



शिक्षा मंत्रालय  
MINISTRY OF  
EDUCATION

INDIAN INSTITUTE OF TECHNOLOGY  
JODHPUR



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING



P M R F

Prime Minister's Research Fellowship

Week 2 - Live Session

# Data Mining

**Swapnil S. Mane**

*mane.1@iitj.ac.in*

PMRF Research Scholar

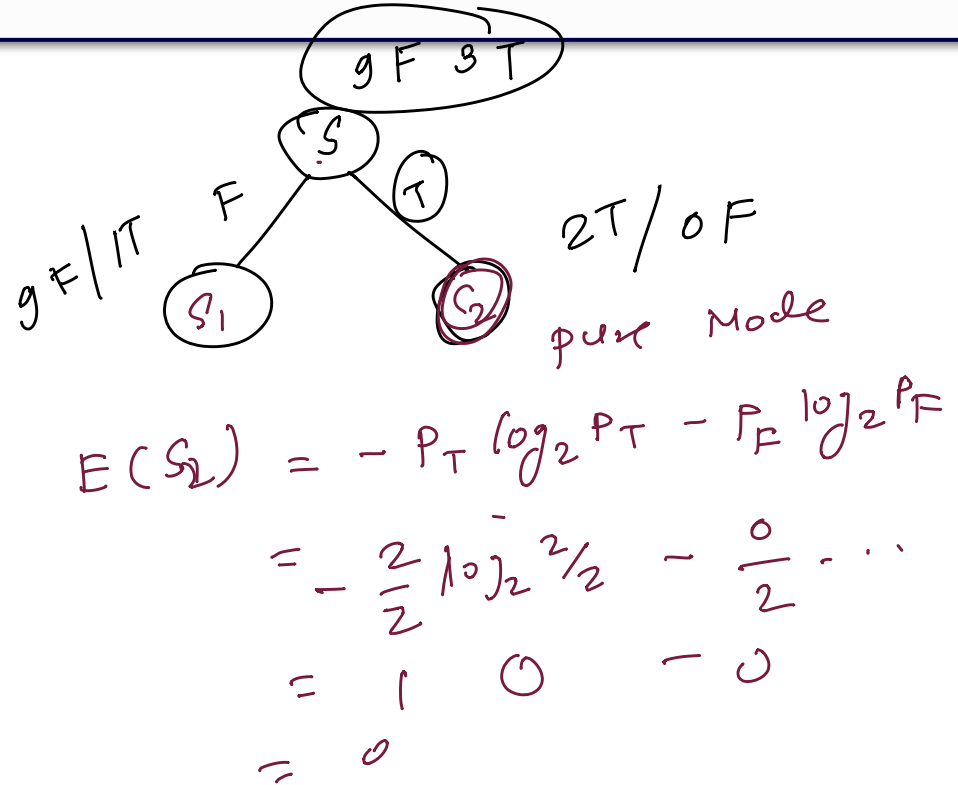
# Week 2

Q1. A decision tree can be used to build models for:

- a) Regression problems
- b) Classification problems
- ☒ c) Both of the above
- d) None of the above

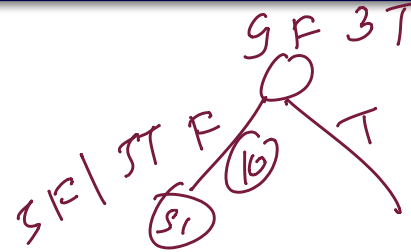
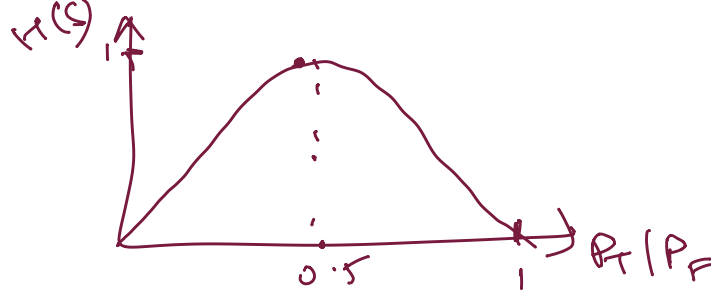
Q2. Entropy value of \_\_\_\_ represents that the data sample is pure or homogenous:

- a) 1
- ☒ b) 0
- c) 0.5
- d) None of the above.



