

Python 实现基于C4.5算法的决策树算法

代码见：<https://github.com/KevinTungs/C4.5-Python.git>

背景信息

什么是决策树？

决策树可以想象成是一个巨大的“如果-那么”规则集合，组织成一个树形结构。它从一个根节点开始，每一次决策就是选择一个分支，直到达到叶节点，叶节点代表了最终的决策结果。你可以将其视为一个问答游戏：每个内部节点代表一个问题，每个分支代表这个问题的一个可能答案，而每个叶节点则代表游戏的最终结果。比如，在判断一个动物是什么的决策树中，一个节点可能问：“它有羽毛吗？”如果答案是“是”，那么你沿着“是”的分支走下去，接下来的问题可能是“它会飞吗？”；如果答案是“否”，你可能就走向了另一个分支，比如进一步询问它是否有鳞片。通过这种方式，决策树能够帮助我们根据一系列问题的答案来分类或者预测事物。

专家系统

专家系统，是人工智能的一个分支，模仿人类专家的决策能力，提供专门领域内的决策支持。它主要由三个部分组成：知识库（存储领域知识和事实）、推理机（模拟专家的推理过程，根据知识库中的规则对问题进行推理）、用户界面（与用户交互，获取问题并展示推理结果）。专家系统依赖于领域专家提供的规则和知识，通过逻辑推理来解决复杂的问题。

决策树和专家系统的决策过程很像。下面是他们的差异：

1. 相似性：

- **模拟决策过程：**决策树和专家系统都旨在模拟（或辅助）人类的决策过程。决策树通过数据学习产生决策路径，而专家系统通过规则和逻辑推理来进行决策。
- **规则表示：**在某种程度上，决策树的路径可以视为决策规则，这与专家系统使用的基于规则的推理相似。

2. 差异：

- **知识来源：**决策树的知识是通过分析训练数据自动学习得到的，而专家系统的知识来自于人类专家的直接输入。
- **构建过程：**决策树的构建是一个自动化的过程，主要依赖于算法从数据中学习；而专家系统的构建需要领域专家明确定义规则和知识。
- **应用领域：**决策树广泛应用于各种分类和回归任务中，专家系统则更多应用于特定领域的问题解决，如医疗诊断、金融分析等。

决策树的最优划分属性选择

决策树学习的关键是如何选择最优划分属性。一般而言，我们希望决策树的分支节点所包含的样本尽可能属于同一类别，也就是说节点的“纯度”越来越高。

最优属性选择的标准主要包括信息增益（Information Gain）、信息增益比（Gain Ratio）、基尼指数（Gini Index）等，它们分别适用于不同的决策树算法，如ID3、C4.5和CART。

信息增益（ID3算法）

信息增益是最常用的选择方法之一，它基于香农的信息论。

信息熵

假定当前样本集合 D 中第 k 类样本所占的比例为 p_k ，则 D 的信息熵定义为：

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

信息熵满足下面三个性质：

1. 单调性，发生概率越高的事件，其携带的信息量越低；
2. 非负性，信息熵可以看作为一种广度量，非负性是一种合理的必然；
3. 累加性，即多随机事件同时发生存在的总不确定性的量度是可以表示为各事件不确定性的量度的和，这也是广度量的一种体现。

信息增益

信息增益是基于熵的概念，熵是一个衡量数据集随机性（或不确定性）的指标。信息增益测量的是在知道属性 A 的信息之后数据集 D 的熵减少了多少。

某个属性在数据集中的信息增益的定义为：

$$Gain(D, A) = Ent(D) - \sum_{v=1}^V \frac{|D_v|}{|D|} Ent(D_v)$$

其中：

1. D 是数据集
2. D_v 是 D 中 A 取值为 V 的子集
3. $\frac{|D_v|}{|D|}$ 相当于该属性的权重
4. A 是某个属性
5. $Values(A)$ 是该属性的所有可能值

信息增益是指父节点的熵与分割后各子节点熵的加权和的差值。信息增益相当于使用这个属性对样本集进行分类所获得的“信息增益”，信息增益越大，意味着使用属性进行分类的结果越好。

信息增益率（C4.5算法）

分裂信息（IV）

分裂信息（Intrinsic Value，IV）衡量的是根据属性 a 对数据集 D 进行划分时所带来的“信息量”或“复杂度”。分裂信息越高，意味着属性 a 的取值越分散。计算公式是：

$$IV(a) = \sum_{v \in Values(a)} -\frac{|D_v|}{|D|} \log_2 \frac{|D_v|}{|D|}$$

这里， \log_2 是以2为底的对数，用于计算信息的传统单位（比特）

增益率

当我们讨论决策树中的增益率（Gain Ratio），我们是在讨论一种特定的属性选择方法，这种方法尤其在C4.5算法中得到了广泛的使用。这种方法的目的 是克服原始信息增益方法偏好选择具有大量值的属性的问题。增益率通过对信息增益（Gain）进行标准化来实现这一点。

信息增益比是在信息增益的基础上进行改进的，它考虑了属性的分裂信息（Split Information），以解决信息增益偏好选择取值多的属性的问题。

增益率比定义为信息增益与分裂信息的比值：

$$Gain_ratio(D, a) = \frac{Gain(D, a)}{IV(a)}$$

属性 a 的固有值 $IV(a)$ 衡量了按属性 a 分裂数据集 D 的信息量。

通过这种方式，增益率既考虑了属性带来的信息增益，又惩罚了属性值过多导致的分裂过度细致。这种方法使得C4.5算法在处理具有不同数量值的属性时表现得更加均衡和有效。

至于为什么C4.5算法要叫C4.5，这其实是算法发展的历史原因。C4.5算法是由Ross Quinlan发展的，作为ID3算法的后续和改进版。ID3（Iterative Dichotomiser 3）是Quinlan在1986年提出的一个决策树算法，它是最早的决策树算法之一，使用信息增益作为属性选择的标准。然而，ID3有一些局限性，比如它不能处理连续属性、缺失值，也没有剪枝过程，容易过拟合。

为了克服这些局限，Quinlan后来提出了C4.5算法，这是对ID3的扩展和改进。C4.5引入了信息增益比来选择属性，能够处理连续属性、缺失值，并且加入了剪枝过程来避免过拟合。C4.5的名字来源于Quinlan的研究和开发过程：它是“C for Classification”的缩写，加上4.5表示这个版本号或者是这个方法的进化版，表明它是在ID3基础上的一次重大改进和迭代。

基尼指数（CART算法）

基尼系数

基尼指数是另一种选择属性的方法，用于CART决策树（既可以用于分类树，也可以用于回归树）。基尼指数衡量的是从数据集中随机选取两个样本，其类别标签不一致的概率。因此，基尼指数越小，数据集的纯度越高。对于二分类问题，基尼指数的计算公式是：

$$Gini(D) = 1 - \sum_{i=1}^m p_i^2$$

这个表达式是在计算数据集 D 的不纯度。对于每一个类别 i ，你先计算其在数据集中的比例 p_i ，然后求其平方 p_i^2 。平方项的意义在于，如果某个类别的样本占比更高，它对整体不纯度的贡献就更小，因为我们希望单一类别的样本占比越高越好。之后，我们将所有类别的 p_i^2 相加，得到的和表示数据集在完全分到一个类别时的“理想”状态。从1中减去这个和，我们就得到了数据集的不纯度度量——基尼指数。

对于选择最优分割属性，会计算每个属性分割后的加权基尼指数，选择使得基尼指数降低最多的属性。

不同的选择标准适用于不同的决策树算法和应用场景，没有绝对的优劣，需要根据实际问题 and 数据特性选择合适的标准。

基尼指数（Gini Index）是决策树算法中用于属性选择的一个重要标准，特别是在CART

（Classification And Regression Trees）算法中。它反映了从数据集中随机选取两个样本，这两个样本属于不同类别的概率。如果我们从一个数据集中随机抽取两次（不放回），基尼指数越低意味着抽取到不同类别样本的概率越小，即数据集的纯度越高。

