

06-NOV-2022

Decision Tree

Mathematical Definition

- { ① Decision Tree classifier
- ② Decision Tree Regressor

Classification Problem

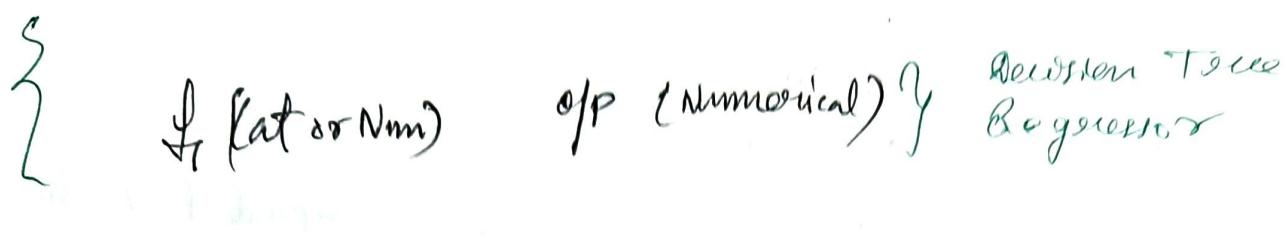
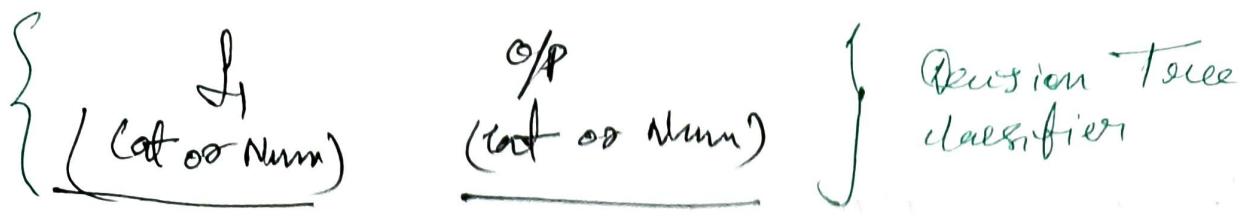
<u>Dataset</u>	f_1 <u>Outlook</u>	f_2 <u>Temperature</u>	f_3 <u>Humidity</u>	f_4 <u>Wind</u>	\downarrow <u>Decision</u>
Bay					
1	Sunny	hot	high	weak	No
2	Sunny	hot	high	strong	No
3	Overcast	hot	high	weak	Yes
4	Rainfall	mild	high	weak	Yes
5	Rainfall	cool	normal	strong	No
6	Rainfall	cool	normal	strong	Yes
7	overcast	cool	normal	weak	No
8	sunny	mild	high	weak	Yes
9	sunny	cool	normal	weak	Yes
10	Rainfall	mild	normal	weak	Yes
11	Sunny	mild	high normal	Strong	Yes
12	overcast	mild	high	strong	Yes
13	overcast	hot	normal	weak	Yes
14	Rainfall	mild	high	strong	No

Dependent Features

Independent Features

- > Here f_1, f_2, f_3, f_4 are categorical features.
- > Based on the output feature we decide whether the problem is about classification or regression.

