

Implementing Minimum Error Rate Classifier

Mayisha Farzana

dept.Computer Science and Engineering
Ahsanullah University of Science and Technology
Dhaka, Bangladesh
160204028@aust.edu

Abstract—This experiment aims to identify specific sample points using the posterior probabilities used to measure the probability probabilities by Gaussian distribution. The goal of this classifier form is to decrease the rate of error during classification. So this classifier takes decisions based on the most posterior probabilities. This classifier is also known as the Bayes classifier with minimum error.

Index Terms—Discriminant functions, pattern recognition, likelihood probabilities ratio, Bayesian classifier, posterior probability

I. INTRODUCTION

A classifier is a minimal error rate classifier, and its function is to decrease the error rate. We are given six sample details in this experiment; we have to identify those. The normal distribution gives a sample's probability probabilities. With two parameters, sigma and mean, any normal distribution can be expressed. Those parameters are given in this experiment. As Bayesian classifier works with posterior probabilities the decision rule is as follows :

If $p(w_1|x) > p(w_2|x)$ then $x \in w_1$

If $p(w_1|x) < p(w_2|x)$ then $x \in w_2$

The posterior probabilities can be calculated with the help of likelihood probabilities. This can be written as:

$$\begin{aligned} P(w_i|x) &= P(x|w_i).P(w_i) \\ \Rightarrow \ln P(w_i|x) &= \ln P(x|w_i).P(w_i) \\ \Rightarrow \ln P(w_i|x) &= \ln P(x|w_i) + \ln P(w_i) \end{aligned}$$

This equation gives the discriminant function for minimum error rate classifier. In here $P(x|w_i)$ is likelihood probabilities and $P(w_i)$ is prior probabilities.

As our dataset is in 2D, so we have to use the following equations:

$$N_k(\mathbf{x}_i | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}_k|}} e^{\left(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_k)\right)}$$

Here, N_k is normal distribution, $\boldsymbol{\Sigma}$ is co-variance matrix, $\boldsymbol{\mu}_k$ is mean, x_i is test data which is a vector, $d=2$ for our experiment because all the data are 2D. We will plot all the test points in these equations, and from this equation we will get the likelihood probability values. After getting the values, if we multiple the values with prior probabilities, we will get our

required value of posterior probability. Posterior probability = likelihood probability * prior probability.

And the decision boundary would be the solution of

$$g_1(x) = g_2(x)$$

$$\Rightarrow p(w_1|x) = p(w_2|x)$$

$$\Rightarrow p(w_1|x) - p(w_2|x) = 0$$

$$\Rightarrow P(x|w_1).P(w_1) - P(x|w_2).P(w_2) = 0$$

Taking ln,

$$\Rightarrow \ln P(x|w_1).P(w_1) - \ln P(x|w_2).P(w_2) = 0$$

$$\Rightarrow \ln P(x|w_1) + \ln P(w_1) - \ln P(x|w_2) - \ln P(w_2) = 0$$

$$\Rightarrow \ln P(x|w_1)/\ln P(x|w_2) - \ln P(w_2)/\ln P(w_1) = 0$$

This is the equation of a decision boundary for minimum error classifier.

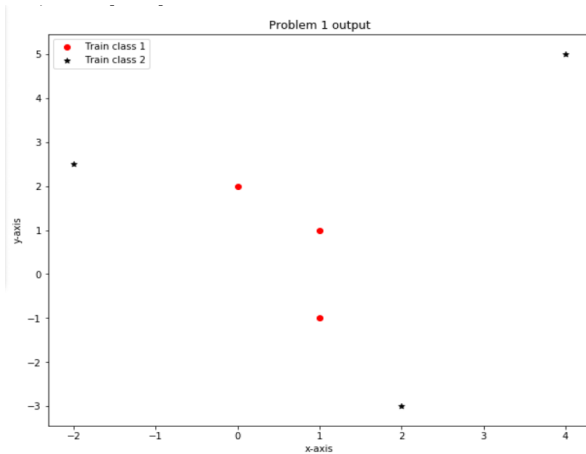
II. EXPERIMENTAL DESIGN / METHODOLOGY

A. Classifying the sample points: At first, we will read the test.txt file. From this file, we will get the x vector values. As we can see that, the classes are not classified in the text dataset. So we need to classify it.

