

**Özyeğin University**

Fall 2020

C.S. 554.A

Homework II

**Submitted by:** M. Furkan Oruc (MSc. Student at CE  
Department)

**Student ID:** S025464

**Submission Date:** 13/11/2020

## Content

1. Nearest Mean Classifier
  - a. Test Data Outcomes
  - b. Training Data Outcomes
2. Images' Visualization (Based on mean pixel values of MNIST Training Data)
3. K – NN Classifier
  - a. Test Data Outcomes
  - b. Training Data Outcomes
4. Complete Source Code
5. References

### 1. Nearest Mean Classifier

Nearest mean classifier operates as a non parametric classification procedure in the following steps.

In this problem, euclidian distance is used due to the fact that the entire range of pixels are provided in the same scale, 0-255 which is extracted from MNIST data. On the other hand, total of 1000 instances are provided for each, training and test purposes.

Euclidian distance, which is a commonly used distance metric can be defined as:

$$E(x, y) = \sqrt{\sum_{i=1}^a d(x_i, y_i)^2}$$

where x and y are two examples, a is the number of attributes and xi refers to the ith attribute value for x.

If the variances weren't assumed to be equal, it would not be possible to act with euclidian distance, we would need the Mahalanobis distance. (Alpaydın, 2014)

In the scope of this question, in order to find the euclidian distance between distinct test instances, mean values for each distinct label and each distinct pixel are calculated. The following code represents the main part of this process.

```

1. for k in range (0,10):
2.     df_pixel_sum[k] = df_train.loc[df_train["label"] == k ,:].sum()
3.     df_pixel_mean[k] = df_train.loc[df_train["label"] == k ,:].mean()
4.     k = k + 1

```

A highlighted part of the source code, where subtraction of a test instance's distinct pixel value from the mean of the related pixel and all labels can be followed as below.

```

1. for i in range(1, 785): # for each pixel. cover all pixels
2.     pixel_difference.append((df_test.iloc[r,i] - df_pixel_mean.iloc[i,
    m])**2)
3.     #For the 0th image in test data, subtract and square each pixel's value from each pixel's mean for each number
4.     i = i + 1
5.     df_pixel_difference[m] = pixel_difference # Store the calculation for each image

```

Right after calculation of distance between each number's mean and storage of that difference for each instance of the test data's each pixel values, the minimum distance to a mean is calculated as below. In this step, also index of the minimum sum value for each instance is also stored. The index value helps storing the assigned instance to the test data based on training data. Of course, while the confusion matrix was being calculated, the instances were belonging to training data.

```

1. #for k in range (0,10): # For each number's mean, find the distance between 0th test data
2.     number_identifier.append(df_pixel_difference.iloc[:, m].sum())
3.     df_number_identifier [m] = number_identifier
4.     m = m + 1
5.     #k = k + 1
6.     number_identifier.clear()
7.
8.     for s in range (0,10): # Finding the minimum distance between 0th test element and number image means
9.         near_mc.append(df_number_identifier.iloc[0,s])
10.        df_near_mc[s] = near_mc
11.        near_mc.clear()
12.        s = s + 1
13.    df_near_mc_transpose = df_near_mc.transpose()
14.
15.    assigned_instance = (df_near_mc_transpose.idxmin())
16.    df_assigned_instance[r] = assigned_instance
17.    #assigned_instance.clear()

```

In order to represent the assigned values of the instances, a confusion matrix is processed.

```
1. for f in range(len(df_test)):
2.
3.     if df_assigned_instance.iloc[f,0] == df_test.iloc[f,0]:
4.         x = df_test.iloc[f,0]
5.         df_confusion.iloc[x,x] = df_confusion.iloc[x,x] + 1
6.
7.     elif df_assigned_instance.iloc[f,0] != df_test.iloc[f,0]:
8.         y = df_test.iloc[f,0]
9.         x = df_assigned_instance.iloc[f,0]
10.        df_confusion.iloc[y,x] = df_confusion.iloc[y,x] + 1
```

In the confusion matrix, each row represents the real label of each test data. Each column represents the assigned label for each test instances. Since there are 100 instances for each label, each row makes up a total of 100 value in the matrix.

Diagonal sums represent the number of correct assignments to each label. Off-diagonal cells (sums) represent the incorrect assignments.

The total number of diagonal values represent the number of correct assignments in the entire dataset, for each case. The accuracy rates for both test dataset and training dataset is provided above. It's important to note that based on the brief of the question, both confusion matrices are based on training dataset means. On the other hand, classification is applied for both test dataset and training dataset.

***Accuracy Rate on Test Data: 70.5%***

***Accuracy Rate on Training Data: 81.2%***

### ***Confusion Matrix for the Test Data***

It can be clearly observed that the most of the incorrect assignments are observed in numbers 2,3 and 8.

The most problematic assignment has been on number 3, from which 31 of them are assigned to the number 5. It can be stated that MNIST dataset provides 5's which may be mistaken as 3 by the nearest mean classifier logic.

| Index | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|----|----|----|----|----|----|----|----|----|
| 0     | 86 | 0  | 1  | 1  | 1  | 6  | 4  | 0  | 0  | 1  |
| 1     | 0  | 95 | 0  | 0  | 0  | 5  | 0  | 0  | 0  | 0  |
| 2     | 1  | 24 | 59 | 3  | 2  | 1  | 0  | 3  | 6  | 1  |
| 3     | 0  | 3  | 1  | 59 | 0  | 31 | 0  | 3  | 1  | 2  |
| 4     | 0  | 2  | 0  | 0  | 69 | 0  | 2  | 0  | 0  | 27 |
| 5     | 3  | 4  | 1  | 9  | 3  | 67 | 0  | 4  | 2  | 7  |
| 6     | 5  | 1  | 7  | 0  | 14 | 9  | 63 | 0  | 1  | 0  |
| 7     | 0  | 10 | 3  | 0  | 2  | 1  | 0  | 73 | 0  | 11 |
| 8     | 2  | 3  | 3  | 8  | 3  | 9  | 2  | 2  | 58 | 10 |
| 9     | 0  | 1  | 0  | 2  | 15 | 0  | 0  | 4  | 2  | 76 |

*Figure 1: Nearest Mean Classification Confusion Matrix for Test Dataset*

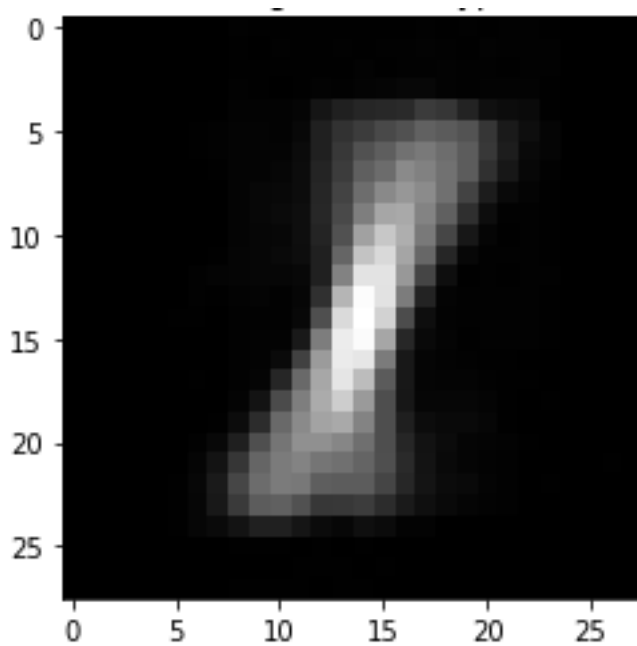
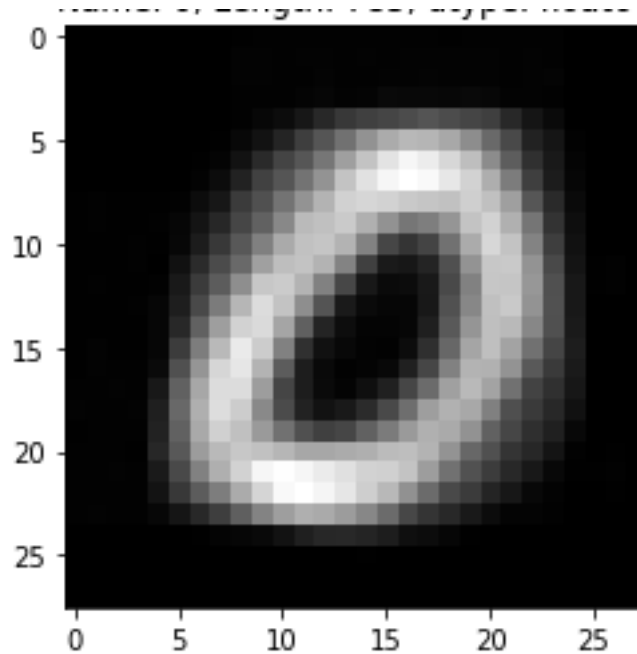
### ***Confusion Matrix for the Training Data***

