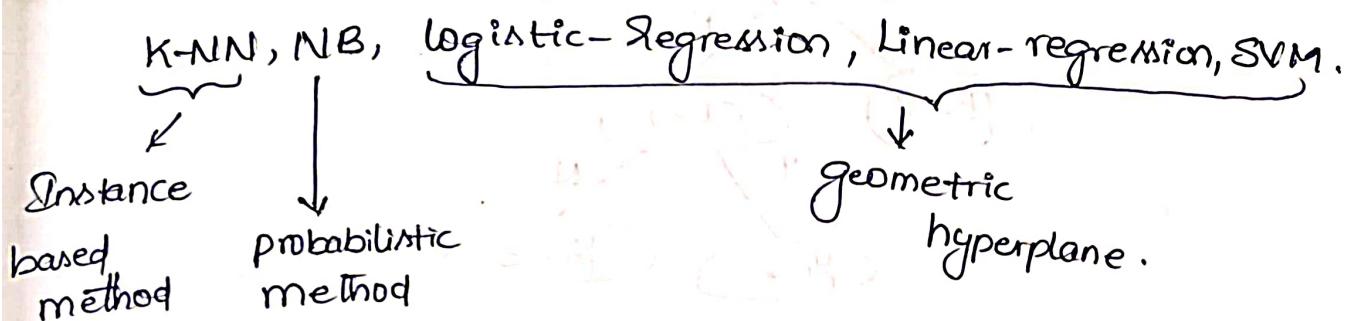


Decision Trees

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 $\text{PL} < a$

$$y_i = 1$$

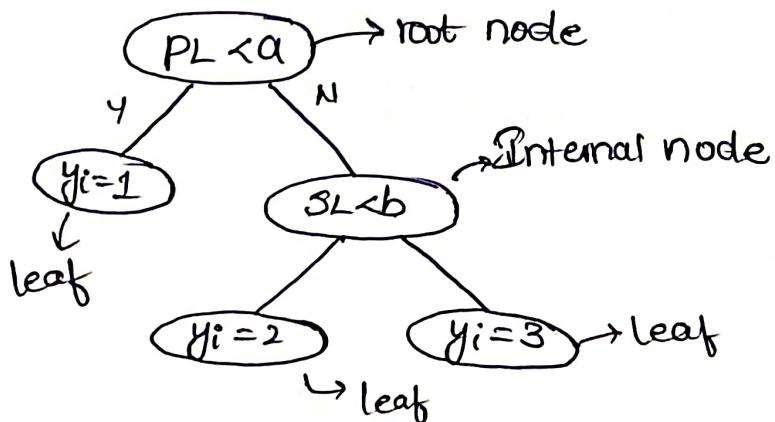
else if $\text{SL} < b$

$$\text{Class} = 2$$

else

$$\text{Class} = 3$$

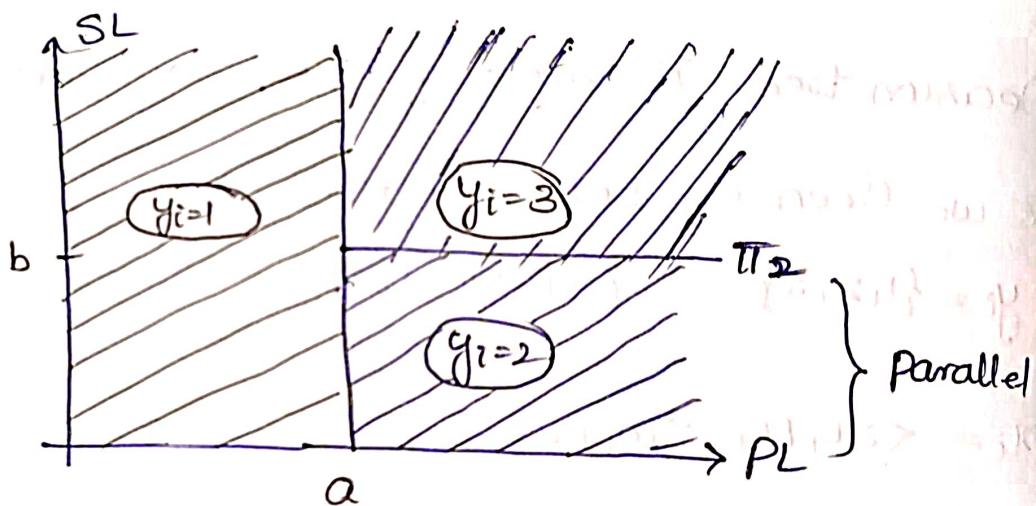
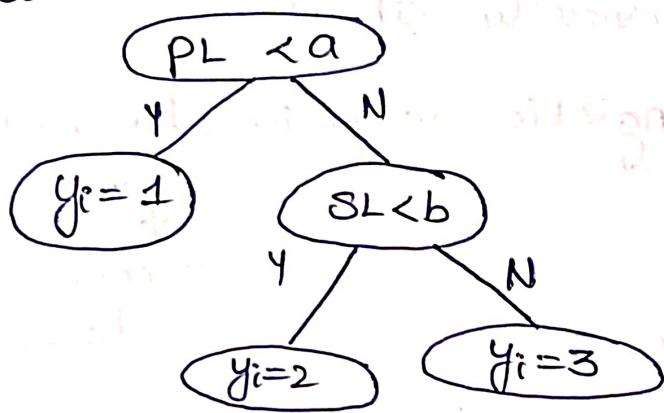
So diagrammatically we can do like



→ 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)

34.2 Sample Decision Tree :-

← Play tennis Example →

34.3 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 Order no. b=2.718

Note :- $\log_2 = \log_{10}/\log_e$

$$\log_e = \ln$$

& $P(Y_i) = P(Y = Y_i)$

If we see play tennis example

Y_1 : Play tennis

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

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

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

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

$$P(Y_+) = \text{No of +ve pts} / \text{Total no of pts}$$

