

Lecture - 04 :

Machine learning algorithms are,

→ Decision Tree Classification

→ Decision Tree Regression

→ Practical Implementation

→ Ensemble Techniques

Decision Tree: [Solving Many Use Cases]

Decision Tree

↳ Regression

↳ Classification

Consider, [If - condition]

if (Age ≤ 18) :

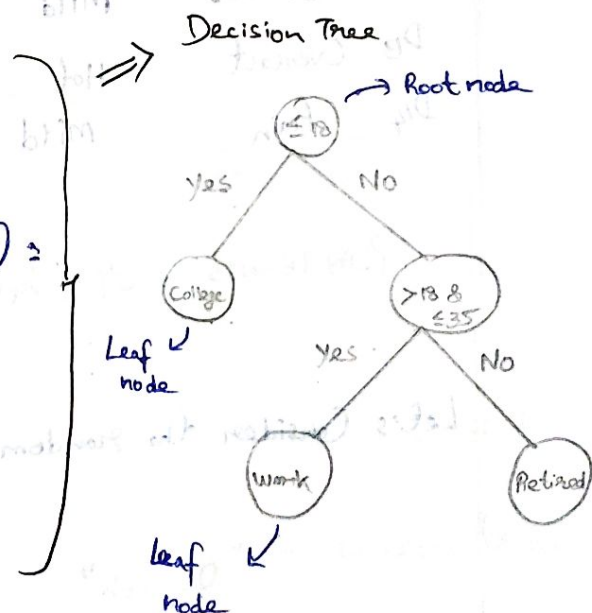
print ("College")

elif (Age > 18 and age ≤ 35) :

print ("Work")

else :

print ("Retired")



Same way, The regression and classification problem is Solved by Decision Tree.

Nested if else \Rightarrow Actually solved with the help of
Decision Tree (visualized manner).

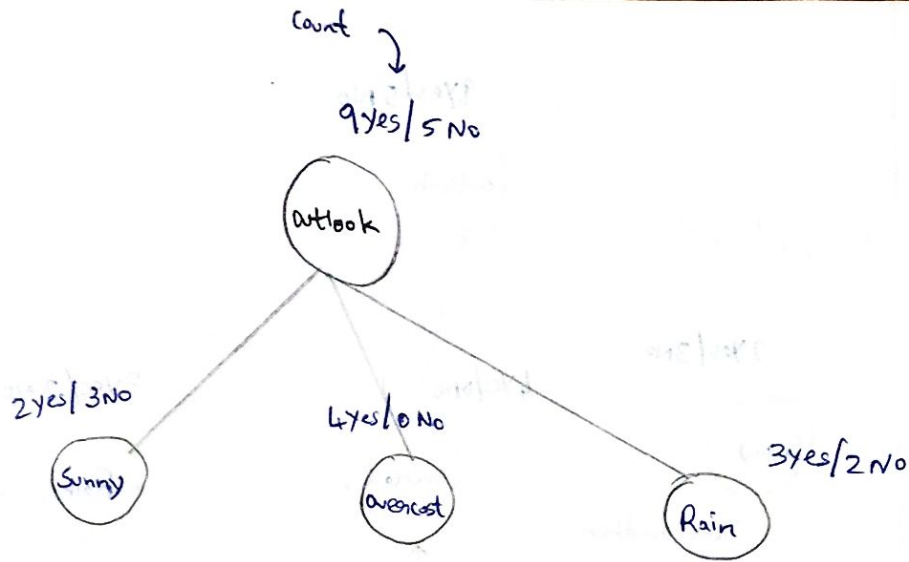
Let's consider the dataset,

DAY	OUTLOOK	TEMPERATURE	HUMIDITY	WIND	PLAY TENNIS
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

PLAY TENNIS \rightarrow O/P [dependent variable]

Let's consider the random feature,

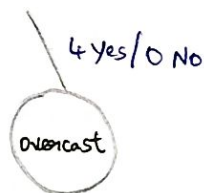
"OUTLOOK".



Pure split - Pure split ensures all the data in Same Category.

Impure Split - Impure Split ensures the data in a group belongs to ~~Impure~~ different Category.

