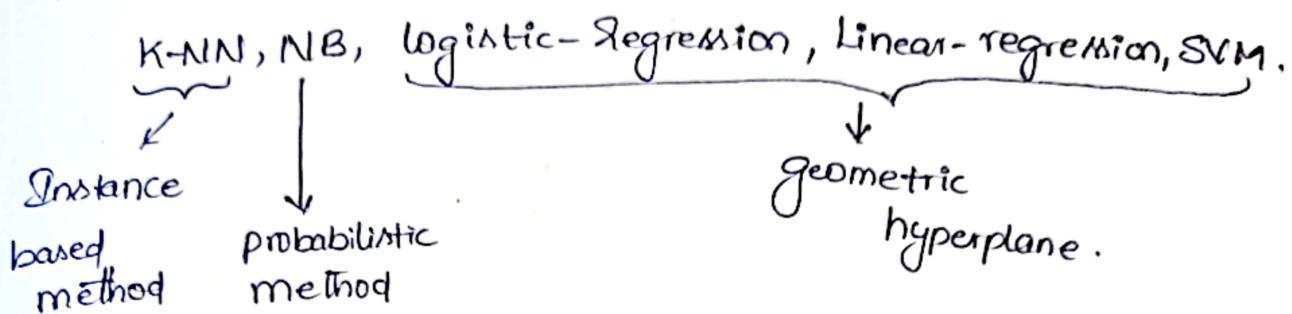


Decision Trees  
Axis parallel

34.1 Geometric Intuition of decision tree: Axis parallel hyperplanes :-

→ As we seen previously about



→ Decision tree is simple if else classifier

# As we seen in IRIS dataset

$$y_i = \{1, 2, 3\} \quad (\text{SL, SW, PL, PW})$$

$$x_i = \langle \text{SL}, \text{PL}, \text{SW}, \text{PW} \rangle$$

if  $\text{PL} < a$

$$y_i = 1$$

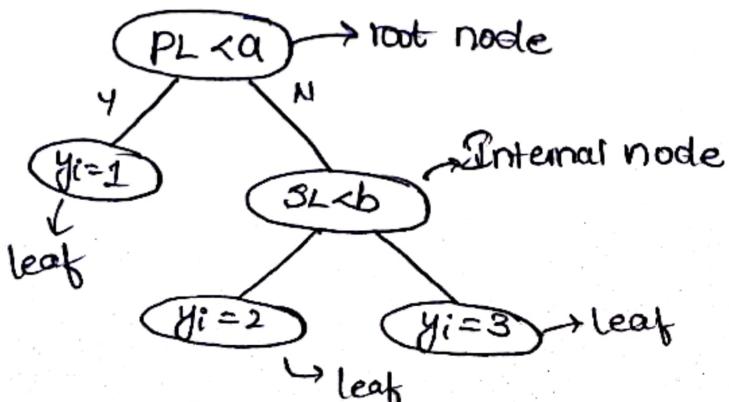
else if  $\text{SL} < b$

$$\text{Class} = 2$$

else

$$\text{Class} = 3$$

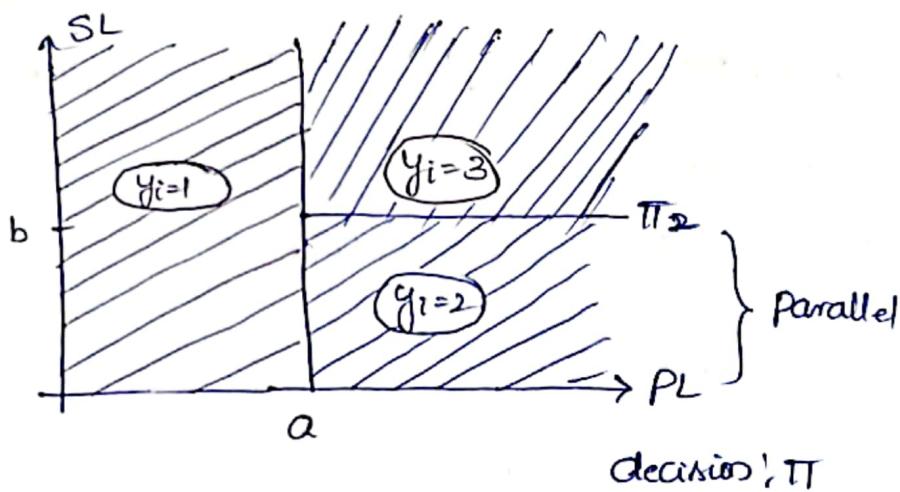
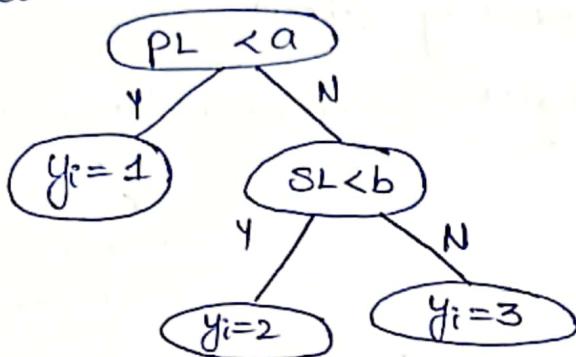
So diagrammatically we can do like