$$\frac{\text{No of +ve pts}}{\text{Total no of pts}} = \text{No of +ve pts in D}$$

best-split (X, Y)

It helps 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 classification Dataset or a Category Dataset.

$$Y \rightarrow Y_+, Y_-$$

Case 1:

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

Case 2:

$$\text{D} \begin{cases} Y_+ \rightarrow 50\% \\ Y_- \rightarrow 50\% \end{cases} \left. \begin{array}{l} H(Y) = -0.5 \log(0.5) - 0.5 \log(0.5) \\ = 1 \end{array} \right.$$

Case 3:

$$\text{D} \begin{cases} Y_+ \rightarrow 0\% \\ Y_- \rightarrow 100\% \end{cases} \left. \begin{array}{l} H(Y) = 0 \end{array} \right.$$

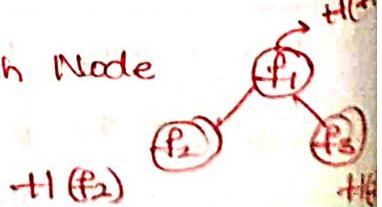
→ from this we can say if both the classes are equally probable then we get maximum entropy ie.

ie. 1. (Impure Subsplit)

→ If one is fully dominant we can say that Entropy is minimum ie. 0. (Pure Subsplit)

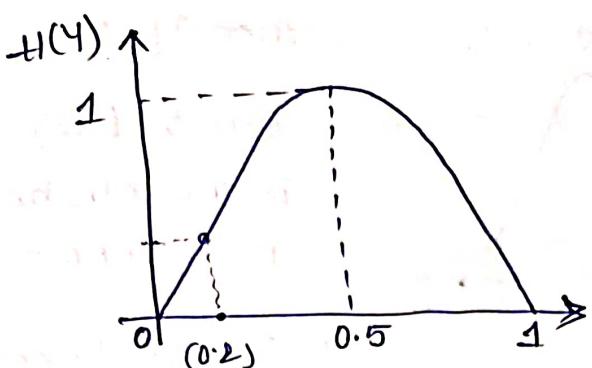
→ As we go from maximum ↓ Our Entropy tends to decrease

