

# Decision Trees

E.g :-

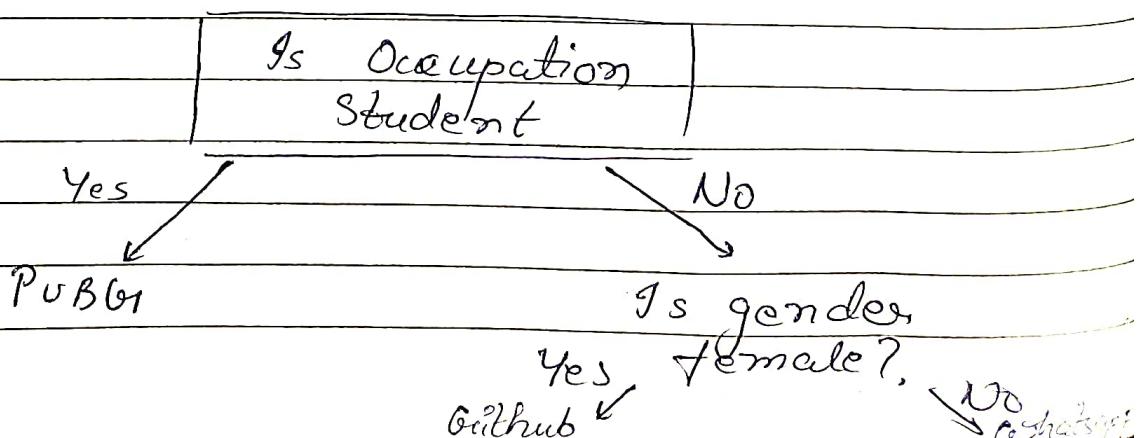
<u>Gender</u>	<u>Occupation</u>	<u>Suggestion</u>
F	Student	PUBG
F	Programmer	Github
M	Programmer	Whatsapp
F	Programmer	Github
M	Student	PUBG
M	Student	PUBG

In python using nested if-else

```

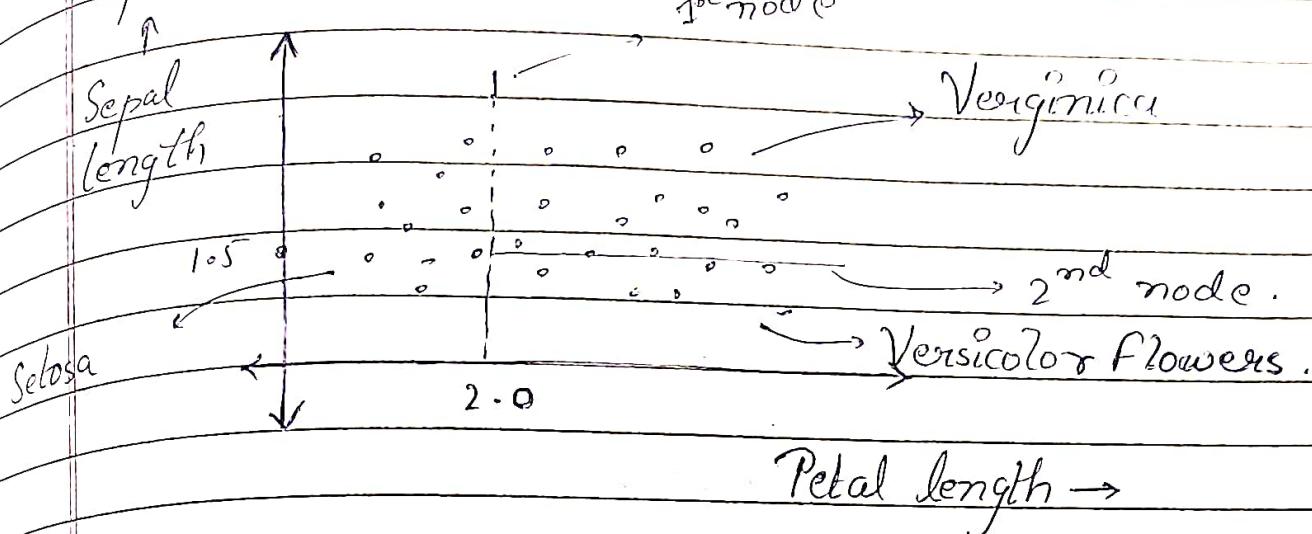
if occupation == student
    print (PUBG)
else
    if gender == female
        print (Github)
    else
        print (Whatsapp).
    
```

↳ Where is the Tree?



