

Homework2

Full name: Rong Zhen zID: z5225226

Q1: PartA

Decision Tree Results

Dataset	Default	0%	25%	50%	75%	
australian	56.52% (2)	81.16% (7)	86.96% (2)	56.52% (2)	20.77% (7)	
labor	66.67% (2)	94.44% (7)	44.44% (7)	66.67% (7)	50.00% (12)	
diabetes	66.23% (2)	67.10% (7)	64.07% (12)	66.23% (2)	35.50% (27)	
ionosphere	66.04% (2)	86.79% (7)	82.08% (27)	71.70% (7)	18.87% (12)	

PartB

---4

PartC

--2

Q2:

Part A:

Answer: Accuracy for training set: 0.8969404186795491

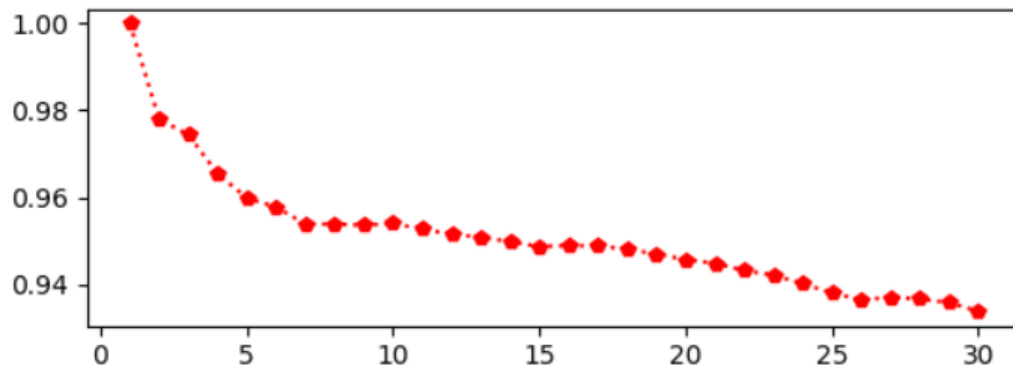
Accuracy for test set: 0.7681159420289855

Part B:

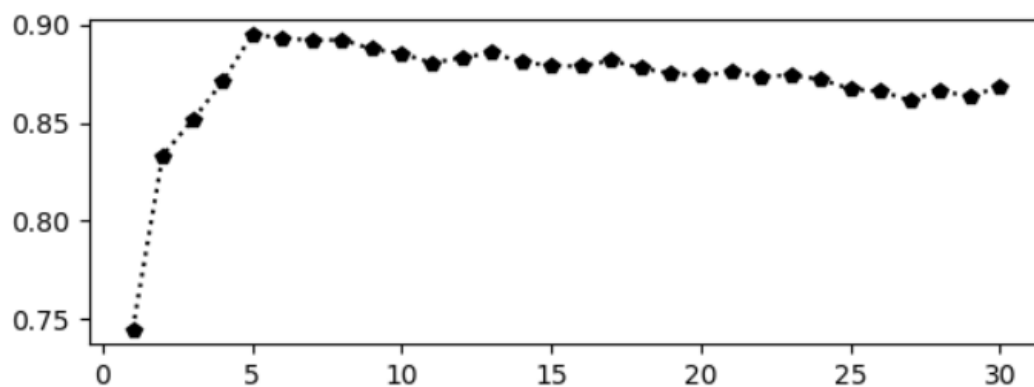
Answer: The optimal number of k is 5.

