# CSL7670 : Fundamentals of Machine Learning

## Lab Report

Name:               **Ganesh Kumar Nagal**
Roll Number:        **M23MEA004**
Program:            **M.Tech Adavance Manufacturing And Design**

2

# Chapter 1

# Lab-2

## 1.1 Objective

The objective of this lab assignment is to get familiarity with the K-Nearest Neighbors (K-NN) algorithm. By completing the problems and exercises in this lab.

## 1.2 Problem-1

1. (Apple vs Orange) You are given a K-NN code for the Apple vs Orange problem. Please read and understand the code. Now perform the following tasks:

1. Synthetically increase the dataset size to 50 samples.

2. Edit the code so that random 80%, 10%, and 10% samples are used for training, testing, and validation respectively.

3. Change the value of $K$ to 3, 5, and 7 and compare the validation set and test set results.

4. Write a code that draws confusion matrices for different $K$. Use the following link to understand about Confusion Matrix.
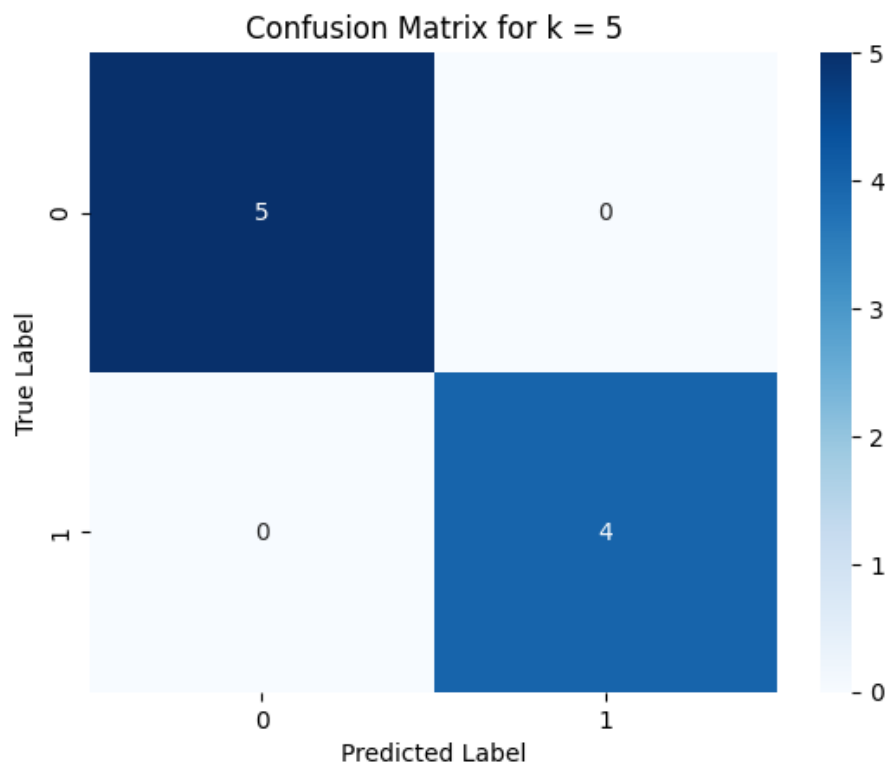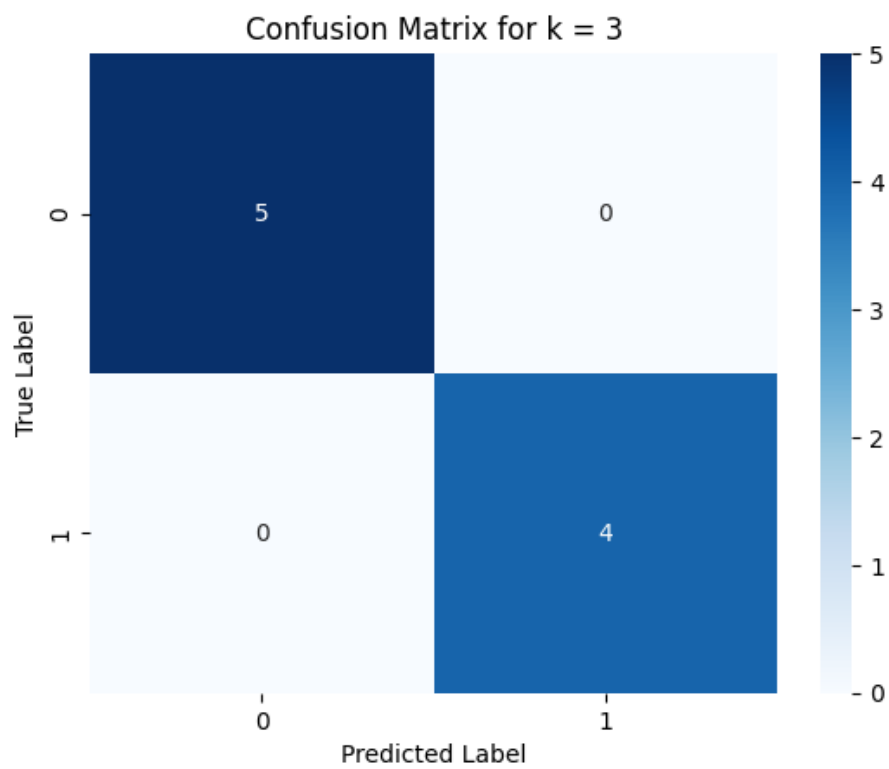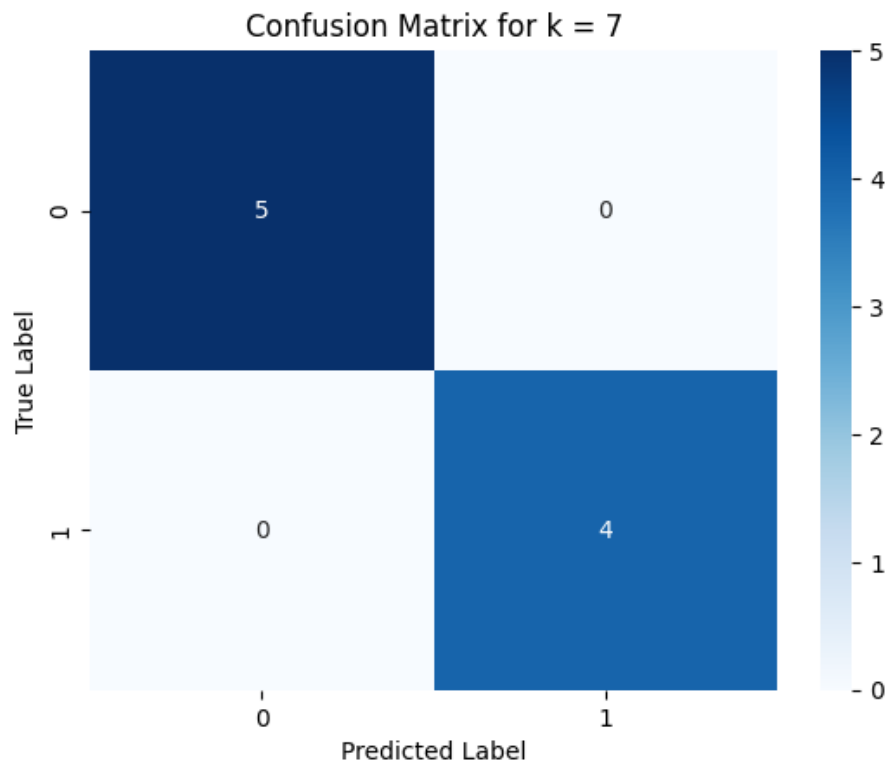
**Solution 1:**

```
1  # -*- coding: utf-8 -*-
2  """prob-1.py
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1
            ↪ Sp9P_8FqwjvNQWMe5biq5010OP1oDWXl
8  """
9
10 import numpy as np
11 import matplotlib.pyplot as plt
12 import seaborn as sns
13 from sklearn.model_selection import train_test_split
14 from sklearn.neighbors import KNeighborsClassifier
15 from sklearn.metrics import accuracy_score, confusion_matrix
16
17 # Sample synthetic dataset, 1=Apple, -1=Orange
18 data_apples = np.array([
```

```
19      [250 , 120 ,  1] ,
20      [240 , 110 ,  1] ,
21      [245 , 130 ,  1] ,
22      [240 , 100 ,  1] ,
23      [220 , 120 ,  1] ,
24      [230 , 125 ,  1] ,
25      [255 , 115 ,  1] ,
26      [215 , 110 ,  1] ,
27      [210 , 130 ,  1] ,
28      [235 , 105 ,  1] ,