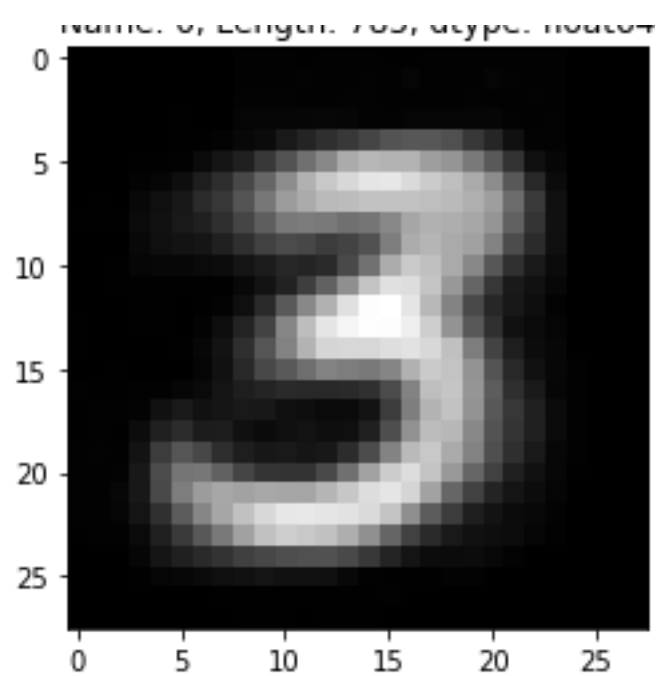
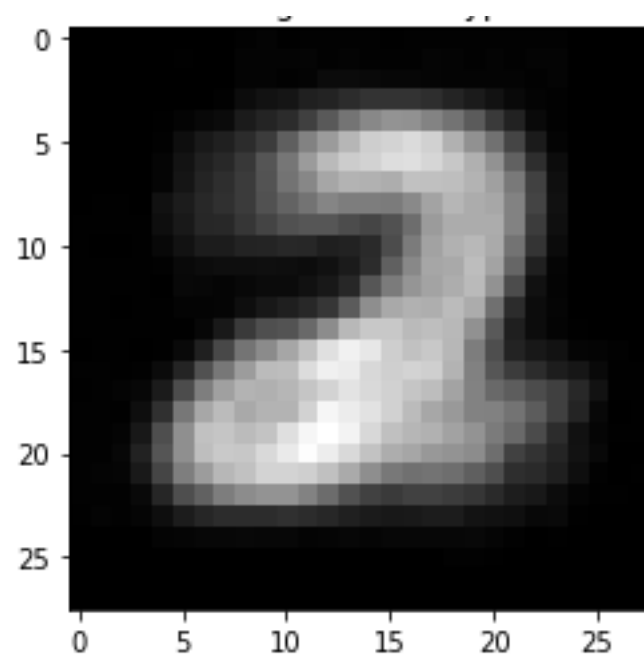
| Index | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|----|----|----|----|----|----|----|----|----|
| 0     | 91 | 0  | 1  | 0  | 0  | 3  | 3  | 1  | 1  | 0  |
| 1     | 0  | 96 | 0  | 0  | 0  | 3  | 0  | 0  | 1  | 0  |
| 2     | 2  | 14 | 69 | 3  | 3  | 1  | 1  | 4  | 1  | 2  |
| 3     | 1  | 1  | 1  | 76 | 0  | 13 | 0  | 3  | 2  | 3  |
| 4     | 0  | 2  | 0  | 0  | 84 | 0  | 2  | 0  | 0  | 12 |
| 5     | 4  | 4  | 1  | 10 | 5  | 73 | 2  | 1  | 0  | 0  |
| 6     | 1  | 5  | 0  | 0  | 4  | 6  | 84 | 0  | 0  | 0  |
| 7     | 0  | 6  | 1  | 0  | 2  | 1  | 0  | 86 | 1  | 3  |
| 8     | 0  | 8  | 2  | 5  | 0  | 8  | 1  | 0  | 73 | 3  |
| 9     | 1  | 2  | 0  | 2  | 7  | 3  | 1  | 4  | 0  | 80 |

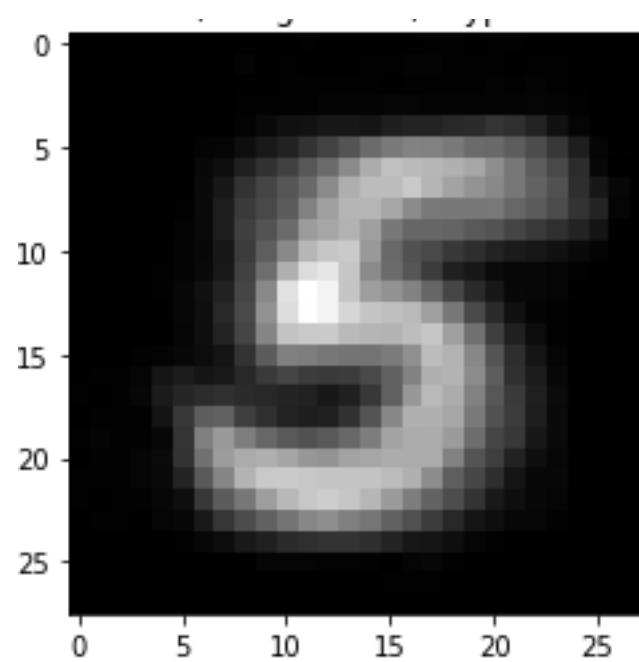
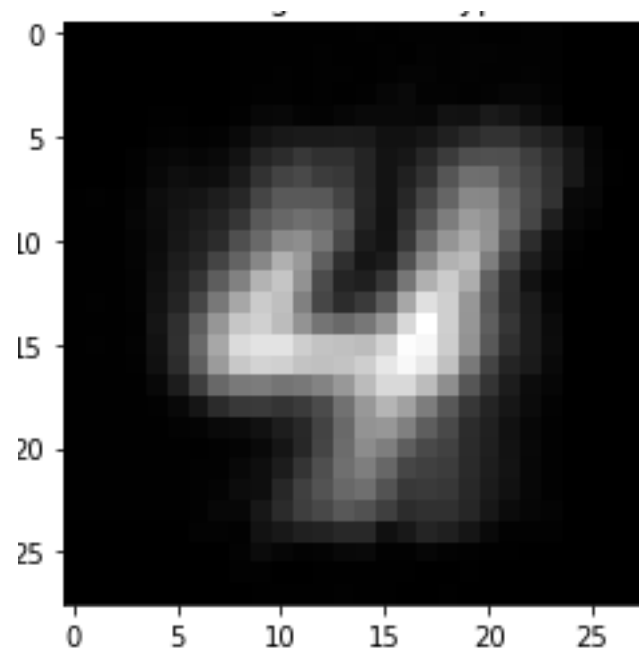
*Figure 2: Nearest Mean Classification Confusion Matrix for Training Dataset*

