

Agenda

→ Decision Tree

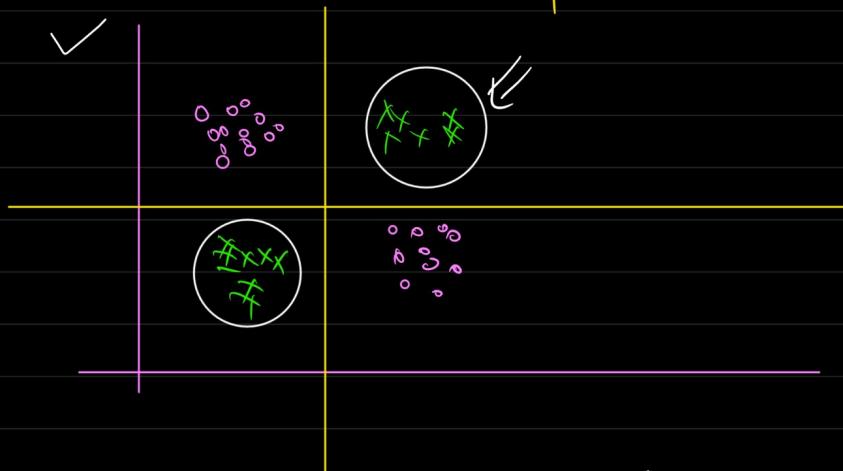
Logistic Regression → Linear Reg + Sigmoid function



Problem wst classification



?



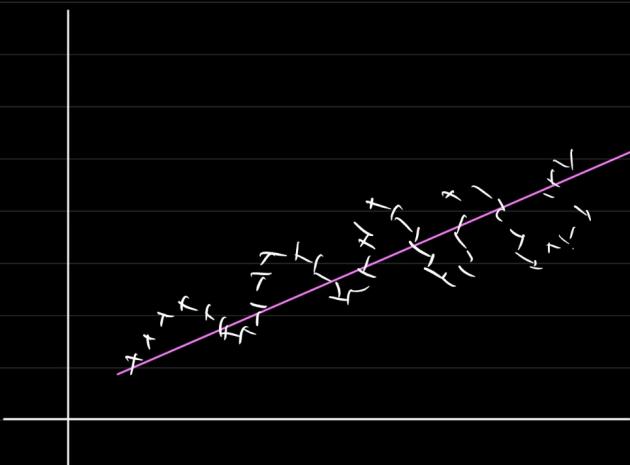
?

No

why? Here data is
non linear

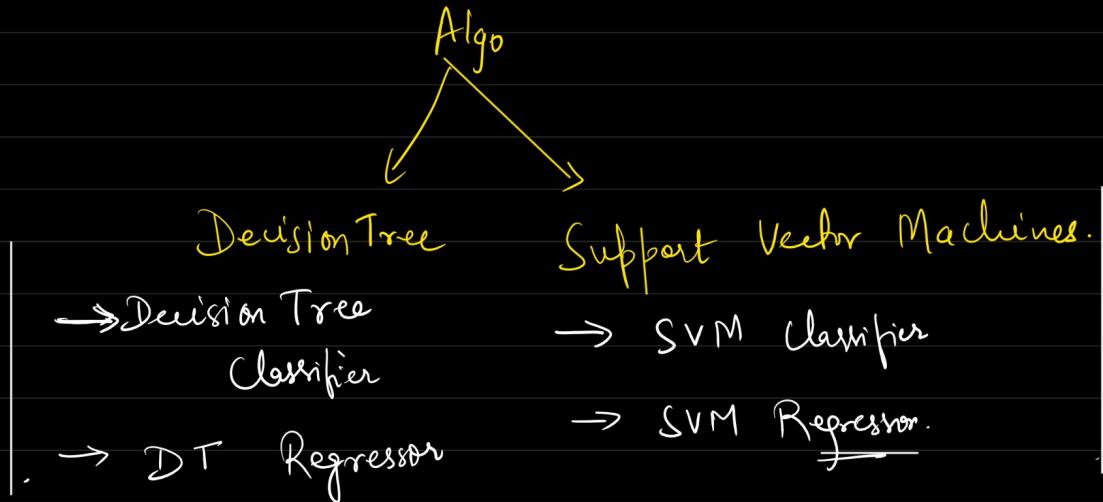
Decision Tree SVM (Support Vector machine)

