

# Decision Tree

①	$X_1$	$X_2$	
$Y_1$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{2}$
$Y_2$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{2}$
	$\frac{1}{2}$	$\frac{1}{2}$	

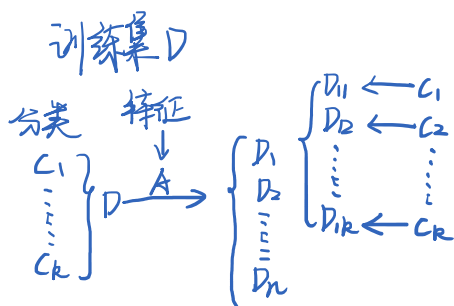
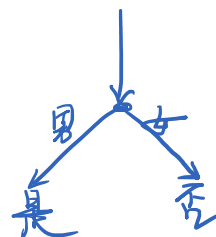
②	$X_1$	$X_2$	
$Y_1$	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{1}{2}$
$Y_2$	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{1}{2}$
	$\frac{1}{4}$	$\frac{3}{4}$	

③	$X_1$	$X_2$	
$Y_1$	$\frac{1}{8}$	$\frac{1}{8}$	$\frac{1}{4}$
$Y_2$	$\frac{3}{8}$	$\frac{3}{8}$	$\frac{3}{4}$
	$\frac{1}{2}$	$\frac{1}{2}$	

$$① H(Y|X) = -(\frac{1}{2} \cdot \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} \log_2 \frac{1}{2}) = 0.5$$

$$② H(Y|X) = -(\frac{1}{4} \cdot \frac{1}{2} \log_2 \frac{1}{2} + \frac{3}{4} \cdot \frac{1}{2} \log_2 \frac{1}{2}) = 0.5$$

$$③ H(Y|X) = -(\frac{1}{2} \cdot \frac{1}{4} \log_2 \frac{1}{4} + \frac{1}{2} \cdot \frac{3}{4} \log_2 \frac{3}{4}) = 0.4056$$



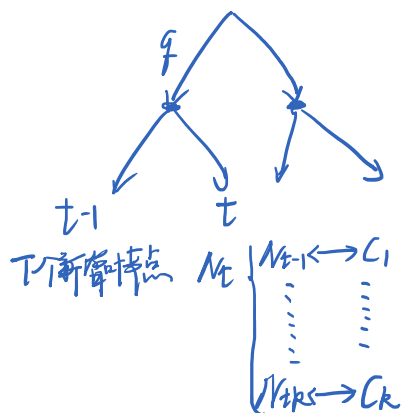
$$① H(D) - H(D|A) = g(D|A) \text{ 信息增益 } ID_3 \begin{cases} \text{阈值} \\ g(D|A) \end{cases}$$

$$② \frac{g(D|A)}{H(D)} = g_R(D|A) \text{ 信息增益比 } C_{4.5} \begin{cases} \text{阈值} \\ g_R(D|A) \end{cases}$$

$$\text{经验熵 } H(D) = -\sum_{k=1}^K \frac{|C_k|}{|D|} \log_2 \frac{|C_k|}{|D|}$$

$$\text{特征A的经验熵 } H(D|A) = -\sum_{i=1}^n \left( \frac{|D_i|}{|D|} \sum_{k=1}^K \frac{|D_{ik}|}{|D_i|} \log_2 \frac{|D_{ik}|}{|D_i|} \right)$$

$$H_A(D) = -\sum_{i=1}^n \frac{|D_i|}{|D|} \log_2 \frac{|D_i|}{|D|}$$



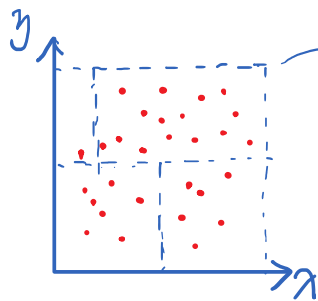
$$Q(T_A) - Q(T_B) = \bar{q}_A - \bar{q}_B$$

