

1/3/23

# # DECISION TREE & RANDOM FOREST

## Topics-

- Introduction
- Entropy & information Gain
- Process kaggle Titanic Dataset
- Implementing Information Gain
- Implementing Decision Tree
- Making predictions
- Decision Tree using sklearn
- Random Forest / Ensembles

## # Entropy & Information Gain

### → Entropy

measure  
of  
randomness

measures the impurity or uncertainty present in the data.

OR

degree of disorder present in the data.

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

Randomness x Entropy

Where,

S = set of all instances

N = number of distinct class values

$p_i$  = event probability

### → Information Gain

It indicates how much information a particular feature / variable gives us about the final outcome



$$\text{Gain}(A, S) = H(S) - \sum |S_i| \cdot H(S_i) = H(S) - H(A, S)$$

where,

$H(S)$  = entropy of whole dataset

$|S_i|$  = no. of instance with  $i$  value

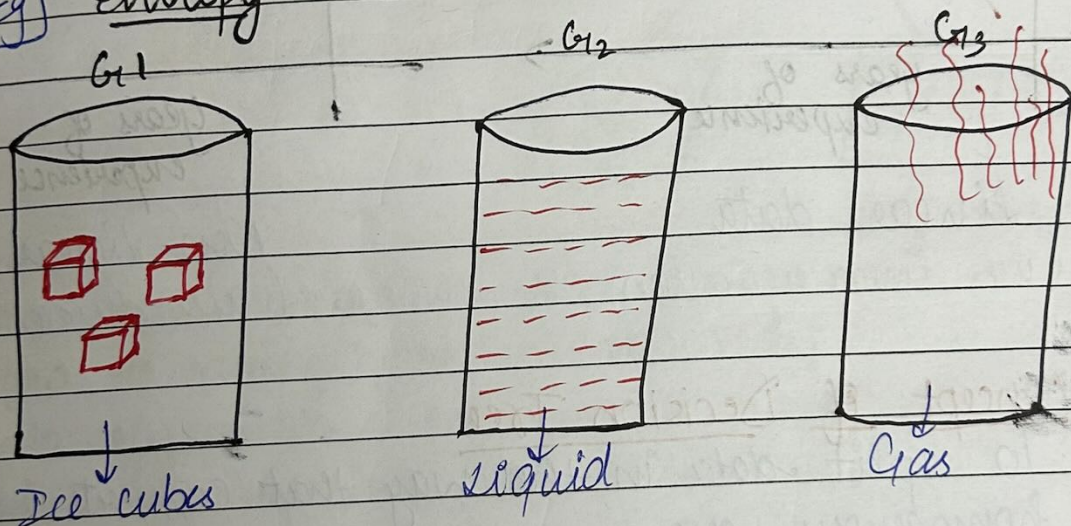
$|S|$  = total no. of instance in dataset

$V$  = set of distinct values of  $A$

$H(S_i)$  = entropy of subset of  $A$

$H(A, S)$  = entropy of an attribute  $A$ .

Eg) Entropy



$G_3 \uparrow$   
 has highest entropy  
 $G_2 \mid$   
 medium  
 $G_1 \downarrow$   
 has lowest entropy

→ High entropy (messy data)  
 Low entropy (clean)

