

MODULE 2

1. Give Decision trees for the following set of training examples

Day	<i>Outlook</i>	<i>Temperature</i>	<i>Humidity</i>	<i>Wind</i>	<i>PlayTennis</i>
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Solution :-

* Entropy(S) = $-P_+ \log_2 P_+ - P_- \log_2 P_-$

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

* Note,

① when all members of S belong to the same class
then, Entropy(S) = 0

② If S contains an equal number of positive and negative examples then,

$$\underline{\text{Entropy}(S) = 1}$$

* The first step is to find the topmost node of the decision tree. ID3 determines the information gain for each attribute, then selects the one with highest information gain.

→ Entropy of S : positive examples = 09
negative examples = 05

$$\text{Entropy}([9+, 5-]) = -\left(\frac{9}{14}\right)\log_2\left(\frac{9}{14}\right) - \left(\frac{5}{14}\right)\log_2\left(\frac{5}{14}\right)$$

$$= \underline{\underline{0.940}}$$

⇒ Information gain of the attribute Outlook is calculated as.

Value(Outlook) = Sunny, Overcast, Rain

$$S_{\text{Sunny}} \leftarrow [2+, 3-]$$

$$S_{\text{Overcast}} \leftarrow [4+, 0]$$

$$S_{\text{Rain}} \leftarrow [3+, 2-]$$

$$\begin{aligned} \text{Gain}(S, \text{Outlook}) &= \text{Entropy}(S) - \left[\left(\frac{5}{14}\right) \text{Entropy}(S_{\text{Sunny}}) \right. \\ &\quad \left. + \left(\frac{4}{14}\right) \text{Entropy}(S_{\text{Overcast}}) + \left(\frac{5}{14}\right) \text{Entropy}(S_{\text{Rain}}) \right] \end{aligned}$$

$$\begin{aligned}
 \textcircled{*} \text{ Entropy}(S_{\text{sunny}}) &= -\left(\frac{2}{5}\right)\log_2\left(\frac{2}{5}\right) - \left(\frac{3}{5}\right)\log_2\left(\frac{3}{5}\right) \\
 &= -(0.4 * (-1.3219)) - (0.6 * (-0.7369)) \\
 &= 0.52876 + 0.44214 \\
 &= 0.970
 \end{aligned}$$

\textcircled{*} Entropy(S_{\text{overcast}}) = 0 (because all members belong to same class)

$$\begin{aligned}
 \textcircled{*} \text{ Entropy}(S_{\text{rain}}) &= -\left(\frac{3}{5}\right)\log_2\left(\frac{3}{5}\right) - \left(\frac{2}{5}\right)\log_2\left(\frac{2}{5}\right) \\
 &= -(0.6 * (-0.7369)) - (0.4 * (-1.3219)) \\
 &= 0.44214 + 0.52876 \\
 &= 0.970
 \end{aligned}$$

$$\begin{aligned}
 &= (0.940) - \left[\left(\frac{5}{14}\right) * 0.9709 + 0 + \left(\frac{5}{14}\right) 0.9709 \right] \\
 &= 0.9409 - [0.3467 + 0.3467] \\
 &= \underline{\underline{0.246}}
 \end{aligned}$$

