



实 验 报 告

(2017 / 2018 学年 第 2 学期)

课程名称	机器学习导论
实验名称	Decision Tree
实验时间	2018 年 5 月 20 日
指导教师	王邦

姓名	游浩然	学号	U201515429
----	-----	----	------------

1 问题重述

- 对课堂上所讲的西瓜数据集，构造决策树。
- 根据用户采集的 WiFi 信息采用决策树预测用户所在房间。

2 西瓜决策树

2.1 Python 代码实现

```
1  -*- coding: utf-8 -*-
2  """
3  @author : Haoran You
4
5  """
6  from math import log
7  import operator
8  import matplotlib
9  import matplotlib.pyplot as plt
10
11  """
12  Module - Create decision tress for datasets
13  :createTree(main)
14  :bestFeature
15  :calcShannonEnt
16  :splitDataset
17  :majorityCnt
18  """
19  def createTree(dataset, label):
20      classlist = [example[-1] for example in dataset]
21      # no subtree
22      if classlist.count(classlist[0]) == len(classlist):
23          return classlist[0]
24      # no feature
25      if len(dataset[0]) == 1:
26          return majorityCnt(classlist)
27      # choose best feature
28      bestFeat = bestFeature(dataset)
29      bestFeatLabel = label[bestFeat]
30      myTree = {bestFeatLabel: {}}
31      del(label[bestFeat])
32      # recursion for each subdataset
33      feat_values = [example[bestFeat] for example in dataset]
34      uniqueVals = set(feat_values)
35      for value in uniqueVals:
36          sub_label = label[:]
37          myTree[bestFeatLabel][value] = createTree(splitDataset(dataset, bestFeat, value), sub_label)
38      return myTree
39
40  def bestFeature(dataset):
41      num_features = len(dataset[0]) - 1
42      baseEntropy = calcShannonEnt(dataset)
43      bestInfoGain = 0.0
44      bestFeat = 0
45      for i in range(0, num_features):
46          featList = [example[i] for example in dataset]
47          uniqueVals = set(featList)
48          newEntropy = 0.0
49          for value in uniqueVals:
50              sub_dataset = splitDataset(dataset, i, value)
51              prob = len(sub_dataset) / float(len(dataset))
52              newEntropy += prob * calcShannonEnt(sub_dataset)
53          infoGain = baseEntropy - newEntropy
54          if infoGain > bestInfoGain:
55              bestInfoGain = infoGain
56              bestFeat = i
57      return bestFeat
58
59  def calcShannonEnt(dataset):
60      num_samples = len(dataset)
61      labelCounts = {}
62      for feat_vector in dataset:
63          current_label = feat_vector[-1]
64          if current_label not in labelCounts.keys():
65              labelCounts[current_label] = 0
66          labelCounts[current_label] += 1
67      shannonEnt = 0.0
68      for key in labelCounts:
69          prob = float(labelCounts[key]) / num_samples
70          shannonEnt -= prob * log(prob, 2)
71      return shannonEnt
```

```

72
73 def splitDataset(dataset, axis, value):
74     sub_dataset = []
75     for feat_vector in dataset:
76         if feat_vector[axis] == value:
77             reduce_feat_vector = feat_vector[:axis]
78             reduce_feat_vector.extend(feat_vector[axis+1:])
79             sub_dataset.append(reduce_feat_vector)
80     return sub_dataset
81
82 def majorityCnt(classlist):
83     classCount = {}
84     for vote in classlist:
