

1.)

ก.
$$\text{gini index} = 1 - \left(\left(\frac{10}{20} \right)^2 + \left(\frac{10}{20} \right)^2 \right) = \frac{1}{2}$$

$$\text{Entropy} = - \left(\frac{5}{10} \right) \log_2 \left(\frac{5}{10} \right) - \left(\frac{5}{10} \right) \log_2 \left(\frac{5}{10} \right) = 1$$

ข. $1 - 20$

$$\text{Gini index} = 1 - (0 + 1^2) = 0$$

$$\text{Information gain} = 1 - 0 = 1$$

$$\text{Entropy} = -0 \log(0) - 1 \log(1) = 0$$

ค. M

$$\text{Gini index} = 1 - \left(\left(\frac{6}{10} \right)^2 + \left(\frac{4}{10} \right)^2 \right) = 0.52$$

$$\text{Information gain} = 1 - 0.97 = 0.03$$

$$\text{Entropy} = - \left(\frac{6}{10} \right) \log_2 \left(\frac{6}{10} \right) - \left(\frac{4}{10} \right) \log_2 \left(\frac{4}{10} \right) = 0.97$$

F

$$\text{Gini index} = 1 - \left(\left(\frac{4}{10} \right)^2 + \left(\frac{6}{10} \right)^2 \right) = 0.52$$

$$\text{Information gain} = 1 - 0.97 = 0.03$$

$$\text{Entropy} = - \left(\frac{4}{10} \right) \log_2 \left(\frac{4}{10} \right) - \left(\frac{6}{10} \right) \log_2 \left(\frac{6}{10} \right) = 0.97$$

ง. Family

$$\text{Gini index} = 1 - \left(\left(\frac{1}{4} \right)^2 + \left(\frac{3}{4} \right)^2 \right) = 0.625$$

$$\text{Information gain} = 1 - 0.81 = 0.19$$

$$\text{Entropy} = - \left(\frac{1}{4} \right) \log_2 \left(\frac{1}{4} \right) - \left(\frac{3}{4} \right) \log_2 \left(\frac{3}{4} \right) = 0.81$$

Sport

$$\text{Gini index} = 1 - \left(\left(\frac{8}{8} \right)^2 + \left(\frac{0}{8} \right)^2 \right) = 0$$

$$\text{Information gain} = 1 - 0 = 1$$

$$\text{Entropy} = - \left(\frac{8}{8} \right) \log_2 \left(\frac{8}{8} \right) - \left(\frac{0}{8} \right) \log_2 \left(\frac{0}{8} \right) = 0$$

Luxury

$$\text{Gini index} = 1 - \left(\left(\frac{7}{8} \right)^2 + \left(\frac{1}{8} \right)^2 \right) = 0.78$$

$$\text{Information gain} = 1 - 0.54 = 0.46$$

$$\text{Entropy} = - \left(\frac{7}{8} \right) \log_2 \left(\frac{7}{8} \right) - \left(\frac{1}{8} \right) \log_2 \left(\frac{1}{8} \right) = 0.54$$

จ. Small

$$\text{Gini index} = 1 - \left(\left(\frac{3}{5} \right)^2 + \left(\frac{2}{5} \right)^2 \right) = 0.52$$

$$\text{Information gain} = 1 - 0.97 = 0.03$$

$$\text{Entropy} = - \left(\frac{3}{5} \right) \log_2 \left(\frac{3}{5} \right) - \left(\frac{2}{5} \right) \log_2 \left(\frac{2}{5} \right) = 0.97$$

Medium

$$\text{Gini index} = 1 - \left(\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right) = 0.51$$

$$\text{Information gain} = 1 - 0.985 = 0.015$$

$$\text{Entropy} = - \left(\frac{3}{7} \right) \log_2 \left(\frac{3}{7} \right) - \left(\frac{4}{7} \right) \log_2 \left(\frac{4}{7} \right) = 0.985$$