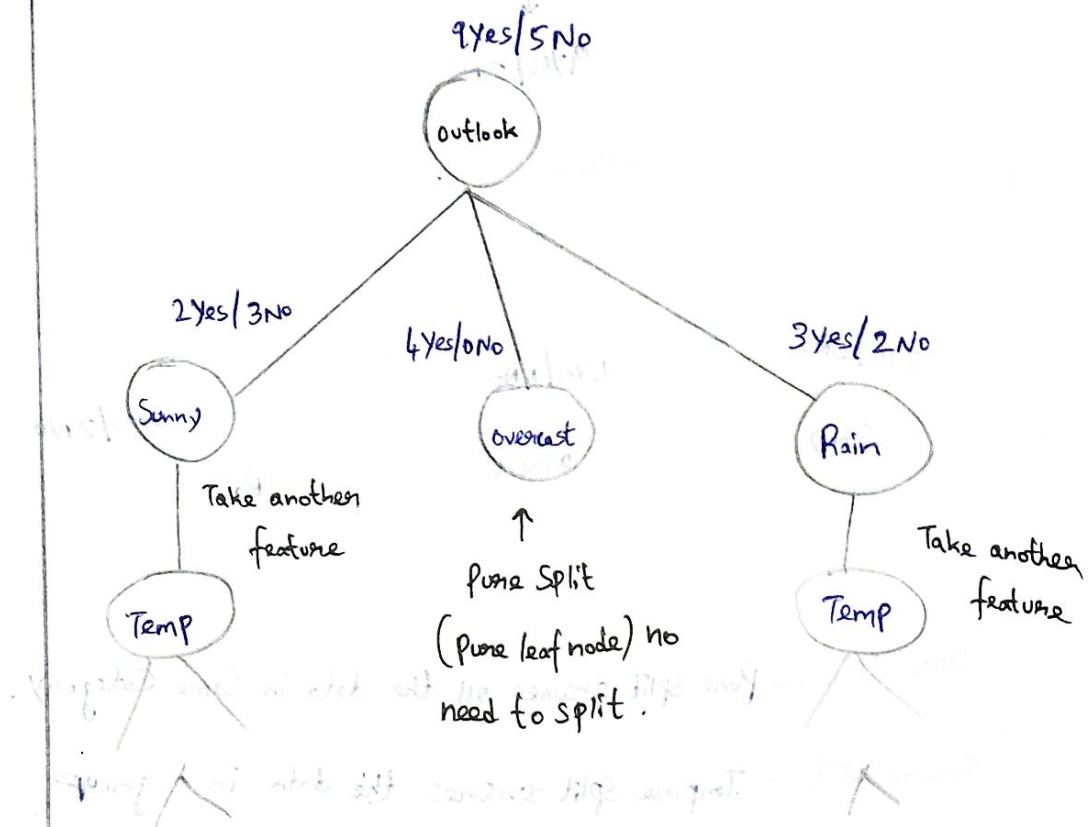
Let's take the Overcast node,



In this node, 4 yes / 0 No, ensures that all are in the Same Category "yes". Suppose the new test data will come, we definitely say Overcast with respect to "yes". You play the Tennis right now.

This type of Split is known as the pure split.

Suppose, The impure Split is Come, The decision tree take the next feature based on the information Gain.



Split will be happen until the Pure Split is Come. The feature Selection is done by information Gain.

Two Questions,

1. How Purity (or) pure Split is known by the Machine?

We are just Count and know the pure Split.
But in Machines have Complex Problems to solve.

Therefore, two techniques used to say the purity of node are,

→ Entropy.

→ Gini Coefficient (or) Gini impurity,

2. How the features are selected?

The feature selection with the help of the information Gain.

Entropy:

Formula,

$$H(s) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

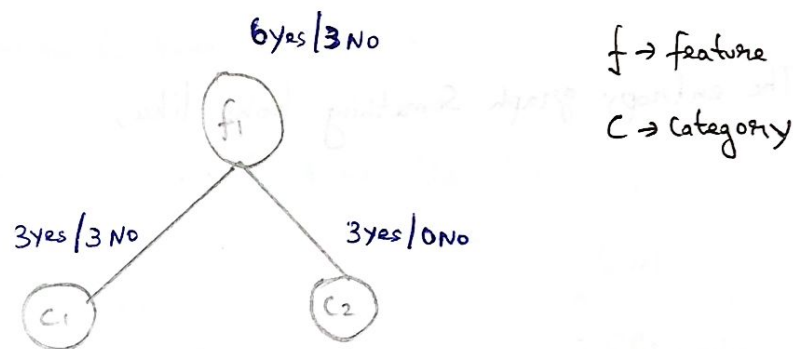
where,

$P_+ \Rightarrow$ probability of yes.