→ Usually we make decisions @ root node and @ Internal nodes.

→ Geometric Intuition :-

If we consider



→ Geometrically Decision trees are the set of axis parallel hyperplane

→ Axis parallel usually can be the parallel lines which are parallel to any axis (either x or y)

## 3.1.2 Sample Decision Tree :-

← Play tennis Example →

## 3.1.3 Building a decision tree : - Entropy

# Entropy: is a metric to measure the impurity in given attribute

$$H(Y) = - \sum_{i=1}^K P(y_i) \log_b(P(y_i))$$

; where b can be  
 $b=2$

or Other no.  $b=2.718$

Note :-  $\log_2 = \lg$

$\log_e = \ln$

&  $P(y_i) = P(Y=y_i)$

# If we see play tennis example

$Y$ : play tennis

$$Y+, Y- \Rightarrow P(Y+) = 9/14$$

$$P(Y-) = 1 - P(Y+) = 5/14$$

$$H(Y) = - \sum_{i=1}^K P(y_i) \log_2(P(y_i))$$

$$H(Y) = - \underbrace{\frac{9}{14} \log_2(\frac{9}{14})}_{\downarrow P(Y+)} - \underbrace{\frac{5}{14} \log_2(\frac{5}{14})}_{\rightarrow P(Y-)} = 0.94$$

No. of +ve pts  
Total no. of pts = % of +ve pts in D

best-split ( $X, Y$ )

# Entropy helps to select the right feature/root node among all the features to construct a decision tree

| OK entropy value  $< 1$  |



## Properties of Entropy

- To understand this let us assume we have a class dataset or a category dataset.

$$\gamma \rightarrow \gamma^+, \gamma^-$$

### case1 {

$$D \left. \begin{array}{l} Y_1 \rightarrow 99\% \\ Y_0 \rightarrow 1\% \end{array} \right\} H(Y) = -0.99 \log 0.99 - 0.01 \log 0.01 = 0.0801$$

### Case 2

$$\text{② } \begin{cases} y_+ \Rightarrow 50\% \\ y_- \Rightarrow 50\% \end{cases} + I(Y) = -0.5 \log(0.5) - 0.5 \log(0.5) = 1$$

se3

$$D \xrightarrow{\quad} Y^+ \rightarrow 0\% \quad \left. \begin{array}{l} \\ \end{array} \right\} H(Y) = 0$$

$$D \xrightarrow{\quad} Y^- \rightarrow 100\%.$$

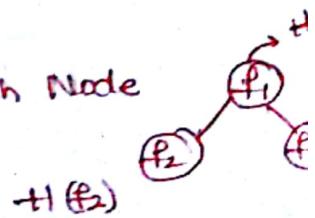
→ from this we can say if both the classes are equally probable then we get maximum entropy

i.e. 1. (Impure Subsplit)

→ If one is fully dominant we can say that Entropy is minimum ie zero. (Pure Subsplit)

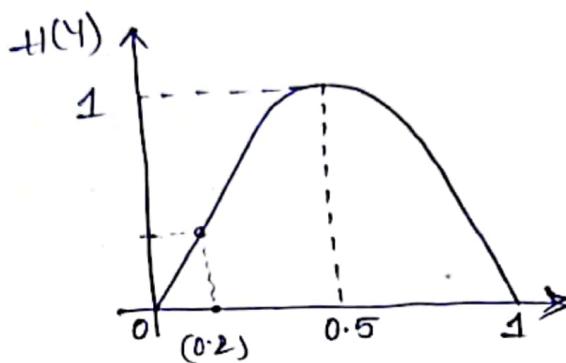
→ As we go from maximum ↓ our Entropy tends to decrease

\* Entropy will be calculated for each Node



$$\star P(Y_t) = 1 - P(Y_i)$$

Graphically we can define as



If  $y_t \rightarrow 0.2$  that means  
 $y \rightarrow 0.8$ .