## 2. Images' Visualization

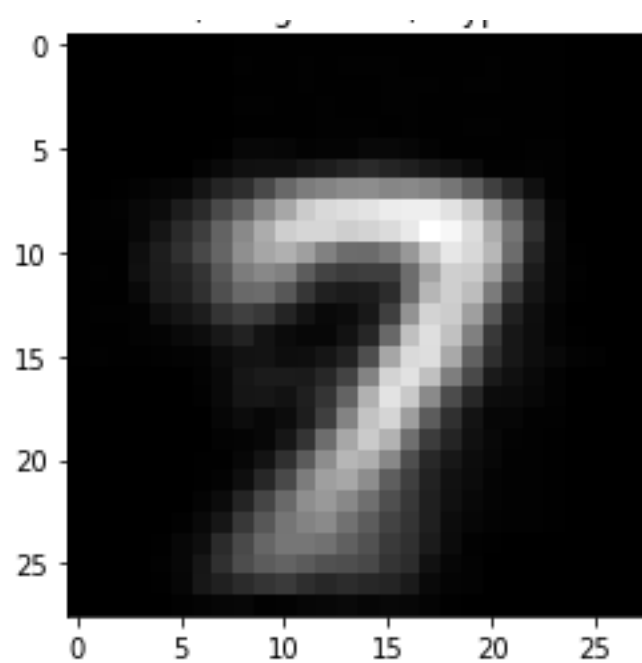
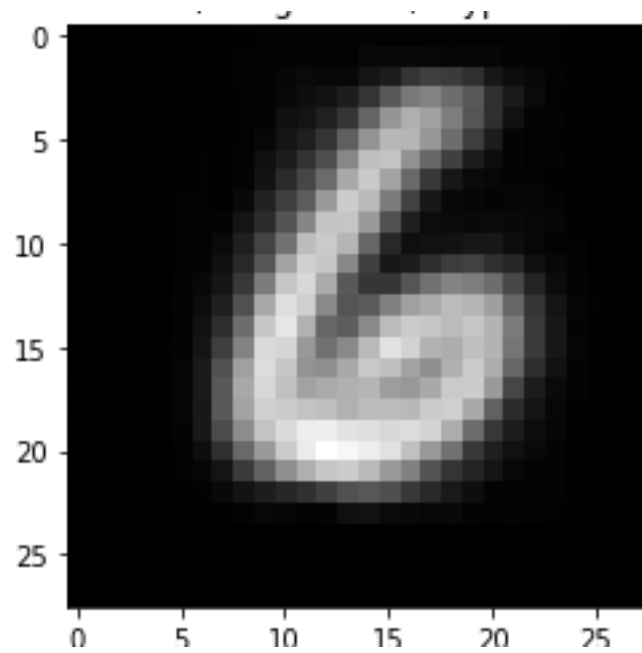
Visualization of MNIST training data's mean values for each number can be observed below.

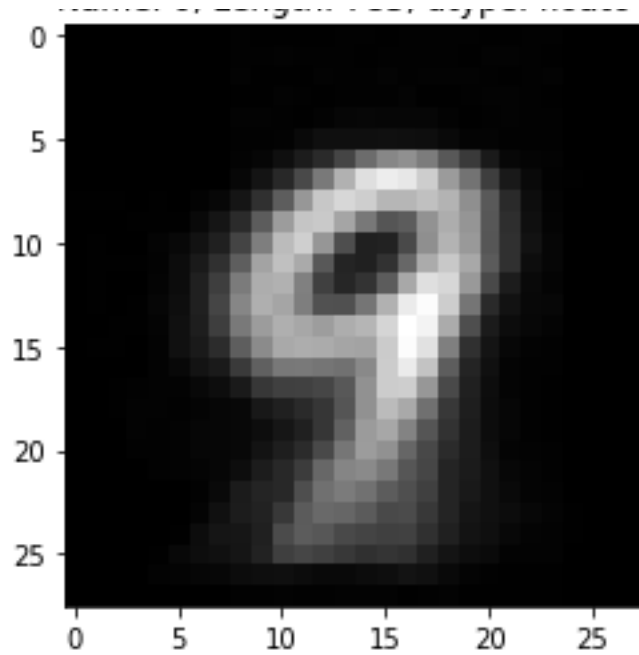
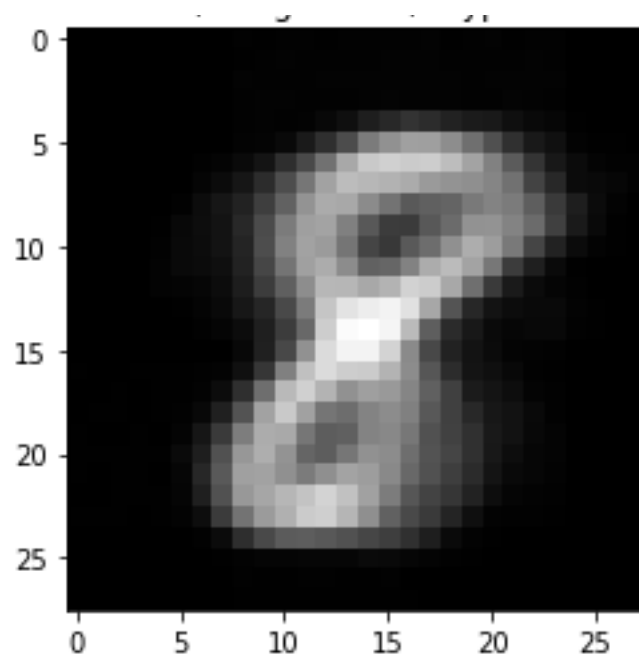












### 3. K-Nearest Neighbor Classification

In the non-parametric classification based on k-nn algorithm, an object is classified based the plurality vote of its neighbors. In other words, based on number of k's, each instance of the dataset is classified based on the number of its closest neighbors inside the “k” set.

Each training instance votes for its class and has no other effect on other classes. Normally, the weight of vote is given to the closest instances. Though, in the k-nn classifier case, the class having the most examples is looked for. It can be also expressed as:

$$\hat{P}(C_i|\mathbf{x}) = \frac{\hat{p}(\mathbf{x}|C_i)\hat{P}(C_i)}{\hat{p}(\mathbf{x})} = \frac{k_i}{k}$$

In this classification's scope, k is valued as 1, 3, 5, 7, 9, 11 for each case.

There has been a quite limited number of ties. Those ties are specifically represented in the accuracy and results section.

A brief part of the source code is represented below, where the Euclidian distance is calculated between each instance of the test data and each pixel with respect to each value of the training dataset.

```
1. for k in range(0, 1000): #1000
2.     summer.clear()
3.     for t in range (1, 785):
4.         distance_knn.clear()
5.         for m in range (len(df_train)):
6.
7.             distance_knn.append((df_test.iloc[k,t] - df_train.iloc[m, t])**2)
8.
9.         df_knn[t] = distance_knn
```

Later on, each Euclidian distance is stored for each test instance and respective training instances. The sum of the pixels are calculated in order to detect the total representative distance between each instance.

```
1.     for n in range(0,1000):
2.
3.         sum_holder = df_knn.iloc[n, :].sum()
4.         summer.append(sum_holder)
5.     df_summer[k] = summer
```

In the next step, for each instance, the values are stored and sorted. While sorting, labels for each instance is calculated.

```
1. for x in range (0,1000): #1000
2.     knn_label.clear()
3.     for y in range (0,11):
4.         knn_label.append(df_train.iloc[df_storella.iloc[y,x],0])
5.     df_knn_label[x] = knn_label
```

Right after the creation of the matrice where for each test data instance, the most repeated index is extracted for each k case in the loop (1, 3, 5, 7, 9, 11); mode specialty is used. By this method, the most repeated cases are determined. This helped storage of the most repeated labels for each k case distinctively.

```
1. for p in range (1,11,2):
2.     for o in range (0,1000): #1000, k = 11
3.         mode_holder.append(df_knn_label.iloc[:p,o].mode())
4.
5.     #df_mode_holder [t] = mode_holder ###
6.
7. df_mode_holder = pd.DataFrame(mode_holder)
```

The confusion matrix is created in a manner where rows' indexes represent each real label already given in the test dataset and each column index represents the assigned label, of course in different matrices for each k case.

Due to the timely cost of nonparametric classification and the brief, the classification is applied on the test data.

On the other hand, one issue which was resolved during the classification was ties. Among 6000 classification of test data, 230 of the classifications occurred with a tie.

In order to break the tie, some specific applications are utilized. Based on the wide scale practice, the closest label is assigned between the multiple members of the tie. In order to make it happen, each “tie” instance is examined and its closest top 6 labels are compared with the tie members. The algorithm suggests that if one of the maximum 5 tie labels are included in the closest 6 label set, then the closest instance among top 6 label is assigned as the tie-breaker label to the related instance.

