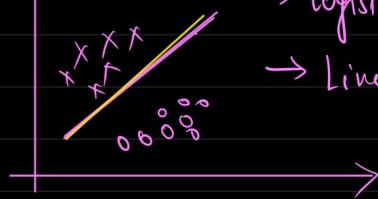
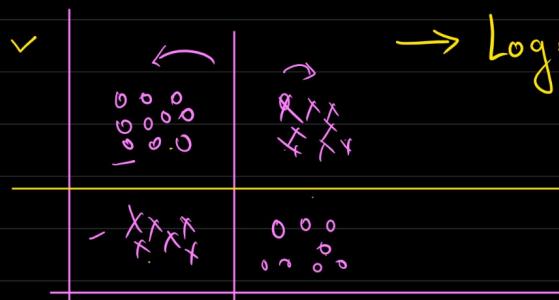


## \* Decision Tree

Log Reg.



→ logistic regression works for linear decision boundary  
→ Linear regression captured linear relationship.



→ Logistic regression | linear regression  
will not work with  
non linear data

## \* Decision tree

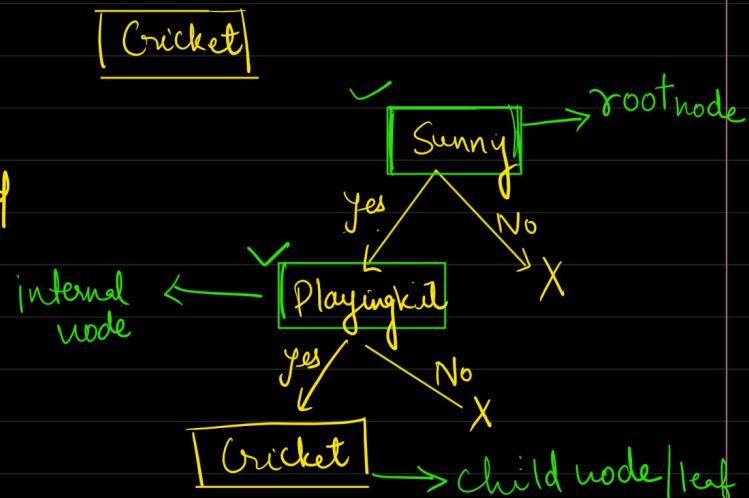
Classification problem (Decision tree classification)  
Regression problem (Decision tree regression)

