

# 实验七 分类 决策树算法

## 【实验目的】

1. 了解 Anaconda 和 python 的使用与设置。
2. 掌握分类决策树算法的基本原理和实现方法。

## 【实验学时】

建议 2 学时

## 【实验环境配置】

- 1、Windows 环境
- 2、Anaconda
- 3、Pandas

## 【实验原理】

数据挖掘中分类决策树算法的实现：

- 1、Pandas 数据导入
- 2、决策树算法
- 3、sk-learn 库的决策树模块使用

## 【实验步骤】

1. 打开 Anaconda 的 jupyter notebook，创建实验七，明确标题和步骤。
2. 导入 pandas 库和 os 库，从 lenses 数据文件中，导入实验训练数据和测试数据。  
(注：可能涉及的函数 os.path.join, os.walk, pandas.read\_csv)

```

1 from math import log
2
3 def calcShannonEnt(dataSet): # 计算数据集的香农熵
4     numEntries = len(dataSet) # 数据集中实例的总数
5     labelCounts = {} # 字典键值记录当前类别出现的次数
6     for featVec in dataSet:
7         currentLabel = featVec[-1] # 当前标签为字典
8         if currentLabel not in labelCounts.keys():
9             labelCounts[currentLabel] = 0
10            labelCounts[currentLabel] += 1
11    shannonEnt = 0.0
12    for key in labelCounts:
13        prob = float(labelCounts[key])/numEntries # 计算每个类别的概率
14        shannonEnt -= prob * log(prob, 2) # 熵值
15    return shannonEnt
16
17 def createDataSet():
18     dataSet = [[1, 1, 'yes'],
19                [1, 1, 'yes'],
20                [1, 0, 'no'],
21                [0, 1, 'no'],
22                [0, 1, 'no']]
23     labels = ['no surfacing', 'flippers']
24     return dataSet, labels
25
26 myDat, labels = createDataSet()
27 myDat

```

```
[[1, 1, 'yes'], [1, 1, 'yes'], [1, 0, 'no'], [0, 1, 'no'], [0, 1, 'no']]
```

```
1 calcShannonEnt(myDat)
```

```
0.9709505944546686
```

- 对导入的实验数据，根据《机器学习实战》教程第三章，完成分类决策树 ID3 算法，并完成对于 lenses 的分类识别。

(1) 按照给定特征分割数据集 (待划分数据集, 划分数据集的特征, 需要返回的特征的值)

```

1 def splitDataSet(dataSet, axis, value): # 按照给定特征分割数据集 (待划分数据集, 划分数据集的特征, 需要返回的特征的值)
2     retDataSet = [] # 新的list对象
3     for featVec in dataSet:
4         if featVec[axis] == value:
5             reducedFeatVec = featVec[:axis] # 接收axis之前的特征
6             reducedFeatVec.extend(featVec[axis+1:]) # 将axis之后的特征进行扩展
7             retDataSet.append(reducedFeatVec)
8     return retDataSet # 得到依据该特征划分后的数据子集
9
10 splitDataSet(myDat, 0, 1)

```

```
[[1, 'yes'], [1, 'yes'], [0, 'no']]
```

(2) 选取最好的数据集的划分方式(选择最好的 feature 来划分数据集)

```

1 def chooseBestFeatureToSplit(dataSet): # 选取最好的数据集的划分方式
2     numFeatures = len(dataSet[0]) - 1 # 特征数量
3     baseEntropy = calcShannonEnt(dataSet) # 计算数据集的香农熵
4     bestInfoGain = 0.0 # 初始化最优的信息增益
5     bestFeature = -1
6     for i in range(numFeatures): # 依据特征循环
7         featList = [example[i] for example in dataSet] # 第i个特征的所有内容
8         uniqueVals = set(featList) # 去重
9         newEntropy = 0.0
10        for value in uniqueVals:
11            subDataSet = splitDataSet(dataSet, i, value) # 数据子集
12            prob = len(subDataSet)/float(len(dataSet)) # 概率
13            newEntropy += prob * calcShannonEnt(subDataSet) # 计算数据子集中每个特征的香农熵
14            infoGain = baseEntropy - newEntropy # 信息增益
15            if infoGain > bestInfoGain: # 比较当前信息增益和上一次的
16                bestInfoGain = infoGain
17                bestFeature = i
18    return bestFeature
19
20 chooseBestFeatureToSplit(myDat)