The assumption is based on both analytical and exploratory analysis on the dataset. It is assumed that if a tie-breaker is present, it would be included between top 6 labels of that instance. In other words, one of the closest 6 label is assumed to be enough to be equal to one of the tie members. The after analysis also suggests that the accuracy on the confusion matrices is presented as 99.99%.

```
1. for t in range (len(nan_2_1)):
2.     if nan_2_1[t] > 1000 and nan_2_1[t] < 2000:
3.         if df_mode_holder_1.iloc[nan_2_1[t],0] == df_knn_label.iloc[0,nan_2_1[
t]-1000]:
4.             counter = counter + 1
5.             counter_list_3.append(nan_2_1[t])
6.         elif nan_2_1[t] > 2000 and nan_2_1[t] < 3000:
7.             if df_mode_holder_1.iloc[nan_2_1[t],4] == df_knn_label.iloc[5,nan_2_1[
t]-2000]:
8.                 counter = counter + 1
9.                 counter_list_6.append(nan_2_1[t])
10.        elif nan_2_1[t] > 3000 and nan_2_1[t] < 4000:
11.            if df_mode_holder_1.iloc[nan_2_1[t],4] == df_knn_label.iloc[5,nan_2_1[
t]-3000]:
12.                counter = counter + 1
13.                counter_list_6.append(nan_2_1[t])
14.        elif nan_2_1[t] > 4000 and nan_2_1[t] < 5000:
15.            if df_mode_holder_1.iloc[nan_2_1[t],4] == df_knn_label.iloc[5,nan_2_1[
t]-4000]:
16.                counter = counter + 1
17.                counter_list_6.append(nan_2_1[t])
18.        elif nan_2_1[t] > 5000 and nan_2_1[t] < 6000:
19.            if df_mode_holder_1.iloc[nan_2_1[t],4] == df_knn_label.iloc[5,nan_2_1[
t]-5000]:
20.                counter = counter + 1
21.                counter_list_6.append(nan_2_1[t])
22. for p in range (5000, 6000): # do it for each k=3,5,7,9,11 for respective 1000
-2000-3000 etc
23.     for z in range (1,5):
```

```

24.         if df_mode_holder_1.iloc[p, z] >= 0: #hold p
25.             if p in counter_list:
26.                 df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[0,p-5000]
27.             elif p in counter_list_2:
28.                 df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[1,p-5000]
29.             elif p in counter_list_3:
30.                 df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[2,p-5000]
31.             elif p in counter_list_4:
32.                 df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[3,p-5000]
33.             elif p in counter_list_5:
34.                 df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[4,p-5000]
35.             elif p in counter_list_6:
36.                 df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[5,p-5000]
37.

```

Below, each confusion matrix is provided with the respective accuracy rate. Also, The total number of diagonal values represent the number of correct assignments in the entire dataset, for each case.

The K-NN Classification is based on both test and training data instances. Following table represents the changing accuracy rate based on different k values, ranging(1,11).

### **Post Reasoning on the Accuracy Rate Difference Between Test & Training Datasets**

It'll be observed in the next page that there is a dramatic accuracy in the training dataset, while the test data performs moderately. After an exploratory post analysis, it can be stated that, in the training case, due to the fact that the closest instance will definitely carry the label of the instance's real label, that's a strong sign. Later on, in the bigger k's such as 9 and 11, it can be stated that training data is made of instances which are ordered by their label and their pixel values (0,255) seem to be closer to each other when they are ordered line by line. This suggests that it was dramatically easier for training data instances to find neighbors which carry the same label with them, since their consecutive instances were clearly similar to each other.

### Test Data Accuracy Rate / K Plot

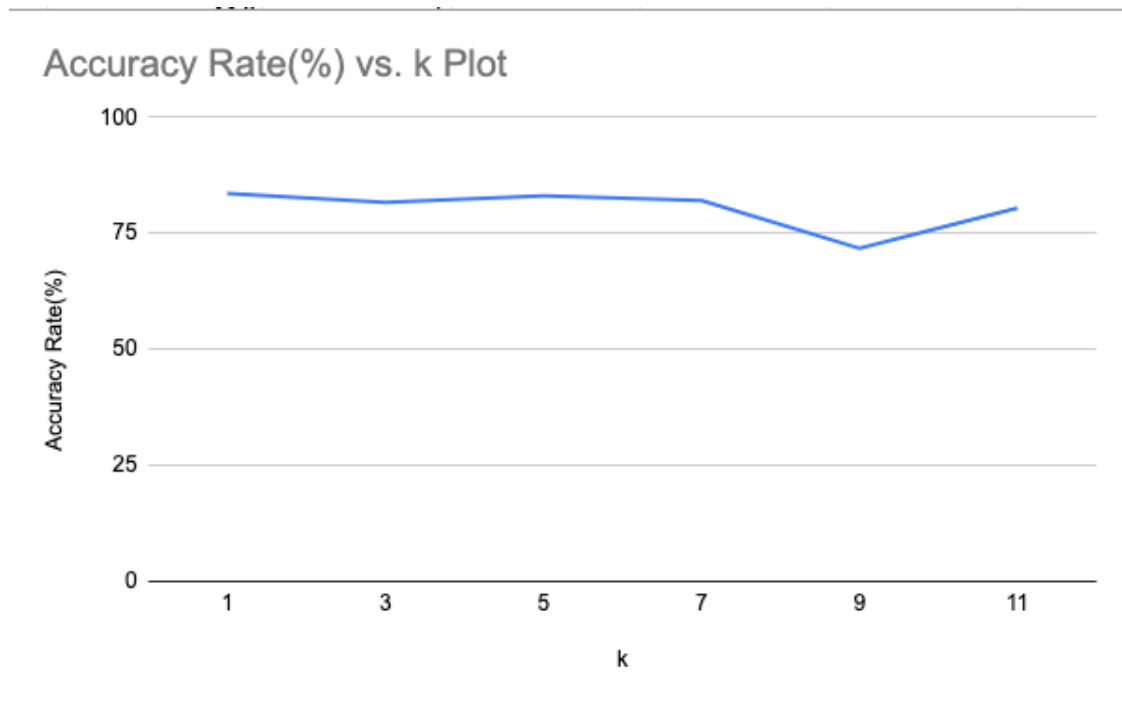


Table 1: Test Data: Plot of Changing Accuracy Rate based on different k values

### Training Data Accuracy Rate / K Plot

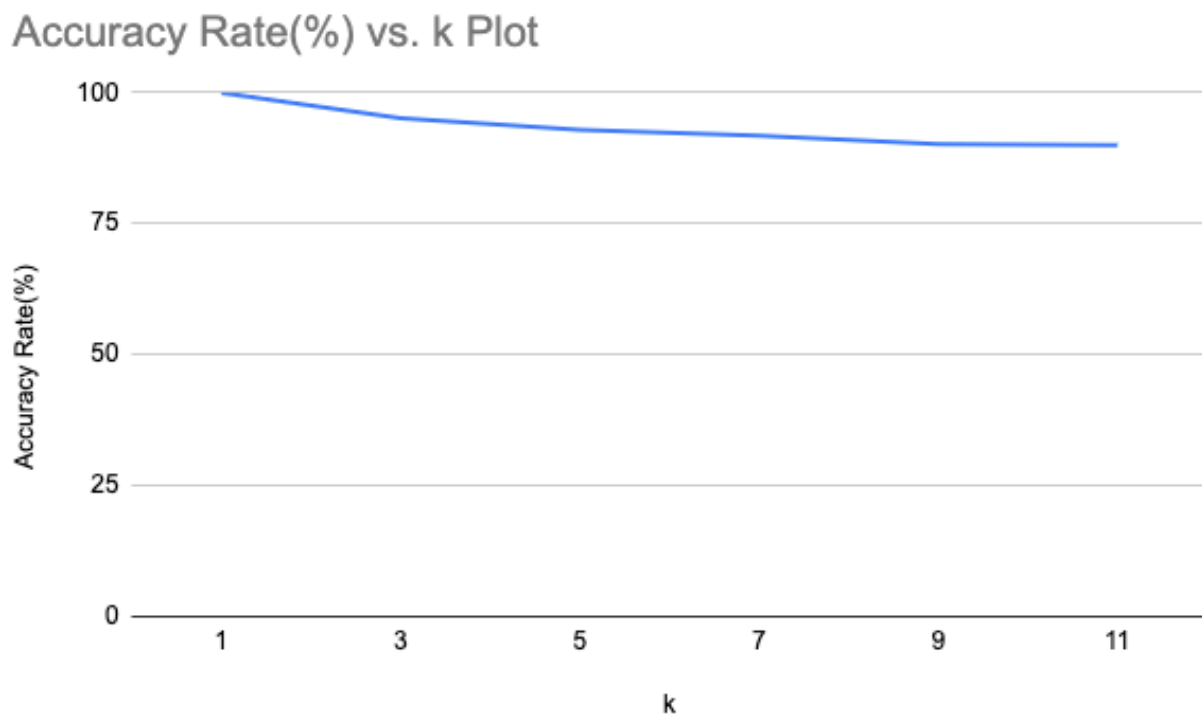


Table 2: Training Data: Plot of Changing Accuracy Rate based on different k value

