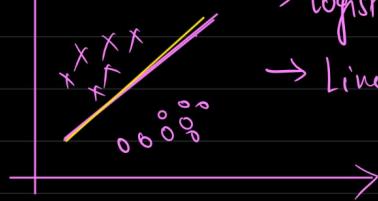
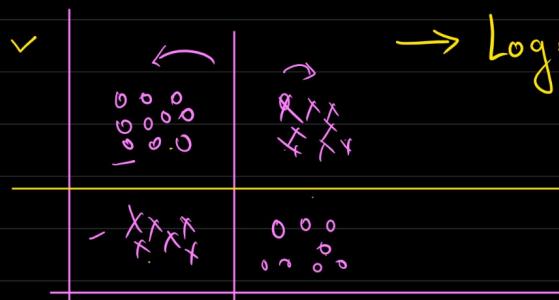


* Decision Tree

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

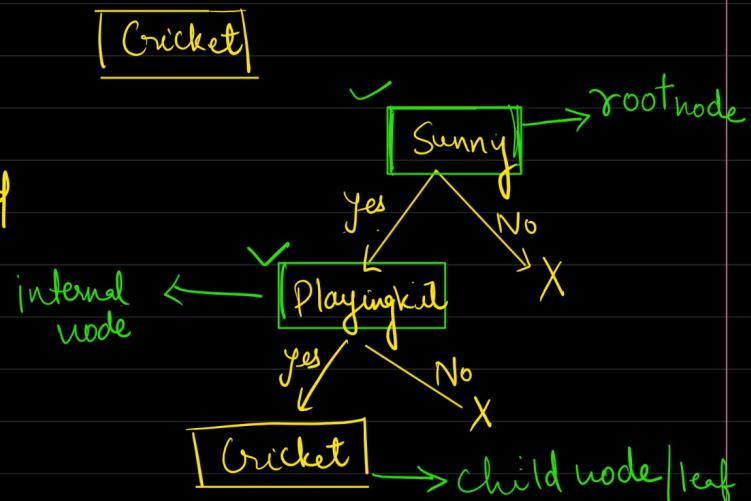
Classification problem (Decision tree classification)
Regression problem (Decision tree regression)

