



# P<sup>6</sup>OORNIMA

## LECTURE NOTES

Campus: PCE Course: BTECH Class/Section: Third CSE-A  
Name of Faculty: Ramesh K. Gade Name of Subject: Machine Learning  
Date (Prep.): 24-02-21 Date (Del.): 25-02-21 Unit No./Topic: 6

OBJECTIVE: To be written before starting the lecture (Pl. write in bullet points the main topics to be covered in the lecture)

Decision Tree and Random Forest.

### IMPORTANT & RELEVANT QUESTIONS:

what is entropy?

### FEED BACK QUESTIONS (AFTER 20 MINUTES):

How we determine the best split for a decision tree.

OUTCOME OF THE DELIVERED LECTURE: To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

good

REFERENCES: Text/Ref. Book with Page No. and relevant Internet Websites:

scikit with learn ebook



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

Decision Tree - Important Terms -

① Entropy - is the measure of randomness or predictability in the dataset.

eg.

② Information gain - it is the measure of decrease in entropy after the dataset is split.

③ Leaf-node - leaf node carries the classification or the decision.

④ Root node - Top most node.

We have to frame the condition in such a way that the information gain is the highest.

$$= \sum_{i=1}^n P(\text{value}_i) \cdot \log_2(P(\text{value}_i))$$

$$\frac{3}{8} \log_2 \frac{3}{8} + \frac{2}{8} \log_2 \frac{2}{8} + \frac{1}{8} \log_2 \frac{1}{8} + \frac{2}{8} \log_2 \left(\frac{4}{8}\right)$$

$$\text{Entropy} = 0.571$$



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

Random Forest Algor<sup>m</sup>

No overfitting - Use of multiple trees, reduces the risk of overfitting.

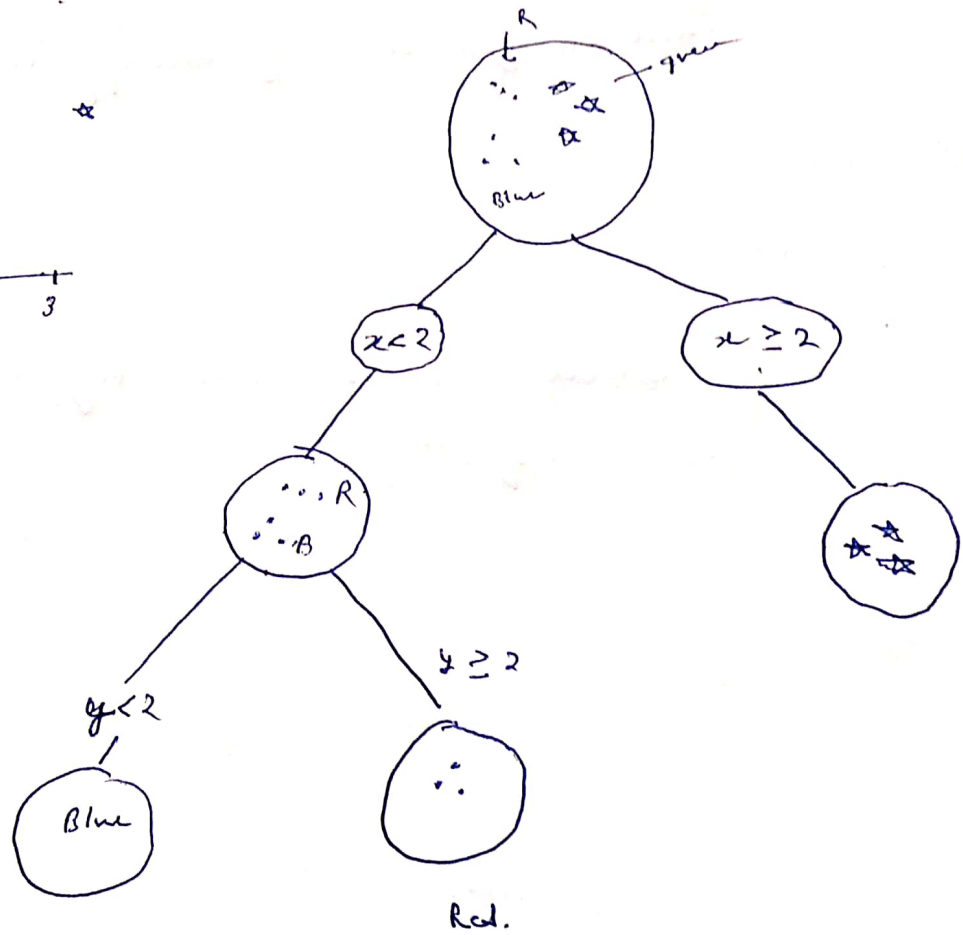
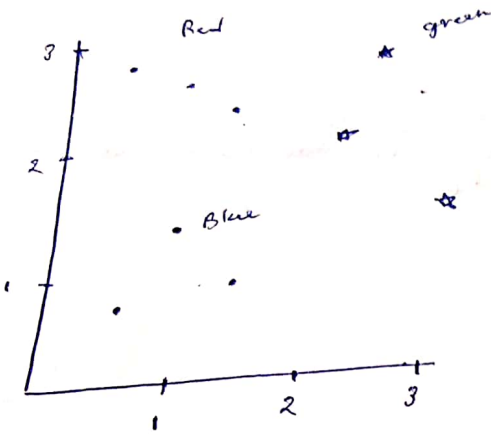
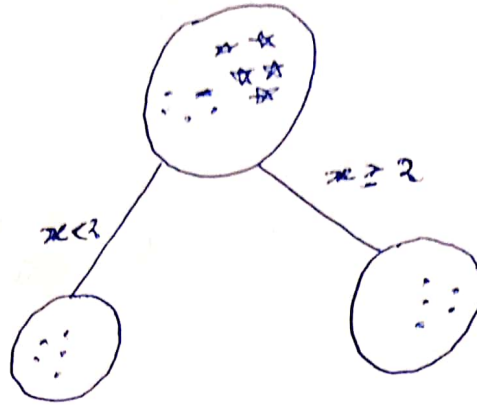
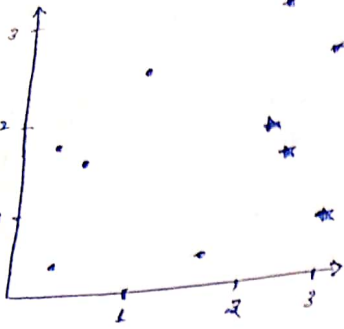
High accuracy - Runs efficiently on large database

