

Bios/CS 534 Project 1

Chenxi Cai

February 21, 2018

1 Problem 1

Class Y = 0:

mean vector:

$$\begin{bmatrix} 1.06242164 & 1.61910524 \end{bmatrix} \quad (1)$$

covariance matrix:

$$\begin{bmatrix} 4.79170095 & 0.90180838 \\ 0.90180838 & 1.2945715 \end{bmatrix} \quad (2)$$

Class Y = 1:

mean vector:

$$\begin{bmatrix} 1.13915258 & -1.18380439 \end{bmatrix} \quad (3)$$

covariance matrix:

$$\begin{bmatrix} 0.7560476 & -0.5093068 \\ -0.5093068 & 3.19387164 \end{bmatrix} \quad (4)$$

1.1 Equal prior case

Classification error rate on the testing data is

$$0.057692 \quad (5)$$

Scatter plots with boundary of the training data is shown as Figure 1.
The s in Class y = 0 and y = 1 is red "o" and blue "x", separately.

1.2 Prior calculated from the data case

Classification error rate on the testing data is

$$0.057692 \quad (6)$$

Scatter plots with boundary of the training data is shown as Figure 2

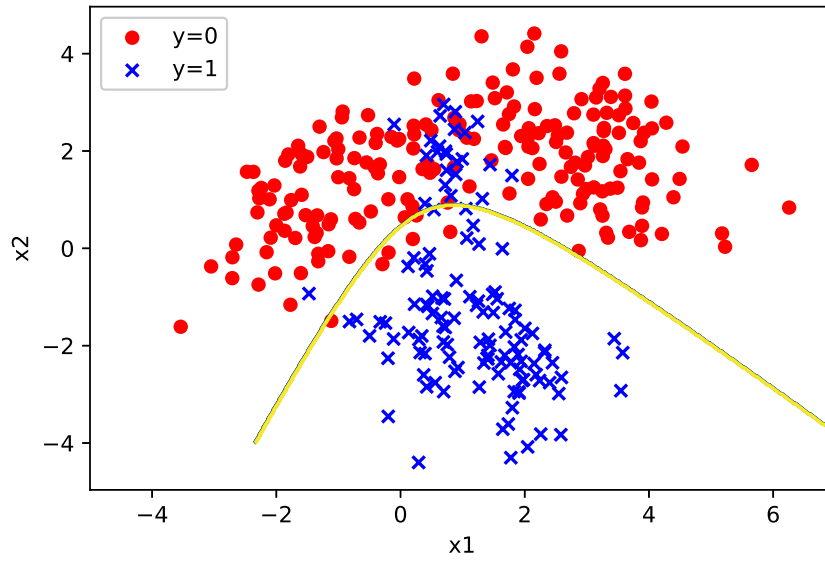


Figure 1: Scatter plots with boundary of equal prior

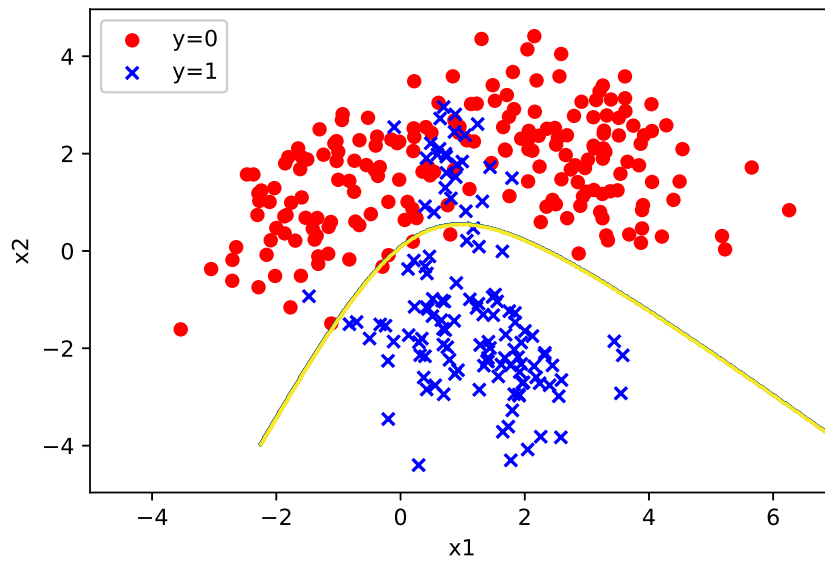


Figure 2: Scatter plots with boundary of prior calculated from the data

2 Problem 2

2.1 Bandwidth=10

Classification error rate on the testing data is

$$0.184615 \quad (7)$$

Scatter plots of the testing data is shown as Figure 3.

