

Homework 4
MNIST Handwritten Digits Recognition
Phys 243
06/01/19

Faddy Sunna

DATA: MNIST Handwritten digits. All 60,000 images were used for training, and 10,000 images for testing. No preprocessing was needed.

At first it was very time consuming to run the whole dataset because sklearn was only using 1 cpu core. I realized that adding, `n_jobs = -1`, allows it to take advantage of all CPU cores, greatly speeding up the runtime.

1. Finding all the 9s, using KNN Sklearn.

First, I found the optimal K value by testing values (1, 3, 5, 7) and measuring the misclassification error. K = 3 had the lowest misclassification error, and thus was used for the rest of the KNN models. Fig 1: Code used to find the optimal K, by iterating k values and measuring misclassification error. I was able to switch from classifying just the 9s and all the digits by adjusting line 42 in figure 3.

Setting Minkowski P=1,2,3 parameters: Fig 3, line 39.

Table 1 and 2 show the final results of the KNN classifier.

FIND NINE	Minkowski	KNN
Mean Accuracy	P	K neighbors
94.50%	1	3
95.90%	2	3
96.80%	3	3

Table 1: Classifying just the 9s. Using KNN Minkowski 1,2,3.

FIND ALL	Minkowski	KNN
Mean Accuracy	P	K neighbors
96.30%	1	3
97.00%	2	3
97.18%	3	3

Table 2: Classifying all digits 0-9. Using KNN Minkowski 1,2,3.

```

9 import struct
10 import numpy as np
11 import matplotlib.pyplot as plt
12 from sklearn.model_selection import train_test_split, cross_val_score
13 from sklearn.neighbors import KNeighborsClassifier
14
15 def read_idx(filename):
16     with open(filename, 'rb') as f:
17         zero, data_type, dims = struct.unpack('>HBB', f.read(4))
18         shape = tuple(struct.unpack('>I', f.read(4))[0] for d in range(dims))
19         return np.fromstring(f.read(), dtype=np.uint8).reshape(shape)
20
21 raw_train = read_idx("train-images-idx3-ubyte")
22 train_data = np.reshape(raw_train, (60000, 28*28))
23 train_label = read_idx("train-labels-idx1-ubyte")
24
25 raw_test = read_idx("t10k-images-idx3-ubyte")
26 test_data = np.reshape(raw_test, (10000, 28*28))
27 test_label = read_idx("t10k-labels-idx1-ubyte")
28
29 idx = (train_label == 0) | (train_label == 1) | (train_label == 2) | (train_label == 3) | (train_label == 4) | (train_label == 5) | (train_label == 6) | (train_label == 7) | (train_label == 8) | (train_label == 9)
30 x_train = train_data[idx]
31 y_train = train_label[idx]
32
33 idx = (test_label == 0) | (test_label == 1) | (test_label == 2) | (test_label == 3) | (test_label == 4) | (test_label == 5) | (test_label == 6) | (test_label == 7) | (test_label == 8) | (test_label == 9)
34 x_test = test_data[idx]
35 y_true = test_label[idx]
36
37 # creating odd list of K for KNN
38 myList = list(range(1,7))
39
40 # subseting just the odd ones
41 neighbors = list(filter(lambda x: x % 2 != 0, myList))
42
43 # empty list that will hold cv scores
44 cv_scores = []
45
46 #15-fold cross validation comparing accuracy of k values
47
48 for k in neighbors:
49     clf = KNeighborsClassifier(n_neighbors=k, n_jobs=-1)
50     scores = cross_val_score(clf, x_test, y_true, cv=3, scoring='accuracy', n_jobs=-1)
51     cv_scores.append(scores.mean())
52
53 # misclassification error
54 MSE = [1 - x for x in cv_scores]
55
56 # determining best k
57 optimal_k = neighbors[MSE.index(min(MSE))]
58 print("The optimal number of neighbors is %d" % optimal_k)
59
60 # plot misclassification error vs k
61 plt.plot(neighbors, MSE)
62 plt.xlabel('Number of Neighbors K')
63 plt.ylabel('Misclassification Error')
64 plt.show()
65
66 read_idx()