- 为什么要用平方：使用平方而非其他幂次是为了加强高比例值的影响（即使数据集更倾向于某一类），并且平方保证了计算过程中数值的非负性和合理性。
- 为什么越小越好：基尼指数越小，意味着数据集中的样本越倾向于属于同一个类别，即数据集的纯度越高。在构建决策树时，我们希望通过选择合适的属性来分割数据集，使得分割后的子数据集具有更低的基尼指数，也就是说，子数据集在类别上更加“纯净”。

在CART算法中的应用

在CART算法中，基尼指数用于评估在特定属性上分割数据集的效果。算法会计算每个属性的每个可能分割点分割后的加权基尼指数，然后选择基尼指数降低最多的分割点作为最佳分割点。通过这种方式，CART算法逐渐构建出能够高效分类或回归的决策树模型。

Python实现

C4.5算法在Python中的实现包括以下步骤：

1. 计算数据集目标属性的信息熵。
2. 对每个属性（除目标属性外），计算其增益比。
3. 选择增益比最高的属性来分割数据集。
4. 对分割产生的每个数据子集递归重复此过程，直到满足停止条件（例如，所有元素属于同一类别，没有更多属性可以分割，或者达到最大树深度）。

代码

```

1 import pandas as pd
2 import math
3 from typing import Any, Dict, List
4 from graphviz import Digraph
5
6
7 # 定义计算熵、条件熵、信息增益、增益率以及选择最优划分属性的函数
8
9 # 计算数据集D的熵
10 def calc_entropy(D: pd.DataFrame) -> float:
11     """
12     计算数据集D的熵
13     :param D: 数据集
14     :return: 数据集D的熵
15     """
16     labels = D.iloc[:, -1]
17     label_counts = labels.value_counts()
18     entropy = -sum((count / len(D)) * math.log2(count / len(D)) for count in
19 label_counts)
19     return entropy
20
21
22 # 计算属性A的分裂信息 (IV)
23 def calc_conditional_entropy(D: pd.DataFrame, A: str) -> float:
24     """
25     计算属性A的条件熵
26     :param D: 数据集
27     :param A: 属性
28     :return: 属性A的分裂信息
29     """
30     A_values = D[A].unique()
31     conditional_entropy = 0.0
32     for v in A_values:
33         sub_D = D[D[A] == v]
34         conditional_entropy += (len(sub_D) / len(D)) * calc_entropy(sub_D)
35     return conditional_entropy
36
37
38 # 计算连续属性在给定划分点下的信息增益
39 def calc_information_gain_continuous(D: pd.DataFrame, attribute: str,
40 split_point: float) -> float:
41     """
42     计算连续属性在给定划分点下的信息增益
43     :param D: 数据集
44     :param attribute: 属性名称
45     :param split_point: 划分点
46     :return: 信息增益

```

```

46     """
47     # 按照划分点分割数据集
48     D1 = D[D[attribute] <= split_point]
49     D2 = D[D[attribute] > split_point]
50
51     # 计算原数据集的熵
52     entropy_before = calc_entropy(D)
53
54     # 计算划分后的加权熵
55     entropy_after = (len(D1) / len(D)) * calc_entropy(D1) + \
56                     (len(D2) / len(D)) * calc_entropy(D2)
57
58     # 计算信息增益
59     information_gain = entropy_before - entropy_after
60     return information_gain
61
62
63 # 计算信息增益
64 def calc_information_gain(D: pd.DataFrame, A: str, split_point: float = None) -
65 > float:
66     """
67     计算信息增益，自动处理连续和分类属性。连续属性需要划分点。
68     :param D: 数据集
69     :param A: 属性
70     :param split_point: 连续属性的划分点，对于分类属性，这个参数不使用
71     :return: 属性A的信息增益
72     """
73     # 检查属性A是连续还是分类的
74     if pd.api.types.is_numeric_dtype(D[A]) and split_point is not None:
75         # 对于连续属性，使用划分点来计算信息增益
76         return calc_information_gain_continuous(D, A, split_point)
77     else:
78         # 对于分类属性，使用原始的信息增益计算方法
79         return calc_entropy(D) - calc_conditional_entropy(D, A)
80
81 # 计算属性A的增益率
82 def calc_gain_ratio(D: pd.DataFrame, A: str) -> float:
83     """
84     计算属性A的增益率
85     :param D: 数据集
86     :param A: 属性
87     :return: 属性A的增益率
88     """
89     information_gain: float = calc_information_gain(D, A)
90     A_values: list = D[A].unique()

```

```

91     iv: float = -sum((len(D[D[A] == v]) / len(D)) * math.log2(len(D[D[A] ==
    v]) / len(D)) for v in A_values)
92     return information_gain / iv if iv != 0 else 0.0
93
94
95 # 选择最优划分属性
96 def choose_best_feature(D: pd.DataFrame) -> str:
97     """
98     选择最优划分属性
99     :param D: 数据集
100    :return: 最优划分属性
101    """
102    features = D.columns[:-1]
103    gain_ratios = {feature: calc_gain_ratio(D, feature) for feature in
    features}
104    return max(gain_ratios, key=gain_ratios.get)
105
106
107 # 寻找连续属性的最优划分点
108 def find_best_split_point_for_continuous_attribute(D: pd.DataFrame, attribute:
    str) -> float:
109     """
110     寻找连续属性的最优划分点
111     :param D: 数据集
112     :param attribute: 属性名称
113     :return: 最优划分点
114     """
115     sorted_values = D[attribute].sort_values().unique()
116     split_points = [(sorted_values[i] + sorted_values[i + 1]) / 2 for i in
    range(len(sorted_values) - 1)]
117     max_gain = -float('inf')
118     best_split = None
119     for split_point in split_points:
120         # 计算每个可能划分点的信息增益
121         gain = calc_information_gain_continuous(D, attribute, split_point)
122         if gain > max_gain:
123             max_gain = gain
124             best_split = split_point
125     return best_split
126
127
128 # 递归构建决策树 (不包含剪枝)
129 def create_decision_tree(D: pd.DataFrame, features: List[str],
    continuous_attributes: List[str]) -> dict:
130     """
131     递归构建决策树, 适应连续变量处理
132     :param D: 数据集

```

```

133     :param features: 特征列表
134     :param continuous_attributes: 连续属性列表
135     :return: 决策树
136     """
137     # 基本终止条件
138     class_counts = D.iloc[:, -1].value_counts()
139     if len(class_counts) == 1:
140         return class_counts.index[0]
141     if not features:
142         return class_counts.idxmax()
143
144     # 选择最优属性及其划分点 (如果是连续的)
145     best_gain = -float('inf')
146     best_feature = None
147     split_point = None
148     for feature in features:
149         if feature in continuous_attributes:
150             # 对于连续属性, 找到最优划分点
151             point = find_best_split_point_for_continuous_attribute(D, feature)
152             gain = calc_information_gain(D, feature, point) # 假设这个函数现在可
以处理连续属性
153             if gain > best_gain:
154                 best_gain = gain
155                 best_feature = feature
156                 split_point = point
157         else:
158             # 对于离散属性, 正常处理
159             gain = calc_information_gain(D, feature) # 离散属性的信息增益计算
160             if gain > best_gain:
161                 best_gain = gain
162                 best_feature = feature
163
164     # 根据选择的最优属性分割数据集
165     tree = {best_feature: {}}
166     if best_feature in continuous_attributes:
167         # 处理连续属性的分割
168         left_D = D[D[best_feature] <= split_point]
169         right_D = D[D[best_feature] > split_point]
170         tree[best_feature]['≤' + str(split_point)] =
create_decision_tree(left_D,
171
172         for f in features if f != best_feature],
continuous_attributes)
173         tree[best_feature]['>' + str(split_point)] =
create_decision_tree(right_D,