85         if vote not in classCount.keys():
86             classCount[vote] = 0
87             classCount[vote] += 1
88     sortedClassCount = sorted(classCount.items(), key=operator.itemgetter(1), reverse=True)
89     return sortedClassCount[0][0]
90
91 """
92 Module - Plot decision tree constructed before
93 :createPlot(main)
94 :plotTree
95 :plotNode
96 :plotMidText
97 :getNumLeaves
98 :getTreeDepth
99 """
100 decisionNode = dict(boxstyle='sawtooth', fc='0.8')
101 leafNode = dict(boxstyle='round4', fc='0.8')
102 arrow_args = dict(arrowstyle='<-')
103
104 def createPlot(myTree):
105     fig = plt.figure(1, facecolor='white')
106     fig.clf()
107     axprops = dict(xticks=[], yticks=[])
108     createPlot.ax1 = plt.subplot(111, frameon=True)
109     plotTree.totalW = float(getNumLeaves(myTree))
110     plotTree.totalD = float(getTreeDepth(myTree))
111     plotTree.xOff = -0.5/plotTree.totalW
112     plotTree.yOff = 1.0
113     plotTree(myTree, (0.5, 1.0), '')
114     plt.show()
115
116 def plotTree(myTree, parentPt, nodeTxt):
117     num_leaves = getNumLeaves(myTree)
118     depth = getTreeDepth(myTree)
119     firstStr = list(myTree.keys())[0]
120     cntrPt = (plotTree.xOff + (1.0 + float(num_leaves))/2.0/plotTree.totalW, plotTree.yOff)
121     plotMidText(cntrPt, parentPt, nodeTxt)
122     plotNode(firstStr, cntrPt, parentPt, decisionNode)
123     secondDict = myTree[firstStr]
124     plotTree.yOff = plotTree.yOff - 1.0/plotTree.totalD
125     for key in secondDict.keys():
126         if type(secondDict[key]).__name__ == 'dict':
127             plotTree(secondDict[key], cntrPt, str(key))
128         else:
129             plotTree.xOff = plotTree.xOff + 1.0/plotTree.totalW
130             plotNode(secondDict[key], (plotTree.xOff, plotTree.yOff), cntrPt, leafNode)
131             plotMidText((plotTree.xOff, plotTree.yOff), cntrPt, str(key))
132     plotTree.yOff = plotTree.yOff + 1.0/plotTree.totalD
133
134 def plotNode(nodeTxt, centerPt, parentPt, nodeType):
135     createPlot.ax1.annotate(nodeTxt, xy=parentPt, xycoords='axes fraction',
136                             xytext=centerPt, textcoords='axes fraction',
137                             va="center", ha="center", bbox=nodeType, arrowprops=arrow_args)
138
139 def plotMidText(cntrPt, parentPt, txtString):
140     # positon for text
141     xMid = (parentPt[0]-cntrPt[0])/2.0 + cntrPt[0]
142     yMid = (parentPt[1]-cntrPt[1])/2.0 + cntrPt[1]
143     createPlot.ax1.text(xMid, yMid, txtString, va="center", ha="center")
144
145 def getNumLeaves(myTree):
146     numLeaves = 0
147     firstStr = list(myTree.keys())[0]
148     secondDict = myTree[firstStr]
149     for key in secondDict.keys():
150         if type(secondDict[key]).__name__ == 'dict':
151             numLeaves += getNumLeaves(secondDict[key])
152         else:
153             numLeaves += 1
154     return numLeaves
155
156 def getTreeDepth(myTree):
157     maxDepth = 0
158     firstStr = list(myTree.keys())[0]
159     secondDict = myTree[firstStr]
160     for key in secondDict.keys():