Right points are black and misclassified points are red. The shape of points in Class $y = 0$ and $y = 1$ is "o" and "x", separately.

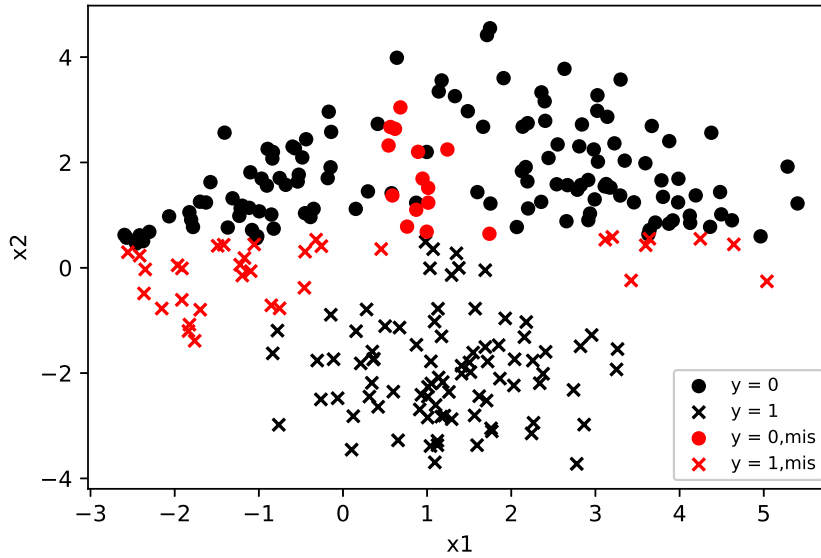


Figure 3: Scatter plots of bandwidth=10

2.2 Bandwidth=1

Classification error rate on the testing data is

$$0.053846 \quad (8)$$

Scatter plots of the testing data is shown as Figure 4.

2.3 Bandwidth=0.1

Classification error rate on the testing data is

$$0.038462 \quad (9)$$

Scatter plots of the training data is shown as Figure 5.

