

# ML Part-3

Decision Tree  
Shiva Reddy

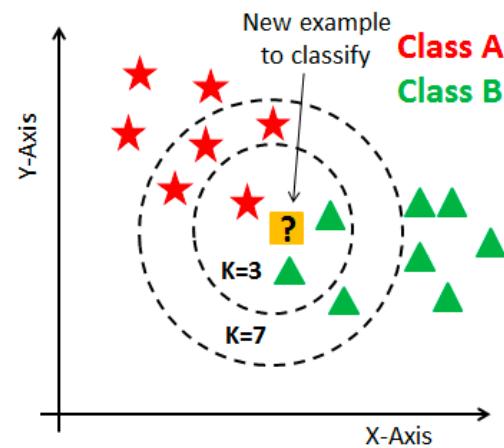
# Imp criteria while choosing best ML Model ?

Scalability

Performance

Interpretability

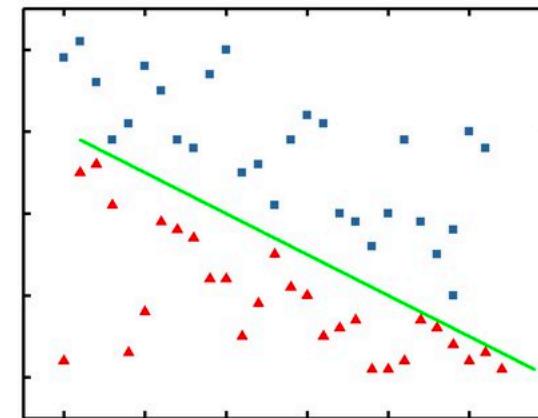
# K-Nearest Neighbor



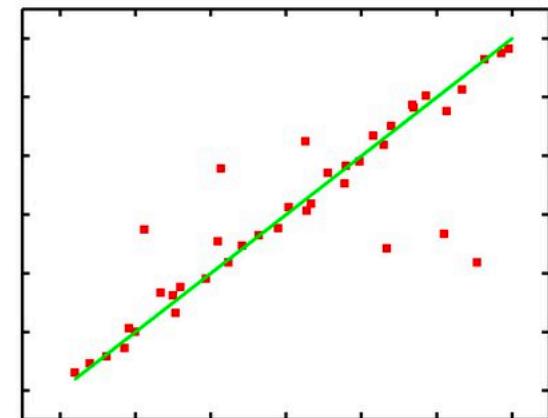
## Imp Pointers on KNN :

Not Scalable

# Linear Models



(a) Logistic Regression



(b) Linear Regression

## Imp Pointers on Linear Models :

Decision boundary : Linear

Highly Interpretable models

E.g. Sales =  $W_1 * \text{Digital Media} + W_2 * \text{TV} + W_3 * \text{Offline Adv}$

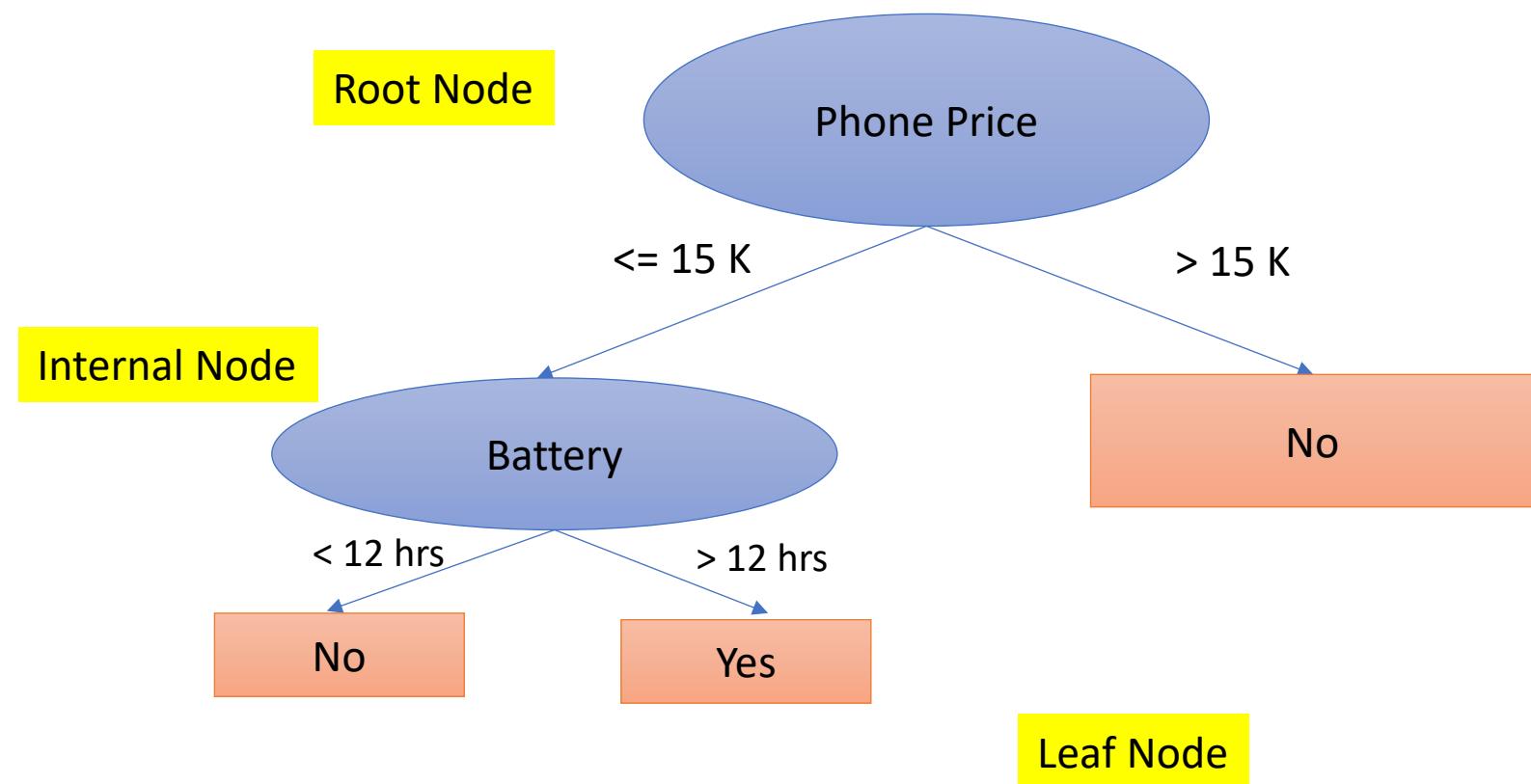
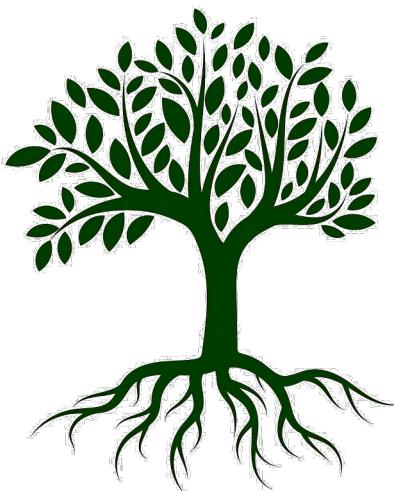
# Decision Tree

- How to take a decision of which phone to buy ?

Each path from the root of the DT to a leaf can be interpreted as a decision rule