## Test Data Confusion Matrices

### Confusion Matrix for the Test Data – Case: $k=1$

| Index | 0  | 1   | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|-----|----|----|----|----|----|----|----|----|
| 0     | 94 | 0   | 0  | 0  | 1  | 0  | 3  | 1  | 0  | 1  |
| 1     | 0  | 100 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 2     | 1  | 4   | 80 | 2  | 1  | 0  | 2  | 7  | 3  | 0  |
| 3     | 0  | 0   | 1  | 76 | 0  | 14 | 2  | 3  | 2  | 2  |
| 4     | 0  | 2   | 0  | 0  | 69 | 1  | 2  | 2  | 0  | 24 |
| 5     | 1  | 3   | 0  | 6  | 4  | 74 | 4  | 2  | 2  | 4  |
| 6     | 4  | 0   | 0  | 0  | 1  | 3  | 92 | 0  | 0  | 0  |
| 7     | 0  | 4   | 0  | 0  | 2  | 1  | 0  | 90 | 0  | 3  |
| 8     | 2  | 2   | 6  | 4  | 1  | 3  | 4  | 1  | 73 | 4  |
| 9     | 0  | 0   | 0  | 1  | 4  | 1  | 0  | 6  | 0  | 88 |

Figure 3: K-nn classification:  $k = 1$

**Accuracy Rate on  $k = 1$ : 83.6%**

### Confusion Matrix for the Test Data – Case: $k=3$

| Index | 0  | 1   | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|-----|----|----|----|----|----|----|----|----|
| 0     | 96 | 0   | 1  | 0  | 0  | 0  | 3  | 0  | 0  | 0  |
| 1     | 0  | 100 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 2     | 2  | 11  | 70 | 1  | 1  | 1  | 2  | 6  | 6  | 0  |
| 3     | 1  | 2   | 6  | 74 | 0  | 9  | 1  | 1  | 3  | 3  |
| 4     | 0  | 2   | 0  | 0  | 70 | 0  | 2  | 1  | 0  | 25 |
| 5     | 1  | 5   | 0  | 6  | 4  | 77 | 4  | 0  | 0  | 3  |
| 6     | 5  | 0   | 0  | 0  | 2  | 3  | 90 | 0  | 0  | 0  |
| 7     | 0  | 6   | 1  | 0  | 2  | 1  | 0  | 87 | 0  | 3  |
| 8     | 5  | 2   | 2  | 6  | 0  | 4  | 3  | 3  | 68 | 7  |
| 9     | 0  | 1   | 0  | 1  | 5  | 0  | 0  | 8  | 0  | 85 |

Figure 4: K-nn classification:  $k = 3$

**Accuracy Rate on  $k = 3$ : 81.7%**



***Confusion Matrix for the Test Data – Case:  $k=5$***

| Index | 0  | 1   | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|-----|----|----|----|----|----|----|----|----|
| 0     | 93 | 0   | 1  | 0  | 0  | 1  | 5  | 0  | 0  | 0  |
| 1     | 0  | 100 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 2     | 1  | 13  | 67 | 1  | 2  | 1  | 2  | 7  | 6  | 0  |
| 3     | 0  | 3   | 1  | 76 | 0  | 11 | 1  | 2  | 3  | 3  |
| 4     | 0  | 2   | 0  | 0  | 70 | 0  | 2  | 1  | 0  | 25 |
| 5     | 1  | 4   | 0  | 3  | 2  | 78 | 5  | 1  | 0  | 6  |
| 6     | 3  | 0   | 0  | 0  | 1  | 3  | 93 | 0  | 0  | 0  |
| 7     | 0  | 7   | 1  | 0  | 1  | 1  | 0  | 89 | 0  | 1  |
| 8     | 3  | 2   | 2  | 5  | 1  | 2  | 4  | 2  | 73 | 6  |
| 9     | 0  | 1   | 0  | 0  | 2  | 0  | 0  | 4  | 1  | 92 |

*Figure 5: K-nn classification:  $k = 5$*

***Accuracy Rate on  $k = 5$ : 83.1%***

***Confusion Matrix for the Test Data – Case:  $k=7$***

| Index | 0  | 1   | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|-----|----|----|----|----|----|----|----|----|
| 0     | 93 | 0   | 1  | 0  | 0  | 1  | 5  | 0  | 0  | 0  |
| 1     | 0  | 100 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 2     | 1  | 18  | 63 | 1  | 2  | 0  | 2  | 7  | 6  | 0  |
| 3     | 0  | 3   | 2  | 77 | 0  | 9  | 1  | 2  | 3  | 3  |
| 4     | 0  | 2   | 0  | 0  | 71 | 0  | 2  | 1  | 0  | 24 |
| 5     | 1  | 4   | 0  | 4  | 1  | 74 | 5  | 1  | 2  | 8  |
| 6     | 3  | 0   | 0  | 0  | 1  | 3  | 93 | 0  | 0  | 0  |
| 7     | 0  | 9   | 1  | 0  | 2  | 0  | 0  | 84 | 0  | 4  |
| 8     | 3  | 3   | 1  | 4  | 1  | 3  | 2  | 3  | 74 | 6  |
| 9     | 0  | 2   | 0  | 0  | 2  | 0  | 0  | 3  | 1  | 92 |

*Figure 6: K-nn classification:  $k = 7$*

***Accuracy Rate on  $k = 7$ : 82.1%***

***Confusion Matrix for the Test Data – Case:  $k=9$***

| Index | 0  | 1   | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|-----|----|----|----|----|----|----|----|----|
| 0     | 93 | 0   | 0  | 0  | 1  | 2  | 3  | 1  | 0  | 0  |
| 1     | 0  | 100 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 2     | 1  | 19  | 61 | 1  | 2  | 0  | 2  | 7  | 6  | 1  |
| 3     | 0  | 5   | 2  | 78 | 0  | 6  | 1  | 2  | 3  | 3  |
| 4     | 0  | 2   | 0  | 0  | 71 | 0  | 2  | 1  | 0  | 24 |
| 5     | 1  | 4   | 0  | 4  | 1  | 76 | 4  | 2  | 0  | 8  |
| 6     | 3  | 0   | 0  | 0  | 2  | 3  | 92 | 0  | 0  | 0  |
| 7     | 0  | 9   | 0  | 1  | 0  | 1  | 0  | 84 | 0  | 5  |
| 8     | 4  | 3   | 1  | 4  | 2  | 3  | 1  | 3  | 70 | 9  |
| 9     | 0  | 1   | 0  | 0  | 3  | 0  | 0  | 2  | 1  | 93 |

*Figure 7: K-nn classification:  $k = 9$*

***Accuracy Rate on  $k = 9$ : 71.8%***

***Confusion Matrix for the Test Data – Case:  $k=11$***

| Index | 0  | 1   | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|-----|----|----|----|----|----|----|----|----|
| 0     | 93 | 0   | 1  | 0  | 1  | 2  | 2  | 1  | 0  | 0  |
| 1     | 0  | 100 | 0  | 0  | 0  | 0  | 0  | 0  | 0  | 0  |
| 2     | 1  | 21  | 58 | 1  | 3  | 0  | 3  | 6  | 6  | 1  |
| 3     | 0  | 6   | 1  | 76 | 0  | 8  | 1  | 2  | 3  | 3  |
| 4     | 0  | 3   | 0  | 0  | 68 | 0  | 2  | 1  | 0  | 26 |
| 5     | 1  | 4   | 0  | 3  | 2  | 75 | 5  | 2  | 0  | 8  |
| 6     | 3  | 0   | 0  | 0  | 4  | 3  | 90 | 0  | 0  | 0  |
| 7     | 0  | 9   | 0  | 1  | 3  | 0  | 0  | 83 | 0  | 4  |
| 8     | 3  | 3   | 2  | 4  | 2  | 4  | 2  | 3  | 69 | 8  |
| 9     | 0  | 1   | 0  | 0  | 2  | 0  | 0  | 3  | 1  | 93 |