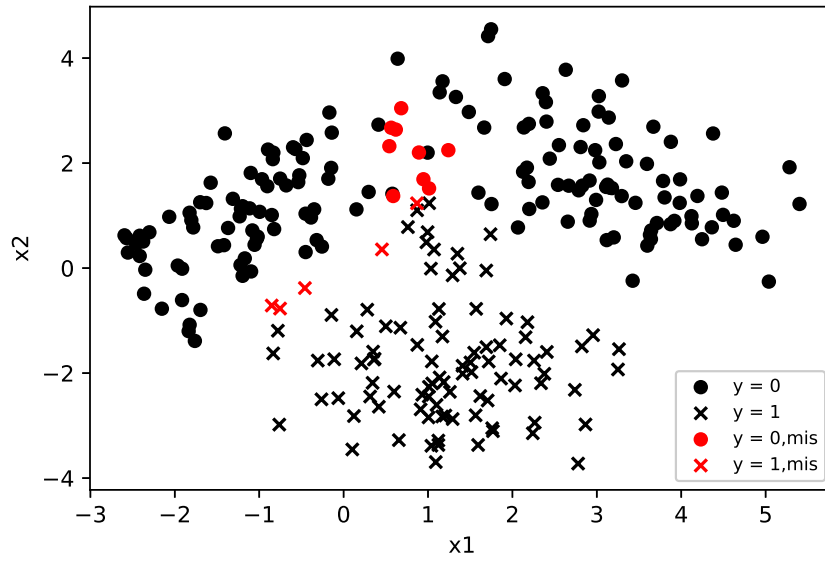


Figure 4: Scatter plots of bandwidth=1

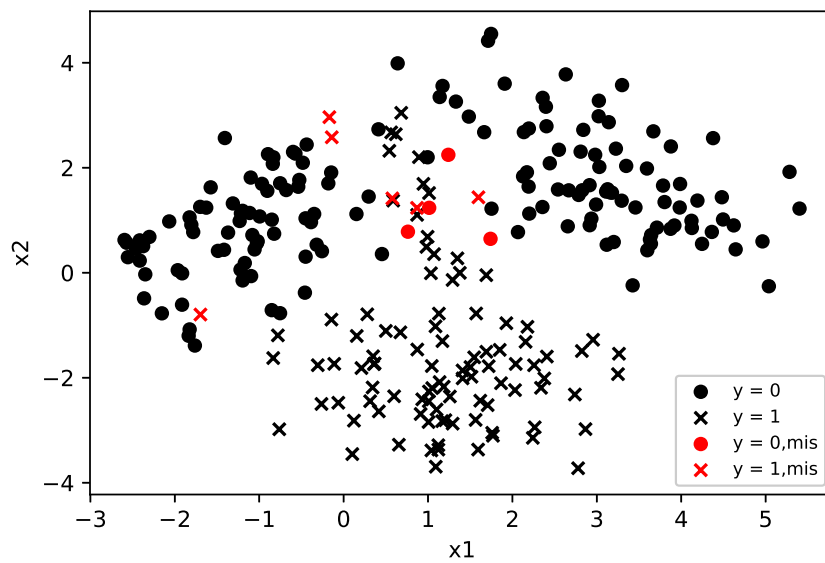


Figure 5: Scatter plots of bandwidth=0.1

3 Problem 3

3.1 K=1

sensitivity = 0.920000, specificity=0.975000, false discovery rate=0.041667

Scatter plots of the testing data is shown as Figure 6. Right points are black and misclassified points are red. The shape of points in Class $y = 0$ and $y = 1$ is "o" and "x", separately.

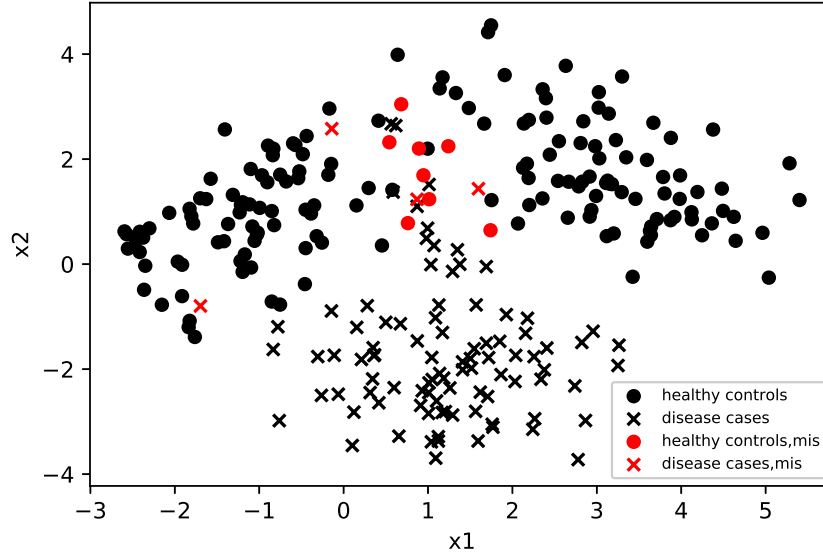


Figure 6: Scatter plots of K=1

3.2 K=5

sensitivity = 0.960000, specificity=0.981250, false discovery rate=0.030303

Scatter plots of the testing data is shown as Figure 7.

3.3 K=10

sensitivity = 0.950000, specificity=0.975000, false discovery rate=0.040404

Scatter plots of the testing data is shown as Figure 8.