E.g Iris Dataset

Geometric Intuition:



Pseudo Code:

- 1. Begin with your training dataset, which should have some feature variables ~~and~~ as classification or regression output.
- 2. Determine the 'best feature' in the dataset to split the data on; more on how we define "best feature" later.
- 3. Split the data into subsets that contain the correct values for this best feature. This splitting point basically defines a node on the tree. i.e each node is a splitting point based on a certain features from our data.
- 4. Recursively generate new tree nodes by using the subset of data created from step 3.

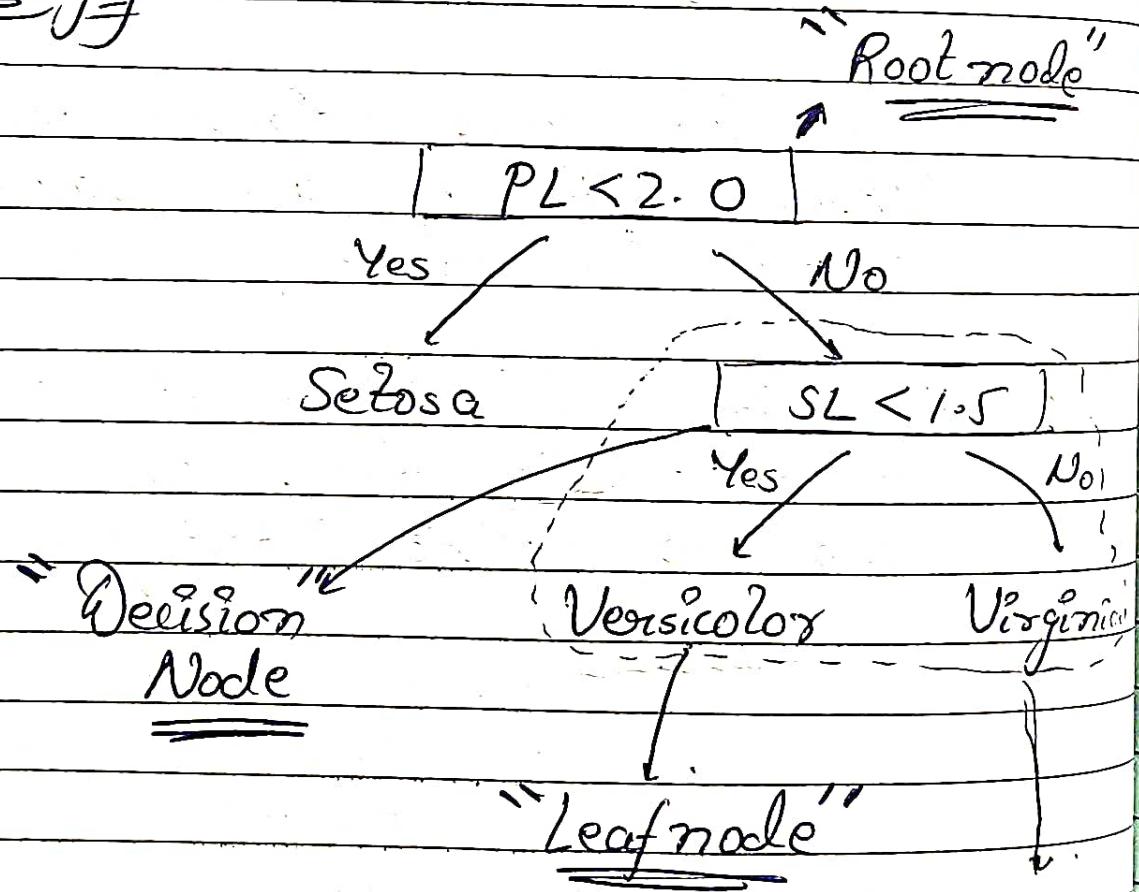
Conclusion :-

Programmatically speaking, Decision trees are nothing but a giant structure of nested if-else condition.

Mathematically speaking, Decision trees use hyperplanes which run parallel to any one of the axes to cut the coordinate system into hyper cuboids.