Large

$$\text{Gini index} = 1 - \left(\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right) = 0.5$$

$$\text{Information gain} = 1 - 1 = 0$$

$$\text{Entropy} = - \left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) - \left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) = 1$$

Extra large

$$\text{Gini index} = 1 - \left(\left(\frac{2}{4} \right)^2 + \left(\frac{2}{4} \right)^2 \right) = 0.5$$

$$\text{Information gain} = 1 - 1 = 0$$

$$\text{Entropy} = - \left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) - \left(\frac{2}{4} \right) \log_2 \left(\frac{2}{4} \right) = 1$$

ฉ. Customer ID

ช. เหมาะ เพราะ มีความบริสุทธิ์มากที่สุด

4. $\text{Cost}(\text{internal}) = \log m$ $\text{Cost}(\text{leaf}) = \log 3$

Tree1

Tree2

$$\text{Cost}(\text{tree1}) = 2 * \log m + 3 * \log 3$$

$$\text{Cost}(\text{tree2}) = 4 * \log m + 5 * \log 3$$

$$= \log m^2 + \log 3^3$$

$$= \log m^4 3^5$$

$$= \log m^2 3^3$$

$$\text{Cost}(D|\text{tree1}) = 7 * \log n$$

$$\text{Cost}(D|\text{tree2}) = 4 * \log n$$

$$\text{Cost}(\text{tree1}, D) = \log m^2 3^3 + 7 * \log n$$

$$\text{Cost}(\text{tree2}, D) = \log m^4 3^5 + 4 * \log n$$

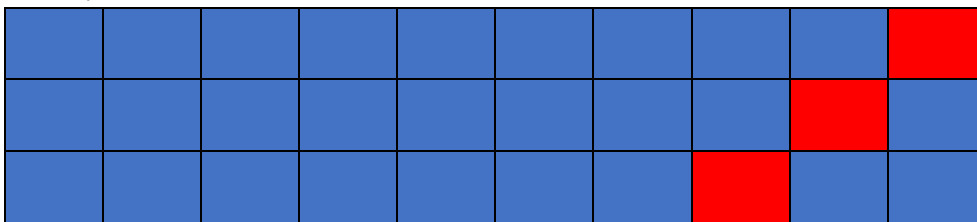
$$= \log m^2 3^3 n^7$$

$$= \log m^4 3^5 n^4$$

$$\text{Cost}(\text{tree1}, D) - \text{Cost}(\text{tree2}, D) = \log m^2 3^3 n^7 - \log m^4 3^5 n^4 = \log m^{-2} 3^{-2} n^3 = 3 \log n - 2 \log 3 - 2 \log m$$

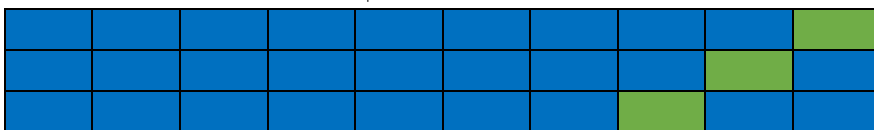
∴ ควรใช้ tree2 เมื่อ $2n \geq 3 + m$ และควรใช้ tree1 เมื่อ $2n \leq 3 + m$

5. แบ่งข้อมูลออกเป็น 10 ส่วน แล้วทำการหาความแม่นยำ 10 รอบโดยแต่ละรอบจะใช้ test set 1 ส่วนที่ไม่ซ้ำกัน



 Train set  Test set

ค่าความแม่นยำจะใช้ค่าเฉลี่ยจากทั้ง 10 รอบ และ train set (กล่องสีฟ้า) จะมีการคัดเลือก model จากการแบ่งเป็น 10 ส่วน โดยใช้ 1 ส่วนที่ไม่ซ้ำกัน 10 รอบ เป็น validation set เพื่อหาค่า error และเลือก model ที่มีค่า error รวมต่ำที่สุด



 Train set  validation set