

决策树

2020年4月24日 星期五 下午7:25

决策树基本算法

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

属性集: $A = \{a_1, a_2, \dots, a_d\}$

过程: 函数 $\text{TreeNode}(D, A)$

1. node;
2. if D 中样本属于同一类 C
→ node 标记为 C : return, endif
3. if $A = \emptyset$ or 样本在 A 上取值相同
→ node 标记为叶结点, 类标为样本中最多的一类
return, endif

4. 从 A 中选择最优划为属性 a_x ;

for a_x 中每一个值 a_x^v do

为 node 生成分支, D_v 表示 D 在 a_x 上取 a_x^v 子集

if $D_v = \emptyset$ → 分支结点标为最多的一类 return

else: $\text{TreeGenerate}(D_v, A \setminus \{a_x\})$ 再生成结点
endif, end for

output: node 为根结点的决策树

剪枝处理: 过拟合

信息熵: $\text{Ent}(D) = -\sum_{k=1}^{|V|} P_k \cdot \log_2 P_k$ ($\text{Ent}(D)$, D 纯度越高

信息熵增益: $\text{Gain}(D, a) = \text{Ent}(D) - \sum_{v=1}^{|V|} \frac{|D^v|}{|D|} \text{Ent}(D^v)$
(越大, 纯度越高)

$a^* = \arg \max_{a \in A} \text{Gain}(D, a)$

权重: 第 v 个分支点取值为 a^v 的样本 D^v

以上数据为例: 西瓜: 好、坏两种, 为划 $\frac{8}{17}$, $\frac{9}{17}$

$$\text{Ent}(D) = -\sum_{k=1}^2 P_k \cdot \log_2 P_k = -(\frac{8}{17} \cdot \log_2 \frac{8}{17} + \frac{9}{17} \cdot \log_2 \frac{9}{17}) = 0.998$$

色泽增益: 三类信息 $0 \log 0 = 0$

$$\text{Ent}(D^1) = 1.000, \text{Ent}(D^2) = 0.918, \text{Ent}(D^3) = 0.722$$

$$\text{Gain}(D, \text{色泽}) = \text{Ent}(D) - \sum_{i=1}^3 \frac{|D^i|}{|D|} \text{Ent}(D^i)$$

$$= 0.998 - (\frac{6}{17} \times 1.0 + \frac{6}{17} \times 0.918 + \frac{5}{17} \times 0.722) = 0.109$$

同理

$$\text{Gain}(D, \text{纹理}) = 0.143 = \text{Gain}(D, \text{纹理}) = 0.38 \dots \text{Gain}(D, \text{形状}) = 0.006$$

纹理 $\text{Gain}(D)$ max → 划为