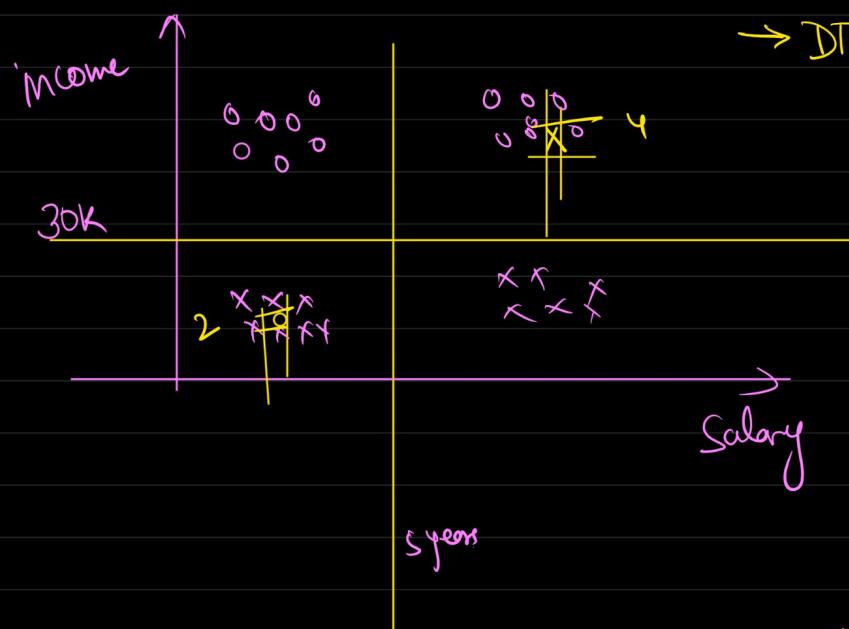
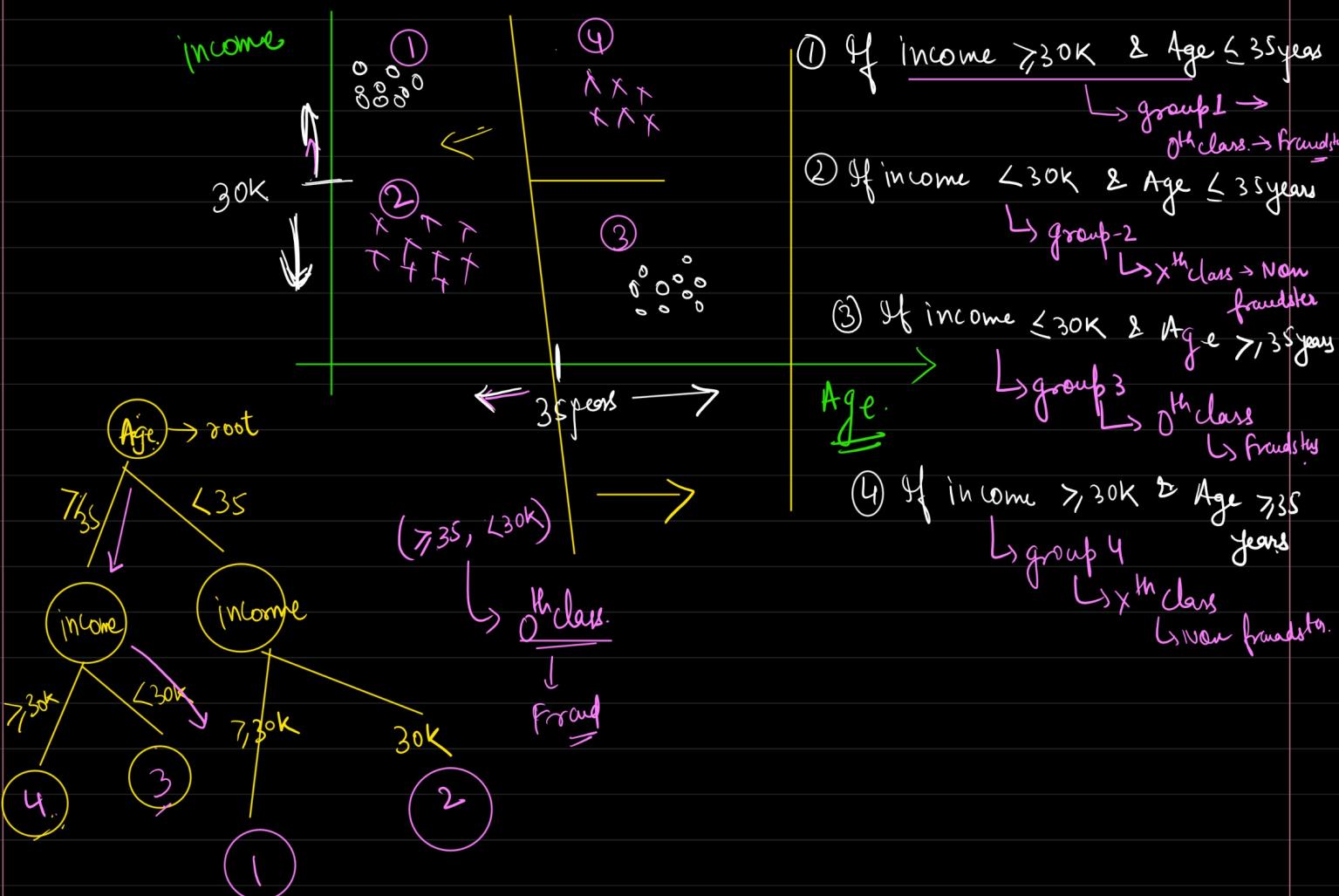
\* Decision Tree works on nested if else condition.

