Introduction:

Classifying hand written letters is both necessary and important to the study of machine learning. This assignment focus on identifying a pool of hand written letters from the MNIST dataset using the Nearest Neighbor Algorithm.
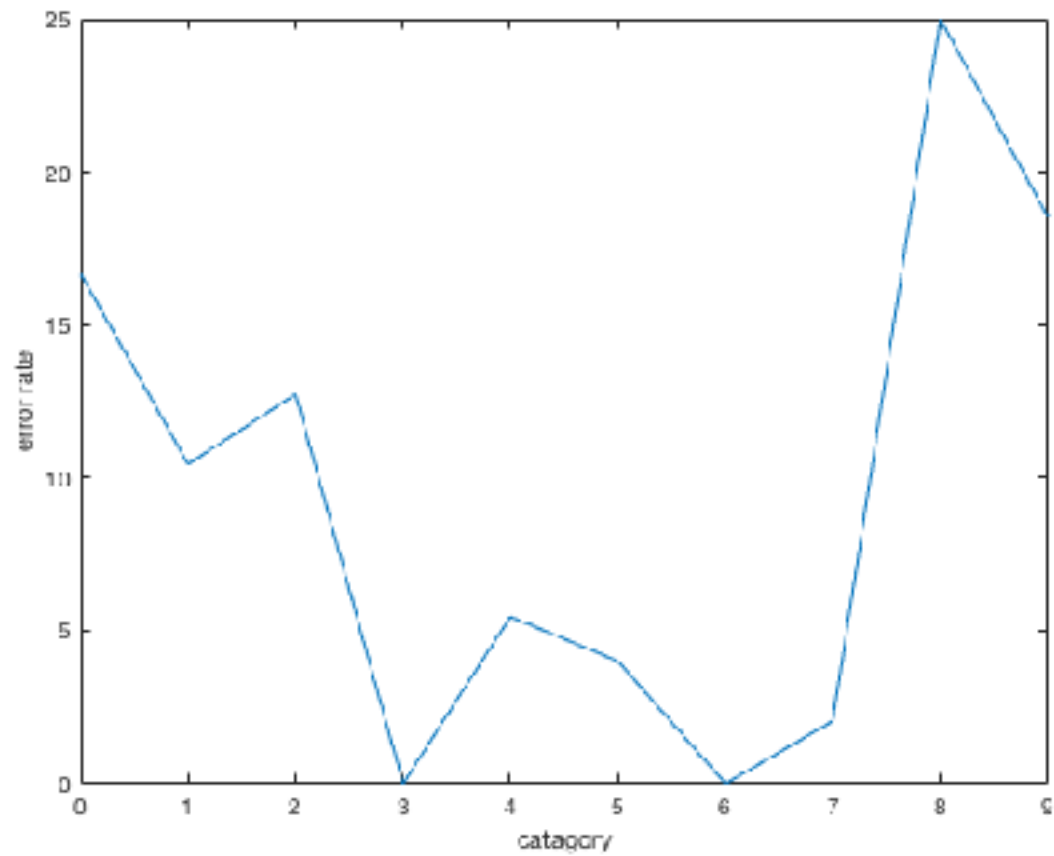
Procedures:

*sidenote: I initially attempted a faster method by averaging all the image with the same truth value into a "trained" model for that specific truth. Then I ran all test images to compare with these 10 "trained" images and yielded about 22% error rate in total. Though this kind of worked, the error rate is too high to be acceptable. I gave up on the attempt to optimize for computing speed and went for the 1-test-compare-all approach, comparing each test image to the entire training set to find a match.

The algorithm simply performed a euclidian distance calculation per pixel between the test and the reference. Each teat image will yield a 5000 element array with distance filled out. We then choose the minimal distance from this list of 5000 elements, then look up the truth of that match in the label_train dataset.

Results:

The algorithm performed as excepted. It made handful of mistakes, which can be characterized in the following manner:

Number 3 and 6 were classified perfectly. While we saw pretty terrible results between 8 and 9. 0 and 2 were also up in the error ally.
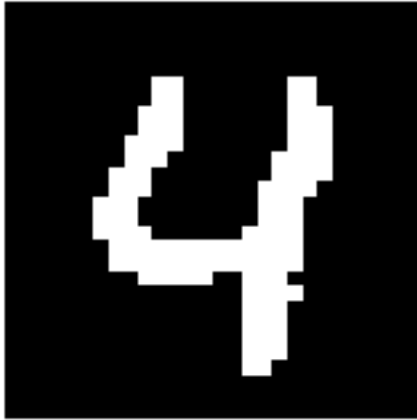
The error rate in the 0-9 order is: (%)

16.6667   10.4478   12.7273      0   5.4545   4.0000      0   2.0408   25.0000   18.5185

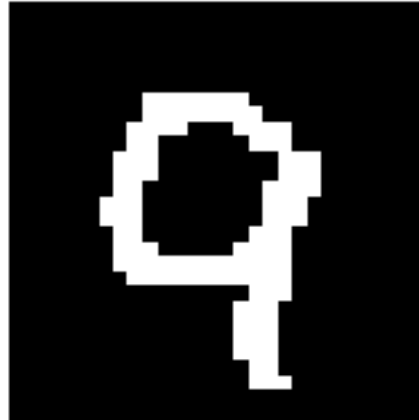The total error rate of the classifier performed at 9.4%, getting 47 images wrong out of 500.

Below is a sample of 5 misclassified images with comments on why it failed embedded in the images.

**Test Image**



the pair matched with dist 339.597899

**Matched Image**



the 4 matches the shape of the 9 rather, the gap, could be confusing for human even

**Test Image**



the pair matched with dist 466.441506

**Matched Image**



another 4-9 mess up, same reason, the general shape is very close

**Test Image**



the pair matched with dist 165.256585

**Matched Image**



yet another one, this one is very bad to human even, just bad writing

| Test Image | Matched Image |
|---|---|



the pair matched with dud 169.762962

Old Guide is the dests of this method, same reason

| Test Image | Matched Image |
|---|---|



the pair matched with dud 606.802485

the first makes of the 2 almost matched 1 perfectly, even the machine to 2nd makes is not enough to

Conclusion:

The algorithm performed as excepted. Due to the simplistic nature of this method, a large number of errors were captured. Further implementations of algorithm improvement can help the classifier perform better.

Appendix, Practice Problems: