

# 文本分类

杨善超

2016 年 4 月 20 日

<sup>1</sup>姓名: 杨善超

<sup>2</sup>学号: 13020510013

<sup>3</sup>邮箱: yangshanchaoysc@gmail.com

# 目 录

<b>1 爬取数据</b>	<b>1</b>
<b>2 中文分词</b>	<b>1</b>
<b>3 实验结果</b>	<b>1</b>
3.1 代码实现 . . . . .	1
3.2 运行结果 . . . . .	2
3.3 结果分析 . . . . .	2
<b>4 参考文献</b>	<b>2</b>
<b>5 附录</b>	<b>3</b>
5.1 程序代码 . . . . .	3
5.2 原始数据 . . . . .	10

## 1 爬取数据

## 2 中文分词

ID3 决策树有两个缺点 [1]:

- (1) 用信息增益选择属性时偏向于选择分枝比较多的属性值, 即取值多的属性
- (2) 不能处理连贯属性

而如表 2 所示, 题目所给数据的两个属性**湿度**和**温度**是连续值, 所以要对其进行离散化处理。此处采用论文中 C4.5 算法处理连续属性的方法 [1]。

Attributes can be either numeric or nominal and this determines the format of the test outcomes. For a numeric attribute A they are  $\{A \leq h, A > h\}$  where the threshold h is found by sorting S on the values of A and choosing the split between successive values that maximizes the criterion above.

经编程计算后得到, 温度最佳分割点为 **83**, 湿度最佳分割点为 **80**, 并将高于阈值的属性值标记为“高”, 不高于阈值的属性值标记为“正常”。

## 3 实验结果

### 3.1 代码实现

本次程序代码中 ID3 算法的核心部分基于《机器学习实战》[2]。另外进行的额外工作为:

- (1) 使用 Pandas 库替代 Numpy 库进行数据分析
- (2) 编写处理连续属性值的函数
- (3) 重新编写决策树的分类函数
- (4) 编写程序读入文件等程序

- (5) 更改原书 2.x 代码，使之在 3.5 版本上可以工作 (绘图程序可以自动绘制决策树)

## 3.2 运行结果

## 3.3 结果分析

正如图??所示, 人们在使用决策树进行决策时, 只需要使用**天气**、**湿度**、**风况**三个属性即可完成最终决策, 不需要使用温度属性。究其原因, 现实生活中除非剧烈高温和低温, 否则温度对人们是否室外活动的影响很小, 而天气当然是最重要的因素, 其次风况对于室外活动则影响很大, 比如羽毛球等运动。湿度的范围倘若不适合运动, 则对人的身体会产生较大影响。

输入新样本属性值, 进行决策后结果如表1所示:

表 1: 不同属性下决策结果

天气	温度	湿度	风况	运动
晴	寒冷	正常	有	适合
晴	寒冷	高	无	不适合

## 4 参考文献

- [1] Top 10 algorithms in data mining. *Knowledge and Information Systems*[M]. 2008(1), Volume 14, 3–4.
- [2] Peter Harringtoo 著, 李悦等译. 机器学习实战 [M]. 人民邮电出版社. 2015(1): 32–52.

## 5 附录

### 5.1 程序代码

```
1 # 主运行函数
2 from math import log
3 import pandas as pd
4 import treePlotter as tP
5
6
7 # 计算湿度和温度的最佳分割点，数据规约后返回 dataSet
8 def getBestSplit():
9     dataSet = pd.read_csv('weatherData.csv', encoding='gbk')
10    labels = list(dataSet.columns[:-1])
11    # 按温度进行升序排序
12    sortedData = dataSet.sort_values(by=dataSet.columns[1])
13    # 计算温度最佳分割点
14    # 只对属性值发生改变的地方才需要切开
15    # 记录信息增益以及最优分割点
16    infGain = 0
17    bestSplit = -1
18    # 初始样本的经验熵
19    HD = calcShannonEnt(dataSet)
20    for pos, item in enumerate(sortedData.values[:-2, :]):
21        prob = (pos + 1) / len(sortedData)
22        newEnt = prob * calcShannonEnt(sortedData.iloc[:pos+1, :]) + (1 -
23            prob) * calcShannonEnt(sortedData.iloc[pos+1:, :])
24        newEnt = HD - newEnt
25        if newEnt > infGain:
26            infGain = newEnt
27            bestSplit = item[1]
28    temperatureSplit = bestSplit
29    print('温度最佳分割点为□%d' % temperatureSplit)
30    # 计算湿度最佳分割点
31    # 按湿度进行升序排序
32    sortedData = dataSet.sort_values(by=dataSet.columns[2])
33    # 计算湿度最佳分割点
34    # 只对属性值发生改变的地方才需要切开
```

```

34     # 记录最优分割点以及信息增益
35     infGain = 0
36     bestSplit = -1
37     # 初始样本的经验熵
38     # HD = calcShannonEnt(dataSet)
39     for pos, item in enumerate(sortedData.values[: -2, :]):
40         # if item[-1] != sortedData.iloc[pos + 1, -1]:
41         prob = (pos + 1) / len(sortedData)
42         newEnt = prob * calcShannonEnt(sortedData.iloc[:pos+1, :]) + (1 -
43             prob) * calcShannonEnt(sortedData.iloc[pos+1:, :])
44         newEnt = HD - newEnt
45         if newEnt > infGain:
46             infGain = newEnt
47             bestSplit = item[2]
48     humiditySplit = bestSplit
49     print('湿度最佳分割点为%d' % humiditySplit)
50     # 将高于阈值的属性值标记为“高”，不高于阈值的属性值标记为“正常”
51     highPos = dataSet.iloc[:, 1] > temperatureSplit
52     lowPos = dataSet.iloc[:, 1] <= temperatureSplit
53     dataSet.ix[highPos, 1] = '高'
54     dataSet.ix[lowPos, 1] = '正常'
55     highPos = dataSet.iloc[:, 2] > humiditySplit
56     lowPos = dataSet.iloc[:, 2] <= humiditySplit
57     dataSet.ix[highPos, 2] = '高'
58     dataSet.ix[lowPos, 2] = '正常'
59     return dataSet, labels
60 # 计算给定的熵DataSetShannnon
61 def calcShannonEnt(dataSet):
62     # 计算类别的分布
63     numEntries = len(dataSet)
64     labelEntries = dataSet.iloc[:, -1]
65     # 计算频度
66     labelCounts = labelEntries.value_counts()
67     # 计算熵
68     shannonEnt = 0
69     for count in labelCounts:
70         prob = count/numEntries

```

```

71         # log base 2
72         shannonEnt -= prob * log(prob, 2)
73     return shannonEnt
74
75 # 对取出dataSet 第轴上值为的子数据aixsvalue
76 def splitDataSet(dataSet, axis, value):
77     retDataSet = dataSet.loc[dataSet.iloc[:, axis] == value]
78     del retDataSet[retDataSet.columns[axis]]
79     return retDataSet
80
81 # 根据信息增益选择最好的划分属性
82 def chooseBestFeatureToSplit(dataSet):
83     # the last column is used for the labels
84     numFeatures = len(dataSet.iloc[0]) - 1
85     baseEntropy = calcShannonEnt(dataSet)
86     bestInfoGain = 0.0
87     bestFeature = -1
88     for i in range(numFeatures):
89         # create a list of all the examples of this feature
90         featList = dataSet.iloc[:, i]
91         # get a set of unique values
92         uniqueVals = set(featList)
93         newEntropy = 0.0
94         for value in uniqueVals:
95             subDataSet = splitDataSet(dataSet, i, value)
96             prob = len(subDataSet)/float(len(dataSet))
97             newEntropy += prob * calcShannonEnt(subDataSet)
98         # calculate the info gain; ie reduction in entropy
99         infoGain = baseEntropy - newEntropy
100        # compare this to the best gain so far
101        if infoGain > bestInfoGain:
102            # if better than current best, set to best
103            bestInfoGain = infoGain
104            bestFeature = i
105    # returns an integer
106    return bestFeature
107
108 # 多数表决

```

```

109 def majorityCnt(classList):
110     classCount={}
111     for vote in classList:
112         if vote not in classCount.keys():
113             classCount[vote] = 0
114             classCount[vote] += 1
115     sortedClassCount = sorted(classCount.iteritems(), key=lambda x: x[1],
                                reverse=True)
116     return sortedClassCount[0][0]
117
118 # 递归创建决策树
119 def createTree(dataSet,labels):
120     classList = list(dataSet.iloc[:, -1])
121     if classList.count(classList[0]) == len(classList):
122         # stop splitting when all of the classes are equal
123         return classList[0]
124     if len(dataSet.iloc[0]) == 1:
125         # stop splitting when there are no more features in dataSet
126         return majorityCnt(classList)
127     bestFeat = chooseBestFeatureToSplit(dataSet)
128     bestFeatLabel = labels[bestFeat]
129     myTree = {bestFeatLabel: {}}
130     del(labels[bestFeat])
131     featValues = list(dataSet.iloc[:, bestFeat])
132     uniqueVals = set(featValues)
133     for value in uniqueVals:
134         myTree[bestFeatLabel][value] = createTree(splitDataSet(dataSet,
                                                                    bestFeat, value), labels[:])
135
136     return myTree
137
138 # given a problem instance, 决策是否适合运动
139 def decision(tree, instance):
140     outcome = tree
141     while type(outcome) == dict:
142         # 得到决策树当前需要决策的属性
143         key = list(outcome.keys())[0]
144         # 得到实际的属性值instance

```



```

145         pos = labels.index(key)
146         keyVal = instance[pos]
147         # get 结果
148         outcome = outcome[key][keyVal]
149     return outcome
150
151 if __name__ == '__main__':
152     # 读入数据
153     dataSet, labels = getBestSplit()
154     # dataSet, labels = createDataSet()
155
156     # 计算出决策树
157     tree = createTree(dataSet, labels[:])
158
159     # test two instances
160     instance = ['晴', '寒冷', '正常', '有']
161     decide = decision(tree, instance)
162     print(', '.join(instance), '->', decide)
163     instance = ['晴', '寒冷', '高', '无']
164     decide = decision(tree, instance)
165     print(', '.join(instance), '->', decide)
166
167     # 绘制决策树
168     tP.createPlot(tree)

```

```

1 # 绘图程序
2 import matplotlib.pyplot as plt
3 import matplotlib as mpl
4 # 指定默认字体
5 mpl.rcParams['font.sans-serif'] = ['SimHei']
6 # 自定义决策节点、叶子节点以及箭头的绘制属性（颜色什么的）
7 decisionNode = dict(boxstyle="sawtooth", fc="c", edgecolor='b')
8 leafNode = dict(boxstyle="round", facecolor="0.8", edgecolor='r')
9 arrow_args = dict(arrowstyle="<-", facecolor='b')
10
11
12 def getNumLeafs(myTree):
13     numLeafs = 0

```

```

14     firstStr = list(myTree.keys())[0]
15     secondDict = myTree[firstStr]
16     for key in secondDict.keys():
17         # test to see if the nodes are dictionnaires, if not they are leaf
            nodes
18         if type(secondDict[key]) == dict:
19             numLeafs += getNumLeafs(secondDict[key])
20         else:
21             numLeafs += 1
22     return numLeafs
23
24 def getTreeDepth(myTree):
25     maxDepth = 0
26     firstStr = list(myTree.keys())[0]
27     secondDict = myTree[firstStr]
28     for key in secondDict.keys():
29         if type(secondDict[key]) == dict:
30             # test to see if the nodes are dictionnaires, if not they are leaf
                nodes
31             thisDepth = 1 + getTreeDepth(secondDict[key])
32         else:
33             thisDepth = 1
34         if thisDepth > maxDepth:
35             maxDepth = thisDepth
36     return maxDepth
37
38 def plotNode(nodeTxt, centerPt, parentPt, nodeType):
39     createPlot.ax1.annotate(nodeTxt, xy=parentPt, xycoords='axes_fraction'
40                             ,
41                             xytext=centerPt, textcoords='axes_fraction', va="
                                center", ha="center", bbox=nodeType,
42                             arrowprops=arrow_args)
43
44 def plotMidText(cntrPt, parentPt, txtString):
45     xMid = (parentPt[0]-cntrPt[0])/2.0 + cntrPt[0]
46     yMid = (parentPt[1]-cntrPt[1])/2.0 + cntrPt[1]
47     createPlot.ax1.text(xMid, yMid, txtString, va="center", ha="center",
48                         rotation=30)

```

```

47
48
49 # if the first key tells you what feat was split on
50 def plotTree(myTree, parentPt, nodeTxt):
51     # this determines the x width of this tree
52     numLeafs = getNumLeafs(myTree)
53     # the text label for this node should be this
54     firstStr = list(myTree.keys())[0]
55     cntrPt = (plotTree.xOff + (1.0 + float(numLeafs))/2.0/plotTree.totalW,
               plotTree.yOff)
56     plotMidText(cntrPt, parentPt, nodeTxt)
57     plotNode(firstStr, cntrPt, parentPt, decisionNode)
58     secondDict = myTree[firstStr]
59     plotTree.yOff = plotTree.yOff + 1.0/plotTree.totalD
60     for key in secondDict.keys():
61         # test to see if the nodes are dictionaries, if not they are leaf
           nodes
62         if type(secondDict[key]) == dict:
63             # recursion
64             plotTree(secondDict[key], cntrPt, str(key))
65         else:
66             # it's a leaf node print the leaf node
67             plotTree.xOff = plotTree.xOff + 1.0/plotTree.totalW
68             plotNode(secondDict[key], (plotTree.xOff, plotTree.yOff), cntrPt
               , leafNode)
69             plotMidText((plotTree.xOff, plotTree.yOff), cntrPt, str(key))
70     plotTree.yOff = plotTree.yOff + 1.0/plotTree.totalD
71
72
73 # if you do get a dictionary you know it's a tree, and the first element
   will be another dict
74 def createPlot(inTree):
75     fig = plt.figure(1, facecolor='white')
76     fig.clf()
77     axprops = dict(xticks=[], yticks=[])
78     # no ticks
79     createPlot.ax1 = plt.subplot(111, frameon=False, **axprops)
80

```

```

81 plotTree.totalW = float(getNumLeafs(inTree))
82 plotTree.totalD = float(getTreeDepth(inTree))
83 plotTree.xOff = -0.5/plotTree.totalW
84 plotTree.yOff = 1.0
85 plotTree(inTree, (0.5, 1.0), '')
86 plt.show()

```

## 5.2 原始数据

表 2: 天气对户外活动影响

天气	温度	湿度	风况	运动
晴	85.0	85.0	无	不适合
晴	85.0	90.0	有	不适合
多云	83.0	78.0	无	适合
有雨	70.0	96.0	无	适合
有雨	68.0	80.0	无	适合
有雨	65.0	70.0	有	不适合
多云	64.0	65.0	有	适合
晴	72.0	95.0	无	不适合
晴	69.0	70.0	无	适合
有雨	75.0	80.0	无	适合
晴	75.0	70.0	有	适合
多云	72.0	90.0	有	适合
多云	81.0	75.0	无	适合
有雨	71.0	80.0	有	不适合