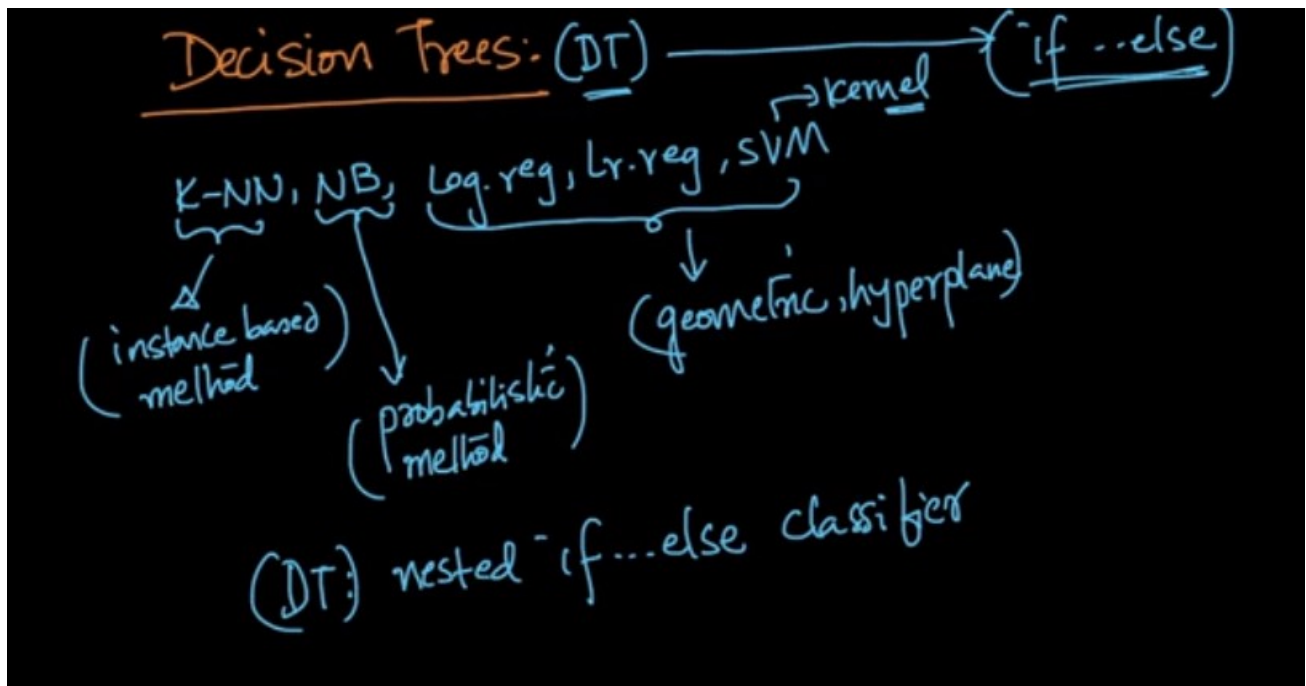


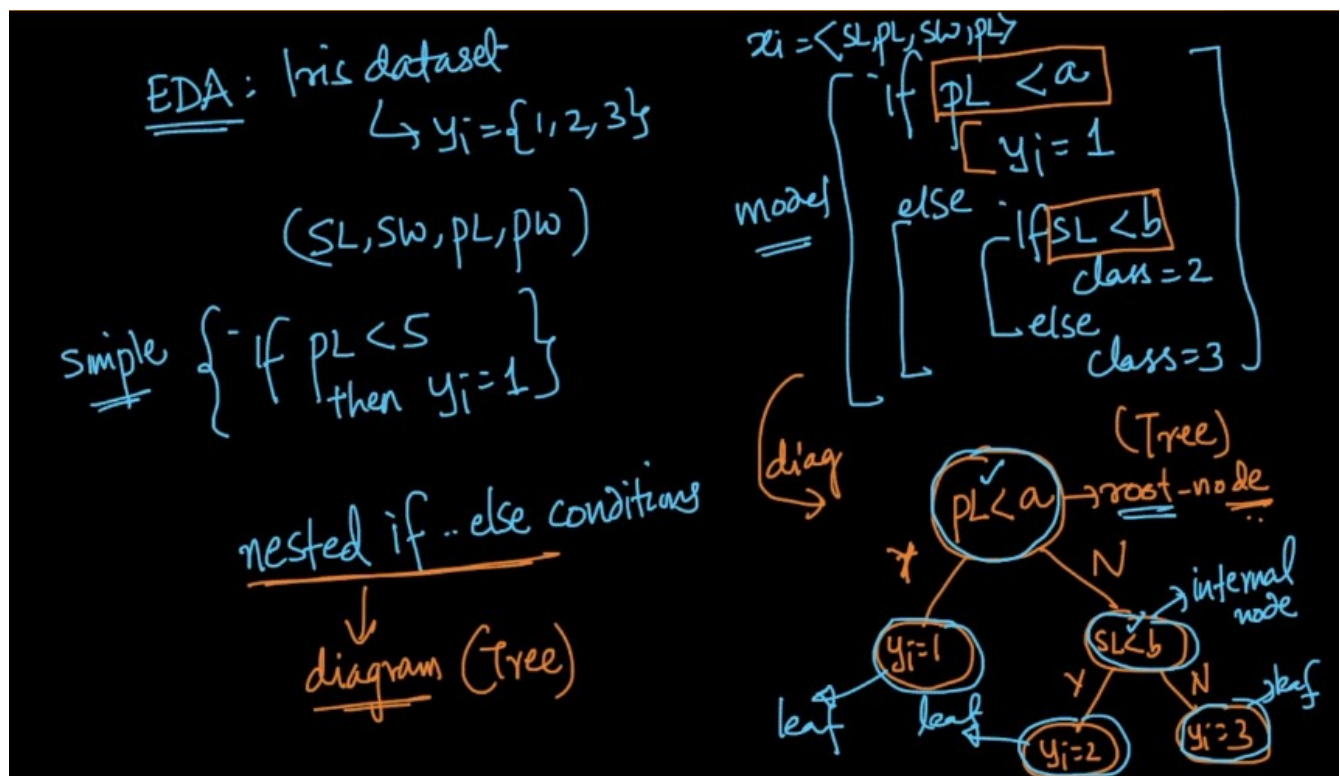
## NOTES

## Different classifiers and there properties.



Nested if – else:

Decision tree starts with a node with conditions and divided into several nodes, this diagram is called a tree.



Root – node: This is the starting node.

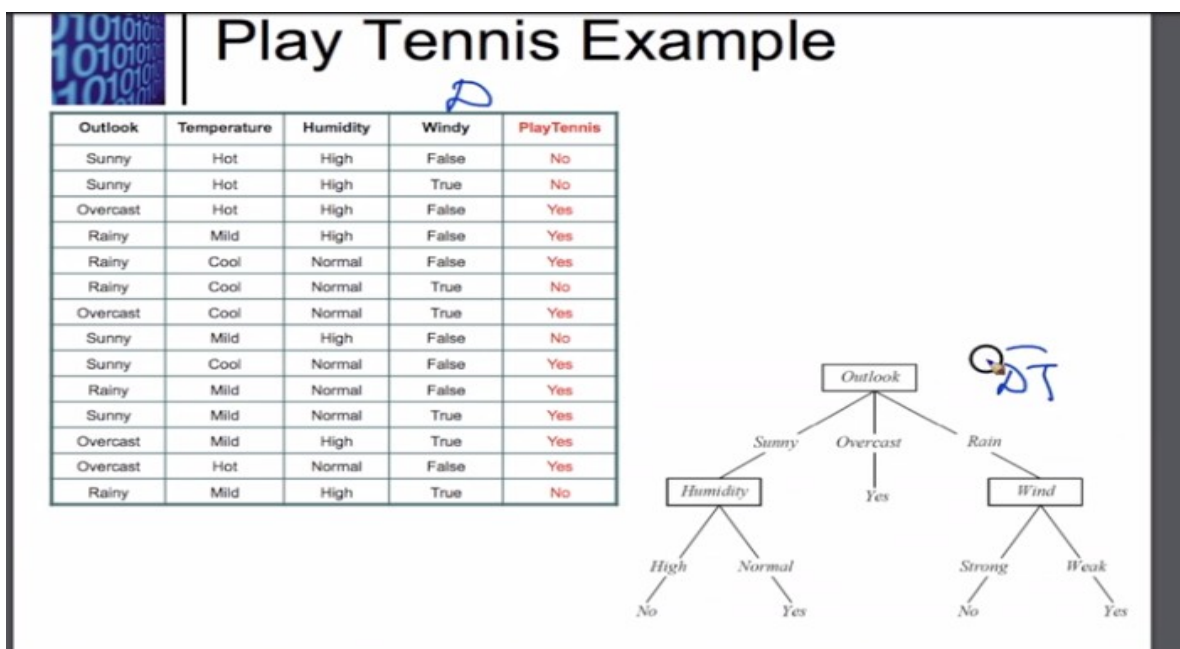
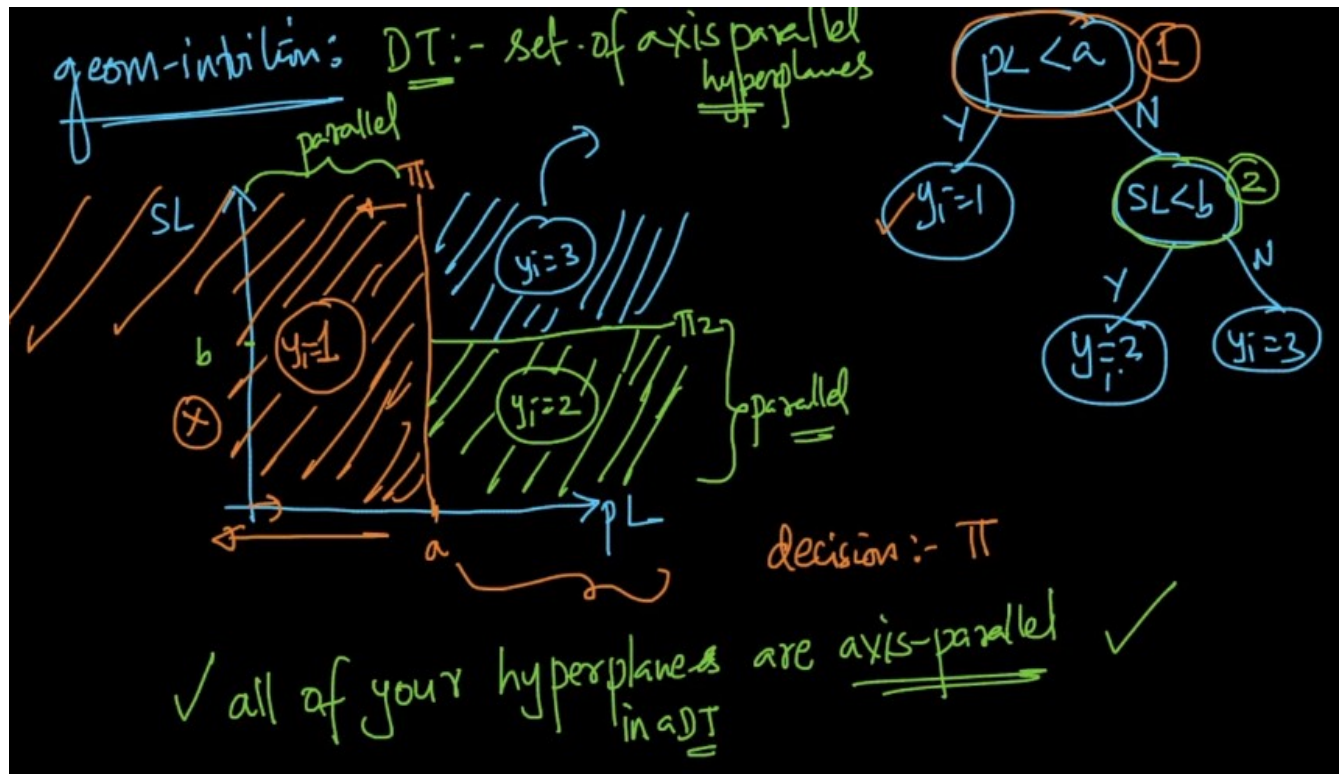
Leaf – node: These are branches of the root node.