Q3. Entropy value of \_\_\_\_\_ represents that the data sample has a 50-50 split belonging to two categories:

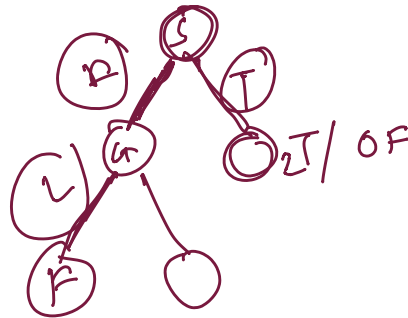
- ☒ a) 1
- ☐ b) 0
- ☐ c) 0.5
- ☐ d) None of the above



$$\begin{aligned}
 E(S_1) &= -P_T \log_2 P_T - P_F \log_2 P_F \\
 &= -\frac{5}{10} \log_2 \frac{5}{10} - \frac{5}{10} \log_2 \frac{5}{10} \\
 &= -0.5 * -1 - 0.5 * -1 \\
 &= 1
 \end{aligned}$$

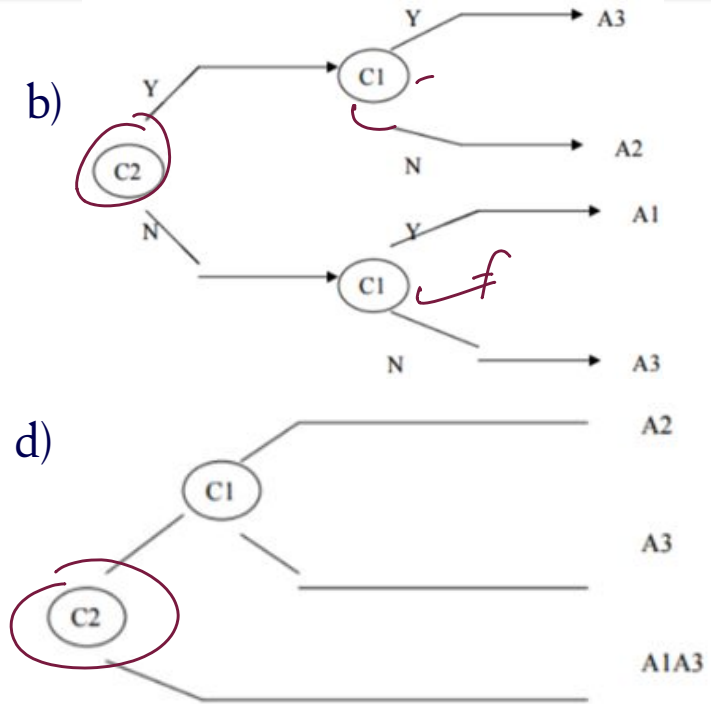
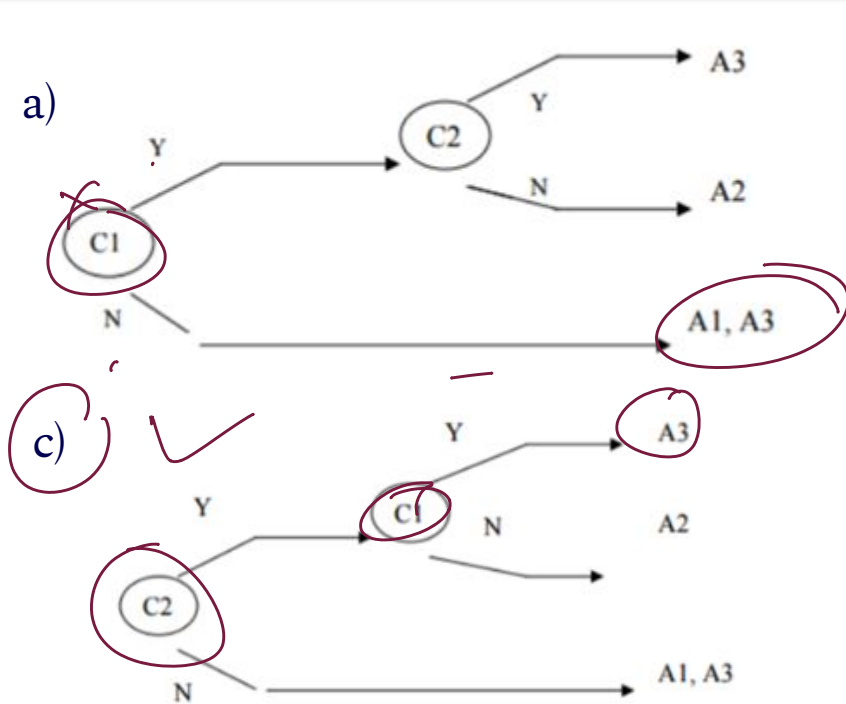
Q4. If a decision tree is expressed as a set of logical rules, then:

