

# Designing a Minimum Distance to Class Mean Classifier

Mayeesha Humaira  
dept. Computer Science and Engineering  
Ahsanullah University of Science and Technology  
Dhaka, Bangladesh  
160204008@aust.edu  
section - A1

**Abstract**—The objective of this assignment is to design a minimum distance to class mean classifier.

**Index Terms**—decision boundary, class mean, linear discriminant function

## I. INTRODUCTION

Minimum distance to the class mean classifier is used to classify an unknown point using the linear discriminant function. When there is more than one class of data clusters minimum mean distance of each class is used to classify unknown points.

## II. EXPERIMENTAL DESIGN / METHODOLOGY

### Task:

Two-class set of prototypes have to be taken from “train.txt” and “test.txt” files.

- 1) Plot all sample points (train data) from both classes, but samples from the same class should have the same color and marker.
- 2) Using a minimum distance classifier with respect to ‘class mean’, classify the test data points by plotting them with the designated class-color but a different marker. Use the Linear Discriminant Function given below. Also, plot the class means.

$$g_i(X) = X^T \bar{Y}_i - \frac{1}{2} \bar{Y}_i^T \bar{Y}_i \quad (1)$$

- 3) Draw the decision boundary between the two classes.
- 4) Find accuracy.

Solution:

- 1) task-1 Plotting Training Data:

Given two-class set of prototypes in a file train.txt as Shown below:

```
2 2 1
3 1 1
-4 3 2
3 3 1
-1 -3 1
2 6 2
```

```
4 2 1
-2 -2 1
0 0 2
-2 2 2
-1 -1 2
-4 2 2
```

The above data are given in the train.txt file where the first column represents the x-coordinates, the second column represents the y-coordinates and the third column represents the class label.

Firstly, I read the data points from the train.txt file. Then using class labels data points were plotted using matplotlib. I employed different markers and colors for different classes. As in the given data there were two classes I used red dot markers to plot all data of class 1 and blue star marker to plot all data of class 2. The image of the plot is shown in Fig. 1.

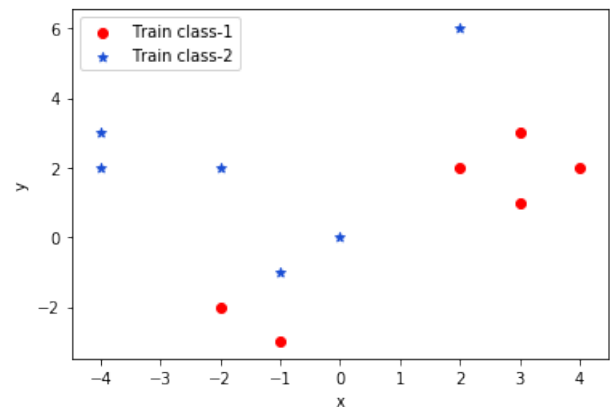


Fig. 1. Graphical representation of training data using matplotlib.

- 2) task-2 Classifying testing with respect to class mean and plotting testing data and mean:

To classify test data I first computed mean of each classes in training data. The mean of class 1 was [1.5, 0.5] and the mean of class 2 was [-1.5, 2.0]. The plot of mean of two classes along with training data using matplotlib is shown in Fig. 2 where red

triangle represents mean of class 1 and blue square represents mean of class 2. Then using the formula  $g_i(X) = X^T \bar{Y}_i - \frac{1}{2} \bar{Y}_i^T \bar{Y}_i$  where  $\bar{Y}_i$  is the mean of class i testing samples were classified. The testing samples were given in a test.txt file. The samples of test.txt file were:

```
-1 -5 1
3 2 1
-2 1 2
8 2 1
6 -1 1
0 2 1
-3 0 2
```

Only the first and the second column of this file were used as x-coordinate and y-coordinate respectively. If for a testing sample  $g_i > g_j$  where i and j represents different classes then that sample belongs to class i else if the testing sample has  $g_i < g_j$  where i and j represents different classes then it belongs to class j. After determining

