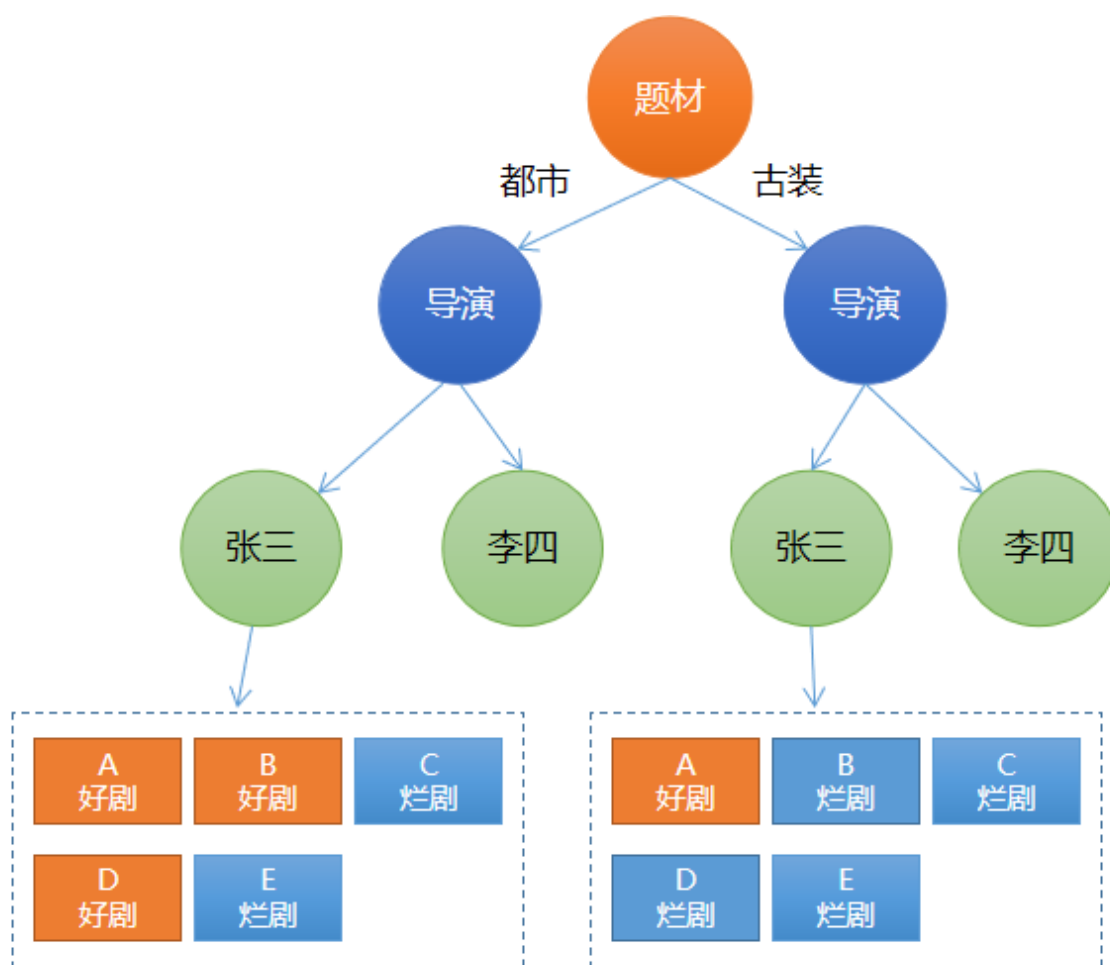


# 一、决策树分类

## 1. 决策树分类原理

决策树可以作为分类器使用，从而实现分类问题。使用决策树实现分类时，先根据不同特征将样本划分到不同叶子节点，再根据投票法（少数服从多数）决定预测结果属于哪个类别，预测类别即同一个子节点下数量最多的样本的类别。例如，预测一部电视剧是好剧还是烂剧，可以用以下决策树来实现：



根据以上决策树模型，可以做这样的分类预测：

- 如果是都市题材，导演为张三，好剧的数量较多，预测结果为好剧；
- 如果是古装题材，导演为张三，烂剧的数据量较多，预测结果为烂剧。

## 2. 决策树分类案例

sklearn中，提供了以下API实现决策树分类器或随机森林分类器：

```

1 sklearn.tree.DecisionTreeClassifier() # 决策树分类器
2 sklearn.ensemble.RandomForestClassifier() # 随机森林分类器

```

该示例中，根据一组小汽车样本数据，该组样本数据统计了小汽车常见特征信息及其分类，特征包括：汽车价格、维修费用、车门数量、载客数、后备箱、安全性，标签为汽车质量。各属性取值如下表所示：

索引-名称	取值范围	含义说明
1-buying	vhigh, high, med, low	购买价格
2-maint	vhigh, high, med, low	维护费用
3-doors	2, 3, 4, 5more	车门数量
4-persons	2, 4, more	载客数
5-lug_boot	small, med, big	后备箱大小
6-safety	low, med, high	安全性
标签	unacc, acc, good, vgood	汽车质量

以下是使用决策树进行分类的示例代码。

```

1 # 决策树分类示例
2 import numpy as np
3 import sklearn.preprocessing as sp
4 import sklearn.ensemble as se
5 import sklearn.model_selection as ms
6
7 raw_samples = [] # 保存一行样本数据
8 with open("../data/car.txt", "r") as f:
9     for line in f.readlines():
10         raw_samples.append(line.replace("\n",
11                                         "").split(","))
12
13 data = np.array(raw_samples).T # 转置

```

```

13 encoders = [] # 记录标签编码器
14 train_x = [] # 编码后的x
15
16 # 对样本数据进行标签编码
17 for row in range(len(data)):
18     encoder = sp.LabelEncoder() # 创建标签编码器
19     encoders.append(encoder)
20     if row < len(data) - 1: # 不是最后一行，为样本特征
21         lbl_code = encoder.fit_transform(data[row]) # 编
22         train_x.append(lbl_code)
23     else: # 最后一行，为样本输出
24         train_y = encoder.fit_transform(data[row])
25
26 train_x = np.array(train_x).T # 转置回来，变为编码后的数组
27 print(train_x)
28
29 # 创建随机森林分类器
30 model = se.RandomForestClassifier(max_depth=8, # 最大树高
31                                   n_estimators=150, # 评
32                                   估系数
33                                   random_state=7) # 随机种
34 # 训练
35 model.fit(train_x, train_y)
36
37 print("accuracy:", model.score(train_x, train_y)) # 打印平
38 均精度
39
40 # 预测
41 ## 待预测数据
42 data = [['high', 'med', '5more', '4', 'big', 'low'],
43         ['high', 'high', '4', '4', 'med', 'med'],
44         ['low', 'low', '2', '2', 'small', 'high'],
45         ['low', 'med', '3', '4', 'med', 'high']]
46 data = np.array(data).T
47 test_x = []
48 for row in range(len(data)):
49     encoder = encoders[row] # 取得每列对应的标签编码器

```

```
47 |     test_x.append(encoder.transform(data[row])) # 待预测数
    |     据编码
48 |
49 | test_x = np.array(test_x).T # 转置回来
50 |
51 | pred_test_y = model.predict(test_x) # 执行预测
52 |
53 | pred_test_y = encoders[-1].inverse_transform(pred_test_y)
    |     # 预测结果反向编码
54 | print("pred_test_y:", pred_test_y) # 预测结果
```

执行结果：

```
1 | accuracy: 0.9704861111111112
2 | pred_test_y: ['unacc' 'acc' 'unacc' 'vgood']
```