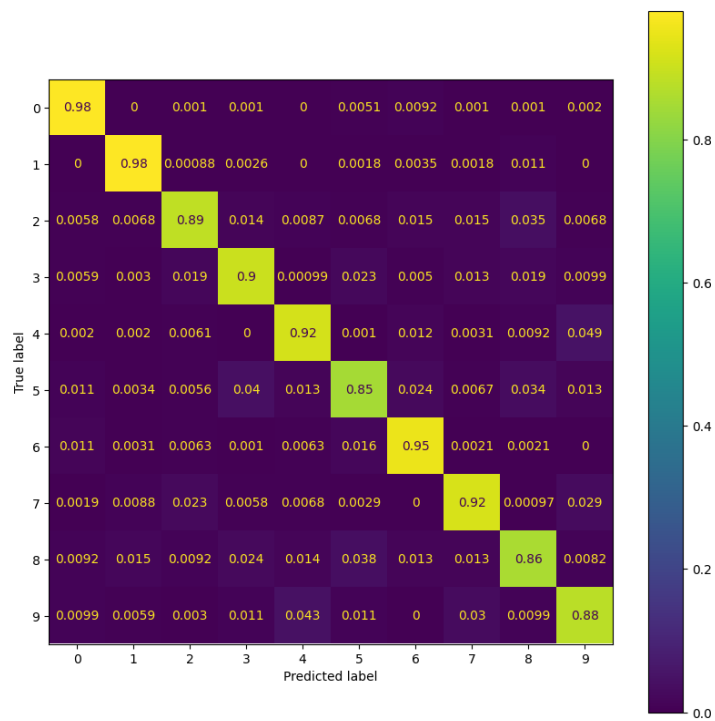
- random classifier
- ROC curve (AUC = 0.997)
- ROC curve (AUC = 0.998)
- ROC curve (AUC = 0.979)
- ROC curve (AUC = 0.978)
- ROC curve (AUC = 0.992)
- ROC curve (AUC = 0.965)
- ROC curve (AUC = 0.995)
- ROC curve (AUC = 0.993)
- ROC curve (AUC = 0.963)
- ROC curve (AUC = 0.976)

Train ROC Plot & AUC Values:



ROC (Test)

True Positive Rate / False Positive Rate

- random classifier
- ROC curve (AUC = 0.998)
- ROC curve (AUC = 0.999)
- ROC curve (AUC = 0.978)
- ROC curve (AUC = 0.985)
- ROC curve (AUC = 0.992)
- ROC curve (AUC = 0.969)
- ROC curve (AUC = 0.996)
- ROC curve (AUC = 0.991)
- ROC curve (AUC = 0.968)
- ROC curve (AUC = 0.973)

**15**

o   ***Analysis and Comparison of Datasets***

The first thing that caught my eyes was that the test accuracy over PCA30 & PCA100 sets are bigger than train accuracy. But this is not the case for Full set.

The next interesting thing is that the classification quality of one and zero is higher than others (Higher precision and recall). Which may happen because of visually distinguishable features of these two digits which make their separation from others easier. Digit 5 has the lowest precision, and from its confusions we can see it mostly mistaken as 8. Which is expected because of these two digits similarities in writing.

We can observe that all the classes got classified with good qualities. But this is no surprise since MNIST is a toy dataset 😊.

o   ***Cost Parameter Effect ([5_models_pca30.ipynb](#))***

First thing first, we loaded the preprocessed dataset, then we separated the test and train data.

By defining a list of C parameters, we can iterate over the list and train a different SVM with each of those Cs in every iteration.

Here, we used the LinearSVC class from sklearn library. We trained this SVM model in an OneVsRest mode, which means it's going to train a single SVM for each class. For example, in case of MNIST dataset, it's going to train 10 different SVM models, one for every class. Every SVM model trained in this mode should predict each corresponding class vs every other class. (It should differentiate between that class and others).