## Decision Tree

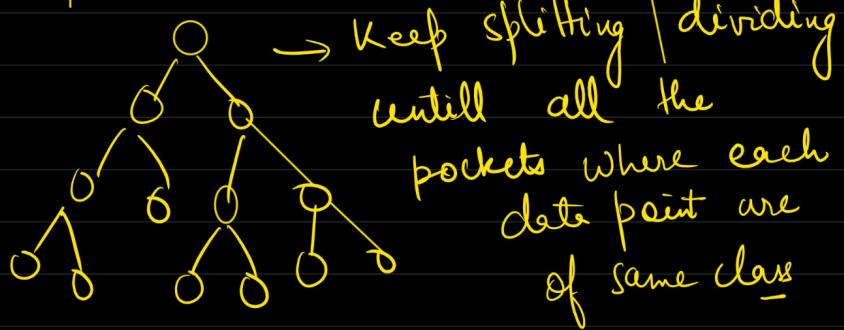
At every node a decision is taken

It looks like a tree



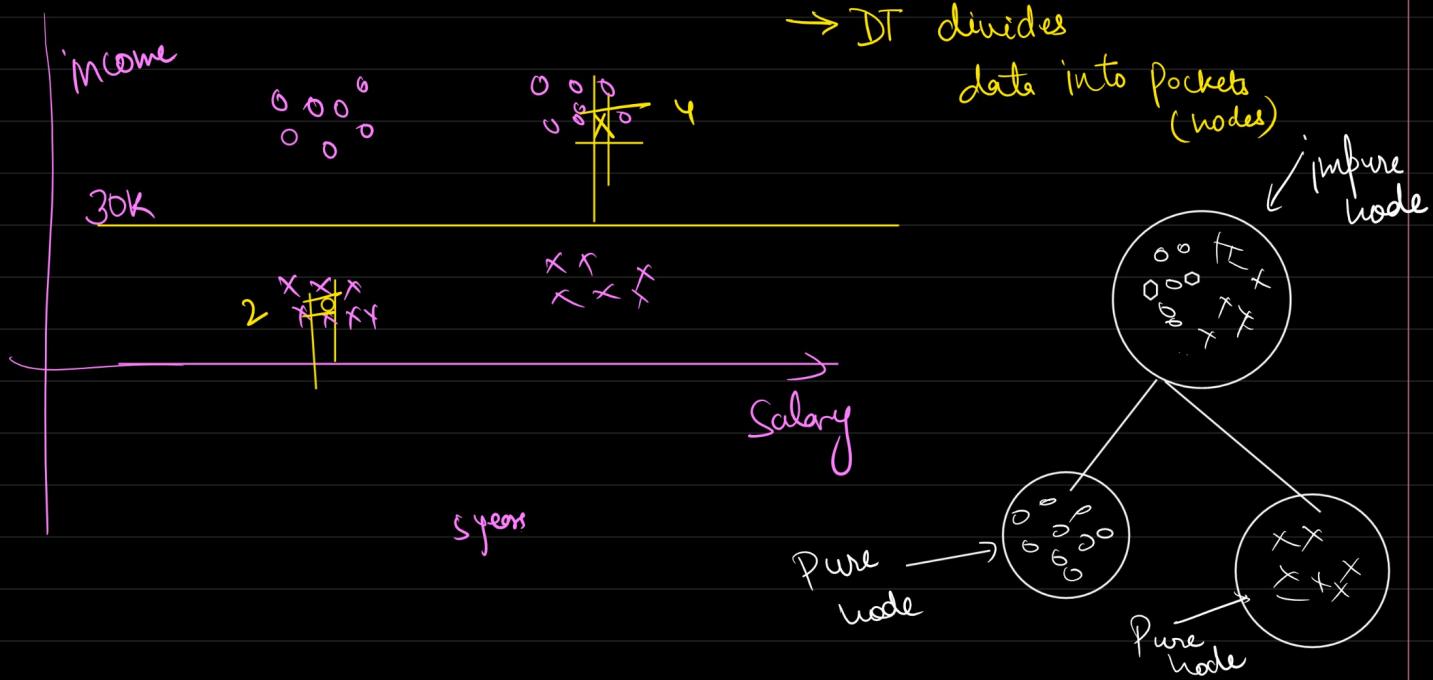


→ DT divides the data into pockets (nodes)