```

0

#### (4) 构建决策树

```

1 def majorityCnt(classList): # 分类次数最多的一类, 参数 (分类名称的列表)
2     classCount = {} # 空字典
3     for vote in classList: # 创建键值为classList中唯一的数据字典
4         if vote in classCount.keys(): # 统计每个类别的出现的频率
5             classCount[vote] += 1
6         classCount[vote] += 1
7     sortedClassCount = sorted(classCount.items(), key=operator.itemgetter(1), reverse=True) # 根据频率排序
8     return sortedClassCount[0][0] # 返回出现次数最多的分类名称
9
10 def createTree(dataSet, labels): # 创建树, 参数 (数据集和包含所有特征的标签列表)
11     classList = [example[-1] for example in dataSet]
12     if classList.count(classList[0]) == len(classList): # 类别完全相同则停止继续划分
13         return classList[0]
14     if len(dataSet[0]) == 1: # 遍历完所有特征返回出现次数最多的
15         return majorityCnt(classList)
16     bestFeat = chooseBestFeatureToSplit(dataSet) # 选取最好的特征
17     bestFeatLabel = labels[bestFeat]
18     myTree = {bestFeatLabel: {}} # 字典myTree存储了树的所有信息
19     del(labels[bestFeat]) # 删除已选的特征
20     featValues = [example[bestFeat] for example in dataSet] # 得到所有属性值
21     uniqueVals = set(featValues) # 去重
22     for value in uniqueVals:
23         subLabels = labels[:] # 特征子集
24         myTree[bestFeatLabel][value] = createTree(splitDataSet(dataSet, bestFeat, value), subLabels)
25     return myTree
26
27 myDat, labels = createDataSet()
28 createTree(myDat, labels)

```

```
{'no surfacing': {0: 'no', 1: {'flippers': {0: 'no', 1: 'yes'}}}}
```

#### (5) 绘图函数

```

7
8
9 def plotNode(nodeTxt, centerPt, parentPt, nodeType): # 定义了实际的绘图功能
10     createplot.ax1.annotate(nodeTxt, xy=parentPt, xycoords='axes fraction',
11                             xytext=centerPt, textcoords='axes fraction',
12                             va="center", ha="center", bbox=nodeType,
13                             arrowprops=arrow_args) # 文本内容, 箭头尖端, 文本位置, coords分别指定坐标系
14
15 def createplot(): # 绘制总体的图
16     fig = plt.figure(1, facecolor='white')
17     fig.clf()
18     createplot.ax1 = plt.subplot(111, frameon=False)
19     plotNode('决策节点', (0.5, 0.1), (0.1, 0.5), decisionNode)
20     plotNode('叶节点', (0.8, 0.1), (0.3, 0.8), leafNode)
21     plt.show()
22
23 def getNumLeafs(myTree): # 获取叶节点的数目
24     numLeafs = 0
25     firstStr = list(myTree.keys())[0]
26     secondDict = myTree[firstStr]
27     for key in secondDict.keys():
28         if type(secondDict[key]).__name__ == 'dict': # 测试节点的数据是否为字典
29             numLeafs += getNumLeafs(secondDict[key]) # 如果是依据子树继续计数
30         else:
31             numLeafs += 1
32     return numLeafs
33
34 def getTreeDepth(myTree):
35     maxDepth = 0
36     firstStr = list(myTree.keys())[0]
37     secondDict = myTree[firstStr]
38     for key in secondDict.keys():
39         if type(secondDict[key]).__name__ == 'dict': # 测试节点的数据是否为字典
40             thisDepth = 1 + getTreeDepth(secondDict[key])
41         else:
42             thisDepth = 1
43         if thisDepth > maxDepth:
44             maxDepth = thisDepth
45     return maxDepth
46
47 def retrieveTree(i): # 输出预先存储的树的信息,用于测试
48     listOfTrees = [{'no surfacing': {0: 'no', 1: {'flippers': {0: 'no', 1: 'yes'}}}},
49                    {'no surfacing': {0: 'no', 1: {'flippers': {0: {'head': {0: 'no', 1: 'yes'}}, 1: 'no'}}}}]
50     return listOfTrees[i]
51
52 myTree = retrieveTree(1)
53 myTree = retrieveTree(0)
54 myTree

```

```
{'no surfacing': {0: 'no', 1: {'flippers': {0: 'no', 1: 'yes'}}}}
```

```
1 getNumLeafs(myTree)
```