→ Entropy will be calculated for each Node



$$* P(y_t) = 1 - P(y_i)$$

Graphically we can define as



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

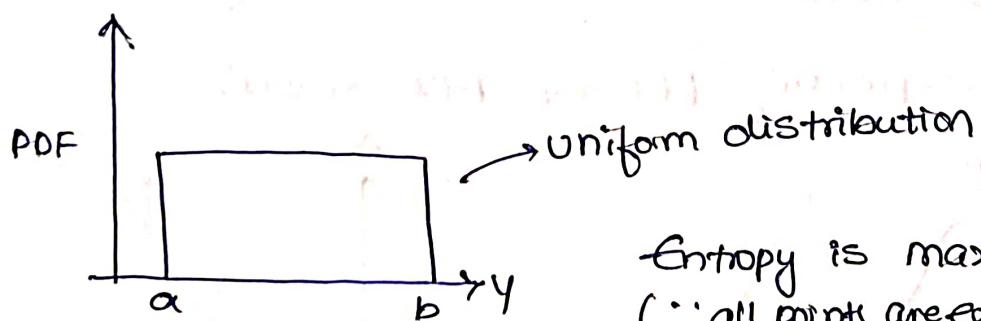
* If we have n class R.V

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

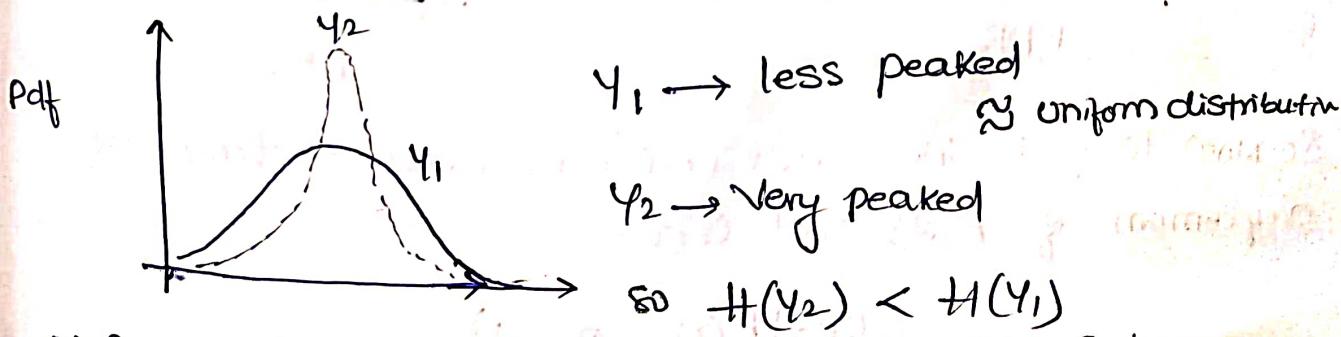
If all of them are equally probable
Entropy is maximum.

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

$$\text{so } H(y_2) < H(y_1)$$

\Rightarrow More peaked distribution is, less is its entropy.