```

[f


```

174                                     [f
    for f in features if f != best_feature],
175
    continuous_attributes)
176     else:
177         # 处理离散属性的分割
178         for value in D[best_feature].unique():
179             sub_D = D[D[best_feature] == value]
180             tree[best_feature][value] = create_decision_tree(sub_D, [f for f in
    features if f != best_feature],
181
    continuous_attributes)
182
183     return tree
184
185
186 # 后剪枝函数
187 def post_pruning(tree: Dict[str, Any], D: pd.DataFrame, features: List[str]) -
    > Dict[str, Any]:
188     """
189     对决策树进行后剪枝
190     :param tree: 当前决策树
191     :param D: 数据集
192     :param features: 特征集
193     :return: 剪枝后的决策树
194     """
195     # 检查树是否是叶节点
196     if not isinstance(tree, dict):
197         return tree
198
199     # 遍历树中的每个节点
200     for feature, branches in tree.items():
201         for value, subtree in branches.items():
202             # 递归剪枝子树
203             subtree = post_pruning(subtree, D[D[feature] == value], [f for f in
    features if f != feature])
204             tree[feature][value] = subtree
205
206     # 尝试剪枝当前节点
207     if all(not isinstance(subtree, dict) for subtree in
    tree[feature].values()):
208         # 计算剪枝前后的准确性
209         accuracy_before_pruning = calc_accuracy(tree, D)
210         # 将当前节点替换为最常见的类
211         most_common_class = D.iloc[:, -1].mode()[0]
212         pruned_tree = most_common_class
213         accuracy_after_pruning = calc_accuracy(pruned_tree, D)

```

```

214         # 如果剪枝后准确性不降低, 则进行剪枝
215         if accuracy_after_pruning >= accuracy_before_pruning:
216             return pruned_tree
217
218     return tree
219
220
221 # 对单个实例进行预测
222 def predict(tree: Dict[str, Any], instance: Dict[str, Any]) -> Any:
223     """
224     对单个实例进行预测
225     :param tree: 决策树
226     :param instance: 单个数据实例
227     :return: 预测结果
228     """
229     if not isinstance(tree, dict):
230         return tree
231     root = next(iter(tree))
232     subtree = tree[root]
233     value = instance[root]
234     if value in subtree:
235         return predict(subtree[value], instance)
236     else:
237         return None
238
239
240 # 计算决策树在数据集D上的准确性
241 def calc_accuracy(tree: Dict[str, Any], D: pd.DataFrame) -> float:
242     """
243     计算决策树在数据集D上的准确性
244     :param tree: 决策树
245     :param D: 数据集
246     :return: 准确性
247     """
248     correct_predictions = 0
249     for _, row in D.iterrows():
250         if predict(tree, row) == row.iloc[-1]:
251             correct_predictions += 1
252     return correct_predictions / len(D)
253
254
255 # 绘制决策树
256 def plot_decision_tree(tree, parent_name=None, edge=None, graph=None):
257     if graph is None:
258         graph = Digraph(comment='Decision Tree', format='png')
259
260     if not isinstance(tree, dict):

```

```

261     node_name = f"Leaf_{tree}"
262     graph.node(node_name, label=str(tree), shape='ellipse')
263     if parent_name is not None:
264         graph.edge(parent_name, node_name, label=str(edge))
265     else:
266         for idx, (feature, branches) in enumerate(tree.items()):
267             node_name = f"Node_{feature}_{idx}"
268             if parent_name is None:
269                 graph.node(node_name, label=str(feature))
270             else:
271                 graph.edge(parent_name, node_name, label=str(edge))
272                 graph.node(node_name, label=str(feature))
273
274             for value, subtree in branches.items():
275                 plot_decision_tree(subtree, node_name, value, graph)
276
277     return graph
278
279
280 # 载入数据
281 data = dict(
282     色泽=['青绿', '乌黑', '乌黑', '青绿', '浅白', '青绿', '乌黑', '乌黑', '青绿', '浅
283         白', '浅白', '青绿', '浅白',
284         '乌黑', '浅白', '青绿'],
285     根蒂=['蜷缩', '蜷缩', '蜷缩', '蜷缩', '蜷缩', '稍蜷', '稍蜷', '稍蜷', '硬挺', '硬
286         挺', '蜷缩', '稍蜷', '稍蜷',
287         '稍蜷', '蜷缩', '蜷缩'],
288     敲声=['浊响', '沉闷', '浊响', '沉闷', '浊响', '浊响', '浊响', '浊响', '清脆', '清
289         脆', '浊响', '浊响', '沉闷',
290         '浊响', '浊响', '沉闷'],
291     纹理=['清晰', '清晰', '清晰', '清晰', '清晰', '清晰', '清晰', '稍糊', '清晰', '清晰', '模
292         糊', '模糊', '稍糊', '稍糊',
293         '清晰', '模糊', '稍糊'],
294     脐部=['凹陷', '凹陷', '凹陷', '凹陷', '凹陷', '稍凹', '稍凹', '稍凹', '平坦', '平
295         坦', '平坦', '凹陷', '凹陷',
296         '稍凹', '平坦', '稍凹'],
297     触感=['硬滑', '硬滑', '硬滑', '硬滑', '硬滑', '软粘', '软粘', '硬滑', '软粘', '硬
298         滑', '软粘', '硬滑', '硬滑',
299         '软粘', '硬滑', '硬滑'],
300     含糖率=[0.460, 0.376, 0.264, 0.318, 0.215, 0.237, 0.149, 0.211, 0.267,
301             0.057, 0.099, 0.161, 0.198, 0.370, 0.042, 0.103],
302     好瓜=['是', '是', '是', '是', '是', '是', '是', '是', '是', '否', '否', '否', '否',
303         '否', '否', '否', '否']
304 )
305
306 df = pd.DataFrame(data)
307

```

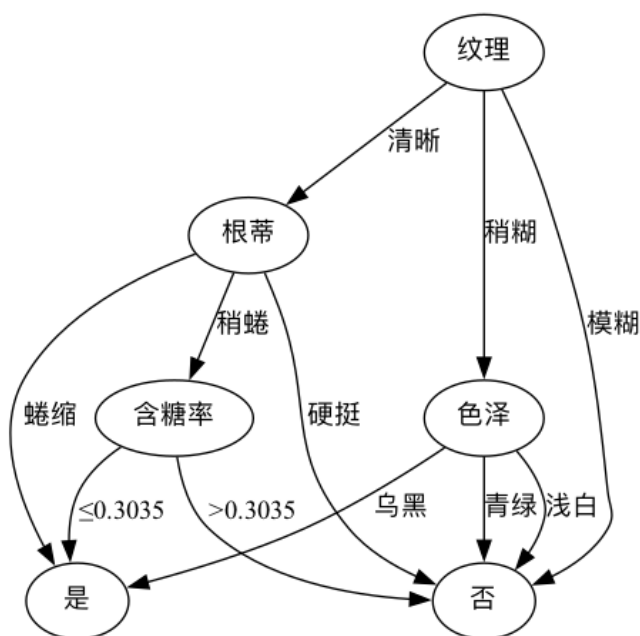
```

300 # 特征列表
301 features = list(df.columns[:-1])
302
303 # 连续属性列表
304 continuous_attributes = ['含糖率']
305
306 # 创建决策树
307 decision_tree = create_decision_tree(df, features, continuous_attributes)
308
309 # 输出剪枝前决策树
310 graph_decision_tree = plot_decision_tree(decision_tree)
311
312 # 对决策树进行后剪枝
313 pruned_tree = post_pruning(decision_tree, df, features)
314
315 # 输出剪枝后决策树
316 graph_pruned_tree = plot_decision_tree(pruned_tree)
317
318 # 保存并显示图像
319 graph_pruned_tree.render(filename='pruned_tree', directory='.', view=True) # 将
    文件保存在当前工作目录
320 graph_decision_tree.render(filename='decision_tree', directory='.', view=True)
    # 将文件保存在当前工作目录

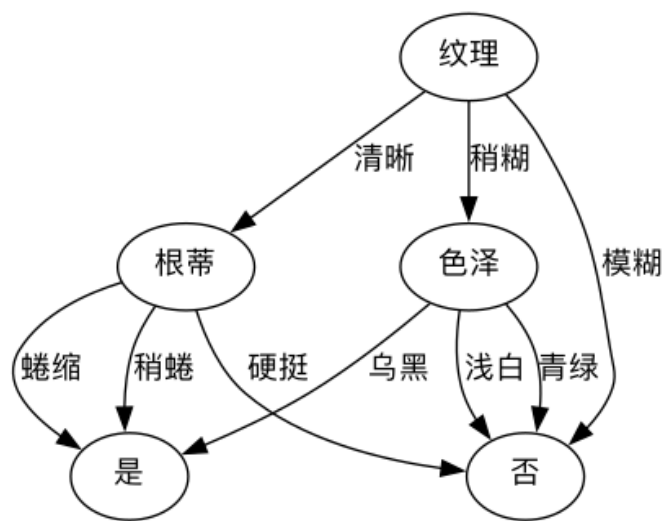
```