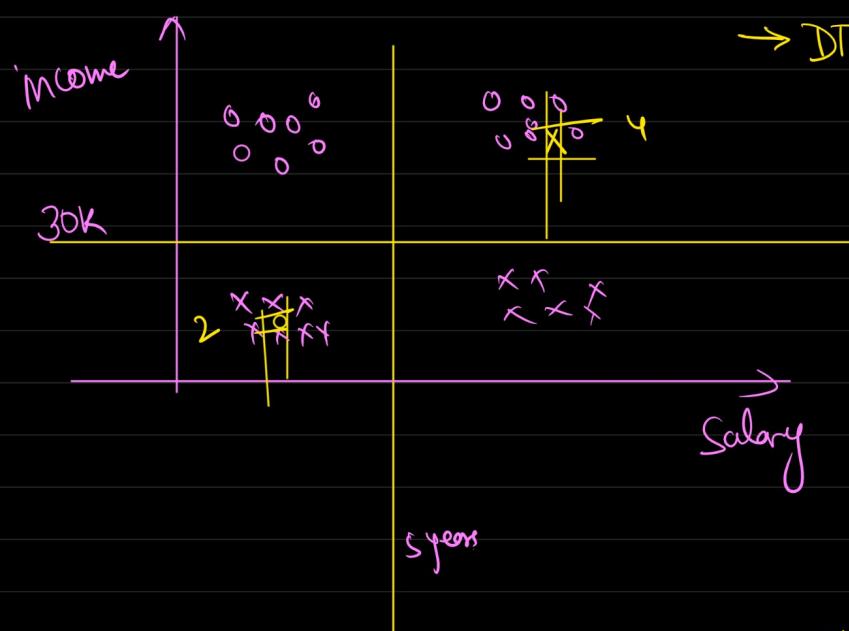
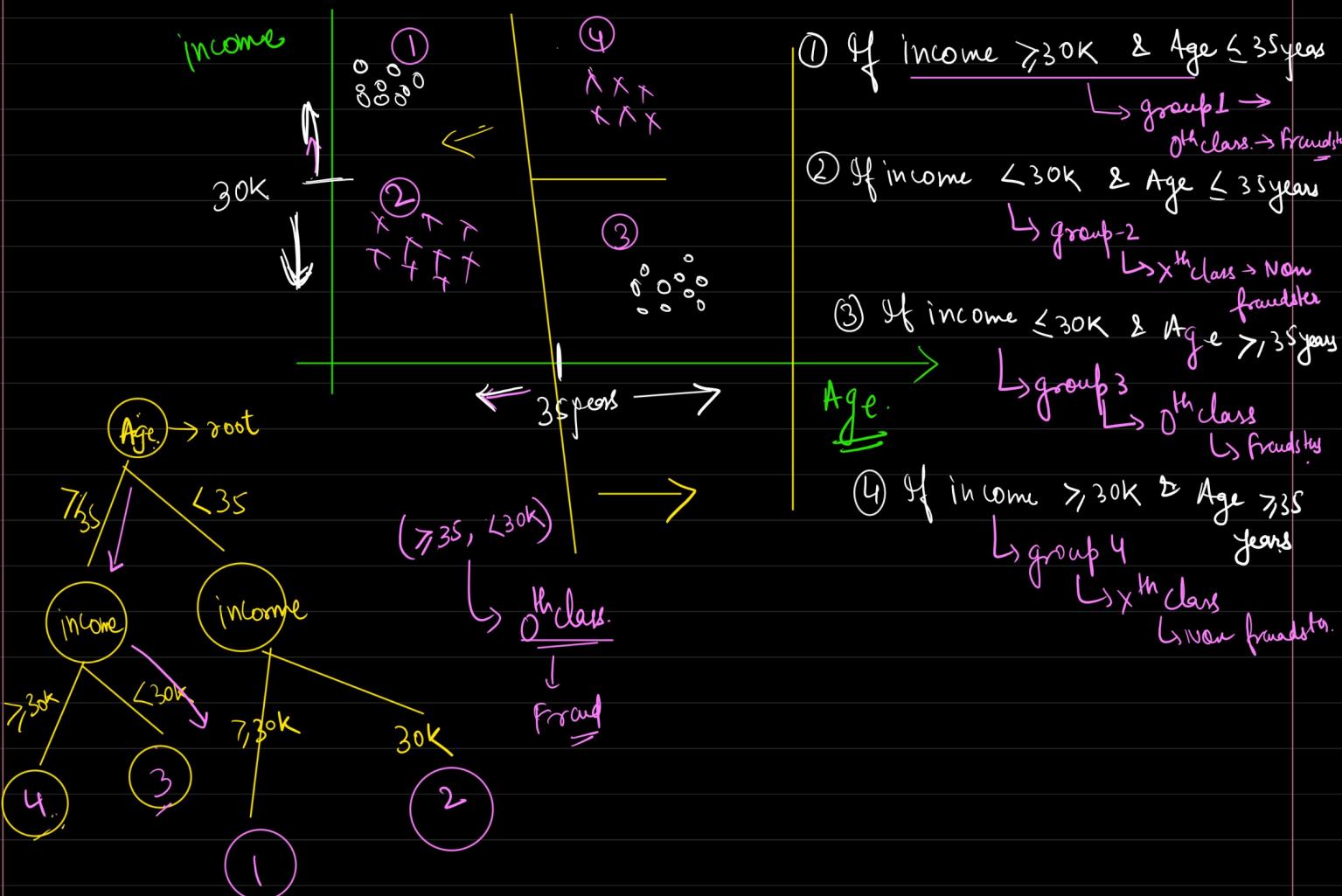
* Decision Tree works on nested if else condition.