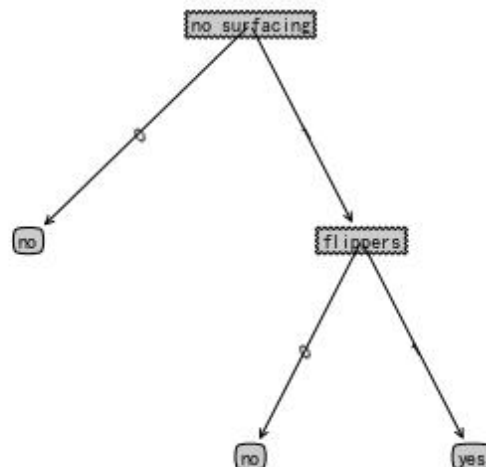
3

(6) 对初始数据使用绘图函数进行测试

```

1 def plotMidText(cntrPt, parentPt, txtString): # 在父子节点间填充文本信息
2     xMid = (parentPt[0]-cntrPt[0])/2.0 + cntrPt[0] # 填充信息的x横坐标
3     yMid = (parentPt[1]-cntrPt[1])/2.0 + cntrPt[1] # 填充信息的y纵坐标
4     createplot.ax1.text(xMid, yMid, txtString, va="center", ha="center", rotation=30)
5
6 def plotTree(myTree, parentPt, nodeTxt):
7     numLeafs = getNumLeafs(myTree) # 计算树叶节点的数量
8     depth = getTreeDepth(myTree) # 计算树深
9     firstStr = list(myTree.keys())[0]
10    cntrPt = (plotTree.xOff + (1.0 + float(numLeafs))/2.0/plotTree.totalW, plotTree.yOff) # 文本位置
11    plotMidText(cntrPt, parentPt, nodeTxt) # 填充父子节点之间的信息
12    plotNode(firstStr, cntrPt, parentPt, decisionNode) # 实际绘图
13    secondDict = myTree[firstStr]
14    plotTree.yOff = plotTree.yOff - 1.0/plotTree.totalD # 减少y的偏移
15    for key in secondDict.keys(): # 遍历决策树
16        if type(secondDict[key]).__name__ == 'dict':
17            plotTree(secondDict[key], cntrPt, str(key))
18        else:
19            plotTree.xOff = plotTree.xOff + 1.0/plotTree.totalW
20            plotNode(secondDict[key], (plotTree.xOff, plotTree.yOff), cntrPt, leafNode)
21            plotMidText((plotTree.xOff, plotTree.yOff), cntrPt, str(key))
22    plotTree.yOff = plotTree.yOff + 1.0/plotTree.totalD
23
24 def createplot(inTree):
25     fig = plt.figure(1, facecolor='white')
26     fig.clf() # 清除去上图数据
27     axprops = dict(xticks=[], yticks=[])
28     createplot.ax1 = plt.subplot(111, frameon=False, **axprops)
29     plotTree.totalW = float(getNumLeafs(inTree)) # 树宽
30     plotTree.totalD = float(getTreeDepth(inTree)) # 树深
31     plotTree.xOff = -0.5/plotTree.totalW; plotTree.yOff = 1.0;
32     plotTree(inTree, (0.5, 1.0), '')
33     plt.show()
34
35
36 myTree = retrieveTree(0)
37 createplot(myTree)

```





## 完成分类函数

```
: 1 def classify(inputTree, featLables, testVec):
2     firstStr = list(inputTree.keys())[0]
3     secondList = inputTree[firstStr]
4     featIndex = featLables.index(firstStr) # 将标签字符串转换为索引
5     for key in secondList.keys():
6         if testVec[featIndex] == key:
7             if type(secondList[key]).__name__ == 'dict':
8                 classLabel = classify(secondList[key], featLables, testVec)
9             else:
10                classLabel = secondList[key]
11     return classLabel
12
13
14 myDat, labels = createDataSet()
15 myTree = retrieveTree(0)
16 classify(myTree, labels, [1, 0])

: 'no'
```

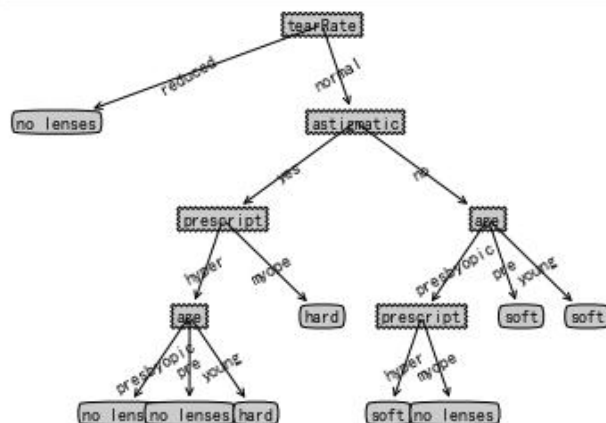
## 储存树函数

```
1 def storeTree(inputTree, filename):
2     import pickle
3     fw = open(filename, 'wb')
4     pickle.dump(inputTree, fw)
5     fw.close()
```

```
1 def grabTree(filename):
2     import pickle
3     fr = open(filename, 'rb')
4     return pickle.load(fr)
5
```

## 从文件读取数据并完成决策树分类算法

```
6 fr = open('lenses.txt')
7 lenses = [inst.strip().split('\t') for inst in fr.readlines()]
8 lensesLabels = ['age', 'prescript', 'astigmatic', 'tearRate']
9 lensesTree = createTree(lenses, lensesLabels)
10 createplot(lensesTree)
```