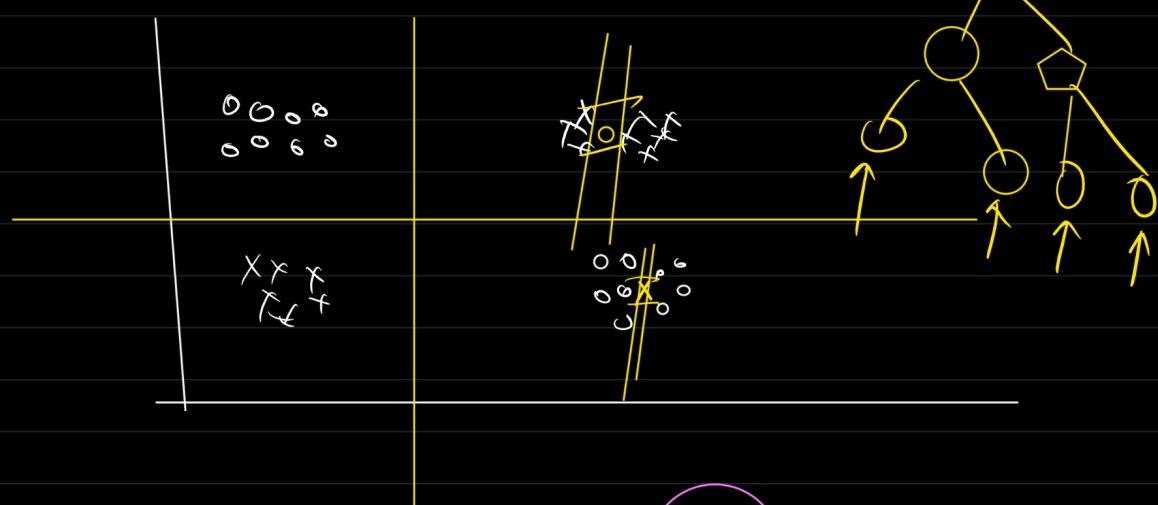


## \* Advantages

- Simple to understand
- Captures non-linear relationship

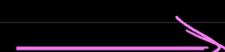


\* The DT intent is to create pure leaf node/pockets where each data point belongs to one class



Depth of DT - 2

Purenode



Some condition

Some cond

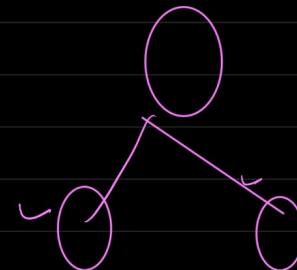
Some condition

Pure node



- ① How this split will happen? (purity / impurity)
- ② Which features to be used for splitting?

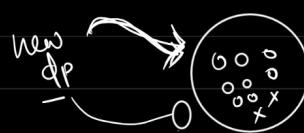
pure - homogeneous node (Same class)  
impure - heterogeneous node



### Impurity Measures

- Classification Error
- Gini
- Entropy.

Salary	Age	Fromed
-	-	0
-	-	1
-	-	0
-	-	1



$$0: 8 \\ X: 2$$

$$\text{Classification Error} = 1 - \max p_i$$

$$P_0 = \frac{8}{8+2} = \frac{8}{10} = 0.8$$

$$C.E = 1 - \max (0.8, 0.2)$$

$$P_X = \frac{2}{8+2} = \frac{2}{10} = 0.2$$

$$= 1 - 0.8 \\ = 0.2$$

\* If I assign everything to majority class then what is the error ratio.

### ② Gini

$$G_I = 1 - \sum_{i=0}^{n-1} (p_i)^2$$

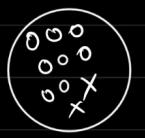
$$(G_I = \sum_{i=0}^{n-1} p_i(1-p_i))$$

$$\text{for two class } p(1-p) + p_X(1-p_X)$$

$$\sum_{i=0}^{n-1} \Rightarrow \text{from class 0 to class } n-1$$

$$= 2p_0(1-p_0) p_X = 1 - p_0^2$$

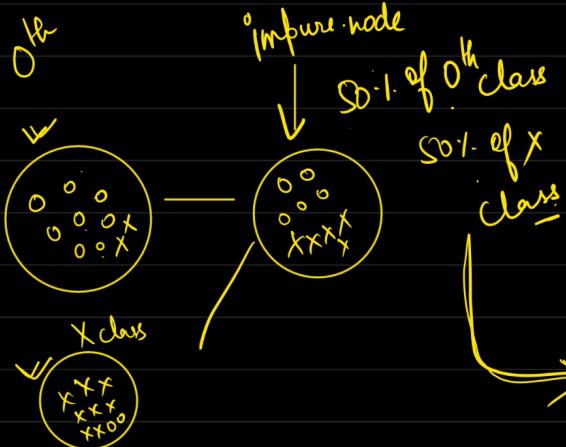
$$= 0.8(1-0.8) + 0.2(1-0.2) \\ = 0.8 \times 0.2 + 0.2 \times 0.8 \\ = 0.32$$