结果

未剪枝结果：



剪枝后结果：



说明含糖率并不能够为判断是否为好瓜提供有用信息

代码解释

包导入

```
1 import pandas as pd
2 import math
3 from typing import Any, Dict, List
4 from graphviz import Digraph
```

Pandas 库

Pandas 是一个流行的 Python 库，用于数据处理和分析。它提供了高性能、易于使用的数据结构和数据操作工具，使得在 Python 中进行数据操作更加简单和高效。Pandas 最核心的数据结构是 DataFrame，它类似于电子表格或数据库中的表格，能够存储和处理二维数据，包括具有标签的行和列。除了 DataFrame，Pandas 还提供了 Series 数据结构，用于存储一维数据，以及各种数据操作功能，如数据的读取、写入、选择、过滤、合并、重塑、分组、聚合等。

一些 Pandas 的主要特点包括：

- 数据结构：**DataFrame 和 Series 是 Pandas 中最重要的数据结构，提供了灵活、高效的数据存储和操作方式。
- 数据导入导出：**Pandas 支持从多种数据源中导入数据，如 CSV 文件、Excel 表格、SQL 数据库、JSON 格式等，并且可以将数据导出到这些格式。
- 数据处理：**Pandas 提供了丰富的数据处理功能，包括数据清洗、处理缺失值、重复值、数据转换、数据合并等。
- 数据分析：**Pandas 提供了统计计算、分组聚合、数据透视表等功能，方便进行数据分析和探索性数据分析（EDA）。

5. 时间序列处理：Pandas 对时间序列数据有着良好的支持，包括日期范围生成、时间索引、时间频率转换等功能。

Pandas 是 Python 数据科学生态系统中不可或缺的一部分，它使得数据分析和数据处理变得更加简单、高效，并且广泛应用于数据科学、机器学习、金融分析、业务分析等领域。

graphviz 库

graphviz 是一个用于创建和可视化图形的 Python 库，而 `Digraph` 则是其中的一个类，用于创建有向图（Directed Graph）对象。

下面是对 `graphviz` 中 `Digraph` 类的一些说明：

- 1. 有向图（Digraph）：**有向图是图论中的一种重要概念，它由一组顶点和一组有向边组成，每条边连接两个顶点，并且有一个方向。`Digraph` 类允许你创建有向图，并且可以在图中添加顶点和有向边。
- 2. 创建图：**要创建一个有向图对象，你需要导入 `graphviz` 库，然后使用 `Digraph` 类来实例化一个对象。例如，`graph = Digraph()` 就创建了一个空的有向图对象。
- 3. 添加节点和边：**在有向图中，你可以通过 `node()` 方法添加节点（顶点），通过 `edge()` 方法添加有向边。例如，`graph.node('A')` 就添加了一个名为 'A' 的节点，`graph.edge('A', 'B')` 就添加了一个从节点 'A' 到节点 'B' 的有向边。
- 4. 可视化：**一旦你添加了节点和边，你可以使用 `render()` 方法将图可视化为各种格式，如图像文件（PNG、PDF 等）或者文本形式。例如，`graph.render('graph')` 将生成一个名为 'graph' 的图形文件。

`Digraph` 类提供了一种方便的方式来创建和操作有向图，而 `graphviz` 库则为我们提供了各种可视化图形的工具和功能，使得我们能够更好地理解 and 展示数据之间的关系。

相关函数

计算数据集D的熵

```
1 # 计算数据集D的熵
2 def calc_entropy(D: pd.DataFrame) -> float:
3     """
4     计算数据集D的熵
5     :param D: 数据集
6     :return: 数据集D的熵
7     """
8     labels = D.iloc[:, -1]
9     label_counts = labels.value_counts()
10    entropy = -sum((count / len(D)) * math.log2(count / len(D)) for count in
11                  label_counts)
12    return entropy
```

1. 提取标签列

```
1 labels = D.iloc[:, -1]
```

这行代码从数据集 `D` 中提取了最后一列，即数据集中的标签列。假设数据集的最后一列包含了类别标签，表示样本所属的分类。

可以使用 -1 这样的负索引

2. 统计了标签列中每个类别出现的次数

```
1 label_counts = labels.value_counts()
```

这行代码统计了标签列中每个类别出现的次数，并将结果存储在 `label_counts` 中。
`value_counts()` 方法返回一个 Series 对象，其中索引是类别标签，值是该类别出现的次数。

3. 计算数据集 `D` 的熵

```
1 entropy = -sum((count/len(D)) * math.log2(count/len(D)) for count in  
label_counts)
```

这行代码计算了数据集 `D` 的熵。具体来说，它使用了熵的定义公式：

$$\text{Entropy}(D) = - \sum_{i=1}^k p_i \log_2(p_i)$$

这行代码使用了 Python 的列表推导式（List Comprehension）和生成器表达式（Generator Expression），结合了 `sum()` 函数和 `math.log2()` 函数来计算数据集 `D` 的熵。

- `(count/len(D)) * math.log2(count/len(D)) for count in label_counts`：这部分是一个生成器表达式，它会遍历 `label_counts` 中的每个元素 `count`，并对每个 `count` 执行 `(count/len(D)) * math.log2(count/len(D))` 这个表达式，得到一个结果。
- `sum(...)`：这是 Python 内置的求和函数 `sum()`，它将会对生成器表达式生成的结果进行求和。

计算属性A的分裂信息（IV）

```
1 # 计算属性A的分裂信息（IV）
```

```

2 def calc_conditional_entropy(D: pd.DataFrame, A: str) -> float:
3     """
4     计算属性A的条件熵
5     :param D: 数据集
6     :param A: 属性
7     :return: 属性A的分裂信息
8     """
9     A_values = D[A].unique()
10    conditional_entropy = 0.0
11    for v in A_values:
12        sub_D = D[D[A] == v]
13        conditional_entropy += (len(sub_D) / len(D)) * calc_entropy(sub_D)
14    return conditional_entropy

```

这段代码是用来计算给定数据集 `D` 中某个属性 `A` 的条件熵的函数。条件熵用于衡量在给定属性 `A` 的情况下，数据集 `D` 的不确定性。下面是对代码的解释：

1. 取出属性 `A` 的所有取值

```
1 A_values = D[A].unique()
```

这行代码从数据集 `D` 中取出属性 `A` 的所有取值，并使用 `unique()` 方法去除重复值。这样，`A_values` 就是属性 `A` 的所有可能取值的列表。

`unique()` 方法是 Pandas 库中 Series 对象的一个方法，用于获取 Series 中的唯一值，并返回一个包含这些唯一值的数组。下面是对 `unique()` 方法的介绍：

