

Decision Tree

[Advance & complex]

classmate

Date _____

Page _____

- It comes under supervised learning
- It is used to solve both Regression & classification problem.

① Decision Tree Regressor.

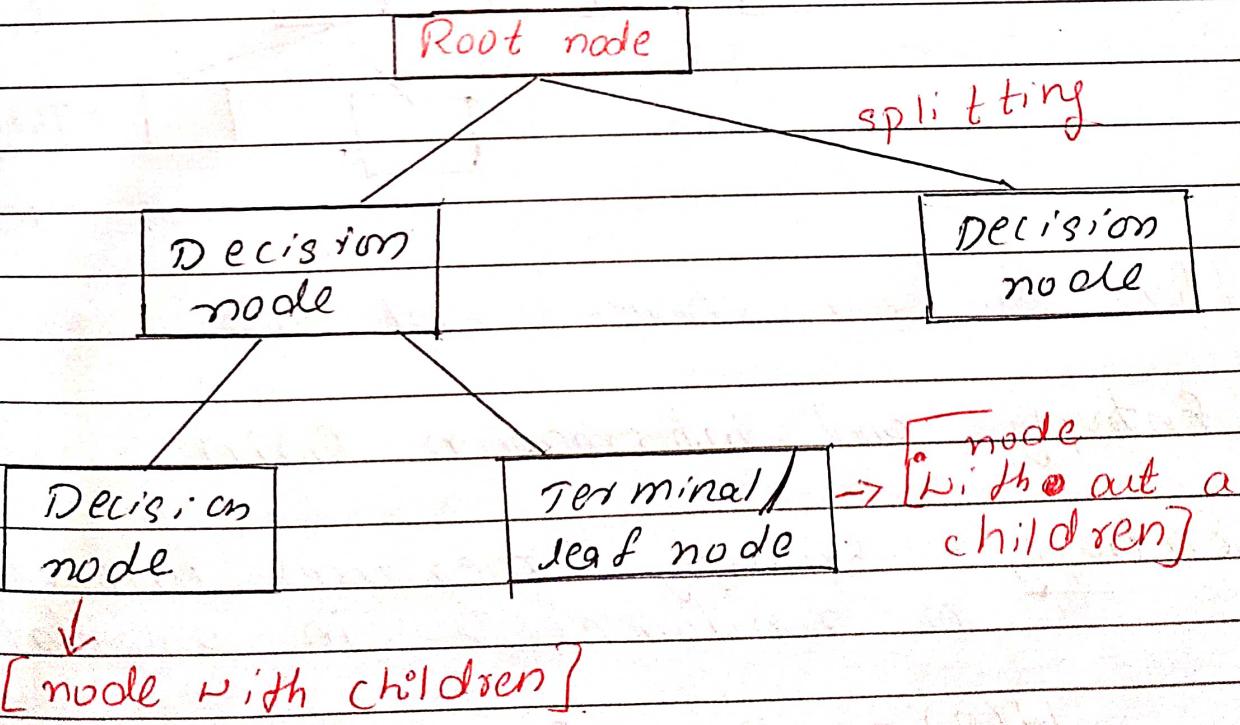
It is used to solve data which has continuous target.

② Decision Tree classifier.

It is used to solve data which has categorical target.