problem for Regression



→ Polynomial
regression

Preparing
the data/
dealing
with
F.E.
LR



* Decision Tree

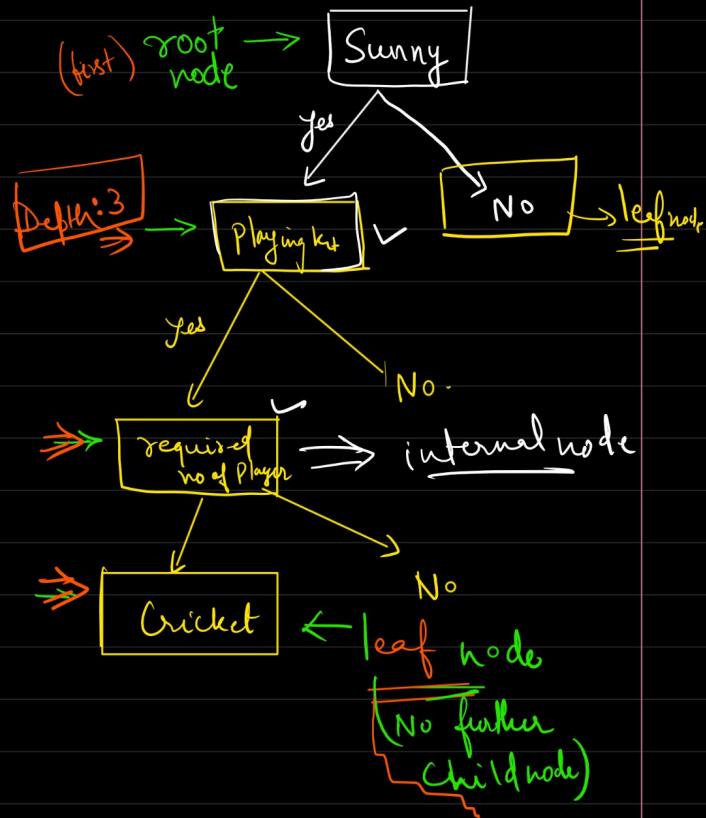
→ DT Cls
→ DT Reg.

Tree

→ At every node, some decision is being made.

→ Decision Tree

Cricket

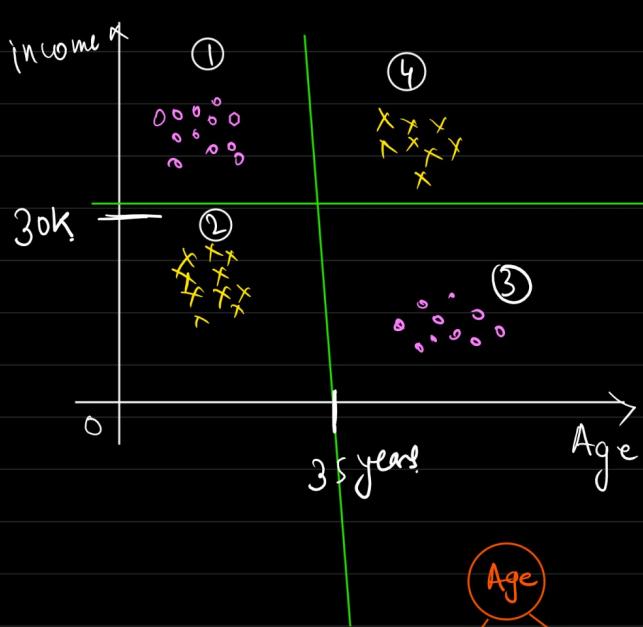


Root node - first node of tree

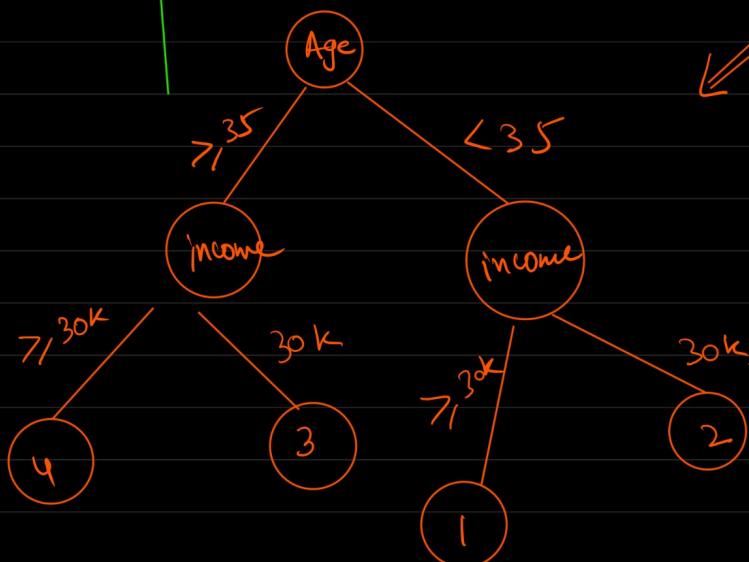
leaf node - No child node

internal node - Node w/ child

Depth → No. of levels.