Internal – nodes: These are neither leaf nodes or root nodes.

Non – leaf nodes: These are nodes that make decisions, these are called the nodes of decision tree.

Geometry:

In decision tree all of the hyper planes are axis parallel.



## Building a Decision Tree: Entropy:

Play Tennis Example

Outlook	Temperature	Humidity	Windy	Play Tennis
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Handwritten notes:

- $Y$ : play Tennis
- $Y_+, Y_-$ 

$$\begin{cases} P(Y_+) = \frac{9}{14} \\ P(Y_-) = 1 - P(Y_+) = \frac{5}{14} \end{cases}$$

Decision Tree Structure:

```

graph TD
    Outlook[Outlook] -- Sunny --> Humidity[Humidity]
    Outlook -- Overcast --> Yes1[Yes]
    Outlook -- Rain --> Wind[Wind]
    Humidity -- High --> No1[No]
    Humidity -- Normal --> Yes2[Yes]
    Wind -- Strong --> No2[No]
    Wind -- Weak --> Yes3[Yes]
  
```

Handwritten note: 2.5

Entropy

Handwritten notes:

- $x.v \quad Y \rightarrow y_1, y_2, y_3, \dots, y_k$
- Entropy:  $H(Y) = - \sum_{i=1}^K p(y_i) \log_b(p(y_i))$
- $p(y_i) =$
- $p(y_i) = P(Y=y_i)$
- $b = 2$  or  $b = e = 2.718$
- $\log_2 = \lg$
- $\log_e = \ln$

## Calculation of the Entropy

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

$$\underline{H(Y)} = - \left[ \frac{9}{14} \log_2\left(\frac{9}{14}\right) + \left(\frac{5}{14}\right) \log_2\left(\frac{5}{14}\right) \right] = \underline{\underline{0.94}} \checkmark$$

Annotations:

- $\frac{9}{14}$  is labeled as  $p(y_+)$  and  $\frac{5}{14}$  is labeled as  $p(y_-)$ .
- A box below the fractions contains the text:  $\frac{\# \text{ +ve pts}}{\text{Total } \# \text{ pts}} = \frac{\% \text{ age of +ve pts in } D$ .
- A circle around  $\frac{5}{14}$  is labeled:  $\% \text{ age of -ve pts in } D$ .

Properties of Entropy:

Various cases:

Properties:  $Y \rightarrow y_+, y_-$  (2 class, 2 category)

Case 1:  $\begin{cases} y_+ \rightarrow 99\% \\ y_- \rightarrow 1\% \end{cases} \quad H(Y) = -0.99 \log_2 0.99 - 0.01 \log_2 0.01 = \underline{\underline{0.0801}}$

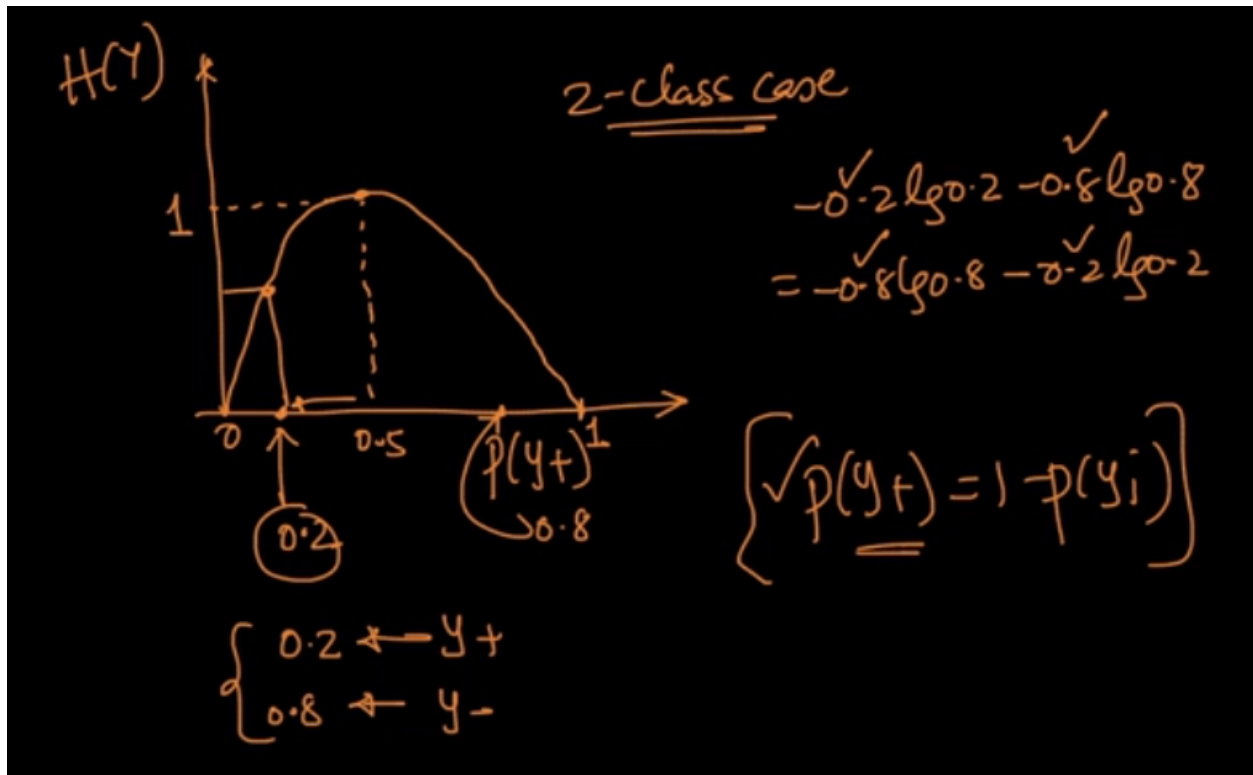
Case 2:  $\begin{cases} y_+ \rightarrow 50\% \\ y_- \rightarrow 50\% \end{cases} \quad H(Y) = -0.5 \log_2 0.5 - 0.5 \log_2 0.5 = \underline{\underline{1}}$

Case 3:  $\begin{cases} \times y_+ \rightarrow 0\% \\ \checkmark y_- \rightarrow 100\% \end{cases} \quad H(Y) = \underline{\underline{0}}$



Entropy curve:

This is a symmetric curve.

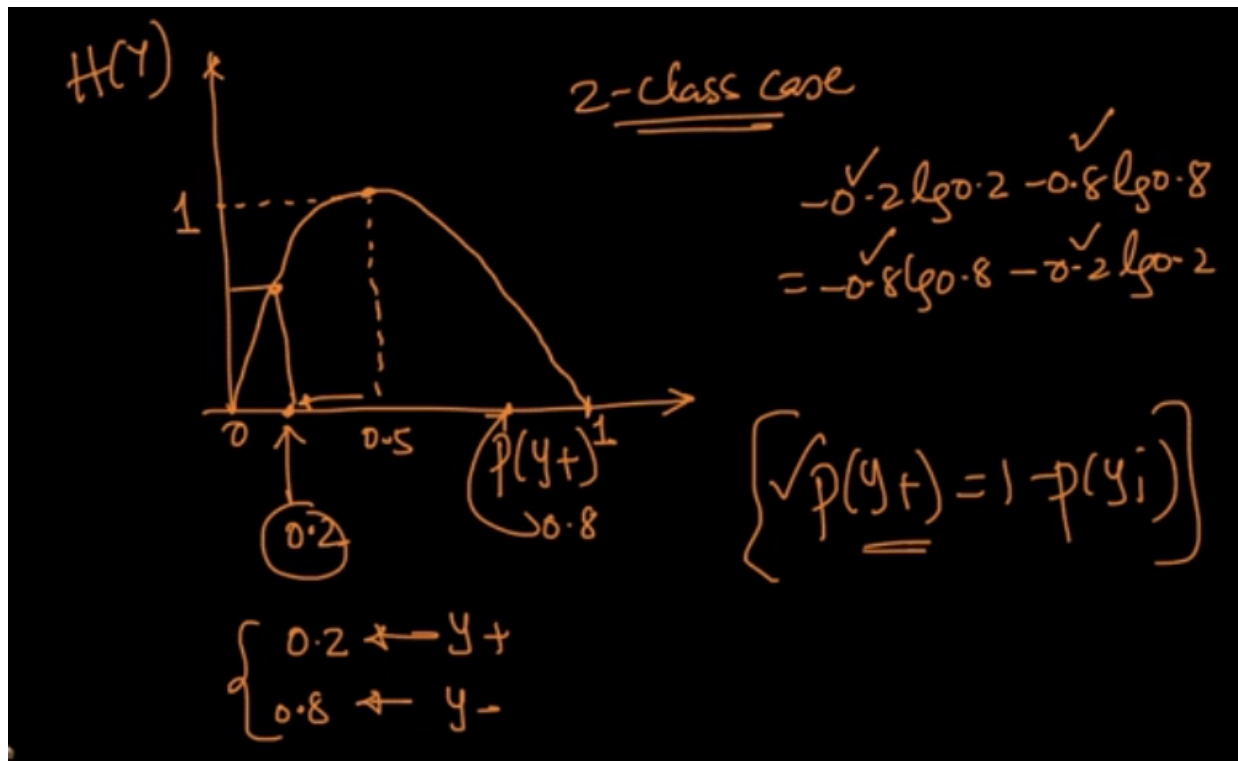


Conclusion:

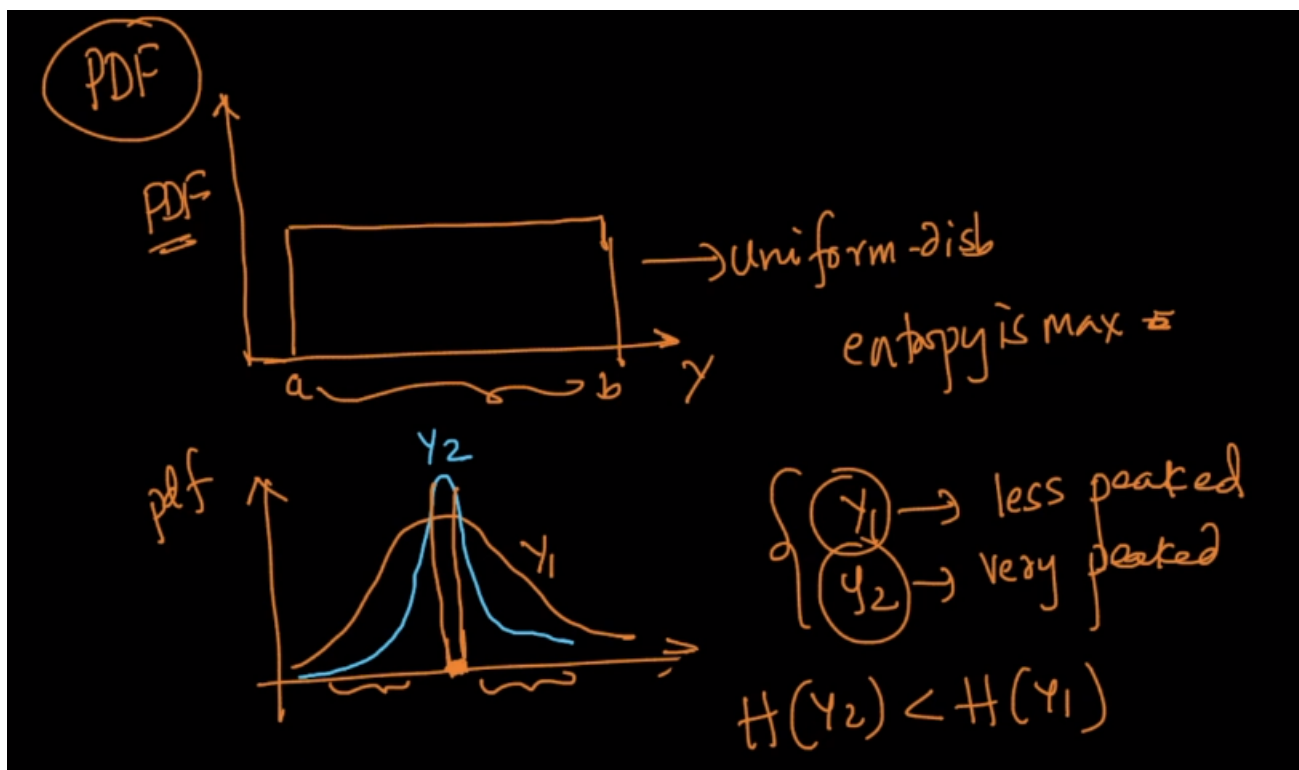
$X \rightarrow y_1, y_2, \dots, y_k$   
 equi-probable  $\rightarrow$  entropy is maximum