Terminology :

Sample  
Decision  
Tree



"Branch  
Subtree"

Some unanswered ques.

b) How to decide which col. should be considered as root node?

b) How to select subsequent decision node?

b) How to decide splitting criteria in case of numerical columns?

Advantages :

- b) Intuitive and easy to understand.
- b) Minimal Data preparation is required.
- b) The cost of using the tree for inference is logistic in the number of data points used to train the tree.

Disadvantages :

- b) Overfitting.
- b) Prone to errors for imbalanced datasets.

Points :

- b) Decision Trees are mostly used for classification problems but can also be applied to regression, hence the name CART.

1.

Classification & Regr Trees

## Working of decision Tree

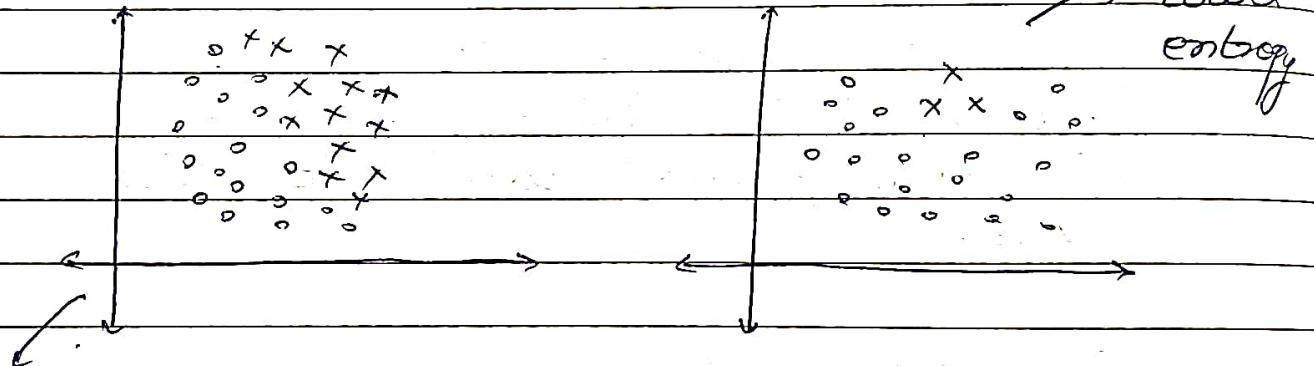
Entropy? (Degree of randomness)

Entropy is nothing but the measure calculated of disorder. Or we can also only on call it the measure of purity / the basis impurity of label.

E.g.:

Higher entropy  $\rightarrow$  ~~impure~~ lower  $\rightarrow$  pure

lower entropy



higher entropy more knowledge less entropy

Mathematical formula of entropy

$$E(S) = \sum_{i=1}^n -p_i \log_2 p_i \quad i \Rightarrow \text{classes in label}$$

Where  $p_i$  is simply the frequentist probability of an element / class "i" in our data.

For e.g.

If our data has only 2 class labels Yes and No.

$$E(D) = -P_{\text{Yes}} \log_2(P_{\text{Yes}}) - P_{\text{No}} \log_2(P_{\text{No}})$$

Eg. Dataset:

D<sub>(1)</sub>

D<sub>(2)</sub>

Salary	Age	Purchase	Salary	Age	Purchase
20000	21	Yes	34000	31	No
10000	45	No	15000	25	No
60000	27	Yes	69000	57	Yes
15000	31	No	25000	21	No
12000	18	No	32000	28	No

$$H(d) = -P_y \log_2(P_y) - P_n \log_2(P_n)$$

$$H(d) = -2 \log_2 \left( \frac{2}{5} \right) - 3 \log_2 \left( \frac{3}{5} \right) \quad H(d) = -\frac{1}{5} \log_2 \left( \frac{1}{5} \right) - \frac{4}{5} \log_2 \left( \frac{4}{5} \right)$$

$$H(d) = 0.97$$

$$H(d) = 0.72$$

→ Higher entropy.

for 3 classes or more classes. (Yes, No, May Be)