语法：

```
1 Series.unique()
```

这里的 `Series` 是一个 Pandas Series 对象。

功能：

`unique()` 方法用于返回 Series 中的唯一值，即去除重复值后的数组。如果 Series 中有多个相同的值，则只保留一个。

返回值：

返回一个包含 Series 中唯一值的数组。数组中的值顺序是根据它们在 Series 中首次出现的顺序排列的。

2. 筛选属性 `A` 取值为 `v` 的子集 `sub_D`


```
1 sub_D = D[D[A] == v]
```

这行代码筛选出了在数据集 `D` 中，属性 `A` 取值为 `v` 的子集 `sub_D`。这是通过 Pandas 的索引操作实现的，`D[A] == v` 返回一个布尔数组，表示属性 `A` 的取值是否等于 `v`，然后将这个数组传给 `D[]`，就可以得到满足条件的子集。

3. 计算属性 `A` 取值为 `v` 的子集 `sub_D` 的熵

```
1 conditional_entropy += (len(sub_D)/len(D)) * calc_entropy(sub_D)
```

这行代码计算了属性 `A` 取值为 `v` 的子集 `sub_D` 的熵，并将其加到条件熵 `conditional_entropy` 中。具体计算过程如下：

- `(len(sub_D)/len(D))` 计算了子集 `sub_D` 在数据集 `D` 中的占比，表示属性 `A` 取值为 `v` 的概率。
- `calc_entropy(sub_D)` 调用了之前定义的 `calc_entropy` 函数，计算了子集 `sub_D` 的熵。
- `(len(sub_D)/len(D)) * calc_entropy(sub_D)` 计算了子集 `sub_D` 的条件熵，并将其累加到条件熵 `conditional_entropy` 中。

计算连续属性在给定划分点下的信息增益

```
1 # 计算连续属性在给定划分点下的信息增益
2 def calc_information_gain_continuous(D: pd.DataFrame, attribute: str,
   split_point: float) -> float:
3     """
4     计算连续属性在给定划分点下的信息增益
5     :param D: 数据集
6     :param attribute: 属性名称
7     :param split_point: 划分点
8     :return: 信息增益
9     """
10    # 按照划分点分割数据集
11    D1 = D[D[attribute] <= split_point]
12    D2 = D[D[attribute] > split_point]
13
14    # 计算原数据集的熵
15    entropy_before = calc_entropy(D)
16
17    # 计算划分后的加权熵
```

```

18     entropy_after = (len(D1) / len(D)) * calc_entropy(D1) + \
19                     (len(D2) / len(D)) * calc_entropy(D2)
20
21     # 计算信息增益
22     information_gain = entropy_before - entropy_after
23     return information_gain

```

1. 数据集划分：

```

1 D1 = D[D[attribute] <= split_point]
2 D2 = D[D[attribute] > split_point]

```

这两行代码将数据集 `D` 按照给定的划分点 `split_point` 划分成两个子集 `D1` 和 `D2`。`D1` 包含了属性 `attribute` 值小于或等于 `split_point` 的样本，而 `D2` 则包含了属性值大于 `split_point` 的样本。

利用pandas库对DataFrame对象进行条件筛选的语法，具体是基于布尔索引进行的。

1. `D[attribute]`：

这部分代码中 `D[attribute]` 是对DataFrame对象 `D` 中的某一列进行访问。`attribute` 是一个字符串变量，应该是DataFrame中的一个列名。通过 `D[attribute]`，我们获取了DataFrame中名为 `attribute` 的列，这个列是一个Series对象。

2. `D[attribute] <= split_point`：

这部分是一个条件表达式，它将返回一个布尔型Series，其中每个元素都是对应位置上 `D[attribute]` 列的值是否满足小于或等于 `split_point` 的结果。换句话说，它会返回一个与DataFrame的索引相同长度的布尔型Series，对应位置上为True或False。

3. `D[D[attribute] <= split_point]`：

这一部分是基于布尔索引进行的筛选操作。`D[condition]` 中的 `condition` 是一个布尔型Series，它表示了对DataFrame的筛选条件。因此，`D[D[attribute] <= split_point]` 会返回一个新的DataFrame对象，其中包含了满足条件的行。

4. `D1 = D[D[attribute] <= split_point]` 和 `D2 = D[D[attribute] > split_point]`：

这两行代码将根据条件 `D[attribute] <= split_point` 和 `D[attribute] > split_point` 分别筛选出DataFrame `D` 中满足条件的行，分别赋值给变量 `D1` 和 `D2`。这样就完成了根据划分点 `split_point` 对数据集进行划分的操作。

2. 计算原数据集的熵：

```

1 entropy_before = calc_entropy(D)

```

这行代码调用了一个叫做 `calc_entropy` 的函数，用于计算整个数据集 `D` 的熵，并将结果保存在变量 `entropy_before` 中。

5. 计算划分后的加权熵：

```
1 entropy_after = (len(D1) / len(D)) * calc_entropy(D1) + \
2                 (len(D2) / len(D)) * calc_entropy(D2)
```

这行代码计算了划分后的加权熵。首先，它计算了划分后每个子集占总数据集的比例（即权重），然后将每个子集的熵乘以其权重，并将结果相加，得到了划分后的加权熵。

6. 计算信息增益：

```
1 information_gain = entropy_before - entropy_after
```

这行代码计算了信息增益，即原数据集的熵与划分后的加权熵之差。最后，函数返回了信息增益的值。

计算信息增益

```
1 # 计算信息增益，自动处理连续和分类属性
2 def calc_information_gain(D: pd.DataFrame, A: str, split_point: float = None) -
  > float:
3     """
4     计算信息增益，自动处理连续和分类属性。连续属性需要划分点。
5     :param D: 数据集
6     :param A: 属性
7     :param split_point: 连续属性的划分点，对于分类属性，这个参数不使用
8     :return: 属性A的信息增益
9     """
10    # 检查属性A是连续还是分类的
11    if pd.api.types.is_numeric_dtype(D[A]) and split_point is not None:
12        # 对于连续属性，使用划分点来计算信息增益
13        return calc_information_gain_continuous(D, A, split_point)
14    else:
15        # 对于分类属性，使用原始的信息增益计算方法
16        return calc_entropy(D) - calc_conditional_entropy(D, A)
```

1. 判断属性类型：

```
1 if pd.api.types.is_numeric_dtype(D[A]) and split_point is not None:
```

这部分代码检查属性 `A` 是否是连续属性，方法是使用 `pd.api.types.is_numeric_dtype()` 函数。如果 `D[A]` 的数据类型是数值型（即连续属性）且给定了划分点 `split_point`，则说明属性 `A` 是连续属性。

这行代码是一个条件语句，它检查属性 `A` 的数据类型是否是数值型（即连续属性），并且判断是否提供了划分点 `split_point`。让我们逐步解释：

1. `pd.api.types.is_numeric_dtype(D[A])`：

- `pd.api.types` 是 Pandas 库的一个模块，提供了用于数据类型检查的函数。
- `is_numeric_dtype()` 是其中的一个函数，用于检查指定对象的数据类型是否为数值型（即整数或浮点数）。
- 在这里，`D[A]` 表示数据集 `D` 中的列 `A`，这是一个 pandas Series 对象，该语句检查该列的数据类型是否为数值型。

2. `split_point is not None`：

- 这是一个简单的条件判断，检查变量 `split_point` 是否不为 `None`。
- 如果 `split_point` 不是 `None`，说明在调用函数时提供了一个划分点，即表示要对连续属性进行处理。

2. 根据属性类型调用不同的函数计算信息增益：

```
1 return calc_information_gain_continuous(D, A, split_point)
```

如果属性 `A` 是连续属性，则调用 `calc_information_gain_continuous()` 函数来计算信息增益；否则，说明属性 `A` 是分类属性，将调用原始的信息增益计算方法。

3. 返回信息增益值：

```
1 return calc_entropy(D) - calc_conditional_entropy(D, A)
```

如果属性 `A` 是分类属性，则调用 `calc_entropy()` 函数计算数据集 `D` 的熵，然后调用 `calc_conditional_entropy()` 函数计算在给定属性 `A` 的条件下的条件熵，并将二者相减得到信息增益的值。