$y_1 \rightarrow$  most probable }  $\rightarrow$  entropy is minimum  
 $y_2, y_3, \dots \rightarrow 0$

Entropy for real – valued feature:

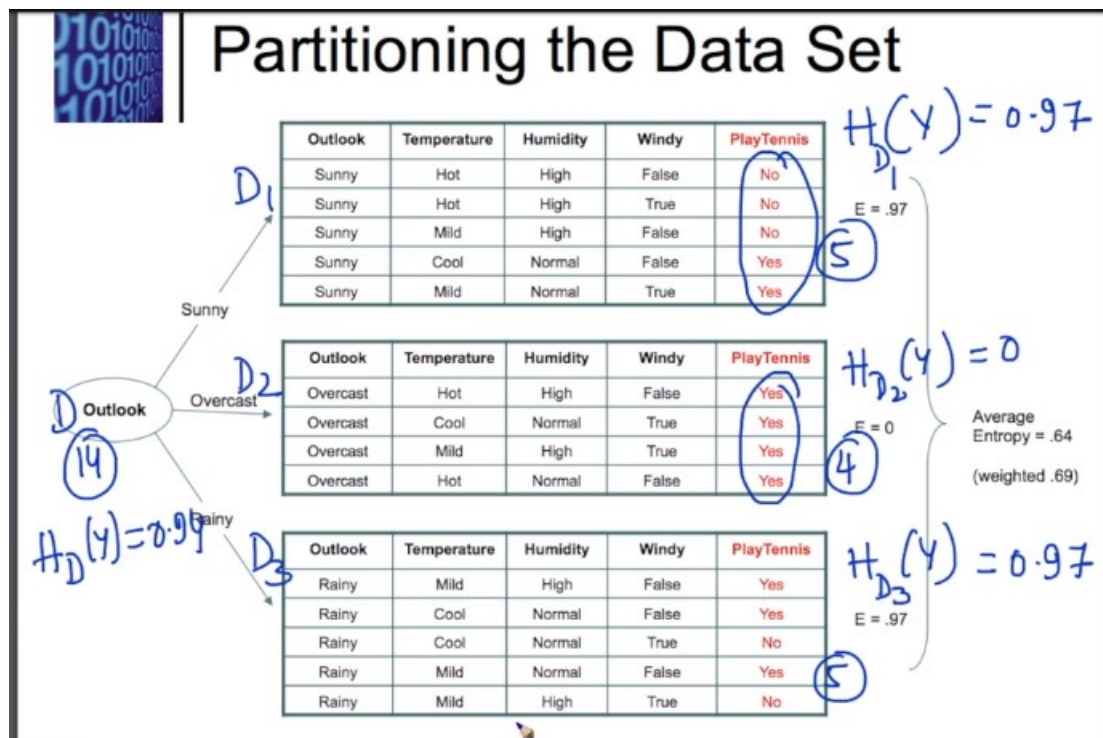


Given a random variable, if all of them are equi – probable, entropy is maximum.  
If the data distribution is more uniform then the entropy is more.



If the data is more like gaussian then the data has less entropy.

Information gain:



Calculation of information gain:

$$IG(Y, outlook) = \left( \frac{5}{7} \times 0.97 \right) - 0.94 = 1G$$

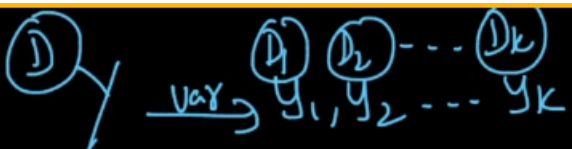
Weighted entropy after  $D_1, D_2, D_3$

$$\left( \frac{5}{14} \times 0.97 \right) + \left( \frac{4}{14} \times 0 \right) + \left( \frac{5}{14} \times 0.97 \right)$$

$$\frac{5}{7} \times 0.97$$

$$\frac{5}{7} \times 0.97 \times 2$$

Formula for Information gain:



The diagram shows a variable  $Y$  in a circle, with an arrow labeled  $var$  pointing to a sequence of outcomes  $y_1, y_2, \dots, y_k$ . Above each outcome is a circled label  $D_1, D_2, \dots, D_k$ , representing a decision node.

$$IG(Y, var) = \sum_{i=1}^k \frac{|D_i|}{|D|} H(Y) - H_D(Y)$$



Gini – Impurity:  
Case 1:

Gini Impurity  $\sim$  similar to Entropy

$$I_G(Y) = 1 - \sum_{i=1}^K (p(y_i))^2$$

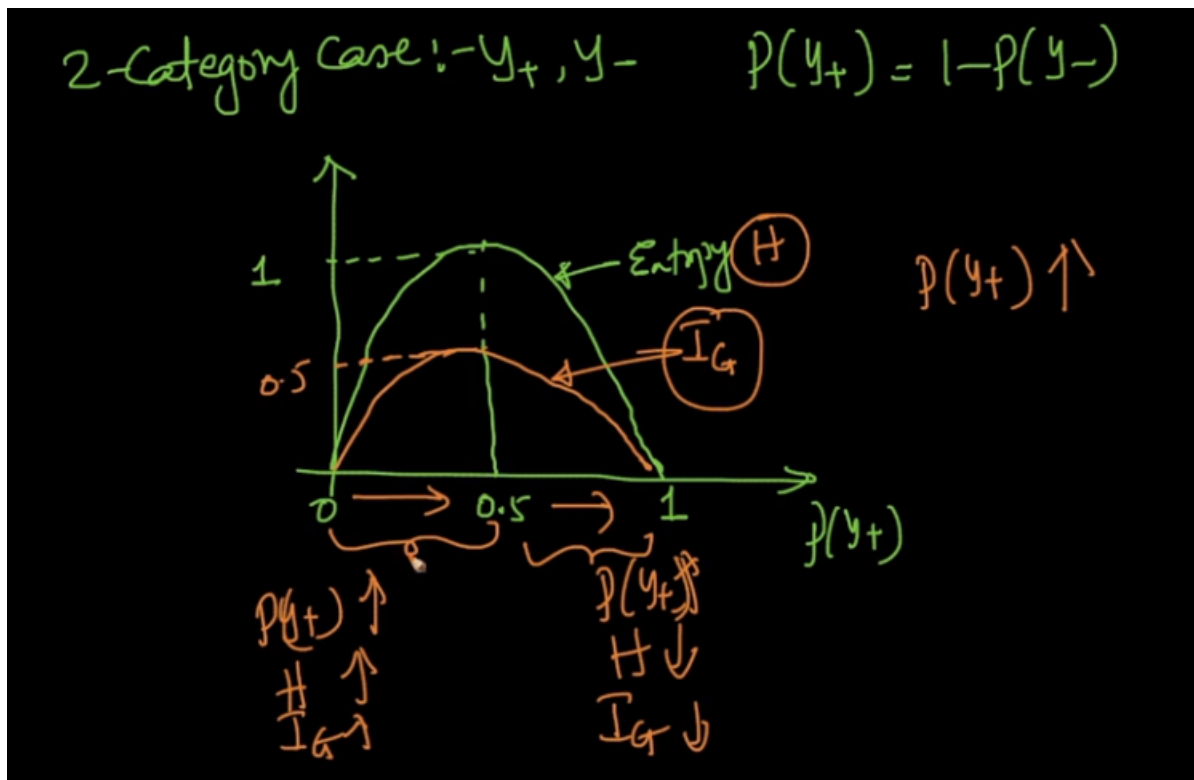
$Y \rightarrow y_+, y_-$

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

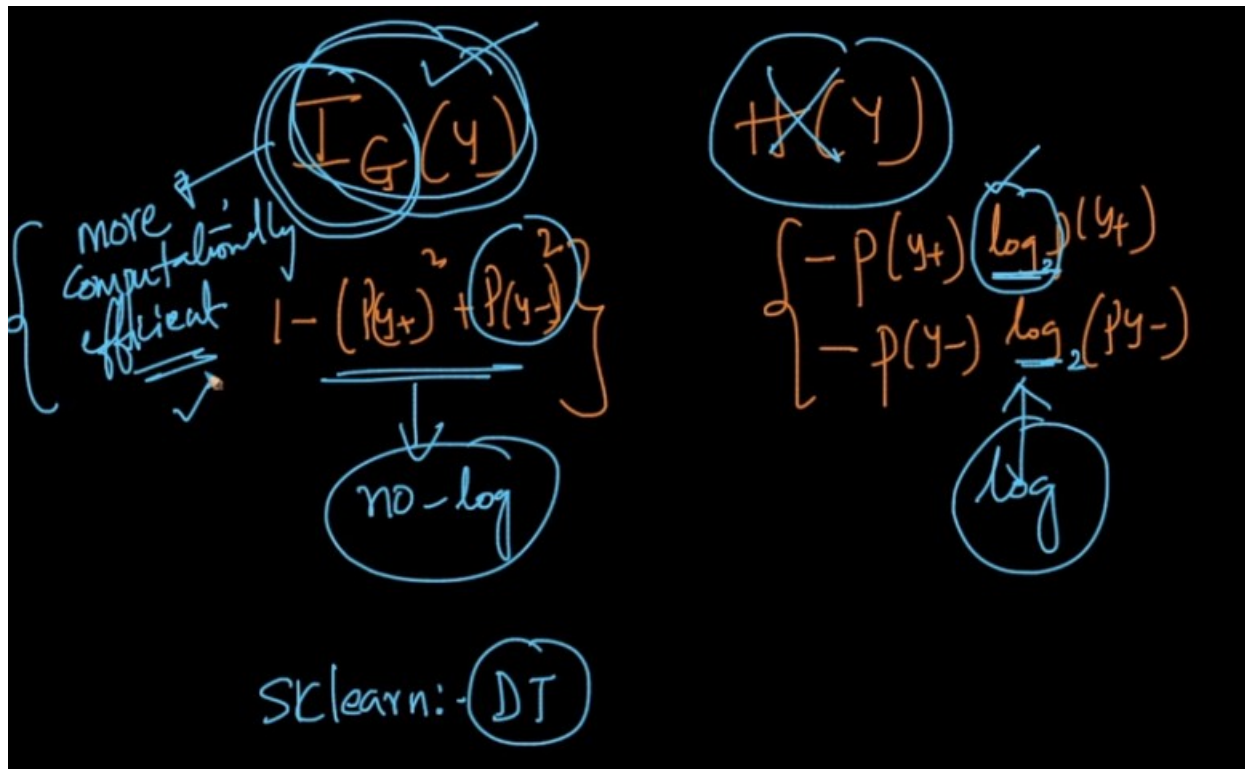
Case 2:  $p(y_+) = 1$   
 $p(y_-) = 0$   
 $I_G(Y) = 1 - (1 + 0) = 0$   
 $H(Y) \uparrow$

Case 1:  $p(y_+) = 0.5$   
 $p(y_-) = 0.5$   
 $I_G(Y) = 1 - (0.25 + 0.25) = 0.5$   
 $H(Y) = 1$