$$= \text{节点数} \times \text{该操作减少的经验熵} - \alpha (\text{该操作增加的节点数} - 1)$$

$$= -|N_t| \sum_{k=1}^K \frac{|N_{tk}|}{|N_t|} \log_2 \frac{|N_{tk}|}{|N_t|} - (-\sum_{t=1}^T |N_t| \sum_{k=1}^K \frac{|N_{tk}|}{|N_t|} \log_2 \frac{|N_{tk}|}{|N_t|}) - \alpha (\text{分支增加的叶节点数} - 1)$$

$$\frac{Q(T_A) - Q(T_B)}{|N_t|} = (\text{父节点的经验熵} - \text{子节点的经验熵期望}) - \alpha \frac{\text{子节点数} - 1}{|N_t|}$$

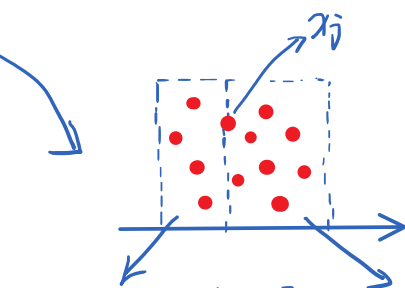
# CART 回归树



分成  $M$  单元  $R_j$   
对应一个值  $c_1, \dots, c_M$

$$f(x) = \sum_{m=1}^M c_m I(x \in R_m)$$

$$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$$



$$R_1(j, s) = \{x | x^{(j)} \leq s\}$$

$$R_2(j, s) = \{x | x^{(j)} > s\}$$

$$\hat{c}_1 = \text{ave}(y_i | x_i \in R_1(j, s))$$

$$\hat{c}_2 = \text{ave}(y_i | x_i \in R_2(j, s))$$

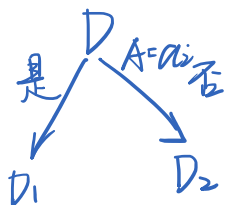
① 找一点  $x_j^{(j)}$  =  
$$\min_{j, s} [\min_{c_1, x \in R_1} \sum (y_i - c_1)^2 + \min_{c_2, x \in R_2} \sum (y_i - c_2)^2]$$

② 其对应值  $s$

$$Gini(D) = 1 - \sum_{k=1}^K p_k^2 = 1 - \sum_{k=1}^K \left(\frac{|C_k|}{|D|}\right)^2$$

$$Gini(D, A) = \frac{|D_1|}{|D|} Gini(D_1) + \frac{|D_2|}{|D|} Gini(D_2)$$

CART = 决策树  $\rightarrow$  二叉树



计算  $Gini(D, A = a_i)$  取最小的特征及对应值

$B = b_j, \dots$

停止方法  $\left\{ \begin{array}{l} \text{结点样本数} < \text{阈值} \\ Gini \text{指数} < \text{阈值} \end{array} \right.$

CART 决策树 = 减枝 =

① 如何理解熵, 基尼指数?

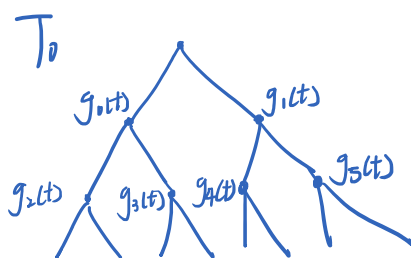
② 回归和拟合区别?

③ 减枝时, 损失函数  $\alpha$  如何定?

( $\alpha > 0$ , 但又取什么最优?)

$$\alpha = g(t) = \frac{c(t) - c(T)}{|T| - 1}$$

(当父节点和子节点损失函数相等时,  $\alpha = 1$ )



重复  $\rightarrow$  减枝节点在  $\{g_i(t) | i=1, \dots, 5\}$   
得到  $T_1, \dots, T_n$

只有根节点和两个子节点

再  $\{T_0, \dots, T_n\}$  中找最优子树  $T_a$

设  $m \geq 2$ ,  $m$  层 = 叉树有  $2^{(m-1)} - 2$  个内部节点

则  $n = [2, 2^{(m-1)} - 2]$ ,  $n$  为整数, 进行了  $[3, 2^{(m-1)} - 1]$  次验证