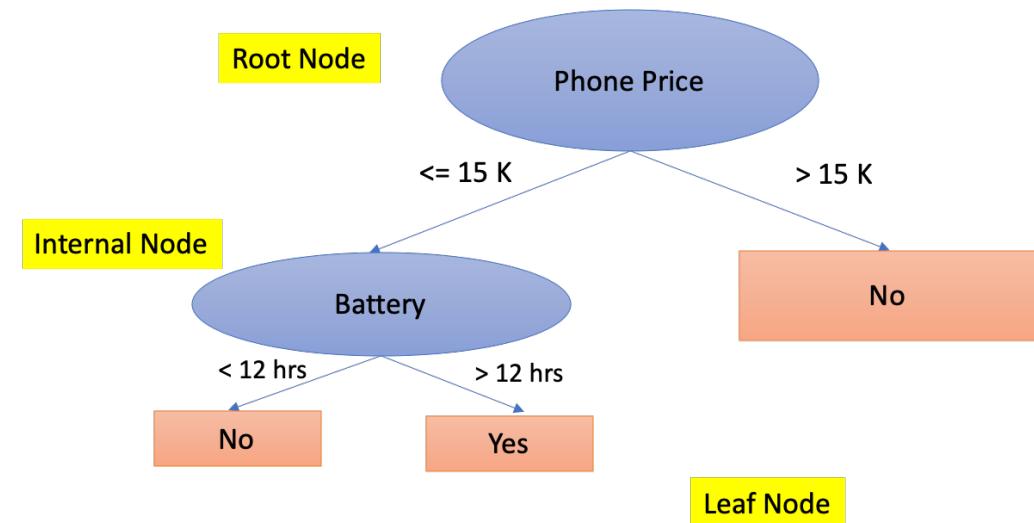
**Imp criteria:**

1. Budget Phone : 15 K
2. Battery : atleast 12 hrs

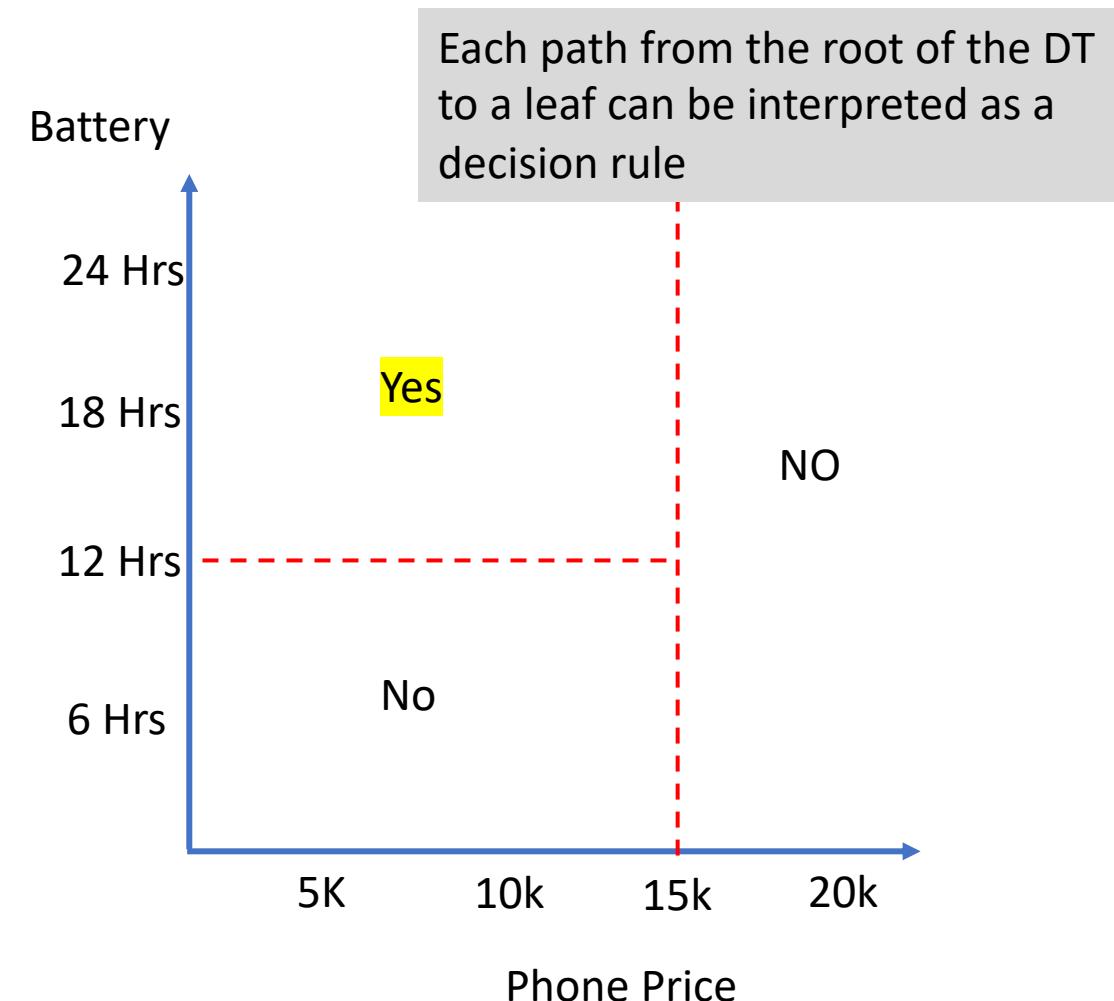
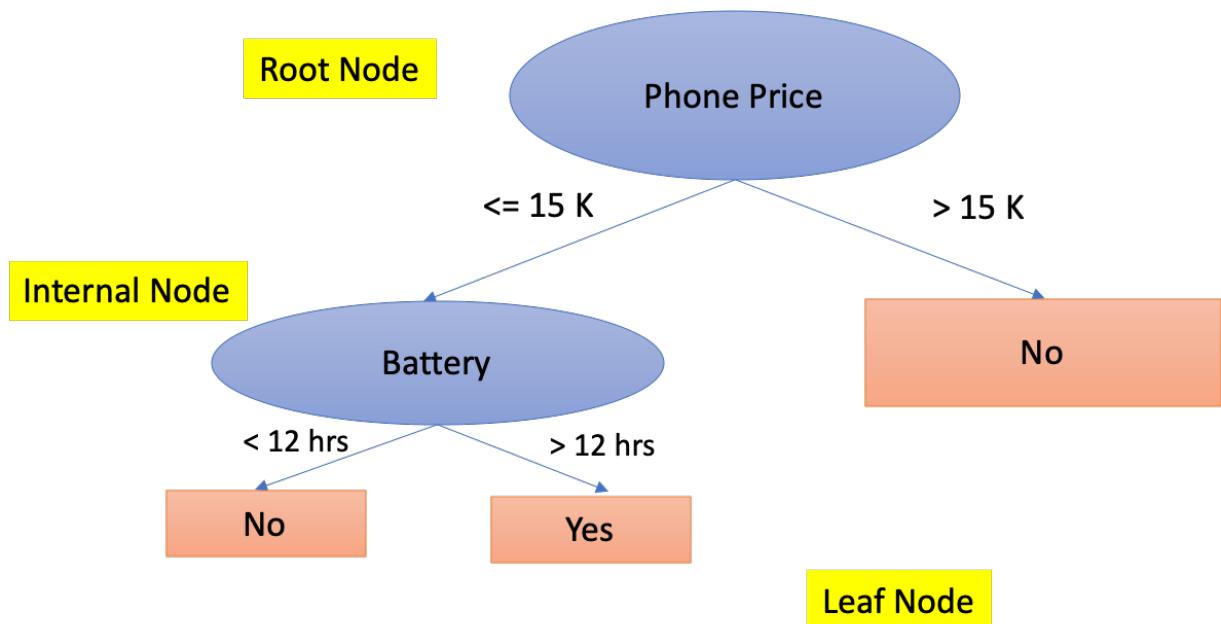


# DT Terminology

- **Root Node :** This is first node while building a DT ,where all data is considered.
- **Internal Node :** All nodes after root node & before leaf node are internal Nodes  
Root Node & Internal nodes are known as '**Decision Node**'.
- **Leaf Node :** Last node is Leaf node/Terminal node.  
At leaf node ,we take final decision /labelling.



# Graphical representation of DT



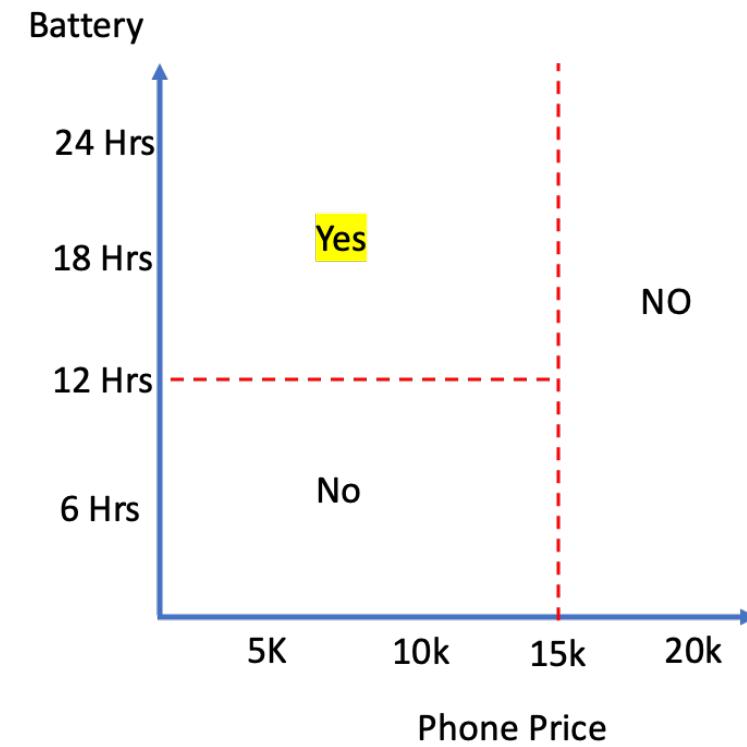
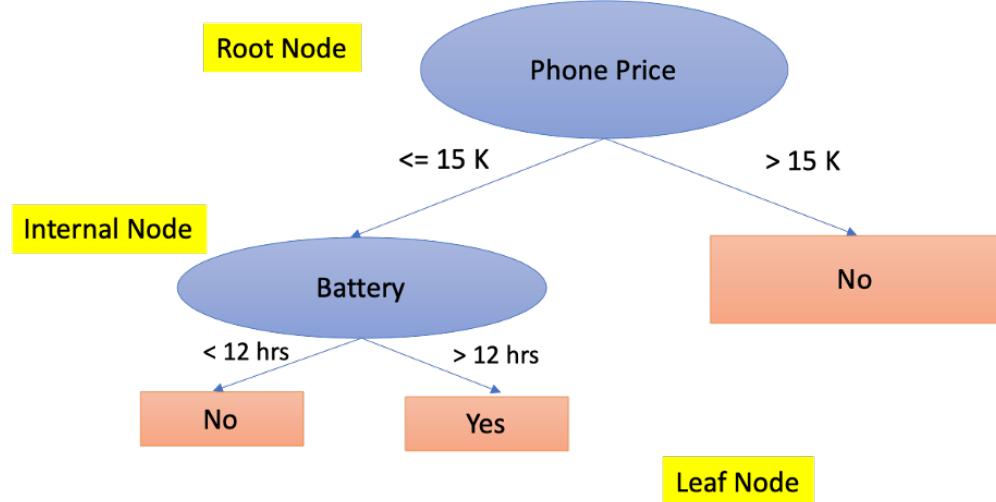
**Decision Boundary :**

Axis Parallel (Non-linear)

Geometric: set of axis parallel hyperplane that divides your whole region into cubes ,cuboid, hypercube.