计算属性A的信息增益率

```

1 # 计算属性A的增益率
2 def calc_gain_ratio(D: pd.DataFrame, A: str) -> float:
3     """
4     计算属性A的增益率
5     :param D: 数据集
6     :param A: 属性
7     :return: 属性A的增益率
8     """
9     information_gain: float = calc_information_gain(D, A)
10    A_values: list = D[A].unique()
11    iv: float = -sum((len(D[D[A] == v]) / len(D)) * math.log2(len(D[D[A] ==
12    v]) / len(D)) for v in A_values)
13    return information_gain / iv if iv != 0 else 0.0

```

1. `A_values = D[A].unique()`：这行代码获取了数据集 `D` 中属性 `A` 的所有取值，并使用 `unique()` 方法去除重复值，得到一个包含所有可能取值的数组 `A_values`。
2. `iv = -sum((len(D[D[A] == v])/len(D)) * math.log2(len(D[D[A] == v])/len(D)) for v in A_values)`：这行代码计算了属性 `A` 的固有值（Intrinsic Value），用于计算增益率。具体计算过程如下：
 - `(len(D[D[A] == v])/len(D))` 计算了属性 `A` 取值为 `v` 的样本在数据集 `D` 中的比例，即该属性值的概率。
 - `math.log2(len(D[D[A] == v])/len(D))` 计算了该属性值的概率的对数。
 - `-sum(...)` 对所有属性值的概率的对数求和，并加上负号，得到固有值 `iv`。
3. `return information_gain / iv if iv != 0 else 0`：最后，函数返回属性 `A` 的增益率。如果固有值 `iv` 不为零，则返回信息增益与固有值的比值；否则返回零，以避免除以零的错误。

选择最优划分属性

```

1 # 选择最优划分属性
2 def choose_best_feature(D: pd.DataFrame) -> str:
3     """
4     选择最优划分属性
5     :param D: 数据集
6     :return: 最优划分属性
7     """
8     features = D.columns[:-1]
9     gain_ratios = {feature: calc_gain_ratio(D, feature) for feature in
10    features}
11    return max(gain_ratios, key=gain_ratios.get)

```

1. `features = D.columns[:-1]`：获取数据集 `D` 中除最后一列外的所有列，即所有属性列。这里使用了 Pandas DataFrame 的 `columns` 属性来获取列标签，并使用切片操作 `[:-1]` 来去除最后一列。
2. `gain_ratios = {feature: calc_gain_ratio(D, feature) for feature in features}`：遍历所有属性，计算每个属性的增益率，并将结果存储在一个字典 `gain_ratios` 中。字典的键是属性名，值是对应属性的增益率。这里使用了字典推导式。
3. `max(gain_ratios, key=gain_ratios.get)`：通过 `max()` 函数找到 `gain_ratios` 中增益率最大的属性。`max()` 函数的 `key` 参数指定了比较的规则，这里使用了 `gain_ratios.get` 函数作为比较的依据，即比较字典的值。最终，该函数返回了具有最大增益率的属性名。

寻找连续属性的最优划分点

```
1 # 寻找连续属性的最优划分点
2 def find_best_split_point_for_continuous_attribute(D: pd.DataFrame, attribute:
   str) -> float:
3     """
4     寻找连续属性的最优划分点
5     :param D: 数据集
6     :param attribute: 属性名称
7     :return: 最优划分点
8     """
9     sorted_values = D[attribute].sort_values().unique()
10    split_points = [(sorted_values[i] + sorted_values[i + 1]) / 2 for i in
   range(len(sorted_values) - 1)]
11    max_gain = -float('inf')
12    best_split = None
13    for split_point in split_points:
14        # 计算每个可能划分点的信息增益
15        gain = calc_information_gain_continuous(D, attribute, split_point)
16        if gain > max_gain:
17            max_gain = gain
18            best_split = split_point
19    return best_split
```

1. 排序并获取唯一值：

```
1 sorted_values = D[attribute].sort_values().unique()
```

这行代码首先通过 `D[attribute]` 获取数据集 `D` 中指定属性 `attribute` 的所有取值，并对这些值进行排序。然后，通过 `unique()` 方法获取排序后的唯一值数组，即属性的所有不同取值。

2. 生成划分点列表：

```
1 split_points = [(sorted_values[i] + sorted_values[i + 1]) / 2 for i in
    range(len(sorted_values) - 1)]
```

这行代码生成了一个划分点的列表 `split_points`。它通过遍历排序后的唯一值数组，并计算相邻两个值的中点来确定划分点。

3. 初始化最大增益和最佳划分点：

```
1 max_gain = -float('inf')
2 best_split = None
```

这两行代码用于初始化最大增益和最佳划分点的变量。`max_gain` 初始化为负无穷，`best_split` 初始化为 `None`。

4. 遍历划分点列表并计算信息增益：

```
1 for split_point in split_points:
2     gain = calc_information_gain_continuous(D, attribute, split_point)
3     if gain > max_gain:
4         max_gain = gain
5         best_split = split_point
```

这部分代码遍历划分点列表 `split_points`，对于每一个划分点，调用 `calc_information_gain_continuous()` 函数计算相应的信息增益，并将结果与当前的最大增益进行比较。如果当前的信息增益大于最大增益，则更新最大增益和最佳划分点。

生成决策树

```
1 # 递归构建决策树（不包含剪枝）
2 def create_decision_tree(D: pd.DataFrame, features: List[str],
    continuous_attributes: List[str]) -> dict:
3     """
4     递归构建决策树，适应连续变量处理
5     :param D: 数据集
6     :param features: 特征列表
7     :param continuous_attributes: 连续属性列表
```

```

8      :return: 决策树
9      """
10     # 基本终止条件
11     class_counts = D.iloc[:, -1].value_counts()
12     if len(class_counts) == 1:
13         return class_counts.index[0]
14     if not features:
15         return class_counts.idxmax()
16
17     # 选择最优属性及其划分点 (如果是连续的)
18     best_gain = -float('inf')
19     best_feature = None
20     split_point = None
21     for feature in features:
22         if feature in continuous_attributes:
23             # 对于连续属性, 找到最优划分点
24             point = find_best_split_point_for_continuous_attribute(D, feature)
25             gain = calc_information_gain(D, feature, point) # 假设这个函数现在可
以处理连续属性
26             if gain > best_gain:
27                 best_gain = gain
28                 best_feature = feature
29                 split_point = point
30         else:
31             # 对于离散属性, 正常处理
32             gain = calc_information_gain(D, feature) # 离散属性的信息增益计算
33             if gain > best_gain:
34                 best_gain = gain
35                 best_feature = feature
36
37     # 根据选择的最优属性分割数据集
38     tree = {best_feature: {}}
39     if best_feature in continuous_attributes:
40         # 处理连续属性的分割
41         left_D = D[D[best_feature] <= split_point]
42         right_D = D[D[best_feature] > split_point]
43         tree[best_feature]['≤' + str(split_point)] =
create_decision_tree(left_D,
44
45
for f in features if f != best_feature],
continuous_attributes)
46         tree[best_feature]['>' + str(split_point)] =
create_decision_tree(right_D,
47
48
for f in features if f != best_feature],

```



```

48     continuous_attributes)
49     else:
50         # 处理离散属性的分割
51         for value in D[best_feature].unique():
52             sub_D = D[D[best_feature] == value]
53             tree[best_feature][value] = create_decision_tree(sub_D, [f for f in
54                 features if f != best_feature],
55                 continuous_attributes)
56     return tree

```

1. 基本终止条件:

```

1 class_counts = D.iloc[:, -1].value_counts()
2 if len(class_counts) == 1:
3     return class_counts.index[0]
4 if not features:
5     return class_counts.idxmax()

