

Experiment No. : 01*Instructor:* Mr. Katkar Atish R.*Name:* Anjali Malav, *Roll number:* 17EEJCS002**AIM : Implementation of k - Nearest Neighbors (KNNs) on Synthetic data using Python****INTRODUCTION TO KNN**

K-Nearest Neighbor (KNN) algorithm is a type of supervised learning algorithm. KNN is used for both regression and classification tasks. When KNN is used for regression problems the prediction is based on the mean or the median of the K-most similar instances. Similarly, when KNN is used for classification, the output can be calculated as the class with the highest frequency from the K-most similar instances. Each instance in essence votes for their class and the class with the most votes is taken as the prediction. Basically, KNN is a prediction algorithm. The predictions are made by training dataset directly. Predictions are made for a new instance (x) by searching through the entire training set for the K most similar instances (the neighbors) and summarizing the output variable for those K instances. For regression this might be the mean output variable, in classification this might be the mode (or most common) class value.

To find out which of the K instances in the training dataset are most similar to a new input a distance measure is used. For real-valued input variables, the most popular distance measure is Euclidean distance. The Euclidean distance is represented by the following formula :

$$EuclideanDistance(x, x_i) = \sqrt{\sum (x_j - x_{ij})^2}$$

In above formula, the Euclidean distance is calculated as the square root of the sum of the squared differences between a new point (x) and an existing point (xi) across all input attributes j. The value for K can be found by algorithm tuning. It can be chosen by trying many different values for K (odd values) depending upon the size of the dataset and selected best value which suits for our problem.

KNN works well with a small number of input variables (p), but struggles when the number of inputs is very large. Each input variable can be considered a dimension of a p-dimensional input space. For example, if you had two input variables x1 and x2, the input space would be 2-dimensional. In high dimensions, points that may be similar may have very large distances. All points will be far away from each other and our intuition for distances in simple 2 and 3-dimensional spaces breaks down. This might feel unintuitive at first, but this general problem is called the “Curse of Dimensionality”.

PSEUDO CODE OF KNN

KNN can be implemented with the help of following steps :

1. Load the data
2. Initialise the value of k
3. For getting the predicted class, iterate from 1 to total number of training data points
 - (a) Calculate the Euclidean distance between test data and each row of training data.
 - (b) Sort the calculated distances in ascending order based on distance values
 - (c) Get top k rows from the sorted array
 - (d) Get the most frequent class of these rows
 - (e) Return the predicted class

DATASET DESCRIPTION

The dataset identified for this experimental study is Synthetic dataset. The dataset is broadly categorized

into 3 parts namely Linearly Separable, Non - Linearly Separable and Overlapping Data. These categories are further divided into several groups where in each group we are given with separate training and testing files. Testing and training data is also divided into different class files. Each class contain two features with many feature vectors and both features contain only numeric value.