# Decision Tree

- Decision tree is a tree based method that partition the feature space into a set of rectangles & then assign a constant value (mean/mode) to every region.



# How to built a Decision tree ?

- To built a decision tree, we start at the tree root and split the data on the feature that results in the largest information gain (IG).  
i.e Select the feature at root node which gives Max IG
- Information Gain : is a way to measure expected reduction in Entropy



# Entropy

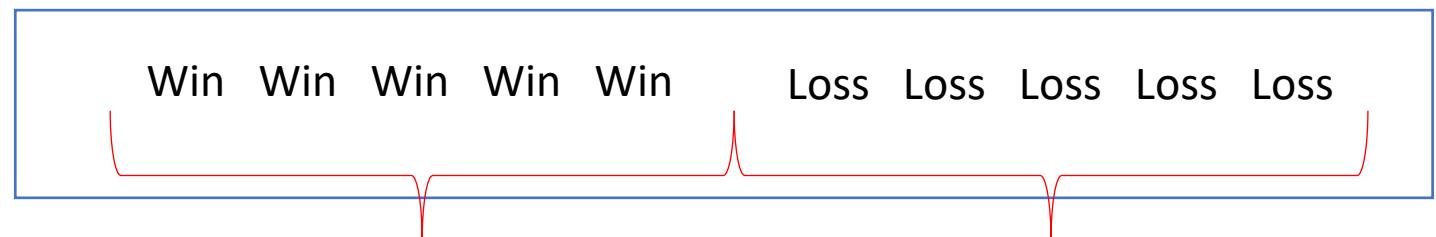
- Entropy is a way to measure impurity/uncertainty in data.
- Low Entropy --- > Good
  - The higher the entropy, the harder it is to draw any conclusions from that information.



# How to calculate Entropy

$$I_H = - \sum_{j=1}^c p_j \log_2(p_j)$$

C : Number of Classes in the label  
P<sub>j</sub>: Count of Class j / Total count



WIN: 5  
Prob 1 = 5/10  
= 0.5

Lost: 5  
Prob 2 = 5/10  
= 0.5

$$= - ( ( \text{prob\_win} * \log_2(\text{prob\_win}) ) + ( \text{prob\_lose} ) * \log_2(\text{prob\_lose}) )$$

$$= -((0.5) * \log_2(0.5) + (0.5) * \log_2(0.5))$$

$$= -(-0.5 + (-0.5))$$

$$= -(-1) = 1$$

# Entropy Max/Min

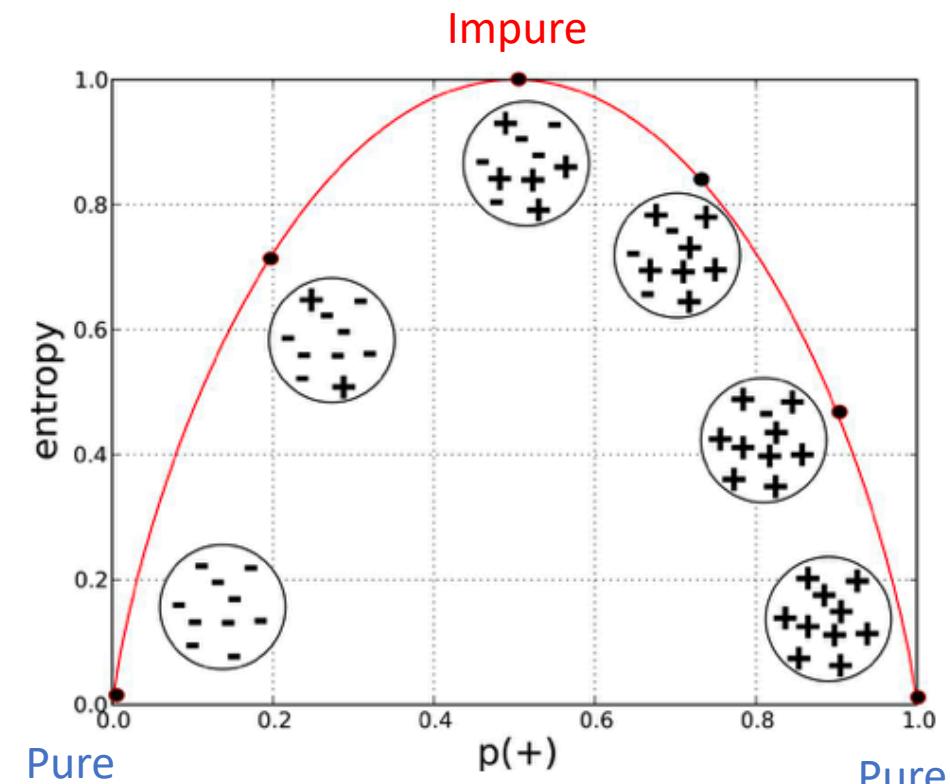
$$I_H = - \sum_{j=1}^c p_j \log_2(p_j)$$

$$= -1 (0 * \log_2(0) + 1 * \log_2(1))$$

Scenarios	#Win	Prob_win	#Lose	Prob_lose	Entropy
1	0	0	10	1	0
2	1	0.1	9	0.9	0.46
3	2	0.2	8	0.8	0.72
4	3	0.3	7	0.7	0.88
5	4	0.4	6	0.6	0.97
6	5	0.5	5	0.5	1
7	6	0.6	4	0.4	0.97
8	7	0.7	3	0.3	0.88
9	8	0.8	2	0.2	0.72
10	9	0.9	1	0.1	0.46
11	10	1	0	0	0

# Entropy Max/Min

Scenarios	#Win	Prob_win	#Lose	Prob_lose	Entropy
1	0	0	10	1	0
2	1	0.1	9	0.9	0.46
3	2	0.2	8	0.8	0.72
4	3	0.3	7	0.7	0.88
5	4	0.4	6	0.6	0.97
6	5	0.5	5	0.5	1
7	6	0.6	4	0.4	0.97
8	7	0.7	3	0.3	0.88
9	8	0.8	2	0.2	0.72
10	9	0.9	1	0.1	0.46
11	10	1	0	0	0



**Entropy is Min**

when all elements belong to one class

**Entropy is Max**

when all elements are equally probable

- Calculate entropy of below variable ‘Play’ :  
[Yes , Yes , Yes , Yes , Yes, No, No, Yes, No, Yes, Yes , Yes, No]

Count of yes	9	9/14
Count of No	5	5/14
Total	14	

$$\text{Entropy.} = -1 * (9/14 * \log_2(9/14) + 5/14 * \log_2(5/14))$$

# Information Gain

- Information Gain : is a way to measure expected reduction in Entropy

**Information Gain** = Parent Entropy - Weighted Average of Child Entropy

Entropy Before Split - Entropy After Split

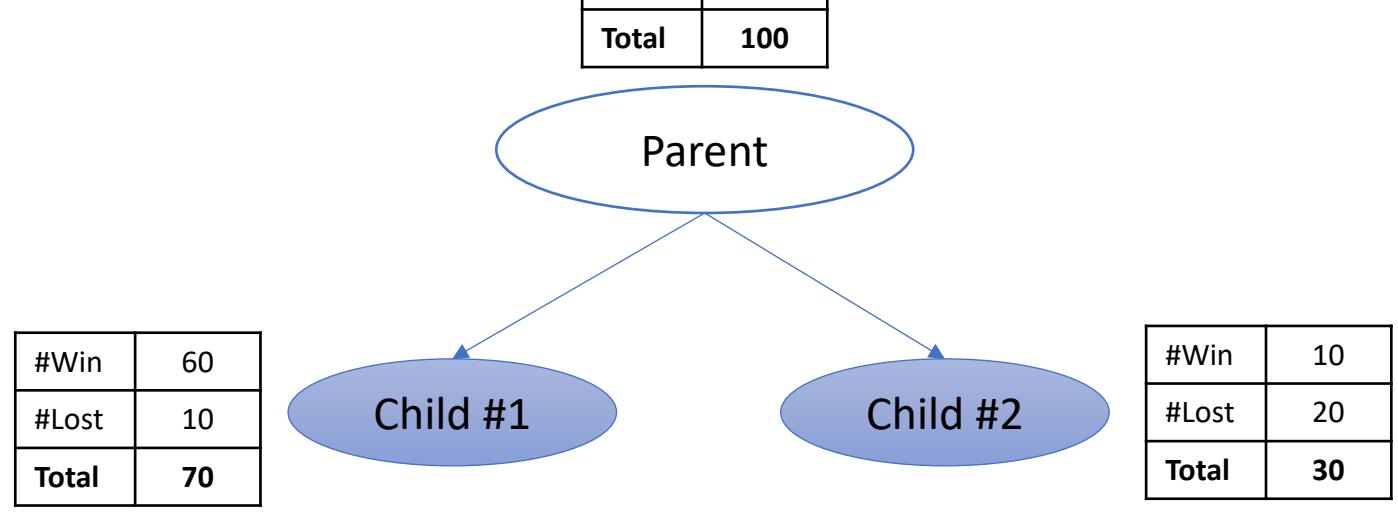
$$= 0.88 - 0.68 = 0.2$$

$$IG(D_p, x_i) = I(D_p) - \frac{N_{left}}{N_p}I(D_{left}) - \frac{N_{right}}{N_p}I(D_{right})$$