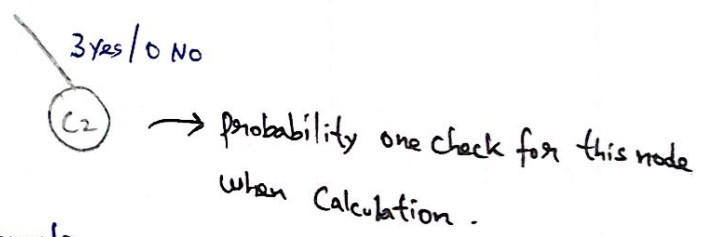
$P_- \Rightarrow$ probability of No.

+ } Binary
- } Classification.

Suppose Consider the node,



Take the node, C_2 to check the pure split or not



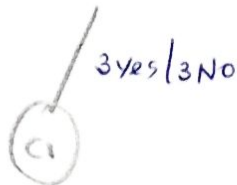
Apply in entropy formula,

$$H(s) = -\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{3} \log_2 \frac{0}{3}$$

$$= -1 \log_2 1$$

$$= -1(0) = 0 \rightarrow \text{Pure Split.}$$

Then check the another node, C_1 .

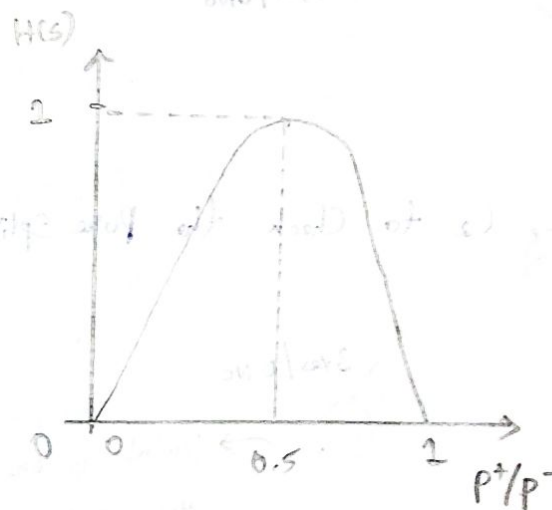


$$H(C_1) = -\frac{3}{6} \log_2 \frac{3}{6} - \frac{3}{6} \log_2 \frac{3}{6}$$

$$= -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$= 1 \rightarrow \text{Impure Split}$$

The entropy graph something looks like,



$$p^+ = 0.5$$

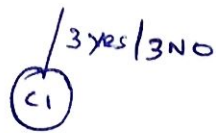
$$p^- = 0.5$$

$$p = 1 - q$$

$$q = 1 - p$$

Suppose, The Consider the Second node we are

take,



The entropy is 1.

The corresponding the P^+/P^- value is 0.5.

It means $0.5 \Rightarrow 50\%$ of P^+ and 50% of P^- .

The entropy values is always the range of 0 to 1.

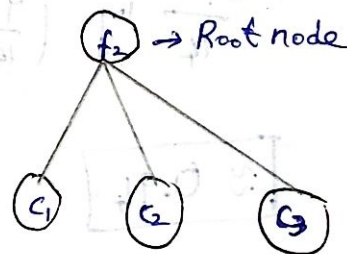
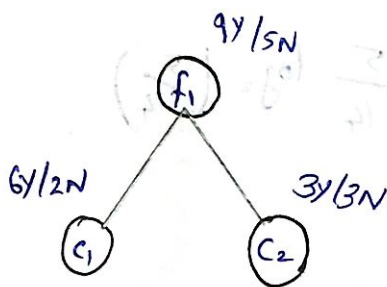
0 \rightarrow Pure Split.

1 \rightarrow Impure Split.

In purity Test, we use the entropy.

Which feature to take to split?

Suppose, we have to take the feature,

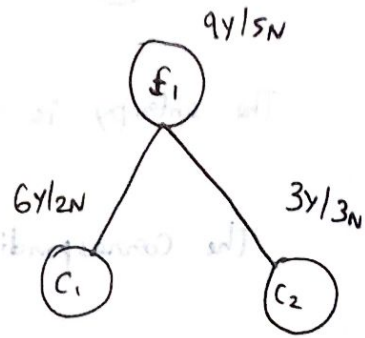


Which feature is to take for impure split?