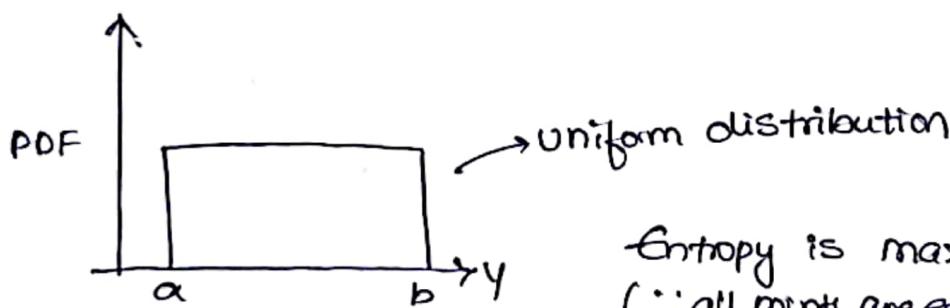
\* By if we have a n class R.V

$$Y_1 = \underbrace{y_1, y_2, \dots, y_k}$$

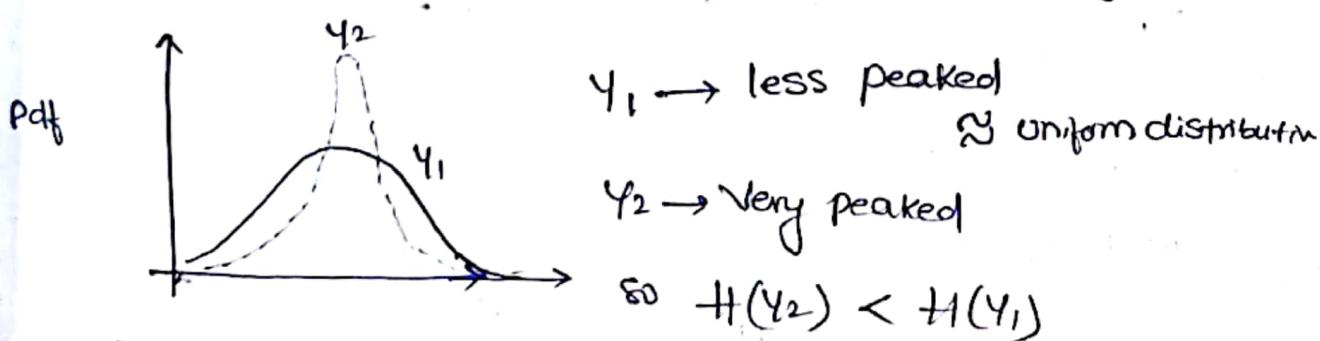
If all of them are equally probable  
Entropy is maximum.

If  $y_1 \rightarrow$  most probable      } Entropy is minimum,  
 $y_2, y_3, \dots, y_k \rightarrow 0$       }

\* If we consider PDF



Entropy is maximum  
( $\because$  all points are equally probable)



$y_1 \rightarrow$  less peaked

$\approx$  uniform distribution

$y_2 \rightarrow$  very peaked

$$\therefore H(y_2) < H(y_1)$$

→ More Peaked distribution is. less is its Entropy

$$D_{KL}(P||Q) = \sum_x P(x) \log \left( \frac{P(x)}{Q(x)} \right)$$

$$= \sum_x P(x) \log P(x) - \sum_x P(x) \log Q(x)$$

- Note:- KS-statistic will be used & Generally in Statistics and KL-divergence much better than KS-Statistic In machine learning.