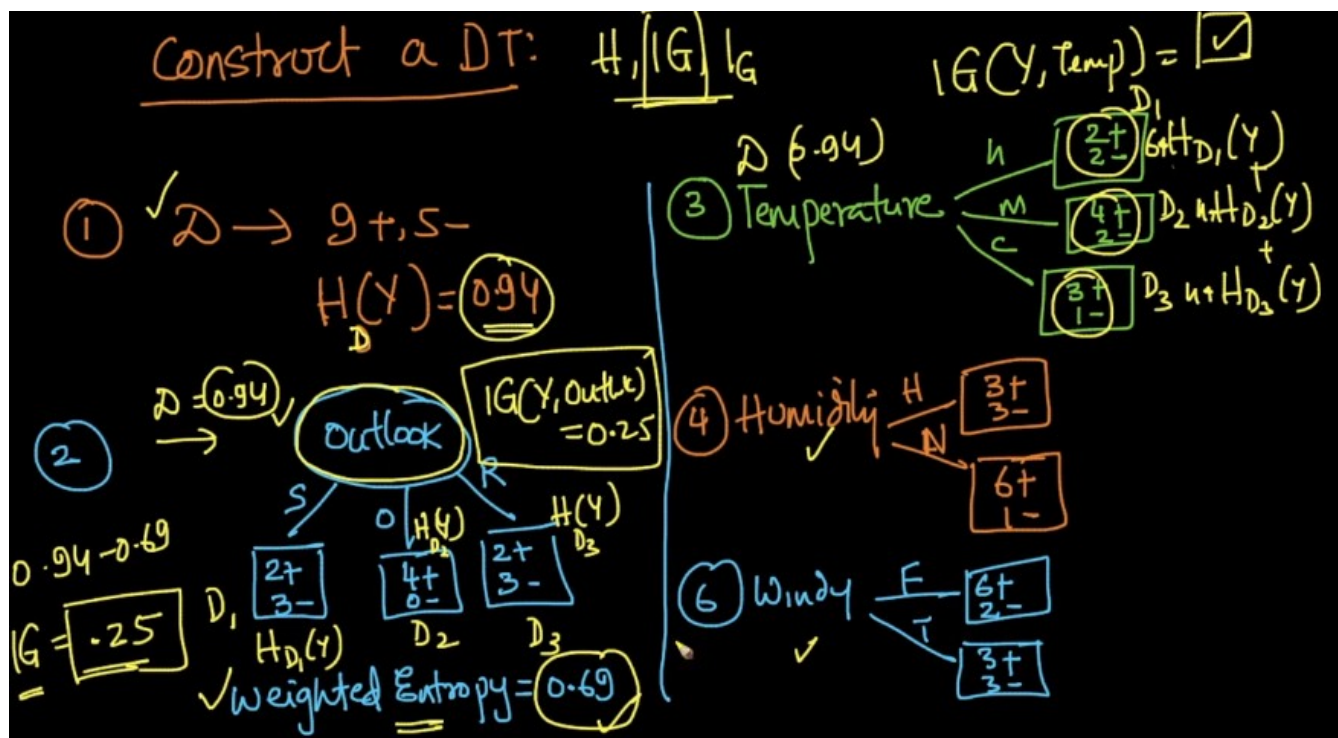
We assume there are only two classes of a feature:



Calculation of information gain is computationally more expensive than entropy, In real world people will take gini - impurity as the measure.



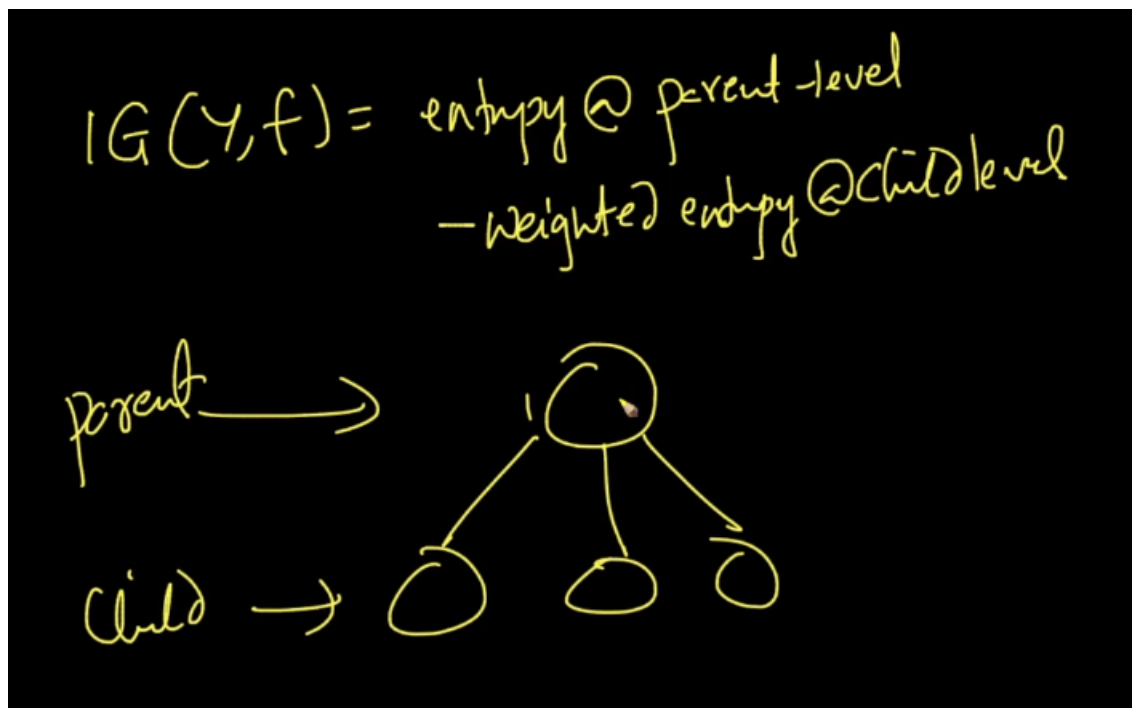
Building a decision Tree: Constructing a DT on a data set.  
Breaking the data set with the feature outlook



$$\checkmark IG(Y, f) = \underbrace{H_D(Y)}_{D \rightarrow} - \underbrace{\left( \sum_{i=1}^k \frac{|D_i|}{|D|} * H_{D_i}(Y) \right)}_{\rightarrow \text{choosing the root node}}$$

$$\left\{ \begin{array}{l} IG(Y, outlook) = 0.25 \\ IG(Y, Temp) = - \\ IG(Y, Humidity) = - \\ IG(Y, Windy) = - \end{array} \right.$$

Information gain on breaking the data set is the (entropy at the parent level) – (weighted entropy at the child level). **The node is chosen which has most information gain.**

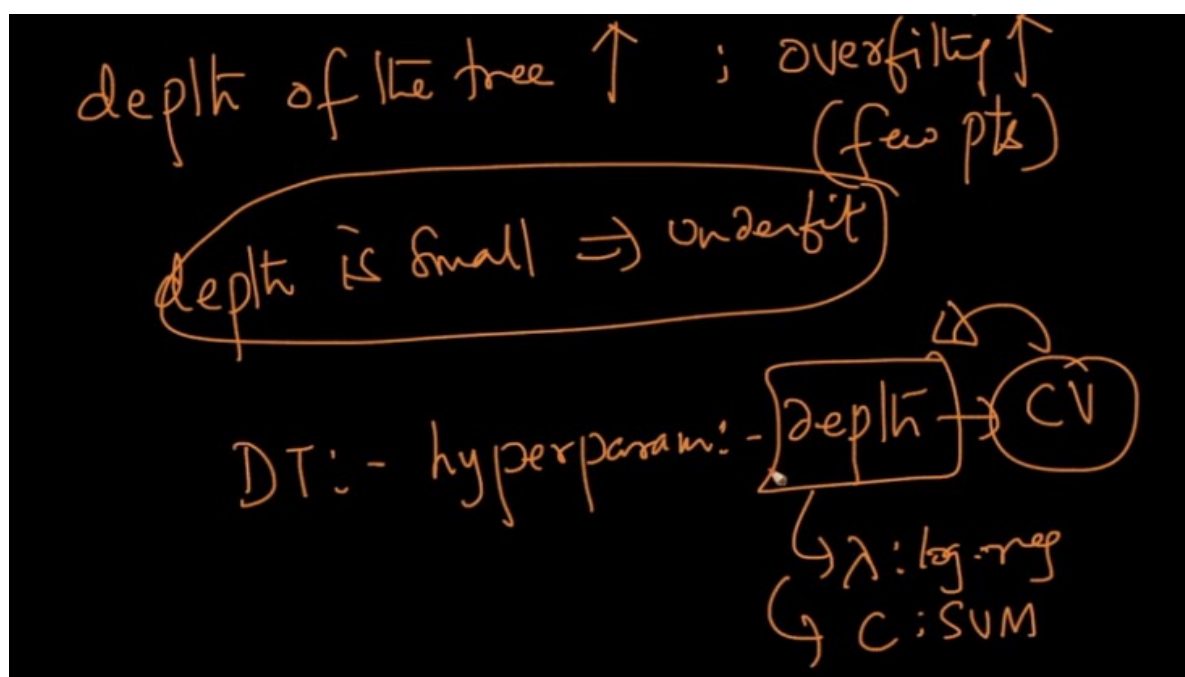
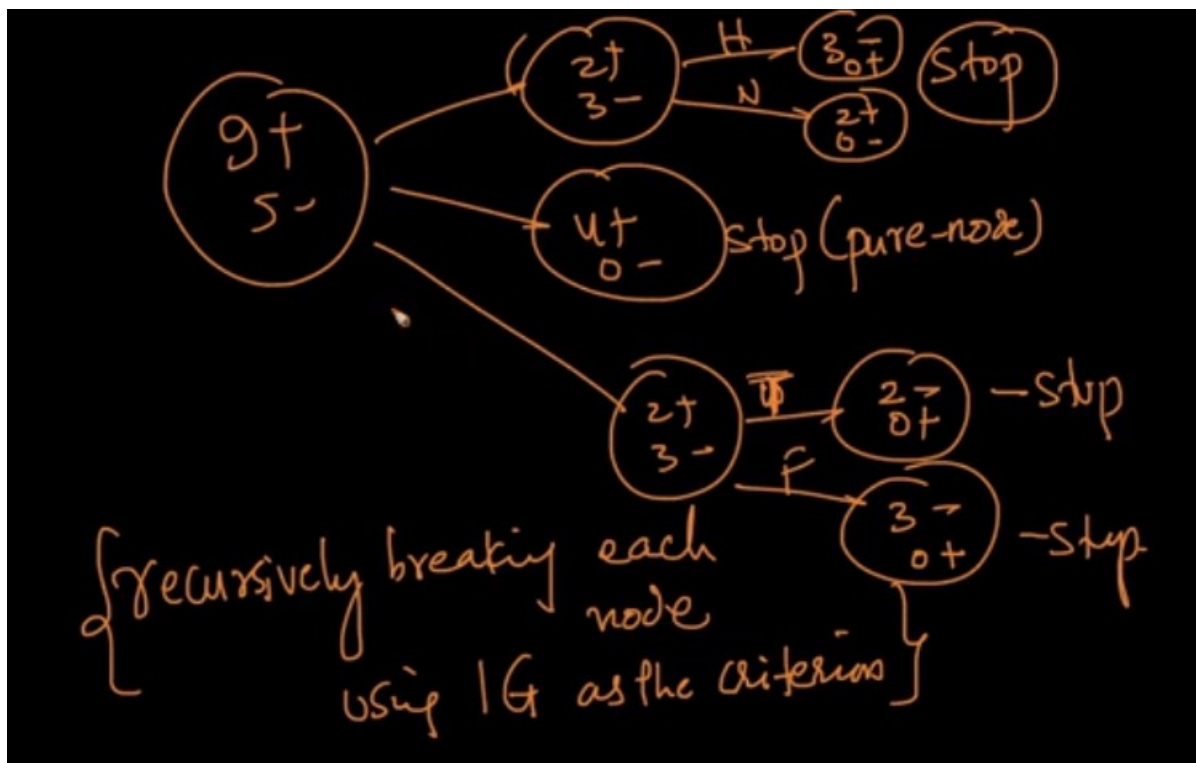




If depth is small then the decision tree tend to under fit.