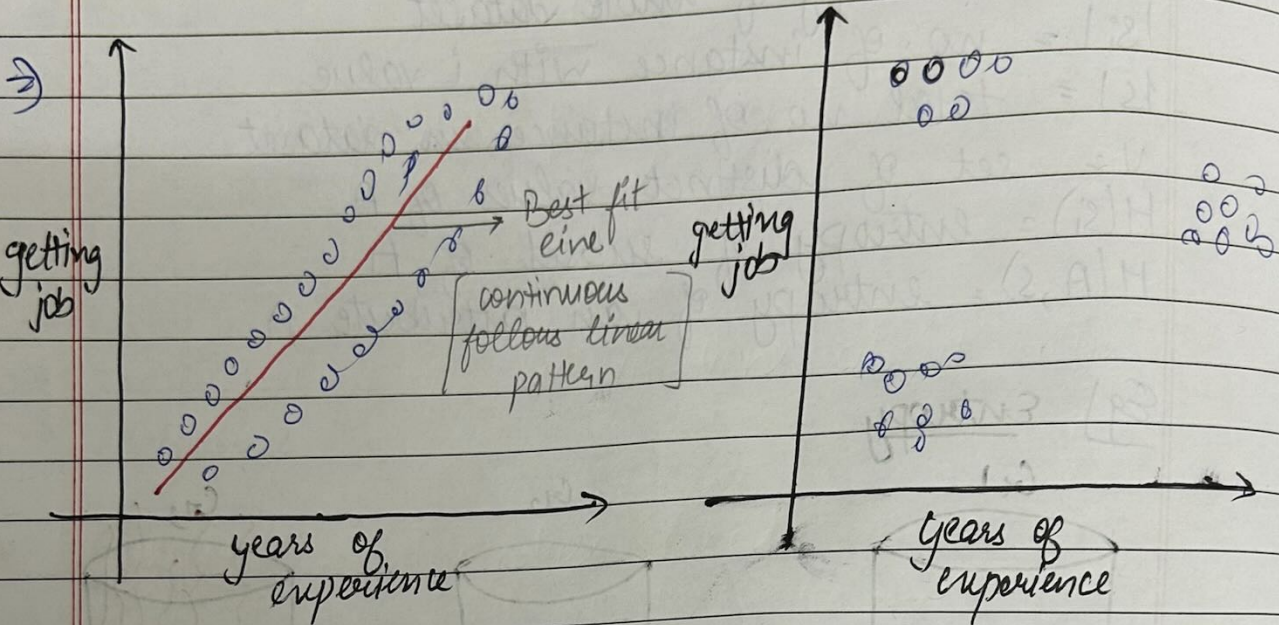
→ Low entropy - high information

\* Entropy  $\propto \frac{1}{\text{Information Gain}}$



\* Decision Tree is used for regression & classification