- $IG$ : Information Gain
- $x_i$ : feature to perform the split
- $N_p$ : number of samples in the parent node
- $N_{left}$ : number of samples in the left child node
- $N_{right}$ : number of samples in the right child node
- $I$ : impurity
- $D_p$ : training subset of the parent node
- $D_{left}$ : training subset of the left child node
- $D_{right}$ : training subset of the right child node

#Win	70
#Lost	30
Total	100

Entropy before split :0.88



Entropy @ left node : 0.59

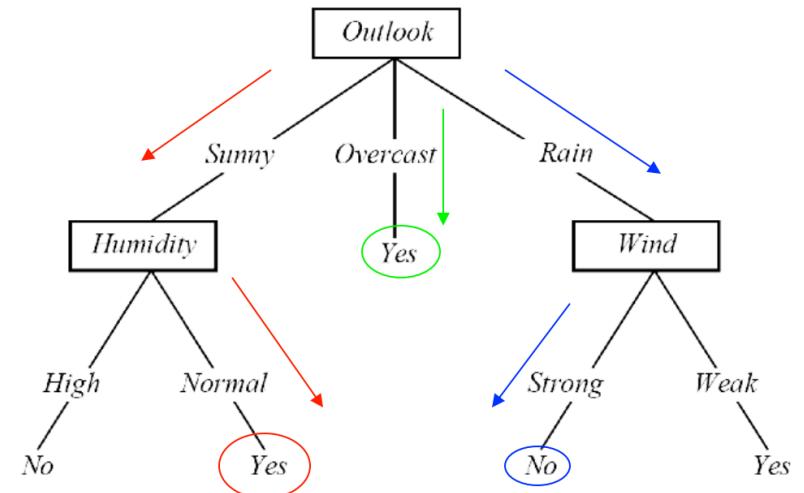
Entropy @ right node : 0.91

$$\begin{aligned} \text{Entropy after split} &= 70/100 * \text{Entropy @ left node} + 30/100 * \text{Entropy @ right node} \\ &= 0.7 * 0.59 + 0.3 * 0.91 = 0.68 \end{aligned}$$

# Play Tennis Dataset

Day	Outlook	Temp.	Humidity	Wind	Play?
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

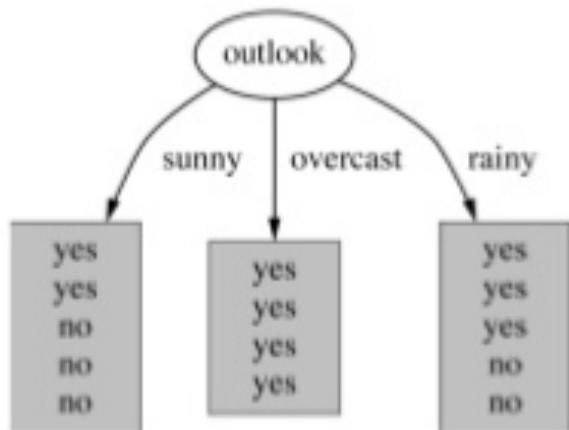
Each path from the root of the DT to a leaf can be interpreted as a decision rule.



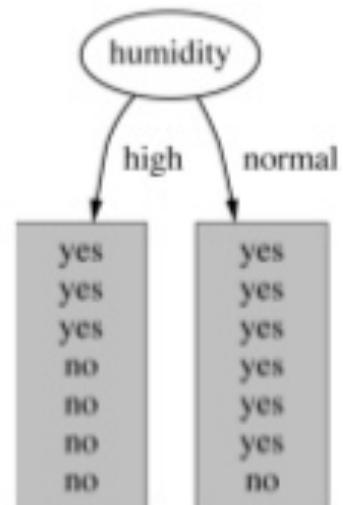
IF *Outlook = Sunny AND Humidity = Normal THEN Playtennis = Yes*  
IF *Outlook = Overcast THEN Playtennis = Yes*  
IF *Outlook = Rain AND Wind = Strong THEN Playtennis = No*