```

Fig 1 : Finding Optimal k , by iterating k=1,3,5,7, and measuring misclassification error

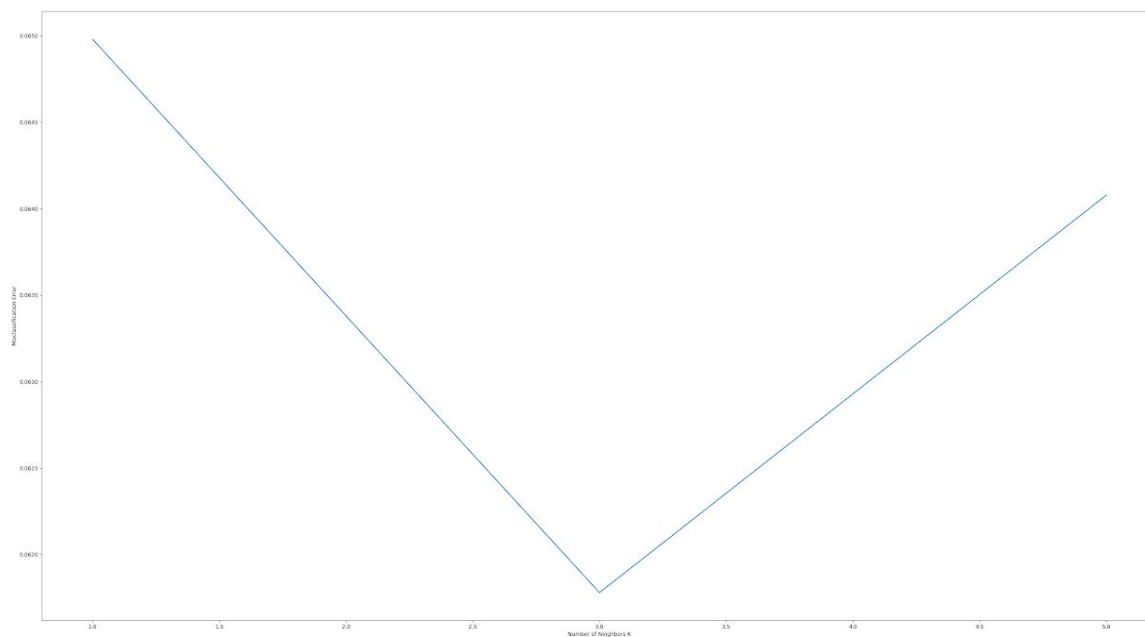


Fig 2: X-axis is k neighbors, Y-axis is misclassification error.

```

1 import struct
2 import numpy as np
3 from sklearn import neighbors, metrics
4 import matplotlib.pyplot as plt
5 from sklearn import preprocessing, cross_decomposition, neighbors
6 from sklearn.neighbors import KNeighborsClassifier
7
8
9 #This code converts the MNIST idx format into numpy arrays
10 """ A function that can read MNIST_s idx file format into numpy arrays.
11 The MNIST data files can be downloaded from here:
12
13 http://yann.lecun.com/exdb/mnist/
14 This relies on the fact that the MNIST dataset consistently uses
15 unsigned char types with their data segments.
16 """
17 def read_idx(filename):
18     with open(filename, 'rb') as f:
19         zero, data_type, dims = struct.unpack('>HBB', f.read(4))
20         shape = tuple(struct.unpack('>I', f.read(4))[0] for d in range(dims))
21         return np.fromstring(f.read(), dtype=np.uint8).reshape(shape)
22
23
24 #import the train and test data and labels. Using all 60,000 images for training and 10,000 images for testing
25 raw_train = read_idx("train-images-idx3-ubyte")
26 train_data = np.reshape(raw_train, (60000, 28*28))
27 train_label = read_idx("train-labels-idx1-ubyte")
28
29 raw_test = read_idx("t10k-images-idx3-ubyte")
30 test_data = np.reshape(raw_test, (10000, 28*28))
31 test_label = read_idx("t10k-labels-idx1-ubyte")
32 #selecting all the data, digits 0-9 for training
33 idx = (train_label == 0) | (train_label == 1) | (train_label == 2) | (train_label == 3) | (train_label == 4) | (train_label == 5) | (train_label == 6) | (train_label == 7) | (train_label == 8) | (train_label == 9)
34 x = train_data[idx]
35 y = train_label[idx]
36
37
38 #Creating our knn classifier and setting parameters to minkowski and p=1,2,3, then fitting the data
39 knn = neighbors.KNeighborsClassifier(n_neighbors=3, metric='minkowski', p=3, n_jobs=-1).fit(x, y)
40
41 #selecting our test data, in this case, it will be all the digits