High wind, Rainfall or not



① ID₃

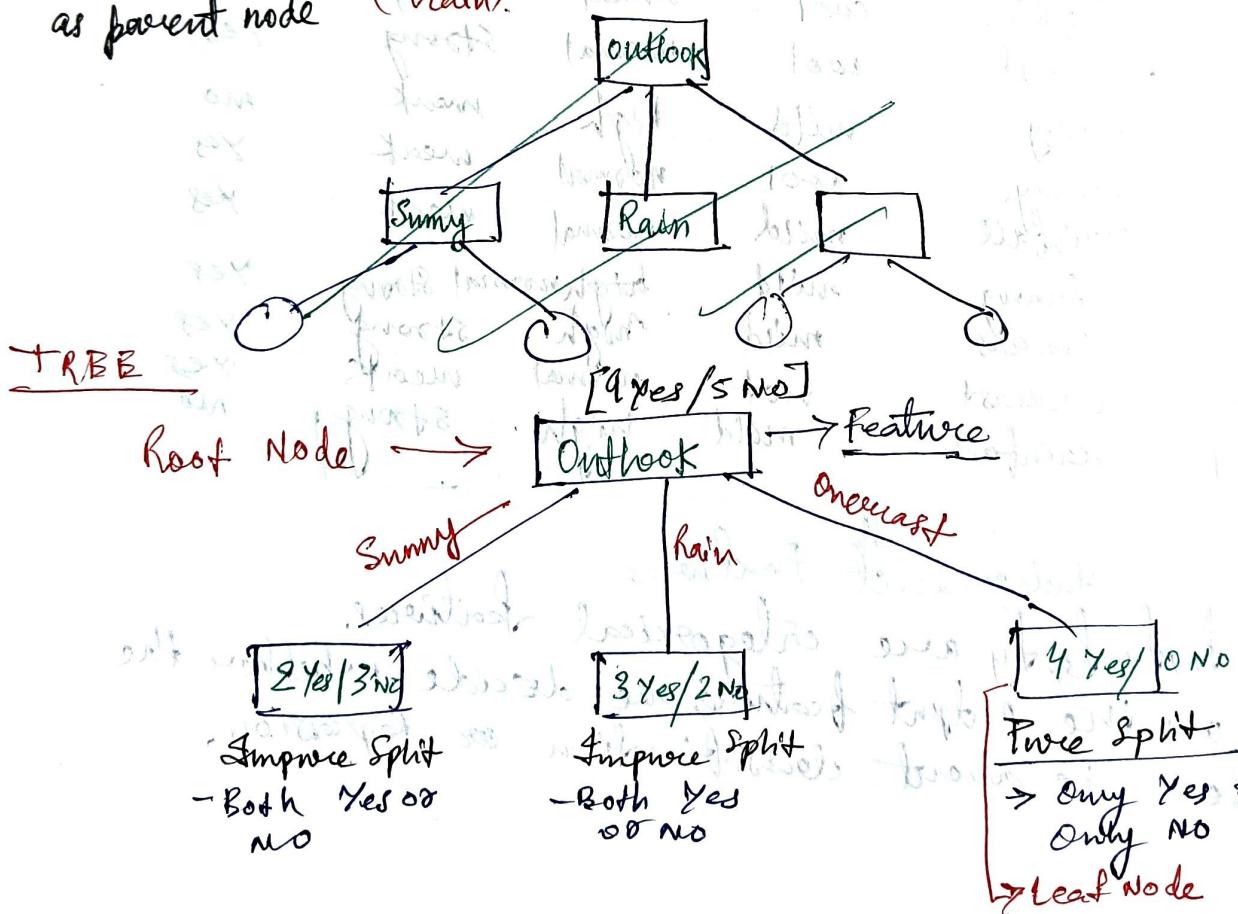
- [Iterative Decomposition]
- Deal with Entropy

② CART

- Classification and Regression Trees.
- Deal with Gini Impurity

→ Choose Outlook as parent node

(We can decide based on entropy or information gain).



→ Try to find out pure split.

Check Purity using below methods

Two methods:

① Entropy

② Gini Impurity or Gini coefficients

Purity

→ How much information one feature is providing.

① Entropy

- Randomness
- we can check purity of the feature using Entropy

$$\text{Entropy} = - \sum_{i=1}^n p_i * \log(p_i)$$

② Gini Impurity

- Check Purity using G.I.

$$G.I. = 1 - \sum_{i=1}^n p_i^2$$

For 2 class

$$\text{Entropy} = -P_Y \log_2(P_Y) - P_N \log_2(P_N)$$

$Y \rightarrow \text{Yes}$ } 2 classes
 $N \rightarrow \text{No}$ }
→ log base 2

For 3 class $\rightarrow (C_1, C_2, C_3)$

$$\text{Entropy} = -P_{C_1} \log_2(C_1) - P_{C_2} \log_2(C_2) - P_{C_3} \log_2(C_3)$$

For 2 class (Yes and No)