# Which attribute to select?

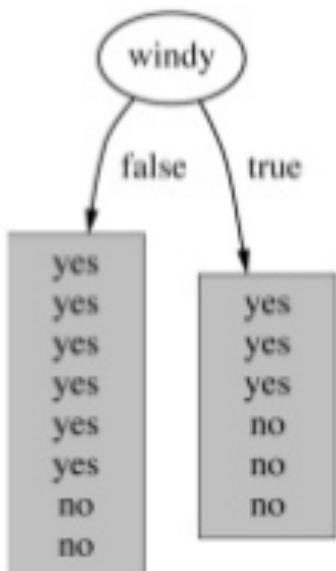
IG = 0.25



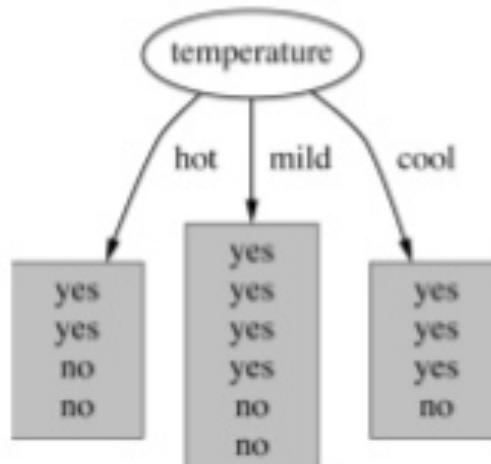
IG = 0.16



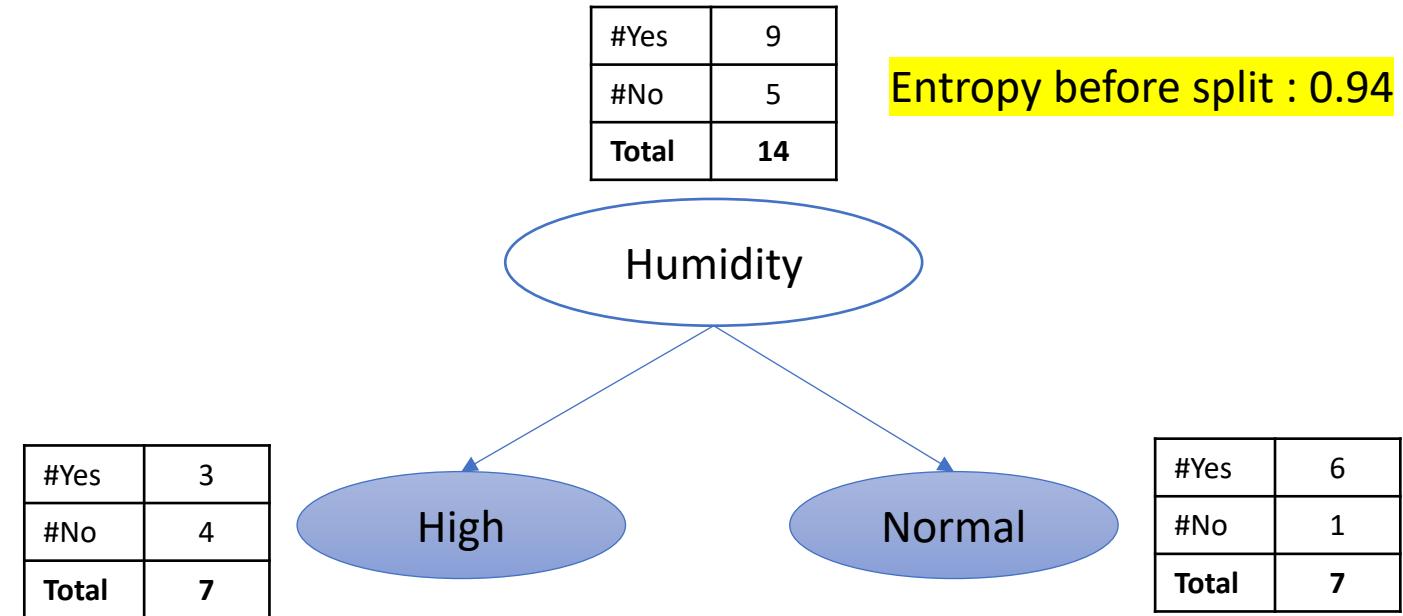
IG = 0.04



IG = 0.051



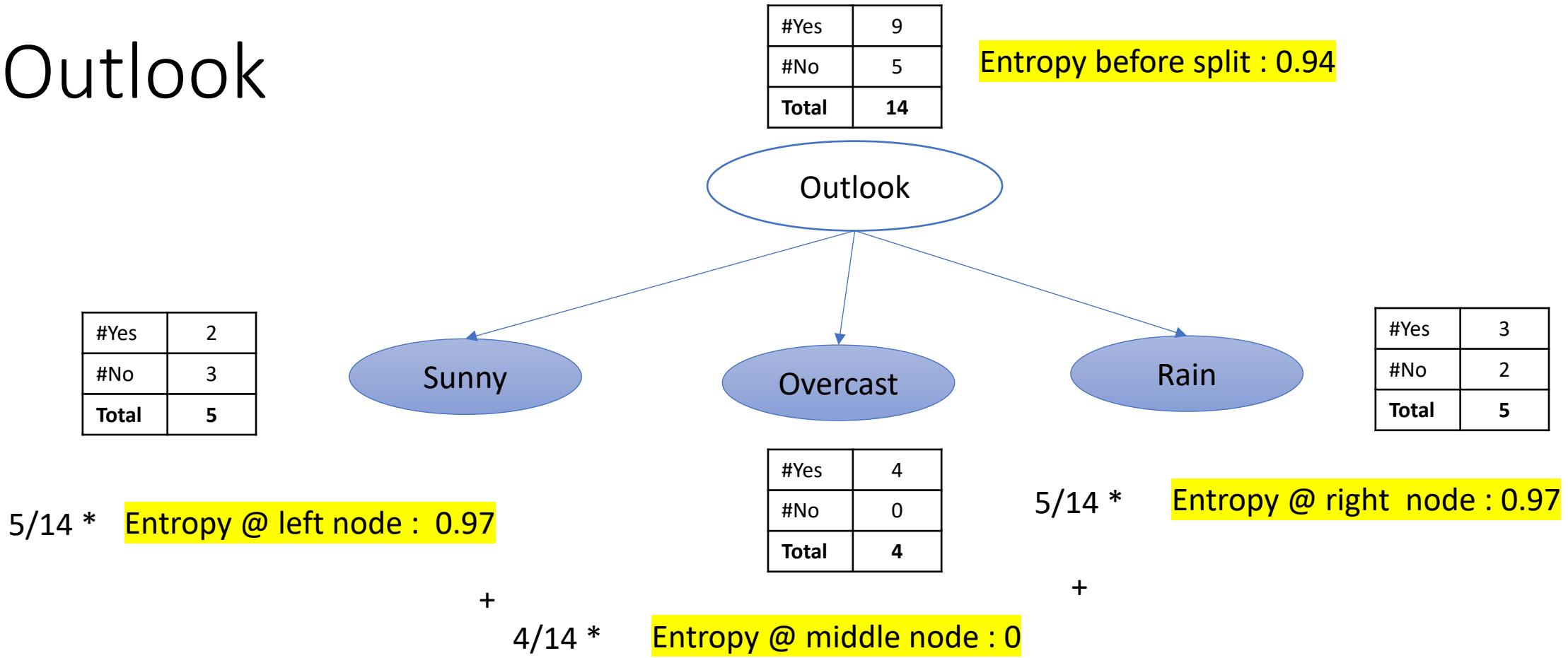
# Humidity



$$\text{Entropy @ after split} = \frac{7}{14} * \text{Entropy @ left node : } 0.98 + \frac{7}{14} * \text{Entropy @ right node : } 0.59$$

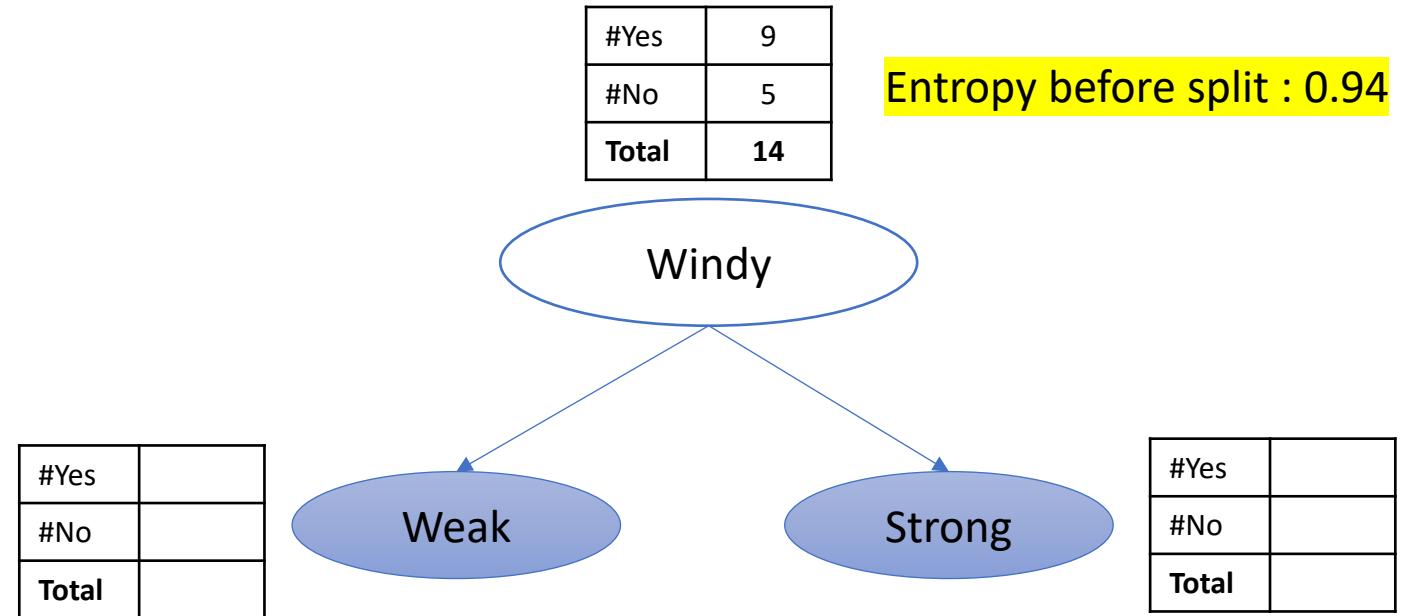
Entropy before split	Entropy after split	Information Gain
0.94	0.78	0.16

# Outlook



Entropy before split	Entropy after split	Information Gain
0.94	0.69	0.25

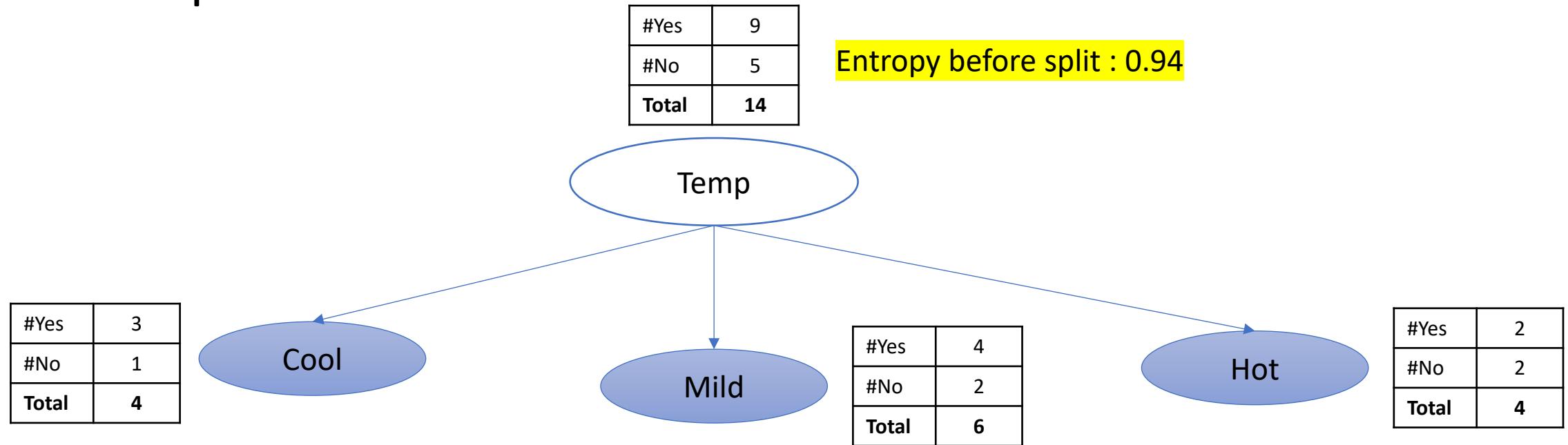
# Windy



	Weight avg @ left node	Entropy @ left node	Weight avg @ right node	Entropy @ right node
Entropy after split				

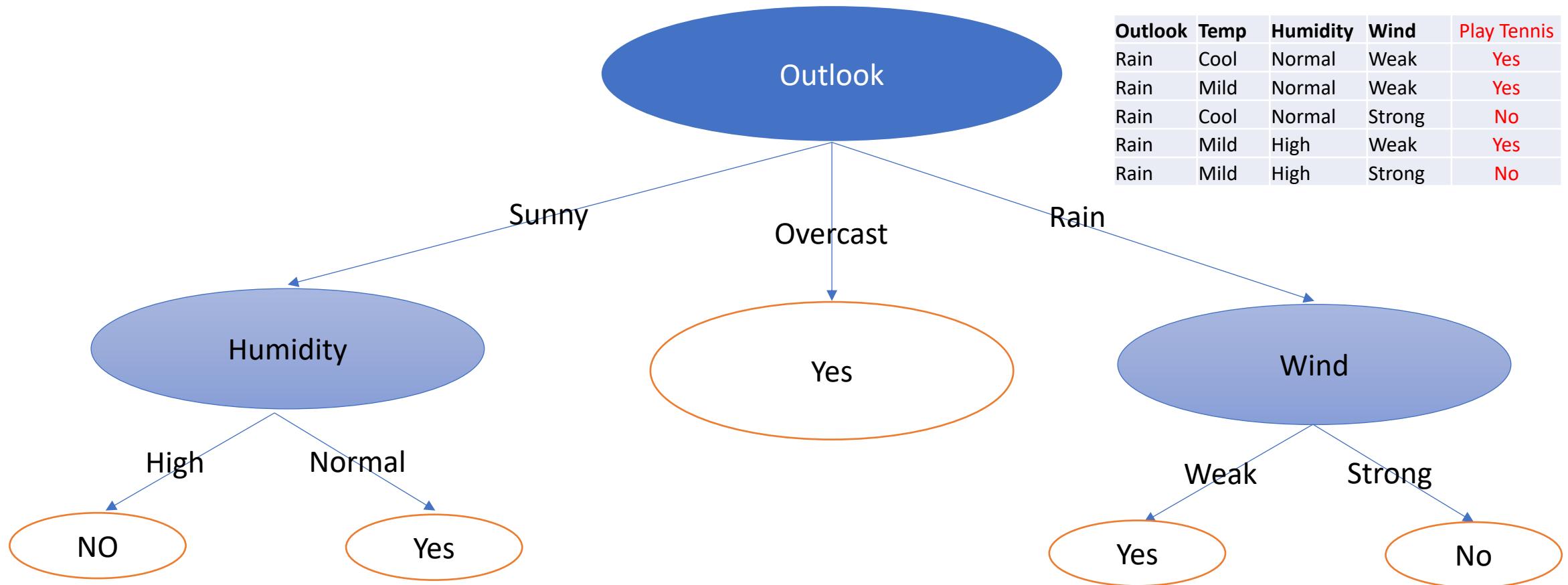
Entropy before split	Entropy after split	Information Gain
0.94		0.04