42 idx = (test_label == 0) | (test_label == 1) | (test_label == 2) | (test_label == 3) | (test_label == 4) | (test_label == 5) | (test_label == 6) | (test_label == 7) | (test_label == 8) | (test_label == 9)
43 x_test = test_data[idx]
44 y_true = test_label[idx]
45 y_pred = knn.predict(x_test)
46
47 import itertools
48
49 #Create confusion matrix
50 def plot_confusion_matrix(cm, classes,
51                           normalize=False,
52                           title='Confusion matrix',
53                           cmap=plt.cm.Blues):
54     """
55     This function prints and plots the confusion matrix.
56     Normalization can be applied by setting 'normalize=True'.
57     """
58     if normalize:
59         cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
60         print("Normalized confusion matrix")
61     else:
62         print('Confusion matrix, without normalization')
63
64     print(cm)
65
66     plt.imshow(cm, interpolation='nearest', cmap=cmap)
67     plt.title(title)
68     plt.colorbar()
69     tick_marks = np.arange(len(classes))
70     plt.xticks(tick_marks, classes, rotation=45)
71     plt.yticks(tick_marks, classes)
72
73     fmt = '.2f' if normalize else 'd'
74     thresh = cm.max() / 2.
75     for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
76         plt.text(j, i, format(cm[i, j], fmt),
77                  horizontalalignment="center",
78                  color="white" if cm[i, j] > thresh else "black")
79
80     plt.tight_layout()
81     plt.ylabel('True label')
82     plt.xlabel('Predicted label')
83
84 #Calculate mean accuracy using built in score function
85 score = KNeighborsClassifier.score(knn, x_test, y_true)
86 scorestr = str(score)
87 print("The mean accuracy score using minkowski metric and p=3 is: " + scorestr)
88 #plot confusion matrix non normalized
89 cm = metrics.confusion_matrix(y_true, y_pred)
90 plot_confusion_matrix(cm, ["0", "1", "2", "3", "4", "5", "6", "7", "8", "9"])
91 plt.show()
92
93 #plot normalized confusion matrix
94 plot_confusion_matrix(cm, ["0", "1", "2", "3", "4", "5", "6", "7", "8", "9"], normalize=True)
95 plt.show()