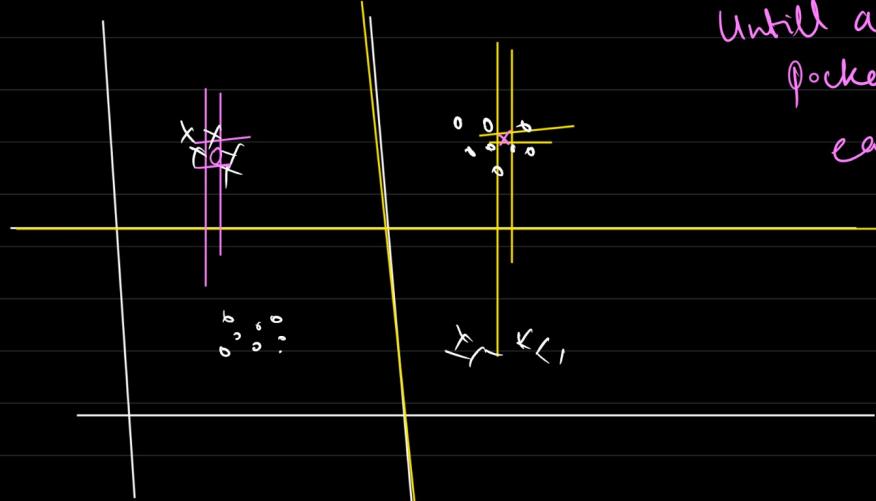


- ① income $> 30K$ & Age ≤ 35 years
↳ group-1 (0^{th} class fraudster)
- ② income $< 30K$ & Age ≤ 35 years
↳ group-2 $\rightarrow x^{th}$ class
↳ Non fraudster.
- ③ if income $\leq 30K$ and Age ≥ 35 years
↳ group-3 $\rightarrow 0^{th}$ class fraudster
- ④ if income $> 30K$ & Age > 35 yrs
↳ group-4
↳ x^{th} class
↳ fraud



\Rightarrow DT divide the data into different Partition or into pockets (nodes)

\rightarrow Keep splitting / dividing until all the pockets where each of the dp is of same class.



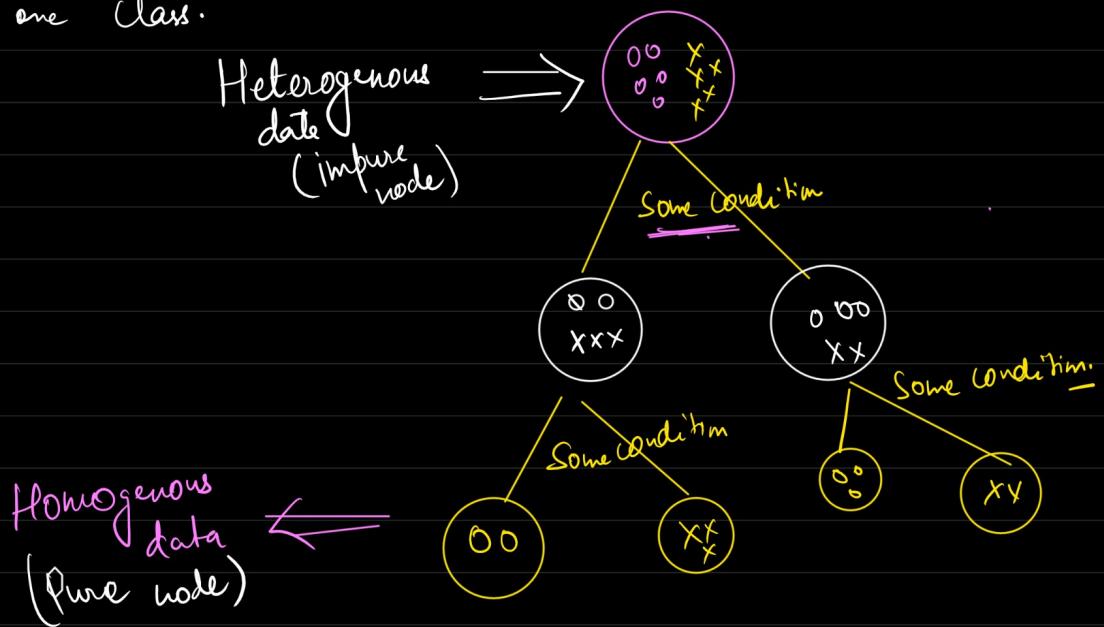
* Advantage of DT

\rightarrow Simple to Understand

\rightarrow Captures the non linear relationship

Intend :-

→ To create leaf node/pocket where each of the data points belong to one class.

