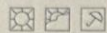


Question 1: Calculate the information gain for splitting CreditScore at 650 in a decision tree classification task, then explain why you would or wouldn't choose this as the root node split.



Mo Tu We Th Fr Sa Su

Memo No.

Date / /

## Practical Problem

Question 1:

- We have 8 records : 4 Low risks
- 4 High risks

$$\begin{aligned} \text{Entropy}(\text{parent}) &= -P(\text{Low}) \log_2 P(\text{Low}) - P(\text{High}) \log_2 P(\text{High}) \\ &= -\left(\frac{4}{8}\right) \log_2 \left(\frac{4}{8}\right) - \left(\frac{4}{8}\right) \log_2 \left(\frac{4}{8}\right) \\ &= 1 \end{aligned}$$

- Split the data based on CreditScore  $\leq 650$   
vs CreditScore  $> 650$

- Left Child (CreditScore  $\leq 650$ ): 3 records,  
all High risk (IDs: 2, 6, 8)
- Right Child (CreditScore  $> 650$ ): 5 records,  
4 Low risk & 1 High risk (IDs: 1, 3, 4, 5, 7)

- Calculate the entropy of each child node

$$\begin{aligned} \text{Entropy}(\text{left}) &= -\frac{0}{3} \log_2 \left(\frac{0}{3}\right) - \left(\frac{3}{3}\right) \log_2 \left(\frac{3}{3}\right) \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{Entropy}(\text{right}) &= -\left(\frac{4}{5}\right) \log_2 \left(\frac{4}{5}\right) - \left(\frac{1}{5}\right) \log_2 \left(\frac{1}{5}\right) \\ &= 0.722 \end{aligned}$$

HẢI TIẾN

- Calculate the information gain

$$\begin{aligned} \text{Info Gain} &= \text{Entropy (parent)} - \text{Weighted Entropy} \\ &= 1 - 0.451 = 0.549 \end{aligned}$$

This split provides an information gain of 0.549, which is quite significant (maximum possible is 1). I would choose this as the root node because:

1. The information gain is substantial (0.549)
2. The split completely separates all high risk cases with Credit Score  $\leq 65$
3. It creates a pure left child node (Entropy 0)
4. The resulting tree would be simple yet effective for this dataset



Question 2: For a regression decision tree predicting CreditScore, calculate the variance reduction when splitting on Age=35, and describe how this splitting criterion differs from information gain in classification trees.

Question 2:

- Calculate the variance of CreditScore for the whole dataset

CreditScore = [720, 650, 750, 600, 780, 630, 710, 640]

$$\text{Mean} = (720 + 650 + 750 + 600 + 780 + 630 + 710 + 640) / 8 = 685$$

$$\begin{aligned} \text{Variance} &= \sum (x - \text{mean})^2 / n \\ &= [(720 - 685)^2 + (650 - 685)^2 + (750 - 685)^2 \\ &\quad + (600 - 685)^2 + (780 - 685)^2 \\ &\quad + (630 - 685)^2 + (710 - 685)^2 \\ &\quad + (640 - 685)^2] / 8 \\ &= 3578 \end{aligned}$$

- Split the data based on Age  $\leq 35$  and Age  $> 35$

• Left child (Age  $\leq 35$ ): [720, 650, 600, 630, 640] (IDs: 1, 2, 4, 6, 8)

• Right child (Age  $> 35$ ): [750, 780, 710] (IDs: 3, 5, 7)

- Calculate the variance for each child node

$$\begin{aligned} \text{Left mean} &= (720 + 650 + 600 + 630 + 640) / 5 \\ &= 648 \end{aligned}$$

HẢI TIẾN



$$\begin{aligned} \bullet \text{ Left variance} &= [(720 - 648)^2 + (650 - 648)^2 \\ &\quad + (600 - 648)^2 + (630 - 648)^2 \\ &\quad + (640 - 648)^2] / 5 \\ &= 1576 \end{aligned}$$

$$\begin{aligned} \bullet \text{ Right mean} &= (750 + 780 + 710) / 3 \\ &= 746.67 \end{aligned}$$

$$\begin{aligned} \bullet \text{ Right variance} &= [(750 - 746.67)^2 \\ &\quad + (780 - 746.67)^2 \\ &\quad + (710 - 746.67)^2] / 3 \\ &= 822.22 \end{aligned}$$

- Calculate the weighted average variance after the split:

$$\begin{aligned} \text{Weight Variance} &= \left(\frac{5}{8}\right) \times \text{Left Variance} \\ &\quad + \left(\frac{3}{8}\right) \times \text{Right Variance} \\ &= 1293.33 \end{aligned}$$

- Calculate the variance reduction

$$\text{Variance Reduction} = \text{Variance (parent)}$$

$$- \text{Weighted Variance}$$

$$= 3575 - 1293.33 = 2281.67$$

HẢI TIẾN



- The difference between variance reduction and information gain:

- Variance reduction is used for regression problems where the target is a continuous value (Credit score)
- Information gain is used for classification problems where the target is a categorical value (Risk Value)
- Variance reduction measures the homogeneity of numeric value after splitting
- Information gain measure the purity of class label labels after splitting

The variance reduction of 2281.67 is significant, showing the Age = 35 split reduces the variance by about 64% of the original variance, making this split effective for a regression tree