

决策树

一、 决策树原理

1. Introduction

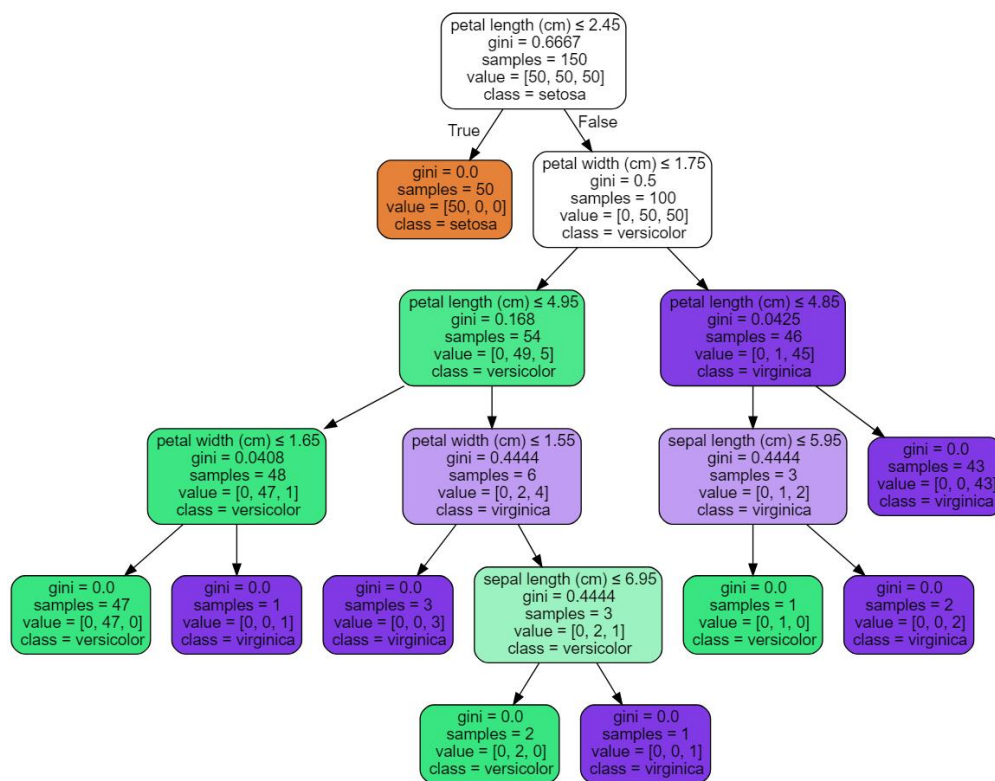


图 1: 决策树

决策树是从训练数据中学习得出一个树状结构的判别模型。决策树通过做出一系列决策来对数据进行划分。

决策树的核心在于如何选择划分字段。一种常见的方法是通过信息增益来进行选择，信息增益越大，则划分后的数据越纯净。

2. 信息熵、条件熵与信息增益

Definition 1 (信息熵). 信息熵用来表示信息量的大小，信息量越大，对应的熵值就越大，分类越不“纯净”。

$$\text{Ent}(D) = - \sum_{k=1}^{|Y|} p_k \log_2 p_k$$

Definition 2 (条件熵). 条件熵表示已知 a 事件发生情况下 D 的信息熵

$$\text{Ent}(D|a) = \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Ent}(D^v)$$

其中， a 有 v 个属性， D 在 a^v 属性上的数据集为 D^v

Definition 3 (信息增益). 信息增益表明在发生事件 A 前后信息熵的差异。信息增益越大，表明 A 对 D 的影响越大，划分后的数据集越纯净。

$$\text{Gain}(D|a) = \text{Ent}(D) - \text{Ent}(D|a)$$

但是通过信息增益划分有一个缺点，即即信息增益会偏向于取值较多的字段。为了克服信息增益指标的缺点，我们可以引入信息增益率：

Definition 4 (信息增益率).

$$\text{Gain_Ratio}_A(D) = \frac{\text{Gain}_A(D)}{\text{IV}(a)}$$

$$\text{IV}(A) = - \sum_{v=1}^V \frac{|D^v|}{|D|} \log_2 \frac{|D^v|}{|D|}$$

其中， $\text{IV}(A)$ 是事件 A 的信息熵

3. 基尼指数

由于信息增益率只能对离散的因变量进行分类，为了克服该缺点，可以引入 Gini 指数来选择划分属性

Definition 5 (Gini 指数).