```

Fig 3: KNN
Sklearn classifier

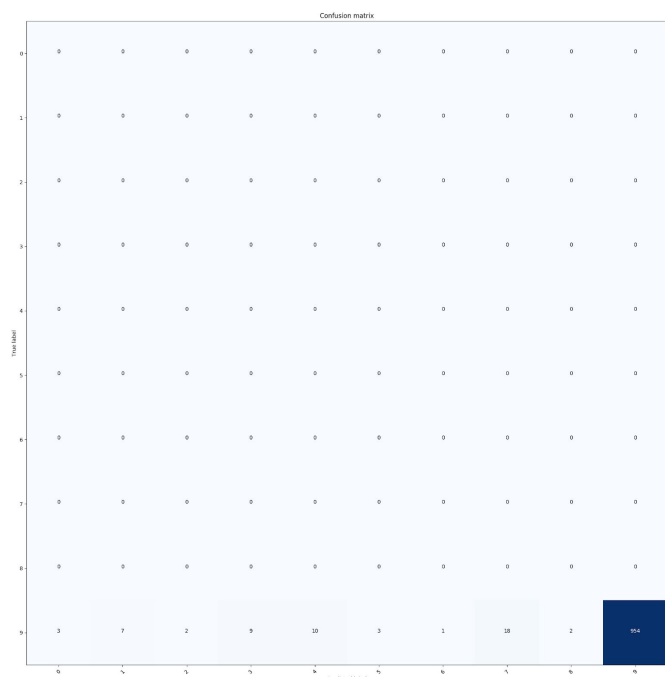


Fig: 4 Confusion matrix, pred 9 , p=1

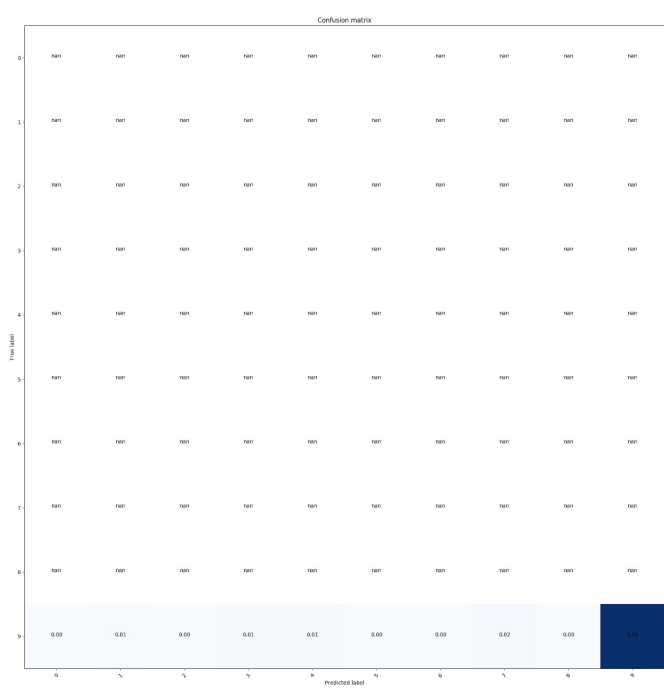


Fig 5: Normalized, classifying 9, P=1

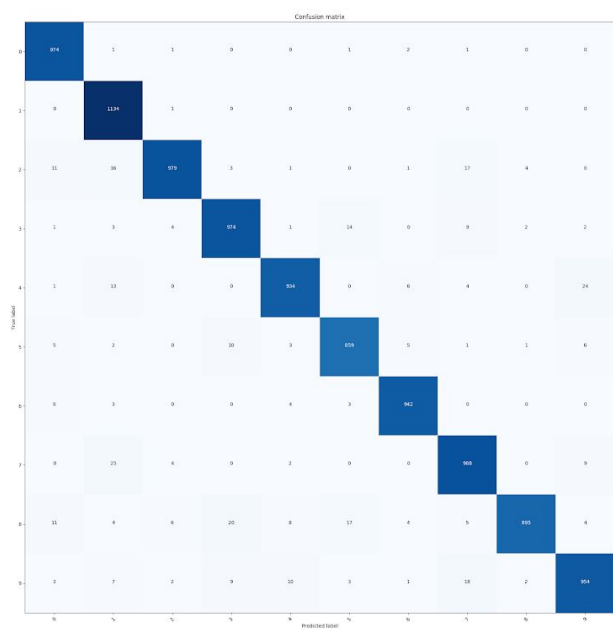


Fig 6: Classifying 0-9, P=1 Minkowski

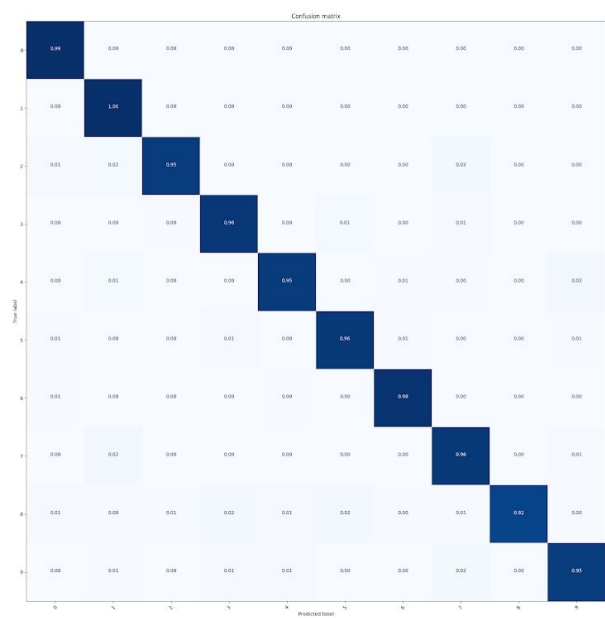


Fig 7: Classifying 0-9, P=1 Normalized

2. Decision Tree: Sklearn was used to run the decision tree classifier. The model was run using entropy and gini criterion for dept =8,16, default. Where the default is to run until the leaf is pure. Mean accuracy was measured using the building sklearn scoring function.

FIND ALL				FIND NINE		
Mean Accuracy	Depth	Criterion		Accuracy	Depth	Criterion
81.80%	8	gini		82.40%	8	gini
88.20%	16	gini		87.00%	16	gini
87.60%	Until pure	gini		84.00%	Until pure	gini
Mean Accuracy	Depth	Criterion		Accuracy	Depth	Criterion
83.60%	8	entropy		80.30%	8	entropy
88.90%	16	entropy		86.20%	16	entropy
88.50%	Until pure	entropy		86.90%	Until pure	entropy

Table 3: Decision tree final results

```

1  from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor
2  import struct
3  import itertools
4  import numpy as np
5  from sklearn import neighbors, metrics
6  import matplotlib.pyplot as plt
7  from sklearn.tree import export_graphviz
8  from sklearn.externals.six import StringIO
9  from IPython.display import Image
10 import pydotplus
11
12
13 #This code converts the MNIST idx format into numpy arrays
14 def read_idx(filename):
15     with open(filename, 'rb') as f:
16         zero, data_type, dims = struct.unpack('>HBB', f.read(4))
17         shape = tuple(struct.unpack('>I', f.read(4))[0] for d in range(dims))
18         return np.fromstring(f.read(), dtype=np.uint8).reshape(shape)