$$G \cdot I = 1 - \sum_{i=0}^{2-1} (p_i)^2$$

$$= 1 - \sum_{i=0}^1 (p_i)^2 = 1 - [(p_0)^2 + (p_X)^2] \\ = 1 - [(0.8)^2 + (0.2)^2]$$

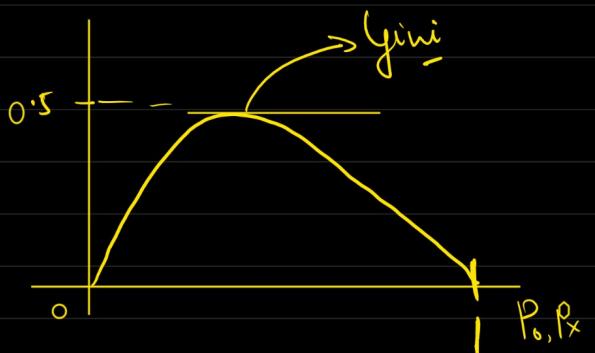
$$= 1 - (0.64 + 0.04) \\ = 1 - 0.68 \Rightarrow 0.32$$



$$G \cdot I = 1 - [(p_0)^2 + (p_X)^2] \\ = 1 - \left(\frac{1}{2}\right)^2 + \left(\frac{1}{2}\right)^2$$

$$= 1 - \left(\frac{1}{4} + \frac{1}{4}\right) = 1 - \frac{2}{4} = \frac{1}{2} = 0.5$$

highest Gini impurity value = 0.5



### ③ Entropy

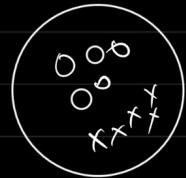
$$H(S) = - \sum_{i=0}^{N-1} p_i \log_2 p_i$$



$$= -p_0 \log_2 p_0 - p_X \log_2 p_X$$

$$= -0.8 \log_2 0.8 - 0.2 \log_2 0.2 \Rightarrow 0.72$$

Highest  $H(S)$   $\rightarrow$   $S01-S01$  of both classes



$$H(S) = -\sum_{i=0}^{N-1} p_i \log_2 p_i$$

$$= -\underbrace{\frac{1}{2} \log_2 \frac{1}{2}}_a - \underbrace{\frac{1}{2} \log_2 \frac{1}{2}}_a$$