It is responsible for information gain.

Information Gain:

$\text{Gain}(S, f)$ where $S \rightarrow$ Sample
 $f \rightarrow$ feature



$$\text{Gain}(S, f) = H(S) - \sum_{v \in \text{val}} \frac{|S_v|}{|S|} H(S'_v)$$

Entropy for
root node

$S_v \rightarrow$ How many samples
each
for split category

Entropy for ~~each~~ how
Many categories found.

$S \rightarrow$ Total Sample Size.

First find $H(S)$,

$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$= -\frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right)$$

$$\approx 0.94$$

Let's find the $H(S'_v)$,

$$H(C_1) = -\frac{6}{8} \log_2 \frac{6}{8} - \frac{2}{8} \log_2 \frac{2}{8}$$

$$H(c_1) = 0.81$$

Same like, Suppose we get

$$H(c_2) = 1$$

$9N/5N$

$$(f_1) \rightarrow \text{Total Sample} = 14 (s)$$

$6N/2N$

$$(c_1) \rightarrow \text{Total Sample} = 8 (s_{v1})$$

$3N/3N$

$$(c_2) \rightarrow \text{Total Sample} = 6 (s_{v2})$$

$$\text{Gain}(s, f_1) = 0.94 - \sum_{i=1}^2 \frac{|s_{vi}|}{|s|} \cdot H(s_{vi})$$

$$= 0.94 - \left[\frac{8}{14} \times 0.81 + \frac{6}{14} \times 1 \right]$$

$$= 0.049$$

$$\text{Gain}(s, f_1) = 0.049$$

Suppose, we get

$$\text{Gain}(s, f_2) = 0.051$$

$$\text{So, } \text{Gain}(s, f_2) >> \text{Gain}(s, f_1)$$

↳ This is more information.

So, take feature (f_2) to split the decision Tree when impure split is happened.

Not only for this two features, Calculate all of the feature and make the decision Tree Split.

GINI IMPURITY:

$n \rightarrow$ no of output

Output $\begin{cases} \text{Yes} \\ \text{No} \end{cases}$

$$G.I = 1 - \sum_{i=1}^n (P_i)^2$$

$$= 1 - [(P_+)^2 + (P_-)^2]$$



$$= 1 - \left[\left(\frac{1}{2} \right)^2 + \left(\frac{1}{2} \right)^2 \right]$$

$P_+ \rightarrow$ probability of yes

$$= 1 - \left[\frac{1}{2} \right]$$

$$P_+ = \frac{2}{4} = \frac{1}{2}$$

$$P_- = \frac{1}{2}$$

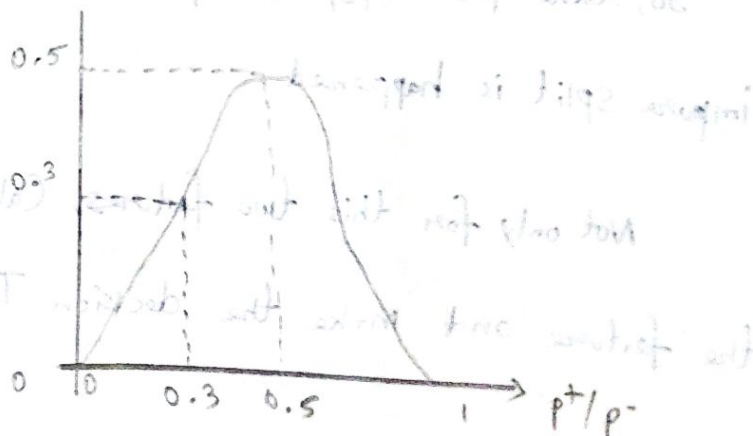
$$G.I = 0.5$$

Gini impurity is faster than entropy. But entropy is also important for calculating information gain.

Large set of datasets \Rightarrow Go Gini impurity.

Minimum or low numbers of dataset records

\hookrightarrow Go Entropy.



Suppose, In Classification problem we have numerical values in features - How decision tree work?

