

人工智能第五次实验报告

决策树模型

课程：人工智能原理	年级专业：19级软件工程
姓名：郑有为	学号：19335286

目录

目录

一、问题背景

1.1 监督学习简介

1.2 决策树简介

二、程序说明

2.1 数据载入

2.2 功能函数

2.3 决策树模型

2.4 决策树可视化

三、程序测试

3.1 数据集说明

3.2 决策树生成和测试

3.3 学习曲线评估算法精度

四、实验总结

附录 - 程序代码

一、问题背景

1.1 监督学习简介

机器学习的形式包括无监督学习，强化学习，监督学习和半监督学习；学习任务有分类、聚类和回归等。

监督学习通过观察“输入—输出”对，学习从输入到输出的映射函数。分类监督学习的训练集为标记数据，每一条数据有对应的“标签”，根据标签可以将数据集分为若干个类别。分类监督学习经训练集生成一个学习模型，可以用来预测一条新数据的标签。

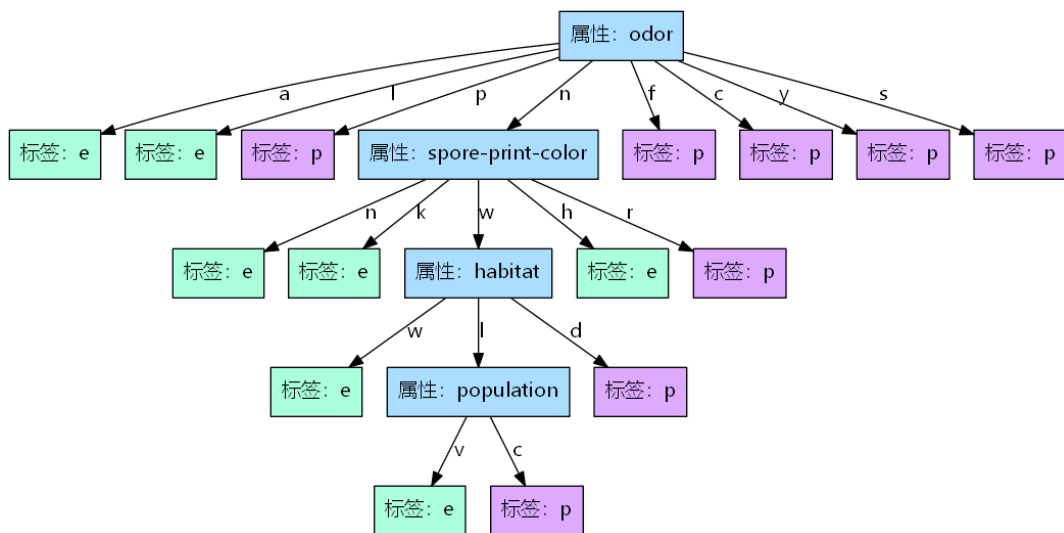
常见的监督学习模型有决策树、KNN算法、朴素贝叶斯和随机森林等。

1.2 决策树简介

决策树归纳是一类简单的机器学习形式，它表示为一个函数，以属性值向量作为输入，返回一个决策。

• 决策树的组成

- 决策树由内节点上的属性值测试、分支上的属性值和叶子节点上的输出值组成。



• 决策树的学习算法

```

1  函数: dt(dataset, par_dataset, attrs)
2  输入: 子数据集 dataset, 父数据集 par_dataset, 属性列表 attrs
3  输出: 决策树
4
5  if dataset 为空
6      return 父数据集 par_dataset 概率最大的标签
7  else if attrs 为空
8      return 子数据集 dataset 概率最大的标签
9  else if 子数据集 dataset 的标签全部相同
10     return 子数据集一致的标签
11 else
12     best_attr <- attrs 在 dataset 中信息收益最高的属性
13     tree <- 值为 best_attr 的节点
14     for best_attr 的所有可选取值 best_attr_value
15         sub_dataset <- dataset中所有属性best_attr的值为best_attr_value的数据
           构成的子数据集
16         sub_tree <- dt(sub_dataset, dataset, attrs - best_attr)
17         添加一条从 tree 指向 sub_tree 的边, 边的值为 best_attr_value
18     return tree

```

• 属性的优先程度判准 - 信息收益

使用信息熵来量化当前状态信息的不确定程度, 对我们来说, 信息熵越小, 信息量就越大, 也就越重要。熵的定义如下:

$$H(V) = - \sum_k P(v_k) \log_2 P(v_k)$$

设有 d 个不同值的属性 A 将训练集划分为 d 个子集 E_k , 使用属性 A 对数据集进行划分所得到的信息收益为:

$$G(A) = H(E) - \sum_{k=1}^d \frac{|E_k|}{|E|} H(E_k)$$

这里的信息收益也就是使用属性 A 划分带来熵的期望减少值。

• 决策树算法的改进

- **增加信息增益的最小阈值:** 如果通过属性 A 划分带来的信息增益小于阈值 ϵ , 则不再继续划分, 而是生成叶子节点, 节点的标签为当前数据集概率最大的标签, 这种改进可以减小决策树的规模, 避免过度拟合, 缩短决策时间。

