

Assignment - 8

Example	A ₁	A ₂	A ₃	Output y
x ₁	1	0	0	0
x ₂	1	0	1	0
x ₃	0	1	0	0
x ₄	1	1	1	1
x ₅	1	1	0	1

There are three binary features A₁, A₂, A₃.

Information gain for A₁.

$$\text{Group 1} = x_3$$

$$\text{Group 2} = x_1, x_2, x_4, x_5$$

$$\text{Entropy} = - \sum_{i=1}^2 p_i \log_2 p_i$$

$$E(\text{Group 1}) = - (1 \log_2 1) = 0$$

$$E(\text{Group 2}) = - \left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right)$$
$$= \frac{1}{2} + \frac{1}{2} = 1$$

$$E(\text{split}) = \frac{1}{5} \times 0 + \frac{4}{5} \times 1$$

$$I = 1 - E(\text{split})$$

$$= \frac{1}{5} = 0.2$$

Information gain for A_2

group 1 - x_1, x_2

group 2 - x_3, x_4, x_5

$$\text{Entropy} = - \sum_{i=1}^n P_i \log_2 P_i$$

$$E(\text{group 1}) = -(0) = 0$$

$$E(\text{group 2}) = - \left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3} \right)$$

$$= 0.528 + 0.389$$

$$= 0.917$$

$$E(\text{split}) = \frac{2}{5} \times 0 + \frac{3}{5} \times 0.917$$

$$= 0.5502$$

$$I = 0.97 - 0.5502$$

$$= 0.4498$$

Information gain for A_3

group 1 - x_1, x_3, x_5

group 2 - x_2, x_4

$$\text{Entropy} = - \sum_{i=1}^n P_i \log_2 P_i$$

$$E(\text{group 1}) = - \left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right)$$

$$= 0.917$$

$$E(\text{group 2}) = - \left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right)$$

$$= 1$$

Q3

$$E(\text{Split}) = \frac{3}{5} \times 0.917 + \frac{2}{5} \times 1$$

$$= 0.5502 + 0.4$$

$$= 0.9502$$

$$I = .97 - 0.9502$$

$$= 0.0498$$

Since A_2 split has the highest information gain we will consider A_2 .

$$\text{left node} = \{x_1, x_2\} \quad \text{Right node} = \{x_3, x_4, x_5\}$$

left node

Information gain A_1

$$\text{group1} = \{\text{null}\} \quad \text{group2} = \{x_1, x_2\}$$

$$E(\text{group1}) = 0$$

$$E(\text{group2}) = -\sum_{i=1}^n p_i \log_2 p_i = -(1 \log_2 1) = 0$$

$$E(\text{Split}) =$$

$$I = 0.97$$

Information gain A_2

$$\text{group1} = \{x_1, x_2\} \quad \text{group2} = \{\text{null}\}$$

$$E(\text{group1}) = -(1 \log_2 1) = 0$$

$$E(\text{group2}) = 0$$

$$E(\text{Split}) = \frac{2}{2} \times 0 + 0 \times 0 = 0$$

$$I = 0.97$$

Information gain A_3

$$\text{group1} = \{x_1\} \quad \text{group2} = \{x_2\}$$

$$\begin{aligned} E(\text{group1}) &= -(1 \log_2 1) = 0 \\ E(\text{group2}) &= -(1 \log_2 1) = 0 \\ E(\text{Split}) &= 0 \\ I &= 0.97 \end{aligned}$$

Best split for left node = A_3 .

Right node

Information gain A_1

$$\text{group1} = \{x_3\} \quad \text{group2} = \{x_4, x_5\}$$

$$\begin{aligned} E(\text{group1}) &= -(1 \log_2 1) = 0 \\ E(\text{group2}) &= -(2 \log_2 1) = 0 \\ E(\text{Split}) &= 0 \\ I &= 0.97 \end{aligned}$$

Information gain for A_2

$$\text{group1} = \{\text{null}\} \quad \text{group2} = \{x_3, x_4, x_5\}$$

$$E(\text{group1}) = 0$$

$$E(\text{group2}) = -\left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3}\right)$$

$$= 0.528 + 0.389$$

$$= 0.917$$

$$E(\text{Split}) = 0 + 1 \times 0.917$$

$$= 0.917$$

$$I = 0.083$$

Information gain A_3
 $\text{group1} = \{x_3, x_5\}$ $\text{group2} = \{x_4\}$

$$E(\text{group1}) = -\left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}\right)$$

$$E(\text{group2}) = -\left(1 \log_2 1\right) = 0$$

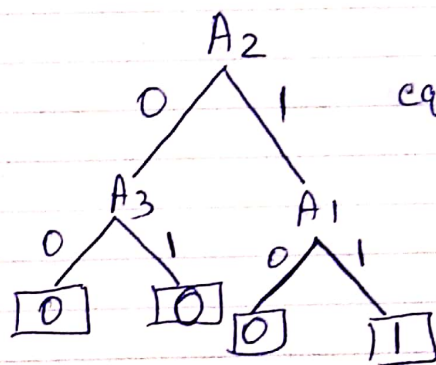
$$E(\text{Split}) = \frac{2}{3} \times 1 + \frac{1}{3} \times 0$$