Suppose we have the Continuous feature,

f_1 O/P

2-3 -

1-3 -

4 -

5 -

7 -

3 -

The decision tree
sort the features

\Rightarrow

f_1 O/P

1-3 -

2-3 -

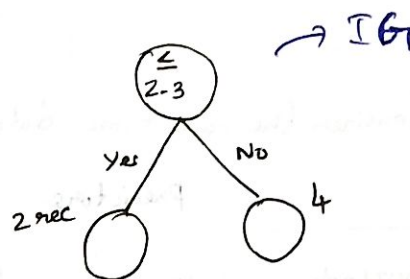
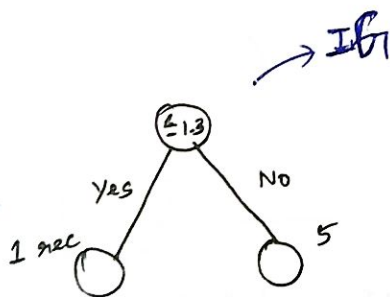
3 -

4 -

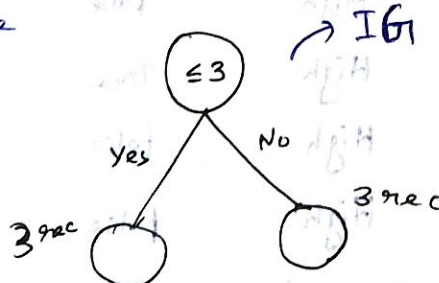
5 -

7 -

Then take first value,



leaf node



Highest IG Selected and
make the split.

IG \rightarrow Information Gain

Likewise, Computing the
all information gain
for every value. Then,

Compare and Select for
Splitting.

The main thing is that how decision tree work for Regression? . Because it is Continuous value . Then we move to,

Decision Tree Regressor:

Suppose,

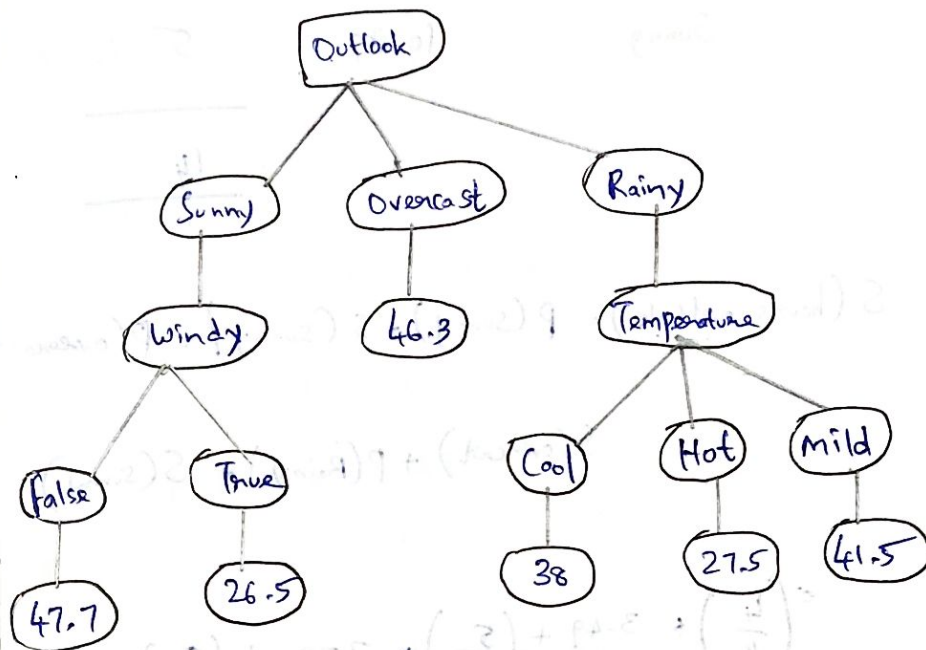
f_1	f_2	O/p	
2-3	5-6	-	} Continuous Values
5-2	7-8	-	
6-9	6-5	-	
4-4	8-4	-	
3-3	4-3	-	

Consider the real time dataset,

Predictors				target
OUTLOOK	TEMP	HUMIDITY	WINDY	HOURS PLAYED
Rainy	Hot	High	false	26
Rainy	Hot	High	True	30
Overcast	Hot	High	false	48
Sunny	mild	High	false	46
Sunny	Cool	Normal	false	62
Sunny	Cool	Normal	True	23
Overcast	Cool	Normal	True	43
Rainy	mild	High	false	36
Rainy	Cool	Normal	false	38
Sunny	mild	Normal	false	48

Rainy	mild	Normal	True	48
overcast	mild	High	True	62
overcast	Hot	Normal	False	44
Sunny	mild	High	True	30

Decision Tree
Converts to



How?

