

## 1 Datasets

ID	Age	CreditScore	Education	RiskLevel
1	35	720	16	Low
2	28	650	14	High
3	45	750	missing	Low
4	31	600	12	High
5	52	780	18	Low
6	29	630	14	High
7	42	710	16	Low
8	33	640	12	High

Table 1: Training Dataset (8 records)

ID	Age	CreditScore	Education
T1	37	705	16
T2	30	645	missing

Table 2: Test Dataset (2 records)

## 2 Question 1

Let us define  $A$  as the initial training dataset before the split, with 8 records: 4 in class *High* and 4 in class *Low*. We can compute the entropy of  $A$ , denoted by  $H(A)$ , as:

$$H(A) = -\frac{4}{8} \log_2 \frac{4}{8} - \frac{4}{8} \log_2 \frac{4}{8} = 1$$

Splitting  $A$  on *CreditScore* at 650, we obtain the following results:

- Left node ( $CreditScore \leq 650$ ): IDs 2, 4, 6, 8  $\rightarrow$  All are *High*
- Right node ( $CreditScore > 650$ ): IDs 1, 3, 5, 7  $\rightarrow$  All are *Low*

Without calculating the entropy of the dataset after the split, we can immediately obtain the information gain of the split as 1. This is because the dataset is completely pure after the split (all records with *CreditScore* less than or equal to 650 have *HighRiskLevel* and all records with *CreditScore* greater than 650 have *LowRiskLevel*), therefore, the information gain equals the initial entropy before the split.

However, we can calculate the entropy of the dataset after the split as follows:

$$H(CreditScore \leq 650) = -\frac{4}{4} \log_2 \frac{4}{4} = 0$$

$$H(CreditScore > 650) = -\frac{4}{4} \log_2 \frac{4}{4} = 0$$

$$InformationGain(CreditScore = 650) = H(A) - \frac{4}{8} \cdot 0 - \frac{4}{8} \cdot 0 = 1 - 0 - 0 = 1$$

Since the information gain is maximized, we would choose this as a root node split.

## 3 Question 2

// To do

## **4 Question 3**

// To do

## **5 Question 4**

// To do

## **6 Question 5**

// To do

## **7 Question 6**

// To do

## **8 Question 7**

// To do

## **9 Question 8**

// To do

## **10 Question 9**

// To do

## **11 Question 10**

// To do

## **12 Question 11**

// To do

## **13 Question 12**

// To do

## **14 Question 13**

// To do

## **15 Question 14**

// To do

## **16 Question 15**

// To do

## **17 Question 16**

// To do

## **18 Question 17**

// To do

## **19 Question 18**

// To do

## **20 Question 19**

// To do

## **21 Question 20**

// To do