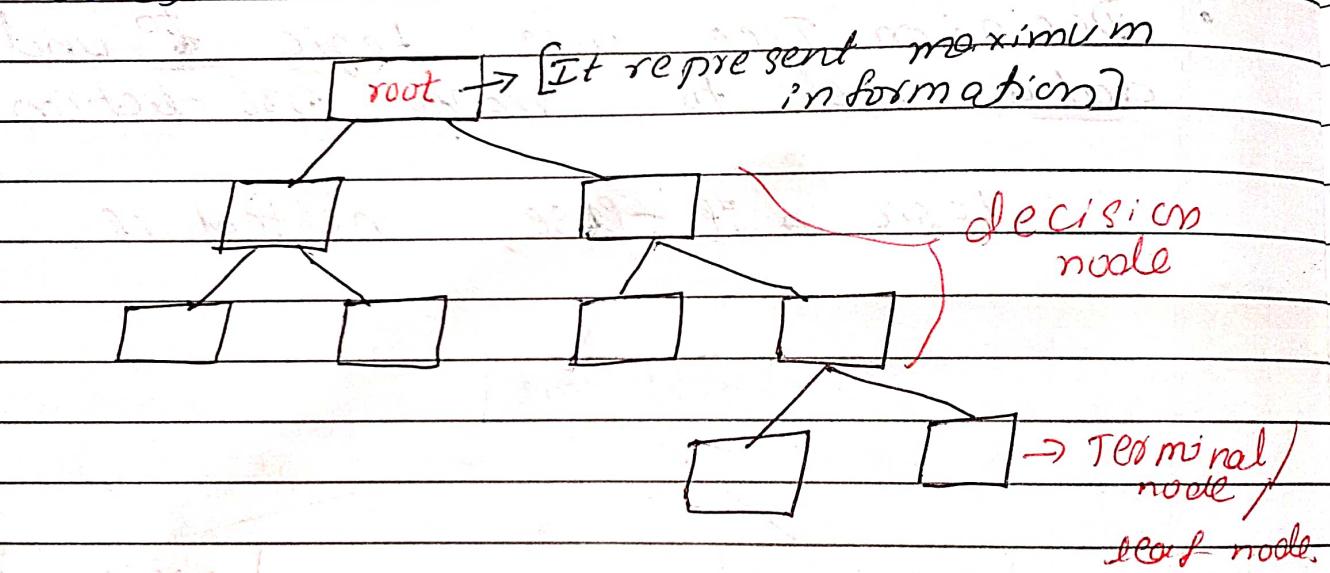
- Decision Tree uses Logic to understand data and to make prediction

Logic \Rightarrow if - else, nested if ...



How decision Tree work?

- Decision Tree builds Regression and classification models in the form of tree.
- It divides entire data into subsets such that an incremental tree will be developed.
- The final result will be the tree which includes decision node and leaf nodes



→ How do you choose root node?

- Entropy and information gain

① Entropy :- Entropy measures uncertainty or randomness in a data.

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

P_+ \rightarrow Probability of positive class
 P_- \rightarrow Probability of negative class.

<u>Purity</u>	<u>Impurity</u>	<u>Impurity</u>
x pure split (3) Yes (0) No	x 3 Yes 2 No	x 3 Yes 3 No
$\text{Entropy} = 0$ $\text{Entropy will be high.}$		$\text{Entropy} = 1$

- Range of Entropy = $[0, 1]$
- Entropy = 0, \rightarrow pure split \rightarrow x
 $(3 \text{ Yes } 0 \text{ No})$
- Entropy > 0 and < 1 , impurity.

Note:-

Entropy is calculated on categories of a column.

② Information Gain:-

- Information Gain tells about how much information a feature is giving.
- Gain is used to select the root node among input variables.
- Root node will be the feature which gives maximum information gain.

$$\text{Gain}(S) = H(S) - \frac{\sum S_v}{|S|} \cdot H(S_v)$$

$H(S)$ \rightarrow Entropy before split

$H(S_v)$ \rightarrow Entropy after split.

S \rightarrow no. of samples before split.

S_v \rightarrow no. of samples after split

\rightarrow Steps to calculate information gain

- ① Find entropy of data $H(S)$
- ② select column
- ③ Find entropy of categories in that column
- ④ calculate the information gain of that column
- ⑤ Repeat above 3 steps for different columns
- ⑥ Select the column which gives maximum information gain.

e.g.: Weather \rightarrow Sunny, rainy, cloudy

$\underbrace{H}_{\text{Entropy}}(\text{sunny}) \quad H(\text{rainy}) \quad H(\text{cloudy})$

Gain (weather)

Weather	Humidity	Wind	play
sunny	high	weak	no
sunny	high	strong	no
cloudy	high	weak	yes
rainy	high	weak	yes
rainy	normal	weak	yes
rainy	normal	strong	no
cloudy	normal	strong	yes
sunny	high	weak	no
sunny	normal	weak	yes
rainy	normal	weak	yes
sunny	normal	strong	yes
cloudy	high	strong	yes
cloudy	normal	weak	yes
rainy	high	strong	no

Eg 1: To select root for above table.