# Temp



	Weight avg @ left node	Entropy @ left node	Weight avg @ middle node	Entropy @ middle node	Weight avg @ right node	Entropy @ right node
Entropy after split	0.28	0.81	0.42	0.91	0.28	1

Entropy before split	Entropy after split	Information Gain
0.94	0.889	0.051



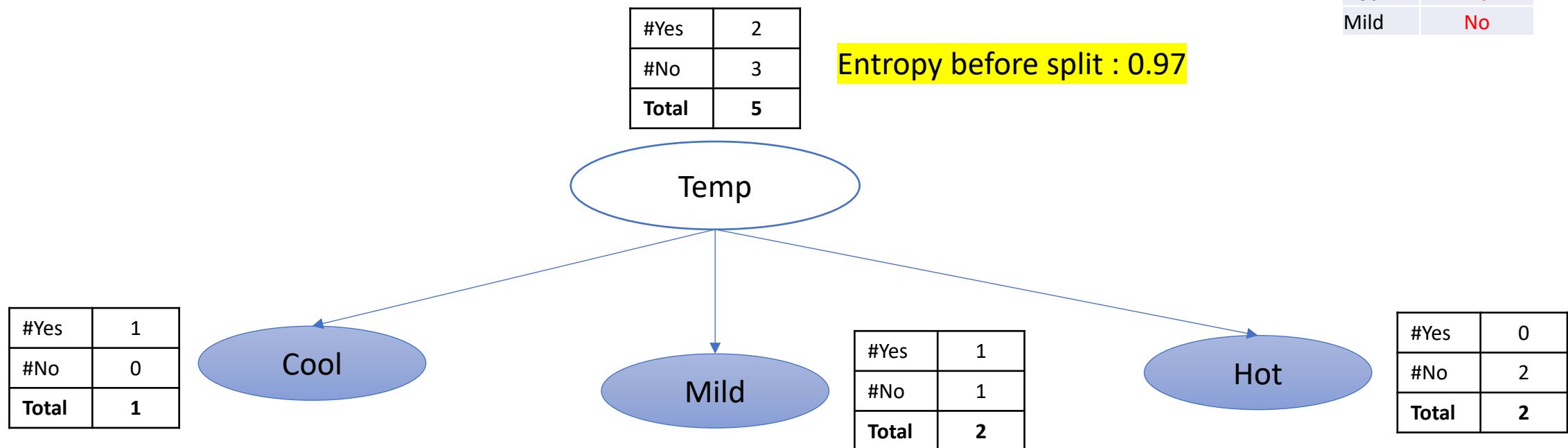
Outlook	Temp	Humidity	Wind	Play Tennis
Sunny	Cool	Normal	Strong	Yes
Rain	Mild	High	Strong	No

Pure node : All elements belongs to one class  
 Entropy = 0

Outlook	Temp	Humidity	Wind	Play Tennis
Rain	Cool	Normal	Weak	Yes
Rain	Mild	Normal	Weak	Yes
Rain	Cool	Normal	Strong	No
Rain	Mild	High	Weak	Yes
Rain	Mild	High	Strong	No

# Temp : 2<sup>nd</sup> Level

Temp	Play Tennis
Cool	Yes
Mild	Yes
Hot	No
Hot	No
Mild	No

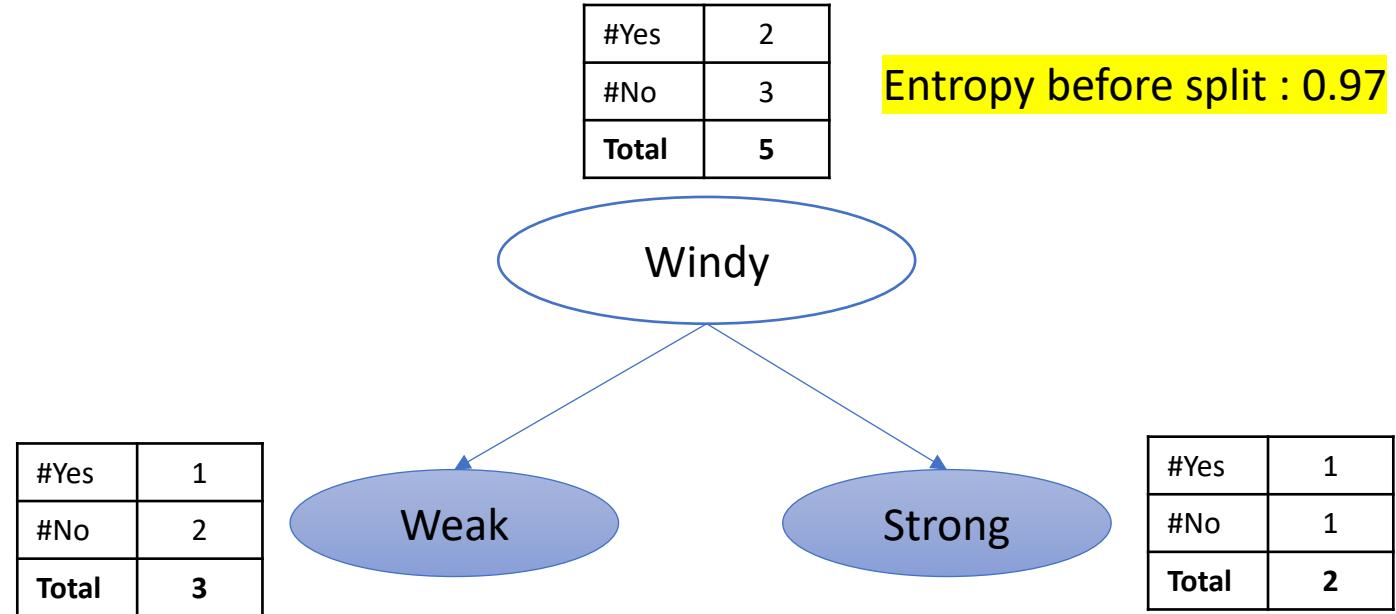


	Weight avg @ left node	Entropy @ left node	Weight avg @ middle node	Entropy @ middle node	Weight avg @ right node	Entropy @ right node
Entropy after split	0.2	0	0.4	1	0.4	0

Entropy before split	Entropy after split	Information Gain
0.97	0.4	0.57

# Windy : 2<sup>nd</sup> Level

Wind	Play Tennis
Weak	Yes
Strong	Yes
Weak	No
Strong	No
Weak	No



	Weight avg @ left node	Entropy @ left node	Weight avg @ right node	Entropy @ right node
Entropy after split	0.6	0.918	0.4	1

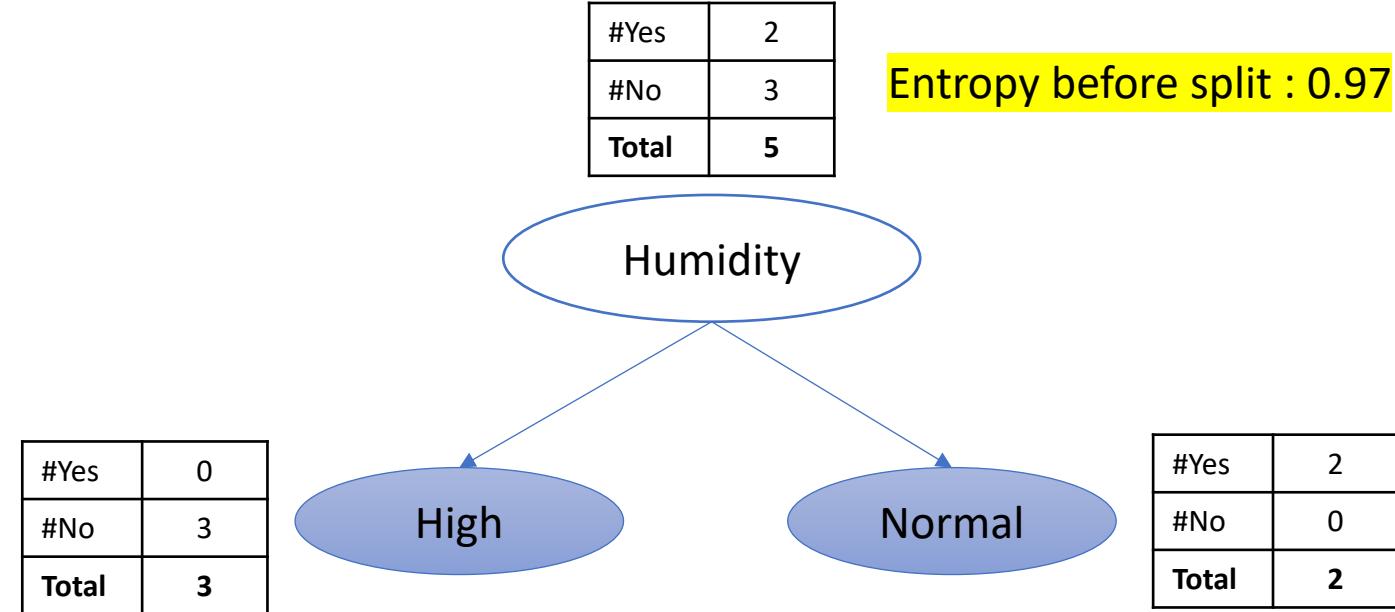
Entropy before split	Entropy after split	Information Gain
0.97	0.9508	0.019