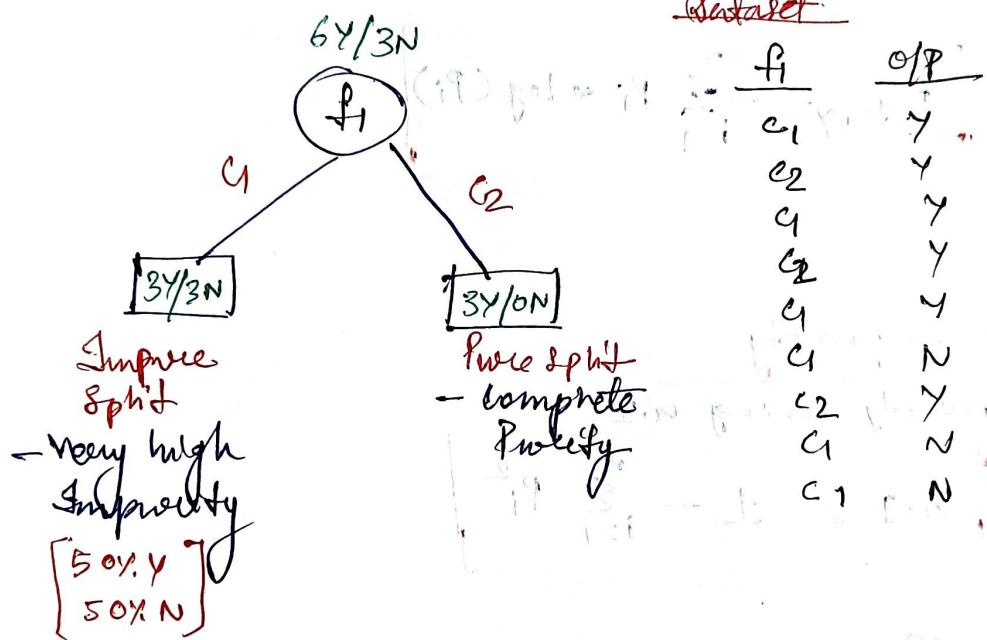
$$H.I = 1 - [P_Y^2 + P_N^2]$$

For 3 class (C1, C2, C3)

$$H.I = 1 - (P_{C_1}^2 + P_{C_2}^2 + P_{C_3}^2)$$

Q why are we calculating the H.I and entropy?
A To check the impurity and randomness of the feature.

Example



Checking impurity for f_1

① Entropy (For 3Y/3N split)

$$(P_Y = \frac{3}{6}, P_N = \frac{3}{6})$$

$$\begin{aligned}
 \text{Entropy}(H(S)) &= - \sum_{i=1}^n P_i \log_2(P_i) \\
 &= - P_Y \log_2(P_Y) - P_N \log_2(P_N) \\
 &= - \frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right) \\
 &= - \log_2\left(\frac{3}{6}\right) = - \log_2(2^{-1}) \\
 &\quad \text{QD} \qquad \qquad \qquad = 1
 \end{aligned}$$

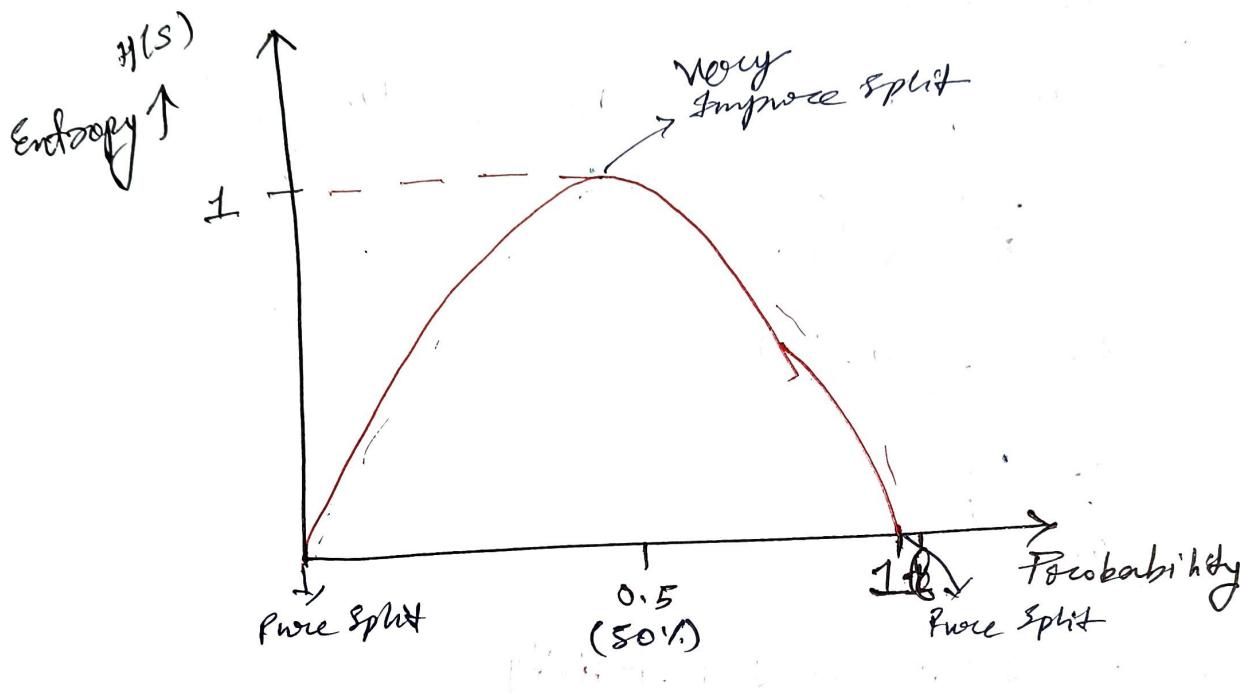
$\left\{ \begin{array}{l} \text{If entropy} = 0 \rightarrow \text{Pure split} \\ \text{If entropy} = 1 \rightarrow \text{Very Improve split} \end{array} \right.$