19
20 #import the train and test data and labels. Using all 60,000 images for training and 10,000 images for testing
21 raw_train = read_idx("train-images-idx3-ubyte")
22 train_data = np.reshape(raw_train, (60000, 28*28))
23 train_label = read_idx("train-labels-idx1-ubyte")
24
25 raw_test = read_idx("t10k-images-idx3-ubyte")
26 test_data = np.reshape(raw_test, (10000, 28*28))
27 test_label = read_idx("t10k-labels-idx1-ubyte")
28 #selecting all the data, digits 0-9 for training
29 idx = (train_label == 0) | (train_label == 1) | (train_label == 2) | (train_label == 3) | (train_label == 4) | (train_label == 5) | (train_label == 6) | (train_label == 7) | (train_label == 8) | (train_label == 9)
30 x = train_data[idx]
31 y = train_label[idx]
32
33
34 #creating the decision tree classifier and setting paramters, either gini or entropy criterion. As well as depth. Default is 5
35 DT = DecisionTreeClassifier(criterion='entropy', max_depth=8).fit(x, y)
36
37 #Setting the test Label
38 idx = (test_label == 9)
39 x_test = test_data[idx]
40 y_true = test_label[idx]
41 y_pred = DT.predict(x_test)
42
43 import itertools
44
45 #Create confusion matrix
46 def plot_confusion_matrix(cm, classes,
47                           normalize=False,
48                           title='Confusion matrix Decision Tree with entropy criterion max depth = 8',
49                           cmap=plt.cm.Blues):
50     """
51     This function prints and plots the confusion matrix.
52     Normalization can be applied by setting `normalize=True`.
53     """
54     if normalize:
55         cm = cm.astype('float') / cm.sum(axis=1)[:, np.newaxis]
56         print("Normalized confusion matrix")