B. Plotting the sample points with different markers: In this part, for classifying the dataset class, we need to calculate the posterior probability values. We know that,

$$P(w_i|x) = P(x|w_i).P(w_i)$$

posterior probability = likelihood probability * prior probability
As we know that, we can get the likelihood values from normal distribution. We will calculate the normal distribution values. Using this values and prior values, we can easily get the posterior probability values. From this values, we need to apply decision rule. We calculate the value of $g(x)$ for each sample point with the two given Gaussian distributions and check for the following conditions $g_1(x) > g_2(x)$. If the above condition is true, then the sample point, x, belongs to the Gaussian distributions' corresponding regions. The value of $g_1(x)$ greater than $g_2(x)$ means the sample point likelihood probabilities close to the used Gaussian distribution to assign this sample point to that region. After plotting all the sample values the output is as follows-



Now that we are plotting all the samples to their intended class, we need to draw a decision limit to split the whole space into two regions.

C. Drawing decision boundary in Contour Plot To draw the decision boundary we have to obtain the equation of the decision boundary $g_1(x) - g_2(x) = 0$. We will draw the decision boundary in contour plotting. To draw the contour plotting, at first we need to take the X,Y dimension to a higher dimension. We need to distribute our 2-dimensional distribution over variables X and Y. We will pack the X and Y dimension into a 3 dimension data. From this data, we can plot the values in contour plotting. For Z values, we will pass the values in multivariate distribution. We will calculate the Z values for different class. From it, we will measure the decision boundary $db=Z-Z1$. After that, we will plot the decision boundary line in the contour plotting. This is how we will draw the decision boundary.

3D surface with 2D contour plot projections

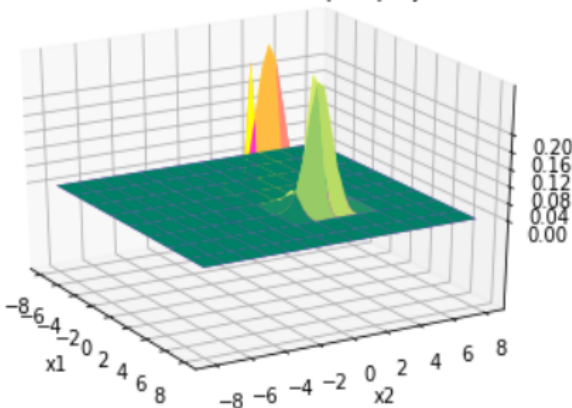


Fig: Contour plotting with probability probability distribution function

3D surface with 2D contour plot projections

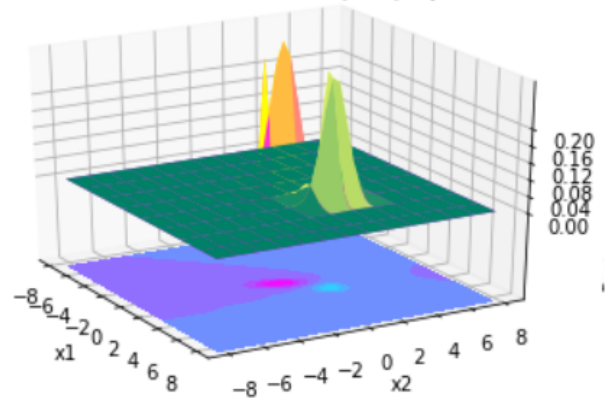


Fig: Contour plotting with decision boundary

III. RESULT ANALYSIS

As we know that minimum error rate classifier tries to minimize the error. So if we change the parameter of Gaussian distribution the sample values also can be shifted to another class because the likelihood probabilities can also be shifted towards another class.

IV. CONCLUSION

In this experiment, we came to understand how a minimum error rate classifier functions and what it means by sigma and mean of Gaussian distribution. There are, however, some drawbacks of this classifier. This classifier relies solely on probability and should be understood regarding the Gaussian distribution.