- χ^2 剪枝：早期终止可能会组织我们发现更好的属性组合，在生成决策树后进行剪枝，即自底向上移除决策树中划分效果不好的节点，替换为叶子节点。 χ^2 剪枝通过量化数据的偏离程度（带 $v - 1$ 个自由度的卡方分布， v 为样例总数），来判断是否剪枝。剪枝也能减小决策树的规模，避免过度拟合，缩短决策时间。

- (类ID3) 决策树的适用范围

上述给出的决策树基于信息熵的度量，根据属性值来划分节点，当属性值不是离散的选项，如布尔值，而是浮点类型时，需要对数值进行划分，否则会生成非常多的分支节点，不能达到很好的效果。故上述算法要求属性值离散或经过预先分组，构建出的决策树也称为离散决策树。

- 常见的决策树算法

常见的决策树算法除了有基于信息收益的ID3算法，还有基于信息增益比的C4.5算法和基于基尼指数的CART算法。

C4.5算法与ID3算法基本一致，只是将ID3算法的信息收益换成了信息增益比：

$$GR(D, A) = \frac{G(A)}{H(D)}$$

其中 A 为属性， D 为训练集。

CART算法可以处理非离散的属性值，每次划分将训练集根据某个特征 A 的一个可能值 a_m 分割为两个部分 ($a < a_m$ 和 $a \geq a_m$)，在所有可能的特征 A 以及他们所有可能的切分点 a_m 中，选择基尼指数最小的特征及其对应可能的切分点作为最有特征与最优切分点，以此来构建出一棵二叉决策树。

- 决策树的优缺点：

- 优点：分类快，具有可读性，易于理解。
- 缺点：可能会出现过度拟合的问题，并非适用于所有分类场景。

二、程序说明

2.1 数据载入

数据集被统一格式化为一个类，它包含四个属性：

- `feature_names`：属性名列表
- `target_names`：标签(分类)名
- `data`：属性数据矩阵, 每行是一个数据, 每个数据是每个属性的对应值的列表
- `target`：目标分类值列表

类 `load_car` 和 `load_mushroom` 分别从文件中导入汽车数据集和蘑菇数据集，类 `new_dataset` 用于快速生成一个满足格式的数据集，使用方法如下：

2.2 功能函数

- `get_h(target)` 计算并返回信息熵
 - 参数：`target` 给定目标分类值列表
- `get_subset(dataset, feature_name, feature_value)` 筛选并返回子数据集
 - 参数：`dataset` 待分析数据集；`feature_name` 属性名字符串；`feature_value` 属性值
 - 选择条件：原数据集中的属性 `feature_name` 值是否等于 `feature_value`
 - 备注：选择后会从数据子集中删去 `feature_name` 属性对应的一列
- `best_split(dataset)` 寻找并返回信息收益最大的属性划分，返回最佳属性
 - 参数：`dataset` 待分析数据集
- `vote_most(dataset)` 返回数据集中一个数据概率上最可能的标签

- 参数: `dataset` 待分析数据集
- `accuracy_rate(predict_result, target_result)` 返回测试的正确率
 - 参数: `predict_result`: 预测标签列表, `target_result`: 实际标签列表

2.3 决策树模型

- 决策树节点** `dt_node` 一个节点有以下属性:

属性名	注释
<code>id</code>	为节点赋予一个全局ID, 目的是方便画图
<code>feature_name</code>	非叶子节点的属性名, 若节点为叶子结点, 则为 <code>None</code>
<code>target_value</code>	叶子节点的标签, 若节点为非叶子节点, 则为 <code>None</code>
<code>vote_most</code>	非叶子节点最可能的标签, 若节点为叶子结点, 则为 <code>None</code>
<code>parent</code>	父亲节点
<code>child</code>	儿子节点, 以当前节点的属性对应的属性值作为键值的字典

初始化: `dt_node(self, content, is_leaf=False, parent=None)`

参数名	注释
<code>content</code>	节点的内容, 对于叶子节点为标签值, 对于非叶子节点为属性名
<code>is_leaf</code>	节点是否为叶子节点, 默认为 <code>False</code>
<code>parent</code>	指定节点的父节点, 默认为 <code>None</code>

- 决策树模型** `dt_tree`

属性名	注释
<code>tree</code>	决策树的根节点
<code>map_str</code>	pydotplus 格式的代码串, 用于做图
<code>color_dir</code>	叶子节点可选颜色, 以标签值为键值的字典

方法名	注释
<code>fit(self, train_set)</code>	训练: 根据提供的训练集生成决策树
<code>predict(self, test_set)</code>	预测: 对测试集的数据进行预测, 返回预测结果数组
<code>show_tree(self, path="demo.png")</code>	可视化: 基于绘图工具 pydotplus 生成图片, 保存在 <code>path</code> 中