### 34.5 Building a decision Tree : Information Gain:

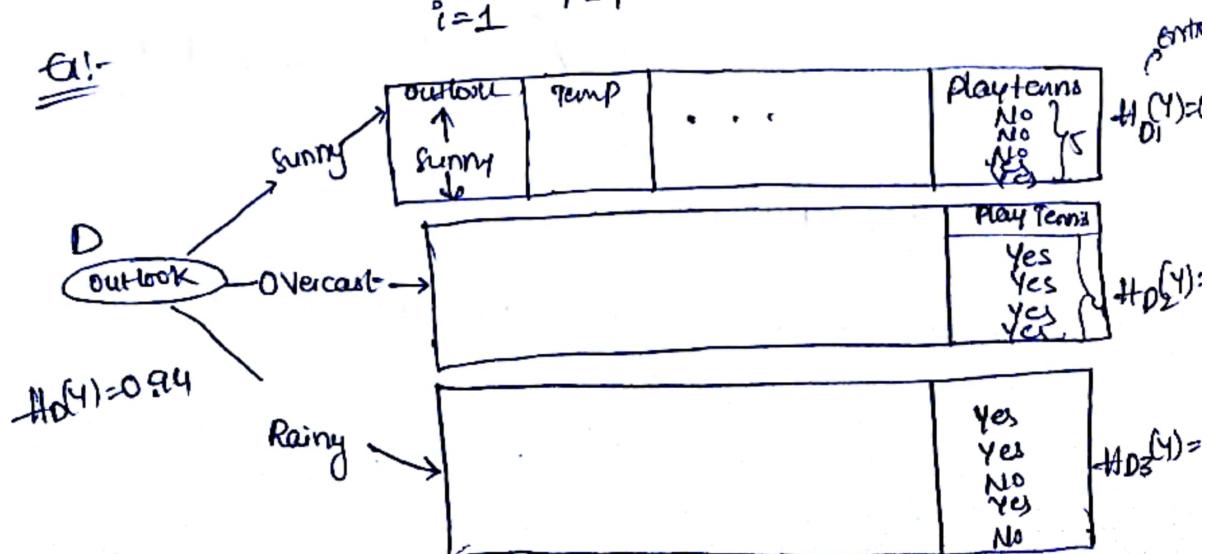
# for a Random Variable  $Y$  if that Random Variable divided into small sub set.

$$Y \xrightarrow{\text{Variation}} Y_1, Y_2, \dots, Y_k$$

or considering a dataset  $D$  if that Dataset is divided into  $k$  Data sub sets

$$D \rightarrow D_1, D_2, \dots, D_k$$

$$IG(Y, \text{Var}) = \sum_{i=1}^k \frac{|D_i|}{|D|} * H_{D_i}(Y) - H_D(Y)$$



# simply if a Dataset  $D$  is given, we need to divide the Datasets into  $n$  Sub Datasets using a feature. ( $D_1, D_2, D_3$ )

- After that find Entropy of the Datasets, (that is for the Dataset  $D$  and for DataSub Set  $D_1, D_2, D_3$ ).
- Then we can get Information gain as

Information Gain = Entropy (Parent) - [Weighted Average Entropy of child nodes].

where

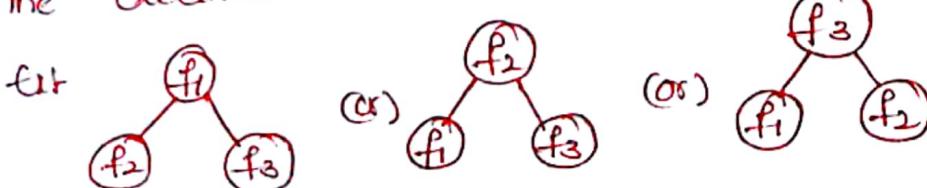
$$\text{Weighted Average Entropy of Child nodes} = \frac{|D_1|}{|D|} \times H_{D_1}(Y) + \frac{|D_2|}{|D|} \times H_{D_2}(Y) + \frac{|D_3|}{|D|} \times H_{D_3}(Y)$$

$$\text{Entropy of parent} = H_D(Y)$$

- for the above example

$$0.97 - \frac{5}{14} \times 0.97 + \frac{4}{14} \times 0 + \frac{5}{14} \times 0.97$$

Note: Information gain helps us to select the pattern how the decision trees to be



The higher the Information gain that particular model / structure need to be considered.

### 34-B Building a decision Tree : Gini Impurity:-

# Gini Impurity is similar to Entropy.

$$I_G(Y) = 1 - \sum_{i=1}^K [P(Y_i)]^2 \quad Y \rightarrow y_1, y_2, y_3, \dots, y_K$$

# Compared to Entropy Gini Impurity is less computation complex and we have less Time complexity as Entropy uses logarithms.