4. 对导入的实验数据，根据 sklearn 的分类决策树算法模块 sklearn.tree.DecisionTreeClassifier，完成对于 lenses 数据的分类识别。

(1) 数据集：

young	myope	no	reduced	no lenses
young	myope	no	normal	soft
young	myope	yes	reduced	no lenses
young	myope	yes	normal	hard
young	hyper	no	reduced	no lenses
young	hyper	no	normal	soft
young	hyper	yes	reduced	no lenses
young	hyper	yes	normal	hard
pre	myope	no	reduced	no lenses
pre	myope	no	normal	soft
pre	myope	yes	reduced	no lenses
pre	myope	yes	normal	hard
pre	hyper	no	reduced	no lenses
pre	hyper	no	normal	soft
pre	hyper	yes	reduced	no lenses
pre	hyper	yes	normal	no lenses
presbyopic	myope	no	reduced	no lenses
presbyopic	myope	no	normal	no lenses
presbyopic	myope	yes	reduced	no lenses
presbyopic	myope	yes	normal	hard

(2) 代码实现

导入数据

```
1 import os
2 import pandas as pd
3 dir_ds=os.path.join(os.path.curdir,'lenses.txt')
4 df_data=pd.read_csv(dir_ds,sep='\t',header=None,names=['age','spectacle','astigmatic','tear','class'])
5 df_data

1 from sklearn.model_selection import train_test_split
2 training_set,testing_set=train_test_split(df_data,test_size=0.2)
3 training_set=training_set.reset_index()
4 training_set
```

数据整理

```
1 testing_set=testing_set.reset_index()
2 testing_set
```

	index	age	spectacle	astigmatic	tear	class
0	2	young	myope	yes	reduced	no lenses
1	9	pre	myope	no	normal	soft
2	20	presbyopic	hyper	no	reduced	no lenses
3	21	presbyopic	hyper	no	normal	soft
4	12	pre	hyper	no	reduced	no lenses

```
1 from sklearn import tree
2 X=df_data.iloc[:,0:4]
3 Y=df_data.iloc[:,4]
4 from sklearn.preprocessing import LabelEncoder
5 ds=LabelEncoder()
6 ds.fit(['no lenses','hard','soft','young','pre','presbyopic','myope','hyper','no','yes','reduced','normal'])
7 Y1=ds.transform(Y.iloc[:])
8 Y1
```

```
array([4, 9, 4, 0, 4, 9, 4, 0, 4, 9, 4, 0, 4, 9, 4, 4, 4, 4, 4, 0, 4, 9,
       4, 4])
```

```
1 X1=[ds.transform(x) for i,r in X.iterrows()]
2 X1
```

```
[array([11, 2, 3, 8]),
 array([11, 2, 3, 5]),
 array([11, 2, 10, 8]),
 array([11, 2, 10, 5]),
 array([11, 1, 3, 8])]
```

使用 sklearn 完成分类算法，并将运行结果存入文档 dt.txt

```
1 model_dt=tree.DecisionTreeClassifier(criterion='entropy')
2 model_dt=model_dt.fit(X1,Y1)
3 tree.export_graphviz(model_dt,out_file='dt.txt')
```

### (3) 运行结果

dt.txt 文档结果：

```
digraph Tree {
node [shape=box];
0 [label="X[3] <= 6.5\nentropy = 1.326\nsamples = 24\nvalue = [4, 15, 5]";
1 [label="X[2] <= 6.5\nentropy = 1.555\nsamples = 12\nvalue = [4, 3, 5]";
0 -> 1 [labeldistance=2.5, labelangle=45, headlabel="True"];
2 [label="X[1] <= 1.5\nentropy = 0.65\nsamples = 6\nvalue = [0, 1, 5]";
1 -> 2;
3 [label="entropy = 0.0\nsamples = 3\nvalue = [0, 0, 3]";
2 -> 3;
4 [label="X[0] <= 6.5\nentropy = 0.918\nsamples = 3\nvalue = [0, 1, 2]";
2 -> 4;
5 [label="entropy = 0.0\nsamples = 1\nvalue = [0, 0, 1]";
4 -> 5;
6 [label="X[0] <= 9.0\nentropy = 1.0\nsamples = 2\nvalue = [0, 1, 1]";
4 -> 6;
7 [label="entropy = 0.0\nsamples = 1\nvalue = [0, 1, 0]";
6 -> 7;
8 [label="entropy = 0.0\nsamples = 1\nvalue = [0, 0, 1]";
6 -> 8;
9 [label="X[1] <= 1.5\nentropy = 0.918\nsamples = 6\nvalue = [4, 2, 0]";
1 -> 9;
```