V. ALGORITHM IMPLEMENTATION / CODE

```
1 '''Problem 1 '''
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 train_dataset = np.loadtxt('test-Minimum-Error-Rate-
6 Classifier.txt', delimiter=",", dtype='float64')
7 print(train_dataset)
8 length=len(train_dataset)
9 print(length)
10 import numpy as np
11 def normal_distribution(x, u, cov):
12     k = len(u) # scalar dimension
13     uu = u.copy() #mu
14     xx = x.copy() #x
15     t1 = (2 * np.pi)**k # scalar (2pie)^k
16     t2 = np.linalg.det(cov) # scalar det(
17     t3 = 1.0 / np.sqrt(t1 * t2) # scalar
18     t4 = np.transpose(xx - uu) # (x-mu)T
19     t5 = np.linalg.inv(cov) # inverse(covariance)
20     t6 = (xx - uu) # (x-mu)
21     t7 = -0.5 * (np.dot(t4,t5).dot(t6))
22     result = t3 * np.exp(t7) # 1x1
23     return result
24 #Initialize
25 mu1=np.array([0,0]).reshape(2,)
26 mu2=np.array([2,2]).reshape(2,)
27 cov1=np.array([.25,.3,.3,1]).reshape(2,2)
28 cov2=np.array([.5,0,0,.5]).reshape(2,2)
```

```

30 p_w1=0.5
31 p_w2=0.5
32
33 x_class1=[]
34 y_class1=[]
35 x_class2=[]
36 y_class2=[]
37 for x in train_dataset:
38     posterior1=p_w1*normal_distribution(np.array([x
39         [0],x[1]]),mu1 , cov1)
40     posterior2=p_w2*normal_distribution(np.array([x
41         [0],x[1]]),mu2 , cov2)
42     print( posterior1, posterior2)
43     if( posterior1> posterior2):
44         print("class 1")
45         x_class1.append(x[0])
46         y_class1.append(x[1])
47     else:
48         print("class 2")
49         x_class2.append(x[0])
50         y_class2.append(x[1])
51 fig,ax=plt.subplots()#to show it in the same figure
52 plt.title("Problem 1 output")
53 plt.xlabel('x-axis', color='black')
54 plt.ylabel('y-axis', color='black')
55 ax.scatter(x_class1,y_class1,marker='o',color='r',
56     label='Train class 1')
57 ax.scatter(x_class2,y_class2,marker='*',color='black',
58     label='Train class 2')
59 fig.set_figheight(8)
60 fig.set_figwidth(10)
61 ax.legend()#show the output figure
62 '''Contour Plotting'''
63
64 import numpy as np
65 import matplotlib.pyplot as plt
66 from matplotlib import cm
67 from mpl_toolkits.mplot3d import Axes3D
68 # Our 2-dimensional distribution will be over
69 variables X and Y
70
71 N = 40
72 X = np.linspace(-8, 8, N)#it divides the -8 to 8
73 range in to 40 piece
74 Y = np.linspace(-8, 8, N)
75 X, Y = np.meshgrid(X, Y)#Make N-D coordinate arrays
76 for vectorized evaluations of N-D scalar
77 #print('X',X)
78 # Mean vector and covariance matrix
79 mu1 = np.array([0., 0.])
80 cov1 = np.array([[.25 , 0.3], [.3, 1.]])
81 mu2 = np.array([2.,2.])
82 cov2 = np.array([[.5 , 0.], [0., .5]])
83 # Pack X and Y into a single 3-dimensional array
84 pos = np.empty(X.shape + (2,))#taking X shape(40,40)
85 and one additional dimension 2D points
86
87 print(pos.shape)
88 pos[:, :, 0] = X
89 pos[:, :, 1] = Y
90
91 def multivariate_gaussian(pos, mu, cov):
92     """Return the multivariate Gaussian distribution
93     on array pos."""
94     dimen = mu.shape[0]
95     cov_det = np.linalg.det(cov)
96     cov_inv = np.linalg.inv(cov)
97     N = np.sqrt((2*np.pi)**dimen * cov_det)
98     # This einsum call calculates (x-mu)T.Sigma-1.(x
99     -mu) in a vectorized
100     # way across all the input variables.
101     ans= np.einsum('...a,ab,...b->...', pos-mu,
102         cov_inv, pos-mu)#EINSTEIN SUMMATION
103
104     return np.exp(-ans / 2) / N
105
106 Z = multivariate_gaussian(pos, mu1, cov1)
107 Z1 = multivariate_gaussian(pos, mu2, cov2)