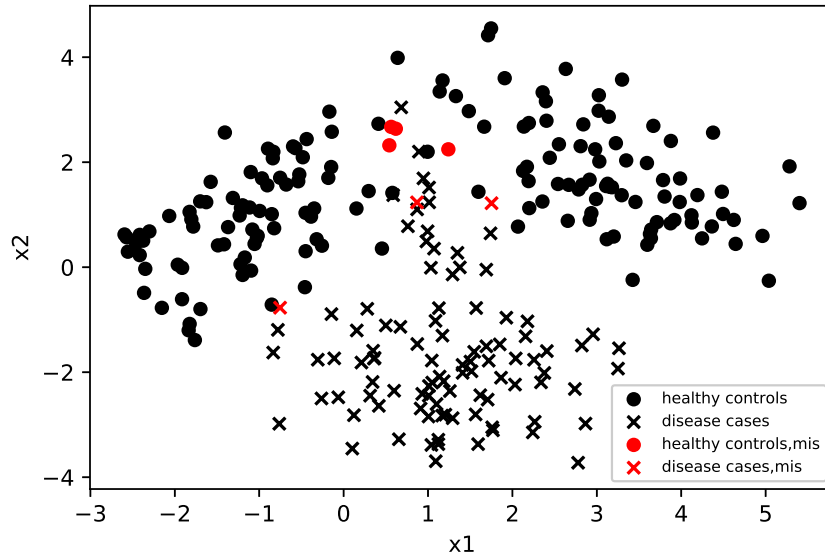


Figure 7: Scatter plots of K=5

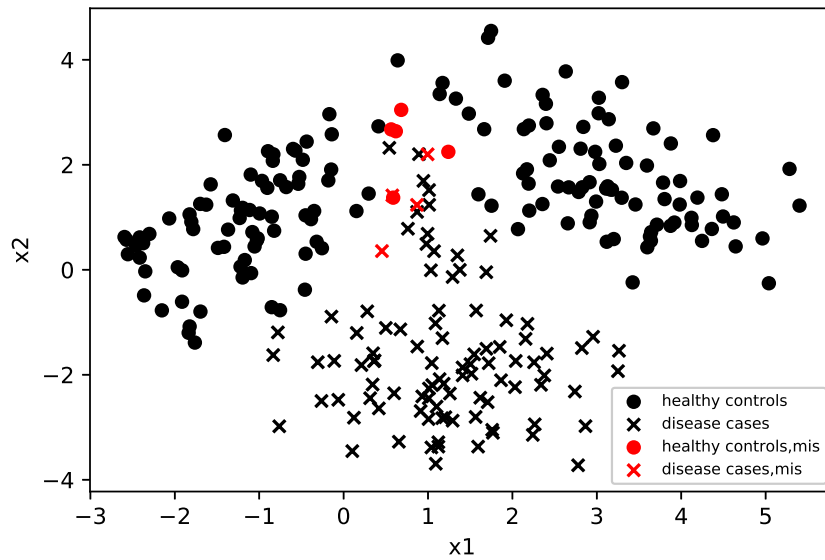


Figure 8: Scatter plots of K=10

Python codes of problem1:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5
6 #define functions of parameters and decision
7 def discrtninant_func(f,y):
```

```

8   mean_vectors = np.zeros((1,2))      #mean vectro and covariance matrix
9   covariance = np.zeros((2,2))
10  mean_vectors = np.mean(f.loc[lambda f : f[2] == y,[0,1]].values, axis=0)
11  covariance = np.cov(f.loc[lambda f : f[2] == y,[0,1]].values.T)
12  return mean_vectors, covariance
13
14 def decision_function(x,x_m,x_c,prior_probability):
15     sig_inv = np.linalg.inv(x_c)      #decision function
16     return -.5*np.log( np.linalg.det(x_c))-.5*np.dot(np.dot((x-x_m),sig_inv),(
17     x-x_m).T) + np.log(prior_probability)
18
19 def predict(X,m1,c1,m2,c2,frac1,frac2):
20     s = X.shape[0]
21     new_Y = np.zeros((s,1))
22     X_new = np.zeros((s,2))
23     j = 0
24     l = s-1
25     for i in range(s):
26         k = []
27         k.append(decision_function(X[i,:],m1,c1,frac1))    #result of decision
28         k.append(decision_function(X[i,:],m2,c2,frac2))    #result of decision
29         #classify
30         new_Y[i] = np.argmax(k)    #find the index that has bigger results(
31         results in which class is the biggest)
32         if np.argmax(k) == 0:
33             X_new[j,:] = X[i,:]    #rearrange testing points accroding to
34             predict class
35             j += 1
36         elif new_Y[i] == 1:
37             #the predict class of points in testing
38             file
39             X_new[l,:] = X[i,:]
40             l -= 1
41
42     return X_new,new_Y,j,l
43
44 def Error_rate(Y,new_Y):
45     s = Y.shape[0]
46     n = 0
47     for i in range(s):
48         if (Y[i] != new_Y[i]):    #error_rate = #points whose predict class
49             is different form original one/#total points
50             n += 1
51     return n/float(s)
52
53 def plot(x,m1,c1,m2,c2,frac1,frac2):
54     #build grid of xx1,xx2
55     xx1,xx2 = np.meshgrid(np.arange(-5, 7,0.02 ),
56     np.arange(-4, 4,0.02))
57     #z is the new_Y of xx1 and xx2 field
58     Z = predict(np.array([xx1.ravel(), xx2.ravel()]).T,m1,c1,m2,c2,frac1,frac2
59     )[1]
60
61     #change Z to a matrix with the same shape of xx1
62     Z = Z.reshape(xx1.shape)
63     #draw the boundary line
64     plt.contour(xx1, xx2, Z, alpha=0.3)
65
66     #plot points
67     plt.scatter(x[:200,0], x[:200,1],30,
68     color='red', marker='o', label='y=0')
69     plt.scatter(x[200:,0], x[200:,1],30,
70     color='blue', marker='x', label='y=1')
71     plt.xlabel('x1')
72     plt.ylabel('x2')
73     plt.legend(loc='upper left')
74     return 0

```