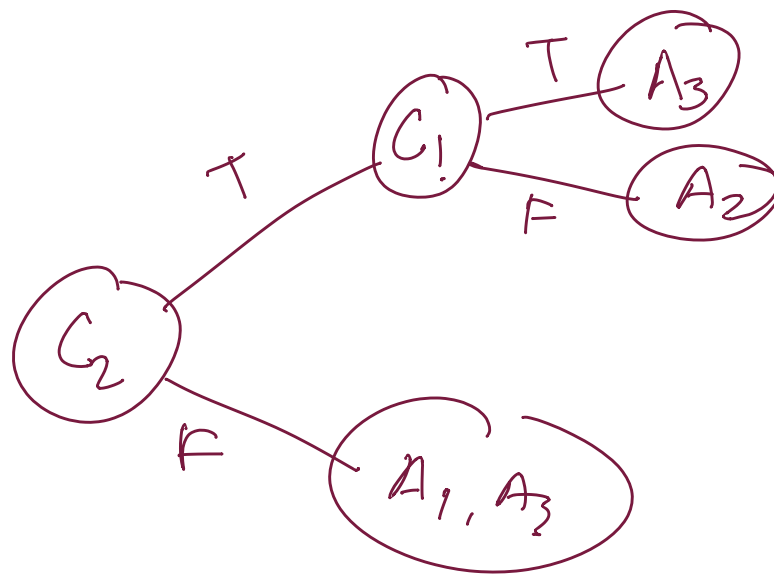
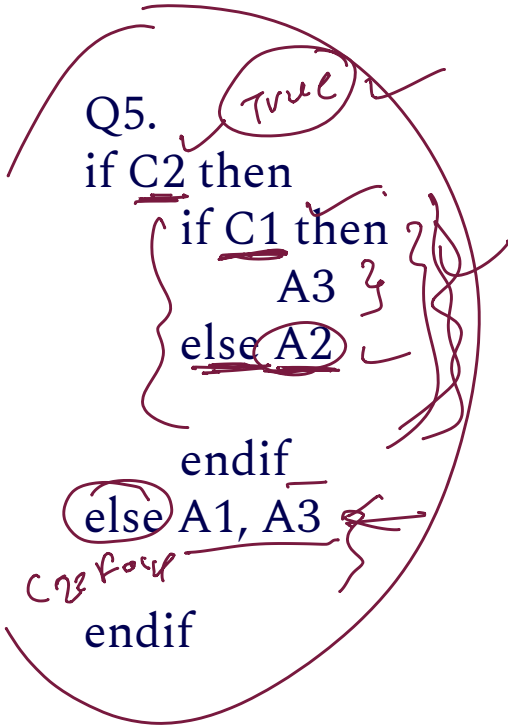
- a) The internal nodes in a branch are connected by AND and the branches by AND
- b) The internal nodes in a branch are connected by OR and the branches by OR
- c) The internal nodes in a branch are connected by AND and the branches by OR
- d) The internal nodes in a branch are connected by OR and the branches by AND



if  $S = F$  OR  $T = 2$   
then  $F$   
end if  
else then  $T$

Q5. The Decision tree corresponding to the following is? (1 Mark) if C2 then if C1 then A3 else A2 endif else A1, A3 endif







Q6. What is the entropy of the dataset?