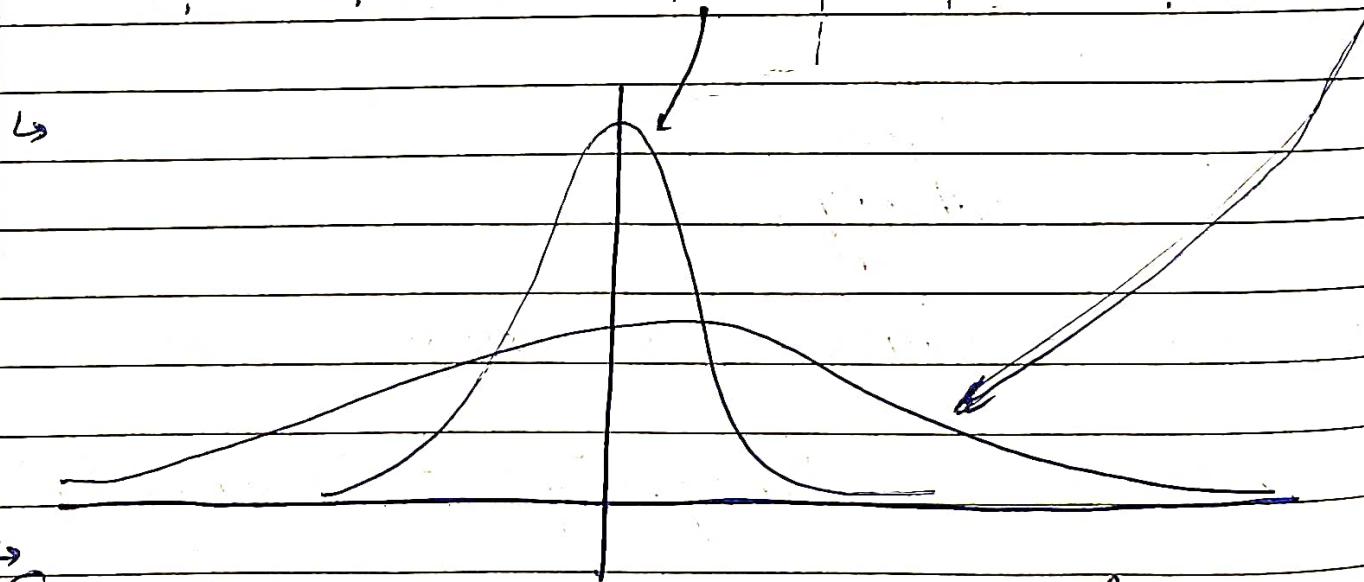
$$H(d) = -P_y \log_2(P_y) - P_n \log_2(P_n) \\ - P_m \log_2(P_m)$$

## Observation :

- ↳ More the uncertainty more is entropy.
- ↳ For a 2 class problem the min entropy is 0 and the max is 1.
- ↳ For more than 2 classes the min entropy is 0 but the max can be greater than 1.
- ↳ Both  $\log_2$  or  $\log_e$  can be used to calculate entropy.

## • Entropy for continuous variable

Eg:	Area	Built in	Price	Area	Built in	Price
	1200	1999	3.5	2200	1989	4.6
	1800	2011	5.6	800	2018	6.5
	1400	2000	7.3	1100	2005	12.8
	:	:	:	:	:	:



↳ So intuitively we can say the

randomness of label in dataset 2 is more. So, its entropy is higher than dataset 1.