(For 3Y/0N split)

$$[P_Y = \frac{3}{3} = 1, P_N = 0]$$

$$\begin{aligned} \text{Entropy } H(S) &= -P_Y \log_2 P_Y - P_N \log_2 P_N \\ &= -\frac{3}{3} \log_2 \left(\frac{3}{3}\right) - \frac{0}{3} \log_2 \left(\frac{0}{3}\right) \\ &= 0 \end{aligned}$$

Entropy Graph w.r.t Probability



(For 2Y/3N)

$$(P_Y = \frac{2}{5}, P_N = \frac{3}{5})$$

$$\begin{aligned} \text{Entropy } H(S) &= -P_Y \log_2 (P_Y) - P_N \log_2 (P_N) \\ &= -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) \\ &= 0.97 \end{aligned}$$

③ Gini coefficients/Gini Impurity

for $(3Y/3N)$ split $\left\{ P_Y = \frac{2}{3}, P_N = \frac{1}{3} \right\}$

$$H.I = 1 - \sum_{i=1}^n (P_i)^2 = 1 - \left(\left(\frac{2}{3}\right)^2 + \left(\frac{1}{3}\right)^2 \right) \\ = 1 - \left(\frac{4}{9} + \frac{1}{9} \right) \\ = \frac{1}{2} = 0.5$$

If $H(S) = 0$ $(3Y/3N)$
How will you decide?
→ we use $H.I$ in
that case

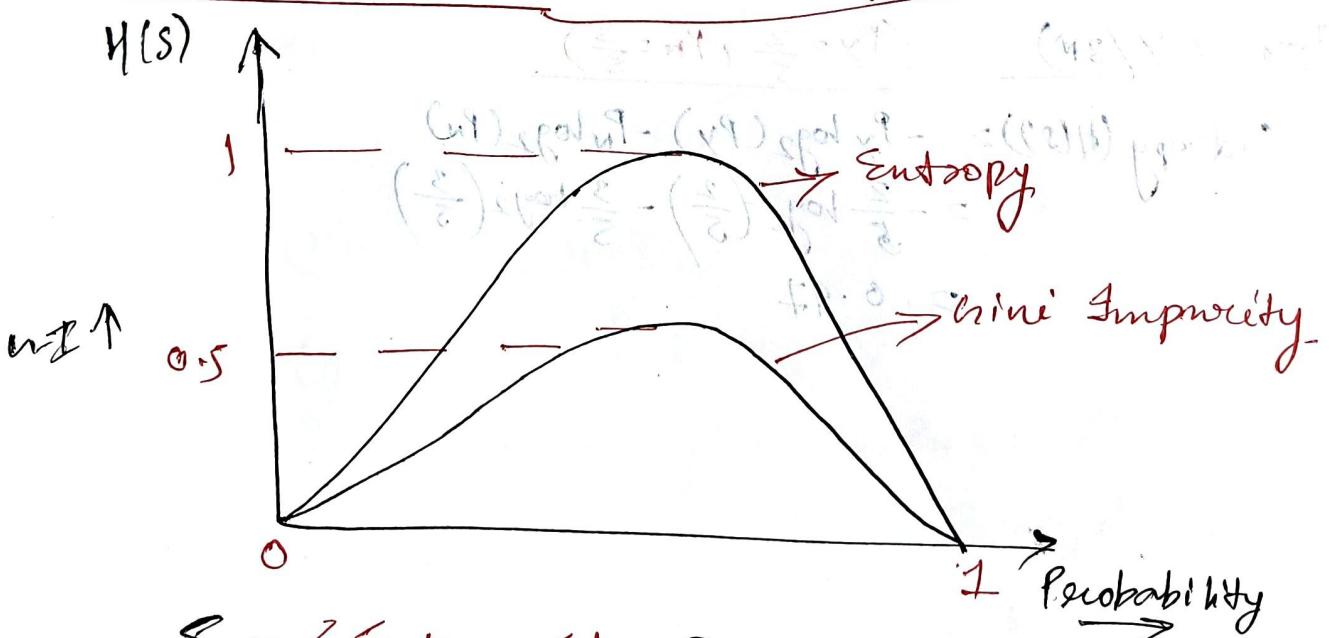
for $(3Y/0N)$ split $\left\{ P_Y = \frac{3}{3}, P_N = 0 \right\}$