29      [240 , 122 ,  1] ,
30      [245 , 125 ,  1] ,
31      [248 , 117 ,  1] ,
32      [243 , 127 ,  1] ,
33      [225 , 112 ,  1] ,
34      [230 , 118 ,  1] ,
35      [238 , 123 ,  1] ,
36      [222 , 115 ,  1] ,
37      [248 , 130 ,  1] ,
38      [210 , 100 ,  1] ,
39      [255 , 121 ,  1] ,
40      [245 , 128 ,  1] ,
41      [240 , 132 ,  1] ,
42      [235 , 120 ,  1] ,
43      [250 , 110 ,  1] ,
44      [243 , 135 ,  1] ,
45      [227 , 125 ,  1] ,
46      [222 , 132 ,  1] ,
47      [232 , 110 ,  1] ,
48      [240 , 105 ,  1] ,
49      [238 , 128 ,  1] ,
50      [215 , 118 ,  1] ,
51      [255 , 112 ,  1] ,
52      [250 , 125 ,  1] ,
53      [240 , 118 ,  1] ,
54      [223 , 130 ,  1] ,
55      [232 , 120 ,  1] ,
56      [235 , 130 ,  1] ,
57      [245 , 112 ,  1] ,
58      [228 , 123 ,  1] ,
59      [240 , 116 ,  1] ,
60      [220 , 128 ,  1] ,
61      [255 , 120 ,  1] ,
62      [230 , 130 ,  1] ,
63      [225 , 122 ,  1] ,
64      [240 , 126 ,  1] ,
65      [233 , 128 ,  1]
66  ])
67
68  data_oranges = np.array ([
69      [20 , 91 ,  -1] ,
70      [30 , 95 ,  -1] ,
71      [14 , 84 ,  -1] ,
72      [32 , 95 ,  -1] ,
```

```
 73        [18, 94, -1],
 74        [23, 90, -1],
 75        [25, 92, -1],
 76        [12, 80, -1],
 77        [29, 87, -1],
 78        [19, 89, -1],
 79        [22, 82, -1],
 80        [30, 88, -1],
 81        [16, 85, -1],
 82        [26, 89, -1],
 83        [28, 86, -1],
 84        [21, 83, -1],
 85        [24, 88, -1],
 86        [15, 81, -1],
 87        [17, 84, -1],
 88        [27, 91, -1],
 89        [23, 84, -1],
 90        [22, 88, -1],
 91        [20, 80, -1],
 92        [25, 87, -1],
 93        [18, 85, -1],
 94        [19, 82, -1],
 95        [28, 83, -1],
 96        [26, 81, -1],
 97        [30, 84, -1],
 98        [29, 82, -1],
 99        [21, 87, -1],
100        [12, 78, -1],
101        [27, 84, -1],
102        [16, 80, -1],