Part C:

Plot for training set:



Plot for test set:



Part D:

Precision and recall for k=5 is:(0.7666666666666667, 0.8518518518518519)

precision and recall for k=2 is:(0.7894736842105263, 0.5555555555555556)

Compare: We can see the precision scores of these two models are almost the same. But when k=5, the recall score is much larger. So the model with k=5 is better in general.

```
1. import csv
2. import math
3. import numpy as np
4. import matplotlib.pyplot as plt
5. from sklearn.metrics import roc_auc_score
6. from sklearn.metrics import recall_score, precision_score
7.
8. x = [[],[],[],[],[],[],[],[],[],[],[],[],[],[]]
9. y = []
10. with open('CreditCards.csv','r') as csvfile:
```

```

11.     r = csv.reader(csvfile)
12.     for i,rows in enumerate(r):
13.         if i==0:
14.             name = rows
15. with open('CreditCards.csv','r') as csvfile:
16.     reader = csv.DictReader(csvfile)
17.     for row in reader:
18.         for i in range(0,14):
19.             x[i].append(float(row[name[i]]))
20.             y.append(float(row['Y']))
21.
22. #pre-processing
23. def normalisation(xlist):
24.     xmin = min(xlist)
25.     xmax = max(xlist)
26.     for i in range(len(xlist)):
27.         xlist[i] = (xlist[i] - xmin)/(xmax - xmin)
28.         #print(xlist)
29.     return xlist
30.
31. list1 = []
32. for i in range(len(x)):
33.     list1.append(normalisation(x[i]))
34. list2 = [[row[i] for row in list1] for i in range(690)]
35.
36. #creating test and training sets
37. x_training = np.array(list2[0:621])
38. y_training = np.array(y[0:621])
39. x_test = np.array(list2[621:690])
40. y_test = np.array(y[621:690])
41.
42. #Part A get two accuracy
43. from sklearn.neighbors import KNeighborsClassifier
44. knn = KNeighborsClassifier(n_neighbors = 2)
45. knn.fit(x_training,y_training)
46. print("Accuracy for training set: ",knn.score(x_training,y_training))
47. print("Accuracy2 for test set: ",knn.score(x_test,y_test))
48. #Accuracy:  0.8969404186795491
49. #Accuracy2:  0.7681159420289855
50.
51. #Part B AUC score for training and test sets
52. #Part C plot them
53. neighbors = np.arange(1,31)
54.

```

```

55. def auclist(rangelist,xlist,ylist,xtest,ytest):
56.     auclist_train = np.empty(len(rangelist))
57.     auclist_test = np.empty(len(rangelist))
58.     for i,k in enumerate(rangelist):
59.         knn_n = KNeighborsClassifier(n_neighbors=k)
60.         knn_n.fit(xlist,ylist)
61.         y_pred = knn_n.predict_proba(xlist)
62.         y_pred2 = knn_n.predict_proba(xtest)
63.         auclist_train[i] = roc_auc_score(ylist,y_pred[:,1])
64.         auclist_test[i] = roc_auc_score(ytest,y_pred2[:,1])
65.     return auclist_train,auclist_test
66. auclist_train,auclist_test = auclist(neighbors,x_training,y_training,x_test,
    y_test)
67. auclist_testlist = auclist_test.tolist()
68.
69. print("the optimal value is: ",auclist_testlist.index(max(auclist_testlist))
    +1)
70.
71.
72. fig = plt.figure()
73. ax1 = fig.add_subplot(2,1,1)
74. ax2 = fig.add_subplot(2,1,2)
75.
76. ax1.plot(range(1,31), auclist_train, "p:",label="training", color='r')
77. ax2.plot(range(1,31), auclist_test, "p:",label="test", color='k')
78. plt.show()
79.
80.
81.
82. #Part D precision and recall for k=5 and k=2
83. def precision_and_recall(k,xlist,ylist,xtest,ytest):
84.     knn = KNeighborsClassifier(n_neighbors=k)
85.     knn.fit(xlist,ylist)
86.     y_pred = knn.predict(xtest)
87.     return precision_score(ytest,y_pred),recall_score(ytest,y_pred)
88.
89. print("precision and recall for k = 5 is: ",precision_and_recall(5,x_trainin
    g,y_training,x_test,y_test))
90. print("precision and recall for k = 2 is: ",precision_and_recall(2,x_trainin
    g,y_training,x_test,y_test))

```