$$I_G(Y)$$

$$1 - P(Y+)^2 - P(Y-)^2$$

$$H(Y)$$

$$- P(Y+) \log_2 P(Y+)$$

$$- P(Y-) \log_2 (P(Y-))$$

# If we consider the difference b/w Entropy & Gini Impurity

let  $Y \begin{cases} \uparrow & Y+ \\ \downarrow & Y- \end{cases}$

Case 1:-  $P(Y+) = 0.5$

$$P(Y-) = 0.5$$

$$I_G(Y) = 1 - (0.25 + 0.25) = 0.5$$

$$H(Y) = 1$$

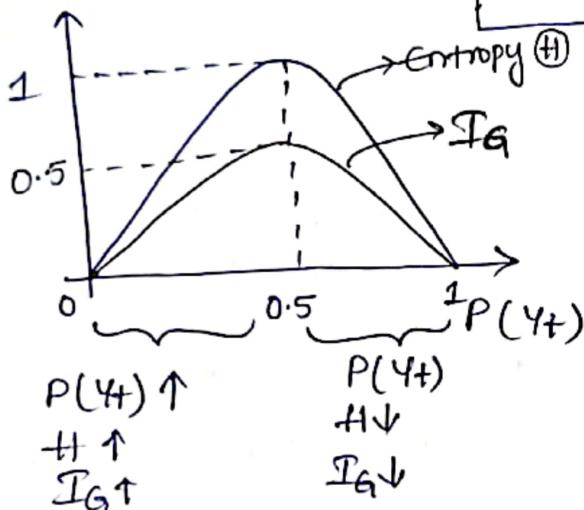
Case 2  $P(Y+) = 1$

$$P(Y-) = 0$$

$$I_G(Y) = 1 - (1+0) = 0 \quad H(Y) = 0$$

# Geometrically.

2-Category Case :-  $y_+$ ,  $y_-$



Gini Impurity (Programmatically):  

$$1 - \sum((\text{np.sum}(y == c) / m) * 2 \text{ for } c \text{ in np.unique}(y))$$



Note:-

Gini Impurity Notation =  $I_G(Y)$

Information Notation =  $I_G(Y)$

347 Building a Decision Tree: Constructing a Decision Tree:-

# Considering the example of Decision tree. play. tennis

Dataset.

# To construct a Decision tree.

Input:-  
 $x, y, \text{max-depth}, \text{min-sample-split} = 2$   
 $\text{min-samples-leaf} = 1$

Step 1 :-

1. Considering any feature in the dataset and find

its Entropy

Outlook

$\Theta \rightarrow 9+, 5-$

$$H_D(Y) = 0.94$$

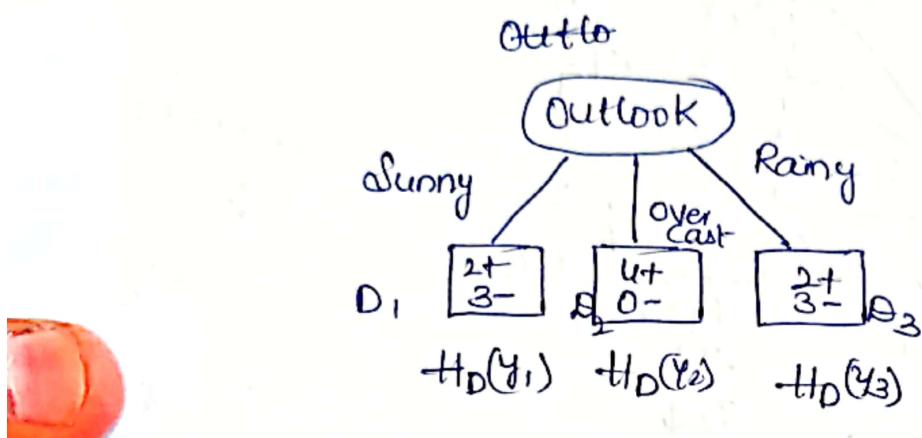
$$\text{num-samples-per-class} =$$

$$\sum_i P(y_i)$$

$$\text{Predicted class} = \arg \max [P(y_i)]$$

### Step 2

Now Divide The feature Data into 'n' sub data using internal features and find their Entropies,



and after that find weighted Entropy  
in This Example = 0.69

### Step 3

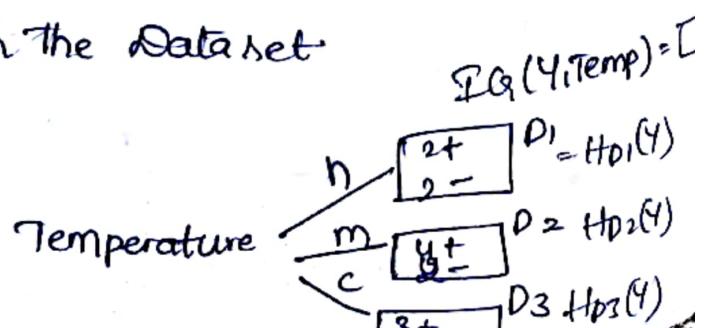
find IG (Y, Outlook) [that means find the Information Gain by using the feature1] using  $\sum H_D(Y_i)$  Weighted Entropy.

