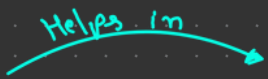



Decision Tree

Types

- ① Decision tree classifier  [Classification Problems]
- ② Decision tree Regressor  [Regression Problems]

Decision tree → It is a multination if else condition.

Decision tree classifier :-

- ID3 [Iterative Dichotomiser 3]
- CART [Classification and Regression tree]

NOTE:- CART in the Sklearn library Solve the classification and Reg. Problems in Decision tree.

➤ ID3 :-

When D.T. designed by ID3. We have a Root Node, we split the Node. split can be binary or multiple split. whenever we split the root node in such a way that it has multiple child Node, more than binary and these further going on. it keeps getting split like this then we can call this technique as ID3 Technique. child Node can be 2-3-4,.....,n in ID3 Technique.

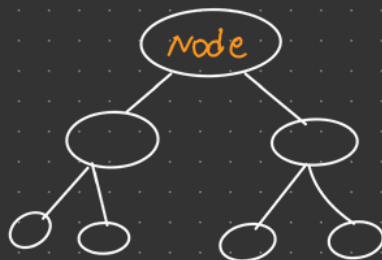
➤ CART :-

CART is called Classification and Regression Tree, the specific split done in it will be binary split.

ID3



CART



ex:-

Person, age = 14

if else condition

Making this nested if else condition into decision Tree

if (age \leq 15)

print ("School")

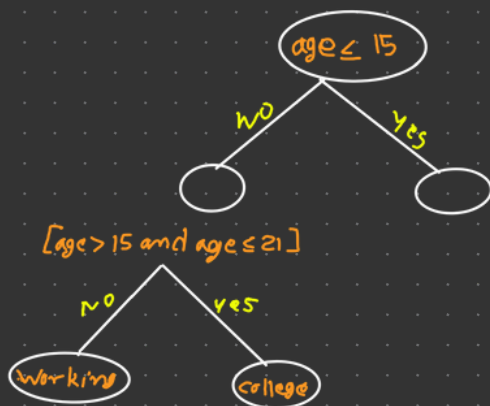
elif (age > 15 and age \leq 21)

print ("college")

else:

Print ("Working")

age = 14



ex, if age = 21 \rightarrow college
if age = 25 \rightarrow Working


How Decision tree works in the dataset *

Dataset → To Predict play Tennis or Not. (Goal)

Day	Outlook	Temp.	Humidity	wind	play Tennis
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


→ We select one feature and after that we create categories in the specific feature present in dataset.

Number of category = No. of Split 

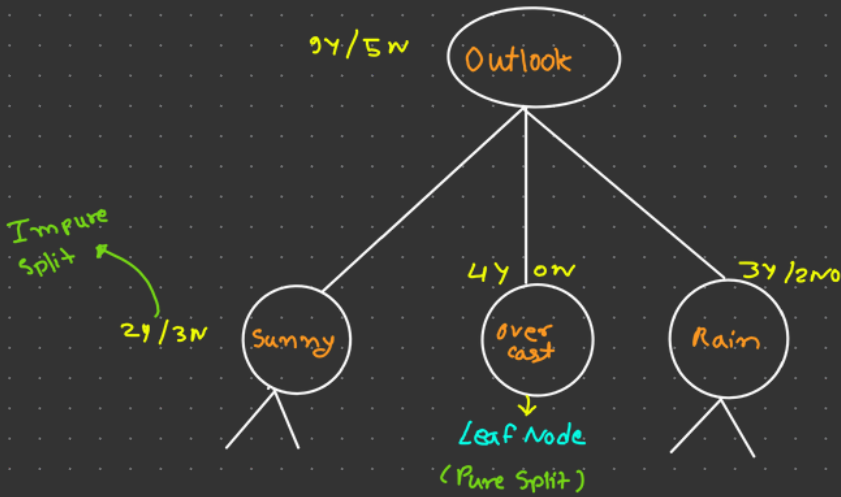
 Impure Split :- In the split of category when we get Yes OR No both it called Impure split. In case of Sunny we got 3 No and 2 Yes, So it is impure split.

Impure split give access to split again. when we got 50% Yes and 50% No, then it is also consider Impure Split.

Impure = 1 → in 50% Yes, 50% No

 Pure Split :- In the split of category when we got only Yes or only No. It called Pure split. In the case of overcast we got Yes and No. So it is a pure Split.

Pure split don't give access to split again, So it is called "Leaf Node".



Entropy Vs Gini Impurity

➤ Whenever dataset is small \Rightarrow Go with Entropy
(1000 - 2000 Records)

➤ Whenever huge is Huge \Rightarrow Go with Gini Impurity
(Millions, 100K Records)

How to Check Purity In ML *

Purity check

Pure Split OR Impure Split

we use two Techniques

➤ Entropy

➤ Gini Impurity

Used for measure of Purity

👉 What feature we need to select to start the split?

↳ For this we can use Information Gain Technique.