## Information Gain:

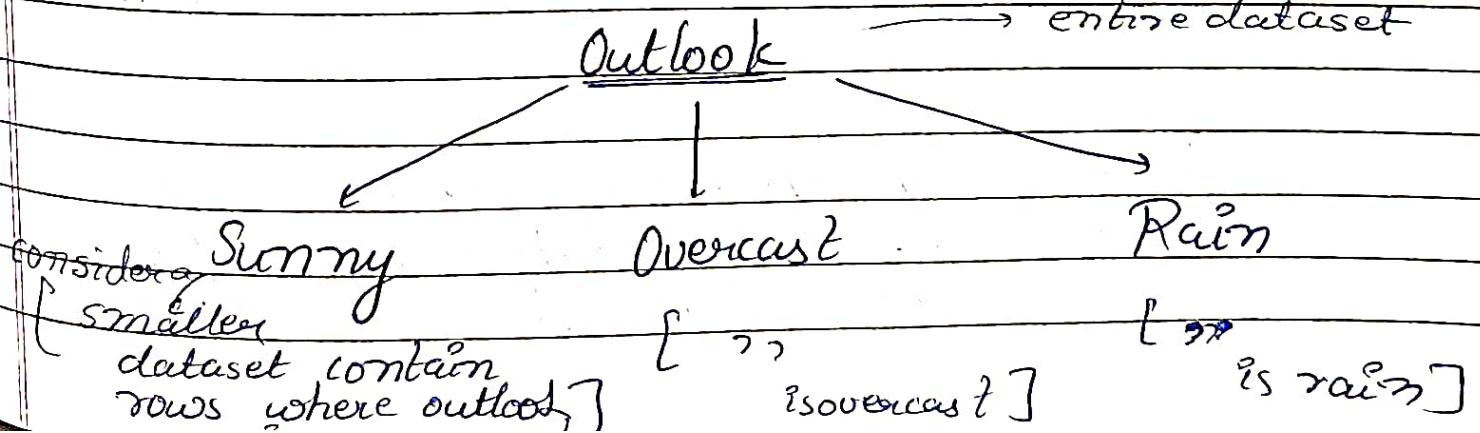
It is a metric used to train Decision Trees. Specifically, this metric measures the quality of a split.

The information gain is based on the decrease in entropy after a data-set is split on an attribute. Constructing a decision tree is all about finding attribute that returns the highest information gain.

$$\text{Information gain} = E(\text{parent}) -$$

$$\sum [\{\text{Weighted Average}\} * E \{\text{Children}\}]$$

Step 1:- entropy of parent. Weighted entropy



Step 2 : Calculate entropy for children.

Each child i.e  
for sunny, overcast & rain

Step 3 : Calculate weighted Entropy of children

$$\text{Weighted Entropy} = \frac{\text{no. of row in sunny dataset}}{\text{Total no. of row in the entire dataset}} \times \text{Entropy of sunny}$$

$$+ \frac{\text{no. of row in overcast}}{\text{Total no. of row in the entire dataset}} \times \text{Entropy of overcast}$$

$$+ \frac{\text{no. of row in rain}}{\text{Total no. of row in the entire dataset}} \times \text{Entropy of rain}$$

Step 4 : Calculate Information Gain

So the information gain when we split this data on the basis of Outlook col. is say 0.28

Step 5 : Calculate Information Gain for all the col.

"Whichever column has the highest Information Gain (maximum decrease in entropy) the algorithm will select that column to split the data"

Gini Impurity :

This is just like entropy, a measure to check the purity/impurity of data column.

Only difference is formula :-

$$G_I = 1 - \sum p_i^2$$

i.e for 2 class classification.

$$G_I = 1 - (P(y)^2 + P(m)^2)$$

Higher Gini  $\rightarrow$  Impure

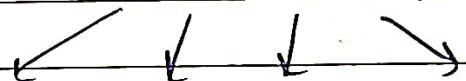
The only benefit is, it is computationally fast.

## • Handling Numerical Data.

E.g.

S.No	User Rating	Downloaded
1	3.5	Yes
2	4.6	Yes
3	2.5	No
4	1.6	Yes
:	:	:

User rating



$n = \text{no. of rows.}$

$\dots \rightarrow n$  number  
of sub trees.

So computationally ~~done~~ very  
much infeasible.

Step 1 :

Sort the data on the basis of  
numerical column.

Step 2 :

Split the entire data on the basis  
of every value of user-rating.