$$\begin{aligned} \text{Gini}(D) &= \sum_{k=1}^{|Y|} p_k(1 - p_k) \\ &= 1 - \sum_{k=1}^{|Y|} p_k^2 \\ &= 1 - \sum_{k=1}^{|Y|} \left(\frac{|C_k|}{|D|} \right)^2 \end{aligned}$$

Definition 6 (条件 Gini 指数). 类似条件熵，条件 Gini 指数可以定义为

$$\text{Gini}(D|a) = \sum_{v=1}^V \frac{|D^v|}{|D|} \text{Gini}(D^v)$$

| 编号 | 色泽 | 根蒂 | 敲声 | 纹理 | 脐部 | 触感 | 好瓜 |
|----|----|----|----|----|----|----|----|
| 1 | 青绿 | 蜷缩 | 浊响 | 清晰 | 凹陷 | 硬滑 | 是 |
| 2 | 乌黑 | 蜷缩 | 沉闷 | 清晰 | 凹陷 | 硬滑 | 是 |
| 3 | 乌黑 | 蜷缩 | 浊响 | 清晰 | 凹陷 | 硬滑 | 是 |
| 4 | 青绿 | 蜷缩 | 沉闷 | 清晰 | 凹陷 | 硬滑 | 是 |
| 5 | 浅白 | 蜷缩 | 浊响 | 清晰 | 凹陷 | 硬滑 | 是 |
| 6 | 青绿 | 稍蜷 | 浊响 | 清晰 | 稍凹 | 软粘 | 是 |
| 7 | 乌黑 | 稍蜷 | 浊响 | 稍糊 | 稍凹 | 软粘 | 是 |
| 8 | 乌黑 | 稍蜷 | 浊响 | 清晰 | 稍凹 | 硬滑 | 是 |
| 9 | 乌黑 | 稍蜷 | 沉闷 | 稍糊 | 稍凹 | 硬滑 | 否 |
| 10 | 青绿 | 硬挺 | 清脆 | 清晰 | 平坦 | 软粘 | 否 |
| 11 | 浅白 | 硬挺 | 清脆 | 模糊 | 平坦 | 硬滑 | 否 |
| 12 | 浅白 | 蜷缩 | 浊响 | 模糊 | 平坦 | 软粘 | 否 |
| 13 | 青绿 | 稍蜷 | 浊响 | 稍糊 | 凹陷 | 硬滑 | 否 |
| 14 | 浅白 | 稍蜷 | 沉闷 | 稍糊 | 凹陷 | 硬滑 | 否 |
| 15 | 乌黑 | 稍蜷 | 浊响 | 清晰 | 稍凹 | 软粘 | 否 |
| 16 | 浅白 | 蜷缩 | 浊响 | 模糊 | 平坦 | 硬滑 | 否 |
| 17 | 青绿 | 蜷缩 | 沉闷 | 稍糊 | 稍凹 | 硬滑 | 否 |

图 2: 西瓜数据集

Example 1. 对于西瓜数据集（图 2），计算色泽的信息增益与基尼指数增益。信息增益

$$\begin{aligned}
 \text{Ent}(D) &= -\frac{8}{17} \log_2\left(\frac{8}{17}\right) - \frac{9}{17} \log_2\left(\frac{9}{17}\right) = 0.9975 \\
 \text{Ent}(D|\text{色泽}) &= \left(\frac{6}{17}\right)\left(-\frac{3}{6} \log \frac{3}{6} - \frac{3}{6} \log \frac{3}{6}\right) \\
 &\quad + \left(\frac{6}{17}\right)\left(-\frac{4}{6} \log \frac{4}{6} - \frac{2}{6} \log \frac{2}{6}\right) \\
 &\quad + \frac{5}{17}\left(-\frac{1}{5} \log \frac{1}{5} - \frac{4}{5} \log \frac{4}{5}\right) \\
 &= 0.353 + 0.324 + 0.212 = 0.889 \\
 \text{Gain}(D, \text{色泽}) &= \text{Ent}(D) - \text{Ent}(D|\text{色泽}) = 0.109
 \end{aligned}$$