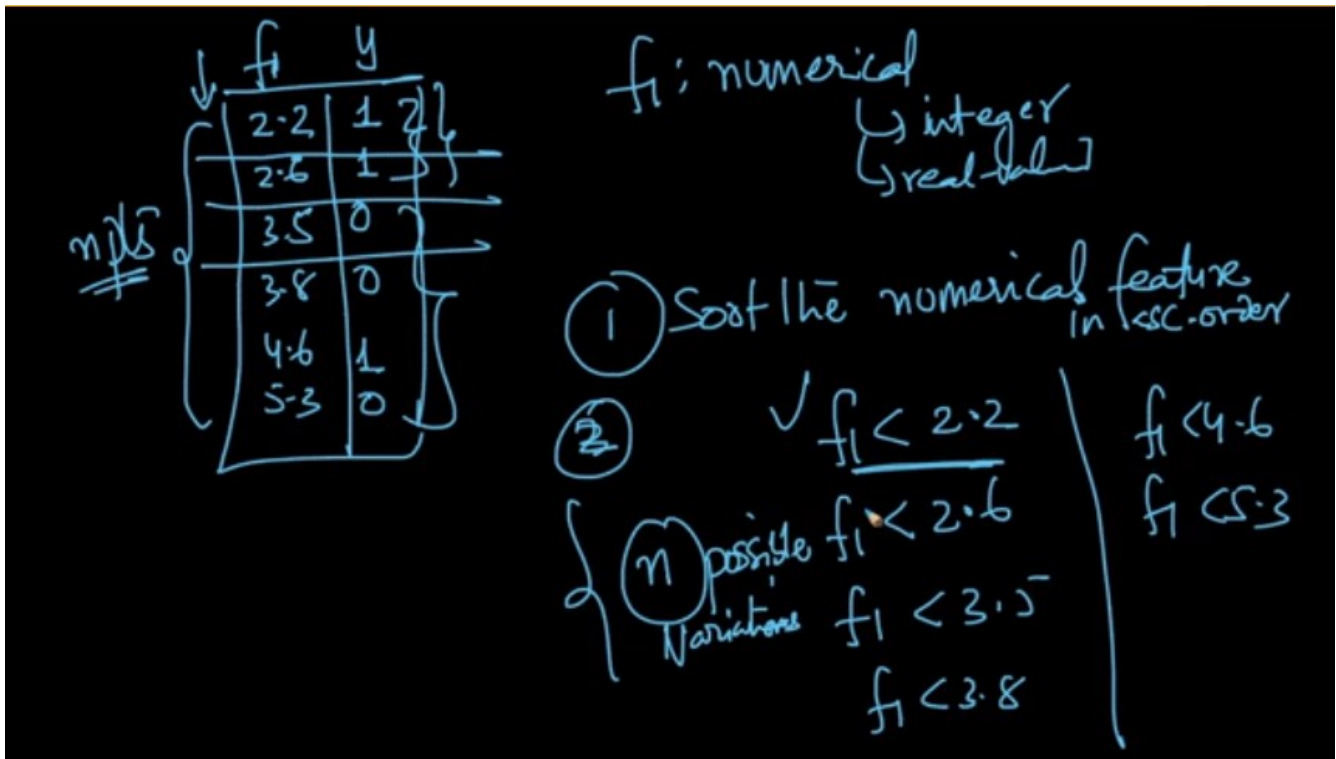
The Hyper parameter in decision tree is the height of the tree.

We use cross validation to choose the depth of the decision tree.



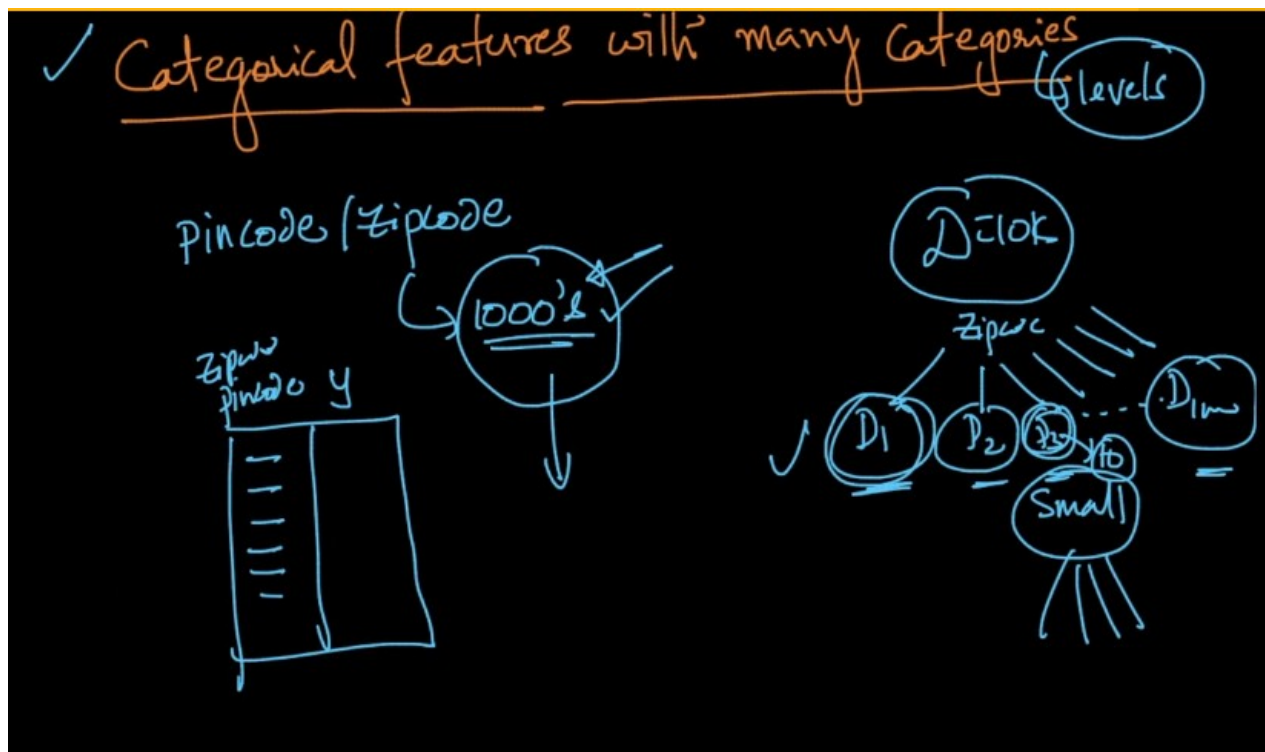


## Building a decision Tree: Splitting numerical features:

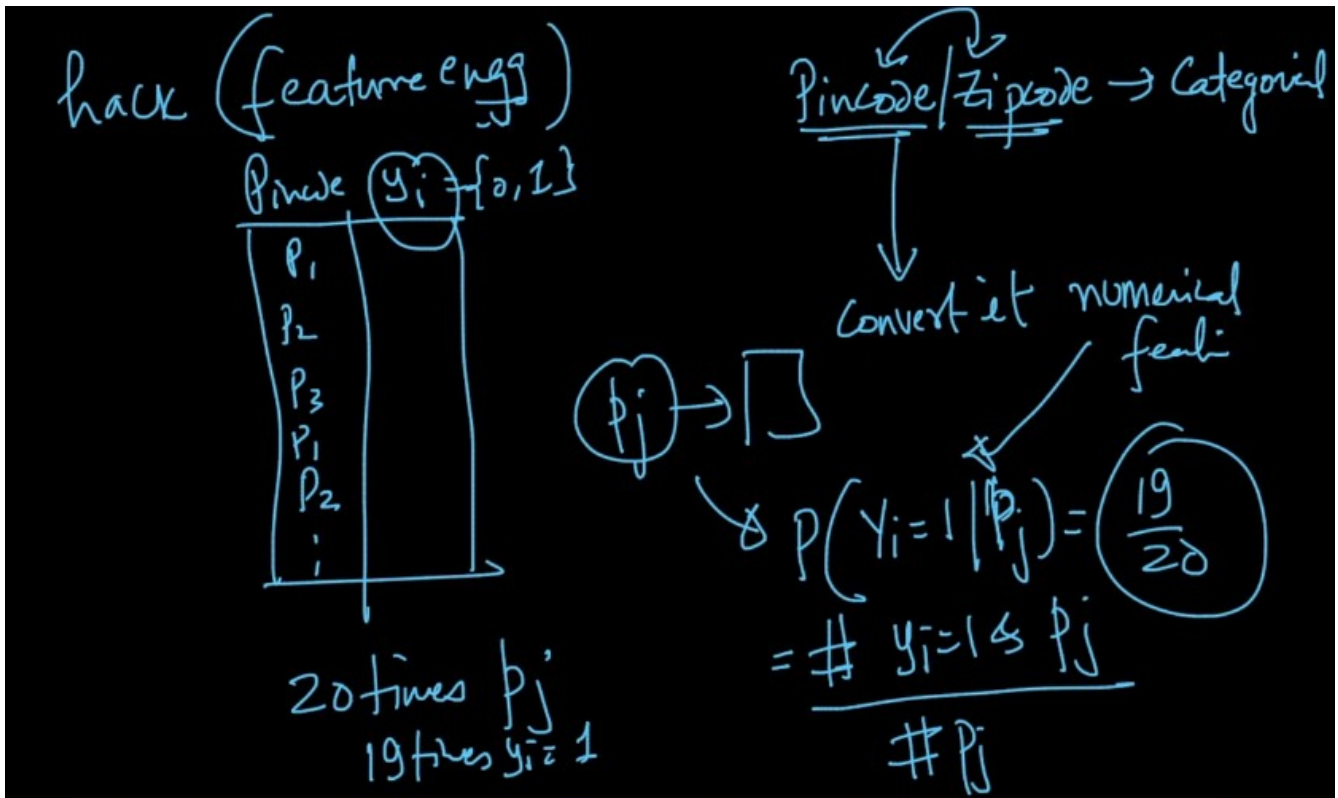


### Building a decision Tree: Categorical features with many possible values:

Example: ZIP code, PIN code.