```

68
69
70 #read data
71 f1 = pd.read_csv('/Users/ccai28/Desktop/HW_1_training.txt',sep='\t', header =
    None, skiprows = 1)
72 f2 = pd.read_csv('/Users/ccai28/Desktop/HW_1_testing.txt',sep='\t', header =
    None, skiprows = 1)
73
74 #mean vector and covariance matrix
75 x = f1.loc[:, [0,1]].values
76 m1,c1 = discrtminant_func(f1,0)
77 m2,c2 = discrtminant_func(f1,1)
78 print m1,m2,c1,c2
79 #read testing data
80 X = f2.loc[:, [0,1]].values
81 Y = f2.loc[:, [2]].values
82
83 #equal prior
84 X_new_e,new_Y_e,j_e,l_e = predict(X,m1,c1,m2,c2,0.5,0.5)
85 error_rate_e = Error_rate(Y,new_Y_e)
86 print('Equal prior: \nerror rate = {:.f} \n'.format(error_rate_e))
87 plot(x,m1,c1,m2,c2,0.5,0.5)
88 plt.savefig('Equal prior .eps',format='eps')
89 plt.show()
90
91 #from data
92 X_new_d,new_Y_d,j_d,l_d = predict(X,m1,c1,m2,c2,200/float(325),125/float(325))
93 error_rate_d = Error_rate(Y,new_Y_d)
94 print('Prior calculated from the data: \nerror rate = {:.f} \n'.format(
    error_rate_e))
95 plot(x,m1,c1,m2,c2,200/float(325),125/float(325))
96 plt.savefig('Prior calculated from the data.eps',format='eps')
97 plt.show()

```

Python codes of problem2:

```

1 from sklearn.neighbors.kde import KernelDensity
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5
6 def kernel(X,x,bandwidth):
7     kde_1 = KernelDensity(kernel='gaussian', bandwidth=bandwidth).fit(x)
8     d = 10**((kde_1.score_samples(X)) #using scleaner to calculate density p(x)
9     return d
10
11 def predict(d1,d2,frac):
12     sum = d1*frac+d2*frac
13     k = []
14     k.append(d1*frac/float(sum)) #calculate posterior according to (density
    * prior)/evidence
15     k.append(d2*frac/float(sum))
16     #classify
17     new_y = np.argmax(k) #find the index with 1 -> predict value
18     return new_y
19
20 def Error_rate(Y,new_Y):
21     s = Y.shape[0]
22     n = 0
23     for i in range(s):
24         if (Y[i] != new_Y[i]): #error-rate = #points of testing
            predict to another class/#totla points
25             n += 1
26     return n/float(s)
27
28 def misclassified_point(X,Y,new_Y):
29     mis_0 = []

```



