

决策树流程

● 输入:

训练集 $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$

属性集 $A = \{a_1, a_2, \dots, a_d\}$

● 生成 `TreeGenerate(D, A)` 的具体实现:

- 若 D 全属于同一类别 C , 则将该分支标记为叶结点, 类别标记为 C
- 若 $A = \emptyset$ (属性集为空) 或者 D 在 A 上的所有属性的取值都相同, 则将该分支标记为叶结点, 类别标记为 D 中最多的类

① 从 A 中选出最优划分属性 a_*

② for (a_* 的每个属性值 a_*^v) :

令 D_v 为 D 中属性 a_* 的取值为 a_*^v 的子集

if D_v 为空 (即该属性值不含任何样本) :

- 将该分支标记为叶结点, 类别标记为 D 中最多的类 —— 停止条件

else (非空)

生成 `TreeGenerate($D_v, A \setminus \{a_*\}$)` —— 递归 (对子集继续进行分支)

说明: 决策树的生成是运用了递归, 在递归式里有3个停止条件。

(1) 信息熵

$$\text{Ent}(D) = - \sum_{k=1}^K p_k \log_2 p_k \quad k=\{1,2,\dots,K\} \text{——最终评判的分类类别}$$

作用：度量样本集的不确定性（区分度）。值越大，区分度越弱。

范围： $[0, \log_2 K]$

当 $p_i=1$ ，其余 $p_j(j \neq i)=0$ 时， $\text{Ent}(D)=0$

当 $p_1 = p_2 = \dots = p_K=1/K$ 时， $\text{Ent}(D)=\log_2 K$

(2) 信息增益

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

解释：使用属性a进行划分，能提升的区分度的增加程度。值越大，说明该属性越能区分样本。

特点：对可取值数目较多的属性有所偏好。

2、增益率

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

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

特点：对可取值数目较少的属性有所偏好。

3、基尼指数

(1) 基尼值

$$\text{Gini}(D) = 1 - \sum_{k=1}^K p_k^2 \quad \text{—— 从训练集D中任意取2个样本，2个样本类别不一致的概率。}$$

基尼值越大，样本区分度越高。

(2) 基尼指数

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