34.5 Building a decision Tree : Information Gain

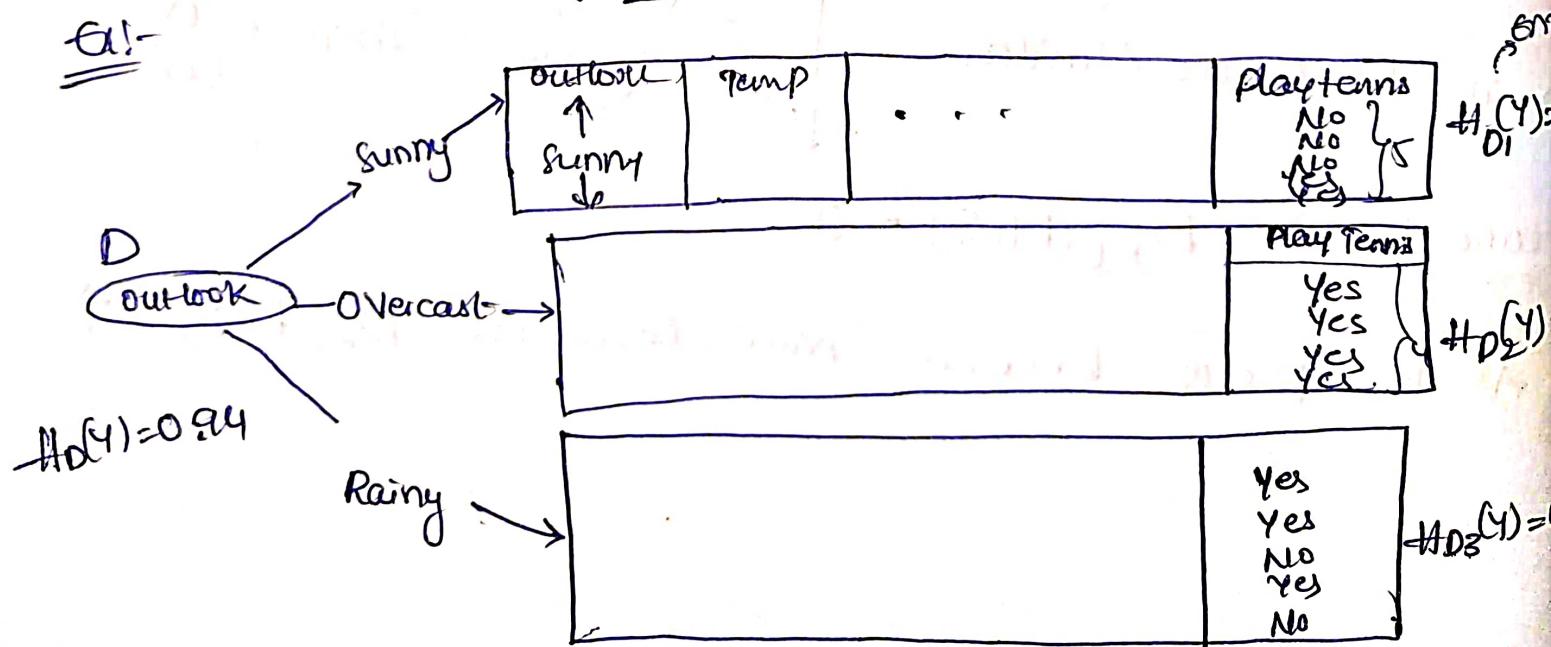
for Random Variable Y if that Random Variable divided into small sub set.

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

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

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

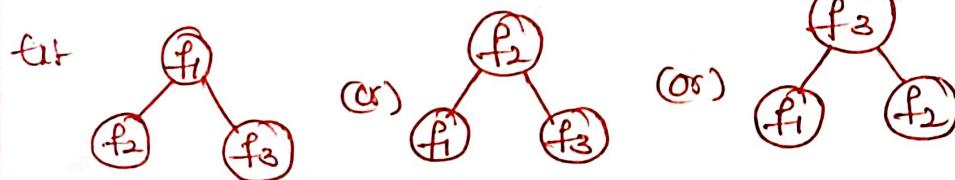
$$\text{Entropy of parent} = H_D(Y)$$

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.

34.6 Building a decision Tree : Gini Impurity

Gini Impurity is similar to entropy.

