# 编程实践：基于决策树和 C4.5 算法进行二分类

**主讲人**　张阳阳

张阳阳，清华大学博士生，研究方向：人工智能、机器学习、深度学习。擅长 Python, MATLAB, Tensorflow, Pytorch 框架，熟悉神经网络，支持向量机，朴素贝叶斯，决策树，聚类分析等AI算法和编程实践。曾作为 Python 课程讲师，有多年教学经验。

# 目录

● 二分类

| outlook | temperature | humidity | windy | result |
|---|---|---|---|---|
| sunny | hot | high | false | N |
| sunny | hot | high | true | N |
| overcast | hot | high | false | Y |
| rain | mild | high | false | Y |
| rain | cool | normal | false | Y |
| rain | cool | normal | true | N |
| overcast | cool | normal | true | Y |

● 训练集

| outlook | temperature | humidity | windy | result |
|---------|-------------|----------|-------|--------|
| sunny | hot | high | false | N |
| sunny | hot | high | true | N |
| overcast | hot | high | false | Y |
| rain | mild | high | false | Y |
| rain | cool | normal | false | Y |
| rain | cool | normal | true | N |
| overcast | cool | normal | true | Y |

● 测试集

| outlook | temperature | humidity | windy |
| --- | --- | --- | --- |
| sunny | mild | high | false |
| sunny | cool | normal | false |
| rain | mild | normal | false |
| sunny | mild | normal | true |
| overcast | mild | high | true |
| overcast | hot | normal | false |
| rain | mild | high | true |

```python
# -*- coding: utf-8 -*-
from math import log
import operator
import treePlotter

def calcShannonEnt(dataSet):
    """
    输入：数据集
    输出：数据集的香农熵
    描述：计算给定数据集的香农熵；熵越大，数据集的混乱程度越大
    """
    numEntries = len(dataSet)
    labelCounts = {}
    for featVec in dataSet:
        currentLabel = featVec[-1]
        if currentLabel not in labelCounts.keys():
            labelCounts[currentLabel] = 0
        labelCounts[currentLabel] += 1
    shannonEnt = 0.0
    for key in labelCounts:
        prob = float(labelCounts[key])/numEntries
        shannonEnt -= prob * log(prob, 2)
    return shannonEnt
```

```python
25  def splitDataSet(dataSet, axis, value):
26      """
27      输入：数据集，选择维度，选择值
28      输出：划分数据集
29      描述：按照给定特征划分数据集；去除选择维度中等于选择值的项
30      """
31      retDataSet = []
32      for featVec in dataSet:
33          if featVec[axis] == value:
34              reduceFeatVec = featVec[:axis]
35              reduceFeatVec.extend(featVec[axis+1:])
36              retDataSet.append(reduceFeatVec)
37      return retDataSet
```

```python
39  def chooseBestFeatureToSplit(dataSet):
40      """
41      输入：数据集
42      输出：最好的划分维度
43      描述：选择最好的数据集划分维度
44      """
45      numFeatures = len(dataSet[0]) - 1
46      baseEntropy = calcShannonEnt(dataSet)
47      bestInfoGainRatio = 0.0
48      bestFeature = -1
49      for i in range(numFeatures):
50          featList = [example[i] for example in dataSet]
51          uniqueVals = set(featList)
52          newEntropy = 0.0
53          splitInfo = 0.0
54          for value in uniqueVals:
55              subDataSet = splitDataSet(dataSet, i, value)
56              prob = len(subDataSet)/float(len(dataSet))
57              newEntropy += prob * calcShannonEnt(subDataSet)
58              splitInfo += -prob * log(prob, 2)
59          infoGain = baseEntropy - newEntropy
60          if (splitInfo == 0): # fix the overflow bug
61              continue
62          infoGainRatio = infoGain / splitInfo
63          if (infoGainRatio > bestInfoGainRatio):
64              bestInfoGainRatio = infoGainRatio
65              bestFeature = i
66      return bestFeature
```

```python
68  def majorityCnt(classList):
69      """
70      输入：分类类别列表
71      输出：子节点的分类
72      描述：数据集已经处理了所有属性，但是类标签依然不是唯一的，
73           采用多数判决的方法决定该子节点的分类
74      """
75      classCount = {}
76      for vote in classList:
77          if vote not in classCount.keys():
78              classCount[vote] = 0
79          classCount[vote] += 1
80      sortedClassCount = sorted(classCount.iteritems(), key=operator.itemgetter(1), reversed=True)
81      return sortedClassCount[0][0]
```