→ Decision Tree is used for non-linear data

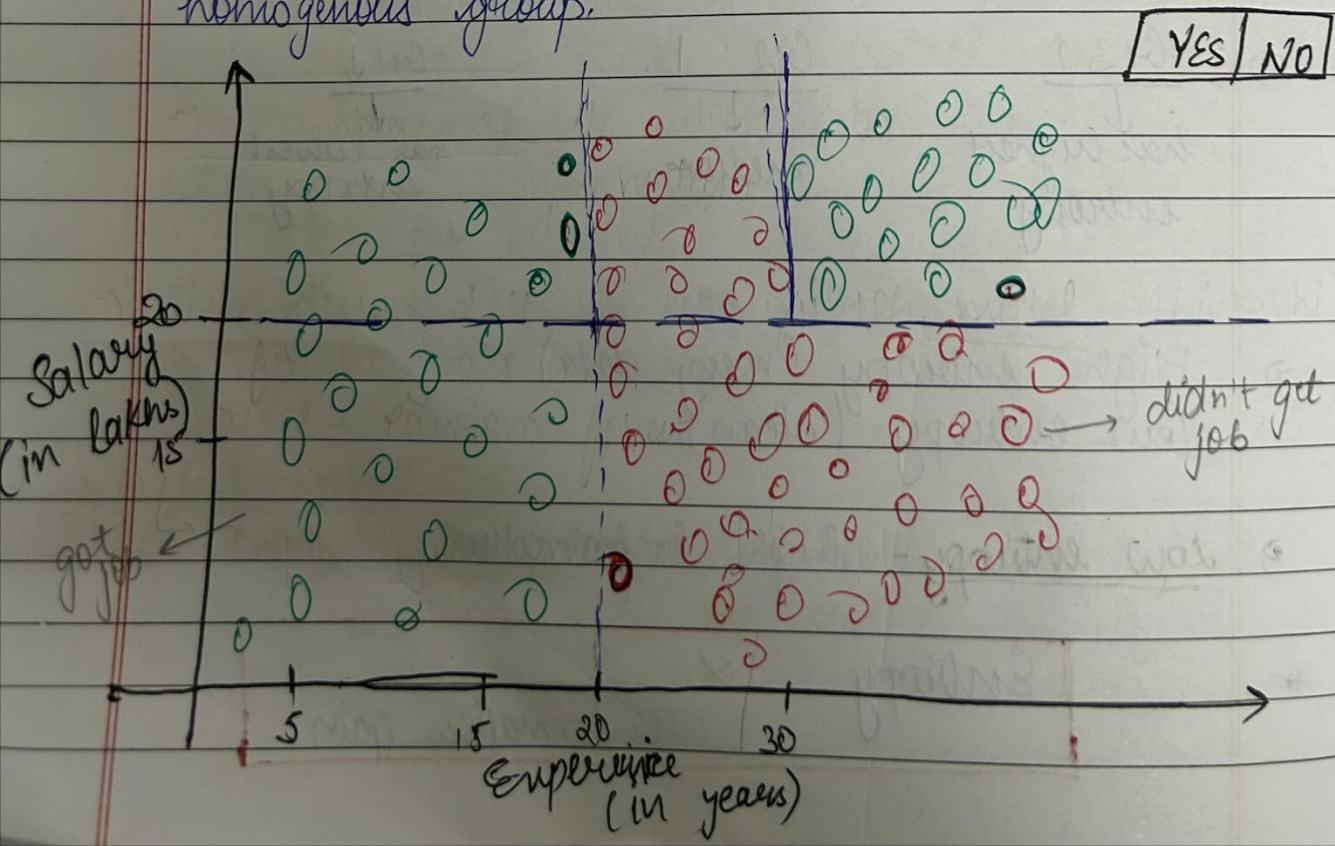


Linear data  
(uses linear regression)

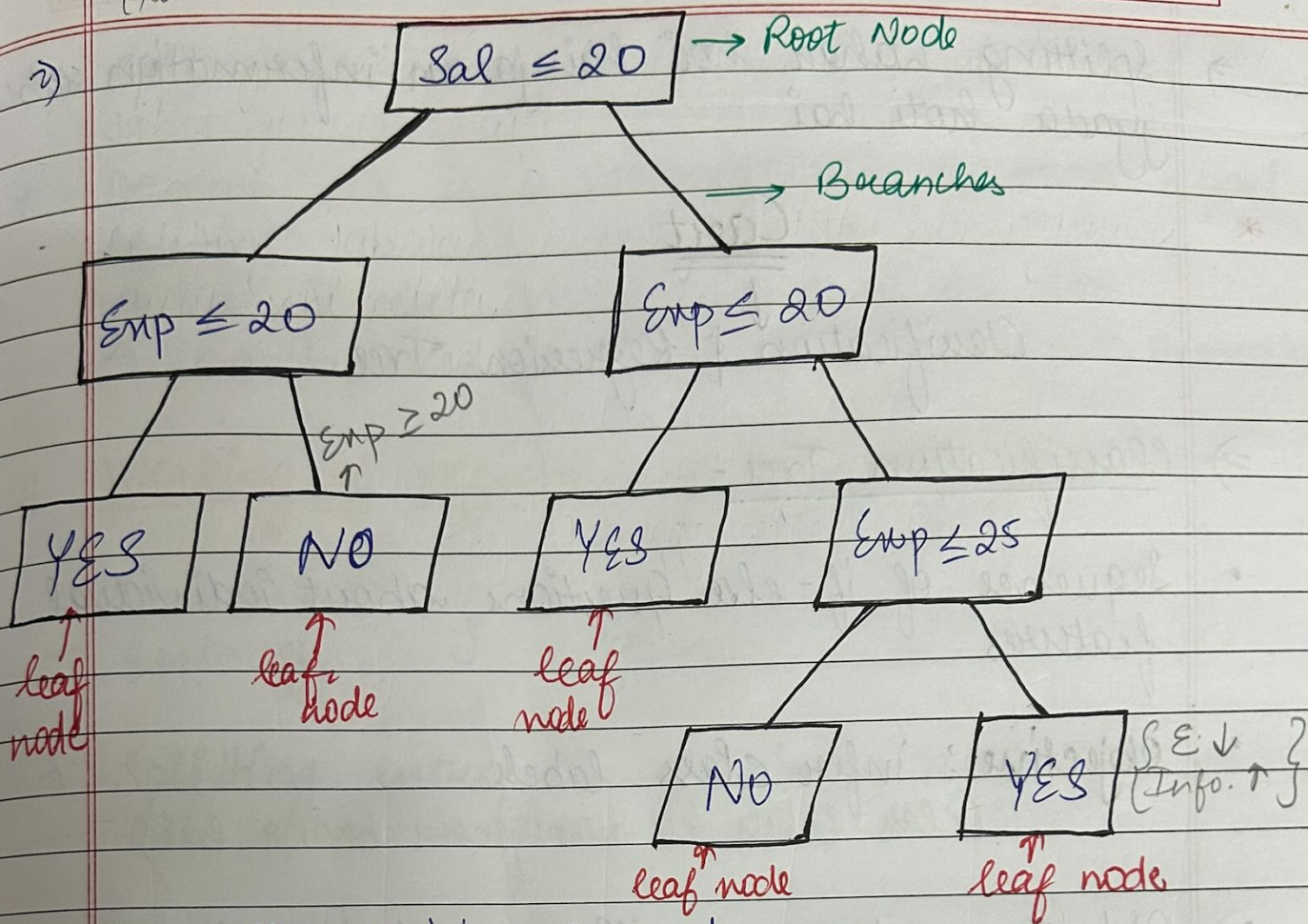
Non-linear data  
(uses decision-tree)