基尼指数增益

$$\begin{aligned}
 \text{Gini}(D) &= 1 - \left(\frac{8^2}{17} + \frac{9^2}{17} \right) = 0.498 \\
 \text{Gini}(D|\text{色泽}) &= \frac{6}{17} \left(1 - \left(\frac{3^2}{6} + \frac{3^2}{6} \right) \right) \\
 &\quad + \frac{6}{17} \left(1 - \left(\frac{4^2}{6} + \frac{2^2}{6} \right) \right) \\
 &\quad + \frac{5}{17} \left(1 - \left(\frac{1^2}{5} + \frac{4^2}{5} \right) \right) \\
 &= 0.427 \\
 \text{Gini}(D, \text{色泽}) &= 0.071
 \end{aligned}$$

4. 剪枝

剪枝是决策树应对过拟合的主要手段。剪枝分为预剪枝和后剪枝。

Defination 7 (预剪枝). 将数据集划分为训练集和验证集，每次选择字段 a 并生成 V 个叶节点后，在验证集上计算精度。若精度大于划分前，则继续划分，否则停止划分

优点

- 降低过拟合风险
- 显著减少训练时间和测试时间开销

缺点

- 无法达到全局最优解

Defination 8 (后剪枝). 将数据集划分为训练集和验证集，从训练集学习到一颗完整的决策树。按照分类节点从外往里的顺序，每次将分类节点替换为叶节点，类别为该节点下最多的类。在验证集中计算精度，如果精度提升，则剪枝，否则不剪枝。

优点

- 欠拟合风险小
- 泛化性能优于预剪枝的决策树

缺点

- 训练时间开销和测试时间开销大于预剪枝决策树

二、 算法

1. 生成决策树

生成决策树的算法如 Algorithm1所示

Algorithm 1 决策树**Input:** $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$; $A = \{a_1, a_2, \dots, a_d\}$

```

1: Def TreeGenerate( $D, A$ )
2: 生成节点 node
3: if  $D$  中样本属于同一类别  $C$  then                                ▷ 已经分类正确
4:   将 node 标记为  $C$  类叶节点; return
5: end if
6: if  $A = \emptyset$  OR  $D$  中样本在  $A$  上取值相同 then                ▷ 不可再分条件
7:   将 node 标记为叶节点, 其类别标记为  $D$  中样本数最多的类; return
8: end if
9: 从  $A$  中选择划分最优属性  $a_*$ 
10: for  $a_*$  的每一个值  $a_*^v$  do
11:   为 node 生成一个分支; 令  $D_v$  表示  $D$  中在  $a_*$  上取值为  $a_*^v$  的样本子集
12:   if  $D_v = \emptyset$  then
13:     将分支节点标记为叶节点, 其类别标记为  $D$  中样本中最多的类; return    ▷ 为了泛化性
14:   else
15:     以 TreeGenerate( $D_v, A \setminus \{a_*\}$ ) 为分支节点
16:   end if
17: end for

```

Output: 以 node 为结点的一颗决策树

表 1: 节点类型

| 类别 | 属性 | | | | |
|------|-------|------------|------|----------|----------|
| | class | class_dict | attr | attrName | children |
| 根节点 | | ✓ | | ✓ | ✓ |
| 分类节点 | | ✓ | ✓ | ✓ | ✓ |
| 叶节点 | ✓ | | ✓ | | |

(1) 程序设计

决策树

共有 3 类结点：根节点，分类节点和叶节点。其属性值有一些差异，如表1所示。所有的节点共用一个 Node 类，但是会在 Node 类中表明该节点属于哪种类型。

其中，class 用于储存叶节点的分类类型，如是好瓜或者坏瓜，并用于决策树的预测。class_dict 用于储存在该节点对应的训练集中，各个分类的数量，如有 3 个好瓜 1 个坏瓜，在进行剪枝操作时，可以较为方便地判断剪枝后该节点的分类类型。attrName 储存属性名称，如“色泽”。attr 储存某一属性的值，如“青绿”，表示从父节点划分到该节点依据的是“色泽” = “青绿”。children 储存该节点的子节点。