Decision Tree

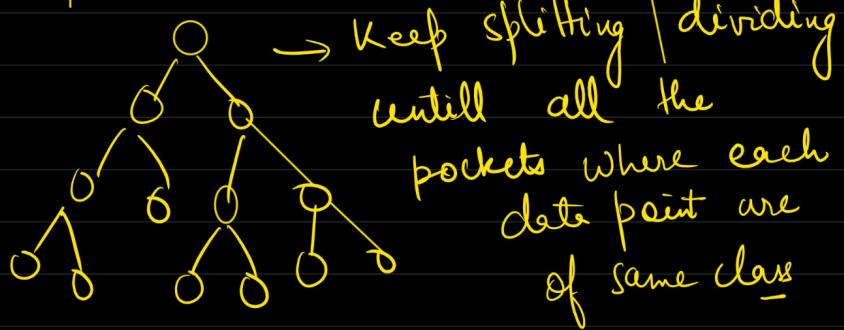
At every node a decision is taken

It looks like a tree



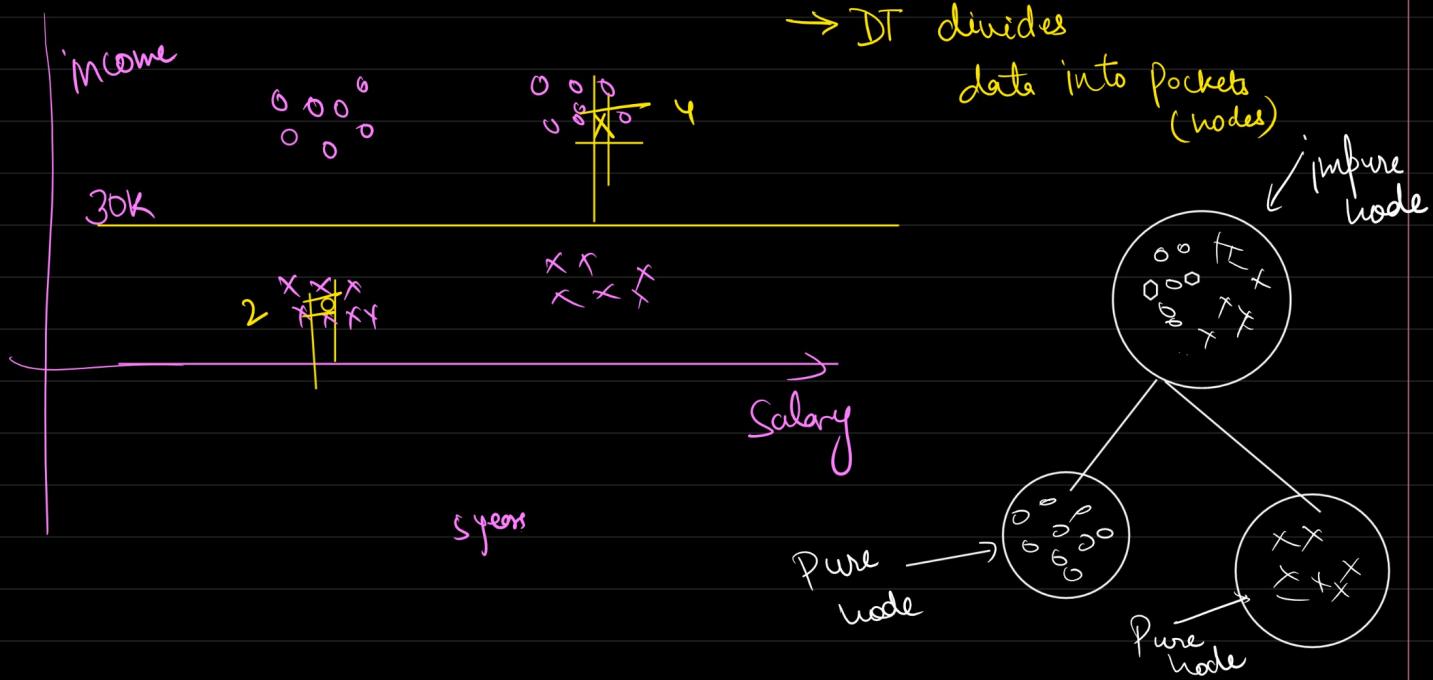


→ DT divides the data into pockets (nodes)

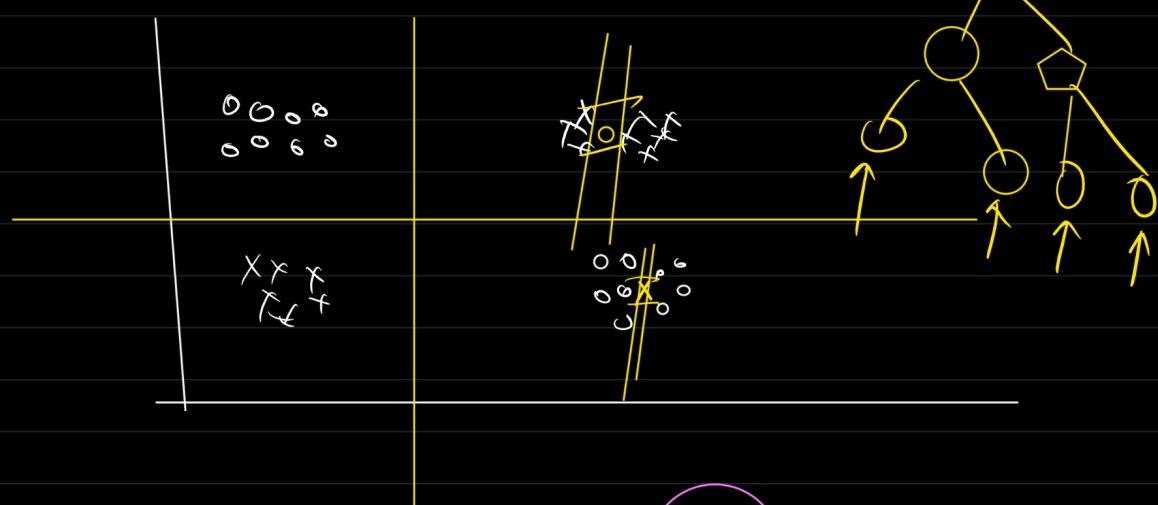


* Advantages

- 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

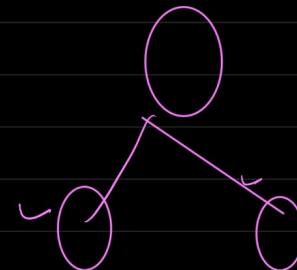
Some condition

Pure node



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

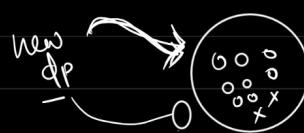
pure - homogeneous node (Same class)
impure - heterogeneous node