*Figure 8: K-nn classification:  $k = 11$*

***Accuracy Rate on  $k = 11$ : 80.5%***

## Training Data Confusion Matrices

*Confusion Matrix for the Training Data – Case:  $k=1$*

| Index | 0   | 1   | 2   | 3   | 4   | 5   | 6   | 7   | 8   | 9   |
|-------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 0     | 100 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   |
| 1     | 0   | 100 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   |
| 2     | 0   | 0   | 100 | 0   | 0   | 0   | 0   | 0   | 0   | 0   |
| 3     | 0   | 0   | 0   | 100 | 0   | 0   | 0   | 0   | 0   | 0   |
| 4     | 0   | 0   | 0   | 0   | 100 | 0   | 0   | 0   | 0   | 0   |
| 5     | 0   | 0   | 0   | 0   | 0   | 100 | 0   | 0   | 0   | 0   |
| 6     | 0   | 0   | 0   | 0   | 0   | 0   | 100 | 0   | 0   | 0   |
| 7     | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 100 | 0   | 0   |
| 8     | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 100 | 0   |
| 9     | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 100 |

*Accuracy Rate on  $k = 1$ : 100%*

*Confusion Matrix for the Training Data – Case:  $k=3$*

| Index | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|----|----|----|----|----|----|----|----|----|
| 0     | 99 | 0  | 0  | 0  | 0  | 0  | 1  | 0  | 0  | 0  |
| 1     | 0  | 99 | 0  | 0  | 0  | 0  | 0  | 1  | 0  | 0  |
| 2     | 0  | 3  | 93 | 0  | 0  | 0  | 1  | 1  | 1  | 1  |
| 3     | 0  | 0  | 1  | 94 | 0  | 2  | 1  | 0  | 0  | 2  |
| 4     | 0  | 2  | 0  | 0  | 92 | 0  | 1  | 0  | 0  | 5  |
| 5     | 0  | 1  | 0  | 2  | 0  | 93 | 2  | 0  | 0  | 2  |
| 6     | 1  | 1  | 0  | 0  | 0  | 0  | 98 | 0  | 0  | 0  |
| 7     | 0  | 4  | 0  | 0  | 1  | 0  | 0  | 93 | 0  | 2  |
| 8     | 0  | 3  | 0  | 1  | 0  | 1  | 1  | 0  | 93 | 1  |
| 9     | 0  | 0  | 0  | 1  | 1  | 0  | 1  | 0  | 0  | 97 |

*Accuracy Rate on  $k = 3$ : 95.1%*

***Confusion Matrix for the Training Data – Case:  $k=5$***

| Index | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|----|----|----|----|----|----|----|----|----|
| 0     | 97 | 0  | 0  | 0  | 0  | 1  | 2  | 0  | 0  | 0  |
| 1     | 0  | 98 | 0  | 1  | 0  | 0  | 0  | 1  | 0  | 0  |
| 2     | 0  | 5  | 88 | 0  | 0  | 0  | 3  | 2  | 1  | 1  |
| 3     | 1  | 1  | 0  | 92 | 0  | 2  | 1  | 1  | 0  | 2  |
| 4     | 0  | 3  | 0  | 0  | 90 | 1  | 1  | 0  | 0  | 5  |
| 5     | 0  | 2  | 1  | 2  | 1  | 89 | 2  | 0  | 1  | 2  |
| 6     | 1  | 2  | 0  | 0  | 0  | 1  | 96 | 0  | 0  | 0  |
| 7     | 0  | 6  | 0  | 0  | 0  | 0  | 0  | 91 | 0  | 3  |
| 8     | 0  | 3  | 0  | 2  | 0  | 2  | 0  | 0  | 92 | 1  |
| 9     | 1  | 0  | 0  | 1  | 2  | 0  | 0  | 0  | 0  | 96 |

***Accuracy Rate on  $k = 5$ : 92.9%***

***Confusion Matrix for the Training Data – Case:  $k=7$***

| Index | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|----|----|----|----|----|----|----|----|----|
| 0     | 98 | 0  | 0  | 0  | 0  | 1  | 1  | 0  | 0  | 0  |
| 1     | 0  | 99 | 0  | 0  | 0  | 0  | 0  | 1  | 0  | 0  |
| 2     | 0  | 8  | 80 | 0  | 1  | 1  | 2  | 6  | 1  | 1  |
| 3     | 1  | 1  | 0  | 92 | 0  | 2  | 1  | 1  | 0  | 2  |
| 4     | 0  | 4  | 0  | 0  | 91 | 0  | 1  | 0  | 0  | 4  |
| 5     | 0  | 2  | 1  | 3  | 1  | 88 | 2  | 0  | 1  | 2  |
| 6     | 1  | 2  | 0  | 0  | 0  | 1  | 96 | 0  | 0  | 0  |
| 7     | 0  | 6  | 0  | 0  | 0  | 0  | 0  | 91 | 0  | 3  |
| 8     | 0  | 4  | 0  | 3  | 0  | 3  | 1  | 0  | 88 | 1  |
| 9     | 1  | 0  | 0  | 2  | 1  | 0  | 0  | 1  | 0  | 95 |

***Accuracy Rate on  $k = 7$ : 91.8%***

***Confusion Matrix for the Training Data – Case:  $k=9$***

| Index | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|----|----|----|----|----|----|----|----|----|
| 0     | 98 | 0  | 0  | 0  | 0  | 1  | 1  | 0  | 0  | 0  |
| 1     | 0  | 98 | 0  | 1  | 0  | 0  | 0  | 1  | 0  | 0  |
| 2     | 1  | 11 | 78 | 0  | 1  | 1  | 1  | 5  | 1  | 1  |
| 3     | 1  | 3  | 1  | 87 | 0  | 2  | 1  | 2  | 1  | 2  |
| 4     | 0  | 4  | 0  | 0  | 91 | 0  | 2  | 0  | 0  | 3  |
| 5     | 0  | 2  | 1  | 2  | 2  | 88 | 2  | 0  | 1  | 2  |
| 6     | 2  | 2  | 0  | 0  | 0  | 1  | 95 | 0  | 0  | 0  |
| 7     | 0  | 6  | 0  | 0  | 1  | 0  | 0  | 89 | 0  | 4  |
| 8     | 0  | 4  | 0  | 4  | 0  | 2  | 1  | 0  | 87 | 2  |
| 9     | 1  | 0  | 0  | 1  | 1  | 0  | 1  | 4  | 1  | 91 |

***Accuracy Rate on  $k = 9$ : 90.2%***

***Confusion Matrix for the Training Data – Case:  $k=11$***

| Index | 0  | 1  | 2  | 3  | 4  | 5  | 6  | 7  | 8  | 9  |
|-------|----|----|----|----|----|----|----|----|----|----|
| 0     | 97 | 0  | 0  | 0  | 0  | 2  | 1  | 0  | 0  | 0  |
| 1     | 0  | 98 | 0  | 1  | 0  | 0  | 0  | 1  | 0  | 0  |
| 2     | 1  | 12 | 78 | 0  | 2  | 1  | 0  | 4  | 1  | 1  |
| 3     | 1  | 3  | 1  | 87 | 0  | 2  | 1  | 2  | 1  | 2  |
| 4     | 0  | 4  | 0  | 0  | 91 | 0  | 2  | 0  | 0  | 3  |
| 5     | 0  | 2  | 1  | 4  | 1  | 87 | 1  | 0  | 2  | 2  |
| 6     | 2  | 3  | 0  | 0  | 0  | 1  | 94 | 0  | 0  | 0  |
| 7     | 0  | 5  | 0  | 0  | 1  | 0  | 0  | 89 | 0  | 5  |
| 8     | 0  | 4  | 0  | 3  | 0  | 2  | 1  | 0  | 88 | 2  |
| 9     | 2  | 0  | 0  | 1  | 3  | 0  | 1  | 2  | 0  | 91 |