2.4 决策树可视化

基于绘图工具 pydotplus，生成绘图代码

- 初始化代码串：

```
1 self.map_str = ""
2 digraph demo{
3 node [shape=box, style="rounded", color="black", fontname="Microsoft
  YaHei"];
4 edge [fontname="Microsoft YaHei"];
5 ""
```

- 创建叶子节点：

```
1 node_content = "标签: " + str(node.target_value)
2 self.map_str += "id" + str(node.id) + "[label=\"" + node_content + "\",
  fillColor=\"" + self.color_dir[node.target_value] + "\", style=filled]\n"
```

- 创建非叶子节点：

```
1 node_content = "属性: " + node.feature_name
2 self.map_str += "id" + str(node.id) + "[label=\"" + node_content + "\",
  fillColor=\"#AADDFF\", style=filled]\n"
```

- 创建边：从节点 `node` 指向一个子节点 `node.child[feature_value]`

```
1 self.map_str += "id" + str(node.id) + " -> " + "id" +
  str(node.child[feature_value].id) + "[label=\"" + str(feature_value) +
  "\"]\n"
```

- 生成的代码举例（部分代码）：

```
1 digraph demo{
2 node [shape=box, style="rounded", color="black", fontname="Microsoft
  YaHei"];
3 edge [fontname="Microsoft YaHei"];
4 id0[label="属性: safety", fillColor="#AADDFF", style=filled]
5 id1[label="标签: 0", fillColor="#AAFFDD", style=filled]
6 id0 -> id1[label="low"]
7 id2[label="属性: persons", fillColor="#AADDFF", style=filled]
8 id3[label="标签: 0", fillColor="#AAFFDD", style=filled]
9 id2 -> id3[label="2"]}
```

三、程序测试

3.1 数据集说明

- 汽车数据 `load_car()`：汽车评价数据库是由一个简单的层次决策模型派生而来，该模型最初是为DEX（M. Bohanec, V. Rajkovic: Expert system for decision making. Sistemica 1(1), pp. 145-157, 1990.）模型根据以下概念结构评估汽车的可接受性。

属性：如下表所示

属性	属性值	注释
buying	v-high, high, med, low	购买价格
maint	v-high, high, med, low	维护花销
doors	2, 3, 4, 5-more	门的数量
persons	2, 4, more	承载人数
lug_boot	small, med, big	后备箱的大小
safety	low, med, high	安全性

标签：0 - 不可接受， 1 - 可接受

数据总数：1728条

- **蘑菇数据** load_mushroom()：这个数据集包括了蘑菇样本的描述，每一种都被确定为绝对可食用；绝对有毒；或未知可食用，不推荐。后一类是与有毒的一类结合起来的。

属性：共有22项，依次是 "cap-shape", "cap-surface", "cap-color", "bruises", "odor", "gill-attachment", "gill-spacing", "gill-size", "gill-color", "stalk-shape", "stalk-root", "stalk-surface-above-ring", "stalk-surface-below-ring", "stalk-color-above-ring", "stalk-color-below-ring", "veil-type", "veil-color", "ring-number", "ring-type", "spore-print-color", "population", "habitat"。

标签：e - 可食用（4208项）， p - 有毒（3916项）

数据总数：8124条

数据来源：<https://archive.ics.uci.edu/ml/datasets/Mushroom>

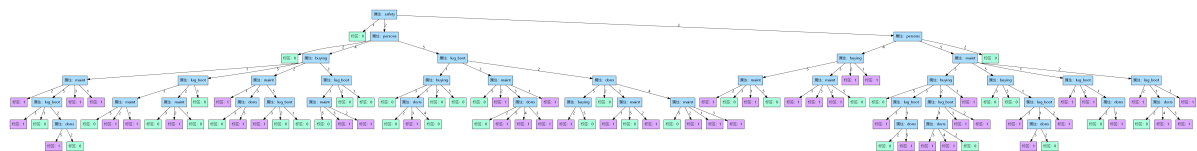
3.2 决策树生成和测试

- **汽车模型**

当训练规模为500时所生成的决策树如下图所示。由于每一个属性都有3-5个可选值，故决策树看起来比较扁。

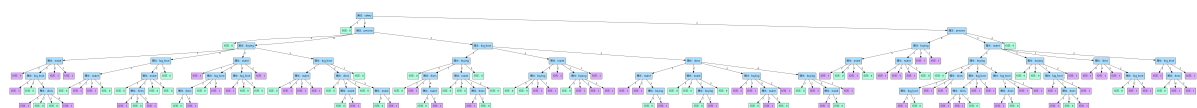
在图中，蓝色的节点为属性节点，绿色和紫色的节点为叶子节点，绿色的节点标签值为0，也就是不购买该汽车，蓝色节点标签值为1，表示购买该汽车，每条从一个节点指向其子节点的边上都有一个值，也就是该节点对应的属性的属性值。