深蓝学院
shenlanxueyuan.com

```python
83  def createTree(dataSet, labels):
84      """
85      输入：数据集，特征标签
86      输出：决策树
87      描述：递归构建决策树，利用上述的函数
88      """
89      classList = [example[-1] for example in dataSet]
90      if classList.count(classList[0]) == len(classList):
91          # 类别完全相同，停止划分
92          return classList[0]
93      if len(dataSet[0]) == 1:
94          # 遍历完所有特征时返回出现次数最多的
95          return majorityCnt(classList)
96      bestFeat = chooseBestFeatureToSplit(dataSet)
97      bestFeatLabel = labels[bestFeat]
98      myTree = {bestFeatLabel:{}}
99      del(labels[bestFeat])
100     # 得到列表包括节点所有的属性值
101     featValues = [example[bestFeat] for example in dataSet]
102     uniqueVals = set(featValues)
103     for value in uniqueVals:
104         subLabels = labels[:]
105         myTree[bestFeatLabel][value] = createTree(splitDataSet(dataSet, bestFeat, value), subLabels)
106     return myTree
```

```python
108  def classify(inputTree, featLabels, testVec):
109      """
110      输入：决策树，分类标签，测试数据
111      输出：决策结果
112      描述：跑决策树
113      """
114      firstStr = list(inputTree.keys())[0]
115      secondDict = inputTree[firstStr]
116      featIndex = featLabels.index(firstStr)
117      for key in secondDict.keys():
118          if testVec[featIndex] == key:
119              if type(secondDict[key]).__name__ == 'dict':
120                  classLabel = classify(secondDict[key], featLabels, testVec)
121              else:
122                  classLabel = secondDict[key]
123      return classLabel
```

深蓝学院
shenlanxueyuan.com

```python
125  def classifyAll(inputTree, featLabels, testDataSet):
126      """
127      输入：决策树，分类标签，测试数据集
128      输出：决策结果
129      描述：跑决策树
130      """
131      classLabelAll = []
132      for testVec in testDataSet:
133          classLabelAll.append(classify(inputTree, featLabels, testVec))
134      return classLabelAll
```