### Converting a categorical feature to a numeric feature:

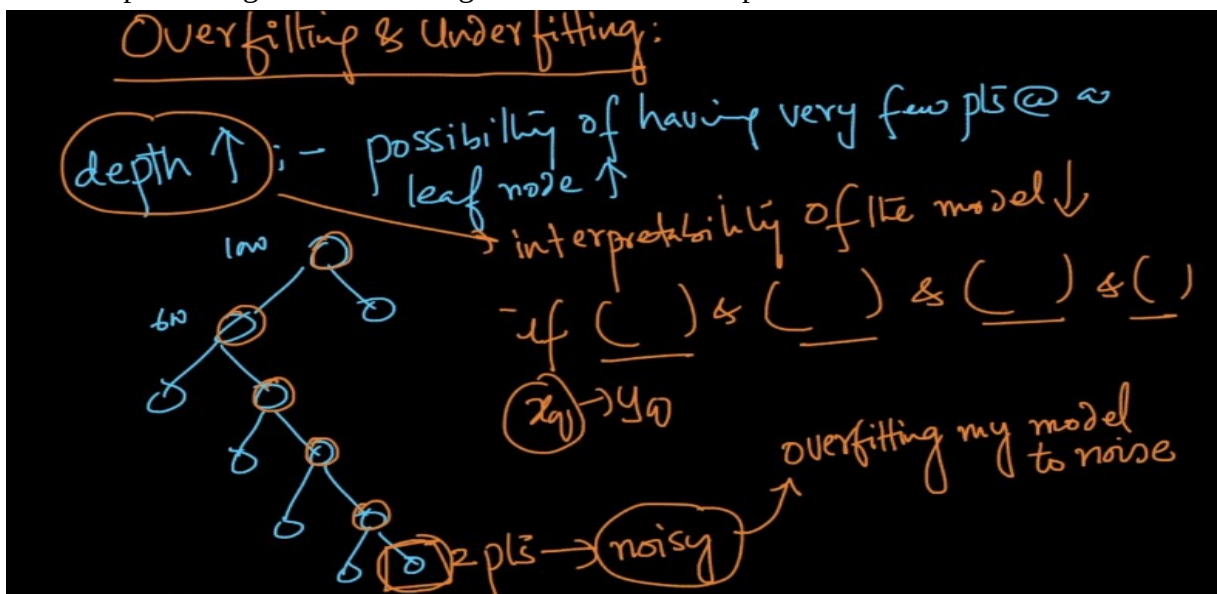


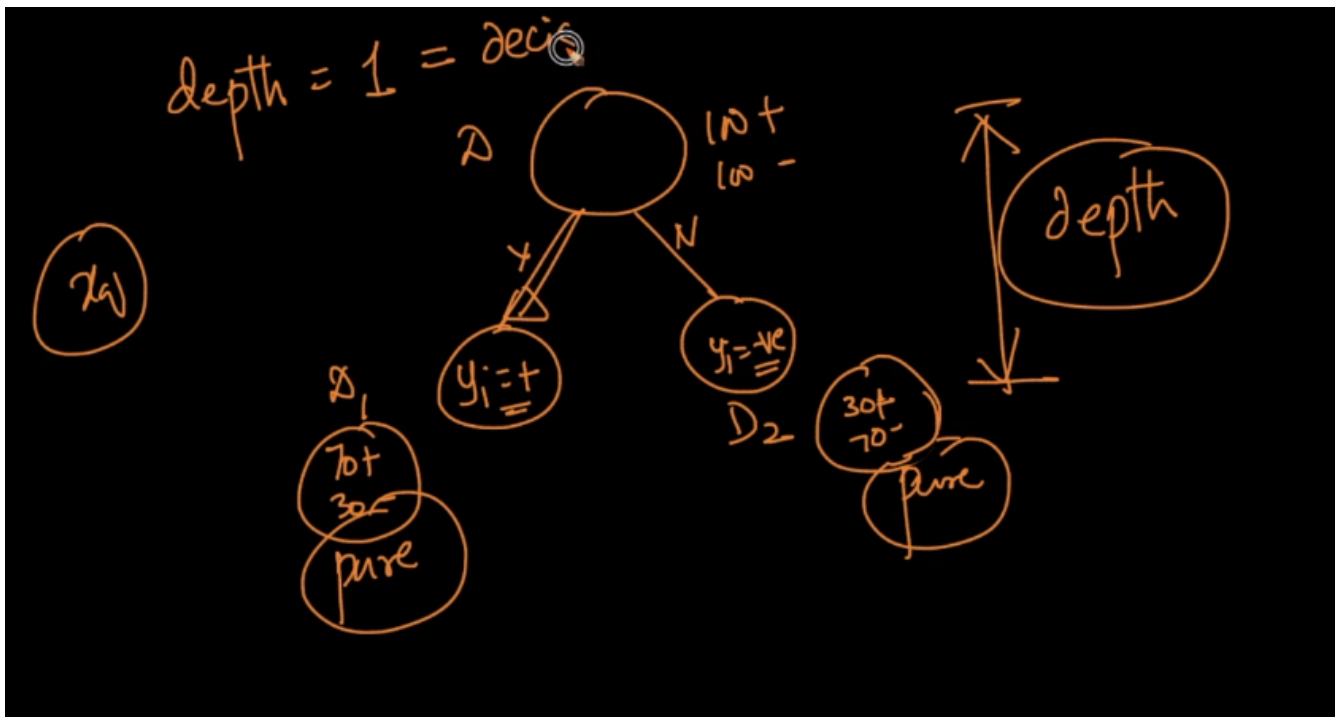
Each feature is converted according to the class of the feature by calculating the probability of occurrence of the feature variables.

Overfitting and Underfitting:

If there are outliers or noise in the data then the decision tree tend to fit these points and make the model to over fit the data.

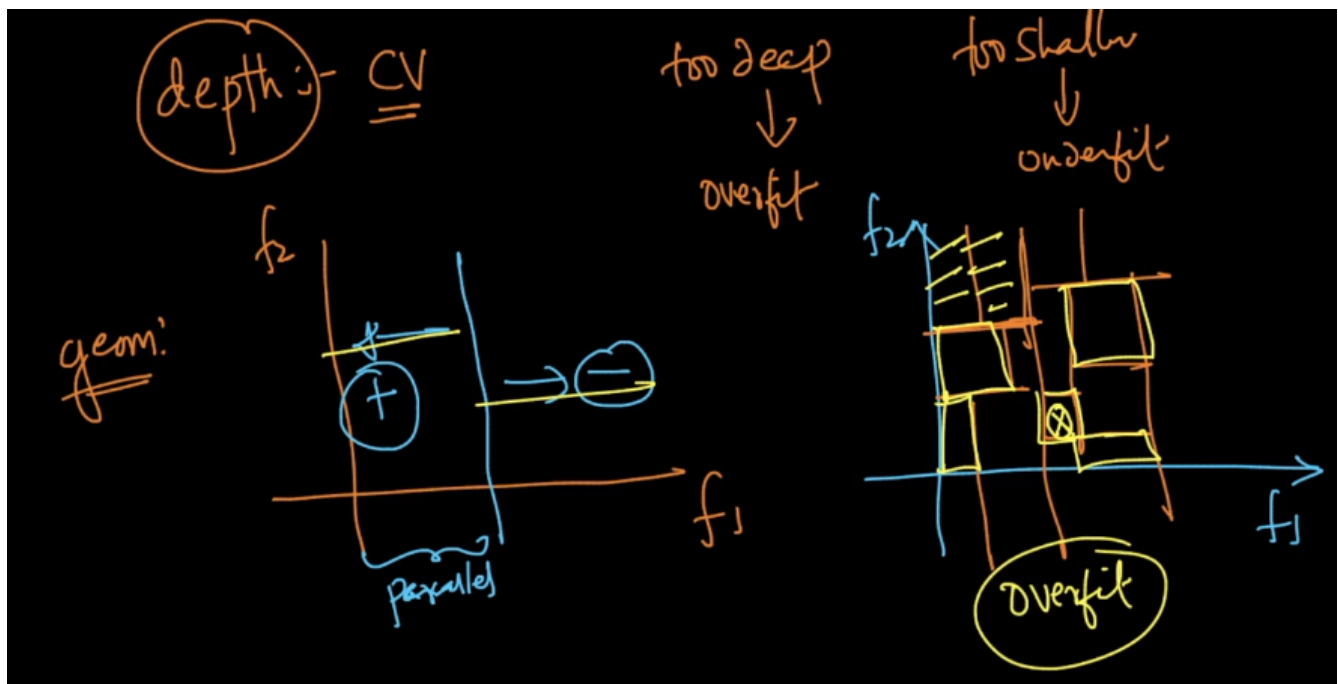
A decision stump is nothing but under fitting the data with less depth.





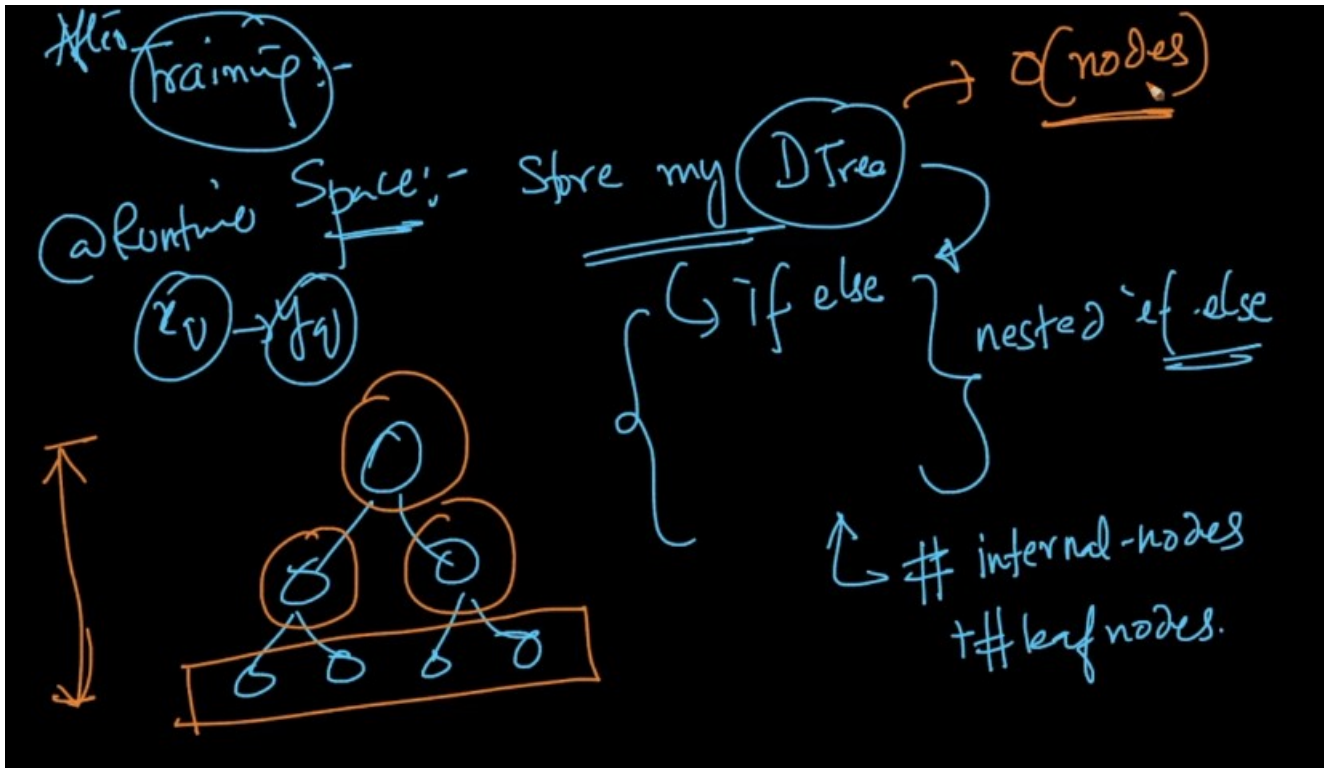
Depth is calculated using cross – validation

Visualizing Over-fitting and under-fitting:

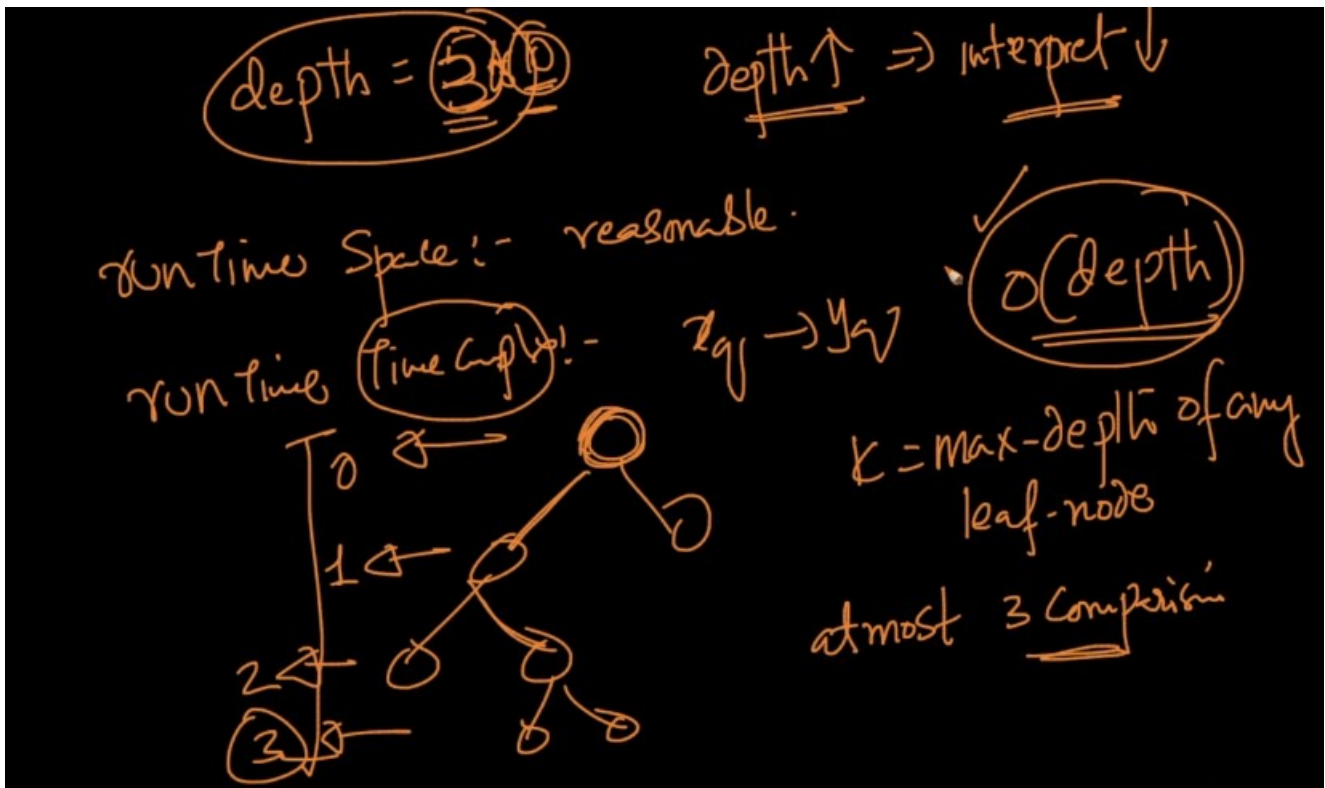


Train and run time complexity:

Converting a decision tree to nested if – else conditions can save space.