对于该决策树，用剩下的数据集作为测试集进行测试，**准确度为：0.9307253463732681。**



当训练规模为1000时所生成的决策树如下图所示，比较上下两图可以看到随着测试规模的增大，决策树的规模也变得更大，但决策树的高度没有变化，只是叶子结点的数目有所增加。

对于该决策树，用剩下的数据集作为测试集进行测试，**准确度为：0.9559834938101788。**

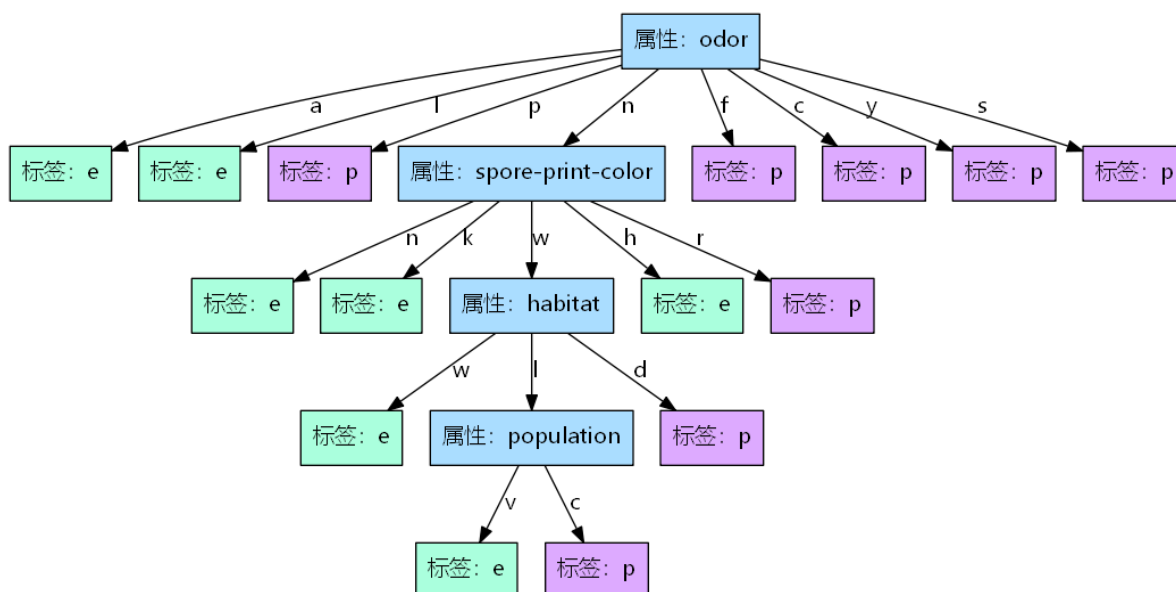


- **蘑菇模型**

虽然蘑菇模型的属性非常多，数据规模也远大于汽车，但其训练得到的决策树模型很简单。

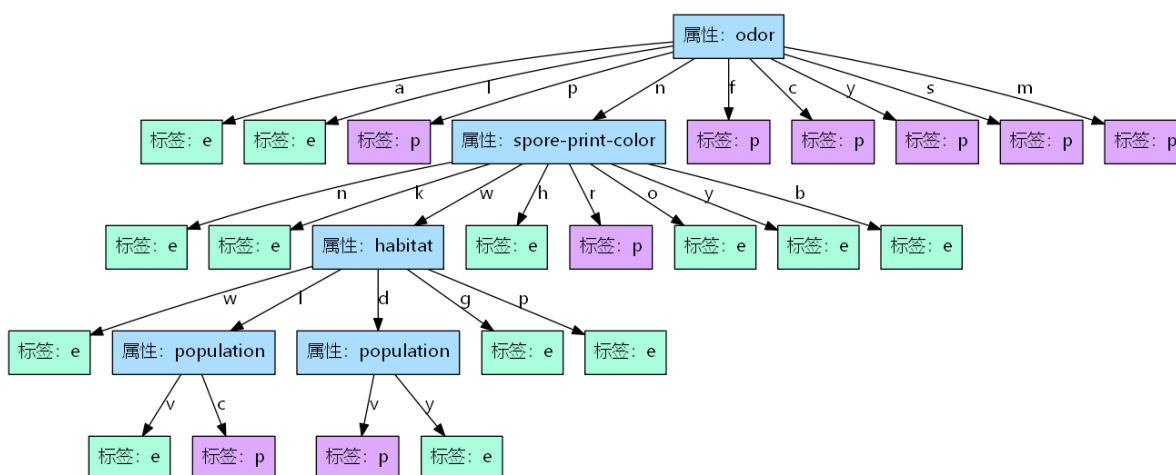
下图是**训练数据规模为5000时的决策树**，该树的高度只有5，也就意味着最多考虑4个属性（odor, spore-print-color, habitat, population）就可以很好的推断出一中蘑菇是否有毒。绿色和紫色的节点表示叶子节点，标签为 e 代表此种蘑菇可食用，标签为 p 代表此种蘑菇有毒。