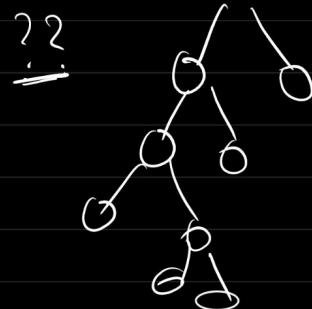
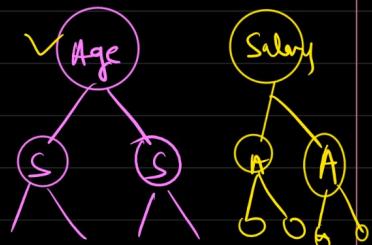


* Pure → homogenous (same class)
impure - heterogenous (mixture)

Q.1. How to split? → On what factor condition should be use?

Q.2. Which feature to be used for splitting?

Q.3. Till when to split?



* How to split

✓ Impurity measures

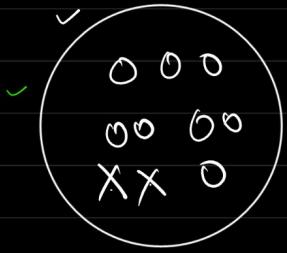
↳ Classification Error

↳ Gini

↳ Entropy

| Salary | Age | Brand |
|--------|-----|-------|
| o | o | |
| o | x | |
| o | x | |
| o | x | |
| x | | |

① Classification error



O : 8

X : 2

$$CE = 1 - \max p_i$$

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

$$P_O = \frac{8}{10} = 0.8$$

$$\Rightarrow = 0.2$$

$$CE = 1 - \max(0.8, 0.2)$$

$$= 1 - 0.8$$

To make pure node, if assign everything to majority class^(o), what is the error rate?

minority misclassified dP's.

* For multiclass



O - 4

X - 2

D - 4

$$P_O = \frac{4}{10}$$

$$P_X = \frac{2}{10}$$

$$D = \frac{4}{10}$$

$$1 - \max(0.4, 0.2, 0.4)$$

$$= \boxed{0.6}$$

$\begin{matrix} 1, 2, 3 \\ \uparrow \quad \uparrow \quad \uparrow r \end{matrix}$

② $Gini$

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

$i \rightarrow$ from class 0 to n-1

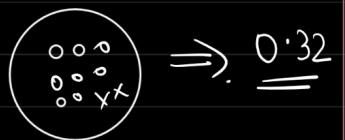


$$\begin{aligned} \Rightarrow 1 - \sum_{i=0}^1 (p_i)^2 &= 1 - \left[(p_0)^2 + (p_X)^2 \right] \\ &= 1 - \left[(0.8)^2 + (0.2)^2 \right] \\ &= 1 - (0.64 + 0.04) \\ &= 1 - 0.68 = 0.32 \end{aligned}$$

$$\begin{aligned} G.I &= \sum_{i=1}^n p_i (1-p_i) \\ &\downarrow \\ &\text{for two clas} \quad p_0 (1-p_0) + p_X (1-p_X) \\ &\downarrow \\ &= 2p_0 (1-p_0) \quad p_0 + p_X = 1 \\ &= 2 \times 0.8 (1-0.8) \quad p_X = 1-p_0 \\ &= 0.32 \end{aligned}$$

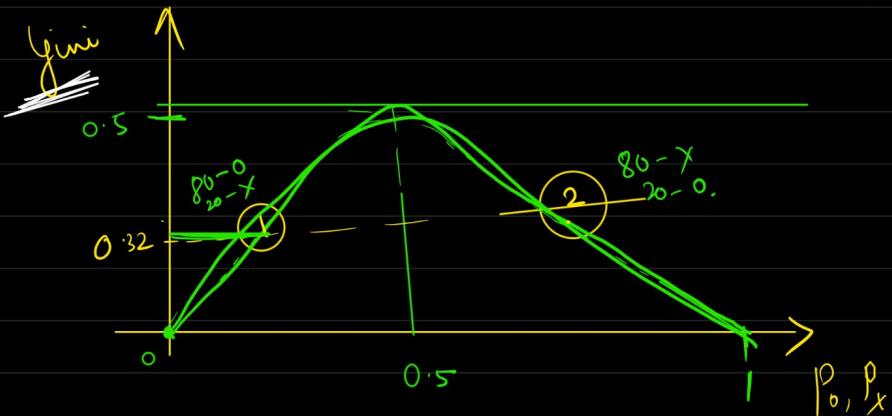
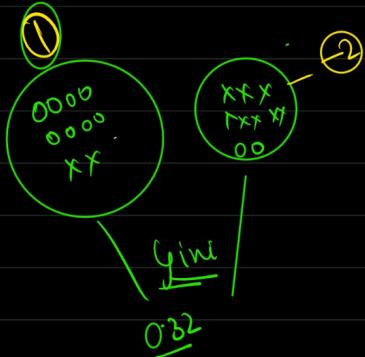
$$\rightarrow \text{Gini} \rightarrow 1 - \sum_{i=1}^n p_i$$