```

Figure 8: Decision tree Sklearn model


```

57     else:
58         print('Confusion matrix, without normalization')
59
60     print(cm)
61
62     plt.imshow(cm, interpolation='nearest', cmap=cmap)
63     plt.title(title)
64     plt.colorbar()
65     tick_marks = np.arange(len(classes))
66     plt.xticks(tick_marks, classes, rotation=45)
67     plt.yticks(tick_marks, classes)
68
69     fmt = '.2f' if normalize else 'd'
70     thresh = cm.max() / 2.
71     for i, j in itertools.product(range(cm.shape[0]), range(cm.shape[1])):
72         plt.text(j, i, format(cm[i, j], fmt),
73                 horizontalalignment="center",
74                 color="white" if cm[i, j] > thresh else "black")
75
76     plt.tight_layout()
77     plt.ylabel('True label')
78     plt.xlabel('Predicted label')
79
80     #create decision tree image
81     dot_data = StringIO()
82     export_graphviz(DT, out_file=dot_data,
83                     filled=True, rounded=True,
84                     special_characters=True, class_names=['0', '1', '2', '3', '4', '5', '6', '7', '8', '9'])
85     graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
86     graph.write_png('tree.png')
87     Image(graph.create_png())
88
89
90     #using the sklearn score function, we find the mean accuracy of the decision tree
91     score = DecisionTreeClassifier.score(DT, x_test, y_true)
92     scorestr = str(score)
93     print("The mean accuracy score using entropy criterion and max depth = 8 is: " + scorestr)
94
95     param = DecisionTreeClassifier.get_params(DT, deep=True)
96
97     print(param)
98
99     #plots confusion matrix non-normalized
100     cm = metrics.confusion_matrix(y_true, y_pred)
101     plot_confusion_matrix(cm, ["0", "1", "2", "3", "4", "5", "6", "7", "8", "9"])
102     plt.show()
103
104     #plots confusion matrix normalized
105     plot_confusion_matrix(cm, ["0", "1", "2", "3", "4", "5", "6", "7", "8", "9"], normalize=True)
106     plt.show()

```

Figure 9 : Decision Tree Sklearn continued

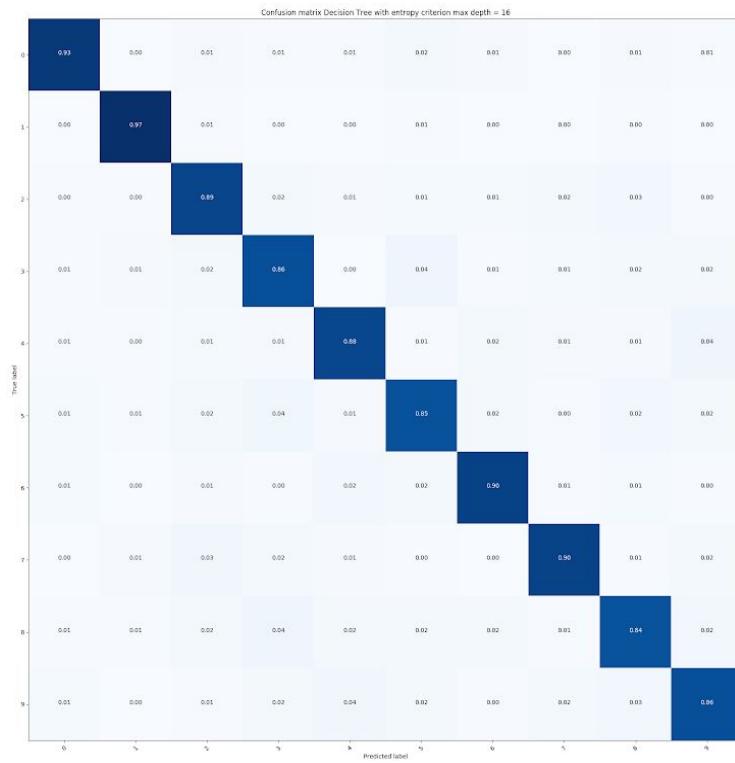


Fig10: Normalized Confusion matrix: Decision Tree, Max depth = 16, Entropy, Digits 0-9