$$H.I = 1 - \left(\left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2 \right) \\ = 1 - 1 = 0$$

For $(2Y/3N)$ split $\left\{ P_Y = \frac{2}{5}, P_N = \frac{3}{5} \right\}$

$$H.I = 1 - \left\{ \left(\frac{2}{5}\right)^2 + \left(\frac{3}{5}\right)^2 \right\} \\ = 1 - \left\{ \frac{4}{25} + \frac{9}{25} \right\} \\ = 1 - \frac{13}{25} \\ = \frac{12}{25}$$

Entropy and G.I. w.r.t. worst probability



$$\left\{ \begin{array}{l} 0 \leq \text{Entropy} \leq 1 \\ 0 \leq H.I \leq 0.5 \end{array} \right\}$$

$$\text{For } (4Y/8N \text{ split}) \quad S_2 P_Y = \frac{4}{12}, P_N = \frac{8}{12} \}$$

$$H.I. = 1 - \left\{ \left(\frac{4}{12} \right)^2 + \left(\frac{8}{12} \right)^2 \right\}$$

$$= 1 - \left\{ \frac{16}{144} + \frac{64}{144} \right\} = 1 - \frac{80}{144}$$

$$= 0.44$$

$$\text{For } (8Y/10N \text{ split}) \quad S_2 P_Y = \frac{8}{10}, P_N = \frac{2}{10} \}$$

$$H.I. = 1 - \left\{ \left(\frac{8}{10} \right)^2 + \left(\frac{2}{10} \right)^2 \right\}$$

$$= 1 - \left\{ \frac{64}{100} + \frac{4}{100} \right\}$$

$$= 1 - \frac{68}{100}$$

$$= \frac{32}{100}$$

$$= 0.32$$

③ Information gain

- Either we consider Entropy or H.I. to find information gain.
- In ID3 algorithm we use entropy to calculate information gain.
- In CART algorithm we use H.I. to calculate information gain.

feature1, feature2, feature3

How to select one feature for root node?

→ Find information gain using entropy or H.I.
 ↓
 ID3 ↑
 CART

→ For big dataset we use H.I.

→ Generally in CART binary split happens.

→ Generally in ID3, more than 2 split happens.

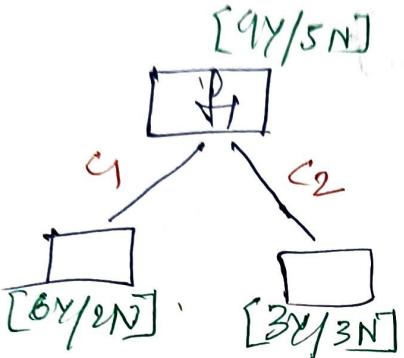
→ Calculations are fast in Gini Impurities.

Information Gain , Root Feature Entropy

$$\text{Gain}(S, F_i) = H(S) - \sum \frac{|S_v|}{|S|} H(S_v) \quad | \text{value}$$

Either we can entropy or Gini Impurities
Sv is each each split, entropy

Example



Q Calculate Information Gain?