i.e. say

User rating  $> 1.6$

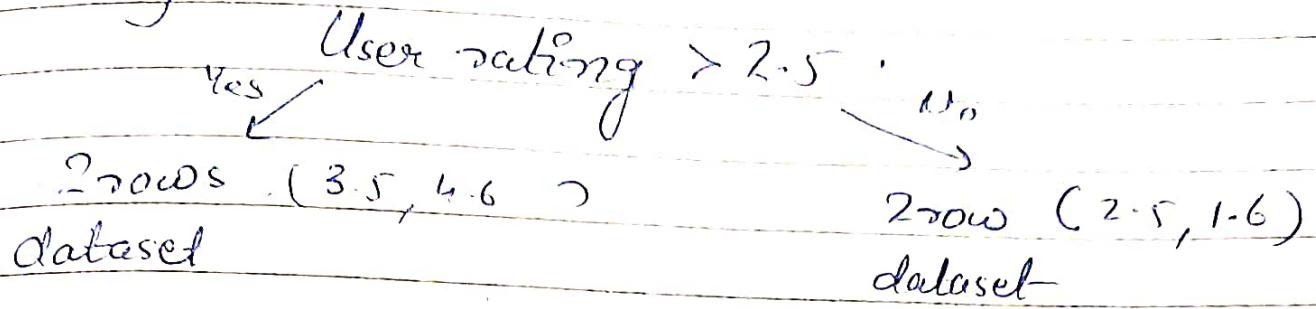
Yes

No

dataset containing  
every row except one

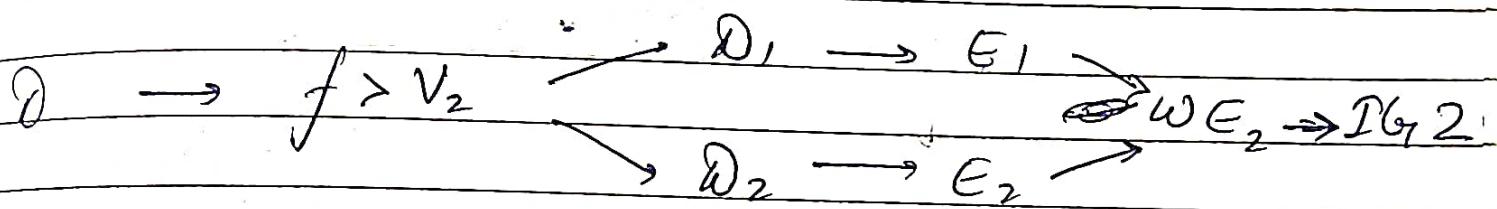
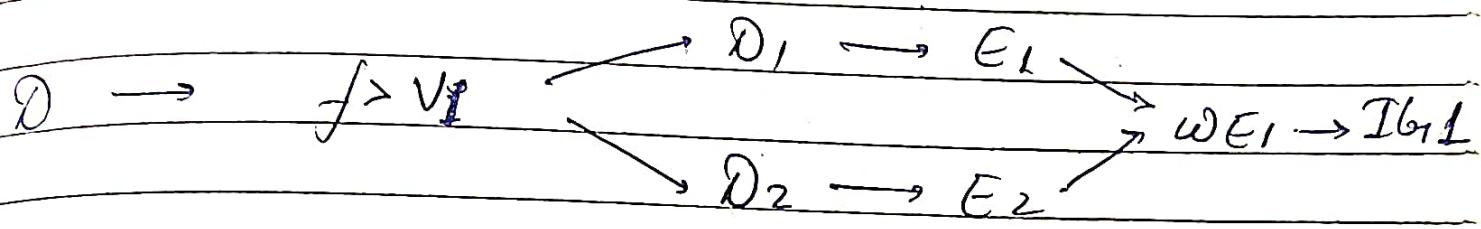
only one  
row logically

again say



and keep on repeating.

So,



Step 3 :-

Max {  $IG_1, IG_2, IG_3, \dots$  }

based onto which  
we select our splitting criteria.

# • Decision Trees Hyperparameter

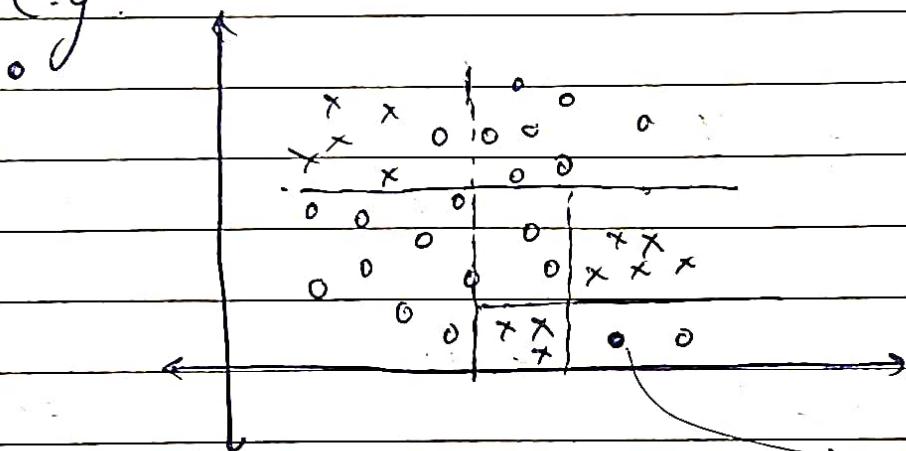
To reduce overfitting

~~Depth of tree~~

Performs well → Training  
Performs badly → Test

↳ Max\_depth

E.g.:

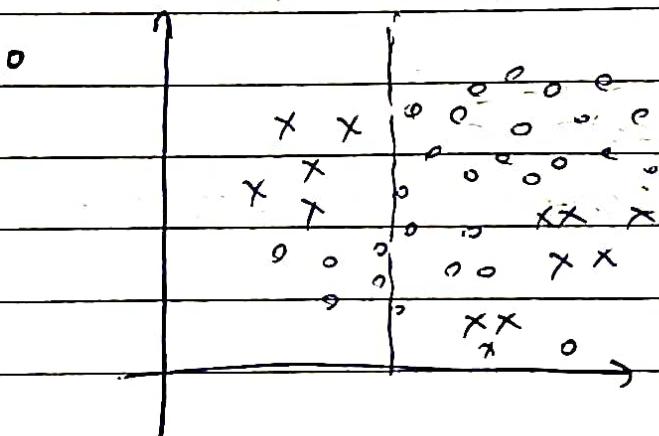


- max\_depth = 'None'

↓  
overfitting

consider this point.

logically it seems that it should be  $x$  but decision tree classified it as  $o$ .



- max\_depth = '1'

↓  
Underfitting

Splitter = 'best' or 'random'

Split best fitting  
criteria i.e.  
highest Info.gain

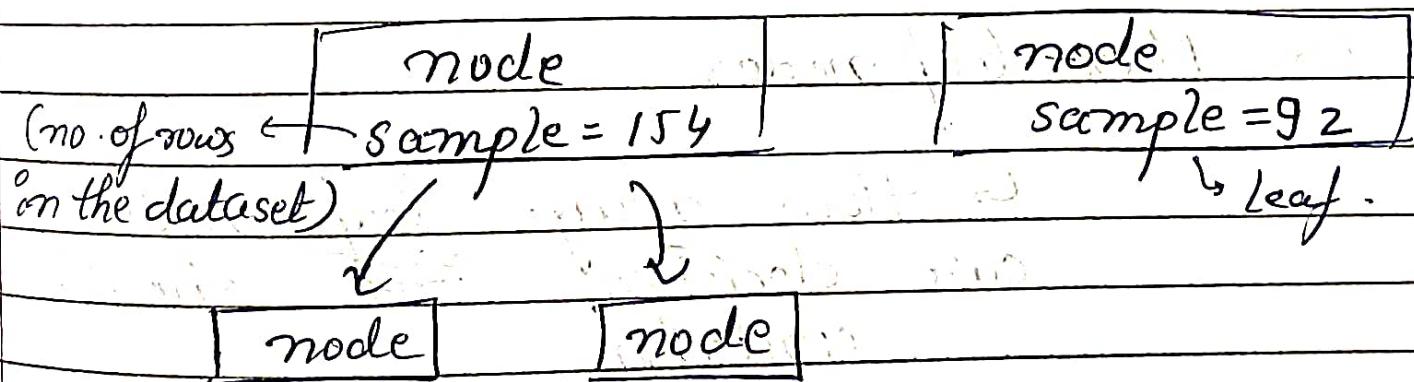
Makes any  
random split  
(best for  
reducing  
overfitting)

Min Samples split :

Using this we can decide that  
if, minimum =  $x$  rows will present  
on a node (means node dataset)  
then an only then the split happen  
rather the node becomes a leaf



Say min Samples split = 100



Min Sample leaf :

means if the leaf node on splitting  
comes down to min Sample leaf  
Value, it will not allow to split.

E.g  $\text{min Sample Leaf Node} = 100$

sample = 116
--------------

↳ no further split

(Because if we split the no. of sample reduces below 100)

Means the leaf node must have 100 no. of samples.

↳ Max features

↳ How many cols. to be given to decision tree

↳ Max leaf nodes

↳ How many leaf node our decision tree should contain.