```

30 mis_1 = []
31 right_0 = []
32 right_1 = []
33 s = Y.shape[0]
34 for i in range(s):
35     if (Y[i] != new_Y[i] and new_Y[i] == 0):           #find original class 1
36         training point -> predict class 0
37         mis_0.append(X[i,:])
38         elif (Y[i] != new_Y[i] and new_Y[i] == 1):       #find original class 0
39             training point -> predict class 1
40             mis_1.append(X[i,:])
41             elif (Y[i] == new_Y[i] and new_Y[i] == 0):    #find original class 0
42                 training point -> predict class 0
43                 right_0.append(X[i,:])                    #find original class 1
44                 training point -> predict class 1
45             else:
46                 right_1.append(X[i,:])
47 mis_0 = np.matrix(mis_0)
48 right_0 = np.matrix(right_0)
49 mis_1 = np.matrix(mis_1)
50 right_1 = np.matrix(right_1)
51 return mis_0, mis_1, right_0, right_1
52
53 def plot(X, x1, x2, bandwidth, Y, frac):
54
55     d1 = kernel(X, x1, bandwidth)
56     d2 = kernel(X, x2, bandwidth)
57
58     #predict
59     s = X.shape[0]
60     new_Y = np.zeros((s,1))
61     for i in range(s):
62         new_Y[i] = predict(d1[i], d2[i], frac)           #calculate predict value for
63         every point in testing file
64
65     #error_rate
66     print('When bandwidth={:f}, \nerror rate = {:f} \n'.format(bandwidth,
67         Error_rate(Y, new_Y)))
68
69     #mis points
70     mis_0, mis_1, right_0, right_1 = misclassified_point(X, Y, new_Y)
71
72     plt.scatter(right_0[:,0], right_0[:,1], 30,
73         color='black', marker='o', label='y = 0')
74     plt.scatter(right_1[:,0], right_1[:,1], 30,
75         color='black', marker='x', label='y = 1')
76     plt.scatter(mis_0[:,0], mis_0[:,1], 30,
77         color='red', marker='o', label='y = 0, mis')
78     plt.scatter(mis_1[:,0], mis_1[:,1], 30,
79         color='red', marker='x', label='y = 1, mis')
80
81     plt.xlabel('x1')
82     plt.ylabel('x2')
83     plt.legend(loc='lower right',
84         fontsize = 7)
85     plt.savefig('bandwidth={:f}.eps'.format(bandwidth), format='eps')
86     plt.show()
87     return 0
88
89 #read files
90 f1 = pd.read_csv('/Users/ccai28/Desktop/HW_1_training.txt', sep='\t', header =
91     None, skiprows = 1)
92 f2 = pd.read_csv('/Users/ccai28/Desktop/HW_1_testing.txt', sep='\t', header =
93     None, skiprows = 1)
94
95 #read testing data
96 X = f2.loc[:, [0,1]].values
97 Y = f2.loc[:, [2]].values

```

```

90
91 #read traing data and kernel density estimator
92 # y = 0
93 x1 = f1.loc[lambda f : f[2] == 0,[0,1]].values
94 #y = 1
95 x2 = f1.loc[lambda f : f[2] == 1,[0,1]].values
96
97 plot(X,x1,x2,10,Y,0.5)
98 plot(X,x1,x2,1,Y,0.5)
99 plot(X,x1,x2,0.1,Y,0.5)

```

Python codes of problem3:

```

1 import numpy as np
2 import pandas as pd
3 import math
4 import matplotlib.pyplot as plt
5 from matplotlib.colors import ListedColormap
6
7 def distance(x,y):
8     d = math.sqrt((x[0]-y[0])**2+(x[1]-y[1])**2) #distance between testing
point and train point
9     return d
10
11 def KNN(X,x,y,K):
12     s = x.shape[0]
13     S = X.shape[0]
14     new_Y = np.zeros((S,1))
15     for i in range(S):
16         d = np.zeros((s,1))
17         for j in range(s):
18             d[j] = distance(X[i,:],x[j,:]) #distance for all the points
respect with one point i in testing file
19             index = np.argsort(d,axis = 0) #sort according to the distance and
index include index of distance
20             count = 0 #from smallest to longest
21             for k in range(K): #find k nearest points in training respect to
testing i
22                 y_1 = y[index[k]] #find the class of testing point i with its
index
23                 if y_1 == 0: #if the nearest test points are in class 0, and
more than half of these points are in
24                     count += 1 # class 0, the i training point is ->predict
into class 0
25                 if count > K*0.5: # else ->predict into class 1
26                     new_y = 0
27                 else:
28                     new_y = 1
29                 new_Y[i] = new_y #build the new class of traing points after
classified
30
31     return new_Y
32
33 def misclassified_point(X,Y,new_Y):
34     mis_0 = []
35     mis_1 = []
36     right_0 = []
37     right_1 = []
38     s = Y.shape[0]
39     for i in range(s):
40         if (Y[i] != new_Y[i] and new_Y[i] == 0): #find original class 1
training point -> predict class 0
41             mis_0.append(X[i,:])
42         elif (Y[i] != new_Y[i] and new_Y[i] == 1): #find original class 0
training point -> predict class 1
43             mis_1.append(X[i,:])
44         elif (Y[i] == new_Y[i] and new_Y[i] == 0): #find original class 0

```