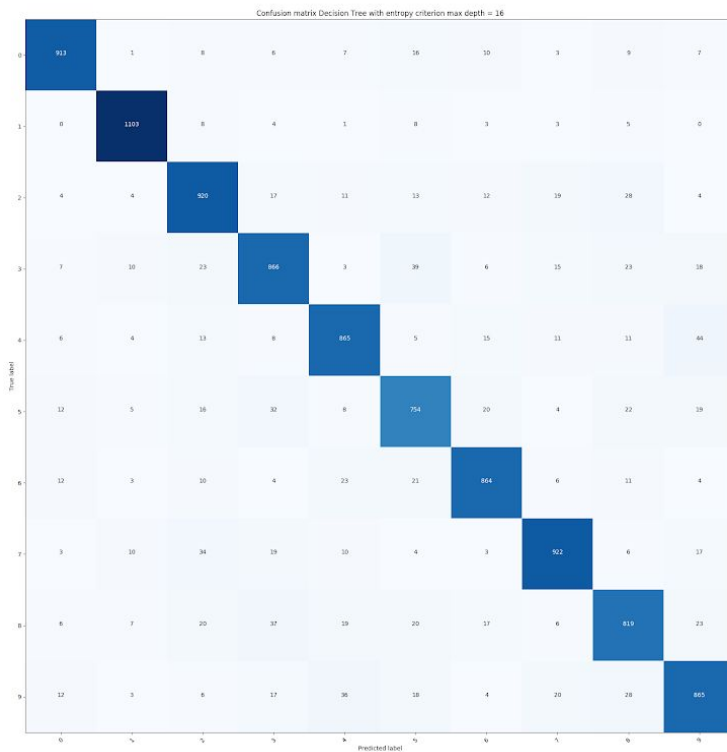


Fig10: Confusion matrix: Decision Tree, Max depth = 16, Entropy, Digits 0-9

3. Random Forest

To run a random forest model, i was able to reuse all of my previous code. As shown in Figure 11, I simply had to import the Random Forest Sklearn module and change the classifier. Parameters are similar to Decision Trees.

FIND ALL			FIND NINE		
Mean Accuracy	Depth	Criterion	Mean Accuracy	Depth	Criterion
90.70%	8	gini	89.90%	8	gini
94.68%	16	gini	92.60%	16	gini
94.32%	Until pure	gini	92.66%	Until pure	gini
Mean Accuracy	Depth	Criterion	Mean Accuracy	Depth	Criterion
90.45%	8	entropy	88.20%	8	entropy
94.84%	16	entropy	92.16%	16	entropy
95.02%	Until pure	entropy	92.17%	Until pure	entropy

Table 3: Random Forest model runs final results

```

38
39 #creating the RANDOM FOREST classifier and setting paramters, either gini or entropy criterion. As well as dep
40 # x, y = make_classification(n_features=784)
41 # DT = RandomForestClassifier(criterion='gini', n_jobs=-1).fit(x, y)
42 DT = RandomForestClassifier(criterion='entropy', max_depth=None, n_jobs=-1)
43 RndTree = DT.fit(x,y)
44 print(RndTree)

```

Fig 11: Code change to use random forest instead of decision tree, with Sklearn.