Ans For this split, entropy has been calculated in entropy section. for $8Y/2N$ and $3Y/3N$ split.

$$\begin{aligned} H(S) &= \text{Root Feature Entropy} = -P_Y \log_2(P_Y) - P_N \log_2(P_N) \\ &= -(9/14) \log_2(9/14) - (5/14) \log_2(5/14) \end{aligned}$$

$$\text{for } 8Y/2N, H(S)_2 = -\frac{6}{8} \log_2\left(\frac{6}{8}\right) - \frac{2}{8} \log_2\left(\frac{2}{8}\right)$$

$$H(S) = 0.81$$

$$\text{for } 3Y/3N, H(S) = -\frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right)$$

$$H(S) = 0.1$$

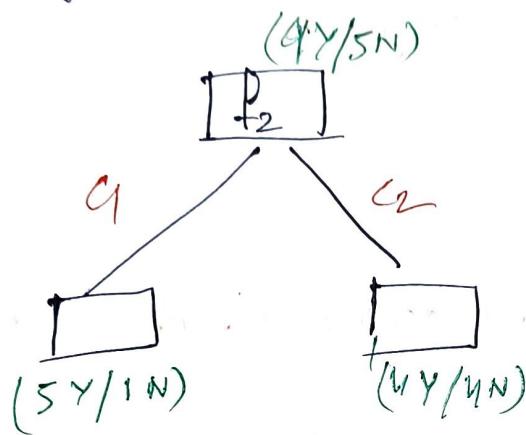
Information gain w.r.t f_1

$$\text{gain}(S, f_1) = 0.94 - \left[\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$\boxed{\text{gain}(S, f_1) = 0.04}$$

Information gain w.r.t f_2

→ Our ultimate goal is to calculate information gain using H.I or entropy.



$$\text{gain}(S, f_2) = 0.94 - \left(\frac{6}{14} \times 0.65 + \frac{8}{14} \times 1 \right)$$

$$\boxed{\text{gain}(S, f_2) = 0.09}$$

NOTE

$$\text{gain}(S, f_2) > \text{gain}(S, f_1)$$

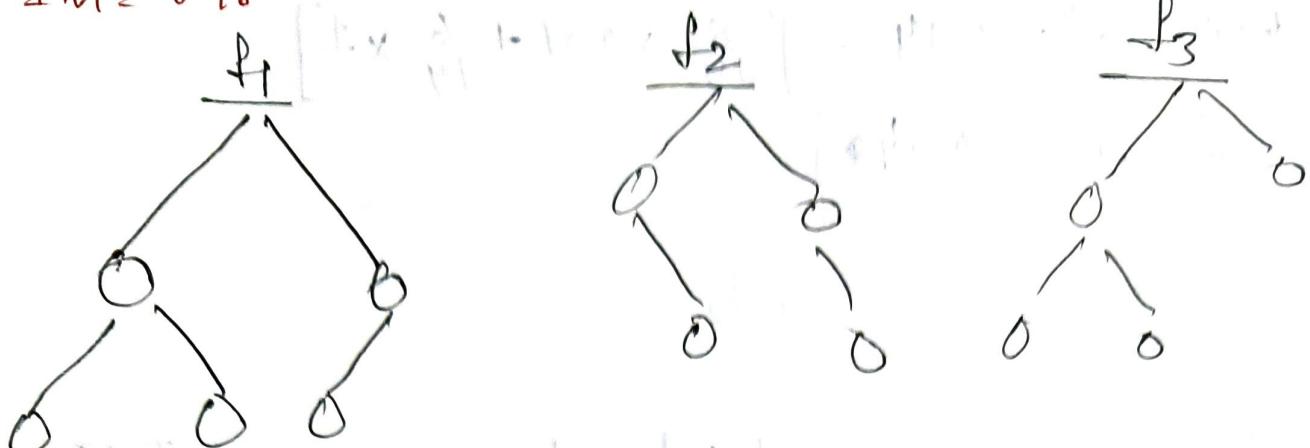
→ we will select f_2 for root node because
it is providing more information.

→ f_2 will be better split

$I(H) = 0.98$

$I(H) = 0.97$

$I(H) = 0.95$



→ Information gain of f_1 is good better than others, so f_1 will be root node.

Example
Example

→ outlook, temp, humidity, wind, decision dataset on 1st page of today's class.

Let's assume outlook has more information gain.

