Sudhamshu B N

IDS 18EC091

7th B

# Introduction to Machine Learning

Apply the best classification technique to determine the highest information among the given attributes.

| Instances | A₁ | A₂ | Target class |
|-----------|----|----|--------------|
| 1 | T | T | + |
| 2 | T | T | + |
| 3 | T | F | — |
| 4 | F | F | + |
| 5 | F | T | — |
| 6 | F | T | — |
| 7 | F | F | — |
| 8 | T | F | + |
| 9 | F | T | — |

$\text{Entropy}(S) = -P_\oplus \log_2 P_\oplus - P_\ominus \log_2 P_\ominus$

$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{Values}(A)} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$

$\text{Entropy}(S) = $ +ve → 4

$\qquad\qquad\quad$ −ve → 5

$\text{Entropy}([4+, 5-]) = -\left(\frac{4}{9}\right) \log_2 \left(\frac{4}{9}\right) - \left(\frac{5}{9}\right) \log_2 \left(\frac{5}{9}\right)$

$\qquad\qquad\qquad = 0.51996 + 0.4711$

$\qquad\qquad\qquad = \underline{0.99107}$

Apply

Entropy [S(A1)]

$$S_T = -\frac{3}{4} \log_2 \left(\frac{3}{4}\right) - \frac{1}{4} \log_2 \left(\frac{1}{4}\right)$$

$$= 0.31127 + 0.5$$

$$= 0.81127$$

$$S_F = -\frac{1}{5} \log_2 \left[\frac{1}{5}\right] - \frac{4}{5} \log_2 \left[\frac{4}{5}\right]$$

$$= 0.4643 + 0.2575$$

$$= 0.7218$$

Gain $(S, A_1) = $ Entropy (s) $- \left\{ \frac{4}{9} \text{ Entropy} (S_T) + \frac{5}{9} \text{ Entropy} (S_F) \right\}$

$$= 0.99107 - \left\{ \frac{4}{9} * 0.81127 + \frac{5}{9} * 0.7218 \right\}$$

$$= 0.2294$$

Entropy [S(A2)]

$$S_T = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right)$$

$$= 0.9709$$

$$S_F = -\frac{2}{4} \log_2 \left(\frac{2}{4}\right) - \frac{2}{4} \log_2 \left(\frac{2}{4}\right)$$

$$= 1$$

Gain $(S_1, A_2) = $ Entropy (S) $- \left\{ \frac{5}{9} \text{ Entropy} (S_T) + \frac{4}{9} * \text{ Entropy} (S_F) \right\}$

$$= 0.99107 - \left\{ \frac{5}{9} * 0.97095 + \frac{4}{9} \times 1 \right\}$$

$$= 0.99107 - 0.98386$$

$$= 0.00721$$

$$\text{Gain}(S, A_1) = 0.2294$$

$$\text{Gain}(S, A_2) = 0.00721$$

$A_1$ has got more information gain. So that will be the root node.