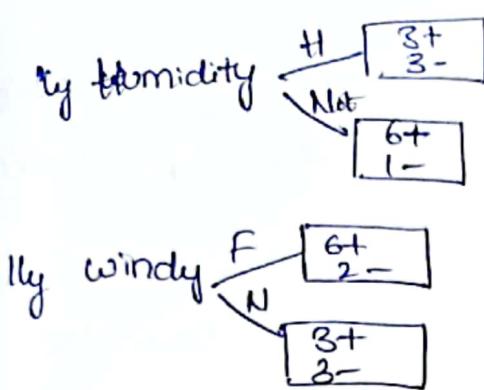
$$\text{In this Example } IG(Y, \text{Outlook}) = 0.94 - 0.69 = 0.25$$

### Step 4

Repeat The Same above steps and find the IG for every feature in the dataset

In This Example





Now once we get the all the IG values

$$IG(Y, \text{outlook}) = \dots$$

$$IG(Y, \text{Temp}) = \dots$$

$$IG(Y, \text{Humidity}) = \dots$$

$$IG(Y, \text{windy}) = \dots$$

using  $IG(Y, f) = H_p(Y) - \sum_{i=1}^K \frac{|D_i|}{|D|} \cdot H_{D_i}(Y)$

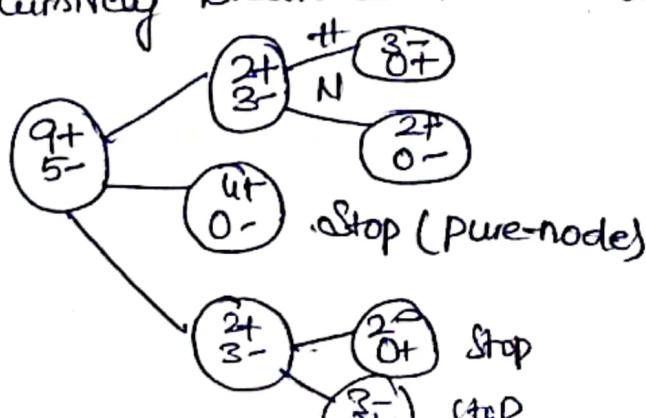
Choose the feature which has highest value and make that as parent node.

$$IG(Y, f) = \text{entropy}@ \text{parent level}$$

- weighted entropy @ child level.

### Step 5:

Recursively break each node using IG as the Criterion



The node which has maximum dominance by the class points in the feature we call that node as pure node and we need to stop @ that node.

# and we also can't grow the tree if we have lack of points.

for example if we have very points @ a node



$n=10K \text{ pts}$

as there are only 2 points  
which are very fewer

∴ They are very fewer and  
If they are noisy points  
we end up with wrong,

# So typically if we got any of the scenario's while constructing a decision tree

- ① pure node ✓ num-samples\_per\_class · Count(o) == 0
- ② few pts @ a node ✓ if depth < maxdepth  
and
- ③ If we are too deep ✓

We need to stop constructing further tree.

\* depth of the tree ↑ : Overfitting ↑ (few pts)

depth is small  $\Rightarrow$  underfit.

# In decision tree hyperparameter is depth and which can findout by Cross Validation (CV).

### 34.8 Building a Decision Tree: Splitting Numerical Features :-

# Construct a Decision tree:-

for the above example we have used Information gain for splitting the node.

# for finding Information Gain we can use Entropy

Instead of Entropy we can also use Gini Impurity which is Computational efficient

IG!  $\rightarrow$  Entropy  
 $\rightarrow$  Gini Impurity.

# till now we have applied Decision tree for the Data which has Categorical features / Discrete RV.

# If we have Numerical featured data.

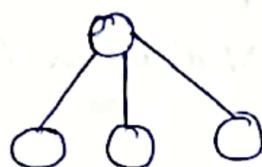
$x_1$	$y$
2.2	1
2.6	1
3.5	0
3.8	0
4.6	1
5.3	0

where  $x_1$ : numerical

$\hookrightarrow$  Integer or real valued

→ Previously we tried to split the data based on Categorical Variable like If we have feature  $f_2$ : 3 categories.

Then we can split into



And If we have numerical data then for every value,



There would be only one node. Which is not prefer.

General way  
→ Hack to Split based on numerical features!

Step 1: Sort the Numerical feature in ascending order

$f_1$	$y$
2.2	1
2.6	1
3.5	0
3.8	0
4.6	1
5.3	0

Step 2 Consider the features based on range

$$f_1 < 2.2$$

$$f_1 < 2.6$$

$$f_1 < 3.5$$

$$f_1 < 3.8$$

$$f_1 < 4.6$$

$$f_1 < 5.3$$

Like this we get n possible variations.