$$= -\frac{1}{2} \cancel{\log_2 \frac{1}{2}} \quad \frac{1}{2} = 2^1$$

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

$$= -1 \times -1 = 1$$

$\rightarrow$  highest Entropy value is 1

$S01-S01$   $\rightarrow$  1

Entropy  $\rightarrow$  from 0 to 1

$y_{ini} \rightarrow$  0 to 0.5

misclassification  $\rightarrow$  0 to 0.5

Error



Class 1

impure  $\rightarrow$  50%

40%  $\times$

60%  $\times$

pure  $\rightarrow$  100%  $\times$

Class 2

50%

60%

40%

0%

Entropy

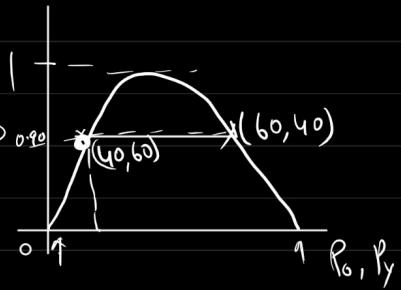
1

0.5

0.90

0.90

0



## Multiclass classification

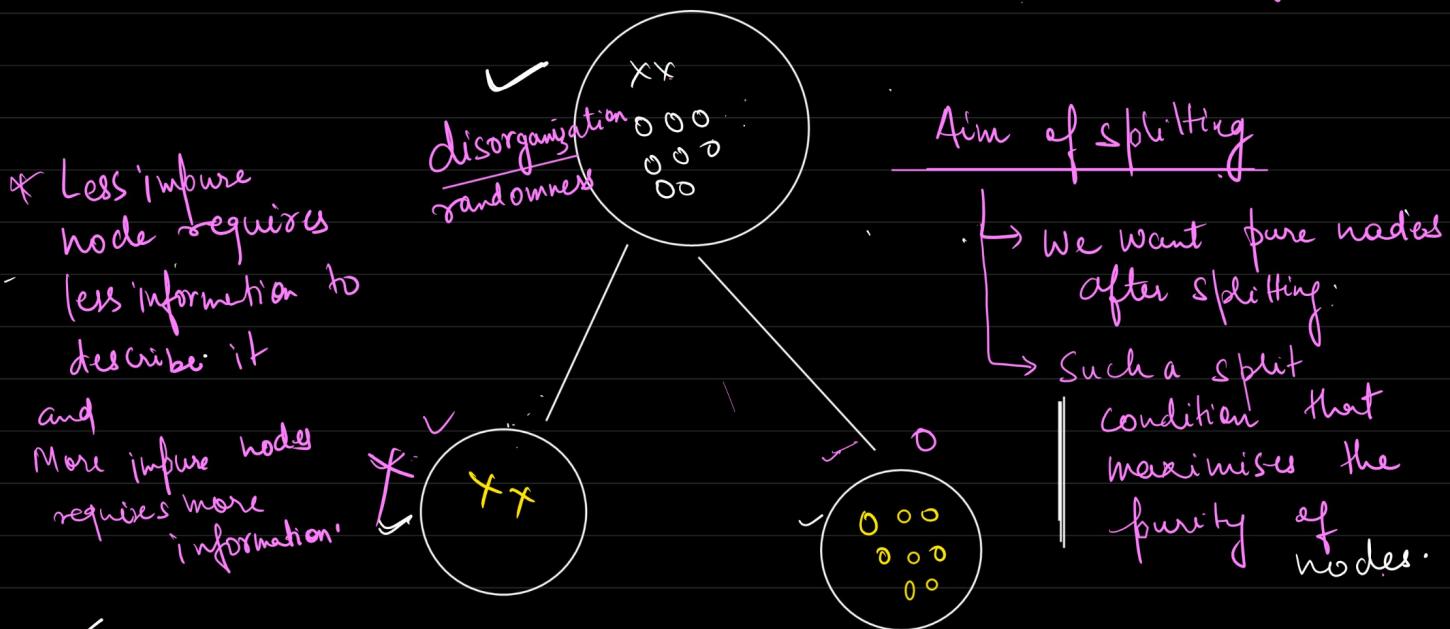
$$\text{Gini Impurity} = 1 - \sum_{i=1}^n p_i^2$$

$$= 1 - [(p_1)^2 + (p_2)^2 + (p_3)^2]$$

$$\text{Entropy}(H(s)) = - \sum_{i=1}^N p_i \log_2 p_i$$

$$= -p_1 \log_2 p_1 - p_2 \log_2 p_2 - p_3 \log_2 p_3$$

Q On what feature to split?  $\Rightarrow$  Information gain



✓  
\*  $\text{impurity}(\text{post split}) < \text{impurity}(\text{pre split})$   
 $\rightarrow$  Split the does this best is chosen as best split

\* This change in impurity is given by the difference of post split impurity and pre split impurity

$$\Delta \text{impurity} = \text{impurity}(\text{pre split}) - \text{impurity}(\text{post split})$$