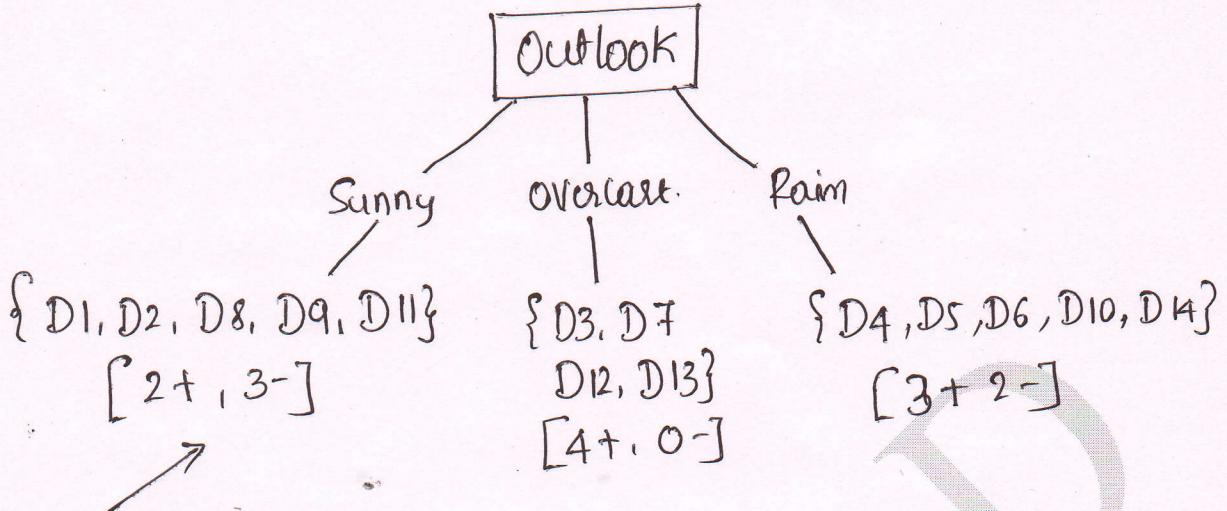
So Gain(S, outlook) = 0.246 ✓

Similarly Gain(S, Temperature) = 0.029

Gain(S, Humidity) = 0.151

Gain(S, wind) = 0.029

So root node will be outlook.



which attribute should be tested here?



$$S_{\text{Sunny}} = \{D_1, D_2, D_8, D_9, D_{11}\}$$

$$\begin{aligned} \text{Gain}(\text{Sunny, Humidity}) &= 0.970 - \left[\left(\frac{3}{5} \right) * 0.0 + \left(\frac{2}{5} \right) * 0.0 \right] \\ &= \underline{\underline{0.970}}. \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{Sunny, Temperature}) &= 0.970 - \left[\left(\frac{2}{5} \right) * 0 + \right. \\ &\quad \left. \left(\frac{2}{5} \right) * 1 + \left(\frac{1}{5} \right) * 0 \right] \\ &= 0.970 - 0.4 \\ &= \underline{\underline{0.570}} \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{Sunny, Wind}) &= 0.970 - \left[\left(\frac{3}{5} \right) * (0.918) + \left(\frac{2}{5} \right) * 1 \right] \\ &= 0.970 - 0.9508 \\ &= \underline{\underline{0.0192}} \end{aligned}$$

So, Attribute Humidity will be descendant node.

$$\Rightarrow S_{\text{Rain}} = \{D_4, D_5, D_6, D_{10}, D_{14}\}$$

$$\begin{aligned} \text{Entropy}(S_{\text{Rain}}) &= -\left(\frac{3}{5}\right)\log_2\left(\frac{3}{5}\right) - \left(\frac{2}{5}\right)\log_2\left(\frac{2}{5}\right) \\ &= \underline{0.970} \end{aligned}$$