a) 0.50

☒ b) 0.92

c) 1

d) 0

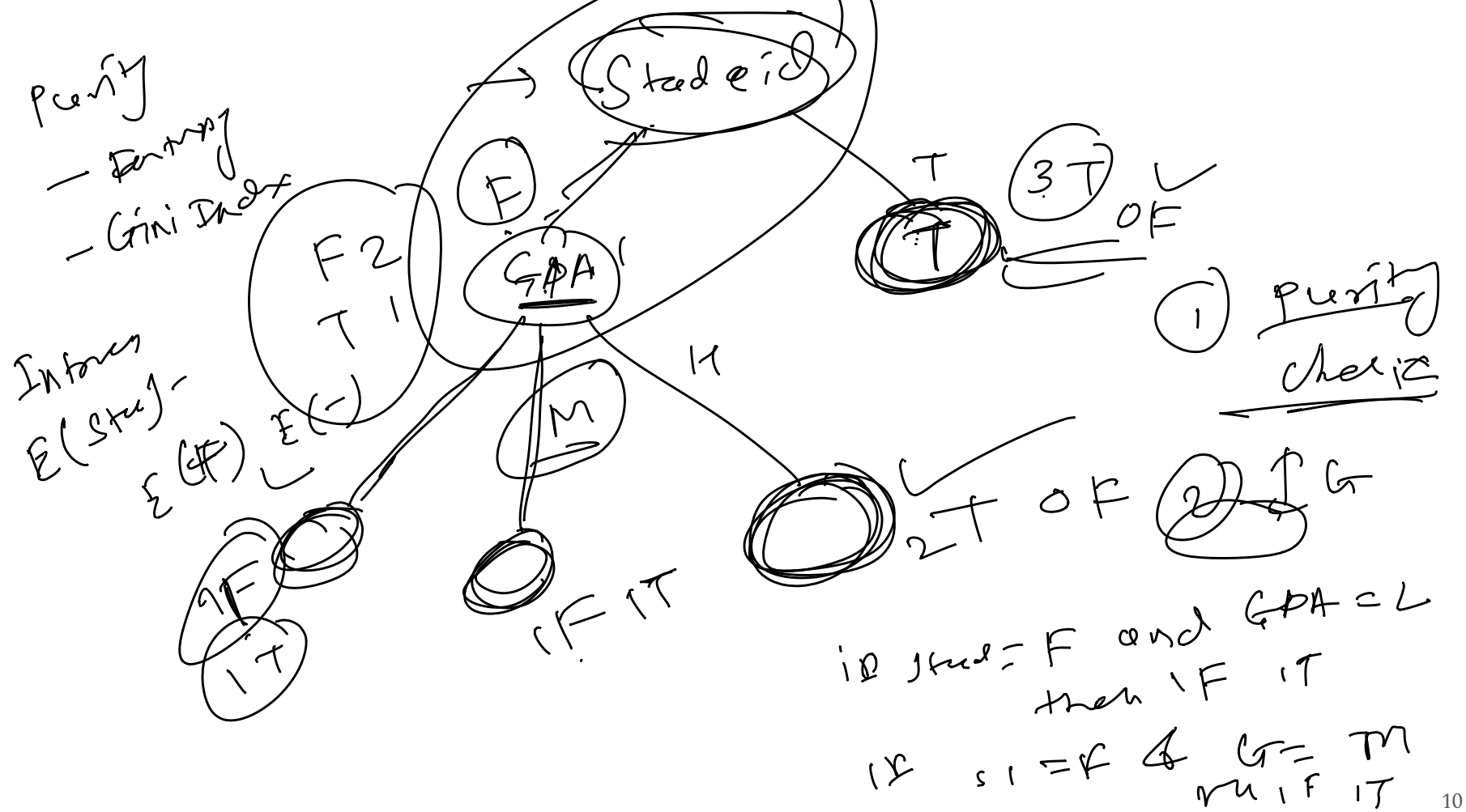
$$T = 4 \\ F = 2$$

$$E(S) = -P_T \log_2 P_T - P_F \log_2 P_F$$

$$E(\text{Passed}) = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6}$$

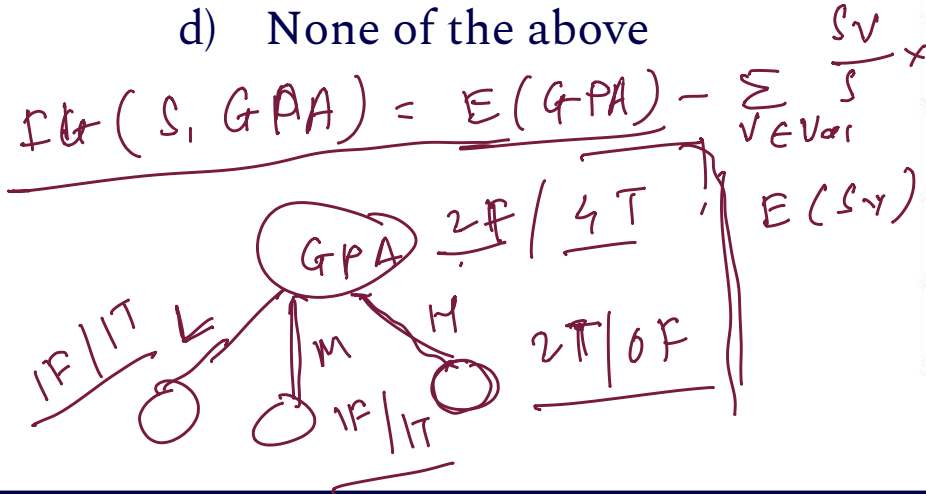
$$= 0.9182 \approx 0.92$$

GPA	Studied	Passed
Low	F	F
Low	T	T
Medium	F	F
Medium	T	T
High	F	T
High	T	T



Q7. Which attribute would information gain choose as the root of the tree?

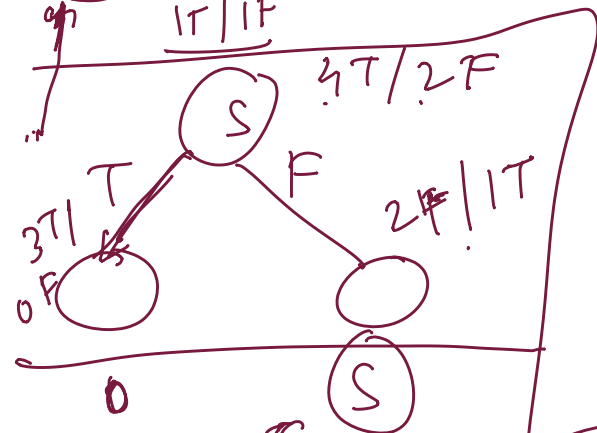
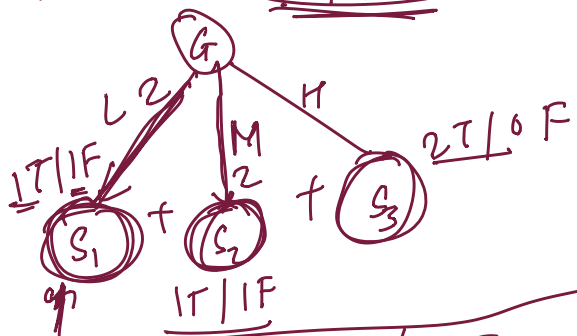
- a) GPA
- ☒ b) Studied
- c) Passed
- d) None of the above



GPA	Studied	Passed
Low	F	F
Low	T	T
Medium	F	F
Medium	T	T
High	F	T
High	T	T

GPA

4T/2F



$$IG(S, GPA) = E(GPA) - \sum_{v \in \text{val}(S)} \frac{1 \cdot \text{val}(S)}{n} \cdot E(\text{val}(S))$$

$$\textcircled{1} E(GPA) = -P_T \log_2 P_T - P_F \log_2 P_F$$

$$= -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = \textcircled{0.92}$$

$$\textcircled{2} E(S_1) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$\textcircled{3} E(S_2) = 1$$

$$\textcircled{4} E(S_3) = 0$$

$$IG(GPA) = 0.92 - \left[ \frac{2}{6} \cdot 1 + \frac{2}{6} \cdot 1 + 0 \right]$$

$$= 0.92 - \left[ \frac{1}{3} + \frac{1}{3} \right]$$

$$IG(GPA) = 0.253,$$

$$= -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3}$$

$$= 0.811$$

$$IG(s) = 0.92 - [(3/2 * 0) + 3/2 * 0.811]$$

$$= 0.5145$$

---


$$IG(s, GA) = 0.2583$$

$$IG(s, s) = 0.545 \quad \checkmark$$