对于该决策树，用剩下的数据集作为测试集进行测试，**准确度为：0.9859109830291386**



下图是当**训练规模增加到8000时的决策树**，只比上图稍微复杂了一些，但高度和考虑的属性都没有改变。

对于该决策树，用剩下的数据集作为测试集进行测试，**准确度为：1.0**



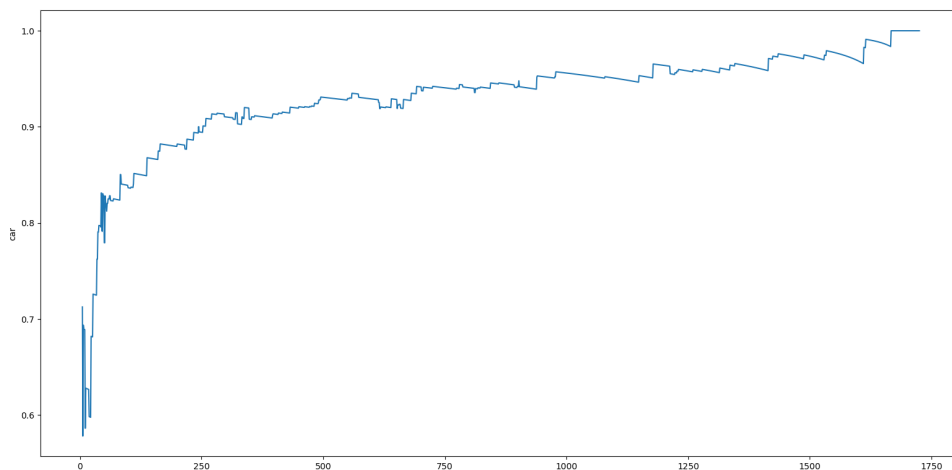
3.3 学习曲线评估算法精度

将所有数据划分为训练集和测试集，用训练集学习生成决策树，再对测试集进行预测，分析预测的准确率。

调用方法 `incremental_train_scale_test(dataset, label, interval=1)` 来生成学习曲线，其中：dataset 为数据训练集, label 为数据名称, interval 为测试规模递增的间隔。

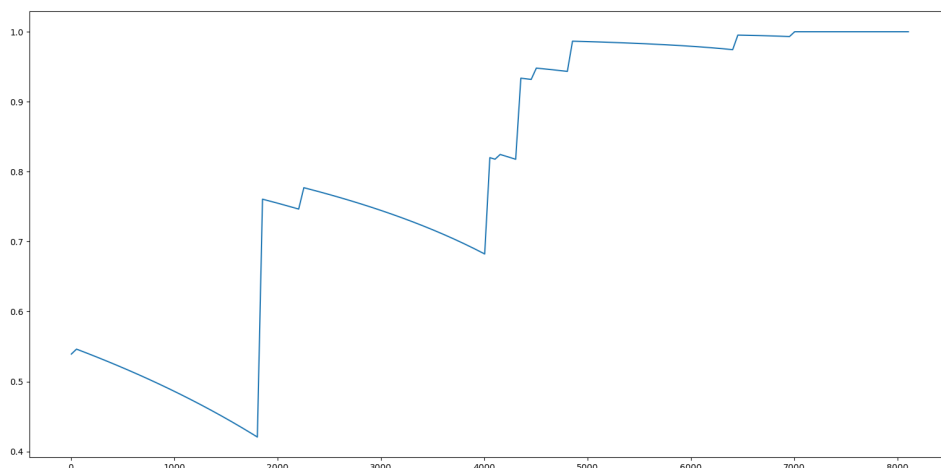
• 汽车模型

首先从规模为5的训练集开始学习，按照1的频率逐步提高训练规模，直到1729。如下图，可以看到随着训练集规模的逐步递增，精度也在不断升高。



- 蘑菇模型

首先从规模为5的训练集开始学习，按照25的频率逐步提高训练规模，直到8124。如下图，可以看到随着训练集规模的逐步递增，精度也在不断升高，在训练集规模大于7000时，测试的准确率已经达到了100%。



四、实验总结

通过本次实验，我们进一步监督学习的基本知识，重点理解决策树的常见算法和改进策略，掌握决策树的基本实现方法，考虑决策树的实现细节，实现了基本的决策树模型并使用汽车模型和蘑菇模型对模型进行测试和可视化，测试效果较好。