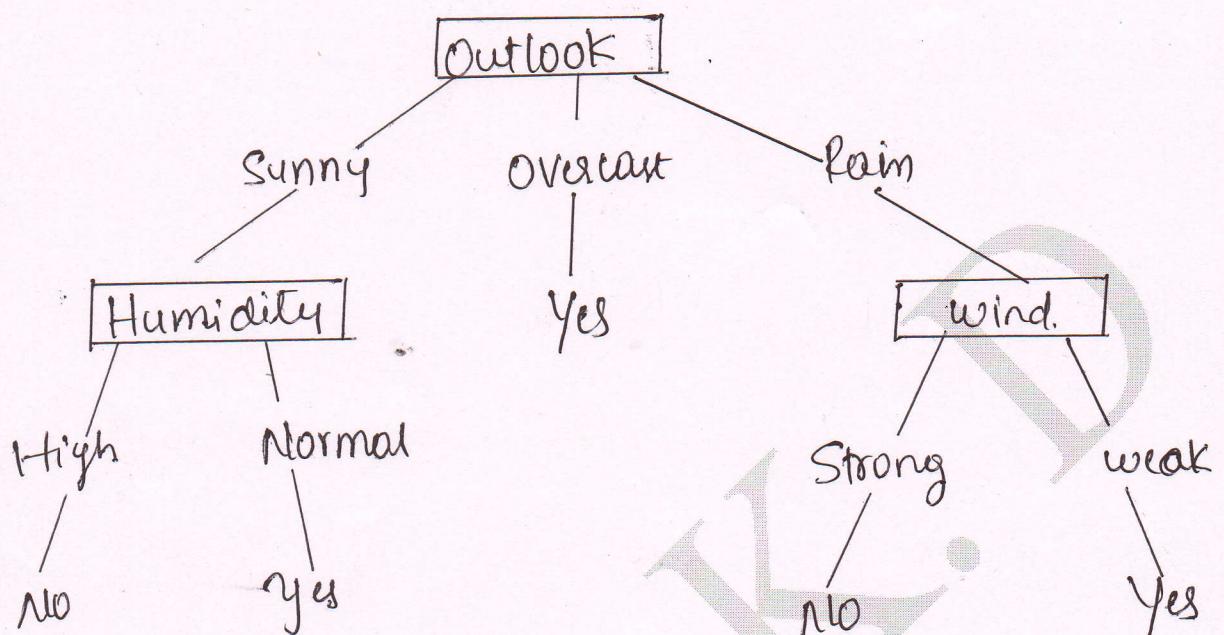
$$\begin{aligned} \text{Gain}(\text{Rain, Temperature}) &= 0.970 - \left[\left(\frac{0}{5}\right)*0 + \left(\frac{3}{5}\right)*0.918 \right. \\ &\quad \left. + \left(\frac{2}{5}\right)*1 \right] \\ &= \underline{\underline{0.0198}} \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{Rain, wind}) &= 0.970 - \left[\left(\frac{3}{5}\right)*0 + \left(\frac{2}{5}\right)*0 \right] \\ &= \underline{\underline{0.970}} \end{aligned}$$

$$\begin{aligned} \text{Gain}(\text{Rain, Humidity}) &= 0.970 - \left[\left(\frac{2}{5}\right)*1 + \left(\frac{3}{5}\right)*0.917 \right] \\ &= \underline{\underline{0.0198}} \end{aligned}$$

So, highest information gain is ~~node~~ attribute
the kind.

So, the final tree is.



2. Consider the following set of training examples.

- a) What is the entropy of this collection of training example with respect to the target function classification?
- b) What is the information gain of a_2 relative to these training examples?

Instance	Classification	a_1	a_2
1	+	T	T
2	+	T	T
3	-	T	F
4	+	F	F
5	-	F	T
6	-	F	T



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

Information gain

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

a)

positive instances = 03

negative instances = 03

$$\text{Entropy}(S) = \left(\frac{3}{6} \right) \log_2 \left(\frac{3}{6} \right) + \left(\frac{3}{6} \right) \log_2 \left(\frac{3}{6} \right)$$

$\text{Entropy}(S) = 1$

when there are equal number of +ve & -ve instances, then $\text{Entropy}(S) = 1$.

$$b) \text{Gain}(S, a_2) = \text{Entropy}(S) - \left[\frac{4}{6} \text{Entropy}(S_T) + \frac{2}{6} \text{Entropy}(S_F) \right]$$

find:

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

$$= \underline{\underline{1}}$$

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

$$= \underline{\underline{1}}$$

$$\Rightarrow \text{Gain}(S, a_2) = 1 - \left[\frac{4}{6}(1) + \frac{2}{6}(1) \right]$$