$$\rightarrow 1 - (p_0^2 + p_x^2) \\ = 1 - (0.5^2 + 0.5^2) \\ = 1 - 0.5 = \underline{\underline{0.5}}$$

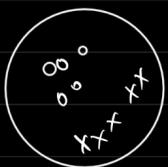


highest gini $\rightarrow \underline{\underline{0.5}}$

$\rightarrow \text{Gini} \rightarrow 0 \rightarrow \text{all of are same} \rightarrow \underline{\underline{\text{no impurity}}},$

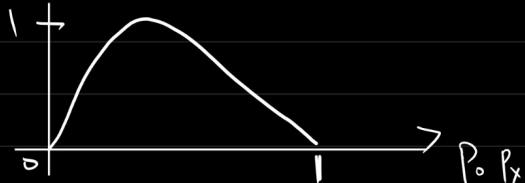


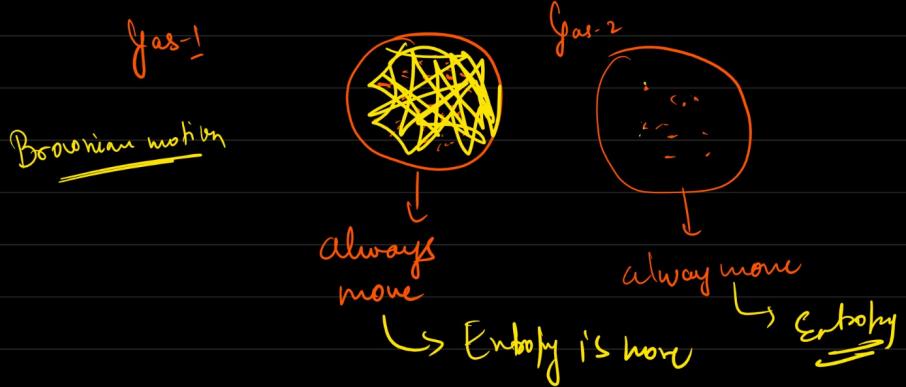
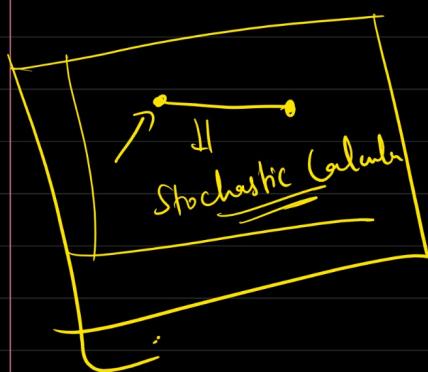
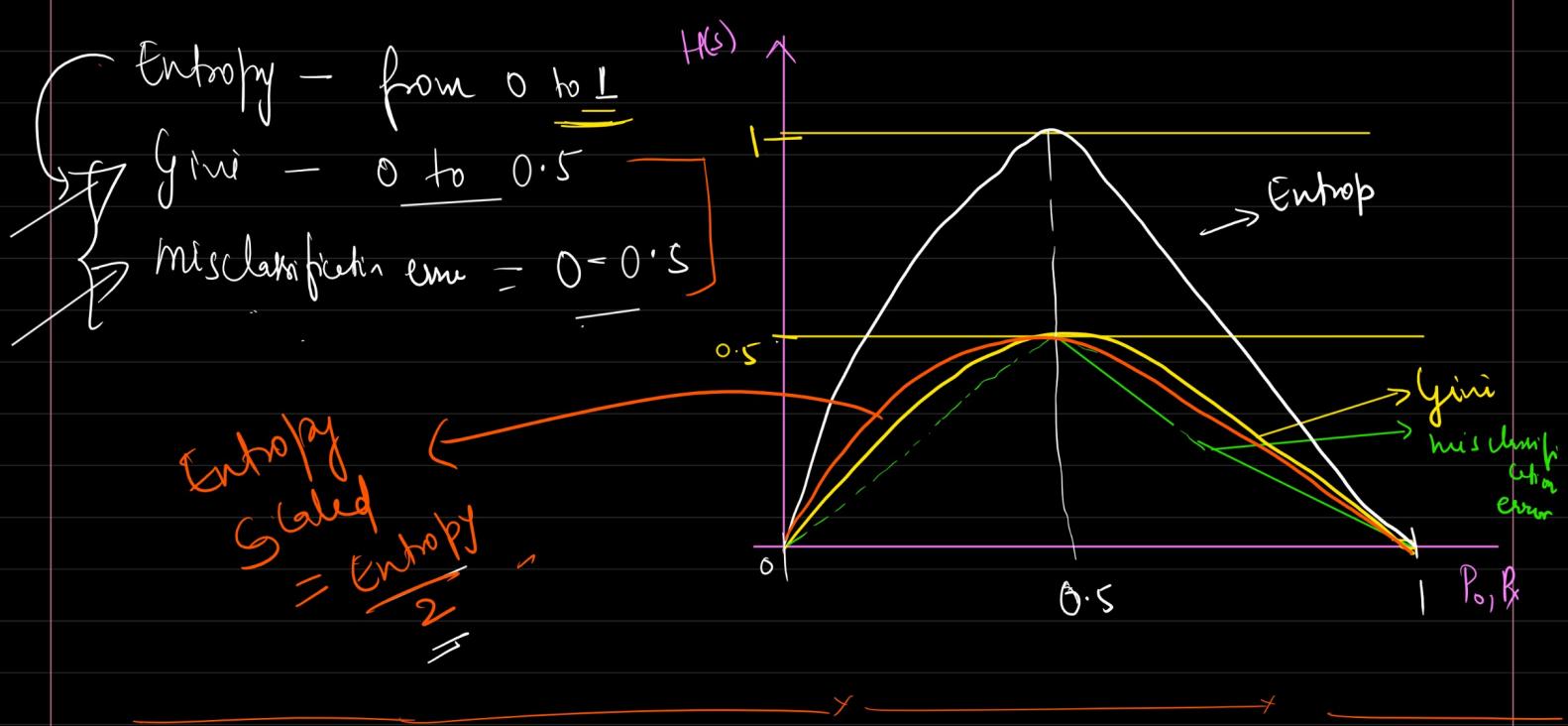
$$③ \text{ Entropy} = - \sum_{i=1}^N p_i \log_2 p_i$$