Highly accurate predictions

Estimate missing data - RF can maintain accuracy when a large proportion of data is missing.

RF or R Decision Forest - is a method that operates by constructing multiple decision trees during training phase.

↳ The decision of the majority of the trees is chosen by the random forest as final decision.







# JOURNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

#### Decision Tree in Machine Learning:-

(Set of nested-if and if else if conditions)  
What is Decision Tree :- Decision tree is a tree shaped

diagram used to determine a course of action. Each branch of the tree represents a possible action/decision occurrence or reaction.

↳ it uses a tree like model to represent decision.  
↳ Decision Tree can be used to solve the classification problems. (For eg- Discriminating the certain type

of vegetables based on certain type of features.)  
↳ set of axis hyperparallel planes are as axis parallel in a DT.

Advantages-1. its simple to understand, interpret and visualize.

2. Little effort required for data preparation.
3. you can handle both numerical and categorical data.
4. Non-linear parameter don't effect its performance.

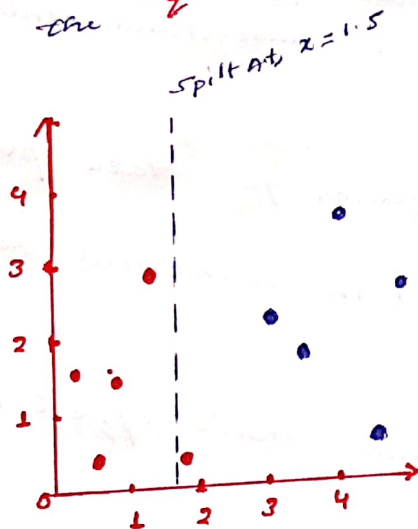
Disadvantages of Decision Tree:-

- 1). Overfitting occurs when the algorithm captures noise in the data.
- 2). The model can get unstable due to small variation in data.
- 3). A highly complicated decision tree tends to have a low bias which makes difficult for the model to work with new data.

Decision Tree - (Important Terms)

1. Entropy:- Entropy is the measure of Randomness or unpredictability in the dataset.

Initially the dataset is having higher entropy. Entropy is a metric used to train decision tree. These metric measures the "quality of a split".



Left branch has ..... (red) dots.  
Right branch has 1 (red) dot and 5 blue dots (.....).

(A data set having only blues would have very low or zero entropy).

(A data set having mixed blues, greens, and red (.....) would have relatively high entropy.)

Information Entropy for a dataset with  $c$  classes are denoted as:-

$$E = - \sum_i^K P_i \log_2 P_i$$

The no of data points /  
or  
 $c$

where  $P_i$  = is the probability of randomly picking  $i^{\text{th}}$  element of class  $c$ .

For eg - consider a dataset with 1 blue, 2 black, 3 red.

Then  $E = - (P_b \log_2 P_b + P_{\text{black}} \log_2 P_{\text{black}} + P_r \log_2 P_r)$

Here  $P_b = \frac{1}{6}$        $P_{\text{black}} = \frac{2}{6}$        $P_r = \frac{3}{6}$

So  $E = - (\frac{1}{6} \log_2 \frac{1}{6}) + \frac{2}{6} \log_2 (\frac{2}{6}) + \frac{3}{6} \log_2 (\frac{3}{6})$   
 $= 1.46$

eg-2. Consider 3 blue ... Entropy would be  
 $E = - (1 \log 1) = 0$

So before split - we have 5 blues and 5 red. So  
 Entropy was  $E_{\text{before}} = - (0.5 \log_2 0.5 + 0.5 \log_2 0.5)$   
 $= 1$

After split, we have two branches - Left branch - 4 red and Right branch - 1 red and 5 blue.  
 So Entropy after split.