108
109 # Create a surface plot and projected filled contour
110 plot under it.
111 fig = plt.figure()
112 ax = fig.gca(projection='3d')
113 #db=Z-Z1#decision boundary
114 z=0
115 ax.set_xlabel('x1')
116 ax.set_ylabel('x2')
117 ax.set_zlabel('Probability density')
118 ax.scatter(x_class1,y_class1, z,color='red')
119 ax.scatter(x_class2,y_class2, z,color='blue')
120 #print(len(pos),pos.shape)
121 ax.plot_surface(X, Y, Z, rstride=3, cstride=5,
122     linewidth=2, antialiased=True,
123     cmap=cm.spring)
124 ax.plot_surface(X, Y, Z1, rstride=3, cstride=5,
125     linewidth=2, antialiased=True,
126     cmap=cm.summer)
127 #ax.contourf(X, Y, db, zdir='z', offset=-0.22, cmap=
128     cm.cool)
129 ax.set_title('3D surface with 2D contour plot
130     projections')
131 # Adjust the limits, ticks and view angle
132 ax.set_zlim(-0.2,0.3)
133 ax.set_zticks(np.linspace(0,0.2,6))
134 ax.view_init(25, -30)
135 plt.show()
136
137 '''Drawing Decision Boundary'''
138 import numpy as np
139 import matplotlib.pyplot as plt
140 from matplotlib import cm
141 from mpl_toolkits.mplot3d import Axes3D
142 # Our 2-dimensional distribution will be over
143 variables X and Y
144
145 N = 40
146 X = np.linspace(-8, 8, N)#it divides the -8 to 8
147 range in to 40 piece
148 Y = np.linspace(-8, 8, N)
149 X, Y = np.meshgrid(X, Y)#Make N-D coordinate arrays
150 for vectorized evaluations of N-D scalar
151 #print('X',X)
152 # Mean vector and covariance matrix
153 mu1 = np.array([0., 0.])
154 cov1 = np.array([[.25 , 0.3], [.3, 1.]])
155 mu2 = np.array([2.,2.])
156 cov2 = np.array([[.5 , 0.], [0., .5]])
157 # Pack X and Y into a single 3-dimensional array
158 pos = np.empty(X.shape + (2,))#taking X shape(40,40)
159 and one additional dimension 2D points
160
161 print(pos.shape)
162 pos[:, :, 0] = X
163 pos[:, :, 1] = Y
164
165 def multivariate_gaussian(pos, mu, cov):
166     """Return the multivariate Gaussian distribution
167     on array pos."""
168     dimen = mu.shape[0]
169     cov_det = np.linalg.det(cov)
170     cov_inv = np.linalg.inv(cov)
171     N = np.sqrt((2*np.pi)**dimen * cov_det)
172     # This einsum call calculates (x-mu)T.Sigma-1.(x
173     -mu) in a vectorized
174     # way across all the input variables.
175     ans= np.einsum('...a,ab,...b->...', pos-mu,
176         cov_inv, pos-mu)#EINSTEIN SUMMATION
177
178     return np.exp(-ans / 2) / N
179
180 Z = multivariate_gaussian(pos, mu1, cov1)
181 Z1 = multivariate_gaussian(pos, mu2, cov2)
182
183 # Create a surface plot and projected filled contour
184 plot under it.
185 fig = plt.figure()
186 ax = fig.gca(projection='3d')
187 db=Z-Z1#decision boundary
188 z=0

```

```
154 ax.set_xlabel('x1')
155 ax.set_ylabel('x2')
156 ax.set_zlabel('Probability density')
157 ax.scatter(x_class1,y_class1, z,color='red')
158 ax.scatter(x_class2,y_class2, z,color='blue')
159 #print(len(pos),pos.shape)
160 ax.plot_surface(X, Y, Z, rstride=3, cstride=5,
161               linewidth=2, antialiased=True,
162               cmap=cm.spring)
163 ax.plot_surface(X, Y, Z1, rstride=3, cstride=5,
164               linewidth=2, antialiased=True,
165               cmap=cm.summer)
166 ax.contourf(X, Y, db, zdir='z', offset=-0.22, cmap=
167             cm.cool)
168 ax.set_title('3D surface with 2D contour plot
169             projections')
170 # Adjust the limits, ticks and view angle
171 ax.set_zlim(-0.2,0.3)
172 ax.set_zticks(np.linspace(0,0.2,6))
173 ax.view_init(25, -30)
174 plt.show()
```