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

$$\text{Highest entropy} = -2 \cdot \frac{1}{2} \log_2 \frac{1}{2} = -\log_2 \frac{1}{2} = \log_2 2 = 1$$





Class 1

50%

40%

60%

(100%)

0%

Class 2

50%

60%

40%

0%

100%

Entropy

1

0.90

? 0.90

0

0

Multiclass classification

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

$$= 1 - \left[p_1^2 + p_2^2 + p_3^2 + \dots + p_n^2 \right]$$

$$\text{Entropy} = - \sum_{i=1}^n p_i \log_2 p_i$$

$$\rightarrow p_1 \log_2 p_1 - p_2 \log_2 p_2 - p_3 \log_2 p_3 - \dots - p_n \log_2 p_n$$

* Which impurity measure to be used?

$$\underline{Gini} \quad \underline{\text{Entropy}}$$

for small dataset \rightarrow Entropy

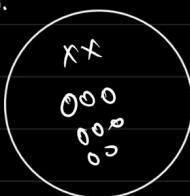
for large dataset \rightarrow Gini

(computation)

* On what features to split

\hookrightarrow Information gain

disorganization/randomness



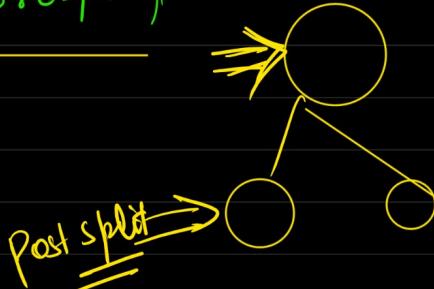
Aim of splitting

We want pure nodes after splitting

Such a split that maximise the purity

impurity (post split) < impurity (pre split)

Split using that feature which does this best



Post split

Impurity (post split) < impurity (pre split)

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

if randomness decreases, purity increase,
then we say Information gain.

impurity
randomness
to be removed
pure

(randomness decreased) \Rightarrow Information Gain
 $(\Delta \text{impurity})$

* Information gain is a measure
to define degree of disorganisation | randomness

In a system \Rightarrow IG \rightarrow GI
IG \rightarrow Entropy

* Whichever feature will have maximum information, that feature will be used for split.

