

## 决策树

决策树是一个树形结构（可以是二叉树，或非二叉树）  
内部节点表示一个属性上的判断，每个分支表示一个判断结果的输出，叶子节点代表分类结果

有监督的多分类模型

决策树的生成步骤：

- 1) 节点的分裂
- 2) 阈值的确立

目的是寻找一个特征，用该特征划分之后，数据集的熵最小。依次选择特征，直至完全分类。

### ★ ID3.5

entropy

信息论中有熵的概念，熵越大越混乱

熵的变化可看做是信息增益

ID3.5的核心思想是以信息增益做为属性选择的标准，选择分裂后信息增益最大的属性分裂

设数据集为  $D$ , 其熵为

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

其中  $p_i$  表示第  $i$  个类别在  $D$  中出现的概率.

将训练数据集  $D$  按照  $A$  属性进行划分

则  $A$  对  $D$  划分的条件熵为

$$\text{Entropy}(D|A) = \sum_{j=1}^v \frac{|D_j|}{|D|} \text{Entropy}(D_j) \quad (2)$$

其中  $v$  表示属性  $A$  的特征值个数,  $D_j$  表示按照第  $j$  个特征值划分的数据集

因此属性  $A$  的信息增益为

$$\text{gain}(A) = \text{Entropy}(D) - \text{Entropy}(D|A) \quad (3)$$

★ C4.5

C4.5 算法实现了分裂信息和信息增益率

属性  $A$  对应的分裂信息为

$$\text{Split}(A) = - \sum_{j=1}^v \frac{|D_j|}{|D|} \log_2 \left( \frac{|D_j|}{|D|} \right) \quad (4)$$

信息增益率为  $\text{gain-rate} = \frac{\text{gain}(A)}{\text{Split}(A)} \quad (5)$

C4.5的核心思想是选择信息增益率最大的属性进行分裂

☆ 实现技巧 某个属性的特征值对应的记录数  
决策树的构造是递归的，所以需要停止条件  
一种可行的办法是当当前节点中的记录数低于  
固定阈值时停止，选择概率最大的类别  
作为当前叶节点的分类。