IMPLEMENTATION CODE FOR KNN

```

1
2 # Package imported for different libraries
3 import math
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import operator
7 import csv
8 import random
9 from pandas import *
10
11
12 # Loading .txt data from computer using numpy
13 loaddata_train1=np.loadtxt('/home/asus/Desktop/prerana_ws/synthetic_data/synthetic_data/
14     nonlinearlySeparable/group4/class1_train.txt')
15 loaddata_train2=np.loadtxt('/home/asus/Desktop/prerana_ws/synthetic_data/synthetic_data/
16     nonlinearlySeparable/group4/class2_train.txt')
17 loaddata_test1=np.loadtxt('/home/asus/Desktop/prerana_ws/synthetic_data/synthetic_data/
18     nonlinearlySeparable/group4/class1_test.txt')
19 loaddata_test2=np.loadtxt('/home/asus/Desktop/prerana_ws/synthetic_data/synthetic_data/
20     nonlinearlySeparable/group4/class2_test.txt')
21 loaddata_train=np.concatenate((loaddata_train1, loaddata_train2), axis=0)
22 loaddata_test=np.concatenate((loaddata_test1, loaddata_test2), axis=0)
23
24
25 # Labeling training data
26 label1=np.ones((loaddata_train1.shape[0], 1))
27 label2=np.ones((loaddata_train2.shape[0], 1)) * -1
28 r1=np.append(label1, loaddata_train1, axis=1)
29 r2=np.append(label2, loaddata_train2, axis=1)
30 train_data=np.concatenate((r1,r2))
31
32
33 # Labeling testing data
34 label3=np.ones((loaddata_test1.shape[0], 1))
35 label4=np.ones((loaddata_test2.shape[0], 1)) * -1
36 r3=np.append(label3, loaddata_test1, axis=1)
37 r4=np.append(label4, loaddata_test2, axis=1)
38 test_data=np.concatenate((r3,r4))
39
40
41 # Plotting Data
42 plt.scatter(loaddata_train1[:,0], loaddata_train1[:,1], marker='o', label='Class1')
43 plt.scatter(loaddata_train2[:,0], loaddata_train2[:,1], marker='x', label='Class2')
44 plt.legend()
45 plt.show()
46
47 # Function to find Euclidean distances
48 def euclideanDistance(instance1, instance2, length):
49     distance = 0
50     for x in range(length):
51         distance += pow((float(instance1[x]) - float(instance2[x])), 2)
52     return math.sqrt(distance)
53
54
55 # Function to find neighbours by sorting Euclidean distances
56 def getKNeighbors(train_data, testInstance, k):
57     distances = []
58     length = len(testInstance)-1
59     for x in range(len(train_data)):
60         dist = euclideanDistance(testInstance, train_data[x], length)
61         distances.append((train_data[x], dist))
62     distances.sort(key=operator.itemgetter(1))
63     neighbors = []

```

```

60     for x in range(k):
61         neighbors.append(distances[x][0])
62     return neighbors
63
64
65 # Calculating label of neighbours and assign it to test instance
66
67
68 def getResponse(neighbors):
69     classVotes = {}
70     for x in range(len(neighbors)):
71         response = neighbors[x][0]
72         if response in classVotes:
73             classVotes[response] += 1
74         else:
75             classVotes[response] = 1
76     sortedVotes = sorted(classVotes.items(), key=operator.itemgetter(1), reverse=True)
77     return sortedVotes[0][0]
78
79
80 # Main function
81 def main():
82     l00 = 0
83     l01 = 0
84     l10 = 0
85     l11 = 0
86     split = 0.70
87     print('Train: ' + repr(len(train_data)))
88     print('Test: ' + repr(len(test_data)))
89     #generate predictions
90     predictions = []
91     k = input("Enter value of k: ")
92     k = int(k)
93     list_pred = []
94     list_act = []
95     for x in range(len(test_data)):
96         neighbors = getKNeighbors(train_data, test_data[x], k)
97         result = getResponse(neighbors)
98         predictions.append(result)
99         list_pred.append(result)
100         list_act.append(test_data[x][0])
101     for i in range(0, len(test_data)):
102         x = list_pred[i]
103         y = list_act[i]
104         if int(list_pred[i]) == -1 and int(list_act[i]) == -1:
105             l00 += 1
106         elif int(list_pred[i]) == -1 and int(list_act[i]) == 1:
107             l01 += 1
108         elif int(list_pred[i]) == 1 and int(list_act[i]) == -1:
109             l10 += 1
110         elif int(list_pred[i]) == 1 and int(list_act[i]) == 1:
111             l11 += 1
112     i+=1
113
114     a = np.array([[l00, l01],
115                  [l10, l11]])
116     print('Confusion Matrix: ')
117     print(DataFrame(a, columns = ['class_0', 'class_1'], index = ['class_0', 'class_1']))
118     prec_0 = (l00/float(l00+l10))
119     prec_1 = (l11/float(l01+l11))
120
121     acc = (l00+l11)*100/(l00+l01+l10+l11)
122     print('Accuracy: ' + repr(acc))
123
124     print('Precision: ')
125     print('Precision for class 0: ' + repr(prec_0))
126     print('Precision for class 1: ' + repr(prec_1))
127     print('Average Precision: ' + repr((prec_0+prec_1)/2))
128
129     rec_0 = (l00/float(l00+l01))
130     rec_1 = (l11/float(l10+l11))
131     print('Recall: ')
132     print('Recall for class 0: ' + repr(rec_0))