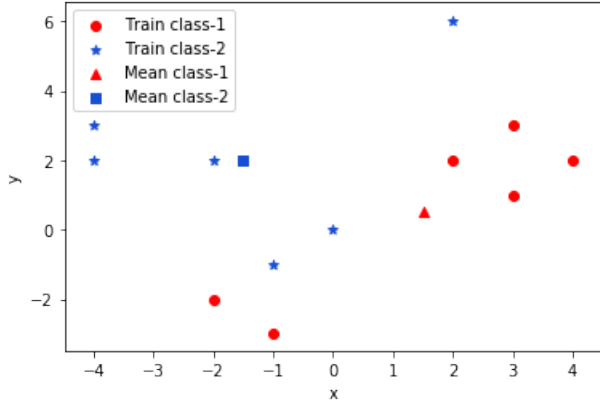


Fig. 2. Graphical representation of training data and mean of each class using matplotlib.

which class each test sample belongs to they were also plotted using matplotlib. Testing sample belonging to class 1 were plotted using red plus marker and samples of class 2 were plotted using blue cross marker. This plot is shown in Fig. 3.

- 3) task-3 Drawing Decision boundary between two classes: To draw the decision boundary I derived an equation using Linear discriminant function of two classes  $g_1$  and  $g_2$ . A point in class-1 is in decision boundary when:

$$g_1 = 0 \quad (2)$$

A point in class-2 is in decision boundary when:

$$g_2 = 0 \quad (3)$$

From Eq. 9 and Eq. 3 it can be said that:

$$g_1 = g_2 \quad (4)$$

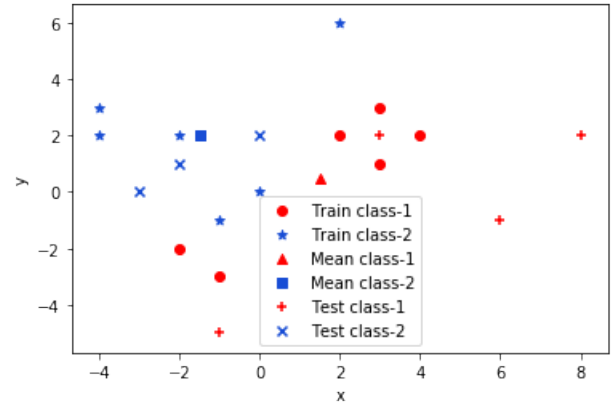


Fig. 3. Graphical representation of training data, mean of each class and classified testing data using matplotlib.

Substituting Eq. 1 in Eq. 4 we get:

$$\bar{\omega}_1^T X - \frac{1}{2} \bar{\omega}_1^T \bar{\omega}_1 = \bar{\omega}_2^T X - \frac{1}{2} \bar{\omega}_2^T \bar{\omega}_2 \quad (5)$$

$$(\bar{\omega}_1^T - \bar{\omega}_2^T) X - \frac{1}{2} \bar{\omega}_1^T \bar{\omega}_1 + \frac{1}{2} \bar{\omega}_2^T \bar{\omega}_2 \quad (6)$$

$$(\bar{\omega}_1^T - \bar{\omega}_2^T) X - \frac{1}{2} (\bar{\omega}_1^T \bar{\omega}_1 - \bar{\omega}_2^T \bar{\omega}_2) = 0 \quad (7)$$

where X in Eq. 7 is:

$$X = \begin{bmatrix} x \\ y \end{bmatrix} \quad (8)$$

Replacing Eq. 8 in Eq. 7 we get:

$$y = \frac{0.5(\bar{\omega}_1^T \bar{\omega}_1 - \bar{\omega}_2^T \bar{\omega}_2) + (\bar{\omega}_1^T - \bar{\omega}_2^T) x}{\bar{\omega}_1^T - \bar{\omega}_2^T} \quad (9)$$

where  $\bar{\omega}_i$  is the mean of class i,  $\bar{\omega}_{ix}$  is the x component of mean of class i and  $\bar{\omega}_{iy}$  is the y component of mean of class i.

Range of x-coordinates in training set was -4 to 4. x-coordinate can be iterated within this range to get the y-coordinates to draw the decision boundary using Eq. 9. I iterated x by 0.3 from range -4 to 4 to get corresponding y values. Training data, classified testing data, mean of each class and the decision boundary in a single plot is shown in Fig. 4.