## # Concept of Decision Tree

To split data in the way that converts into homogenous group.







Now the decision is made.

Ques) How to verify the decision tree?

→ No. of regions = No. of leaf nodes  
then the decision tree is right

Ques) Predict if the person gets the job or not, if the person has salary = 15, emp = 30.

→ Salary = 15  
Emp = 30

Output = NO

→ Decision tree makes different boundaries.

→ Splitting tab tak hoti hai job tak homogeneous groups nhi milta.



→ Splitting wahan hoti hai jahan information sabse jyada hoti hai

\*

## Cost

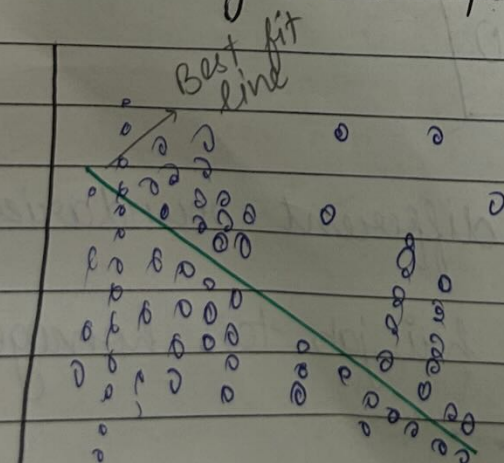
↓  
Classification & Regression Tree

### ⇒ Classification Tree -

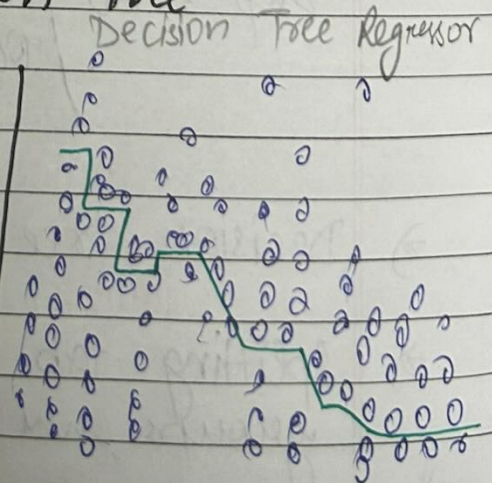
- Sequence of if-else questions about individual features
- Objective: infer class labels  
(class labels ko homogenous banana hai)
- Able to capture non-linear relationships b/w features & labels
- Don't require feature scaling  
(eg - standardization, normalization - etc.)

### # Regression Tree

Linear Regression v/s Regression Tree



↳ Gives more error



↳ Gives less error



\* leaf node is the answer

\* Decision Tree is a data structure to make model of data.

→ Uses if else at every node of the tree

→ can be used for both classification & regression.

## # Working of Algorithm of Decision Tree

→ Gini Index } Information  
→ Entropy } gain

→ Splitting entropy or gini index apne aap karta hai.

→ Splitting is done by

Gini Index      Entropy      ID3

decide where to split & where not to

\* Random Forest → Combination of Decision Trees.