103        [24, 82, -1],
104        [14, 88, -1],
105        [17, 82, -1],
106        [22, 80, -1],
107        [23, 86, -1],
108        [13, 79, -1],
109        [19, 85, -1],
110        [21, 89, -1],
111        [15, 84, -1],
112        [26, 87, -1],
113        [25, 85, -1],
114        [27, 82, -1],
115        [28, 89, -1]
116 ])
117
118 # Concatenate apple and orange data
119 data = np.vstack((data_apples, data_oranges))
120
121 # Separate features (redness and weight) and labels (fruit type)
122 X = data[:, :2]
123 y = data[:, 2]
124
125 # Split dataset into training, testing, and validation sets
126 X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.2,
```

```
          ↪ random_state=42)
127  X_test, X_val, y_test, y_val = train_test_split(X_temp, y_temp, test_size
          ↪ =0.5, random_state=42)
128
129  # List of k values to try
130  k_values = [3, 5, 7]
131
132  for k in k_values:
133      print(f"Testing␣k␣=␣{k}")
134
135      # Create a KNN classifier with the current k value and L1 distance
              ↪ metric
136      knn_classifier = KNeighborsClassifier(n_neighbors=k, p=1)
137
138      # Train the classifier on the training data
139      knn_classifier.fit(X_train, y_train)
140
141      # Predict the fruit types for the test data
142      y_test_pred = knn_classifier.predict(X_test)
143
144      # Create a confusion matrix
145      cm = confusion_matrix(y_test, y_test_pred)
146
147      # Plot the confusion matrix
148      plt.figure()
149      sns.heatmap(cm, annot=True, fmt="d", cmap="Blues")
150      plt.title(f"Confusion␣Matrix␣for␣k␣=␣{k}")
151      plt.xlabel("Predicted␣Label")
152      plt.ylabel("True␣Label")
153      plt.show()
154
155      # Calculate the accuracy of the model on test data
156      accuracy_test = accuracy_score(y_test, y_test_pred)
157      print("Test␣Accuracy:", accuracy_test)
158      print()
```