函数

为了方便程序编写，共用到了 7 个函数，具体功能如表2所示

表 2: 函数

| 函数名称 | 注释 | |
|------------|---------|----------------------|
| | 对应伪代码行数 | 功能 |
| isSame | 3 | 判断 D 中的数据是否是同一个类别 |
| sameValue | 6 | 判断 D 中是否所有的样本都相同 |
| countClass | 7 | 统计 D 中每个类别的个数 |
| countAttr | 10 | 统计 D 中指定分类属性的各个属性的个数 |
| IF | 9 | 计算信息熵 |
| CE | 9 | 计算每个属性的信息增益 |
| IFD | 9 | 计算 D 上的信息熵 |

2. 后剪枝

Algorithm 2 后剪枝

Input: 根节点 root

```

1: 遍历决策树，获取所有分类节点和对应深度 nodeList
2: while nodeList 不为空 do
3:   选出 nodeList 中深度最深的节点 node
4:   计算剪枝前精度 acc1
5:   复制 node 中的内容                                ▷ 用于恢复节点内容
6:   剪枝: 将 class 赋值为 class_dict 中数量最大的类，将类型修改为叶节点，清空子节点
7:   计算剪枝后精度 acc2
8:   if acc2 ≤ acc1 then
9:     恢复节点内容
10:  end if
11:   从 nodeList 中删除 node 节点
12: end while

```

Output: root 节点

后剪枝的算法如 Algorithm2所示

3. 预剪枝

预剪枝的代码如 Algorithm3所示，预剪枝是在决策树的基础上进行修改的，在递归产生分支之前将程序截断，只往下分支一层，并比较精度。如果精度提升，则返回 True, 继续递归产生分支，否则返回 False, 程序结束。Algorithm3的代码在 Algorism1的 9-10 行之间。

三、 运行结果

在西瓜数据集上运行了决策树，其结果如图3所示

Algorithm 3 预剪枝

```

1: 计算分支前的精度 acc1
2: for  $a_*$  的每一个值  $a_*^v$  do
3:   为 node 生成一个分支; 令  $D_v$  表示  $D$  中在  $a_*$  上取值为  $a_*^v$  的样本子集
4:   if  $D_v = \emptyset$  then
5:     将分支节点标记为叶节点, 其类别标记为  $D$  中样本中最多的类; return    ▷ 为了泛化性
6:   else
7:     生成一个叶节点 node, 属性赋值为  $a_*^v$ , 类别标记为样本中最多的类别
8:   end if
9: end for
10: 计算分之后精度 acc2
11: if acc2 > acc1 then
12:   return True
13: else
14:   return False
15: end if