At max. a decision tree is trained to be 5 – 10 levels of depth.



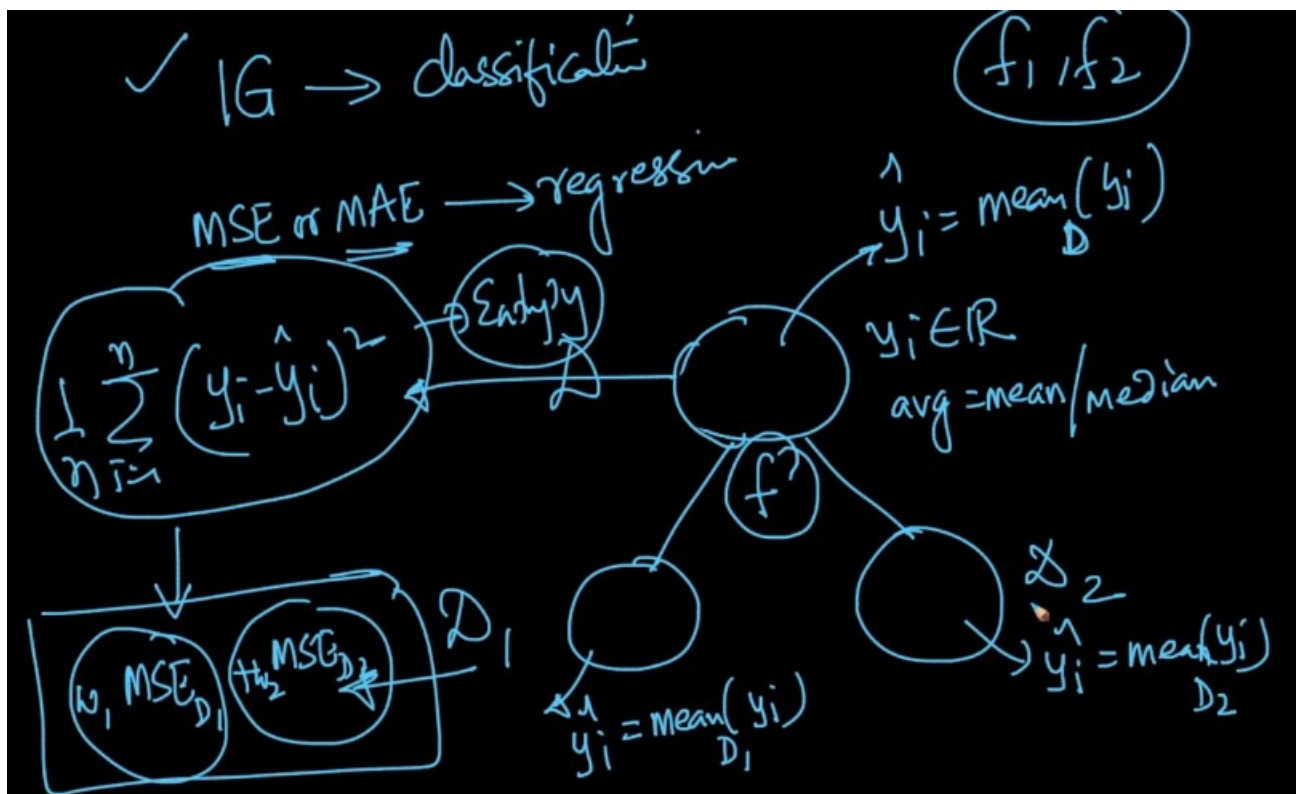
Decision tree can handle Large data, dimensionality is small or reasonable, low latency.



DT :- large data, dim is small  
 low latency  $\rightarrow O(\text{depth})$   
RF, GBDT  
 popular in internet applications

Regression using Decision Trees:

Instead of using Information gain we use Mean square error (or) mean absolute deviation is used to make regression trees.





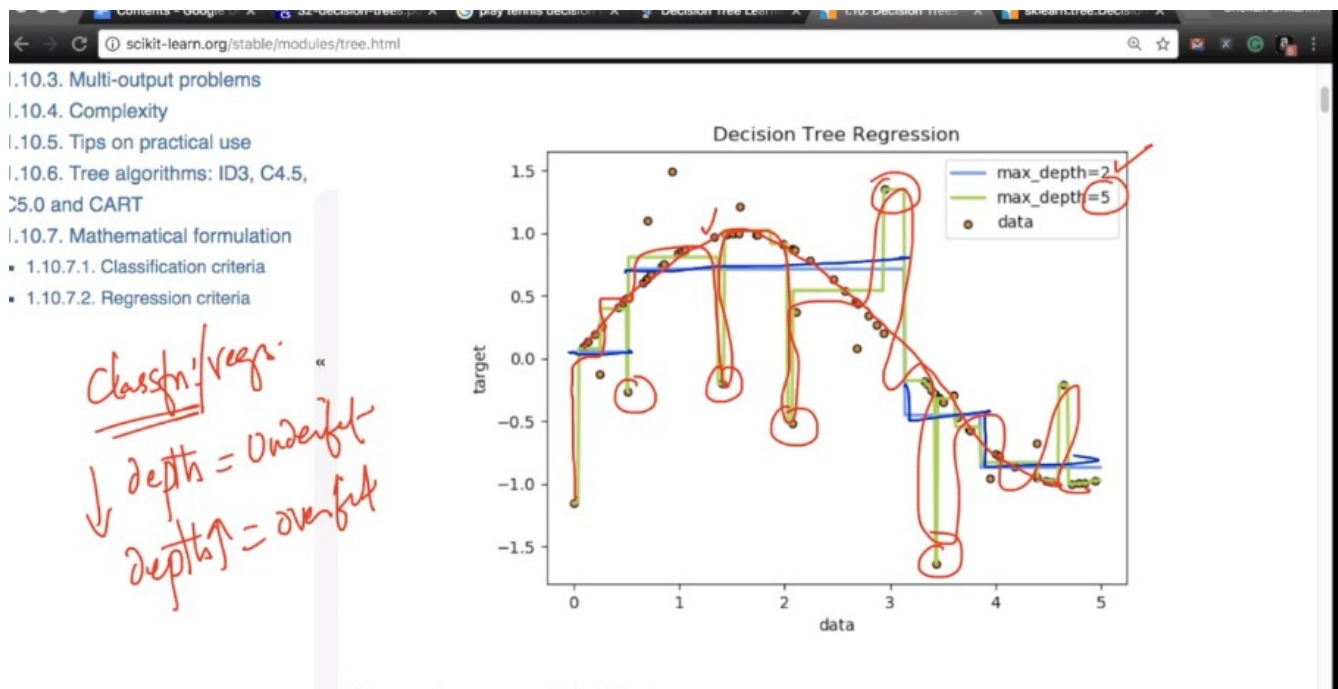
Classif:  $f_i$  :- reducing entropy :- "0"

regress:  $f_i$  :- reduce MSE  $\rightarrow$  "0"

$$\text{MSE}(y_i) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

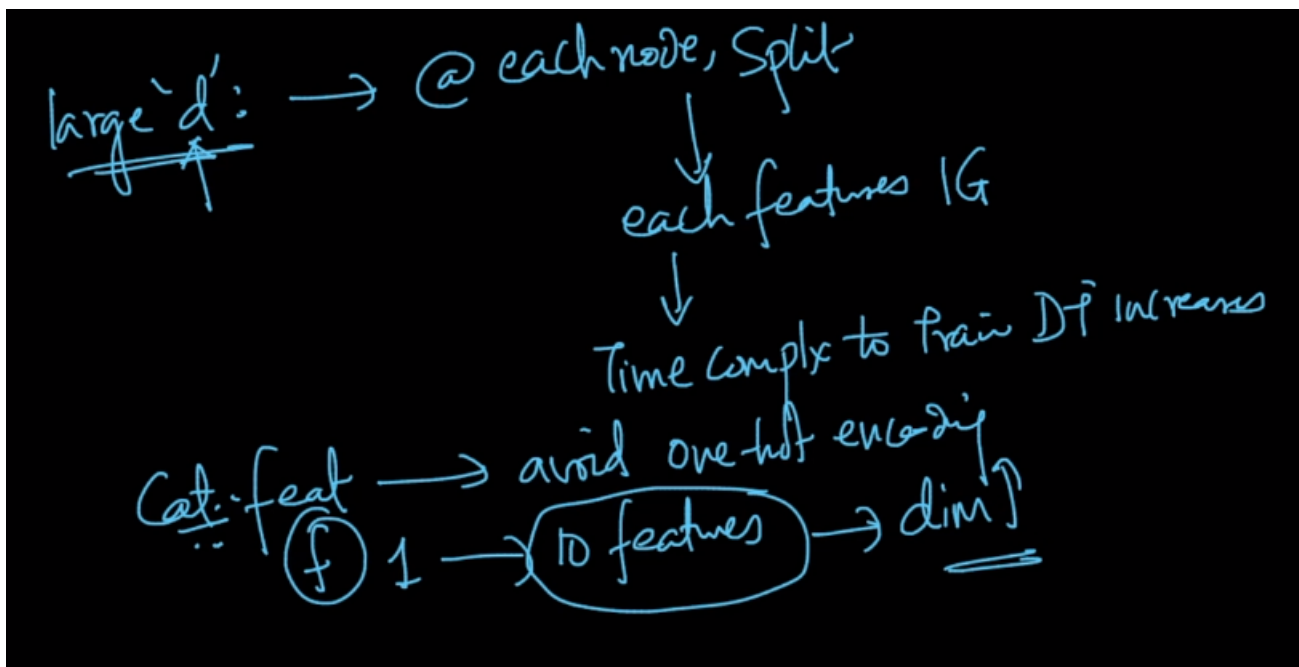
$$\text{MAE} = \text{MAI} \rightarrow "0"$$

$$\text{Median}(|y_i - \hat{y}_i|)$$



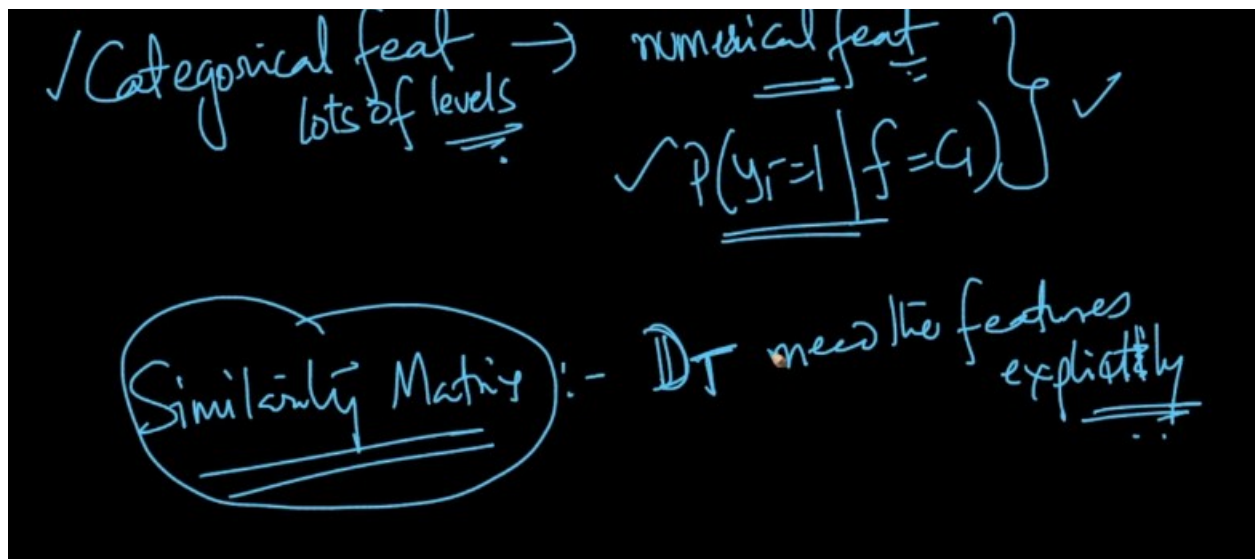
All the lines are axis parallel in decision tree model for regression / classification.