```

```

133 print('Recall for class 1: ' + repr(rec_1))
134 print('Average Recall: ' + repr((rec_0+rec_1)/2))
135
136 f0 = (2*(prec_0*rec_0)/(prec_0+rec_0))
137 f1 = (2*(prec_1*rec_1)/(prec_1+rec_1))
138 print('F1 Score: ')
139 print('F1 Score for class 0: ' + repr(f0))
140 print('F1 Score for class 1: ' + repr(f1))
141 print('Average F1 Score: ' + repr((f0+f1)/2))
142
143 main()

```

OUTPUT

1. Output for Linearly Separable Data

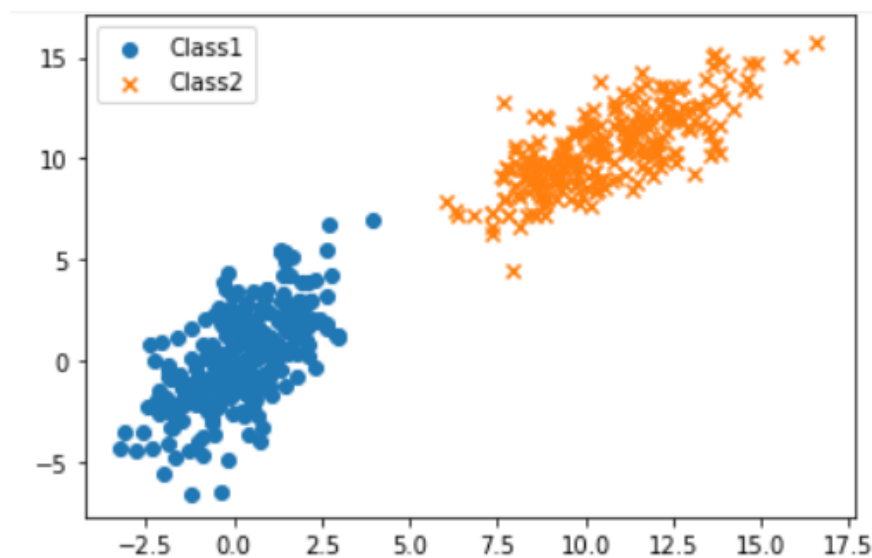


FIGURE 1: Plotting of Non-Linearly Separable Data

```

1 Anjali_Malav/Machine_Learning \$
2 Train: 500
3 Test: 200
4 Enter value of k: 9
5 Confusion Matrix:
6      class_0  class_1
7 class_0      100      0
8 class_1       0      100
9 Accuracy: 100.0
10 Precision:
11 Precision for class 0: 1.0
12 Precision for class 1: 1.0
13 Average Precision: 1.0
14 Recall:
15 Recall for class 0: 1.0
16 Recall for class 1: 1.0
17 Average Recall: 1.0
18 F1 Score:
19 F1 Score for class 0: 1.0
20 F1 Score for class 1: 1.0
21 Average F1 Score: 1.0