$$I_G(Y) = 1 - \sum_{i=1}^K [P(Y_i)]^2$$

$Y \rightarrow y_1, y_2, y_3, \dots, y_K$

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

$$I_G(Y)$$

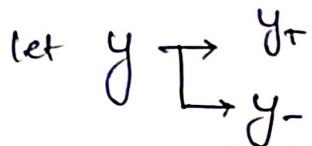
$$1 - P(Y+)^2 + P(Y-)^2$$

$$H(Y)$$

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

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

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



$$\text{Case 1: } P(Y+) = 0.5$$

$$P(Y-) = 0.5$$

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

$$H(Y) = 1$$

$$\text{Case 2: } P(Y+) = 1$$

$$P(Y-) = 0$$

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

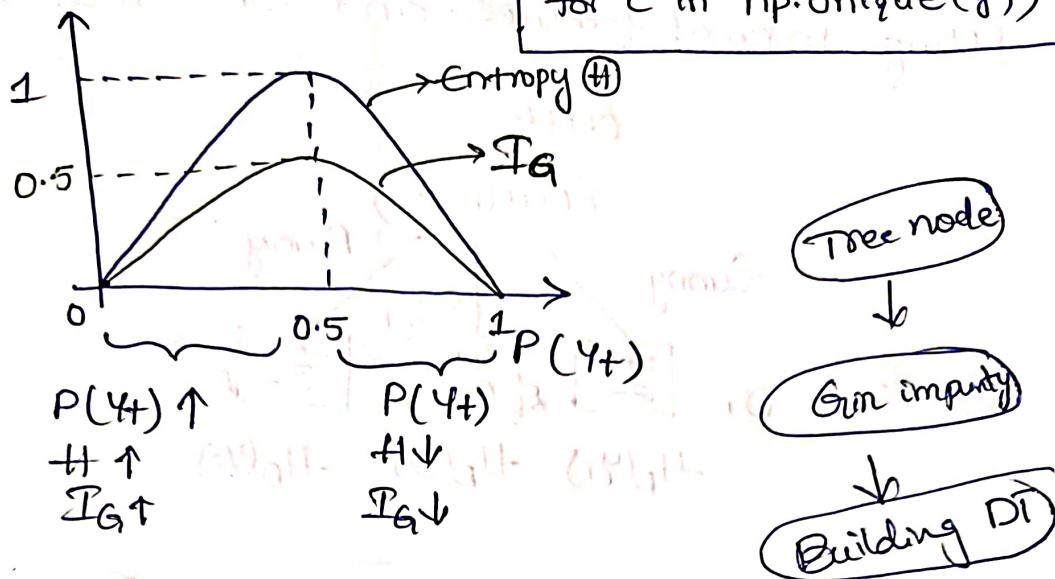
Geometrically,

2-Category Case :- y_+ , y_-

Gini Impurity (Programmatically):

$$1 - \sum ((\text{np.sum}(y == c) / m))^2$$

for c in $\text{np.unique}(y)$



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

Data set.

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

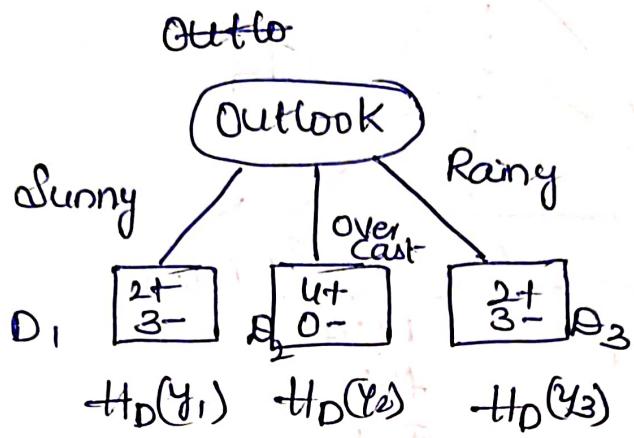
$\Theta \rightarrow 9+, 5-$

$$H_D(y) = 0.94$$

$\text{num-samples-per-class} =$
 $\frac{\text{left}}{\text{right}} P(y_i)$
 $\text{predicted class} = \arg \max [P(y_i)]$

Step 2

Now Divide The feature Data into 'n' Sub datasets using internal features. and find their Entropies.



and after that find either weighted Entropy

In This Example = 0.69

Step 3

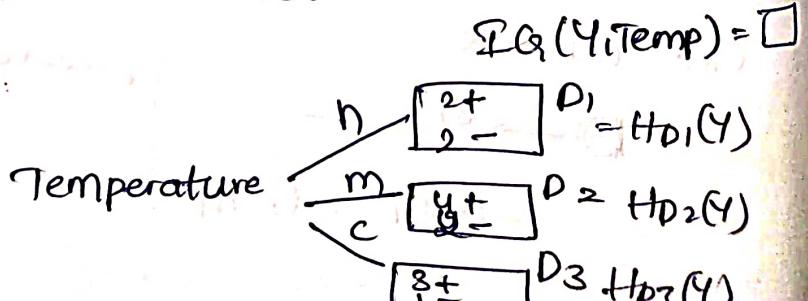
find $I_G(Y, \text{outlook})$ [that means find the Information By using the feature 1] using $\sum H_D(Y_i)$ & weighted Entropy.