```

这部分代码是函数的基本终止条件。如果数据集中的样本都属于同一个类别，或者特征列表为空，则直接返回该类别；否则，继续选择最优属性进行划分。

`class_counts.index[0]` 和 `class_counts.idxmax()` 都是返回索引

2. 选择最优属性:

```

1 best_gain = -float('inf')
2 best_feature = None
3 split_point = None
4 for feature in features:
5     if feature in continuous_attributes:
6         # 对于连续属性，找到最优划分点
7         point = find_best_split_point_for_continuous_attribute(D, feature)
8         gain = calc_information_gain(D, feature, point)
9         if gain > best_gain:
10             best_gain = gain
11             best_feature = feature
12             split_point = point
13     else:
14         # 对于离散属性，正常处理
15         gain = calc_information_gain(D, feature) # 离散属性的信息增益计算
16         if gain > best_gain:

```

```

17         best_gain = gain
18         best_feature = feature

```

这部分代码是选择最优属性的过程。对于每个特征，如果它是连续属性，则调用

`find_best_split_point_for_continuous_attribute()` 函数找到最优划分点，并计算信息增益；如果是离散属性，则直接计算信息增益。最终选择信息增益最大的属性作为最优属性，并记录最优划分点（如果是连续属性）。

`-float('inf')` 表示负无穷大。在 Python 中，`float('inf')` 表示正无穷大，而 `-float('inf')` 则表示负无穷大。

5. 根据最优属性分割数据集：

```

1  tree = {best_feature: {}}
2  if best_feature in continuous_attributes:
3      # 处理连续属性的分割
4      left_D = D[D[best_feature] <= split_point]
5      right_D = D[D[best_feature] > split_point]
6      tree[best_feature]['≤' + str(split_point)] = create_decision_tree(left_D,
7                                                                           [f for f
8      in features if f != best_feature],
9      continuous_attributes)
10     tree[best_feature]['>' + str(split_point)] = create_decision_tree(right_D,
11                                                                           [f for f
12     in features if f != best_feature],
13     continuous_attributes)
14 else:
15     # 处理离散属性的分割
16     for value in D[best_feature].unique():
17         sub_D = D[D[best_feature] == value]
18         tree[best_feature][value] = create_decision_tree(sub_D, [f for f in
19     features if f != best_feature],
20     continuous_attributes)

```

这部分代码根据最优属性将数据集进行分割。对于连续属性，根据最优划分点将数据集分为左右两个子集，并分别递归调用 `create_decision_tree()` 函数；对于离散属性，根据属性值进行分割，并为每个属性值递归调用 `create_decision_tree()` 函数。

```

1  tree[best_feature]['≤' + str(split_point)] = create_decision_tree(left_D,
2                                                                           [f for f
3  in features if f != best_feature],

```

这句代码涉及到字典的嵌套、列表推导式和函数调用：

- `tree[best_feature]` : 这部分访问了名为 `tree` 的字典中键为 `best_feature` 的值。由于 `tree` 是一个字典，因此 `tree[best_feature]` 得到的是一个子字典或者空字典（如果之前未定义过）。
- `['<' + str(split_point)]` : 在上述子字典中，将键设置为 `'<' + str(split_point)`。这表明此键用于表示连续属性的子分支，其键的形式是 `"<"` 加上划分点的字符串形式，例如 `"<0.5"`。
- `create_decision_tree(left_D, [f for f in features if f != best_feature], continuous_attributes)` : 这是一个函数调用，调用了名为 `create_decision_tree` 的函数。它的参数是 `left_D`、一个特征列表和一个连续属性列表。特征列表通过列表推导式 `[f for f in features if f != best_feature]` 构建，该列表表示去除了当前最佳特征 `best_feature` 的特征列表。

`[f for f in features if f != best_feature]` 是一个列表推导式，用于从列表 `features` 中过滤出不等于 `best_feature` 的元素，并构建一个新的列表。让我们逐步解释：

- `for f in features` : 这部分定义了一个循环，遍历列表 `features` 中的每个元素，并将其赋值给变量 `f`。
- `if f != best_feature` : 这是一个条件语句，它筛选出 `features` 中不等于 `best_feature` 的元素。
- `[f for f in features if f != best_feature]` : 综合起来，这个列表推导式遍历 `features` 中的每个元素，当元素不等于 `best_feature` 时，将其加入新列表中。

该列表推导式返回一个新的列表，其中包含了 `features` 中除去 `best_feature` 外的所有元素。

后剪枝

```

1 # 后剪枝函数
2 def post_pruning(tree: Dict[str, Any], D: pd.DataFrame, features: List[str]) -
  > Dict[str, Any]:
3     """
4     对决策树进行后剪枝
5     :param tree: 当前决策树
6     :param D: 数据集
7     :param features: 特征集
8     :return: 剪枝后的决策树
9     """

```

```

10     # 检查树是否是叶节点
11     if not isinstance(tree, dict):
12         return tree
13
14     # 遍历树中的每个节点
15     for feature, branches in tree.items():
16         for value, subtree in branches.items():
17             # 递归剪枝子树
18             subtree = post_pruning(subtree, D[D[feature] == value], [f for f in
features if f != feature])
19             tree[feature][value] = subtree
20
21     # 尝试剪枝当前节点
22     if all(not isinstance(subtree, dict) for subtree in
tree[feature].values()):
23         # 计算剪枝前后的准确性
24         accuracy_before_pruning = calc_accuracy(tree, D)
25         # 将当前节点替换为最常见的类
26         most_common_class = D.iloc[:, -1].mode()[0]
27         pruned_tree = most_common_class
28         accuracy_after_pruning = calc_accuracy(pruned_tree, D)
29         # 如果剪枝后准确性不降低, 则进行剪枝
30         if accuracy_after_pruning >= accuracy_before_pruning:
31             return pruned_tree
32
33     return tree

```

后剪枝是一种决策树优化技术, 旨在通过删除树的某些分支来减少过拟合, 同时尽可能保持或提高模型的准确性。

1. 检查是否是叶节点:

- `if not isinstance(tree, dict): return tree`

这一行检查当前节点 (`tree`) 是否是一个叶节点。在决策树中, 叶节点通常不再是一个字典类型, 而是直接存储类别标签。如果当前节点是叶节点, 就直接返回这个节点, 因为叶节点无需剪枝。

2. 遍历树中的每个节点:

- 循环遍历当前决策树中的每个节点。由于决策树是以嵌套字典的形式存储的, 每个非叶节点的键 (`feature`) 对应一个特征, 值 (`branches`) 是另一个字典, 表示基于该特征的分支。
- 对于每个分支的子树 (`subtree`), 递归地调用 `post_pruning` 函数进行剪枝。

3. 尝试剪枝当前节点:

- `if all(not isinstance(subtree, dict) for subtree in tree[feature].values()):`

这行代码检查当前节点的所有子节点是否都是叶节点。如果是，考虑将当前节点转换为叶节点，即进行剪枝。

4. 计算剪枝前后的准确性:

- 使用一个名为 `calc_accuracy` 的函数来计算剪枝前后的准确率。首先计算当前树的准确率，然后假设将当前节点替换为最常见类别，计算这种情况下的准确率。

5. 执行剪枝:

- 如果将当前节点替换为数据集中最常见的类别（即叶节点）不会降低决策树的准确性，则进行剪枝，将当前节点替换为该类别。

6. 返回剪枝后的决策树:

- 如果对任何节点都没有执行剪枝操作，或者剪枝后的树已经是最优的，则返回当前的决策树。

`post_pruning` 函数通过递归地检查决策树的每个节点，评估是否通过将非叶节点替换为叶节点（最常见的类别）可以维持或提高模型的准确性。如果可以，就执行剪枝操作。

对单个实例进行预测

```
1 # 对单个实例进行预测
2 def predict(tree: Dict[str, Any], instance: Dict[str, Any]) -> Any:
3     """
4     对单个实例进行预测
5     :param tree: 决策树
6     :param instance: 单个数据实例
7     :return: 预测结果
8     """
9     if not isinstance(tree, dict):
10         return tree
11     root = next(iter(tree))
12     subtree = tree[root]
13     value = instance[root]
14     if value in subtree:
15         return predict(subtree[value], instance)
16     else:
17         return None
```