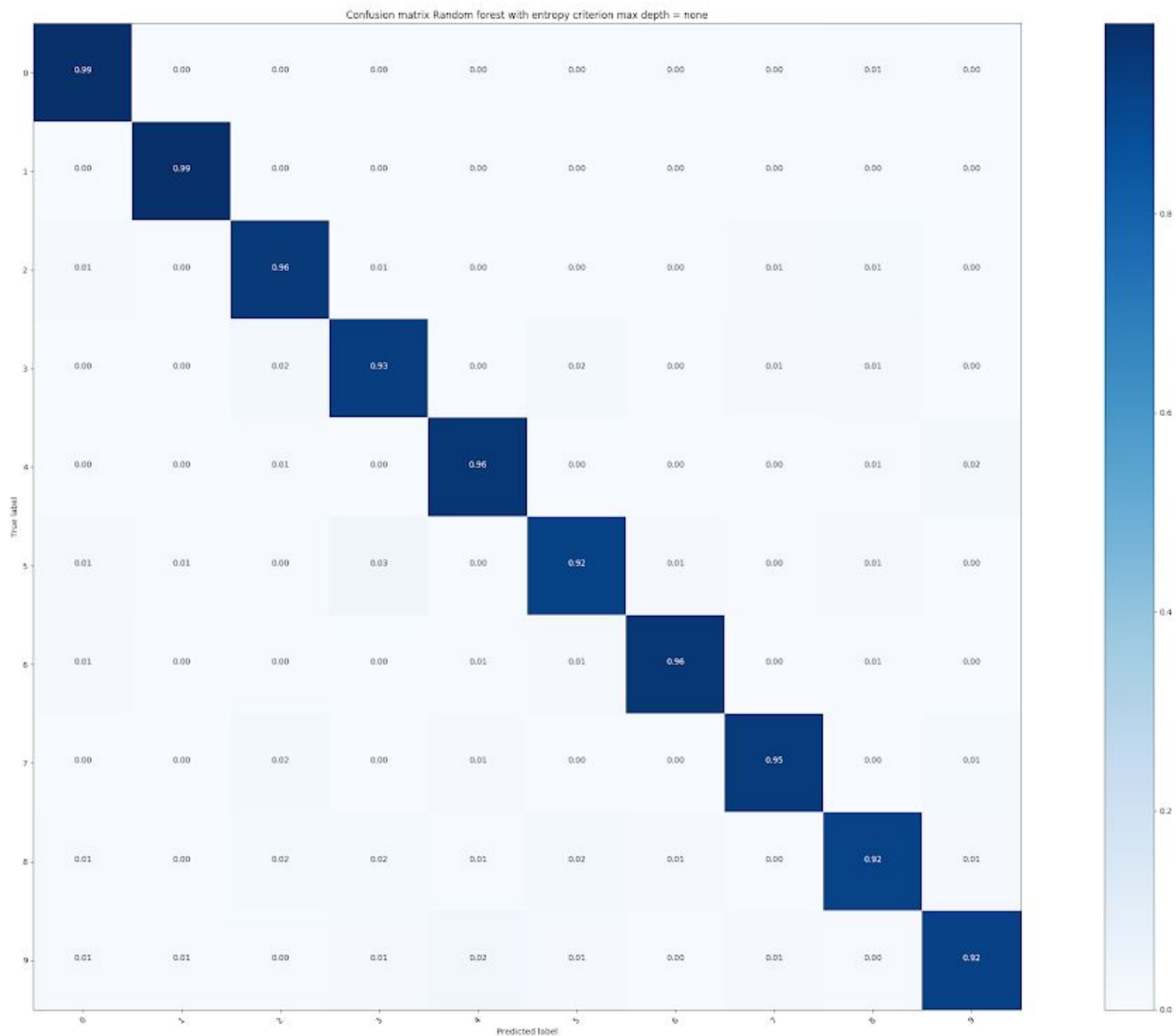


Fig 12: Normalized confusion matrix, random forest, Entropy , Max depth= None.

Final Results: Finding all digits 0-9

KNN- Mean Accuracy **97.3%** - K neighbors = 3, Minkowski P=3

Decision Tree- Mean Accuracy **88.9%** -Max Depth = 16, Entropy

Random Forest- Mean Accuracy **95.0%** - Max Depth = None, Entropy

Final Results: Finding just the 9s

KNN- Mean Accuracy **96.8%** - K neighbors = 3, Minkowski P=3

Decision Tree- Mean Accuracy **87.0%** -Max Depth = 16, Gini

Random Forest- Mean Accuracy **92.66%** - Max Depth = None, Gini

The Multiclass classification was more accurate in all three models, when compared to the binary classification. This is probably due to the fact that there is more data available to fit the clusters, nodes, and leaves.

It is also shown in the confusion matrices, that the digits 0 and 1, have a higher classification accuracy.