

一、计算题

对下表分别用 ID3 和 CART 算法生成决策树（要求写出详细的计算步骤）。

ID	年龄	有工作	有自己的房子	信贷情况	类别
1	青年	否	否	一般	否
2	青年	否	否	好	否
3	青年	是	否	好	是
4	青年	是	是	好	是
5	青年	是	是	非常好	是
6	青年	是	是	一般	是
7	青年	否	否	一般	否
8	中年	否	否	一般	否
9	中年	否	否	好	否
10	中年	是	是	好	是
11	中年	否	是	非常好	是
12	中年	否	是	非常好	是
13	中年	是	是	好	是
14	老年	否	是	非常好	是
15	老年	否	是	好	是
16	老年	是	否	好	是
17	老年	是	否	非常好	是
18	老年	否	否	一般	否

答案：

分别以 A_1 、 A_2 、 A_3 、 A_4 表示特征“年龄”、“有工作”、“有房子”、“信贷情况”，

ID3 决策树生成过程如下：

计算各特征对数据集 D 的信息增益：

$$\text{经验熵为： } H(D) = -\left(\frac{12}{18} \log_2 \frac{12}{18} + \frac{6}{18} \log_2 \frac{6}{18}\right) = 0.9183$$

$$G(D, A_1) = H(D) - H(D, A_1)$$

$$= 0.9183 - \left[\frac{7}{18} \times \left(-\frac{4}{7} \log_2 \frac{4}{7} - \frac{3}{7} \log_2 \frac{3}{7} \right) + \frac{6}{18} \times \left(-\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right) + \frac{5}{18} \times \left(-\frac{4}{5} \log_2 \frac{4}{5} - \frac{1}{5} \log_2 \frac{1}{5} \right) \right]$$

$$= 0.028$$

$$G(D, A_2) = H(D) - H(D, A_2)$$

$$= 0.9183 - \left[\frac{8}{18} \times 0 + \frac{10}{18} \times \left(-\frac{4}{10} \log_2 \frac{4}{10} - \frac{6}{10} \log_2 \frac{6}{10} \right) \right]$$

$$= 0.379$$

$$G(D, A_3) = H(D) - H(D, A_3)$$

$$= 0.9183 - \left[\frac{9}{18} \times 0 + \frac{9}{18} \times \left(-\frac{3}{9} \log_2 \frac{3}{9} - \frac{6}{9} \log_2 \frac{6}{9} \right) \right]$$

$$= 0.459 \quad (\text{信息增益最大})$$

$$\begin{aligned}
G(D, A_4) &= H(D) - H(D, A_4) \\
&= 0.9183 - \left[\frac{5}{18} \times \left(-\frac{1}{5} \log_2 \frac{1}{5} - \frac{4}{5} \log_2 \frac{4}{5} \right) + \frac{8}{18} \times \left(-\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8} \right) + \frac{5}{18} \times 0 \right] \\
&= 0.357
\end{aligned}$$

$G(D, A_3)$ 最大，因此选择特征 A_3 = “有房子” 作为根节点

(1) 分支 1: “有房子”=是时，类别也全部为“是”，因此可以作为叶子，

(2) 分支 2: “有房子”=否时，继续计算后续的信息增益：

$$H(D_2) = -\left(\frac{3}{9} \log_2 \frac{3}{9} + \frac{6}{9} \log_2 \frac{6}{9} \right) = 0.918$$

$$H(D_2, A_1) = -\frac{4}{9} \times \left(\frac{1}{4} \log_2 \frac{1}{4} + \frac{3}{4} \log_2 \frac{3}{4} \right) - \frac{2}{9} \times 0 - \frac{3}{9} \times \left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right) = 0.667$$

$$H(D_2, A_2) = -\frac{3}{9} \times 0 - \frac{6}{9} \times 0 = 0$$

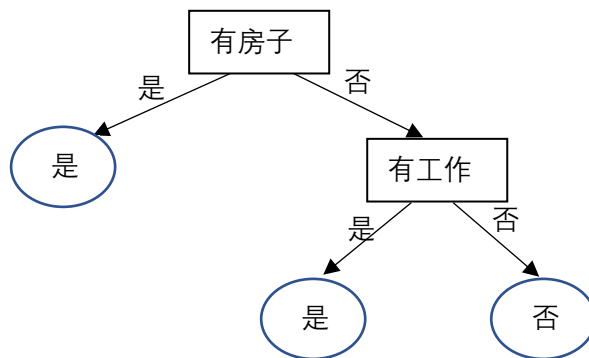
$$H(D_2, A_4) = -\frac{4}{9} \times 0 - \frac{4}{9} \times \left(\frac{2}{4} \log_2 \frac{2}{4} + \frac{2}{4} \log_2 \frac{2}{4} \right) - \frac{1}{9} \times 0 = 0.444$$

$$G(D_2, A_1) = H(D_2) - H(D_2, A_1) = 0.918 - 0.667 = 0.251$$

$$G(D_2, A_2) = H(D_2) - H(D_2, A_2) = 0.918 - 0 = 0.918 \quad (\text{信息增益最大})$$

$$G(D_2, A_4) = H(D_2) - H(D_2, A_4) = 0.918 - 0.444 = 0.474$$

因此 A_2 = “有工作” 为最优特征，作为下一个划分节点：此时，“有工作”=是时，所有样本类别均为“是”；“有工作”=否时，所有样本类别均为“否”。后续全部为叶子。因此构建的 ID3 决策树如下：