深蓝学院
shenlanxueyuan.com

```python
136  def storeTree(inputTree, filename):
137      """
138          输入：决策树，保存文件路径
139          输出：
140          描述：保存决策树到文件
141      """
142      import pickle
143      fw = open(filename, 'wb')
144      pickle.dump(inputTree, fw)
145      fw.close()
```
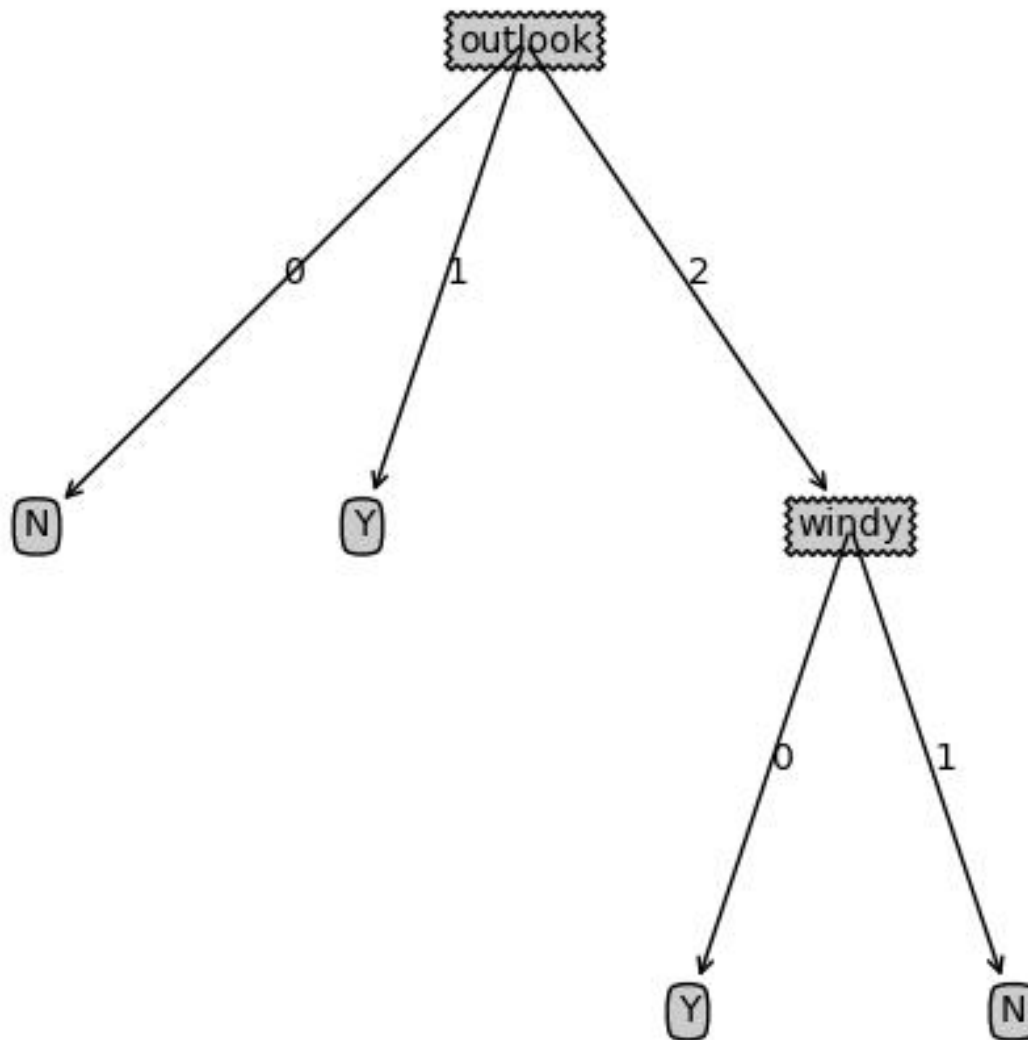
深蓝学院
shenlanxueyuan.com

```python
147  def grabTree(filename):
148      """
149      输入：文件路径名
150      输出：决策树
151      描述：从文件读取决策树
152      """
153      import pickle
154      fr = open(filename, 'rb')
155      return pickle.load(fr)
```

```python
def createDataSet():
    """
    outlook->  0: sunny | 1: overcast | 2: rain
    temperature-> 0: hot | 1: mild | 2: cool
    humidity-> 0: high | 1: normal
    windy-> 0: false | 1: true
    """
    dataSet = [[0, 0, 0, 0, 'N'],
               [0, 0, 0, 1, 'N'],
               [1, 0, 0, 0, 'Y'],
               [2, 1, 0, 0, 'Y'],
               [2, 2, 1, 0, 'Y'],
               [2, 2, 1, 1, 'N'],
               [1, 2, 1, 1, 'Y']]
    labels = ['outlook', 'temperature', 'humidity', 'windy']
    return dataSet, labels
```

```python
174  def createTestSet():
175      """
176      outlook->  0: sunny | 1: overcast | 2: rain
177      temperature-> 0: hot | 1: mild | 2: cool
178      humidity-> 0: high | 1: normal
179      windy-> 0: false | 1: true
180      """
181      testSet = [[0, 1, 0, 0],
182                 [0, 2, 1, 0],
183                 [2, 1, 1, 0],
184                 [0, 1, 1, 1],
185                 [1, 1, 0, 1],
186                 [1, 0, 1, 0],
187                 [2, 1, 0, 1]]
188      return testSet
```

```python
190  def main():
191      dataSet, labels = createDataSet()
192      labels_tmp = labels[:] # 拷贝, createTree会改变labels
193      desicionTree = createTree(dataSet, labels_tmp)
194      #storeTree(desicionTree, 'classifierStorage.txt')
195      #desicionTree = grabTree('classifierStorage.txt')
196      print('desicionTree:\n', desicionTree)
197      treePlotter.createPlot(desicionTree)
198      testSet = createTestSet()
199      print('classifyResult:\n', classifyAll(desicionTree, labels, testSet))
200
201  if __name__ == '__main__':
202      main()
```

深蓝学院
shenlanxueyuan.com

这些年来很多人因为癌症逝世，其中不乏娱乐圈中的一线明星们，而乳腺癌也成了很多女星不幸离世的原因。这种病症已经引起了来自社会的关注，定期检查变得很有必要。Wisconsin医学院的william H.Wolberg博士提供乳腺癌数据样本。所欲数据来自真实临床案例，每个案例有9个属性。



数据来源：http://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/

# Attribute                Domain

-- --------------------------------------

1. Sample code number          id number

2. Clump Thickness          1 - 10

3. Uniformity of Cell Size      1 - 10

4. Uniformity of Cell Shape      1 - 10

5. Marginal Adhesion          1 - 10

6. Single Epithelial Cell Size   1 - 10

7. Bare Nuclei          1 - 10

8. Bland Chromatin          1 - 10

9. Normal Nucleoli          1 - 10

10. Mitoses          1 - 10

**11. Class:**          **(2 for benign, 4 for malignant)**

**Class distribution:**

Benign: 458 (65.5%)

Malignant: 241 (34.5%)

请你使用附件中提供的数据：Breastdata.txt

使用决策树和 C4.5 算法，选择 80% 的数据作为训练集，剩余 80% 的数据作为测试集，然后预测下面这位患者【5,7,6,10,8,4,6,5,4】（9个属性的检测结果，分别用1-10表示）是患良性还是恶性肿瘤。提示：重难点是原始数据的读取，切割，后面稍微调整。



数据来源：http://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/

# Q&A

结语

深蓝学院
shenlanxueyuan.com

感谢各位聆听

请批评指正