→ Based on Threshold we divide Then features

$f_1 < \gamma_1 \leftarrow D_1$   
 $f_1 < \gamma_1 \leftarrow D_2$   
 $f_2 < \gamma_2 \leftarrow D_3$   
 $f_2 < \gamma_2 \leftarrow D_2$   
 $f_3 < \gamma_3 \leftarrow D_1$   
 $f_3 < \gamma_3 \leftarrow D_2$   
⋮  
⋮  
 $f_n < \gamma_n \leftarrow D_1$   
 $f_n < \gamma_n \leftarrow D_2$

(252)

numerical		Y	(a)
$f_1$	$\gamma$		
0.2	1	D <sub>1</sub>	D <sub>1</sub>
0.6	1		D <sub>1</sub>
3.5	0	D <sub>2</sub>	D <sub>2</sub>
3.8	0		D <sub>2</sub>
4.6	1	D <sub>2</sub>	D <sub>2</sub>
5.3	0		D <sub>2</sub>

like this we n possible splitting based on Criteria and  
then we finalize nodes based on Information Gain  
Value, the feature which has highest IG will  
be considered to be the node

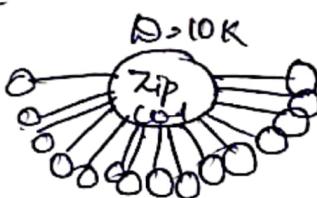
But This General process is the time taking process  
so we have some hacks to work around.

34.10 Building a decision tree: Categorical features with many possible values:-

# let us consider we have a feature Pincode / Zipcode with 1000's of values.

Zip code	1
=	
=	
=	
=	

If we have like this we end up having thousands of node under a node



So we end up having thousands of node @ a level.

# So to overcome this

One of the hack is to do feature engineering the categorical data and convert it into Numerical data.

Pincode	$Y_i = \{0, 1\}$
P <sub>1</sub>	
P <sub>2</sub>	
P <sub>3</sub>	
P <sub>4</sub>	
P <sub>1</sub>	
P <sub>2</sub>	
P <sub>2</sub>	

To convert it into numerical feature

$$P(Y_i=1 | P_j)$$

This means no of time we get  $Y_i=1$  for a particular Pincode ( $P_j$ ) (let be  $P_2$ ) / total no of time  $P_j$  occurred.

### 34.9 Feature Standardization :-

# In the case of logistic regression, SVM, etc. NN as all are distance based algorithms we need to perform feature standardization to make all the values under a feature into same scale.

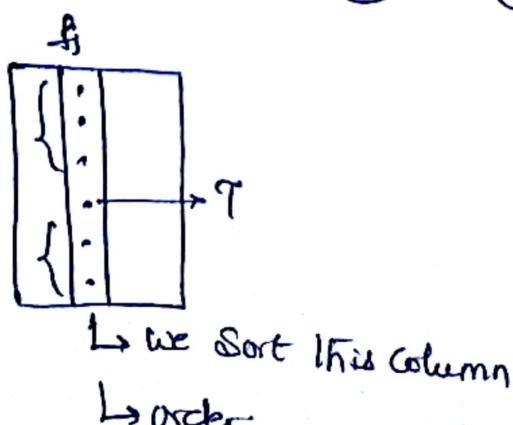
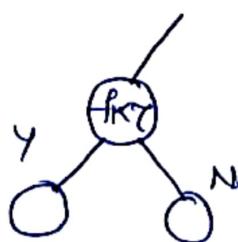
In logistic regression

$$x_{ij}' = \frac{x_{ij} - \bar{x}_j}{\sigma_j}$$

The diagram shows a vertical column of the matrix  $f_j$  with entries  $x_{ij}$ . An arrow points from the original entry  $x_{ij}$  down to the mean  $\bar{x}_j$  and standard deviation  $\sigma_j$ , indicating the transformation to the standardized value  $x_{ij}'$ .

# whereas Decision trees are not distance based methods.

Based on Threshold value we define a node and split it based on  $<$ ,  $>$ ,  $=$  values



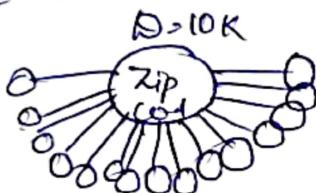
So we do not need to perform feature standardization, as it does not depend on scale rather it only considers the value not on distances that are located on scales.

34.10 Building a decision tree: Categorical features with many possible values:-

# let us consider we have a feature Pincode / Zipcode with 1000's of values.

Zip code	Y
=	
=	
=	
=	

If we have like this we end up having thousands of node under a node



So we end up having thousands of node @ a level.

# How to overcome this

One of the hack is to do feature engineering the categorical data and convert it into Numerical data.