Confusion Matrix for k = 3


Confusion Matrix for k = 5

Confusion Matrix for k = 7

```
1   OUTPUT:
2
3   Testing k = 3
4
5   Test Accuracy: 1.0
6
7   Testing k = 5
8
9   Test Accuracy: 1.0
10
11  Testing k = 7
12
13  Test Accuracy: 1.0
```

# 1.3   Problem-2

2. (Handwritten Digit Classification) Use the provided code for classifying handwritten digits from the MNIST dataset. First, read and understand the code. Then, perform the following tasks:

   1. Modify the code so that it uses L1-distance instead of the default L2-distance (Euclidean).

   2. Determine the value of $K$ that gives better performance.

   3. Report the accuracy achieved with the modified code.

   4. Display the results by showing the image, actual label, and predicted label. Identify some samples where the predicted label is incorrect.

**Solution 2:**

```
1  # -*- coding: utf-8 -*-
2  """prob-1b.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7      https://colab.research.google.com/drive/1VO0qxSQB624n4QwWHn3-
        ↪ nqJn1KDRPHxX
8  """
9
10 warnings.simplefilter(action='ignore', category=FutureWarning)
11 import numpy as np
12 from sklearn.datasets import fetch_openml
13 from sklearn.model_selection import train_test_split
14 from sklearn.neighbors import KNeighborsClassifier
15 from sklearn.metrics import accuracy_score
16 import matplotlib.pyplot as plt
17
18 # Load MNIST dataset
19 mnist = fetch_openml('mnist_784')
20 X = mnist.data
21 y = mnist.target
22
23 # Limit to a subset of the data
24 num_samples = 5000
25 X_subset = X[:num_samples]
26 y_subset = y[:num_samples]
27
28 # Split dataset into training and testing sets
29 X_train, X_test, y_train, y_test = train_test_split(X_subset, y_subset,
     ↪ test_size=0.2, random_state=42)
30
31 # Find the optimal value of K using cross-validation
32 best_accuracy = 0
33 best_k = 0
34 k_range = range(1, 6)  # Limiting the search range
35
36 for k in k_range:
37     knn_classifier = KNeighborsClassifier(n_neighbors=k, p=1)  # p=1
        ↪ indicates L1-distance
38     knn_classifier.fit(X_train, y_train)
39     y_pred = knn_classifier.predict(X_test)
40     accuracy = accuracy_score(y_test, y_pred)
41
42     if accuracy > best_accuracy:
43         best_accuracy = accuracy
44         best_k = k
45
46 print("Best␣K:", best_k)
47 print("Best␣Accuracy:", best_accuracy)
48
49 # Train the final classifier with the best K on the entire training data
50 final_knn_classifier = KNeighborsClassifier(n_neighbors=best_k, p=1)
51 final_knn_classifier.fit(X_train, y_train)
```

```
52
53  # Predict the labels for the test data
54  y_pred = final_knn_classifier.predict(X_test)
55
56  # Calculate the accuracy of the final model
57  accuracy = accuracy_score(y_test, y_pred)
58  print("Final Accuracy:", accuracy)
59
60  # Display images, actual labels, and predicted labels for incorrect
        ↪ predictions
61  incorrect_indices = np.where(y_pred != y_test)[0]
62
63  for index in incorrect_indices[:5]:  # Display the first 5 incorrect
        ↪ predictions
64      image = X_test[index, :].reshape(28, 28)
65      actual_label = y_test[index]
66      predicted_label = y_pred[index]
67
68      plt.figure(figsize=(4, 4))
69      plt.imshow(image, cmap='gray')
70      plt.title(f"Actual: {actual_label}, Predicted: {predicted_label}")
71      plt.axis('off')
72      plt.show()
```

```
1  OUTPUT:
2
3  Best K: 1
4  Best Accuracy: 0.939
5  Final Accuracy: 0.939
```