There are two ways that we can incorporate when it comes to training our SVM model. Primal and Dual. The documentation of sklearn

provided a simple criterion which is going to help to decide which method to use. The criteria are as follows:

If the number of samples are bigger than the number of features, you should use primal method, other wise use dual method.
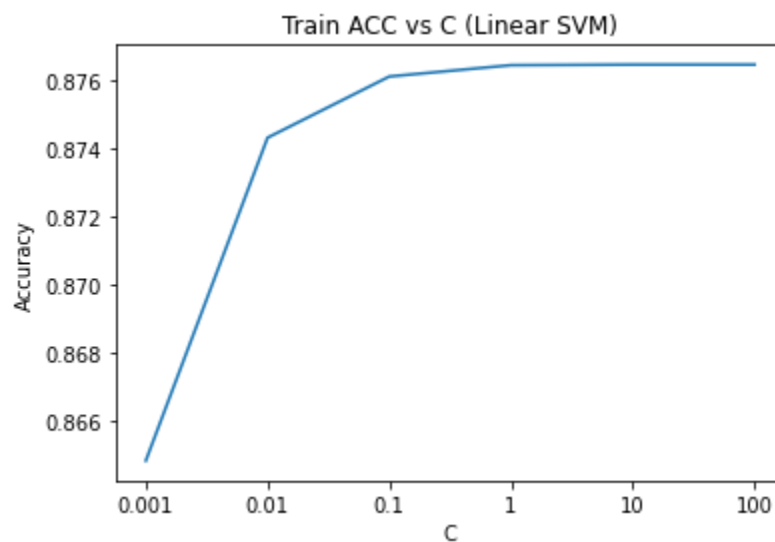
In our case since we know the number of samples are bigger than number of features, we turn dual option off.
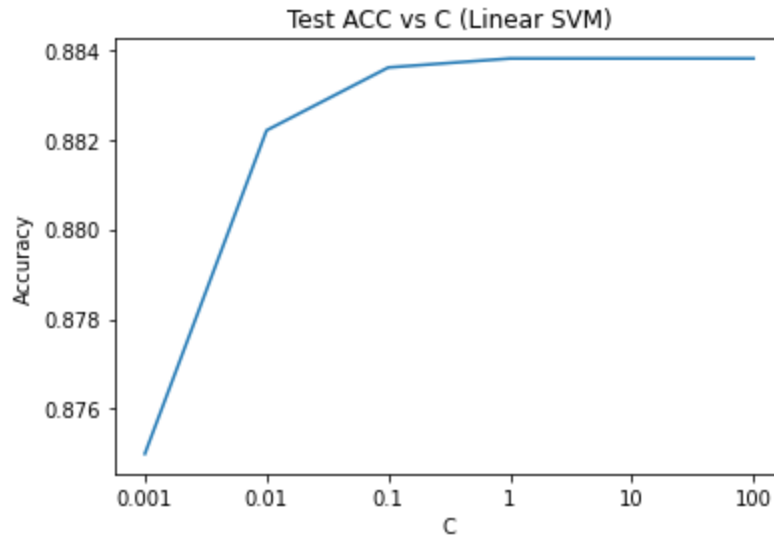
After training each model, the training and test accuracy calculated and stored to be compared with others.

So, here's the result of this experiment:

|  | 0.001 | 0.010 | 0.100 | 1.000 | 10.000 | 100.000 |
|---|---|---|---|---|---|---|
| **TRAIN ACC** | 0.864833 | 0.874333 | 0.876133 | 0.876467 | 0.876483 | 0.876483 |
| **TEST ACC** | 0.875 | 0.8822 | 0.8836 | 0.8838 | 0.8838 | 0.8838 |

We can see that by increasing the C parameter, the accuracy of train and test got increased too. To see this effect more clearly, we plotted the train and test accuracy against the C parameter separately:



Train ACC vs C (Linear SVM)

Test ACC vs C (Linear SVM)

We can observe from these plots that from one point (C=0.1) the effect of C parameters is negligible, hence it increased the training time a lot.