```

45     training_point -> predict class 0
        right_0.append(X[i,:]) #find original class 1
46     training_point -> predict class 1
        else:
47         right_1.append(X[i,:])
48     mis_0 = np.matrix(mis_0)
49     right_0 = np.matrix(right_0)
50     mis_1 = np.matrix(mis_1)
51     right_1 = np.matrix(right_1)
52     return mis_0, mis_1, right_0, right_1
53
54 def plot(X,Y,x,y,K):
55     new_Y = KNN(X,x,y,K)
56
57     #mis points
58     mis_0, mis_1, right_0, right_1 = misclassified_point(X,Y,new_Y)
59
60
61     plt.scatter(right_0[:,0], right_0[:,1], 30,
62                 color='black', marker='o', label='healthy controls')
63     plt.scatter(right_1[:,0], right_1[:,1], 30,
64                 color='black', marker='x', label='disease cases')
65     plt.scatter(mis_0[:,0], mis_0[:,1], 30,
66                 color='red', marker='o', label='healthy controls, mis')
67     plt.scatter(mis_1[:,0], mis_1[:,1], 30,
68                 color='red', marker='x', label='disease cases, mis')
69
70     plt.xlabel('x1')
71     plt.ylabel('x2')
72     plt.legend(scatterpoints=1,
73                loc='lower right',
74                ncol=1,
75                fontsize=7)
76     plt.savefig('k={:d}.eps'.format(K), format='eps')
77     plt.show()
78
79     return 0
80
81 def parameters(X,Y,x,y,K):
82     new_Y = KNN(X,x,y,K)
83     right0 = [] #True Negative, TN
84     right1 = [] #True Positive, TP
85     mis0 = [] #class 0 -> predict class 1 #False Negative, FN
86     mis1 = [] #class 1 -> predict class 0 #False Positive, FP
87
88     for i in range(Y.shape[0]):
89         if np.allclose(new_Y[i], 0):
90             # Class 0
91             if np.allclose(new_Y[i], Y[i]): #Correct class0
92                 right0.append(new_Y[i])
93             else:
94                 mis0.append(new_Y[i])
95         else:
96             if np.allclose(new_Y[i], Y[i]):
97                 right1.append(new_Y[i])
98             else:
99                 mis1.append(new_Y[i])
100     TN = len(right0)
101     TP = len(right1)
102     FN = len(mis0)
103     FP = len(mis1)
104     sensitivity = TP / float((TP + FN))
105     specificity = TN / float((TN + FP))
106     false_discovery_rate = FP / float((FP + TP))
107
108     print('When k={:d}, \nsensitivity = {:f}, \nspecificity={:f}, \nfalse
        discovery rate={:f} \n'.format(K, sensitivity, specificity,
        false_discovery_rate))

```

```

109
110
111
112
113 #read data
114 f1 = pd.read_csv('/Users/ccai28/Desktop/HW_1_training.txt',sep='\t', header =
      None, skiprows = 1)
115 f2 = pd.read_csv('/Users/ccai28/Desktop/HW_1_testing.txt',sep='\t', header =
      None, skiprows = 1)
116
117 #read testing data
118 X = f2.loc[:, [0,1]].values
119 Y = f2.loc[:, [2]].values
120
121 #read traing data
122 x = f1.loc[:, [0,1]].values
123 y = f1.loc[:, [2]].values
124
125 #parameters,k=1
126 parameters(X,Y,x,y,1)
127 #plot,k=1
128 plot(X,Y,x,y,1)
129
130 #k=5
131 parameters(X,Y,x,y,5)
132 plot(X,Y,x,y,5)
133
134 #k=10
135 parameters(X,Y,x,y,10)
136 plot(X,Y,x,y,10)

```