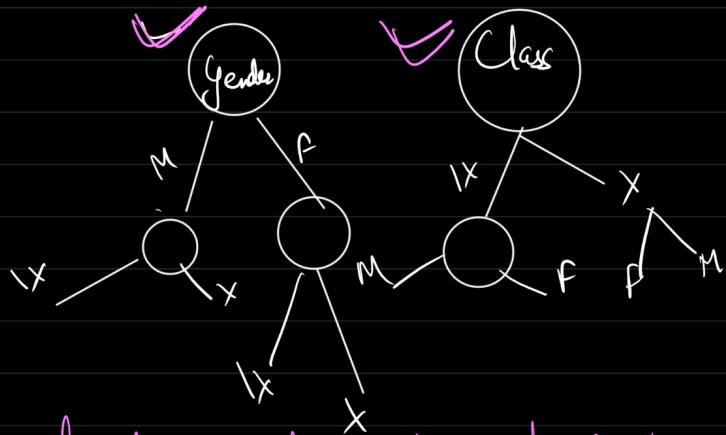
↑  
Information gain

- \* Whichever feature has the information gain, that will be used for splitting.
  - \* If disorganization | randomness decreases, then we say Information gain  $\Delta$  Impurity
  - \* Information theory is a measure to define the degree of disorganization in a system  $\Rightarrow I(G) \leftarrow H$  Entropy

Gender	Class	Cricket (%)
M	IX	Yes
F	X	No
M	X	Yes
-	IX	-
-	-	-
-	-	-
-	-	-
-	-	-
-	-	-
-	-	-

Gender  $\rightarrow 30 = 50\% \text{ play cricket}$  (out of 10 female 2 play cricket)  
 Class  $\rightarrow 10/2 = F \Rightarrow \text{play cricket}$   
 Class  $\rightarrow 20/13 = M \Rightarrow \text{play cricket}$  (out of 20 male, 13 play cricket)

Class  $\rightarrow IX - 14/43\% = \text{out of } 14 \text{ class IX students, } 14\% \text{ play cricket}$   
 Class  $\rightarrow X - 16/57\% = \text{out of } 16 \text{ class X students, } 57\% \text{ play cricket.}$



- \* whichever feature has the highest information gain, that will be used for splitting

$$\begin{aligned} IG &= \text{Impurity split} - \text{Impurity Post split} \\ &= E_n(\text{parent}) - E_n(\text{child combined}) \end{aligned}$$

$$E_n(\text{parent}) = -\frac{15}{30} \log \frac{15}{30} - \frac{15}{30} \log \frac{15}{30} = 1$$

$$IG = En(\text{parent}) - En(\text{child})$$

$$E_n(\text{Parent}) = -\frac{15}{30} \log_2 \frac{15}{30} - \frac{15}{30} \log_2 \frac{15}{30} = 1$$

$$E_{\text{child}}(x) \rightarrow -\frac{43}{100} \log_2 \frac{43}{100} - \frac{57}{100} \log_2 \frac{57}{100} - x_1$$

En child Gender

$$E_{\text{female child}} = -\frac{2}{10} \log \frac{2}{10} - \frac{8}{10} \log \frac{8}{10}$$

$$= 0.72$$

$$E_{\text{male child}} = -\frac{13}{20} \log \frac{13}{20} - \frac{7}{20} \log \frac{7}{20} = 0.92$$

$$E_n(\text{child combined}) = \frac{\text{No of females}}{\text{Total student}} \times E_{\text{female child}} + \frac{\text{No of males}}{\text{Total student}} \times E_{\text{male child}}$$

$$= \frac{10}{30} \times 0.72 + \frac{20}{30} \times 0.92$$

$$= \underline{\underline{0.85}}$$

$$IG_{\text{Gender}} = E_p - E_{\text{cc}}$$

$$= 1 - 0.85 = \underline{\underline{0.15}}$$

$$E_n(x) = -\frac{57}{100} \log_2 \frac{57}{100} - \frac{43}{100} \log_2 \frac{43}{100} - x_2$$

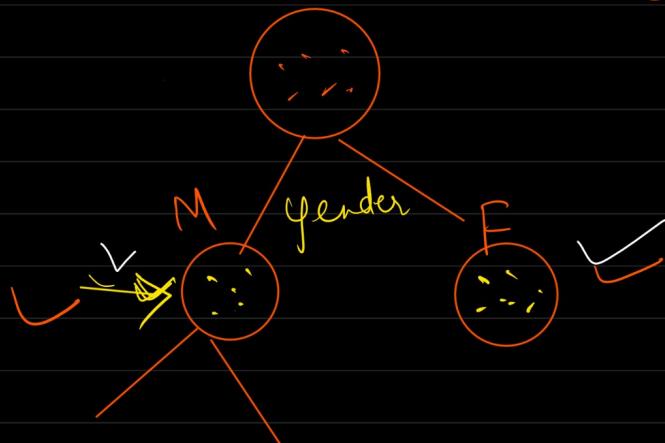
$$E_{\text{combined}} = \frac{14}{30} x_1 + \frac{16}{30} x_2 = 0.99$$

$$IG_{\text{class}} = 1 - 0.99 = 0.01$$

$$= \underline{\underline{0.01}}$$

Since  $IG_{\text{Gender}} > IG_{\text{class}}$

$\therefore$  We will Gender to split.