附录 - 程序代码

电子版报告和代码地址：<https://gitee.com/WondrousWisdomcard/ai-homework>

```
1 import numpy as np
2 from matplotlib import pyplot as plt
3 from math import log
4 import pandas as pd
5 import pydotplus as pdp
6
7 """
8 19335286 郑有为
```



```

9  人工智能作业 - 实现ID3决策树
10  """
11
12  nonce = 0 # 用来给节点一个全局ID
13  color_i = 0
14  # 绘图时节点可选的颜色, 非叶子节点是蓝色的, 叶子节点根据分类被赋予不同的颜色
15  color_set = ["#AAFFDD", "#DDAAFF", "#DDFFAA", "#FFAADD", "#FFDDAA"]
16
17  # 载入汽车数据, 判断顾客要不要买
18  class load_car:
19      # 在表格中, 最后一列是分类结果
20      # feature_names: 属性名列表
21      # target_names: 标签(分类)名
22      # data: 属性数据矩阵, 每行是一个数据, 每个数据是每个属性的对应值的列表
23      # target: 目标分类值列表
24      def __init__(self):
25          df = pd.read_csv('../dataset/car/car_train.csv')
26          labels = df.columns.values
27          data_array = np.array(df[1:])
28          self.feature_names = labels[0:-1]
29          self.target_names = labels[-1]
30          self.data = data_array[0:,0:-1]
31          self.target = data_array[0:,-1]
32
33  # 载入蘑菇数据, 鉴别蘑菇是否有毒
34  class load_mushroom:
35      # 在表格中, 第一列是分类结果: e 可食用; p 有毒.
36      # feature_names: 属性名列表
37      # target_names: 标签(分类)名
38      # data: 属性数据矩阵, 每行是一个数据, 每个数据是每个属性的对应值的列表
39      # target: 目标分类值列表
40      def __init__(self):
41          df = pd.read_csv('../dataset/mushroom/agaricus-lepiota.data')
42          data_array = np.array(df)
43          labels = ["edible/poisonous", "cap-shape", "cap-surface", "cap-
44                  color", "bruises", "odor", "gill-attachment",
45                  "gill-spacing", "gill-size", "gill-color", "stalk-shape",
46                  "stalk-root", "stalk-surface-above-ring",
47                  "stalk-surface-below-ring", "stalk-color-above-ring",
48                  "stalk-color-below-ring",
49                  "veil-type", "veil-color", "ring-number", "ring-type",
50                  "spore-print-color", "population", "habitat"]
51          self.feature_names = labels[1:]
52          self.target_names = labels[0]
53          self.data = data_array[0:,1:]
54          self.target = data_array[0:,0]
55
56  # 创建一个临时的子数据集, 在划分测试集和训练集时使用
57  class new_dataset:
58      # feature_names: 属性名列表
59      # target_names: 标签(分类)名
60      # data: 属性数据矩阵, 每行是一个数据, 每个数据是每个属性的对应值的列表
61      # target: 目标分类值列表
62      def __init__(self, f_n, t_n, d, t):
63          self.feature_names = f_n
64          self.target_names = t_n
65          self.data = d
66          self.target = t

```

```

63
64 # 计算熵，熵的数学公式为： $H(V) = - \sum_{k} P(v_k) \log_2 P(v_k)$ 
65 # 其中  $P(v_k)$  是随机变量  $V$  具有值  $v_k$  的概率
66 # target: 分类结果的列表，return: 信息熵
67 def get_h(target):
68     target_count = {}
69     for i in range(len(target)):
70         label = target[i]
71         if label not in target_count.keys():
72             target_count[label] = 1.0
73         else:
74             target_count[label] += 1.0
75     h = 0.0
76     for k in target_count:
77         p = target_count[k] / len(target)
78         h -= p * log(p, 2)
79     return h
80
81 # 取数据子集，选择条件是原数据集中的属性 feature_name 值是否等于 feature_value
82 # 注：选择后会从数据子集中删去 feature_name 属性对应的一列
83 def get_subset(dataset, feature_name, feature_value):
84     sub_data = []
85     sub_target = []
86     f_index = -1
87     for i in range(len(dataset.feature_names)):
88         if dataset.feature_names[i] == feature_name:
89             f_index = i
90             break
91
92     for i in range(len(dataset.data)):
93         if dataset.data[i][f_index] == feature_value:
94             l = list(dataset.data[i][:f_index])
95             l.extend(dataset.data[i][f_index+1:])
96             sub_data.append(l)
97             sub_target.append(dataset.target[i])
98
99     sub_feature_names = list(dataset.feature_names[:f_index])
100     sub_feature_names.extend(dataset.feature_names[f_index+1:])
101     return new_dataset(sub_feature_names, dataset.target_names, sub_data,
102                        sub_target)
103
104 # 寻找并返回信息收益最大的属性划分
105 # 信息收益值划分该数据集前后的熵减
106 # 计算公式为： $Gain(A) = get\_h(ori\_target) - \sum(|sub\_target| / |ori\_target|$ 
107 #  $* get\_h(sub\_target))$ 
108 def best_spilt(dataset):
109     base_h = get_h(dataset.target)
110     best_gain = 0.0
111     best_feature = None
112     for i in range(len(dataset.feature_names)):
113         feature_range = []
114         for j in range(len(dataset.data)):
115             if dataset.data[j][i] not in feature_range:
116                 feature_range.append(dataset.data[j][i])
117
118         spilt_h = 0.0
119         for feature_value in feature_range:

```

```

119         subset = get_subset(dataset, dataset.feature_names[i],
120                               feature_value)
121         spilt_h += len(subset.target) / len(dataset.target) *
122         get_h(subset.target)
123
124         if best_gain <= base_h - spilt_h:
125             best_gain = base_h - spilt_h
126             best_feature = dataset.feature_names[i]
127
128         return best_feature
129
130 # 返回数据集中一个数据最可能的标签
131 def vote_most(dataset):
132     target_range = {}
133     best_target = None
134     best_vote = 0
135
136     for t in dataset.target:
137         if t not in target_range.keys():
138             target_range[t] = 1
139         else:
140             target_range[t] += 1
141
142     for t in target_range.keys():
143         if target_range[t] > best_vote:
144             best_vote = target_range[t]
145             best_target = t
146
147     return best_target
148
149 # 返回测试的正确率
150 # predict_result: 预测标签列表, target_result: 实际标签列表
151 def accuracy_rate(predict_result, target_result):
152     # print("Predict Result: ", predict_result)
153     # print("Target Result: ", target_result)
154     accuracy_score = 0
155     for i in range(len(predict_result)):
156         if predict_result[i] == target_result[i]:
157             accuracy_score += 1
158     return accuracy_score / len(predict_result)
159
160 # 决策树的节点结构
161 class dt_node:
162
163     def __init__(self, content, is_leaf=False, parent=None):
164         global nonce
165         self.id = nonce # 为节点赋予一个全局ID, 目的是方便画图
166         nonce += 1
167         self.feature_name = None
168         self.target_value = None
169         self.vote_most = None # 记录当前节点最可能的标签
170         if not is_leaf:
171             self.feature_name = content # 非叶子节点的属性名
172         else:
173             self.target_value = content # 叶子节点的标签
174
175         self.parent = parent
176         self.child = {} # 以当前节点的属性对应的属性值作为键值

```

```

175
176 # 决策树模型
177 class dt_tree:
178
179     def __init__(self):
180         self.tree = None # 决策树的根节点
181         self.map_str = ""
182         digraph demo{
183             node [shape=box, style="rounded", color="black",
fontname="Microsoft YaHei"];
184             edge [fontname="Microsoft YaHei"];
185             """ # 用于作图: pydotplus 格式的树图生成代码结构
186             self.color_dir = {} # 用于作图: 叶子节点可选颜色, 以标签值为键值
187
188         # 训练模型, train_set: 训练集
189         def fit(self, train_set):
190
191             if len(train_set.target) <= 0: # 如果测试集数据为空, 则返回空节点, 结束
递归
192                 return None
193
194             target_all_same = True
195             for i in train_set.target:
196                 if i != train_set.target[0]:
197                     target_all_same = False
198                     break
199
200             if target_all_same: # 如果测试集数据中所有数据的标签相同, 则构造叶子节点,
结束递归
201                 node = dt_node(train_set.target[0], is_leaf=True)
202                 if self.tree == None: # 如果根节点为空, 则让该节点成为根节点
203                     self.tree = node
204
205                 # 用于作图, 更新 map_str 内容, 为树图增加一个内容为标签值的叶子节点
206                 node_content = "标签: " + str(node.target_value)
207                 self.map_str += "id" + str(node.id) + "[label=\"" +
node_content + "\", fillcolor=\"" + self.color_dir[node.target_value] +
"\", style=filled]\n"
208
209                 return node
210             elif len(train_set.feature_names) == 0: # 如果测试集待考虑属性为空, 则
构造叶子节点, 结束递归
211                 node = dt_node(vote_most(train_set), is_leaf=True) # 这里让叶子
结点的标签为概率上最可能的标签
212                 if self.tree == None: # 如果根节点为空, 则让该节点成为根节点
213                     self.color_dir[vote_most(train_set)] = color_set[0]
214                     self.tree = node
215
216                 # 用于作图, 更新 map_str 内容, 为树图增加一个内容为标签值的叶子节点
217                 node_content = "标签: " + str(node.target_value)
218                 self.map_str += "id" + str(node.id) + "[label=\"" +
node_content + "\", fillcolor=\"" + self.color_dir[node.target_value] +
"\", style=filled]\n"
219
220                 return node
221             else: # 普通情况, 构建一个内容为属性的非叶子节点
222                 best_feature = best_spilt(train_set) # 寻找最优划分属性, 作为该结点
的值