① find Entropy of data $H(S)$

→ Data
1
9 Yes 5 No

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

$$H(S) = -P(\text{Yes}) \cdot \log_2 P(\text{Yes}) - P(\text{No}) \cdot \log_2 P(\text{No})$$

$$P(\text{Yes}) = 9/14 \quad P(\text{No}) = 5/14$$

$$H(S) = (9/14) \cdot \log_2 (9/14) - (5/14) \cdot \log_2 (5/14)$$

$$H(S) = 0.94$$

② find entropy of weather. & gain of weather

→ sunny cloudy rainy
 | | |
 2 Yes 3 No 4 Yes 0 No 3 Yes 2 No

$$\bullet H(S) = -P(\text{Yes}) \cdot \log_2 P(\text{Yes}) - P(\text{No}) \cdot \log_2 P(\text{No})$$

→ calculate $H(\text{sunny})$

$$P(\text{Yes}) = 2/5$$

$$P(\text{No}) = 3/5$$

$$H(\text{sunny}) = -2/5 \log_2 (2/5) - 3/5 \log_2 (3/5)$$

$$H(\text{sunny}) = 0.9711$$

→ calculate entropy of cloudy ($H(\text{cloudy})$)

$$P(\text{yes}) = \frac{4}{4} \quad P(\text{no}) = 0$$

$$H(\text{cloudy}) = -\left(\frac{4}{4}\right) \log_2 \left(\frac{4}{4}\right) - \left(0\right) \log_2 \left(0\right)$$

$$H(\text{cloudy}) = 0.11$$

→ calculate $H(\text{Rainy})$

$$P(\text{yes}) = \frac{3}{5} \quad P(\text{no}) = \frac{2}{5}$$

$$H(\text{Rainy}) = -\left(\frac{3}{5}\right) \log_2 \left(\frac{3}{5}\right) - \left(\frac{2}{5}\right) \log_2 \left(\frac{2}{5}\right)$$

$$H(\text{Rainy}) = 0.97$$

$$\rightarrow H(S) = 0.94 \quad H(\text{sunny}) = 0.97$$

$$H(\text{cloudy}) = 0 \quad H(\text{Rainy}) = 0.97$$

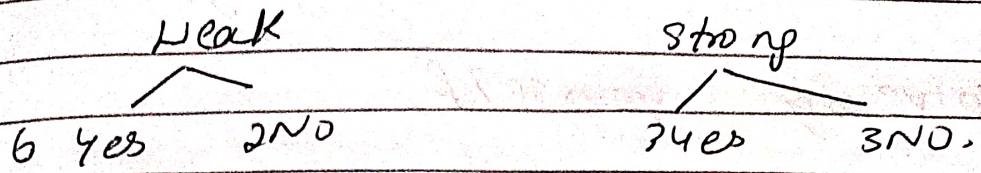
$$\rightarrow \text{Gain}(\text{Weather}) = H(S) - \frac{1}{|S_v|} \cdot H(S_v)$$

$$\text{Gain} = H(S) - \left(\frac{|S_v|}{S}\right) H(\text{sunny}) - \left(\frac{|S_v|}{S}\right) H(\text{cloudy}) \\ - \left(\frac{|S_v|}{S}\right) H(\text{rainy})$$

$$\text{Gain} = 0.94 - \left(\frac{1}{4}\right)(0.94) - \left(\frac{4}{4}\right)(0) - \left(\frac{1}{4}\right)(0.97)$$

$$\text{Gain} = 0.24$$

③ Find entropy of wind. and gain of wind



→ calculate $H(\text{weak})$:-

$$\bar{P}(\text{Yes}) = 6/8, P(\text{No}) = 2/8$$

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

$$H(\text{weak}) = 0.81$$

~~$$H(\text{strong}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$~~

→ calculate $H(\text{strong})$:-

$$P(\text{Yes}) = \frac{1}{2}, P(\text{No}) = \frac{1}{2} = \frac{3}{6} = \frac{1}{2}$$