```

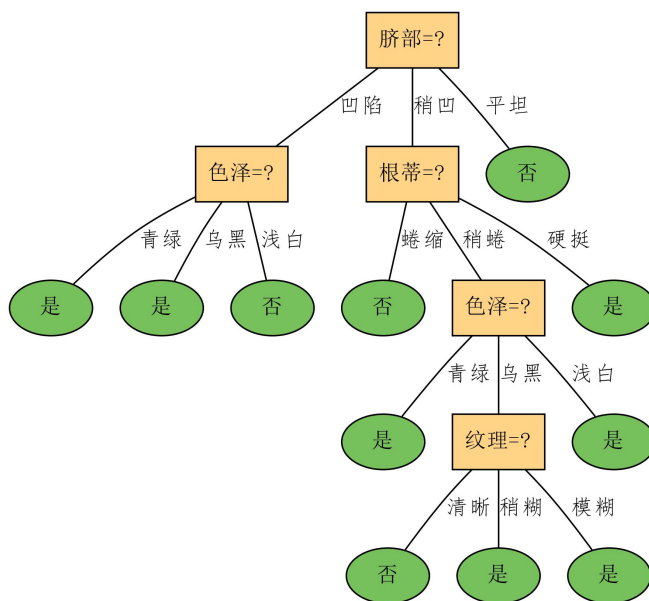


图 3: 不剪枝

附录

A 代码

```
1 import pandas as pd
```

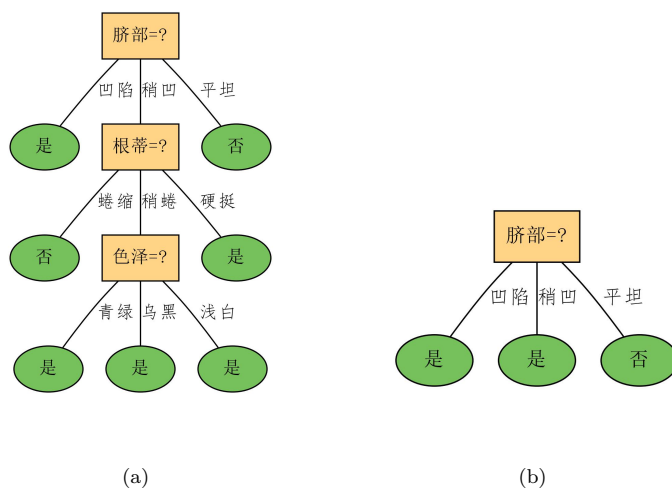


图 4: 预剪枝与后剪枝. (a) 后剪枝. (b) 预剪枝.

```

2 import numpy as np
3 import graphviz as gz
4
5
6 class Node:
7     def __init__(self, class_=None, type=None, attr=None, attr_name=None):
8         self.children = []      # 子节点
9         self.class_ = class_    # 叶节点类别
10        self.type = type        # 节点类型
11        self.attr = attr        # 父节点到该节点的属性
12        self.class_dict = {}    # 节点类别字典
13        self.attrName = attr_name
14    def copy(self, node):
15        node.children = self.children
16        node.class_ = self.class_
17        node.type = self.type
18        node.attr = self.attr
19        node.class_dict = self.class_dict
20        node.attrName = self.attrName
21    def print(self):
22        print(self.class_,self.class_dict,self.attrName,self.type,self.children)
23
24
25 class Decision_Tree:
26     def __init__(self, data, val, A, plot=True):
27         self.data = data
28         self.val = val
29         self.attr_num = {}
30         self.plot = plot
31         self.node = Node()
32         for i in range(len(A)):