# Humidity: 2<sup>nd</sup> Level

Outlook	Temp	Humidity	Wind	Play Tennis
Sunny	Cool	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Sunny	Mild	High	Weak	No

#Yes	2
#No	3
Total	5

Entropy before split : 0.97



Entropy @ after split =

Entropy @ left node :

+

Entropy @ right node :

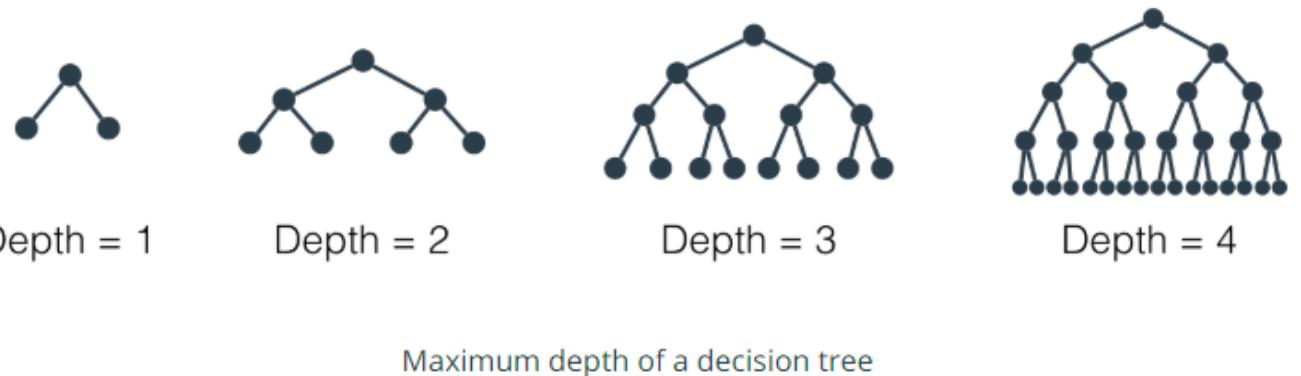
Entropy before split	Entropy after split	Information Gain
0.97	0	0.97

# DT Algorithm

1. Start at the root node : Select Feature which gives Max IG
2. Split the data at root node into different nodes using the feature  $x_i$  which has maximize information gain.
3. Assign the subset data to all new child nodes after splitting.
4. Repeat steps 1 and 2 for each new child node
  - Till : Stop if leaf nodes are pure or early stopping criteria is satisfied

# Stopping criteria

- ❖ **Leaf nodes are Pure** : All elements belong to one class
- ❖ **Depth of Tree** : A maximal node depth is reached
  - Depth of Tree is small → Model Underfit
  - Depth of Tree is large → Model Overfits
  - if last node is not a pure node then majority is taken as decision



## ❖ **min\_samples\_split :**

The minimum number of samples required to split an internal node

## ❖ **min\_samples\_leaf :**

The minimum number of samples required to be at a leaf node.

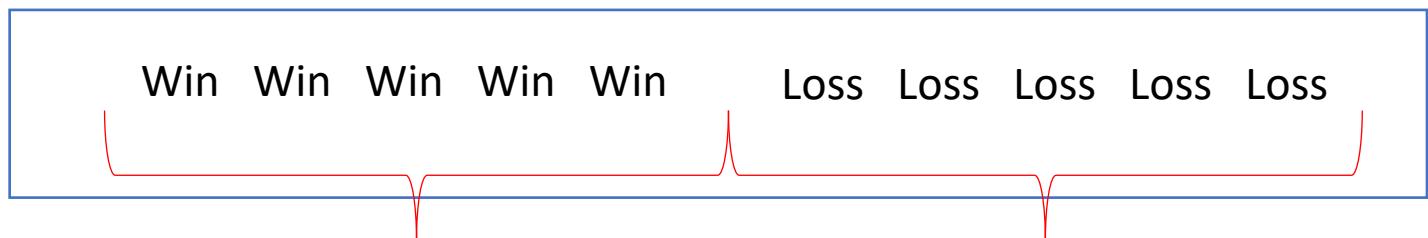
# Gini Index

Gini Impurity tells us what is the probability of misclassifying an observation.  
Gini impurity can be considered as an alternative for the entropy method.

Note that **the lower the Gini the better the split**. In other words the lower the likelihood of misclassification.

$$Gini = 1 - \sum_{i=1}^C (p_i)^2$$

C : Number of Classes in the label  
Pi: Count of Class j / Total count



WIN: 5  
Prob 1 = 5/10  
= 0.5

Lost: 5  
Prob 2 = 5/10  
= 0.5

$$\begin{aligned} \text{Gini Index} &= 1 - (\text{prob(win)}^2 + \text{prob(loss)}^2) \\ &= 1 - (0.25 + 0.25) \\ &= 1 - 0.5 = 0.5 \end{aligned}$$

# Which one to use : Entropy vs Gini Index

Entropy

$$I_H = - \sum_{j=1}^c p_j \log_2(p_j)$$

C : Number of Classes in the label  
Pj: Count of Class j / Total count

log computation

Takes more time ,more complex

Gini Index

$$Gini = 1 - \sum_{i=1}^C (p_i)^2$$

C : Number of Classes in the label  
Pi: Count of Class j / Total count

square computation

- Since the computing square is cheaper than the logarithmic function we prefer Gini impurity over entropy.

# Variants of Decision Tree

- CART Classification & Regression Tree) :
  - Uses Gini Index(Classification) as metrics
  - Uses MAE (Regression)
- ID3 (Iterative Dichotomiser 3) :
  - Uses Entropy as metrics

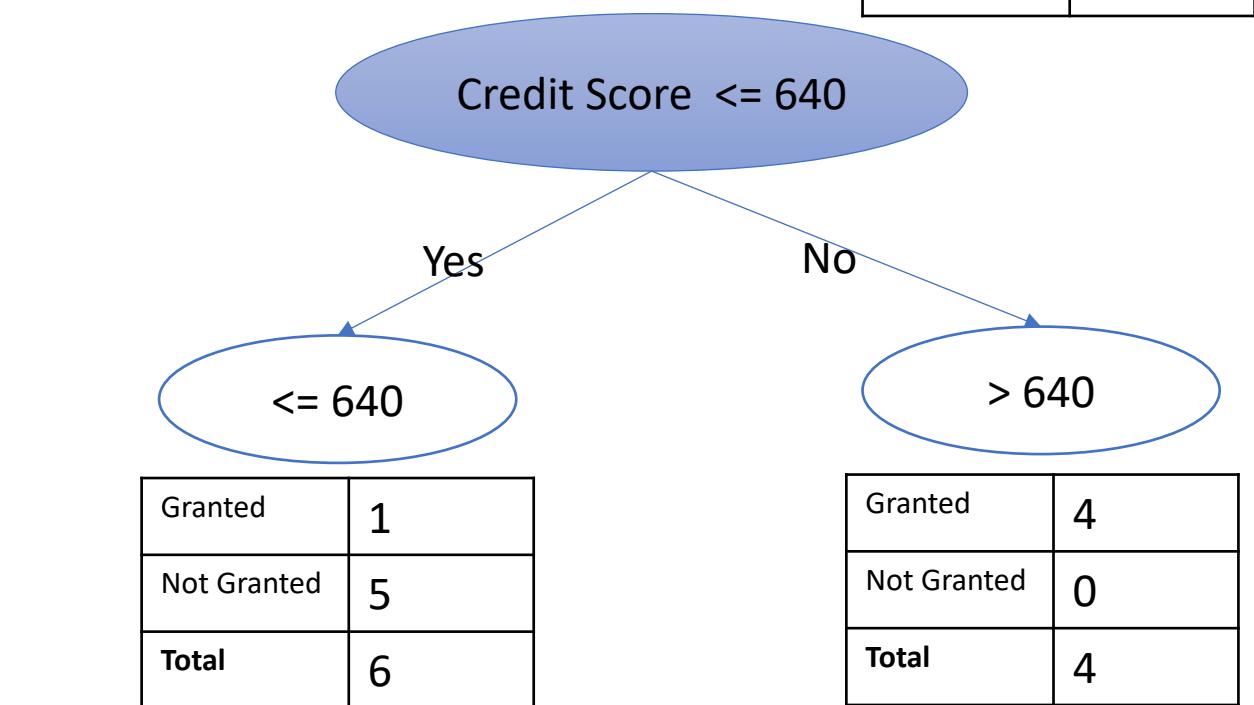
# How to deal with Numeric Feature ?