Firstly, calculate mean and sd for target variable,

$$\text{Mean} = 39.8$$

$$\text{sd} = 9.32$$

$$\text{Coefficient of Variation } (CV) = \frac{S}{\bar{x}} \times 100 = 23\%$$

↳ variation
in percentage.

Secondly, Take one feature and Calculate ~~mean~~ Sd along with forget variable,

		Hours played (Std)	Count
Outlook	Overcast	3.49	4
	Rainy	7.78	5
	Sunny	10.87	5
			<hr/>
			14
			<hr/>

$$S(\text{hours, outlook}) = P(\text{Sunny}) * S(\text{Sunny}) + P(\text{overcast}) * S(\text{overcast}) + P(\text{Rainy}) * S(\text{Rainy})$$

$$= \left(\frac{4}{14}\right) * 3.49 + \left(\frac{5}{14}\right) * 7.78 + \left(\frac{5}{14}\right) * 10.87$$

$$= 7.66$$

Standard deviation is high, So reduce it,

Something called as Standard deviation reduction.

		Hours played (Std)	
Outlook	Overcast	3.49	<div style="border: 1px solid black; padding: 5px; display: inline-block;">SDR = 1.66</div>
	Rainy	7.78	
	Sunny	10.87	

$$SDR(T, x) = S(T) - S(T, x)$$

$$= 9.32 - 7.66 = 1.66$$



$$SDR(\text{Hours, outlook}) = S(\text{Hours}) - S(\text{Hours, Outlook})$$

outlook \Rightarrow

$$\boxed{SDR = 1.66}$$

Same like,

$$\text{Humidity } SDR = 0.28$$

$$\text{Temp } SDR = 0.48$$

$$\text{Mild } SDR = 0.29$$

Then set one threshold using hypermeter,

and for many processes.

Refer: <https://www.Saadsayed.com/decision-tree-reg-htm>

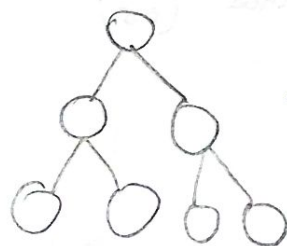
Another method (or)

Decision Tree Regression:

$$f_2(x, r) = 0 \text{ IP} = (x, r) \cdot 10^2$$
[illegible]

85.0 = 202 - Fibinul

In f_1 feature, we assign mean value and calculate MSE



Finally, which mse is low, that

Corresponding mean is Output.

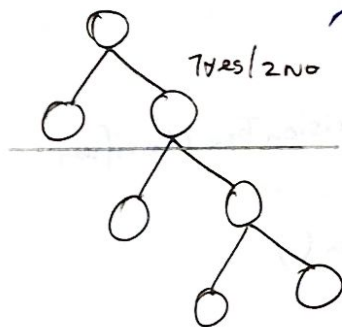
The decision tree have more depth (or) split, The model becomes complex. It will work for training data and fit. But unable to work testing data because we have only some amount of test data so it may be unfit and causes Overfitting. In order to avoid that we use, two techniques,

→ Post-pruning.

→ Pre-pruning.

Post pruning method set the after construction of decision tree.

Example,



→ In this case we know 7 Yes it means 80% of yes category. So, we decide it. So, cut the remaining split.

The advantages are time complexity is reduced and

Prevent Overfitting.

Pre-pruning method pruning the tree before

Constructing Using hyperparameter.

Hyperparameter:

Max depth, Max-leaf } Inside Grid Search CV.

Simple Tree fit Example Problem:

```
import pandas as pd
```

```
import matplotlib.pyplot as plt
```

```
% matplotlib inline
```

```
from sklearn.datasets import load_iris
```

```
iris = load_iris()
```

```
iris.data
```

```
from sklearn.tree import DecisionTreeClassifier
```

```
Classifier = DecisionTreeClassifier()
```

```
Classifier.fit(iris.data, iris.target)
```

```
from sklearn import tree
```

```
plt.figure(figsize=(15,10))
```

```
tree.plot_tree(Classifier, filled=True)
```

output: visualize tree
for get more
information.