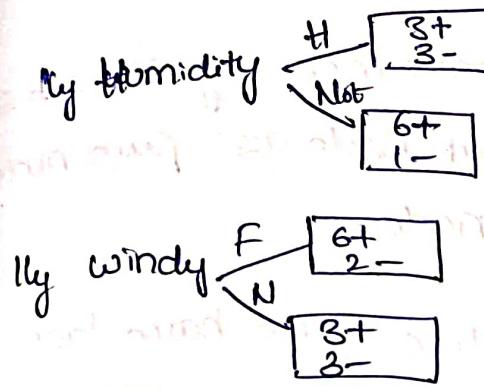
In this Example $I_G(Y, \text{outlook}) = 0.94 - 0.69 = 0.25$

Step 4

Repeat the same above steps and find the I_G for every feature in the Data set.

In This Example





Now once we get the all the IG values

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

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

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

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

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

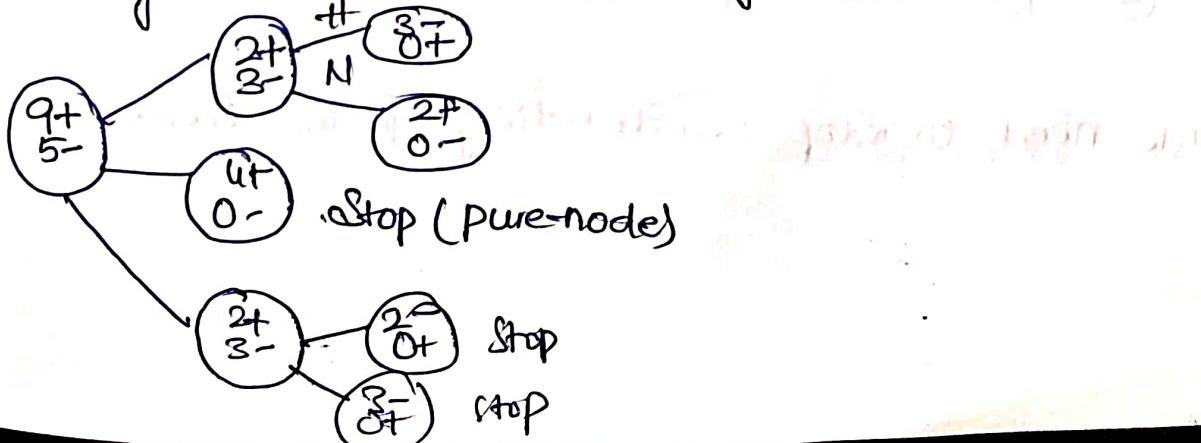
Choose the feature which has highest value and make that as parent node.

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

- weighted entropy @ Child level.

Step 5:

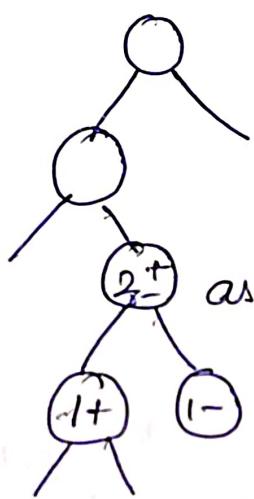
Recursively break each node using IG as the Criterion



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



$n=10K$ 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 op.

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

① pure node ✓

num-samples_per_class.Count(0) == 0

② few pts @ a node ✓

if depth < maxDepth
exit

③ If we are too deep ✓

We need to stop constructing further tree.

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

depth is small \Rightarrow underfit.

In decision tree hyperparameter is depth and which can find out 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.

EG: \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.