Credit Score	Income	Loan Status
150	25000	Not Granted
250	50000	Not Granted
360	60000	Not Granted
450	27000	Not Granted
540	40000	Not Granted
640	30000	Granted
730	45000	Granted
780	25000	Granted
800	55000	Granted
850	42000	Granted



# How to deal with Numeric Feature ?

Credit Score	Income	Loan Status	
150	25000	Not Granted	0.11
250	50000	Not Granted	
360	60000	Not Granted	0.4
450	27000	Not Granted	
540	40000	Not Granted	1
640	30000	Granted	0.61
730	45000	Granted	
780	25000	Granted	
800	55000	Granted	
850	42000	Granted	



0.6 \* Entropy @ left node : 0.65

0.7 \* Entropy @ right node : 0

Entropy before split	Entropy after split	Information Gain
1	0.39	0.61

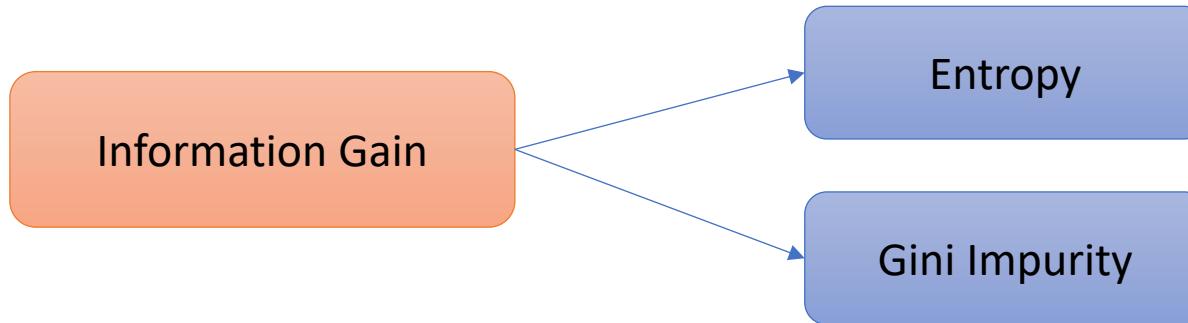
Credit Score	Education	Married	Loan Status
	Graduate	Yes	
	PG	No	

Features	IG
Credit Score	1
Education	0.86
Married	0.67



# Classification

For Classification problem , in decision tree we use Information Gain to build the model.



# Regression

For Regression problem , we can use Mean Square Error (MSE) or Mean Absolute Error (MAE).

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

MSE = mean squared error

$n$  = number of data points

$Y_i$  = observed values

$\hat{Y}_i$  = predicted values (mean value)

N=1000

Person A : Actual Loan Amt = 50000

Predicted Loan Amt = 40000

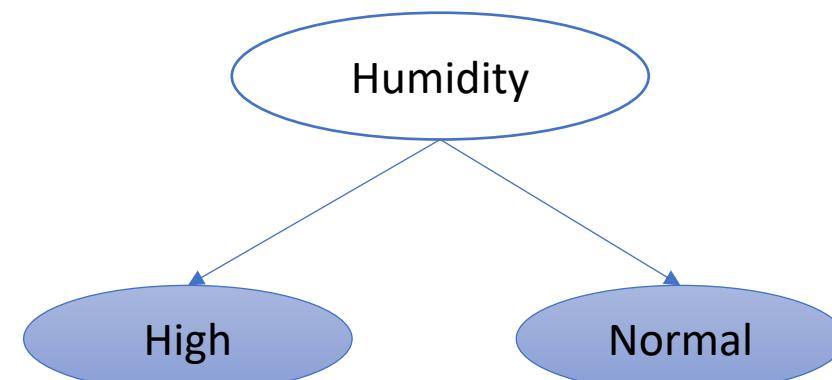
$$= (50000 - 40000)^2 + (60000 - 30000)^2$$

# Regression problem

Features				Target : Numeric
Outlook	Temperature	Humidity	Wind	Hours Played
Overcast	Mild	High	Strong	55
Rain	Mild	High	Strong	15
Sunny	Hot	High	Strong	15
Overcast	Hot	High	Weak	45
Rain	Mild	High	Weak	30
Sunny	Hot	High	Weak	10
Sunny	Mild	High	Weak	10
Overcast	Cool	Normal	Strong	50
Rain	Cool	Normal	Strong	10
Sunny	Mild	Normal	Strong	30
Overcast	Hot	Normal	Weak	40
Rain	Cool	Normal	Weak	30
Rain	Mild	Normal	Weak	20
Sunny	Cool	Normal	Weak	40

**How to select best Features:**

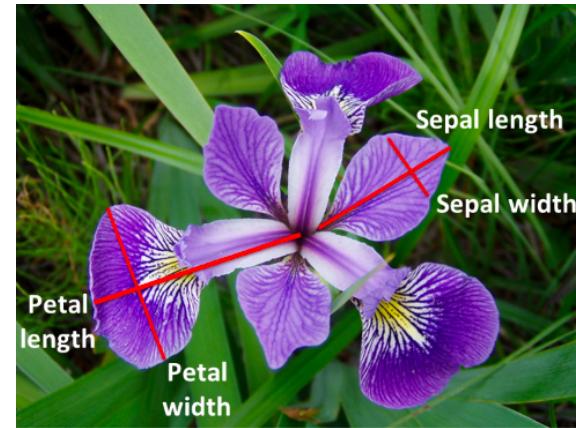
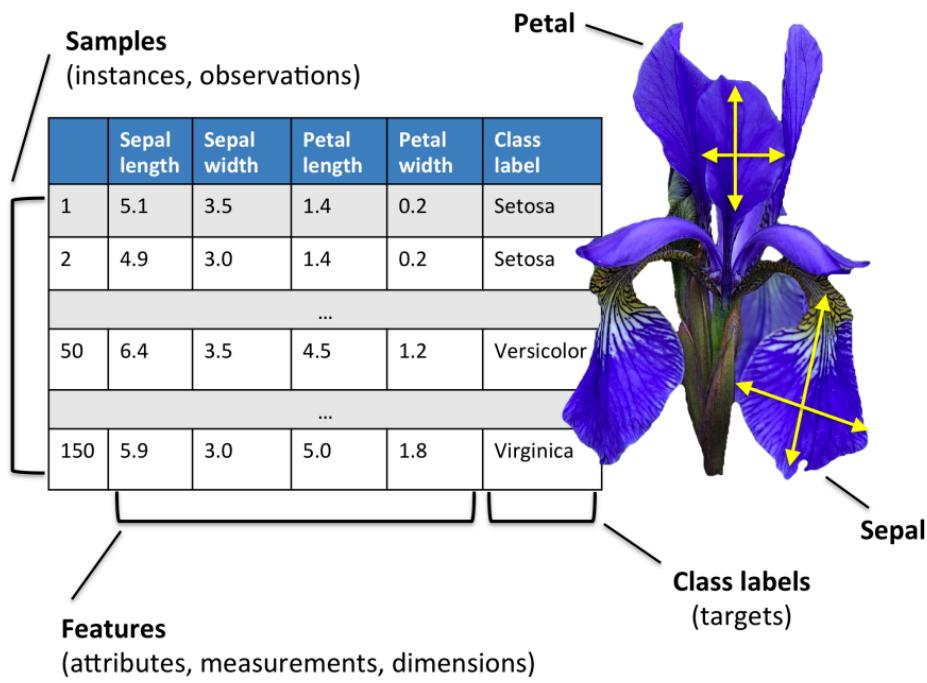
Error Reduction = MSE before split - MSE after split



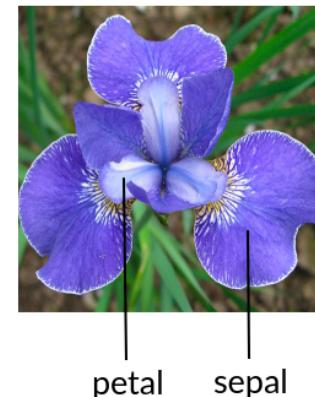
Features	Error Reduction
Outlook	25
Temp	12
Humidity	8.61
Wind	10

# DT Demo

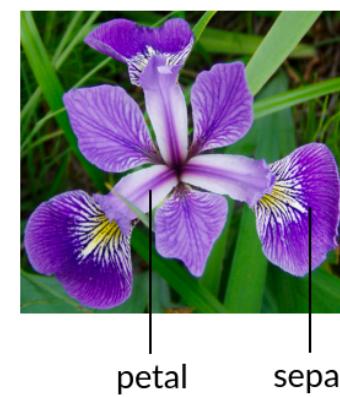
- Decision Tree on IRIS dataset



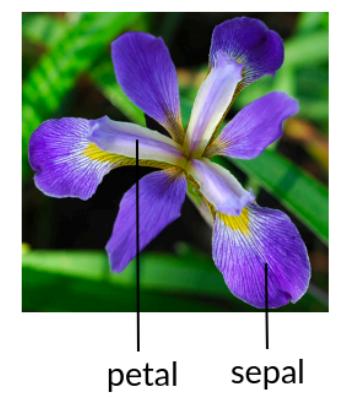
**iris setosa**



**iris versicolor**



**iris virginica**



# Summary of DT

- Interpretable:
  - Decision Tree models are highly interpretable as decision boundaries are axis parallel
- Feature Selection
- No Data Pre-processing is needed.
- It can handle both numeric & Categorical variables
- Easy to Overfit :
  - Same change in data impacts model badly
- Hyperparameter tuning is important to get good results

Confused about social distancing?  
Use this decision tree.



**DATA SCIENTIST**

**DEEP LEARNING**

**DECISION TREE**

Thank you

- Hope you have enjoyed the sessions