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

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

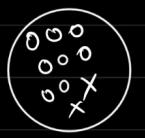
$$(G_I = \sum_{i=0}^{n-1} p_i(1-p_i))$$

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

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

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

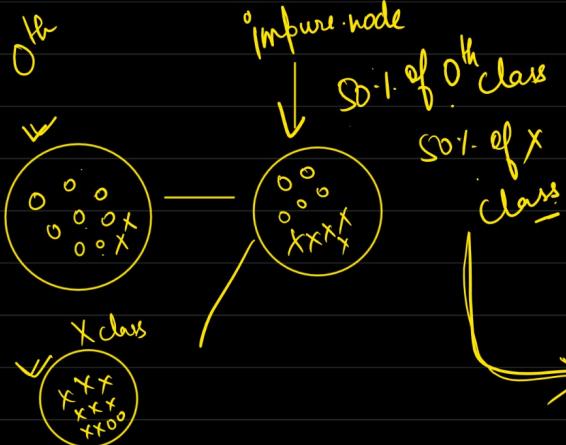
$$= 0.8(1-0.8) + 0.2(1-0.2) \\ = 0.8 \times 0.2 + 0.2 \times 0.8 \\ = 0.32$$



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

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

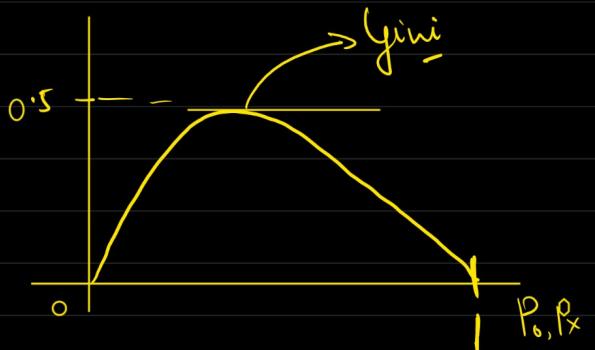
$$= 1 - (0.64 + 0.04) \\ = 1 - 0.68 \Rightarrow 0.32$$



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

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

highest Gini impurity value = 0.5



③ Entropy

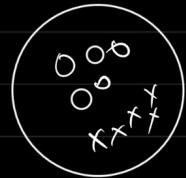
$$H(S) = - \sum_{i=0}^{N-1} p_i \log_2 p_i$$



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

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

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

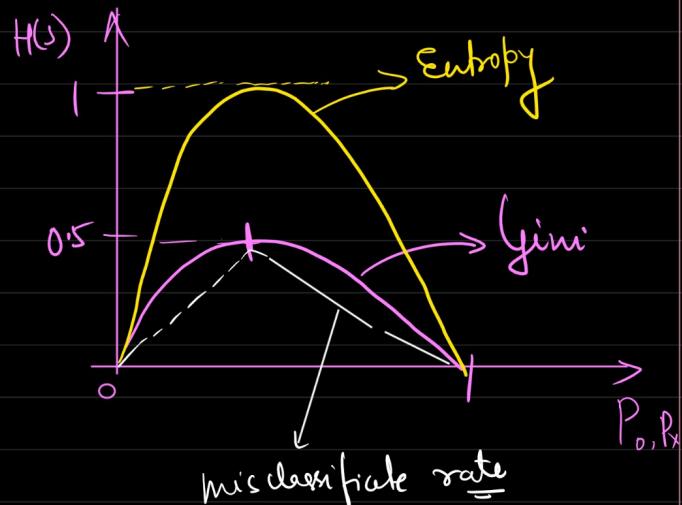
$S01-S01$ \rightarrow 1

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

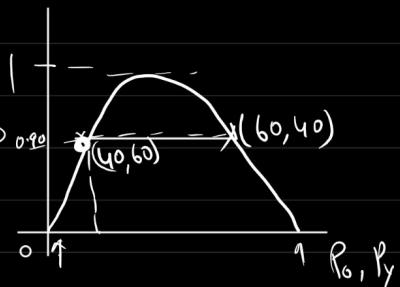
1

0.5

0.90

0.40

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