### III. RESULT ANALYSIS

With respect to the decision boundary two training point were not correctly classified and one test sample was not correctly classified. Hence, I obtained training accuracy of 83.33% and test accuracy of 85.71%.

### IV. CONCLUSION

This algorithm is easy to implement. As the procedures and procedures are simple hence computations are faster. However, this algorithm has higher chance of misclassification as the decision boundary is linear.

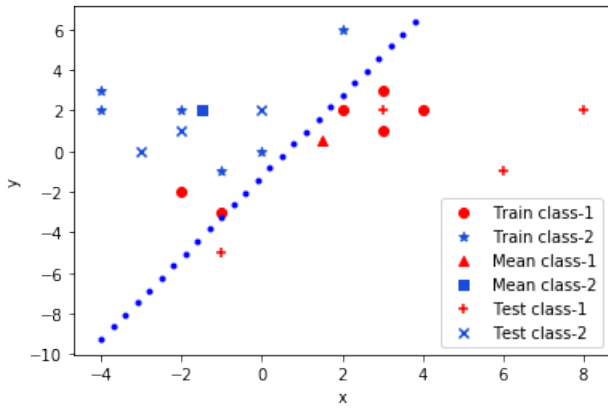


Fig. 4. Graphical representation of decision boundary along with training data, mean of each class and classified testing data using matplotlib.