```

```

160         if type(secondDict[key]).__name__=='dict':
161             thisDepth = 1 + getTreeDepth(secondDict[key])
162         else:
163             thisDepth = 1
164         if thisDepth > maxDepth: maxDepth = thisDepth
165         return maxDepth
166
167 if __name__ == '__main__':
168     # load dataset
169     f = open('watermelon.txt', 'r', encoding='utf-8')
170     label = f.readline()
171     label = label.strip().split(' ')
172     dataset = [inst.strip().split(' ') for inst in f.readlines()]
173     print('feature ', label)
174     print('dataset ', dataset)
175     # tree structure
176     myTree = createTree(dataset, label)
177     print('myTree ', myTree)
178     # draw tree
179     from pylab import *
180     mpl.rcParams['font.sans-serif'] = ['SimHei']
181     createPlot(myTree)

```

2.2 结果图

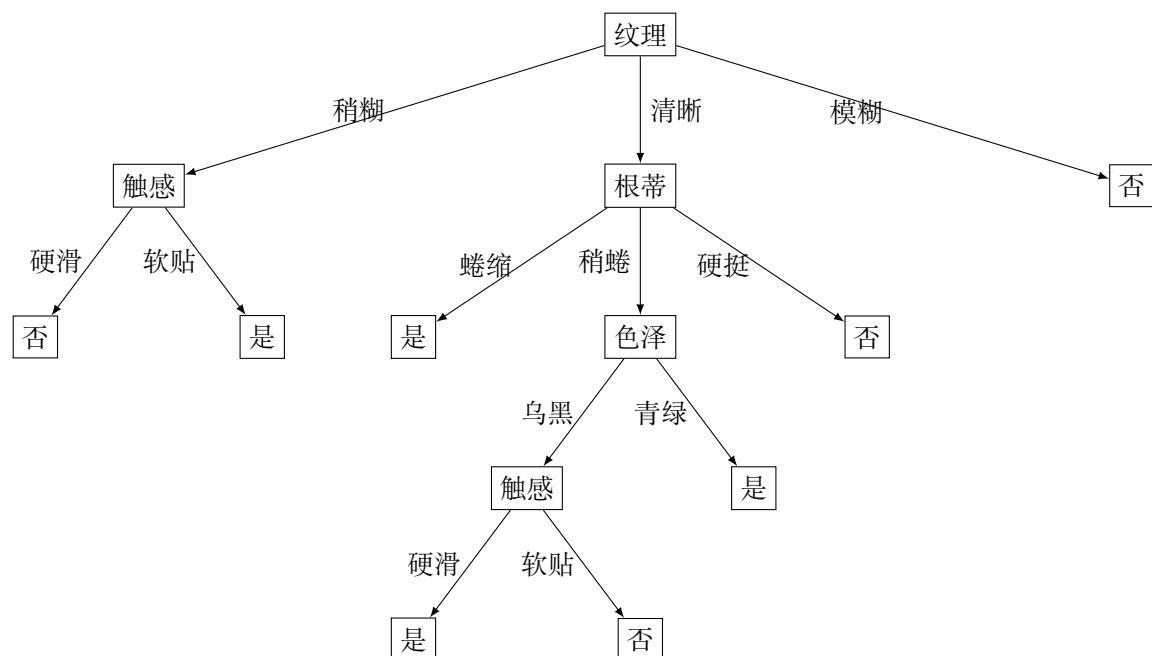


Figure 1: A decision tree for evaluating watermelon.

3 房间决策树

3.1 Python 代码实现

```

1  #-*- coding: utf-8 -*-
2  """
3  @author : Haoran You
4
5  """
6  from math import log
7  import csv
8  import os
9  import random
10 import matplotlib.pyplot as plt
11 from Task1 import *
12
13 def devidetraincsv():
14     trainDT = csv.reader(open('TrainDT.csv', 'r'))
15     dataset = []

```

```

16     for line in trainDT:
17         dataset.append(line)
18     del(dataset[0])
19     time_list = list(set([int(example[-1]) for example in dataset]))
20     BSSID_list = list(set([example[0] for example in dataset]))
21     label_dict = {}
22     for item in dataset:
23         label_dict[item[-1]] = item[2]
24     pre_whole_train = []
25     for i in range(0, len(time_list)):
26         new_dict = {}
27         for j in range(0, len(BSSID_list)):
28             new_dict[BSSID_list[j]] = 0
29         pre_whole_train.append(new_dict)
30     for item in dataset:
31         pre_whole_train[int(item[-1]) - 1][item[0]] = float(item[1])
32     whole_train = []
33     fin_label = 0
34     for item in pre_whole_train:
35         fin_label += 1
36         new_list = []
37         for keys in item.keys():
38             new_list.append(item[keys])
39         new_list.append(label_dict[str(fin_label)])
40         whole_train.append(new_list)
41     num_train = int(0.9 * len(whole_train))
42     num_val = len(whole_train) - num_train
43     train = random.sample(whole_train, num_train)
44     val = random.sample(whole_train, num_val)
45     return train, val, BSSID_list
46
47 def dividetestcsv(feature):
48     testDT = csv.reader(open('TestDT.csv', 'r'))
49     dataset = []
50     for line in testDT:
51         dataset.append(line)
52     del(dataset[0])
53     time_list = list(set([int(example[-1]) for example in dataset]))
54     BSSID_list = feature
55     pre_whole_test = []
56     for i in range(0, len(time_list)):
57         new_dict = {}
58         for j in range(0, len(BSSID_list)):
59             new_dict[BSSID_list[j]] = 0
60         pre_whole_test.append(new_dict)
61     for item in dataset:
62         pre_whole_test[int(item[-1]) - 1][item[0]] = float(item[1])
63     whole_test = []
64     fin_label = 0
65     for item in pre_whole_test:
66         fin_label += 1
67         new_list = []
68         for keys in item.keys():
69             new_list.append(item[keys])
70         whole_test.append(new_list)
71     return whole_test
72
73 def discretization(dataset):
74     for i in range(0, len(dataset[0]) - 1):
75         rss_list = [example[i] for example in dataset]
76         num = len(rss_list)
77         for k in range(0, num):
78             if rss_list[k] == 0.0:
79                 dataset[k][i] = 0
80             else:
81                 dataset[k][i] = 1
82     return dataset
83
84 def predict_val(dataset, feature, tree):
85     label = [example[-1] for example in dataset]
86     i = 0
87     total_num = len(label)
88     corr_num = 0
89     for item in dataset:
90         for keys in tree.keys():
91             index = feature.index(keys)
92             sub_tree = tree[keys][item[index]]
93             if isinstance(sub_tree, dict):
94                 pred = find_leaf(item, feature, sub_tree)
95             else:
96                 pred = sub_tree
97             if pred == label[i]:
98                 corr_num += 1
99             i += 1
100     acc = corr_num / total_num
101     return acc
102
103 def predict(dataset, feature, tree):