$$H(\text{strong}) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2}$$

$$H(\text{strong}) = 1/1$$

$$\rightarrow \text{Gain} = H(S) - \left(\frac{S_V}{S}\right) H(\text{weak}) - \left(\frac{S_V}{S}\right) H(\text{strong})$$

$$\text{Gain} = 0.94 - \left(\frac{8}{14}\right)(0.81) - \left(\frac{6}{14}\right)(1)$$

$$\text{Gain} = 0.0481$$

(4) calculate Entropy of Humidity and Gain of Humidity.

→ solve by yourself.

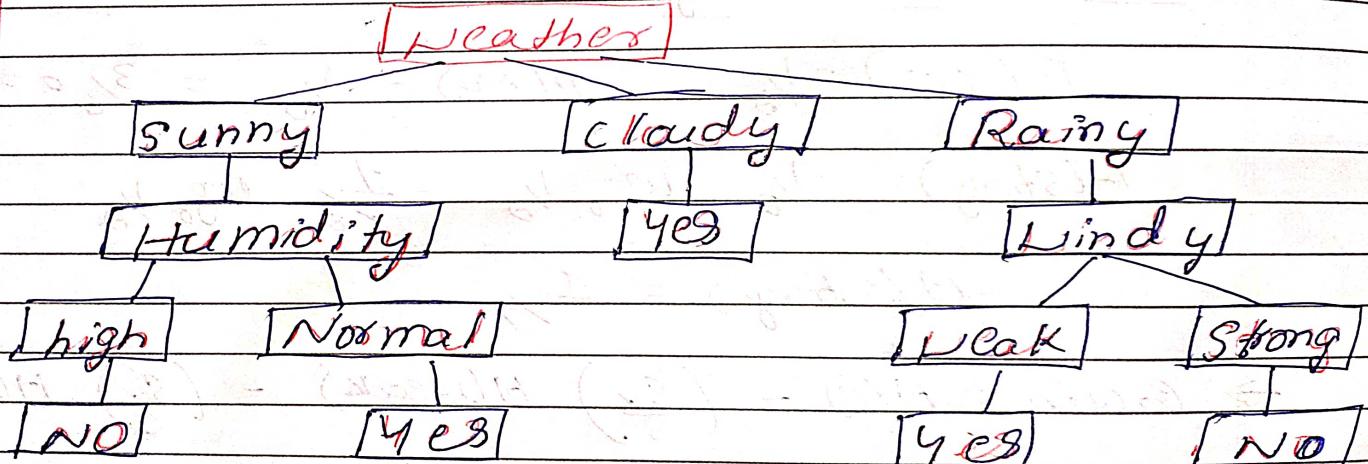
$$\rightarrow \text{Gain (Humidity)} = 0.12$$

$$(5) \text{Gain (Weather)} = 0.24$$

$$\text{Gain (Wind)} = 0.02$$

$$\text{Gain (Humidity)} = 0.12$$

Weather is selected as root node because it is giving maximum information gain.



Question:-

(1) If ~~cloudy~~, ~~high~~, ~~weak~~
Question?

Cloudy, high, weak → Yes or no

Sunny, high, weak → No

Rain, normal, weak → Yes.

~~IMP~~
Note: scaling and handling outliers are not required in tree based algorithms.

→ Limitation of Gain & Entropy:-

- Time taken to compute is very high because of log function.

→ How to overcome the limitation of Entropy and Information gain

- 1) Gini Index
- 2) Weighted Gini

1) Gini :-

- It measures impurity of the data
- It is used to select root node.
- Select the column which has less Weighted Gini.

$$\text{Gini} = 1 - \sum P_i^2 = 1 - [P(\text{Yes})]^2 - [P(\text{No})]^2$$

2) Weighted Gini :-

$$\text{Weighted Gini} = \frac{\sum S_v}{\sum S_i} \text{Gini}(S_v)$$

$S_i \Rightarrow$ no. of samples before split

$S_v \Rightarrow$ no. of samples after split

Note: Select the column which gives less weighted gini as root node.