$$= 1 - 1$$

$$\text{Gain}(S, a_2) = \underline{\underline{0}}$$

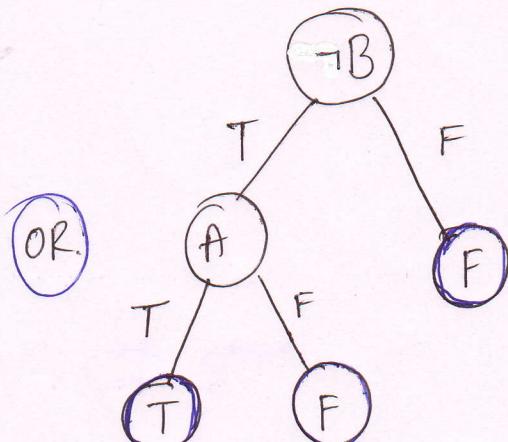
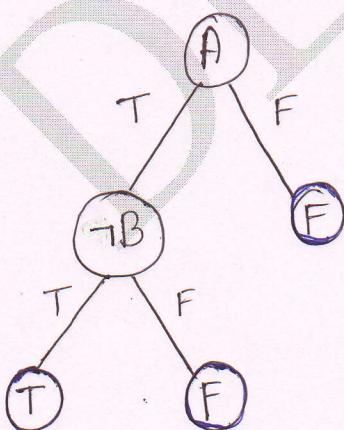
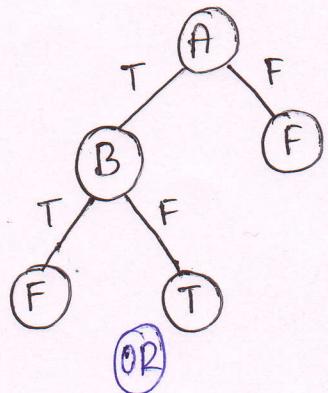
3. Give decision trees to represent the following Boolean functions.

- i) $A \& \neg B$
- ii) $A \vee [B \& C]$
- iii) $A \oplus B$
- iv) $[A \& B] \vee [C \& D]$

Solution

- i) $A \& \neg B$

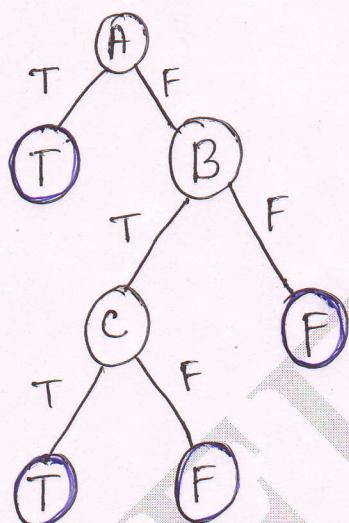
A	B	$\neg B$	$A \& \neg B$
T	T	F	F (-)
T	F	T	T (+)
F	T	F	F (-)
F	F	T	F (-)



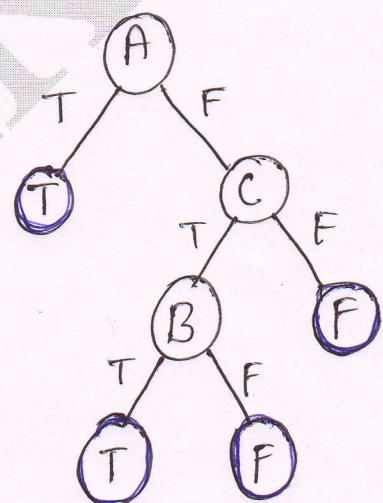
ii) $A \vee [B \wedge C]$

A	B	C	$B \wedge C$	$A \vee [B \wedge C]$
T	T	T	T	T
T	T	F	F	T
T	F	T	F	T
F	T	T	T	T
F	T	F	F	F
F	F	T	F	F
F	F	F	F	F

D
A
?

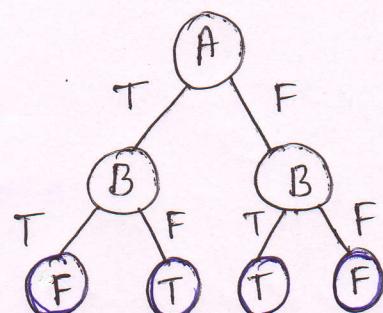


OR.



iii) $A \oplus B$

A	B	$A \oplus B$
T	T	F
T	F	T
F	T	T
F	F	F



iv) $[A \And B] \Or [C \And D]$

A	B	C	D	A $\And\And$ B	C $\And\And$ D	(A $\And\And$ B) \Or (C $\And\And$ D)
T	T	T	T	T	T	T
T	T	T	F	T	F	T
T	T	F	T	T	F	T
T	T	F	F	T	F	T
T	F	T	T	F	T	T
T	F	T	F	F	F	F
T	F	F	T	F	F	F
F	T	T	T	F	T	T
F	T	T	F	F	F	F
F	T	F	F	F	F	F
F	F	T	T	F	T	T
F	F	T	F	F	F	F
F	F	F	T	F	F	F
F	F	F	F	F	F	F

