

## 5.决策树

### 运行效果：

我建立了一个数据集，里边有两类Stu.Meng和Dr.Jiang来代指我和老师，前边的feat1\2\3\4是区分我们两个人的变量（随便从网上找的数据集，略微改了改，原来是根据四个变量区分花的种类的，只用了数据集，代码是自己写的），通过这四个变量区分我和老师,为了方便区分这两类，我把两个人的feat4调的差距大了些，一个是0~0.5，一个是1~2.5，这样我们（人，不是程序）看到的时候就知道是不是分错了

### 示例

```
164         dataSet = loadCSV()
165         decisionTree = buildTree(dataSet)
166         pruneTree(decisionTree, 0.4)
167         pre_name = [6.7, 3.1, 4.4, 2.2]
168
169         print(classify(pre_name, decisionTree))
170
calculateDiffCount()
```

CART ×

C:\Users\mziyu\anaconda3\python.exe C:/Users/mziyu/Desktop/孟子喻作业/机器学习/CART.py  
{'Dr.Jiang': 16}

```
164         dataSet = loadCSV()
165         decisionTree = buildTree(dataSet)
166         pruneTree(decisionTree, 0.4)
167         pre_name = [5.1, 3.5, 1.4, 0.2]
168
169         print(classify(pre_name, decisionTree))
170
```

CART ×

C:\Users\mziyu\anaconda3\python.exe C:/Users/mziyu/Desktop/孟子喻作业/机器学习/CART.py  
{'Stu.Meng': 23}

调节阈值后基本都能分类正确

完整代码见CART.py

写这个挺有意思的，因为之前只是听说，但是自己手写一遍真真让人学到很多，比如之前的分离点的计算，对算法的理解更透彻了

### 5.1构建树类

因为对构建树的方式不熟悉，所以查阅了比较多的资料，以下这种标记 TrueBranch 和 FalseBranch 的方式是比较容易理解的，所以使用此种方式

```
1 class CARTTree:
2     def __init__(self, value=None, trueBranch=None, falseBranch=None,
3         results=None, col=-1, summary=None, data=None):
4         self.value = value
```

```

4         self.trueBranch = trueBranch
5         self.falseBranch = falseBranch
6         self.results = results
7         self.col = col
8         self.summary = summary
9         self.data = data
10
11     def __str__(self):
12         print(self.col, self.value)
13         print(self.results)
14         print(self.summary)
15         return ""

```

## 5.2基尼指数的计算

计算基尼指数，此处按课本公式来即可，比较简单

```

1  def gini(dataset):
2      # 计算基尼指数
3
4      data_num = len(dataset)
5      # 以下部分用于统计每个类别出现的个数，并组成一个字典，存在results中
6      results = {}
7      for data in dataset:
8          # data[-1] means dataType
9          if data[-1] not in results:
10             results.setdefault(data[-1], 1)
11          else:
12             results[data[-1]] += 1
13
14      gini_result = 0
15      gini_result = float(gini_result)
16      for i in results:
17          gini_result += (results[i] / data_num) * (results[i] / data_num)
18      return 1 - gini_result

```

## 5.3建树

建立CART二叉树，即针对不同的切分点，计算整体的基尼指数，取能使基尼指数最小的点做父节点建立CART决策树

```

1  def buildTree(rows):
2      # 递归建立决策树， 当gain=0, 时停止回归
3      # build decision tree bu recursive function
4      # stop recursive function when gain = 0
5      # return tree
6      currentGain = gini(rows)
7      column_lenght = len(rows[0])
8      rows_length = len(rows)
9
10     best_gain = 0.0
11     best_value = None
12     best_set = None
13
14     # choose the best gain
15     for col in range(column_lenght - 1):

```

```

16         col_value_set = set([x[col] for x in rows])
17         for value in col_value_set:
18             list1, list2 = chooseSplitData(rows, value, col)
19             p = len(list1) / rows_length
20             gain = currentGain - p * gini(list1) - (1 - p) * gini(list2)
21             if gain > best_gain:
22                 best_gain = gain
23                 best_value = (col, value)
24                 best_set = (list1, list2)
25         dcY = {'impurity': '%.3f' % currentGain, 'sample': '%d' % rows_length}
26         #
27         # stop or not stop
28
29         if best_gain > 0:
30             trueBranch = buildTree(best_set[0])
31             falseBranch = buildTree(best_set[1])
32             return CARTTree(col=best_value[0], value=best_value[1],
trueBranch=trueBranch, falseBranch=falseBranch, summary=dcY)
33         else:
34             return CARTTree(results=calculatedDiffCount(rows), summary=dcY,
data=rows)

```

## 5.4分离点的计算

根据相对分离点大小将实例点分到左右两个子树中

```

1  def chooseSplitData(dataset, value, column):
2      lefttree = []
3      righttree = []
4
5      if isinstance(value, int) or isinstance(value, float):
6          for row in dataset:
7              if row[column] >= value:
8                  lefttree.append(row)
9              else:
10                 righttree.append(row)
11     else:
12         for row in dataset:
13             if row[column] == value:
14                 lefttree.append(row)
15             else:
16                 righttree.append(row)
17     return lefttree, righttree

```

## 5.5剪枝

计算是否取得损失函数最小，否则进行剪枝

```

1  def pruneTree(tree, miniGain):
2      if tree.trueBranch.results == None:
3          pruneTree(tree.trueBranch, miniGain)
4      if tree.falseBranch.results == None:
5          pruneTree(tree.falseBranch, miniGain)
6
7      if tree.trueBranch.results != None and tree.falseBranch.results !=
None:

```

```

8         len1 = len(tree.trueBranch.data)
9         len2 = len(tree.falseBranch.data)
10        len3 = len(tree.trueBranch.data + tree.falseBranch.data)
11
12        p = float(len1) / (len1 + len2)
13
14        gain = gini(tree.trueBranch.data + tree.falseBranch.data) - p *
gini(tree.trueBranch.data) - (1 - p) * gini(tree.falseBranch.data)
15
16        if gain < miniGain:
17            tree.data = tree.trueBranch.data + tree.falseBranch.data
18            tree.results = calculateDiffCount(tree.data)
19            tree.trueBranch = None
20            tree.falseBranch = None

```

## 5.6分类预测

```

1  def classify(data, tree):
2      if tree.results != None:
3          return tree.results
4      else:
5          branch = None
6          v = data[tree.col]
7          if isinstance(v, int) or isinstance(v, float):
8              if v >= tree.value:
9                  branch = tree.trueBranch
10             else:
11                 branch = tree.falseBranch
12         else:
13             if v == tree.value:
14                 branch = tree.trueBranch
15             else:
16                 branch = tree.falseBranch
17         return classify(data, branch)

```