If dimensions are large, then the time training the data.  
One should avoid One - hot encoding.

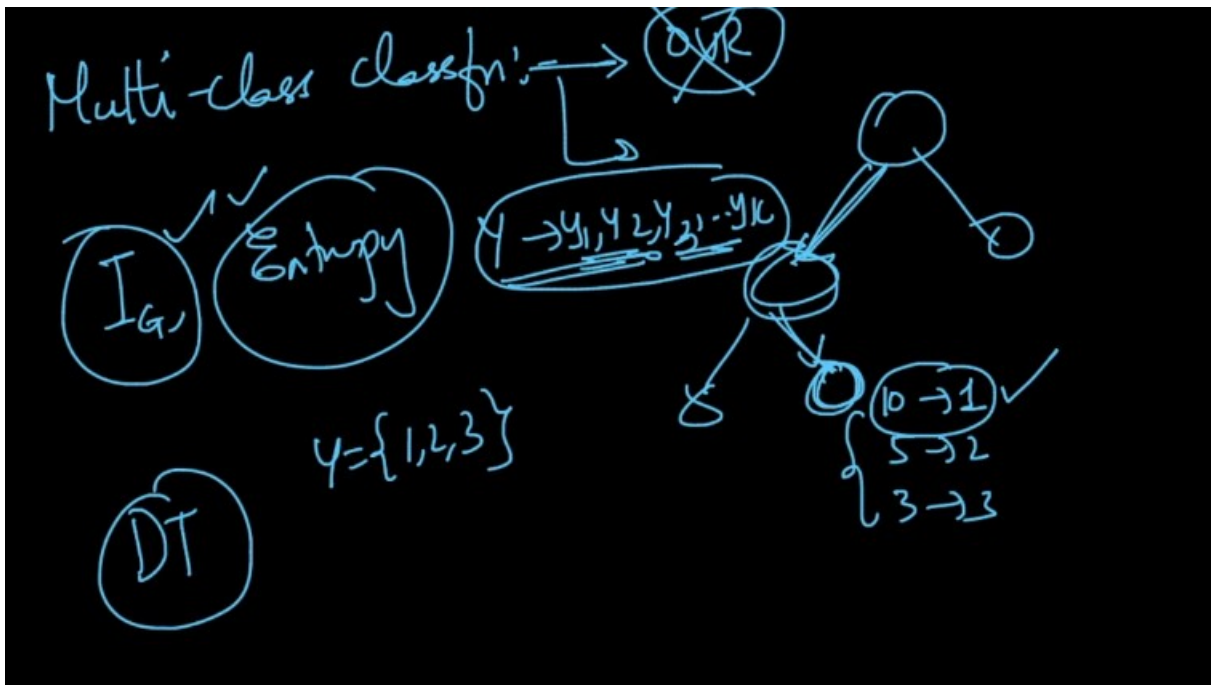


Converting into numerical features will save a lot of time.

Decision trees can read the data explicitly not in case is similarity matrix.

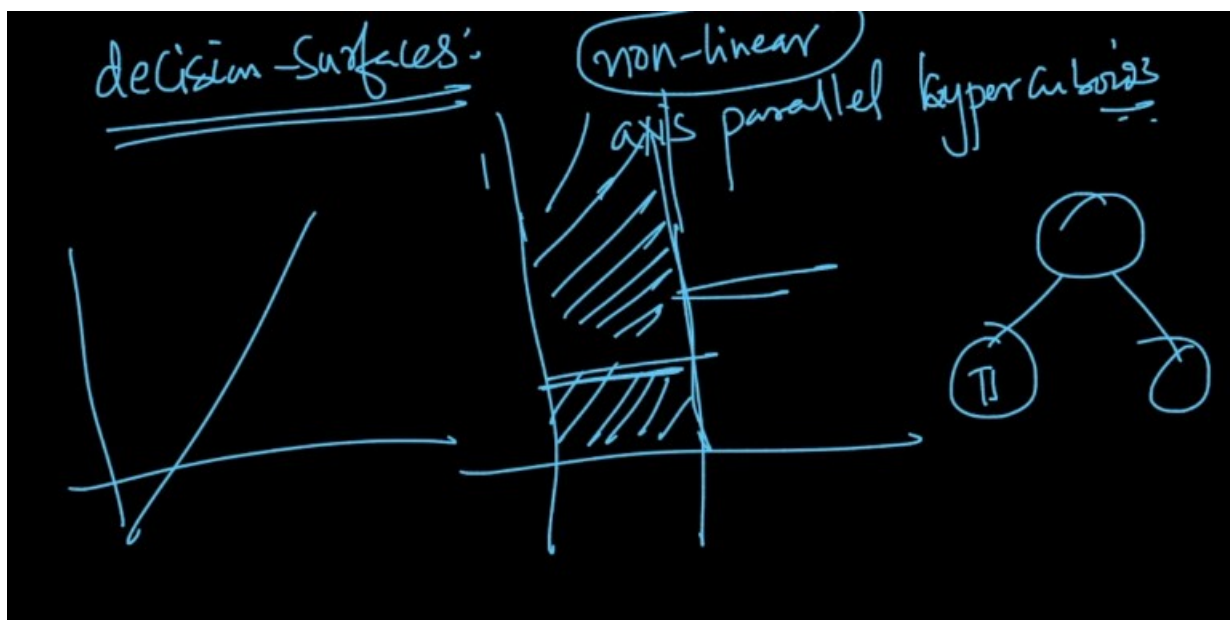


Decision trees can naturally be extended to multi – class classification.



Decision surface:

The decision that we get are non – linear. It basically divides the data into axis – parallel planes / hyper planes.



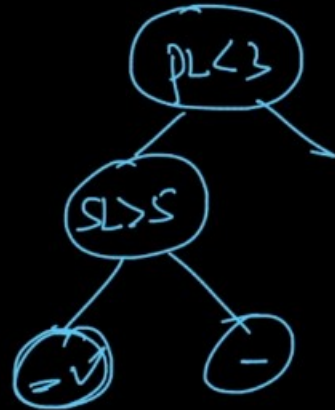
Feature interactions and decision trees:

There are feature interactions in decision trees to decide the class of the query point.  
 There are logical feature interactions in decision trees.

feature-interactions: (logical)

 $S_L, S_W, p_L, p_W$ 

(PLC3) AND (SLYS)



log-yes  
FT: →

$$f_i * f_j$$

$f_i \text{ AND } f_j$

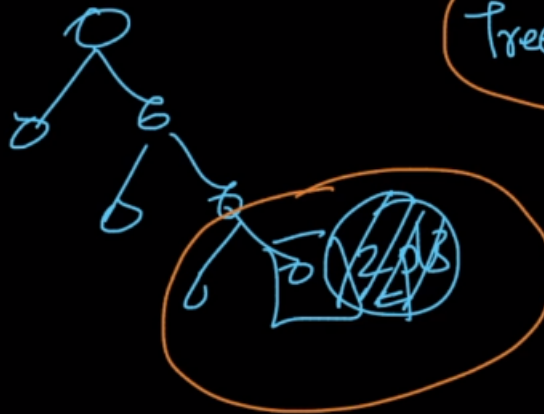
Outliers:

When depth is large then the model is prone to outliers.

## Outliers:

depth ↑ :- outlier will impact

## Tree unstable



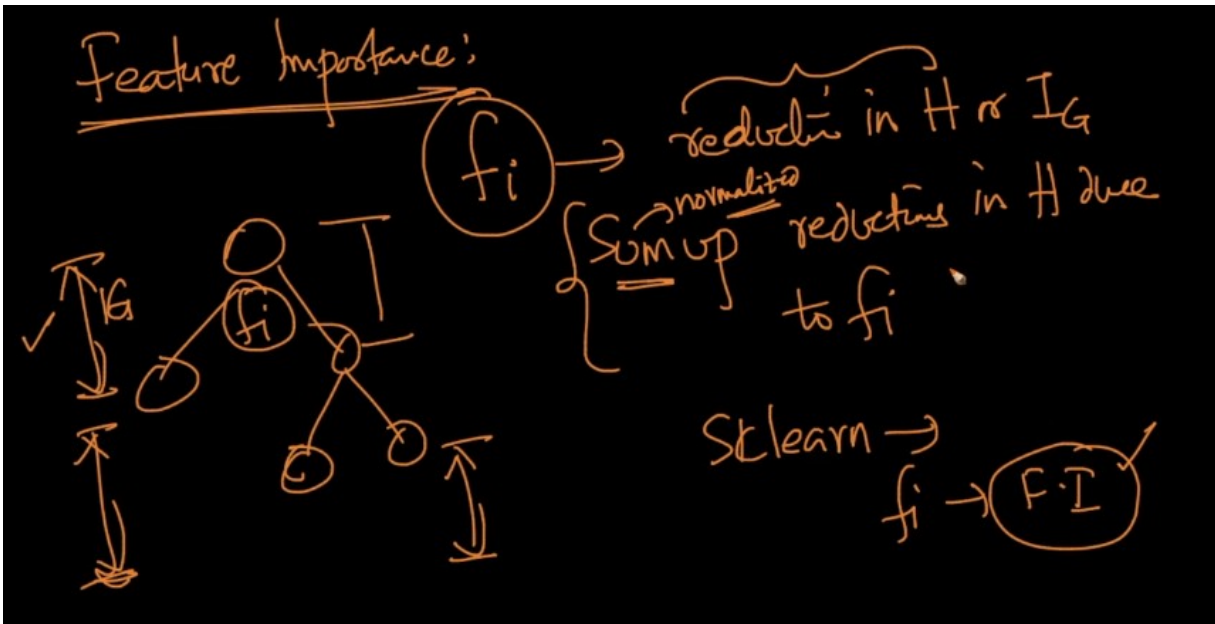
Interpret – ability:

Interpret – ability is very easy in decision trees.

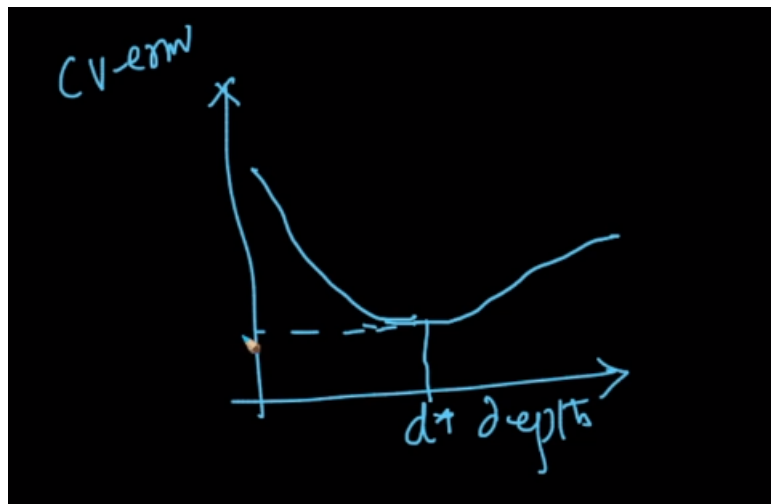
Feature importance:

We can sum up the reductions in entropy of each feature based in the importance of the feature.

If one feature occurs more than one time then we can conclude that feature is more important.



Exercise:





Ex: DT on Amazon Food reviews:

Try:

$\left\{ \begin{array}{l} \text{BoW} \\ \text{tfidf} \end{array} \right\} \rightarrow \underline{d}$  is large (10's of thousand)

$\left\{ \begin{array}{l} \underline{\text{w2v}} \rightarrow \text{avg} \\ \text{tfidf} \end{array} \right\} \rightarrow \text{200 dim}$

DT on w2v  $\rightarrow$  depth  $\rightarrow$  CV