```



```
33         self.attr_num[A[i]]=i
34
35     def train(self, prune = None):
36         if prune is None:
37             self.node = self.TreeGenerate(self.data,list(self.attr_num.keys()),
38                                             attr=None,root=True)
39
40         elif prune == 'post':
41             self.node = self.TreeGenerate(self.data,list(self.attr_num.keys()),
42                                             attr=None,root=True)
43
44             self.postPruning()
45         elif prune == 'pre':
46             self.node = self.TreeGenerate(self.data,list(self.attr_num.keys()),
47                                             attr=None,root=True,pre=True)
48
49         if self.plot:
50             graph = gz.Graph()
51             self.draw_DT(graph, self.node, 0)
52             graph.view()
53
54     def predict_node(self, node, x):
55         attrName = node.attrName
56         if len(node.children)==0:
57             return node.class_
58         for child in node.children:
59             if child.attr == x[self.attr_num[attrName]]:
60                 class_ = self.predict_node(child,x)
61                 return class_
62
63     def predict(self, x):
64         class_ = self.predict_node(self.node, x)
65         return class_
66
67     def accuracy(self):
68         cnt = 0
69         for i in range(len(self.val)):
70             if self.predict(self.val[i,:-1]) == self.val[i,-1]:
71                 cnt+=1
72         return cnt/len(self.val)
73
74     def postPruning(self):
75         divideList = self.getDivideNode(self.node,0)
76         print(divideList)
77         while len(divideList)>0:
78             depth = 0
79             index = 0
80             for i in range(len(divideList)):
81                 if divideList[i][1]>depth:
82                     depth = divideList[i][1]
83                     index = i
84             acc1 = self.accuracy()
85             node = divideList[index][0]
```

```
83     node_copy = Node()
84     node.copy(node_copy)
85     # ----- #
86     # 剪枝
87     # ----- #
88     node.type = 'leaf'
89     class_dict = node.class_dict
90     class_ = max(class_dict, key=class_dict.get)
91     node.class_ = class_
92     node.children = []
93     # ----- #
94     # 比较剪枝前后的准确率,如果效果没有变好,则恢复
95     # ----- #
96     acc2 = self.accuracy()
97     if acc2 <= acc1:
98         node_copy.copy(node)
99         print(f'acc:{acc1},不剪枝,节点信息: 字段 {node.attr},深度
100             {divideList[index][1]}')
101     else:
102         print(f'acc:{acc2},剪枝,节点信息: 字段 {node.attr},深度
103             {divideList[index][1]}')
104         divideList.pop(index)
105     # graph = gz.Graph()
106     # self.draw_DT(graph,self.node,0)
107     # graph.view()
108     return self.node
109
110 def getDivideNode(self, node, depth):
111     devideNodeList = []
112     if len(node.children)>0:
113         for child in node.children:
114             devideNodeList.extend(self.getDivideNode(child, depth+1))
115     if node.type == 'divide':
116         devideNodeList.append((node,depth))
117     return devideNodeList
118
119 def TreeGenerate(self, D, A, attr, root, pre=False):
120     node = Node()
121     if root:
122         self.node=node
123     node.attr = attr
124     # ----- #
125     # 判断是否都是同一个类,如果都是同一个类,分类完成,return
126     # ----- #
127     if self.isSame(D):
128         node.class_ = D[0][-1]
129         node.type = 'leaf'
130         return node
```

```

131     # ----- #
132     # 判断是否可以再分，如果不能再分，分类完成，return
133     # ----- #
134     if len(A)==0 or self.sameValue(D):
135         node.type = 'leaf'
136         count = self.countClass(D)
137         class_ = max(count, key=count.get)
138         node.class_ = class_
139         return node
140     # ----- #
141     # 计算类别的分布 用于后剪枝的时候判断divide节点
142     # 变成叶节点的时候应该归属到哪一个类
143     # ----- #
144     node.class_dict = self.countClass(D)
145     # ----- #
146     # 计算样本D上的信息熵
147     # ----- #
148     ifd = self.IFD(D)
149     # ----- #
150     # 挑选信息增益最大的属性a*
151     # ----- #
152     best_attr = ['',0]
153     for attr in A:
154         gain = ifd - self.CE(D,self.attr_num[attr])
155         if gain > best_attr[1]:
156             best_attr = [attr,gain]
157     # ----- #
158     # 统计D中a*每个类别的数量
159     # 注：a*的类别应该从整个数据集中统计，而非D中，但是数量要在D中统计
160     # ----- #
161     attr_dict = self.countAttr(D,self.attr_num[best_attr[0]])
162     print(attr_dict)
163     # ----- #
164     # 设置node的属性（分类节点，按照什么属性分类）
165     # ----- #
166     node.type = 'divide'
167     node.attrName = best_attr[0]
168     # ----- #
169     # 对于分类属性的每一个类，判断是否为空
170     # 如果为空则新分支标记为D中最多的类，否则重复上述步骤
171     # ----- #
172     ifContinue = True
173     if pre:
174         ifContinue = self.prePruning(node,attr_dict,best_attr,D)
175     if ifContinue:
176         for attr in attr_dict.keys():
177             if attr_dict[attr] == 0:
178                 count = self.countClass(D)
179                 class_ = max(count, key=count.get)
180                 node.children.append(Node(class_=class_,type='leaf',attr=attr))

```

```

181         else:
182             Dv = D[np.where(D[:,self.attr_num[best_attr[0]]]==attr),:][0]
183             A2 = A.copy()
184             A2.remove(best_attr[0])
185             print(best_attr[0],attr)
186             print(A2)
187             print('# ----- #')
188             node.children.append(self.TreeGenerate(Dv,A2,attr,False,pre=pre))
189         return node
190
191     def prePruning(self, node, attr_dict, best_attr, D):
192         if node.class_ is None:
193             node.class_ = max(node.class_dict, key=node.class_dict.get)
194         acc1 = self.accuracy()
195         for attr in attr_dict.keys():
196             if attr_dict[attr] == 0:
197                 count = self.countClass(D)
198                 class_ = max(count, key=count.get)
199                 node.children.append(Node(class_=class_,type='leaf',attr=attr))
200             else:
201                 Dv = D[np.where(D[:,self.attr_num[best_attr[0]]]==attr),:][0]
202                 class_dict = self.countClass(Dv)
203                 class_ = max(class_dict, key=class_dict.get)
204                 child = Node(type='leaf',class_=class_,attr=attr)
205                 node.children.append(child)
206             # graph = gz.Graph()
207             # self.draw_DT(graph, self.node, 0)
208             # graph.view()
209             acc2 = self.accuracy()
210             node.children = []
211             if acc2>acc1:
212                 print(f'划分前: {acc1}, 划分后: {acc2}, 继续划分')
213                 return True
214             else:
215                 print(f'划分前: {acc1}, 划分后: {acc2}, 禁止划分')
216                 return False
217
218     def isSame(self, D):
219         '''
220         判断D中的数据是否是同一个类别
221         :param D: 数据集的一个子集
222         :return: bool
223         '''
224         class_ = D[0][-1]
225         for i in range(len(D)):
226             if D[i][-1] != class_:
227                 return False
228         return True
229
230     def countClass(self, D):