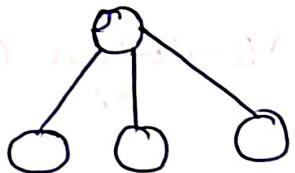
x_1	y
2.2	1
2.6	1
3.5	0
3.8	0
4.6	1
5.3	0

where x_1 : numerical

\hookrightarrow Integer(α) 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 preferable

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

$$f_1 < 2.2$$

Like this we get n possible variations.

$$f_1 < 2.6$$

$$f_1 < 3.5$$

$$f_1 < 3.8$$

$$f_1 < 4.6$$

$$f_1 < 5.3$$

→ Based on Threshold we divide the features

$$\begin{aligned}
 f_1 < T_1 &\quad D_1 \\
 &\quad D_2 \\
 f_2 < T_2 &\quad D_1 \\
 &\quad D_2 \\
 f_3 < T_3 &\quad D_1 \\
 &\quad D_2 \\
 \vdots & \\
 f_n < T_n &\quad D_1 \\
 &\quad D_2
 \end{aligned}$$

f_1	y	
2.2	1	D_1
2.6	1	
35	0	D_2
3.8	0	
4.6	1	D_2
5.3	0	

Like this we have 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 time taking process so we have some hacks to work around.

It first starts with

then it moves the length

towards the minimum

and start comparing

with minimum length

and so on until

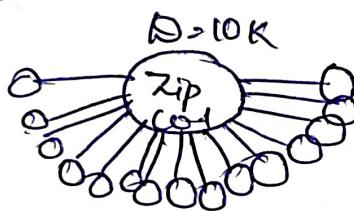
it reaches the length

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.

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 ₁	
P ₂	
P ₃	
P ₄	
P ₁	
P ₂	
P ₃	

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.

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

$$= \frac{\text{no of } \# Y_i=1 \text{ fragmen } P_j}{\# P_j}$$

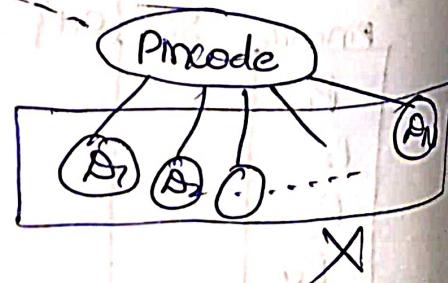
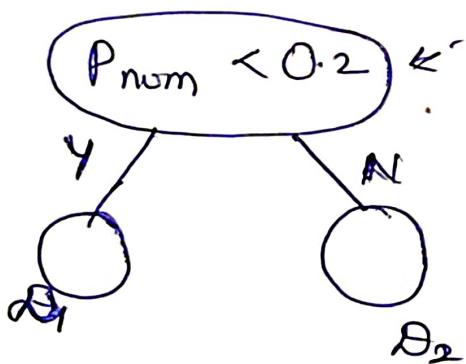
let we have 20 times P_j and in that 19 times $y_i=1 \Rightarrow 19/20$

$$\begin{array}{c} P_j \xrightarrow{\text{Categorical feature}} \text{Numerical feature} = P(Y_i=1 | P_j) \\ \text{Numerical feature} \end{array}$$

$D \rightarrow$ remove
pincode
features

numerical
features

so by this we can make decision tree like



34.11 Overfitting and Underfitting

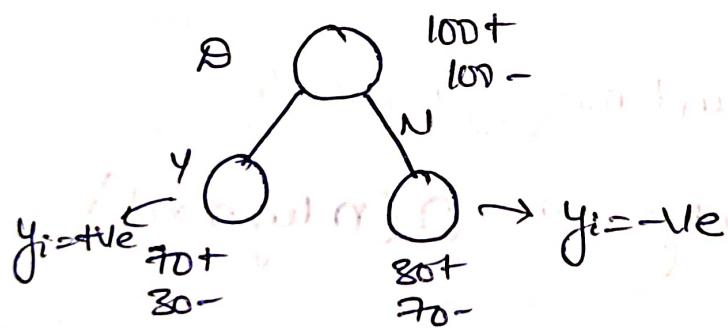
If the depth \uparrow :
possibility of having very few points @ a leaf node \uparrow

By this whole model becomes noisy, Overfitting my model to noise.

