

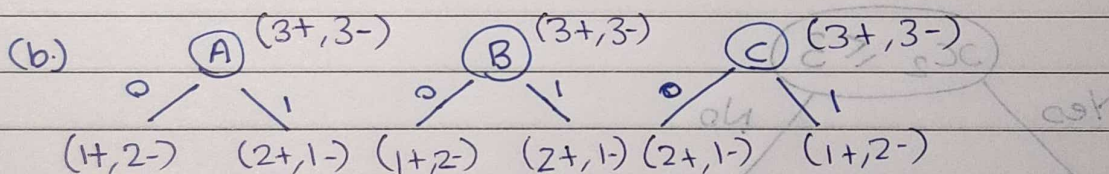
Machine Learning

HOMEWORK-II

1.) The training dataset given is as follows

A	B	C	Y
0	1	0	Yes
1	0	1	Yes
0	0	0	No
1	0	0	No
0	1	1	No
1	1	0	Yes

(a.) The given dataset has instances with same values of A, B, C but different outcomes. Thus, a decision tree with a 100% accuracy on this training set cannot be formed, as no form of tree would satisfy both the instances



$$\begin{aligned} \text{Entropy}(S) &= -(P(y=0) \log_2 P(y=0) + P(y=1) \log_2 P(y=1)) \\ &= -\frac{3}{6} \log_2 \frac{1}{2} - \frac{3}{6} \log_2 \frac{1}{2} \\ &= 1 \end{aligned}$$

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum_{v \in \text{values}(A)} P(v) \text{Entropy}(S_v)$$

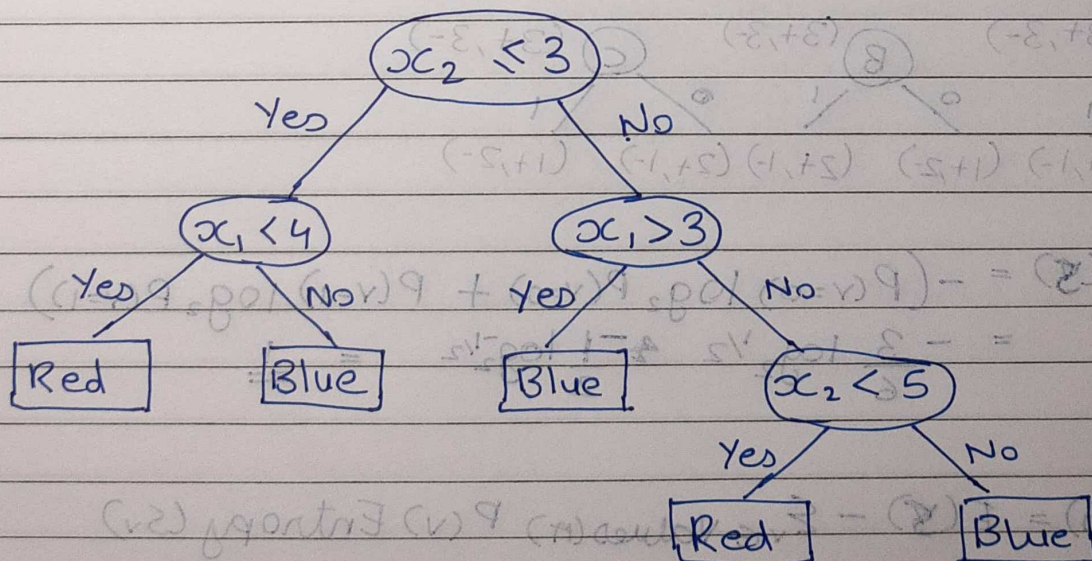
$$\text{Gain}(S, A) = 1 - \left[\frac{3}{6} E(1, 2) + \frac{3}{6} E(2, 1) \right] = 0.545$$

$$\text{Gain}(S, B) = 1 - \left[\frac{3}{6} E(1, 2) + \frac{3}{6} E(2, 1) \right] = 0.545$$

$$\text{Gain}(S, C) = 1 - \left[\frac{3}{6} E(2, 1) + \frac{3}{6} E(1, 2) \right] = 0.545$$

The attributes A, B and C all have equal information gain. Each attribute splits the training dataset in (3+, 3-) into two: (1+, 2-) & (2+, 1-).

2) Decision tree for the given decision boundary is



3.)

	C	D	F
X_2			
15	A	B	E
	5	10	25
	X_1		

4.) Let,

$T_p \rightarrow$ Test results are positive

$D_p \rightarrow$ Patient has disease

$T_n \rightarrow$ Test results are negative

$D_n \rightarrow$ Patient does not have disease

$$P(T_p | D_p) = 0.99$$

$$P(T_n | D_p) = 1 - P(T_p | D_p) = 0.01$$

$$P(D_p) = \frac{1}{10000} = 0.0001$$

$$P(T_n | D_n) = 0.99$$

$$P(T_p | D_n) = 1 - P(T_n | D_n) = 0.01$$

$$P(D_n) = 1 - P(D_p) = 0.9999$$

We need,

$$P(D_p | T_p) = \frac{P(T_p | D_p) \cdot P(D_p)}{P(T_p)} = \frac{(0.99)(0.0001)}{0.010098}$$

$$P(T_p) = P(T_p | D_p) \cdot P(D_p) + P(T_p | D_n) \cdot P(D_n) = (0.99)(0.0001) + (0.01)(0.9999) = 0.010098$$

$$P(D_p | T_p) = 0.00980$$

5.) The mutual information between random variables X & Y with joint probability mass function $P(x, y)$ and marginal probability mass functions $P(x)$ & $P(y)$ is defined as

$$I(x, y) = \sum_{x \in X} \sum_{y \in Y} P(x, y) \log \frac{P(x, y)}{P(x) \cdot P(y)}$$

$$= \sum_{x \in X} \sum_{y \in Y} P(x, y) \log \frac{P(x|y)}{P(x)}$$

$$= - \sum_{x \in X} \sum_{y \in Y} P(x, y) \log P(x) + \sum_{x \in X} \sum_{y \in Y} P(x, y) \log P(x|y)$$

$$= H(x) - H(x|y)$$

OR

$$H(y) - H(y|x)$$

We can infer that X says as much as about Y as Y says about X .

Thus,

$$I(x, y) = H(x) - H(x|y) = H(y) - H(y|x)$$