For Binary Classification

① Entropy :

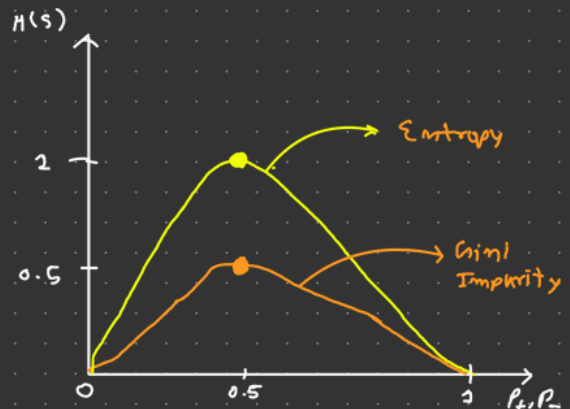
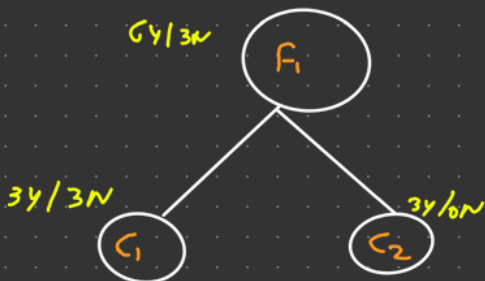
$$H(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

② Gini Impurity

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

Here $p_+ \rightarrow$ Probability of +ve category
 $p_- \rightarrow$ Probability of -ve category

ex



Entropy

$$① -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$H(c_1) = -\left(\frac{3}{6}\right) \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right)$$

$$H(c_1) = 1$$

Impure split

$$② -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$H(c_2) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

$$H(c_2) = \left(-\frac{3}{6}\right) \log_2 \left(-\frac{3}{6}\right) - \left(\frac{0}{6}\right) \log_2 \left(\frac{0}{6}\right)$$

$$H(c_2) = 0$$

Pure split

for Multiclass classification

c_1, c_2, c_3
yes / no / maybe

Entropy

$$H(S) = -P_{c_1} \log_2 P_{c_1} - P_{c_2} \log_2 P_{c_2} - P_{c_3} \log_2 P_{c_3}$$

Gini Impurity

$$G.I. = 1 - \sum_{j=1}^m P_j^2$$

$$G.I.(c_1) = 1 - [(P_+)^2 + (P_-)^2]$$

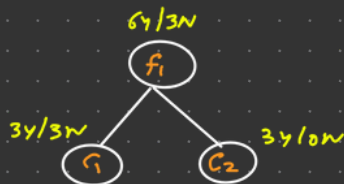
$$1 - \left[\left(\frac{3}{6}\right)^2 + \left(\frac{3}{6}\right)^2\right]$$

$$1 - \frac{1}{4} + \frac{1}{4}$$

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

$$G.I.(c_1) = 0.5$$

Impure split



$$G.I.(c_1) = 1 - [(P_+)^2 + (P_-)^2]$$

$$1 - \left[\left(\frac{3}{3}\right)^2 + \left(\frac{0}{3}\right)^2\right]$$

$$1 - 1$$

$$G.I.(c_2) = 0$$

Pure Split

Information Gain *

$$Gain(S, f_i) = H(S) - \sum_{v \in Val} \frac{|S_v|}{|S|} H(S_v)$$

ex,



⇒ Entropy of Root Node

$$\begin{aligned} H(S) &= -P_+ \log_2 P_+ - P_- \log_2 P_- \\ &= -\left(\frac{9}{16}\right)^2 \log_2 \left(\frac{9}{16}\right) - \left(\frac{5}{16}\right) \log_2 \left(\frac{5}{16}\right) \end{aligned}$$

$$H(S) \cong 0.94$$

$$\triangleright H(c_1) = -\left(\frac{6}{8}\right)^2 \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \frac{2}{8}$$

$$H(c_1) \cong 0.81$$

$$\triangleright H(c_2) = 50\% \text{ Yes and } 50\% \text{ NO}$$

$$H(c_2) = 1$$

→ Impure split

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

Here \rightarrow

$|S_v|$ = In the category total output

$|S|$ = Root Node total output value.

$H(S_v)$ = Entropy

$$\text{Gain}(S, f_1) = 0.94 - \left[\frac{8}{14} \times 0.91 + \frac{6}{14} \times 1 \right]$$

$$\text{Gain} = 0.049$$

f_1 Feature's total Gain = 0.049

Let's \rightarrow



\Rightarrow

f_2 Feature's total Gain = 0.051

To choosing what feature to select for split we take a Gain comparison of each feature. The feature with more information Gain, Selected for the splitting.

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

$>$

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

Information Gain is more when we split Using f_2 , So we go with f_2 feature for the Root Node and splitting.

Decision tree for Numerical Split.

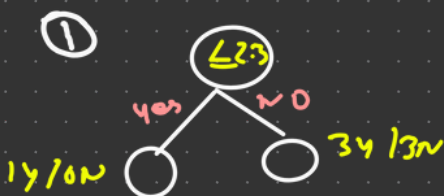
→ Sorting the feature: (in ascending order value set)

Suppose

<u>F₁</u>	<u>O/P</u>
2.3	yes
3.6	yes
4	no
5.2	no
6.7	yes
7.8	no
9.0	yes

→ Threshold Set

Sorting with Threshold = 2.3



② Threshold = 3.6



After that 2.3 Purity check with Entropy then check Information gain, After that next value with same step followed till the end.

③ Threshold = 4



④, ⑤, ⑥, ... ⑦

Same with all feature

At the end feature Selected whose information Gain is Highest.

 Disadvantage :-

when we have Huge Dataset (Millions) then \Rightarrow Time complexity $\uparrow\uparrow\uparrow$