

第4章 决策树

4.1 决策树算法(递归)

递归终止:

- ① 当前结点样本属于同一类, 无需划分
- ② 属性为空, 或所有样本在所有属性取值相同: 无法划分
- ③ 当前结点的样本集合为空, 不须划分
- ④ 利用当前结点的经验分布
- ⑤ 父结点的样本分布作为当前结点的后验分布

4.2 划分选择

关键在于选择划分属性 (结点纯度变高)

4.2.1 信息增益

D中k类占 p_k , 则 $Ent(D) = - \sum_{k=1}^{|Y|} p_k \log_2 p_k$
 $Ent(D) \downarrow$ 纯度 \uparrow

设属性a有V个取值 $\{a^1, a^2, \dots, a^V\}$
 共有V个结点. 每个结点计算 D^v
 的信息熵 (D^v 含 a^v 取值的样本)
 然后加上权重 $\frac{|D^v|}{|D|}$

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

Gain越大 "纯度提升" 越大

4.2.2 增益率:

注意: 信息增益相对对取值数目多的属性有偏好.

C4.5 使用信息增益率

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

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

通常而言, a^m 取值数目越多, 则 $IV(a)$ 通常越大
 所以增益率对数目较多的属性有偏好

所以C4.5使用启发式, 先选信息增益高于平均水平 \bar{g} , 再从这些中选信息增益率较高的
 为步 $\left\{ \begin{array}{l} \text{为选高于平均水平的信息增益} \\ \text{后选信息增益率较高的} \end{array} \right.$

4.2.3 基尼指数

$$\text{基尼指数: } Gini(D) = 1 - \sum_{k=1}^{|Y|} p_k^2$$

随机抽两个样本, 不一致的概率

$Gini(D) \downarrow \rightarrow$ 纯度越高

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

选择使得划分后 Gini 最小的作为最优属性, $a^* = \arg \min Gini_index(D, a)$

4.3 剪枝处理

防止过拟合

贪心: 按划分前估计, 划分能否提升泛化性能

后剪枝: 先得到完整的决策树, 自底向上, 将结点替换成叶结点能否提升泛化性能

性能

如何估计泛化性能: 留一法 $\left\{ \begin{array}{l} \text{训练} \\ \text{验证} \end{array} \right.$

为: 后剪枝优于前剪枝: 保留了更多的分枝, 欠拟合风险更小, 但训练时间开销更大

4.4 连续属性决策值

"二分法"

连续属性a确定了n个取值 $\{a^1, \dots, a^n\}$

n-1个候选划分点

$$T_a = \left\{ \frac{a_i + a_{i+1}}{2} \mid 1 \leq i \leq n-1 \right\}$$

每次基于划分点二分

$$Gain(D, a) = \max_{t \in T_a} Gain(D, a, t) = \max_{t \in T_a} Ent(D) - \sum_{i \in \{1, 2\}} \frac{|D_i|}{|D|} Ent(D_i)$$

与离散属性不同, 当前候选划分点为连续属性, 之后还可使用