① How to calculate weighted gini for above table

A:

→ ① for weather column

• Boxes

sunny	cloudy	rainy
2 yes 3 no	4 yes 0 no	3 yes 2 no
$P(\text{yes}) = 2/5$ $P(\text{no}) = 3/5$	$P(\text{yes}) = 4/4 = 1$ $P(\text{no}) = 0$	$P(\text{yes}) = 3/5$ $P(\text{no}) = 2/5$

$$\rightarrow \text{Gini}(\text{sunny}) = 1 - [P(\text{yes})]^2 - [P(\text{no})]^2$$

$$= 1 - (2/5)^2 - (3/5)^2$$

$$\text{Gini}(\text{sunny}) = \underline{\underline{0.48}}$$

$$\bullet \text{Gini}(\text{cloudy}) = 1 - (1)^2 - (0)^2 \\ = 0$$

$$\bullet \text{Gini}(\text{rainy}) = 1 - (3/5)^2 - (2/5)^2 \\ = \underline{\underline{0.48}}$$

$$\rightarrow \text{Weighted Gini} = \frac{1/5}{1/5} \text{Gini}(\text{sunny}) + \frac{1/5}{1/5} \text{Gini}(\text{cloudy}) \\ + \frac{1/5}{1/5} \text{Gini}(\text{rainy})$$

$$= \left(\frac{5}{14} \right) \cdot (0.48) + \left(\frac{5}{14} \right) (0.48)$$

$$= 0.34$$

Similarly calculate for Wind & humidity and compare all these result.

Weighted Gini (Weather) = 0.34

Weighted Gini (Wind) = 0.48

Weighted Gini (Humidity) = 0.52

→ From above we need to select which has less weighted gini.

→ Weather has less weighted gini so it is considered as root node.

Note: Gini range = [0, 0.5]

→ Converting Categorical column to numerical

① Label Encoder:

- When data is ordinal

- use `LabelEncoder` when data has 2 categories

Eg: Yes, No, Positive, Negative

② One Hot Encoder:

- use when data is nominal [no order]

shift + Tab → after Import algorithm name
from sklearn.ensemble import
gradientboostingclassifier
Date _____
Page _____
Shift + Tab

- Use this one hot encoding when you have ≥ 3 categories.

(3) manual Encoder:

- Use this when column has more than 3 categorical.
- Check the categories relationship with target and assign values according to that.
- It is used to encode target column.

Overfitting

- training score is high
- testing score is less
- Low bias, high variance

Underfitting

- Training score is less
- testing score is less
- high bias, high variance

→ Bias → Training Error / Error made on Train data;

→ Variance → Testing Error / Error made on Test data.

Note In general, if we find a model with low bias & low variance. It is the most suitable model to select.

Hyper parameter tuning in Decision tree :-

- 1) criterion :- It decides which method is used to measure the purity of the split.
→ criterion : [gini, entropy].
↓
It selects any one.
- 2) splitter :- It is the strategy or a technique used to select variables to split at node. (e.g.: - like root, sub node)
 $\text{splitter} = (\text{best}, \text{random})$
- 3) max depth :-
It tells about depth of the tree (levels) (height of a tree).
- 4) min sample split :-
 - minimum number of samples required to split internal node & (sub node)
 - An internal node will have further splits (also called children).
- 5) min sample leaf :-
 - minimum no of sample required to split at leaf node.
 - A leaf node is node without any children (without any further splits).

Note:- building model

- $n - \text{Jobs} = 1$ \Rightarrow It will use only one processor.
- $n - \text{Jobs} = -1$ \Rightarrow It will make use of all the processor.
- scoring = Select "accuracy" or "f1 score"
- verbose = To print message.
- CrossValidation [CV], \Rightarrow Cross validation. ~~the no~~
~~number of folds~~. Define how many CV required.