

Question :- 1 Decision - tree for weather forecasting

Total :- 14

Sunny : 8

Rainy 6

⇒ Before splitting:

$$\text{Probability } P(S) = \frac{8}{14}$$

$$P(R) = \frac{6}{14}$$

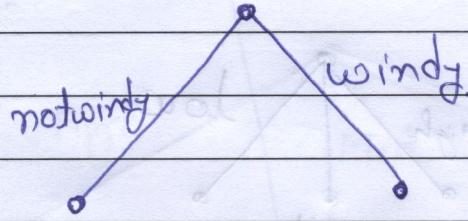
$$\text{Entropy: } H(x) = - \sum_{k=1}^K P(x=x_k) \log_2 P(x=x_k)$$

$$= - \left[\frac{8}{14} \log_2 \left(\frac{8}{14} \right) + \frac{6}{14} \log_2 \left(\frac{6}{14} \right) \right]$$

$$= - \left[\frac{8}{14} (-0.2430) + \frac{6}{14} (-0.36798) \right]$$

$$= 0.2965 \text{ bit}$$

\Rightarrow After splitting
 \Rightarrow windy :-



$\text{sunny} = 6$

$\text{Rainy} = 2$

$\text{sunny} = 2$

$\text{Rainy} = 4$

$$\text{windy} = H(x) = - \left[\frac{2}{6} \log_2 \left(\frac{2}{6} \right) + \left(\frac{4}{6} \right) \log_2 \left(\frac{4}{6} \right) \right]$$

$$= 0.2764 \text{ bit}$$

$$\text{Notwindy} = p(x) = - \left[\frac{6}{8} \log_2 \left(\frac{6}{8} \right) + \left(\frac{2}{8} \right) \log_2 \left(\frac{2}{8} \right) \right]$$

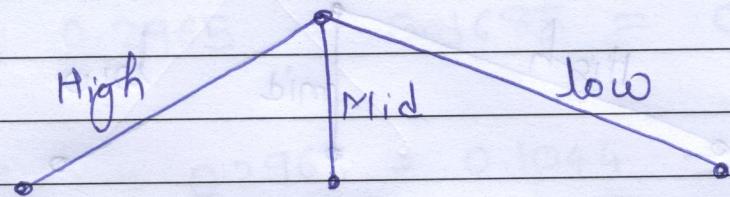
$$= 0.0843 \text{ bit}$$

Weighted Average

$$= E(H(x)) = \left[\frac{6}{14} (0.2764) + \frac{8}{14} (0.0843) \right]$$

$$= 0.1666 \text{ bit}$$

\Rightarrow Humidity :-

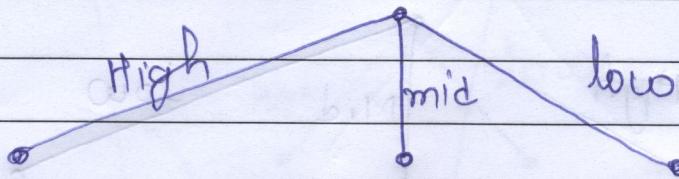


$$\begin{array}{lll} \text{High: } H(x) = - \left[\left(\frac{1}{5} \right) \log_2 \left(\frac{1}{5} \right) + \left(\frac{4}{5} \right) \log_2 \left(\frac{4}{5} \right) \right] \\ \quad \quad \quad = 0.2335 \text{ bit} \\ \text{Mid: } H(x) = - \left[\left(\frac{4}{6} \right) \log_2 \left(\frac{4}{6} \right) + \left(\frac{2}{6} \right) \log_2 \left(\frac{2}{6} \right) \right] \\ \quad \quad \quad = 0.2764 \text{ bit} \\ \text{Low: } H(x) = - \left[\left(\frac{3}{3} \right) \log_2 \left(\frac{3}{3} \right) + \left(\frac{0}{3} \right) \log_2 \left(\frac{0}{3} \right) \right] \\ \quad \quad \quad = 0 \text{ bit} \end{array}$$

$$\begin{array}{l} \text{Weighted Average} = \left[\left(\frac{3}{14} \right) (0.2335) + \left(\frac{6}{14} \right) (0.2764) \right. \\ \quad \quad \quad \left. + \left(\frac{5}{14} \right) (0) \right] \end{array}$$

$$= 0.1685 \text{ bit}$$

\Rightarrow temperature:



Sunny:-5

Rainy:-0

Sunny:3

Rainy:2

Sunny:0

Rainy:4

$$\text{High: } H(x) = - \left[\frac{5}{5} \log_2 \left(\frac{5}{5} \right) + \left(\frac{0}{5} \right) \log_2 \left(\frac{0}{5} \right) \right]$$

$$= 0$$

$$\text{mid: } H(x) = - \left[\frac{3}{5} \log_2 \left(\frac{3}{5} \right) + \left(\frac{2}{5} \right) \log_2 \left(\frac{2}{5} \right) \right]$$

$$= 0.2923$$

$$\text{low: } H(x) = - \left[\frac{0}{4} \log_2 \left(\frac{0}{4} \right) + \left(\frac{4}{4} \right) \log_2 \left(\frac{4}{4} \right) \right]$$

$$= 0$$

$$\text{Weighted Average} = \left[\frac{5}{14}(0) + \frac{5}{14}(0.2923) + \frac{4}{14}(0) \right]$$

$$= 0.1044 \text{ bit.}$$

$$\text{windy} = 0.2965 - 0.1666 = 0.1299 \text{ bit}$$

$$\text{Humidity} = 0.2965 - 0.1685 = 0.128 \text{ bit}$$

$$\text{Temperature} = 0.2965 - 0.1044 = \boxed{0.1921 \text{ bit}}$$

⇒ In this example, splitting the data samples based on the temperature will minimize the entropy (unpredictability) compare to other features

⇒ In other word, splitting based on temperature provide the maximum amount of information gain at this level

⇒ so that best feature to split the data samples at the top of the decision tree is temperature.