## V. ALGORITHM IMPLEMENTATION / CODE

Code:

```
1 #Plotting training data
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import numpy as np
5 #loading training data
6 #Plotting training data
7 import matplotlib.pyplot as plt
8 import pandas as pd
9 import numpy as np
10
11 #Plotting training data
12 x,y,z = np.loadtxt('train.txt',unpack=True,
13     delimiter=' ')
14 plt.xlabel('x')
15 plt.ylabel('y')
16 for m in range(len(z)):
17     #print(m + 1, days[m])
18     if z[m]==1:
19         xc1=plt.scatter(x[m], y[m], color='r')
20     elif z[m]==2:
21         xc2=plt.scatter(x[m], y[m], marker='*',
22             color='#184DD5')
23 plt.legend([xc1, xc2], ["Train class-1", "Train
24     class-2"])
25 plt.show()
26
27 #x component of mean1
28 mlx=[]
29 for m in range(len(z)):
30     if z[m]==1:
31         mlx.extend([x[m]])
32 print(mlx)
33 ux1=np.mean(mlx)
34 print(ux1)
35 #y component of mean1
36 mly=[]
37 for m in range(len(z)):
38     if z[m]==1:
39         mly.extend([y[m]])
40 print(mly)
41 uy1=np.mean(mly)
42 print(uy1)
43 #x component of mean2
```

```
44 m2x=[]
45 for m in range(len(z)):
46     if z[m]==2:
47         m2x.extend([x[m]])
48 print(m2x)
49 ux2=np.mean(m2x)
50 print(ux2)
51 #y component of mean2
52 m2y=[]
53 for m in range(len(z)):
54     if z[m]==2:
55         m2y.extend([y[m]])
56 print(m2y)
57 uy2=np.mean(m2y)
58 print(uy2)
59
60 #making list of x y components of mean of two
61     classes
62 mean1=[ux1,uy1]
63 print(mean1)
64 mean2=[ux2,uy2]
65 print(mean2)
66
67 #plotting mean of two classes along with training
68     data
69 plt.xlabel('x')
70 plt.ylabel('y')
71 for m in range(len(z)):
72     if z[m]==1:
73         xc1=plt.scatter(x[m], y[m], color='r')
74     elif z[m]==2:
75         xc2=plt.scatter(x[m], y[m], marker='*',
76             color='#184DD5')
77 xm1=plt.scatter(ux1, uy1, marker='^', color='r')
78 xm2=plt.scatter(ux2,uy2, marker='s', color='#184DD5')
79 plt.legend([xc1, xc2, xm1, xm2], ["Train class-1", "
80     Train class-2", "Mean class-1", "Mean class-2"])
81 plt.show()
82
83 #classifying testing data using minimum distance to
84     class mean classifier
85 w = np.loadtxt('test.txt',unpack=False, delimiter='
86     ')
87 plt.xlabel('x')
88 plt.ylabel('y')
89 for m in range(len(z)):
90     if z[m]==1:
91         xc1=plt.scatter(x[m], y[m], color='r')
92     elif z[m]==2:
93         xc2=plt.scatter(x[m], y[m], marker='*',
94             color='#184DD5')
95 xm1=plt.scatter(ux1, uy1, marker='^', color='r')
96 xm2=plt.scatter(ux2,uy2, marker='s', color='#184DD5')
97
98 for i in range(len(w)):
99     g1=pd.Series(w[i, 0:2]).dot(pd.Series(mean1))
100     -0.5*pd.Series(mean1).dot(pd.Series(mean1))
101     g2=pd.Series(w[i, 0:2]).dot(pd.Series(mean2))
102     -0.5*pd.Series(mean2).dot(pd.Series(mean2))
103     if g1>g2:
104         xt1=plt.scatter(w[i,0], w[i,1], marker='+',
105             color='r')
106     elif g1<g2:
107         xt2=plt.scatter(w[i,0],w[i,1], marker='x',
108             color='#184DD5')
109 plt.legend([xc1, xc2, xm1, xm2,xt1,xt2], ["Train
110     class-1", "Train class-2", "Mean class-1", "Mean
111     class-2",
```

```

102         "Test
        class-1", "Test class-2"]])
103 plt.show()
104
105 #y intercept
106 minx=min(x)
107 print(minx)
108 maxx=max(x)
109 print(maxx)
110
111 for j in range(int(minx),int(maxx),1):
112     print(j)
113 s1=pd.Series(mean1).dot(pd.Series(mean1))
114 s2=pd.Series(mean2).dot(pd.Series(mean2))
115
116 intercept=0.5*(s1-s2)
117 print(intercept)
118 #x-coefficient and y coefficient
119 coeff=np.array(mean1)-np.array(mean2)
120 print(coeff)
121 print(coeff[1])
122 #list of range of x from -4 to 4
123 import decimal
124 def float_range(start, stop, step):
125     while start < stop:
126         yield float(start)
127         start += decimal.Decimal(step)
128
129 rangex=list(float_range(int(minx), int(maxx), '0.3')
130 )
131 print(rangex)
132
133 #classifying testing data using minimum distance to
134 #class mean classifier
135 w = np.loadtxt('test.txt',unpack=False, delimiter='
136 ')
137 plt.xlabel('x')
138 plt.ylabel('y')
139 pred=[]
140 for m in range(len(z)):
141     #print(m + 1, days[m])
142     if z[m]==1:
143         plt.scatter(x[m], y[m], color='r')
144     elif z[m]==2:
145         plt.scatter(x[m], y[m], marker='*', color='
146 #081F59')
147
148 plt.scatter(ux1, uy1, marker='^', color='r')
149 plt.scatter(ux2,uy2, marker='s', color='#081F59')
150 for i in range(len(w)):
151     g1=pd.Series(w[i, 0:2]).dot(pd.Series(mean1))
152     -0.5*pd.Series(mean1).dot(pd.Series(mean1))
153     g2=pd.Series(w[i, 0:2]).dot(pd.Series(mean2))
154     -0.5*pd.Series(mean2).dot(pd.Series(mean2))
155     if g1>g2:
156         plt.scatter(w[i,0], w[i,1], marker='+',
157 color='r')
158         pred.extend([1])
159     elif g1<g2:
160         plt.scatter(w[i,0],w[i,1], marker='x', color=
161 '#081F59')
162         pred.extend([2])
163
164 for j in range(len(rangex)):
165     yval=-(coeff[0]*(rangex[j])+intercept)/coeff[1]
166     plt.plot(rangex[j],yval, '.b')
167
168 print(pred)
169
170 #Testing Accuracy
171 cnt=0
172 for k in range(len(w)):
173     if pred[k]==w[k][2]:
174         cnt=cnt+1
175
176 print(cnt)
177 print(len(w))
178 test_acc=cnt/len(w)*100
179 print("Test accuracy: ", test_acc, "%")
180
181 #Training Accuracy
182 train_acc=10/12*100
183 print("Train accuracy: ", train_acc, "%")

```