***Accuracy Rate on  $k = 11$ : 90.0%***

## 4. Source Code

```
1. # -*- coding: utf-8 -*-
2. """
3. Created on Mon Nov  9 18:32:52 2020
4.
5. @author: User
6. """
7.
8. #hw2 with mnist dataset, image classifier
9.
10. import pandas as pd
11.
12. df_train = pd.read_csv("train1k.csv", sep=",")
13. df_test = pd.read_csv("test1k.csv", sep=",")
14.
15.
16. df_train.iloc[0, 0]
17.
18.
19. #Creating a dataframe in order to store the means of each pixels for each numbers in th
   e training data
20.
21.
22. pixel_sum = []
23. df_pixel_sum = pd.DataFrame(pixel_sum)
24. pixel_mean = []
25. df_pixel_mean = pd.DataFrame(pixel_mean)
26.
27. label_hold = []
28. k = 0
29. label_holder = df_train["label"]
30.
31. for k in range (0,10):
32.     df_pixel_sum[k] = df_train.loc[df_train["label"] == k ,:].sum()
33.     df_pixel_mean[k] = df_train.loc[df_train["label"] == k ,:].mean()
34.     k = k + 1
35.
36. t=0
37. indexer = []
38. for t in range(785):
39.     indexer.append(t)
40.
41. #Set mean pixel's index to numbers for easy access
42.
43. df_pixel_mean.set_index([indexer], inplace=True)
44.
45. #Nearest Mean Classifier
46.
47. #Right Now: Only for 0 Image & 1st Row of Test Data
48.
49. i = 1
50. pixel_difference = []
51. number_identififer = []
52. near_mc = []
53. #assigned_instance
54. df_assigned_instance = pd.DataFrame()
55. df_pixel_difference = pd.DataFrame(pixel_difference)
56. df_number_identififer = pd.DataFrame(number_identififer)
```

```

57. df_near_mc = pd.DataFrame(near_mc)
58. m = 0
59. k = 0
60. r = 0
61. s = 0
62. l = 0
63.
64. for r in range (0,1000):
65.     for m in range (0, 10): # for each mean of each number
66.         pixel_difference.clear()
67.         for i in range(1, 785): # for each pixel. cover all pixels
68.             pixel_difference.append((df_test.iloc[r,i] - df_pixel_mean.iloc[i, m])**2)
69.         #For the 0th image in test data, subtract and square each pixel's value from each pixel's mean for each number
70.         i = i + 1
71.         df_pixel_difference[m] = pixel_difference # Store the calculation for each image
72.
73.     #for k in range (0,10): # For each number's mean, find the distance between 0th test data
74.         number_identifier.append(df_pixel_difference.iloc[:, m].sum())
75.         df_number_identifier [m] = number_identifier
76.         m = m + 1
77.         #k = k + 1
78.         number_identifier.clear()
79.
80.     for s in range (0,10): # Finding the minimum distance between 0th test element and number image means
81.         near_mc.append(df_number_identifier.iloc[0,s])
82.         df_near_mc[s] = near_mc
83.         near_mc.clear()
84.         s = s + 1
85.     df_near_mc_transpose = df_near_mc.transpose()
86.
87.     assigned_instance = (df_near_mc_transpose.idxmin())
88.     df_assigned_instance[r] = assigned_instance
89.     #assigned_instance.clear()
90.
91. df_assigned_instance = df_assigned_instance.transpose()
92.
93.
94. # Self Test for the Accuracy
95.
96. q = 0
97. rank = 0
98.
99. for q in range (0,1000):
100.     if df_assigned_instance.iloc[q,0] == df_test.iloc[q,0]:
101.         rank = rank + 1
102.     print(rank)
103.
104.
105.     # Confusion matrix as 0s
106.
107.     confusion = [0] * 10
108.     df_confusion = pd.DataFrame(confusion)
109.     f = 0
110.     for f in range (1,10):
111.         df_confusion[f] = df_confusion[0]
112.

```

```

113.     # Create the Confusion Matrix
114.
115.     f = 0
116.     x = 0
117.     y = 0
118.
119.     for f in range(len(df_test)):
120.
121.         if df_assigned_instance.iloc[f,0] == df_test.iloc[f,0]:
122.             x = df_test.iloc[f,0]
123.             df_confusion.iloc[x,x] = df_confusion.iloc[x,x] + 1
124.
125.         elif df_assigned_instance.iloc[f,0] != df_test.iloc[f,0]:
126.             y = df_test.iloc[f,0]
127.             x = df_assigned_instance.iloc[f,0]
128.             df_confusion.iloc[y,x] = df_confusion.iloc[y,x] + 1
129.
130.
131.
132.
133.     # knn time
134.     knn = []
135.     df_knn = pd.DataFrame(knn)
136.
137.     dataframe_holder_list = {}
138.     k = 0
139.     m = 0
140.     n = 0
141.     y = 0
142.     distance_knn = []
143.     sum_holder = 0
144.     min_holder = 0
145.     minner = []
146.     summer = []
147.     trial = []
148.     df_summer = pd.DataFrame(summer)
149.     df_minner = pd.DataFrame(minner)
150.     df_trial = pd.DataFrame()
151.
152.
153.
154.     for k in range(0, 1000): #1000
155.         summer.clear()
156.         for t in range (1, 785):
157.             distance_knn.clear()
158.             for m in range (len(df_train)):
159.
160.                 distance_knn.append((df_test.iloc[k,t] - df_train.iloc[m, t])**2)
161.
162.                 df_knn[t] = distance_knn
163.
164.             for n in range(0,1000):
165.
166.                 sum_holder = df_knn.iloc[n, :].sum()
167.                 summer.append(sum_holder)
168.                 df_summer[k] = summer
169.
170.             #for y in range(0,1000):
171.
172.                 #min_holder = df_summer.iloc[:,k].idxmin()
173.                 #minner.append(min_holder)

```



```

174.         #df_minner[k] = minner
175.
176.
177.     df_trial = df_summer.copy()
178.     #df_trial = df_trial.sort_values(by = [0], ascending = True)
179.     #df_trial = df_trial.set_index(0)
180.     storella = []
181.     df_storella = pd.DataFrame(storella)
182.
183.     #storella.append(df_trial.index[0:11])
184.
185.     storella.clear()
186.     t = 0
187.     for t in range (0,1000): #1000
188.         df_trial = df_trial.sort_values(by = [t], ascending = True)
189.         storella.append(df_trial.index[0:11])
190.
191.
192.     #####
193.     df_storella = pd.DataFrame(storella)
194.     df_storella = df_storella.transpose()
195.
196.     x = 0
197.     y = 0
198.     u = 0
199.     o = 0
200.     knn_label = []
201.     df_knn_label = pd.DataFrame(knn_label)
202.     mode_holder = []
203.     df_mode_holder = pd.DataFrame(mode_holder)
204.
205.
206.     for x in range (0,1000): #1000
207.         knn_label.clear()
208.         for y in range (0,11):
209.             knn_label.append(df_train.iloc[df_storella.iloc[y,x],0])
210.             df_knn_label[x] = knn_label
211.
212.
213.     #mode
214.
215.     mode_holder.clear()
216.     p = 1
217.     t = 0
218.     o = 0
219.
220.     mode_holder_1 = []
221.     df_mode_holder_1 = pd.DataFrame(mode_holder_1)
222.
223.     for p in range (1,13,2):
224.         for o in range (0,1000): #1000, k = 11
225.             mode_holder.append(df_knn_label.iloc[:,p,o].mode())
226.
227.             #df_mode_holder [t] = mode_holder ###
228.
229.     df_mode_holder_1 = pd.DataFrame(mode_holder)
230.
231.
232.     # raw confusion matrix
233.
234.     knn_holder_list = []