```

```

223         best_feature_index = -1
224         for i in range(len(train_set.feature_names)):
225             if train_set.feature_names[i] == best_feature:
226                 best_feature_index = i
227                 break
228
229         node = dt_node(best_feature)
230         node.vote_most = vote_most(train_set)
231         if self.tree == None: # 如果根节点为空,则让该节点成为根节点
232             self.tree = node
233             # 用于作图, 初始化叶子节点可选颜色
234             for i in range(len(train_set.target)):
235                 if train_set.target[i] not in self.color_dir:
236                     global color_i
237                     self.color_dir[train_set.target[i]] =
color_set[color_i]
238                     color_i += 1
239                     color_i %= len(color_set)
240
241             feature_range = [] # 获取该属性出现在数据集中的可选属性值
242             for t in train_set.data:
243                 if t[best_feature_index] not in feature_range:
244                     feature_range.append(t[best_feature_index])
245
246             # 用于做图, 创建一个内容为属性的非叶子节点
247             node_content = "属性: " + node.feature_name
248             self.map_str += "id" + str(node.id) + "[label=\"" +
node_content + "\", fillcolor=\"#AADDFF\", style=filled]\n"
249
250             for feature_value in feature_range:
251                 subset = get_subset(train_set, best_feature, feature_value)
252                 # 获取每一个子集
253                 node.child[feature_value] = self.fit(subset) # 递归调用 fit
函数生成子节点
254                 if node.child[feature_value] == None:
255                     # 如果创建的子节点为空, 则创建一个叶子节点作为其子节点, 其中标签
值为概率上最可能的标签
256                     node.child[feature_value] =
dt_node(vote_most(train_set), is_leaf=True)
257                     node.child[feature_value].parent = node
258
259                     # 用于做图, 创建当前节点到所有子节点的连线
260                     self.map_str += "id" + str(node.id) + " -> " + "id" +
str(node.child[feature_value].id) + "[label=\"" + str(feature_value) +
"\"]\n"
261
262                     # print("Rest Feature: ", train_set.feature_names)
263                     # print("Best Feature: ", best_feature_index, best_feature,
"Feature Range: ", feature_range)
264                     # for feature_value in feature_range:
265                     #     print("Child[" + feature_value + "]: ",
node.child[feature_value].feature_name,
node.child[feature_value].target_value)
266                 return node
267
268             # 测试模型, 对测试集 test_set 进行预测
269             def predict(self, test_set):
270                 test_result = []

```

```

270         for test in test_set.data:
271             node = self.tree # 从根节点一只往下找，知道到达叶子节点
272             while node.target_value == None:
273                 feature_name_index = -1
274                 for i in range(len(test_set.feature_names)):
275                     if test_set.feature_names[i] == node.feature_name:
276                         feature_name_index = i
277                         break
278                 if test[feature_name_index] not in node.child.keys():
279                     break
280                 else:
281                     node = node.child[test[feature_name_index]]
282
283             if node.target_value == None:
284                 test_result.append(node.vote_most)
285             else: # 如果没有到达叶子节点，则取最后到达节点概率上最可能的标签为目标值
286                 test_result.append(node.target_value)
287
288         return test_result
289
290     # 输出树，生成图片，path: 图片的位置
291     def show_tree(self, path="demo.png"):
292         map = self.map_str + "}"
293         # print(map)
294         graph = pdp.graph_from_dot_data(map)
295         graph.write_png(path)
296
297     # 学习曲线评估算法精度 dataset: 数据练集, label: 纵轴的标签, interval: 测试规模递增
    的间隔
298     def incremental_train_scale_test(dataset, label, interval=1):
299         c = dataset
300         r = range(5, len(c.data) - 1, interval)
301         rates = []
302         for train_num in r:
303             print(train_num)
304             train_set = new_dataset(c.feature_names, c.target_names,
305                                     c.data[:train_num], c.target[:train_num])
306             test_set = new_dataset(c.feature_names, c.target_names,
307                                    c.data[train_num:], c.target[train_num:])
308             dt = dt_tree()
309             dt.fit(train_set)
310             rates.append(accuracy_rate(dt.predict(test_set),
311                                           list(test_set.target)))
312
313         print(rates)
314         plt.plot(r, rates)
315         plt.ylabel(label)
316         plt.show()
317
318     if __name__ == '__main__':
319
320         c = load_car() # 载入汽车数据集
321         # c = load_mushroom() # 载入蘑菇数据集
322         train_num = 1000 # 训练集规模(剩下的数据就放到测试集)
323         train_set = new_dataset(c.feature_names, c.target_names,
324                                 c.data[:train_num], c.target[:train_num])
325         test_set = new_dataset(c.feature_names, c.target_names,
326                                c.data[train_num:], c.target[train_num:])

```

```
322
323     dt = dt_tree() # 初始化决策树模型
324     dt.fit(train_set) # 训练
325     dt.show_tree("../image/demo.png") # 输出决策树图片
326     print(accuracy_rate(dt.predict(test_set), list(test_set.target))) # 进
    行测试，并计算准确率吧
327
328     # incremental_train_scale_test(load_car(), "car")
329     # incremental_train_scale_test(load_mushroom(), "mushroom",
    interval=20)
```