* decrease in randomness \Rightarrow Information Gain.

| Gender | Class | Cricket (Yes/No) |
|--------|-------|------------------|
| M | 1X | Yes |
| F | X | No |
| M | 1X | Yes |
| - | - | - |
| - | - | - |
| - | - | - |
| - | - | - |
| - | - | - |

total students: 30 → 50% play cricket Yes - 15
No - 15

Gender $\rightarrow 10/2 = F$ - play cricket (out of 10 females, 2 play cricket)

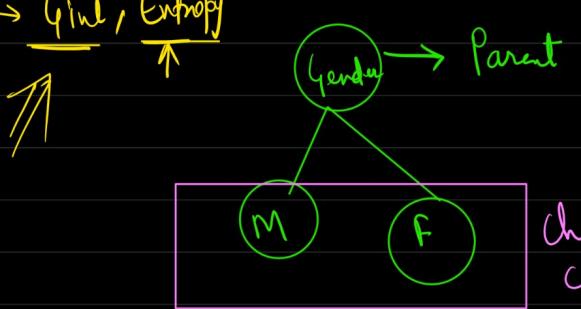
Class $\rightarrow 14/16 = M$ - play cricket (out of 14 class 1x students, 13 play cricket)

Gender $\rightarrow 16/16 = X$ - out of 16 class x student 5 + 1 play cricket



* Whichever feature has the highest information gain, used for splitting.

Impurity \rightarrow Gini, Entropy



$$IG = \text{Imp. Pre-split} - \text{Imp. Post-split}$$

$$= \text{En}(Parent) - \text{En}(\text{Child combined})$$

$$\text{En Parent} \Rightarrow -\frac{15}{30} \log_2 \frac{15}{30} - \frac{15}{30} \log_2 \frac{15}{30} = 1$$

En(Child)

$$\left\{ \begin{array}{l} \text{En(Female child)} = -\frac{2}{10} \log_2 \frac{2}{10} - \frac{8}{10} \log_2 \frac{8}{10} \\ = 0.72 \end{array} \right.$$

$$\left. \begin{array}{l} \text{En(Male child)} = -\frac{13}{20} \log_2 \frac{13}{20} - \frac{7}{20} \log_2 \frac{7}{20} \\ = 0.92 \end{array} \right.$$

$$\text{En}(\text{Child combined}) = \frac{\text{No of females}}{\text{Total Students}} \times \text{En Female child} + \frac{\text{No of males}}{\text{Total Students}} \times \text{En Male child}$$

$$\text{Weighted Entropy} = \frac{10}{30} \times 0.72 + \frac{20}{30} \times 0.92 \Rightarrow 0.85$$

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

En Child class

$$\text{En}(1x) = \frac{43}{100} \log_2 \frac{43}{100} - \frac{57}{100} \log_2 \frac{57}{100} = x_1$$

$$\text{En}(X) = -\frac{57}{100} \log_2 \frac{57}{100} - \frac{43}{100} \log_2 \frac{43}{100} = x_2$$

$$\text{En Combined} = \frac{14}{30} x_1 + \frac{16}{30} x_2$$

$$= 0.99$$

$$IG = E_p - E_c = 1 - 0.99 = 0.01$$

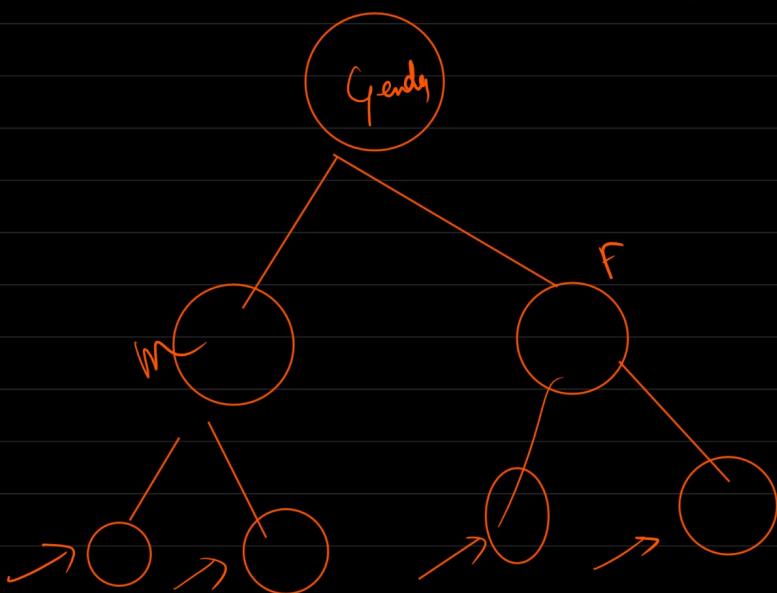
$$\text{IG} = E_p - E_c$$

$$= 1 - 0.85$$

$$= \underline{\underline{0.15}}$$

Since $\text{IG}_{\text{Gender}} > \text{IG}_{\text{Class}}$

\therefore we will split based on Gender



Algorithm of DT

\rightarrow Step 1 \rightarrow Recursive binary splitting | Partitioning the data into smaller subsets

