

决策树: 根结点  $x_1$ , 内部结点  $x_n$ , 叶结点  $x_n$

属性测试

决策结果

## 生成决策树, 递归过程

- base case:
- ① 当前结点所包含样本属于同一类型 无需划分
  - ② 属性集为空  $\rightarrow$  majority vote 利用当前结点的后验分布
  - ③ 样本集为空  $\rightarrow$  父节点样本集 majority vote 当前结点的先验分布.

## 划分选择

### 1. 信息增益

信息熵

$$\text{Entropy: Ent}(D) = - \sum_{k=1}^{|D|} p_k \log_2 p_k$$

$p_k$ : 第  $k$  类样本所占比例

Gain  $\uparrow$  使用  $a$  来划分所获得纯度提升个

信息增益

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

### 2. 增益率

$$\text{Gain-ratio} = \frac{\text{Gain}(D, a)}{\text{IV}(a)}$$

对可取值数目较少的属性有偏好

候选划分属性中找出信息增益高于均值的属性

再从中选增益率最高的

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

### 3. 基尼指数

$$\begin{aligned} \text{Gini}(D) &= \sum_{k=1}^{|D|} \sum_{k'=1}^{|D|} p_k p_{k'} \\ &= 1 - \sum_{k=1}^{|D|} p_k^2 \end{aligned}$$

Gini  $\downarrow$  选

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

## 剪枝处理

- 预剪枝  $\rightarrow$  欠拟合风险 (贪心策略)
- 后剪枝  $\rightarrow$  开销大

## 连续值 & 缺失值

连续值

将该属性的值排序, 选取不同点计算信息增益, 选择最优点作为划分点

缺失值

为那些无缺失值样本赋权重.

## 多变量决策树

单变量: 轴平行 分类边界由若干个与坐标轴平行的分段组成

多变量: 斜划分  $\sum_{i=1}^d w_i a_i = c$  的线性分类器