Interpretability of the model also \downarrow
because

if $() \& () \& () \& () \dots$ get \uparrow
for a $x_2 \rightarrow y_2$ by this interpretability decreases.

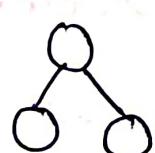
Similarly If the depth is very low (min depth would be 1)



so by this we end up getting none of the node is pure

so If we given a Query point, The class label y_2 would take majority among them.

Note This type of tree called decision stump.

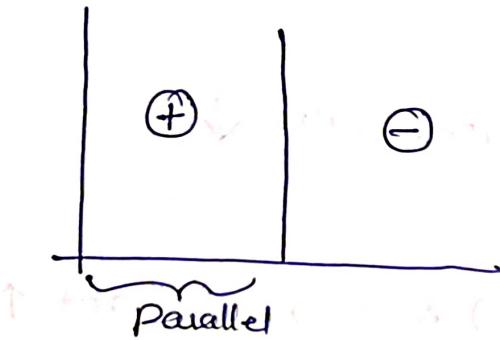


So if we go too deep \rightarrow Overfit

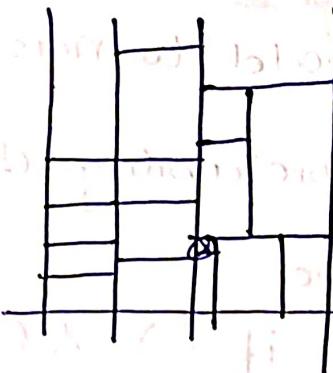
too shallow \rightarrow Underfit

Geometrically

Underfit



Overfit



So we use Cross Validation 'CV' to better determine the depth to be taken into consideration.

34.12 Train and Runtime complexity :-

Train time Complexity $\sim O(n \log n * d)$

n = no of points D_{train}

d = dimensions

When we have numerical features It becomes more computational complex to solve, but there are many algorithmic methods or hacks to solve.

$n \log n$ in the complexity is mostly used for sorting

After training @ runtime the space needed to store my decision tree by a simple nested if else is small.

is small. no of

Simply In nested if else we use 1 internal nodes + no of leaf nodes.

Typically we train decision tree with depth 5 to 10.

as depth $\uparrow \Rightarrow$ Interpretability \downarrow

So the runtime complexity is reasonable.

whereas run time complexity would do K decisions where K is max depth of any leaf node.

run time complexity: $O(\text{depth})$.

Decision trees are good

\rightarrow when we have large data, and dim is small

\rightarrow and also when we have to make low latency

(as) the time complexity for DT is $O(\text{depth})$

34.9 Feature Standardization :-

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

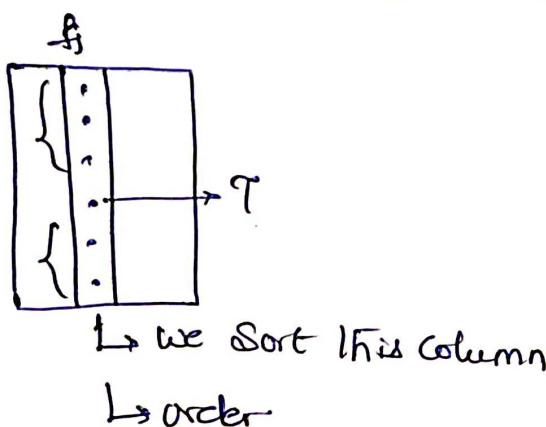
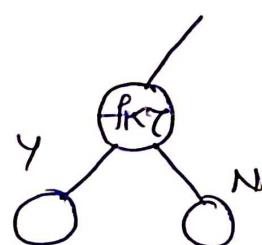
In logistic regression

$$\begin{array}{|c|c|c|} \hline & f_j & \\ \hline x_i & x_{ij} & \\ \hline & \downarrow & \\ & u_j, \sigma_j & \end{array}$$

$$x'_{ij} = \frac{x_{ij} - u_j}{\sigma_j}$$

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 values not on distances that differ based on scales.