1. 检查 `tree` 是否是字典类型

```
1 if not isinstance(tree, dict):
```

这行代码检查 `tree` 是否是字典类型。如果不是字典类型，则直接返回 `tree`，即返回决策树的叶节点（预测结果）。

2. 获取决策树的根节点

```
1 root = next(iter(tree))
```

这行代码获取了决策树的根节点。`iter(tree)` 创建了一个迭代器对象，`next()` 函数返回迭代器的下一个元素，即字典中的第一个键，也就是根节点。

3. 获取了根节点为键的子树

```
1 subtree = tree[root]
```

这行代码获取了以根节点为键的子树，即根节点的所有子节点及其对应的子树。

4. 获取数据实例中根节点对应的特征值

```
1 value = instance[root]
```

5. 返回结果

```
1     if value in subtree:
2         return predict(subtree[value], instance)
3     else:
4         return None
```

检查数据实例的特征值是否在子树中存在。

如果特征值在子树中存在，递归调用 `predict` 函数，并传入对应的子树和数据实例。这样可以继续在子树上进行预测。递归过程将一直持续，直到到达叶节点为止，然后返回叶节点的值作为最终预测结果。

如果特征值不在子树中存在，表示决策树无法根据当前特征对数据实例进行分类，因此返回 `None`。

载入数据

```
1 # 载入数据
2 data = {
```

```

3     '色泽': ['青绿', '乌黑', '乌黑', '青绿', '浅白', '青绿', '乌黑', '乌黑', '青
    绿', '浅白', '浅白', '青绿', '浅白', '乌黑', '浅白', '青绿'],
4     '根蒂': ['蜷缩', '蜷缩', '蜷缩', '蜷缩', '蜷缩', '稍蜷', '稍蜷', '稍蜷', '硬
    挺', '硬挺', '蜷缩', '稍蜷', '稍蜷', '稍蜷', '蜷缩', '蜷缩'],
5     '敲声': ['浊响', '沉闷', '浊响', '沉闷', '浊响', '浊响', '浊响', '浊响', '清
    脆', '清脆', '浊响', '浊响', '沉闷', '浊响', '浊响', '沉闷'],
6     '纹理': ['清晰', '清晰', '清晰', '清晰', '清晰', '清晰', '清晰', '稍糊', '清晰', '清
    晰', '模糊', '模糊', '稍糊', '稍糊', '清晰', '模糊', '稍糊'],
7     '脐部': ['凹陷', '凹陷', '凹陷', '凹陷', '凹陷', '稍凹', '稍凹', '稍凹', '平
    坦', '平坦', '平坦', '凹陷', '凹陷', '稍凹', '平坦', '稍凹'],
8     '触感': ['硬滑', '硬滑', '硬滑', '硬滑', '硬滑', '硬滑', '软粘', '软粘', '硬滑', '软
    粘', '硬滑', '软粘', '硬滑', '硬滑', '软粘', '硬滑', '硬滑'],
9     '含糖率': [0.460, 0.376, 0.264, 0.318, 0.215, 0.237, 0.149, 0.211, 0.267,
    0.057, 0.099, 0.161, 0.198, 0.370, 0.042, 0.103],
10    '好瓜': ['是', '是', '是', '是', '是', '是', '是', '是', '否', '否', '否',
    '否', '否', '否', '否', '否', '否']
11 }
12
13 df = pd.DataFrame(data)

```

1. `data = {...}`: 这里定义了一个 Python 字典 `data`，其中包含了一些数据。这些数据看起来像是描述了一些水果（可能是西瓜）的特征和标签信息，如颜色、根蒂、敲声等等，以及是否为好瓜。
2. `df = pd.DataFrame(data)`: 这行代码使用 Pandas 中的 `DataFrame` 函数将字典 `data` 转换为了一个数据帧（DataFrame），并将其赋值给变量 `df`。

- 在 Pandas 中创建 DataFrame 的主要函数是 `pd.DataFrame()`。这个函数接受多种不同的输入形式，包括字典、列表、数组、Series 对象等，然后将其转换为 DataFrame 数据结构。

下面是对 `pd.DataFrame()` 函数的一些重要参数和使用方式的介绍：

- 数据输入:** `data` 参数是必需的，它可以是各种数据类型，包括字典、列表、数组、Series 对象等。这些数据将会被转换成 DataFrame 的形式。
 - 字典：字典的键将成为 DataFrame 的列标签，字典的值可以是列表、数组或其他序列类型。
 - 列表：列表中的元素可以是字典、列表、数组等，它们将会被转换成 DataFrame 的行。
 - 数组：数组可以是二维的，每一行对应 DataFrame 中的一行数据。
 - Series 对象：单个 Series 对象也可以作为数据输入，它将会成为 DataFrame 的一列。
- 索引和列标签:** `index` 和 `columns` 参数分别用于指定 DataFrame 的行索引和列标签。如果不指定，行索引和列标签将会自动生成。
- 数据类型:** `dtype` 参数用于指定 DataFrame 中每一列的数据类型。如果不指定，Pandas 将会自动推断数据类型。

- d. **复制数据**: `copy` 参数用于控制是否复制输入数据。默认情况下, Pandas 会尝试在可能的情况下复制输入数据, 以避免数据共享问题。
- e. **其他参数**: 除了上述参数之外, `pd.DataFrame()` 函数还接受一些其他参数, 如 `index_col`、`dtype`、`copy` 等, 用于控制 DataFrame 的创建和转换过程。
- 数据帧是 Pandas 中用于处理二维数据的主要数据结构, 类似于电子表格或数据库表格。这个数据帧将包含了字典 `data` 中的数据, 并且每一列的名称将会根据字典中键的名称自动命名。

剩余操作

```
1 # 创建数据集
2 df = pd.DataFrame(data)
3
4 # 特征列表
5 features = list(df.columns[:-1])
6
7 # 连续属性列表
8 continuous_attributes = ['含糖率']
9
10 # 创建决策树
11 decision_tree = create_decision_tree(df, features, continuous_attributes)
12
13 # 输出剪枝前决策树
14 graph_decision_tree = plot_decision_tree(decision_tree)
15
16 # 对决策树进行后剪枝
17 pruned_tree = post_pruning(decision_tree, df, features)
18
19 # 输出剪枝后决策树
20 graph_pruned_tree = plot_decision_tree(pruned_tree)
21
22 # 保存并显示图像
23 graph_pruned_tree.render(filename='pruned_tree', directory='.', view=True) # 将
    文件保存在当前工作目录
24 graph_decision_tree.render(filename='decision_tree', directory='.', view=True)
    # 将文件保存在当前工作目录
```

```
1 graph_pruned_tree.render(filename='pruned_tree', directory='.', view=True)
```

这句代码调用了方法 `render()`, 该方法通常用于绘制图形或图像。

- `graph_pruned_tree` 是一个 `Digraph` 对象，它是由 Graphviz 库提供的一个对象，用于创建图形可视化。
- `render()` 是 `Digraph` 对象的一个方法，用于将图形渲染为一个图像文件，并可选地将其显示或保存到文件中。
- 在这个语法中，`filename='pruned_tree'` 指定了要保存的文件名为 `'pruned_tree'`，这意味着生成的图像文件将被命名为 `pruned_tree.xxx`，其中 `.xxx` 取决于指定的格式（在这里默认为 PNG 格式）。
- `directory='.'` 指定了保存文件的目录。在这里，`.` 表示当前工作目录，即脚本所在的文件夹。
- `view=True` 指定了在渲染后是否立即查看图像。如果设置为 `True`，则会在渲染后自动打开默认的图像查看器来显示生成的图像；如果设置为 `False`，则不会自动显示图像。