Pincode	$y_i = \{0, 1\}$
P <sub>1</sub>	
P <sub>2</sub>	
P <sub>3</sub>	
P <sub>4</sub>	
P <sub>1</sub>	
P <sub>2</sub>	
P <sub>3</sub>	

To convert it into numerical feature

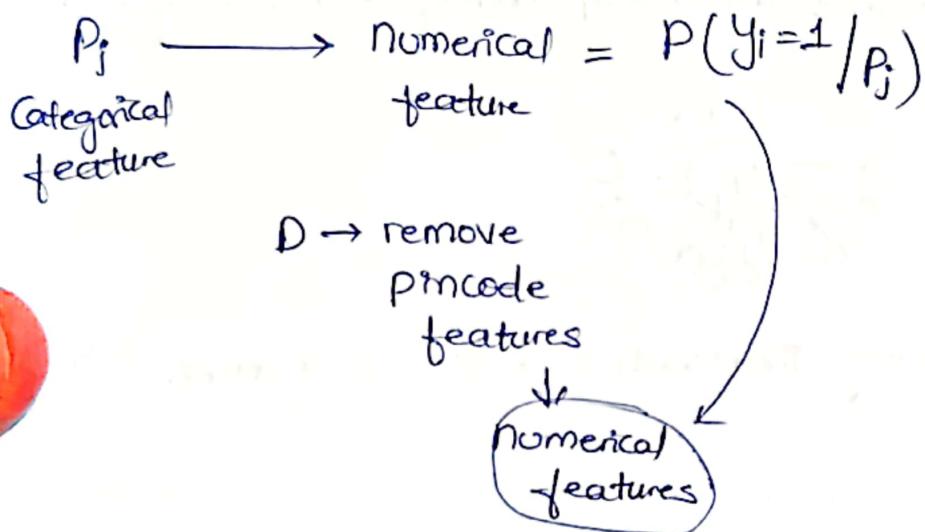
$$P(y_i=1 | P_j)$$

This means no of time we get  $y_i=1$  for a particular Pincode ( $P_j$ ) (let be  $P_2$ ) / total no of time  $P_j$  occurred.

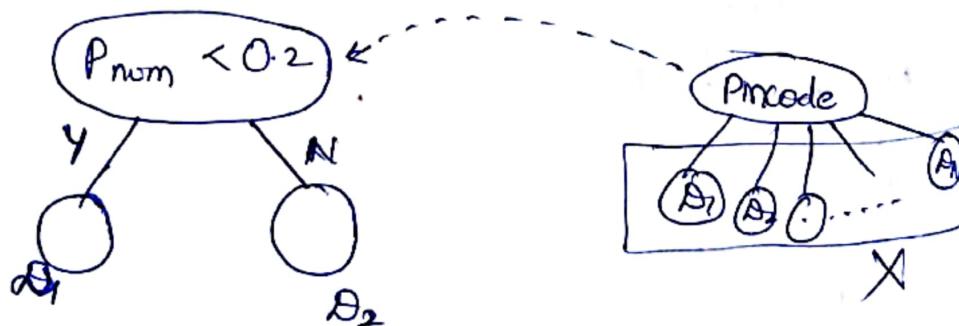
$$= P(Y_i=1 | P_j)$$

$$= \frac{\text{no. of } Y_i=1 \text{ for a given } P_j}{\# P_j}$$

let we have 20-times  $P_j$  and in that 19-times  
 $Y_i=1 \Rightarrow 19/20$



so by this we can make decision tree like



## Decision trees

- \* It is a supervised learning technique. Can be used for both classification and regression.

### Internal nodes

- \* Pruning is the process of removing the unwanted branches from the trees.

- \* To select the best attributes for the Root node and for the sub-trees there is a technique called Attribute Selection measure (ASM).

They are two popular techniques for ASM

- Information Gain
- Gini Impurity

### Information Gain

- Information gain is the measurement of changes in entropy after the segmentation of dataset as on attribute.
- It calculates how much information a feature provides us about a class.
- According to the value of Information gain we split the node and build the Decision tree.

Decision tree algorithm always tries to maximize the value of Information gain and a node / attribute

having the highest information Gain is split first

Information Gain:- Entropy(S) -  $[(\text{weightedavg} \times \text{Ent. of child nodes})]$

Entropy is a metric to measure the Impurity in an attribute. It specifies randomness in data.

Entropy

\* pruning: Getting an Optimal decision tree.

Pruning is a process of deleting the unnecessary nodes from a tree in order to get the optimal decision tree.

→ usually a too-large tree increases the risk of overfitting. And a small tree may not capture all the important features of dataset. Therefore a technique that decreases the size of the learning tree without reducing accuracy is known as pruning.

There are mainly two types of tree pruning techniques used:

- Cost complexity Pruning.
- Reduced Error Pruning.