$E_{\text{left}} = - (0.5 \log 0.5) = 0$   
 $E_{\text{right}} = - (\frac{1}{6} \log_2 (\frac{1}{6})) + \frac{5}{6} \log_2 (\frac{5}{6})$   
 $= 0.65$





# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

$$E_{\text{split}} = 0.4 * 0 + 0.6 * 0.65$$
$$= 0.39$$

we started with  $E_{\text{before}} = 1$  entropy before the split and now are down to 0.39.

$$\text{So Information Gain} = 1 - 0.39 = 0.61$$

higher Information Gain = More Entropy removed.

**Information Gain:-** It is the measure of decrease in entropy after the dataset is split. Information gain is calculated for a split by subtracting weighted entropies of each branch from the original entropy.

**LEAF NODE:-** Leaf node carries the classification.

**ROOT NODE:-** The top most decision node is known as the Root node.

**Note-** we have to frame the conditions such that it split the data in such way that the information gain is the highest.

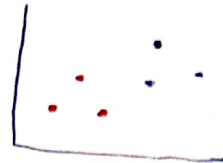
when training a decision tree using these metrics, the best split is chosen by maximizing the Information Gain.



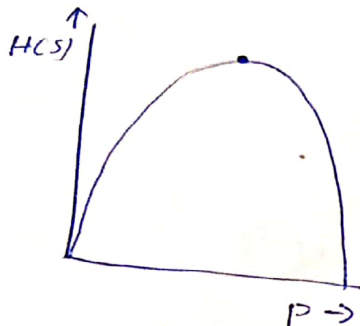
Build a Decision Tree

$$H(S) = - \sum_{c=1}^C P_c \log P_c$$

$P_c$  is prob. of class  $c$ .



$$H(S) = - \left( \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} \right) = 1$$



$P=0.5$  highest entropy.



$$H(S) = 1 \log 1 = 0$$

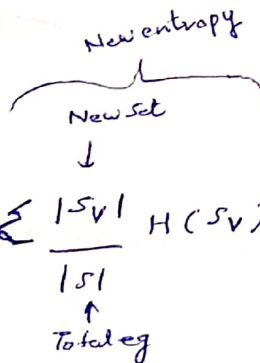
classess  $\begin{cases} \text{yes} & 9 \\ \text{No} & 5 \end{cases}$

$$H(S) = - \left( \frac{9}{14} \log \frac{9}{14} + \frac{5}{14} \log \frac{5}{14} \right) = 0.94$$

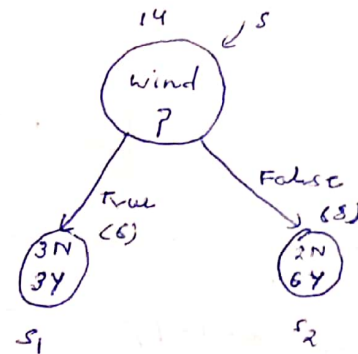
CAH)

$$IG(S, A) = H(S) - \sum \frac{|S_v|}{|S|} H(S_v)$$

$\downarrow$   
 (old entropy)



(Ratio of no of eg in new set over all sets we have)



Maximize the information gain.

$$H(S) = - \sum \frac{8}{14} \left[ -\frac{6}{8} \log \frac{6}{8} - \frac{2}{8} \log \frac{2}{8} \right] + \frac{6}{14} \left[ -\frac{3}{6} \log \frac{3}{6} - \frac{3}{6} \log \frac{3}{6} \right]$$

$$= H(S) - \frac{8}{14} (0.81) + \frac{6}{14}$$

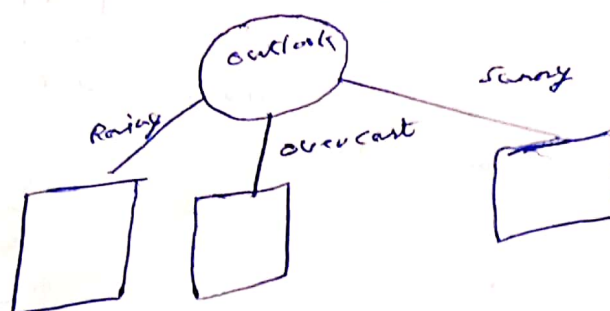
$$= 0.94 - 0.892$$

$$= 0.0048 \leftarrow$$

$$\begin{aligned}
 IG(S, \text{wind}) &= 0.048 \\
 IG(S, \text{outlook}) &= 0.247 \\
 IG(S, \text{Humidity}) &= 0.15 \\
 IG(S, \text{Temp}) &= 0.029
 \end{aligned}$$

Do yourself.

So Here we choose  $IG(S, \text{outlook})$



Repeat this process for other attributes recursively.  
until entropy is zero..

=