$$= 0.6667$$

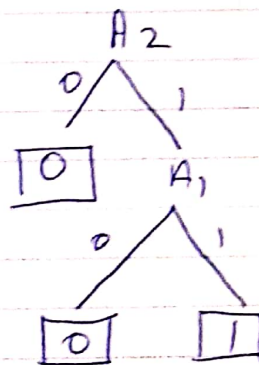
$$I = 0.34$$

Best split for right node = A_1

Tree



equivalent to



2b)

$$x_4 = (A_1 = 1, A_2 = ?, A_3 = 1)$$

Attribute A_1

$$\text{Group 1} = x_3$$

$$\text{Group 2} = x_1, x_4, x_2, x_5$$

$$E(\text{group 1}) = 0$$

$$E(\text{group 2}) = \frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}$$

$$= \frac{1}{2} + \frac{1}{2}$$

$$E(\text{split}) = \frac{1}{5} \times 0 + \frac{4}{5} \times 1$$

$$= \frac{4}{5}$$

$$I = 0.97 - \frac{4}{5} = 0.17$$

Attribute A_2

Case 1

$$x_4 = 1 \ 0 \ 1$$

$$\text{group 1} = x_1, x_2, x_4$$

$$\text{group 2} = x_3, x_5$$

$$E(\text{group 1}) = \left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3} \right) = 0.917$$

$$E(\text{group 2}) = - \left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right) = 1$$

$$E(\text{split}) = \frac{3}{5} \times 0.917 + \frac{2}{5} \times 1$$

$$I = 0.9502$$

$$I = 0.0198$$

Case 2

$x_4 = 1 \ 1 \ 1$

group 1 - x_1, x_2

group 2 - x_3, x_4, x_5

$$E(\text{group 1}) = 0$$

$$E(\text{group 2}) = -\left(\frac{1}{3} \log_2 \frac{1}{3} + \frac{2}{3} \log_2 \frac{2}{3}\right) \\ = 0.917$$

$$E(\text{Split}) = \frac{2}{5} \times 0 + \frac{3}{5} \times 0.917$$

$$= 0.5502 \quad I = 0.97 - 0.5502 = 0.4198$$

Attribute A_3

group 1 - x_1, x_3, x_5

group 2 - x_2, x_4

$$E(\text{group 1}) = -\left(\frac{2}{3} \log_2 \frac{2}{3} + \frac{1}{3} \log_2 \frac{1}{3}\right) \\ = 0.917$$

$$E(\text{group 2}) = -\left(\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2}\right) \\ = 1$$

$$E(\text{Split}) = \frac{3}{5} \times 0.917 + \frac{2}{5} \times 1$$

$$= 0.9502$$

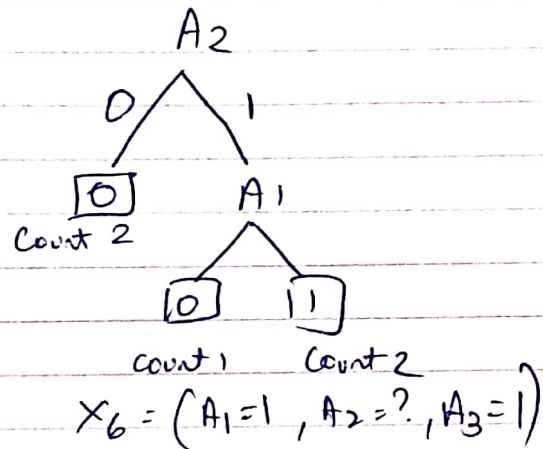
$$I = 0.97 - 0.9502$$

$$= 0.019$$

Since the information gain of A_2 (weighted) is 0.2198 is the highest among all 3 A_2 can be considered as the root.

2a

Decision Tree



$X_6 = 101$ goes into the left subtree with output label 0

$X_6 = 111$ goes into the right subtree with output label 1

The weighted avg is $w_0 \times p_0 + w_1 \times p_1$

$$w_0 = \frac{2}{4} \quad w_1 = \frac{2}{4} \quad p_0 = 0 \quad p_1 = 1$$

$$= \frac{1}{2} \times 0 + \frac{1}{2} \times 1 = 0.5$$

$$A_2 = \begin{cases} 0 & \text{weighted avg} < 0.5 \\ 1 & \text{weighted avg} \geq 0.5 \end{cases}$$

Since weighted avg = 0.5 $A_2 = 1$

$$X_6 = A_1 = 1, A_2 = 1, A_3 = 1$$

The prediction will be 1.