\rightarrow Step 2 - Select the feature to split (Information gain)

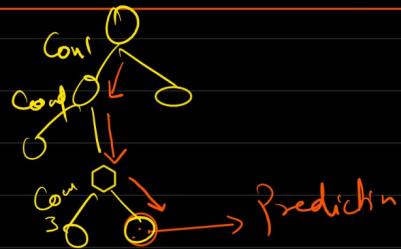
Step 3 - Apply the split

Step 4 - Repeat the process for the subset obtained

Step 5. \rightarrow Continue the process until every node is a purenode / Stopping criteria is reached.

\Rightarrow Step 6 - Majority value in the leaf nodes will be prediction

Age=100, Salary 50k,
Gender=M



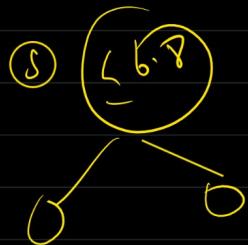
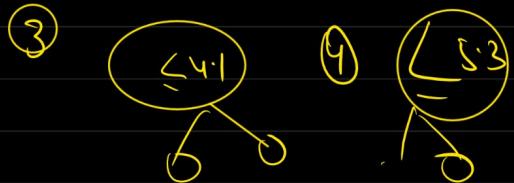
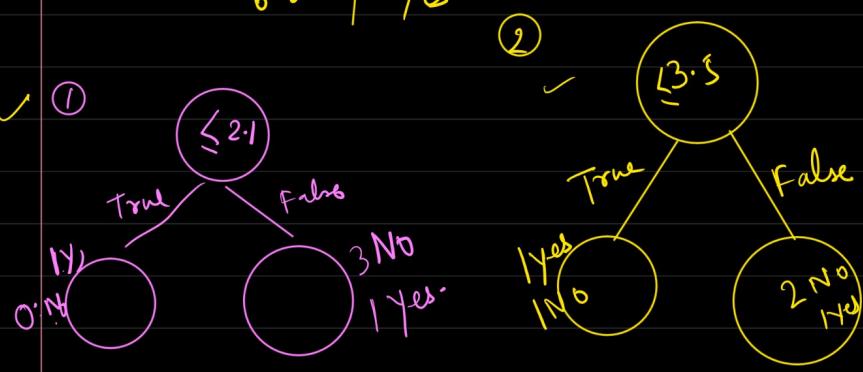
* For categorical feature



* What if feature is continuous

| f_1 | output |
|-------|--------|
| 2.1 | Yes |
| 3.5 | No |
| 4.1 | No |
| 5.3 | No |
| 6.8 | Yes |

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

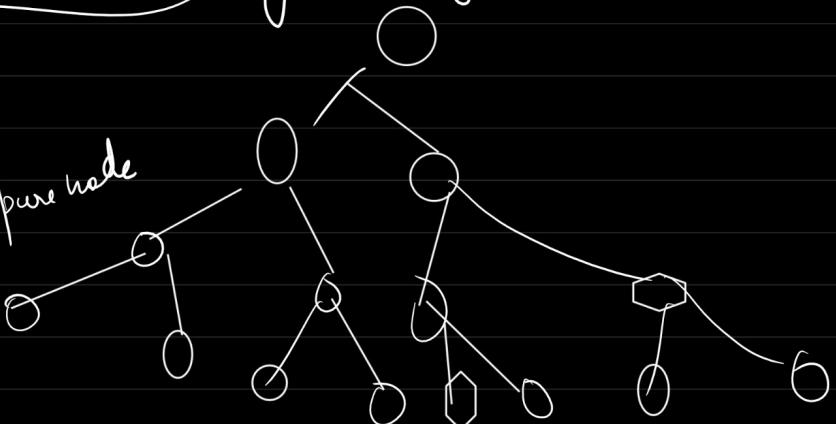


→ whichever split
 gives the highest IG
 that will be used
 for splitting. (feature/value)

* Decision tree is a greedy algorithm.

↳ it keeps on
 splitting until
 every leaf node is a pure node

↓
 memorise the data



Overfitting.



Overfitting → training acc ↑
testing acc ↓