```

```

235.
236.     knn_confusion = [0] * 10 # number??
237.     df_knn_confusion = pd.DataFrame(knn_confusion)
238.     f = 0
239.     for f in range (1,10): # range??
240.         df_knn_confusion[f] = df_knn_confusion[0]
241.
242.
243.
244.     #k - nn confusion matrix
245.
246.     d = 0
247.     s = 0
248.     key = 0
249.     x = 0
250.     u = 0
251.     #for b in range (1,7):
252.
253.     for s in range (1000, 2000):
254.
255.         if df_test.iloc[d, 0] == df_mode_holder_1.iloc[s,0]:
256.             x = df_test.iloc[d,0]
257.             df_knn_confusion.iloc[x,x] = df_knn_confusion.iloc[x,x] + 1
258.
259.         elif (df_test.iloc[d, 0]) != (df_mode_holder_1.iloc[s,0]):
260.             x = df_test.iloc[d,0]
261.             u = int(df_mode_holder_1.iloc[s,0])
262.             df_knn_confusion.iloc[x,u] = df_knn_confusion.iloc[x,u] + 1
263.             d = d + 1
264.             s = s + 1
265.         conf_holder_df = df_knn_confusion.copy()
266.         knn_holder_list.append(conf_holder_df)
267.         key = key + 1
268.
269.     r = 0
270.     z = 0
271.     m = 0
272.     e = 0
273.     nan_hunter = []
274.     for r in range (1000,2000):
275.         for z in range (1,5):
276.             if df_mode_holder_1.iloc[r,z] >= 0:
277.                 nan_hunter.append(r)
278.
279.
280.
281.     len(set(nan_hunter))
282.
283.     ""
284.
285.
286.     #mode
287.
288.     mode_holder.clear()
289.     p = 1
290.     t = 0
291.     o = 0
292.
293.     mode_holder_1 = []
294.     df_mode_holder_1 = pd.DataFrame(mode_holder_1)
295.

```

```

296.     for p in range (1,13,2):
297.         for o in range (0,1000): #1000, k = 11
298.             mode_holder.append(df_knn_label.iloc[:p,o].mode())
299.
300.             #df_mode_holder [t] = mode_holder ###
301.
302.     df_mode_holder_1 = pd.DataFrame(mode_holder)
303.
304.
305.     # raw confusion matrix
306.
307.     knn_holder_list = []
308.
309.     knn_confusion = [0] * 10 # number??
310.     df_knn_confusion = pd.DataFrame(knn_confusion)
311.     f = 0
312.     for f in range (1,10): # range??
313.         df_knn_confusion[f] = df_knn_confusion[0]
314.
315.
316.
317.     #k - nn confusion matrix
318.
319.     d = 0
320.     s = 0
321.     key = 0
322.     x = 0
323.     u = 0
324.     z=0
325.     #for b in range (1,7):
326.     for s in range (5000, 6000):
327.
328.         if df_test.iloc[d, 0] == df_mode_holder_1.iloc[s,0]:
329.             x = df_test.iloc[d,0]
330.             df_knn_confusion.iloc[x,x] = df_knn_confusion.iloc[x,x] + 1
331.
332.         elif (df_test.iloc[d, 0]) != (df_mode_holder_1.iloc[s,0]):
333.             x = df_test.iloc[d,0]
334.             u = int(df_mode_holder_1.iloc[s,0])
335.             df_knn_confusion.iloc[x,u] = df_knn_confusion.iloc[x,u] + 1
336.             d = d + 1
337.             s = s + 1
338.         conf_holder_df = df_knn_confusion.copy()
339.         knn_holder_list.append(conf_holder_df)
340.         key = key + 1
341.
342.     r = 0
343.     z = 0
344.     m = 0
345.     e = 0
346.     g = 0
347.     nan_hunter = []
348.     nan_2 = []
349.     saglama = 0
350.     saglama_2 = []
351.     for r in range (5000,6000):
352.         for z in range (1,5):
353.             if df_mode_holder_1.iloc[r,z] >= 0:
354.                 for g in range (0,5):
355.                     if df_mode_holder_1.iloc[r,g] == df_knn_label.iloc[0,r-
5000]: #Check!

```

```

356.                 saglama = saglama + 1
357.                 break
358.             else:
359.                 nan_2.append(r)
360.                 nan_hunter.append(r)
361.
362.         len(set(nan_hunter))
363.         nan_2_1 = set(nan_hunter)
364.
365.
366.         # Equal Vote Salvation: Problem Solved by manipulating modes
367.         # Do again now.
368.         p = 0
369.         pi = 0
370.         z = 0
371.         for p in range (5000, 6000): # do it for each k=3,5,7,9,11 for respective 1000-
            2000-3000 etc
372.             for z in range (1,5):
373.                 if df_mode_holder_1.iloc[p, z] >= 0: #hold p
374.                     if p in counter_list:
375.                         df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[0,p-5000]
376.                     elif p in counter_list_2:
377.                         df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[1,p-5000]
378.                     elif p in counter_list_3:
379.                         df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[2,p-5000]
380.                     elif p in counter_list_4:
381.                         df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[3,p-5000]
382.                     elif p in counter_list_5:
383.                         df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[4,p-5000]
384.                     elif p in counter_list_6:
385.                         df_mode_holder_1.iloc[p,0] = df_knn_label.iloc[5,p-5000]
386.
387.
388.
389.
390.
391.
392.
393.
394.         # nan_2_1 check
395.         nan_2_1 = list(nan_2_1)
396.         t = 0
397.         counter = 0
398.         counter_list = []
399.         counter_list_2 = []
400.         counter_list_3 = []
401.         counter_list_4 = []
402.         counter_list_5 = []
403.         counter_list_6 = []
404.
405.
406.
407.         for t in range (len(nan_2_1)):
408.             if nan_2_1[t] > 1000 and nan_2_1[t] < 2000:
409.                 if df_mode_holder_1.iloc[nan_2_1[t],0] == df_knn_label.iloc[0,nan_2_1[t]
-1000]:
410.                     counter = counter + 1
411.                     counter_list_3.append(nan_2_1[t])
412.                 elif nan_2_1[t] > 2000 and nan_2_1[t] < 3000:
413.                     if df_mode_holder_1.iloc[nan_2_1[t],4] == df_knn_label.iloc[5,nan_2_1[t]
-2000]:

```

```

414.         counter = counter + 1
415.         counter_list_6.append(nan_2_1[t])
416.         elif nan_2_1[t] > 3000 and nan_2_1[t] < 4000:
417.             if df_mode_holder_1.iloc[nan_2_1[t],4] == df_knn_label.iloc[5,nan_2_1[t]
-3000]:
418.                 counter = counter + 1
419.                 counter_list_6.append(nan_2_1[t])
420.                 elif nan_2_1[t] > 4000 and nan_2_1[t] < 5000:
421.                     if df_mode_holder_1.iloc[nan_2_1[t],4] == df_knn_label.iloc[5,nan_2_1[t]
-4000]:
422.                         counter = counter + 1
423.                         counter_list_6.append(nan_2_1[t])
424.                         elif nan_2_1[t] > 5000 and nan_2_1[t] < 6000:
425.                             if df_mode_holder_1.iloc[nan_2_1[t],4] == df_knn_label.iloc[5,nan_2_1[t]
-5000]:
426.                                 counter = counter + 1
427.                                 counter_list_6.append(nan_2_1[t])
428.
429.         len(set(counter_list_2))
430.
431.         in_counter_list = set(counter_list)
432.         in_counter_list_2 = set(counter_list_2)
433.         in_counter_list_3 = set(counter_list_3)
434.
435.         insecondnotfirst = in_counter_list_2 - in_counter_list
436.         insecondnotthird = in_counter_list_2 - in_counter_list_3
437.         infirstnotsecond = in_counter_list - in_counter_list_2
438.         infirstnotstthird = in_counter_list - in_counter_list_3
439.         inthirdnotfirst = in_counter_list_3 - in_counter_list
440.
441.         result_list =

```

## Image Visualization

```

1. from matplotlib import pyplot
2. import numpy as np
3. import csv
4. import matplotlib.pyplot as plt
5. k = 0
6. for k in (0,9):
7.
8.     label = df_pixel_mean_transpose[k]
9.
10.    pixels_mean = df_pixel_mean_transpose[1:]
11.
12.    pixels_mean = np.array(pixels_mean, dtype='uint8')
13.
14.    deneme = pixels_mean[:,k].reshape(28, 28)    #CHANGE k based on numbers
15.
16.    plt.title('Label is {label}'.format(label=label))
17.    plt.imshow(deneme, cmap='gray')
18.    plt.show()
19.    k = k + 1
20.
21.
22.
23. deneme = []

```

## **5. References**

1. Alpaydın, 2014.