```

```

104     if os.path.exists('results.csv'):
105         os.remove('results.csv')
106     f = open('results.csv', 'a', newline='')
107     csv_write = csv.writer(f, dialect='excel')
108     i = 1
109     for item in dataset:
110         for keys in tree.keys():
111             index = feature.index(keys)
112             sub_tree = tree[keys][item[index]]
113             if isinstance(sub_tree, dict):
114                 pred = find_leaf(item, feature, sub_tree)
115             else:
116                 pred = sub_tree
117             i += 1
118             result = []
119             result.append(i)
120             result.append(pred)
121             csv_write.writerow(result)
122
123 def find_leaf(item, feature, tree):
124     for keys in tree.keys():
125         index = feature.index(keys)
126         sub_tree = tree[keys][item[index]]
127         if isinstance(sub_tree, dict):
128             pred = find_leaf(item, feature, sub_tree)
129         else:
130             pred = sub_tree
131     return pred
132
133 if __name__ == '__main__':
134     train_dataset, val_dataset, feature = devidetraincsv()
135     test_dataset = dividetestcsv(feature)
136     print('feature ', feature)
137     print('num_train ', len(train_dataset))
138     print('num_val ', len(val_dataset))
139     train_dataset = discretization(train_dataset)
140     myTree = createTree(train_dataset, feature.copy())
141     print(myTree)
142     createPlot(myTree)
143     val_dataset = discretization(val_dataset)
144     acc = predict_val(val_dataset, feature.copy(), myTree)
145     print('accuracy ', acc)
146     test_dataset = discretization(test_dataset)
147     predict(test_dataset, feature.copy(), myTree)

```

3.2 结果图

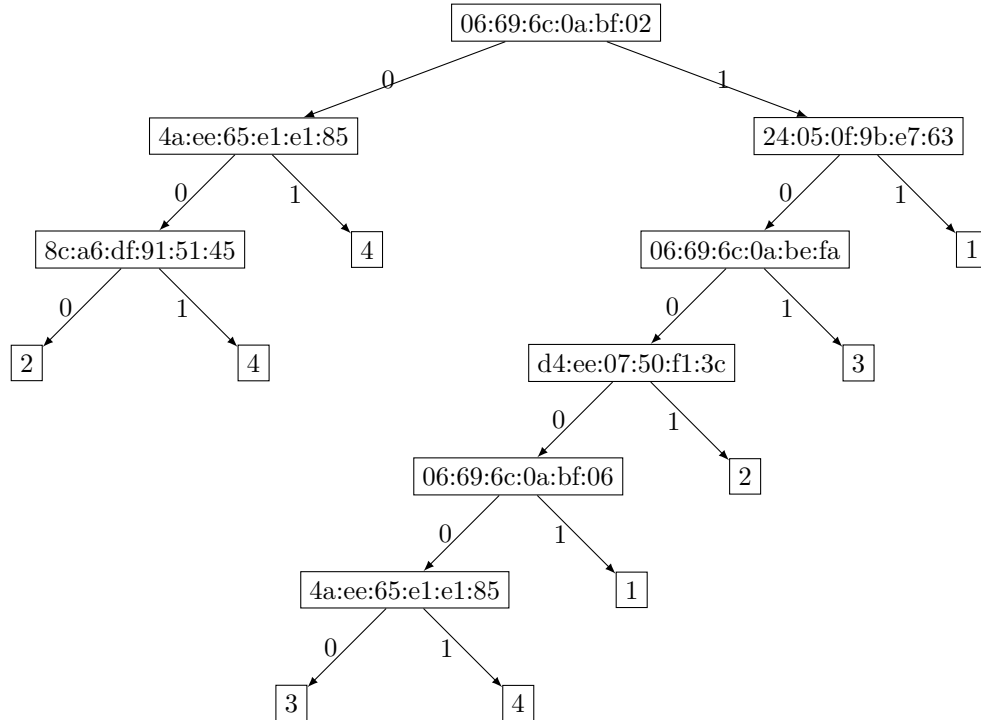


Figure 2: A decision tree for evaluating roomlabel. 0 presents no signal received, 1 means inversely.