```

2. Output for Non-Linearly Separable Data

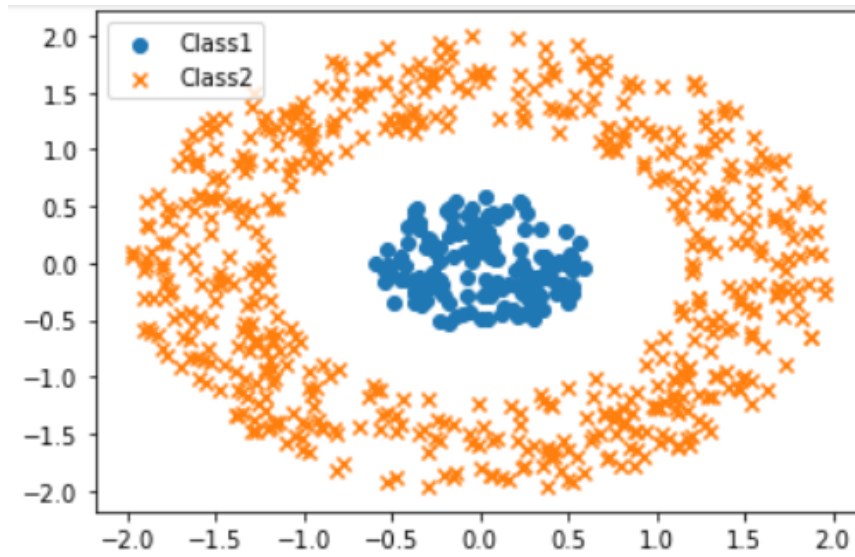


FIGURE 2: Plotting of Non-Linearly Separable Data

```

1 Anjali_Malav/Machine_Learning \$
2 Train: 750
3 Test: 300
4 Enter value of k: 7
5 Confusion Matrix:
6      class_0  class_1
7 class_0      240      0
8 class_1       0      60
9 Accuracy: 100.0
10 Precision:
11 Precision for class 0: 1.0
12 Precision for class 1: 1.0
13 Average Precision: 1.0
14 Recall:
15 Recall for class 0: 1.0
16 Recall for class 1: 1.0
17 Average Recall: 1.0
18 F1 Score:
19 F1 Score for class 0: 1.0
20 F1 Score for class 1: 1.0
21 Average F1 Score: 1.0

```

3. Output for Over-lapping Data

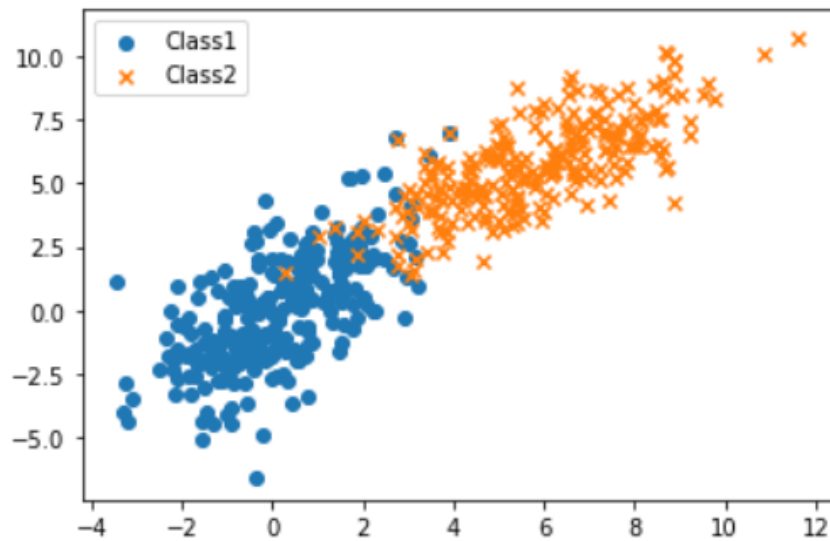


FIGURE 3: Plotting of Over-Lapping Data

```

1 Anjali_Malav/Machine_Learning \$
2 Train: 500
3 Test: 200
4 Enter value of k: 11
5 Confusion Matrix:
6      class_0  class_1
7 class_0      100      0
8 class_1       0      100
9 Accuracy: 100.0
10 Precision:
11 Precision for class 0: 1.0
12 Precision for class 1: 1.0
13 Average Precision: 1.0
14 Recall:
15 Recall for class 0: 1.0
16 Recall for class 1: 1.0
17 Average Recall: 1.0
18 F1 Score:
19 F1 Score for class 0: 1.0
20 F1 Score for class 1: 1.0
21 Average F1 Score: 1.0

```

CONCLUSION

The KNN algorithm outperforms for the above cases. The implementation shows that KNN gives **100%** testing accuracy for linearly separable, non-linearly separable and over-lapping data respectively for used dataset.