CART 决策树生成过程如下：

计算各特征对数据集 D 的基尼指数：

$$\text{Gini}(D, A_1 = \text{“青年”}) = \frac{7}{18} \times 2 \times \frac{4}{7} \times \left(1 - \frac{4}{7} \right) + \frac{11}{18} \times 2 \times \frac{8}{11} \times \left(1 - \frac{8}{11} \right) = 0.423$$

$$\text{Gini}(D, A_1 = \text{“中年”}) = \frac{6}{18} \times 2 \times \frac{4}{6} \times \left(1 - \frac{4}{6} \right) + \frac{12}{18} \times 2 \times \frac{8}{12} \times \left(1 - \frac{8}{12} \right) = 0.444$$

$$\text{Gini}(D, A_1 = \text{“老年”}) = \frac{5}{18} \times 2 \times \frac{4}{5} \times \left(1 - \frac{4}{5} \right) + \frac{13}{18} \times 2 \times \frac{8}{13} \times \left(1 - \frac{8}{13} \right) = 0.431$$

$$\text{Gini}(D, A_2) = \frac{8}{18} \times 2 \times \frac{8}{8} \times (1 - \frac{8}{8}) + \frac{10}{18} \times 2 \times \frac{4}{10} \times (1 - \frac{4}{10}) = 0.267$$

$$\text{Gini}(D, A_3) = \frac{9}{18} \times 2 \times \frac{9}{9} \times (1 - \frac{9}{9}) + \frac{9}{18} \times 2 \times \frac{3}{9} \times (1 - \frac{3}{9}) = 0.222 \quad (\text{基尼指数最小})$$

$$\text{Gini}(D, A_4 = \text{“一般”}) = \frac{5}{18} \times 2 \times \frac{1}{5} \times (1 - \frac{1}{5}) + \frac{13}{18} \times 2 \times \frac{11}{13} \times (1 - \frac{11}{13}) = 0.277$$

$$\text{Gini}(D, A_4 = \text{“好”}) = \frac{8}{18} \times 2 \times \frac{6}{8} \times (1 - \frac{6}{8}) + \frac{10}{18} \times 2 \times \frac{6}{10} \times (1 - \frac{6}{10}) = 0.433$$

$$\text{Gini}(D, A_4 = \text{“非常好”}) = \frac{5}{18} \times 2 \times \frac{5}{5} \times (1 - \frac{5}{5}) + \frac{13}{18} \times 2 \times \frac{7}{13} \times (1 - \frac{7}{13}) = 0.359$$

$\text{Gini}(D, A_3)$ 最小，因此选择特征 $A_3 = \text{“有房子”}$ 作为根节点。

(1) 分支 1: “有房子”=是时，类别也全部为“是”，因此可以作为叶子；

(2) 分支 2: “有房子”=否时，继续计算后续的基尼指数：

$$\text{Gini}(D_2, A_1 = \text{“青年”}) = \frac{4}{9} \times 2 \times \frac{1}{4} \times (1 - \frac{1}{4}) + \frac{5}{9} \times 2 \times \frac{2}{5} \times (1 - \frac{2}{5}) = 0.433$$

$$\text{Gini}(D_2, A_1 = \text{“中年”}) = \frac{2}{9} \times 2 \times \frac{2}{2} \times (1 - \frac{2}{2}) + \frac{7}{9} \times 2 \times \frac{3}{7} \times (1 - \frac{3}{7}) = 0.381$$

$$\text{Gini}(D_2, A_1 = \text{“老年”}) = \frac{3}{9} \times 2 \times \frac{2}{3} \times (1 - \frac{2}{3}) + \frac{6}{9} \times 2 \times \frac{1}{6} \times (1 - \frac{1}{6}) = 0.333$$

$$\text{Gini}(D_2, A_2) = \frac{3}{9} \times 2 \times \frac{3}{3} \times (1 - \frac{3}{3}) + \frac{6}{9} \times 2 \times \frac{6}{6} \times (1 - \frac{6}{6}) = 0 \quad (\text{基尼指数最小})$$

此时，因为 $\text{Gini}(D_2, A_2)$ 最小，达到可能取到的最小值，因此直接选择 $A_2 = \text{“有工作”}$ 作为

下一个最优特征，并且由 $A_2 = \text{“有工作”}$ 作为中间节点，后续分支全部为叶子，从而构建的

CART 决策树与上述 ID3 决策树完全一致。