Algorithm for DT

Step-1 Recursive binary splitting | Partitioning the data into smaller subset.

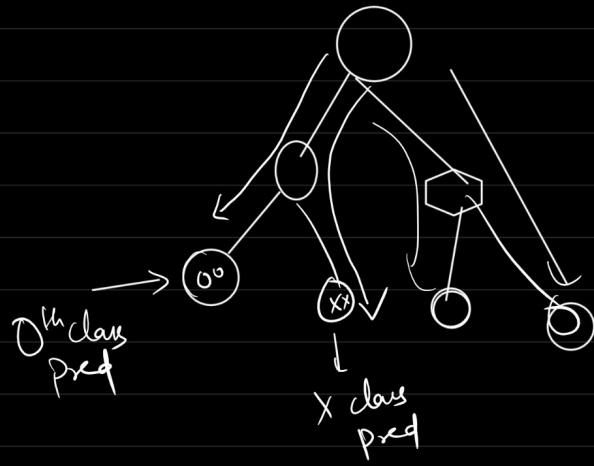
Step-2 Select the best variable for split (Information Gain)

Step-3 Apply the split

Step-4 - Repeat the process for the subset obtained.

Step-5 - Continue the process until every node is a pure node | Stopping criteria is reached.

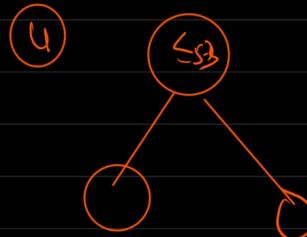
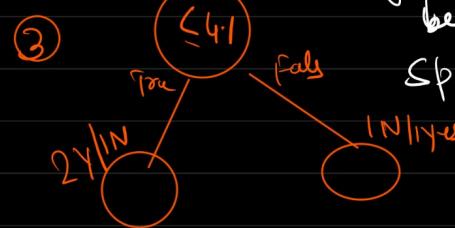
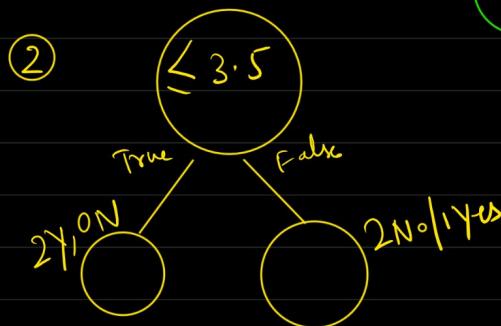
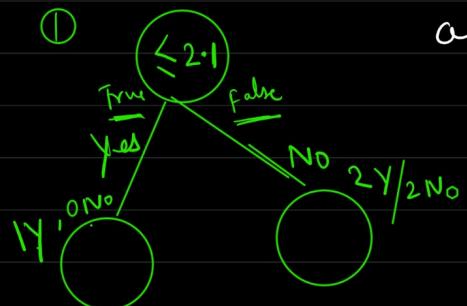
Step-6 → Majority value in the leaf node will be prediction.



\* What if the feature is continuous

$f_1$ (No of hours studied)	Output (Cricket or not)
→ 2.1	Yes
→ 3.5	Yes
→ 4.1	No
→ 5.3	No
→ 6.8	Yes

- ① Sort the feature ( $f_1$ )
- ② treat each value / row as a cut off / threshold and try to build a decision tree and whichever threshold gives you the best information gain that will be used for splitting.



\* Disadvantage

Time-complexity  $\uparrow$

\* Decision tree is a Greedy Algorithm

↓  
It keeps on splitting  
until every  
leaf node is a pure node

↓  
memorise the data

↓  
Overfitting