```

```
231     '''
232     统计D中每个类别的个数
233     :param D: 数据集的一个子集
234     :return: dict{class:cnt}
235     '''
236     class_dict = {}
237     for i in range(len(D)):
238         class_ = D[i][-1]
239         class_dict.setdefault(class_,0)
240         class_dict[class_] += 1
241     return class_dict
242
243 def sameValue(self, D):
244     '''
245     判断D中是否所有的样本都相同
246     :param D: 数据集的一个子集
247     :return: bool
248     '''
249     v = D[0]
250     for i in range(D.shape[0]):
251         if (v != D[i]).any():
252             return False
253     return True
254
255 def countAttr(self, D, attr_idx):
256     '''
257     统计D中指定分类属性的各个属性的个数
258     :param D: 数据集的一个子集
259     :param attr_idx: 属性编号，输入时生成
260     :return: dict{attr:cnt}
261     '''
262     attr_dict = {}
263     for i in range(len(self.data)):
264         attr = self.data[i][attr_idx]
265         attr_dict.setdefault(attr,0)
266     for i in range(len(D)):
267         attr = D[i][attr_idx]
268         # attr_dict.setdefault(attr, 0)
269         attr_dict[attr] += 1
270     return attr_dict
271
272 def IF(self, x):
273     '''
274     计算信息熵
275     :param x: ndarray
276     :return: 信息熵
277     '''
278     x[np.where(x==0)] = 1
279     return -np.sum(x*np.log2(x))
280
```

```
281 def CE(self, D, attr_idx):
282     '''
283     计算每个属性的信息增益
284     :param D: 数据集的一个子集
285     :param attr_idx: 属性编号，输入时生成
286     :return: 信息增益
287     '''
288     attr_dict = self.countAttr(D, attr_idx)
289     EF = 0
290     for attr in attr_dict.keys():
291         Dv = D[np.where(D[:, attr_idx] == attr), :] [0]
292         ifd = self.IFD(Dv)
293         EF += len(Dv) / len(D) * ifd
294     return EF
295 def IFD(self, D):
296     '''
297     计算D上的信息熵
298     :param D: 数据集的一个子集
299     :return: 信息熵
300     '''
301     class_dict = self.countClass(D)
302     class_dict = np.array(list(class_dict.values()))
303     class_dict = class_dict / np.sum(class_dict)
304     return self.IF(class_dict)
305
306 def draw_DT(self, graph, node, nodeid):
307     '''
308     通过graphviz进行决策树可视化（递归）
309     :param graph: 图
310     :param node: 上一级节点（决策树的节点）
311     :param nodeid: 上一级节点编号（graph的节点编号）
312     :return: nodeid
313     '''
314     if len(node.children) == 0:
315         graph.node(str(nodeid), node.class_, fontname='SimSun',
316                    style='filled', color='black', fillcolor='#78C25E')
317     else:
318         graph.node(str(nodeid), node.attrName + '=?',
319                    fontname='SimSun', shape='box',
320                    style='filled', color='black', fillcolor='#FFD588')
321     node.print()
322     curr_nodeid = nodeid
323     # print(node.class_, nodeid)
324     if node.children:
325         for child in node.children:
326             child_id = nodeid + 1
327             nodeid = self.draw_DT(graph, child, child_id)
328             graph.edge(str(curr_nodeid), str(child_id), label=child.attr,
329                        fontname='FangSong')
329     return nodeid
```