   o **Gamma Parameter Effect**

Again, for this assignment, we loaded the preprocessed training and test data, then we trained 3 different SVM models with RBF kernels and each time we used a different gamma value.
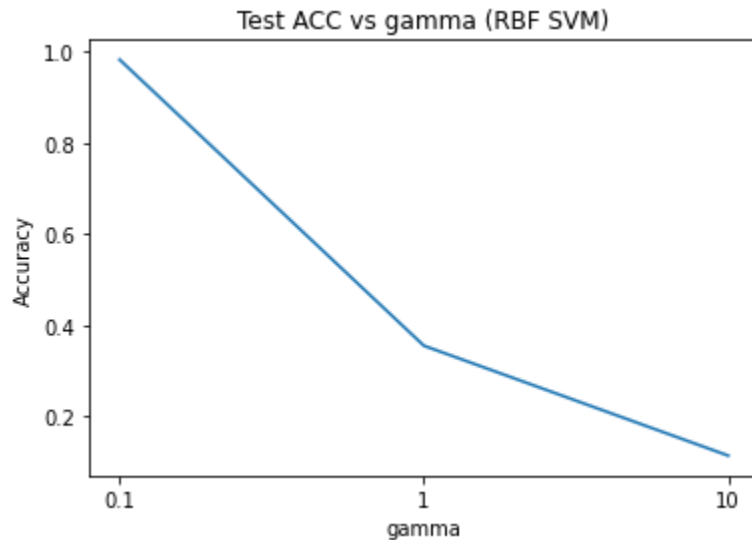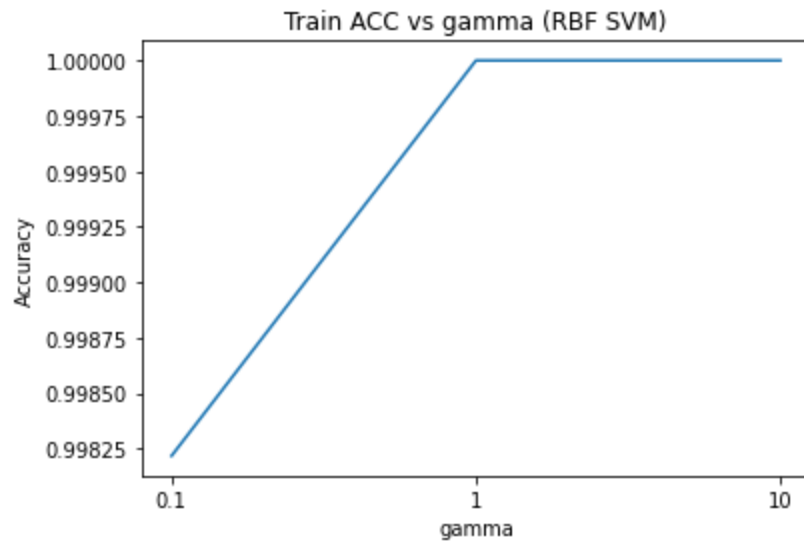
All the models, used the same C=1 value and OneVsRest method for handling the multiclass problem.

Here's the result of this experiment:

|  | 0.1 | 1.0 | 10.0 |
| --- | --- | --- | --- |
| **TRAIN ACC** | 0.998217 | 1.0 | 1.0 |
| **TEST ACC** | 0.9833 | 0.3549 | 0.1135 |

Here's we can observe a direct relation between gamma parameter and overfitting. Which means by increasing the gamma value, the accuracy over training data is maximized but the test accuracy reduced.

Here's some plot for accuracy against gamma value:

**18**

Train ACC vs gamma (RBF SVM)



Test ACC vs gamma (RBF SVM)

- o *Crammer & Singer method*
  *(6_crammer_singer_pca100.ipynb)*

Traditional SVM is one of the most used and most famous among classifiers but the problem with this classifier is it's limitations when it comes to multiclass problems. Traditional SVM can only separate two classes which means it's suitable for binary classification problems.

People usually use some tricks to get SVM works on multiclass problems, for example they are going to train different independent SVM models for each class. The problem with this method is its lack of ability when it comes to capturing inter-class correlations.

To address this issue, crammer and singer proposed a new method of multiclass SVM in which all classes are trained together.

Basically, what Crammer and Singer proposed is to transform the problem to an optimization problem and now the task is to solve this optimization problem.

This is the optimization problem:

$$
\begin{aligned}
\text{maximize} \quad & -(1/2)\mathbf{Tr}(U^T Q U) + \mathbf{Tr}(E^T U) \\
\text{subject to} \quad & U \preceq \gamma E \\
& U\mathbf{1}_m = 0
\end{aligned}
$$

There's an option to use C&S method in LinearSVC of sklearn library.

Here's the result of the algorithm:

Train Accuracy: 92.56%

Test Accuracy: 92.6%

Head to the notebook to see Confusion Matrix and ROC, AUC.

- ○ ***Crammer and Singer Implementation
  ([7_c&s_svm_final.ipynb](7_c&s_svm_final.ipynb))***

For this part of the assignment, we implemented and corrected an implementation of Crammer and Singer method for training Multiclass SVM.

Test Accuracy: 92.642