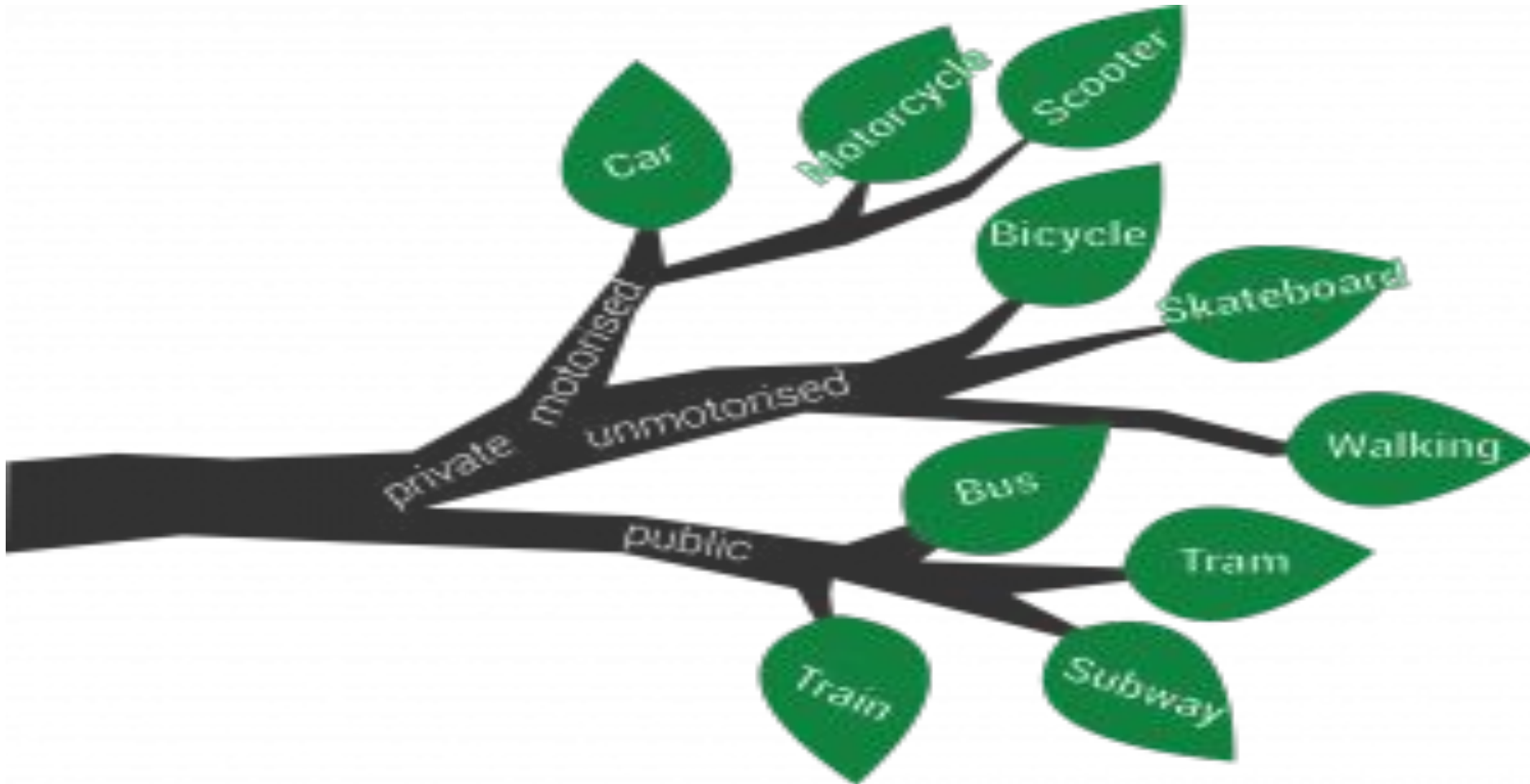
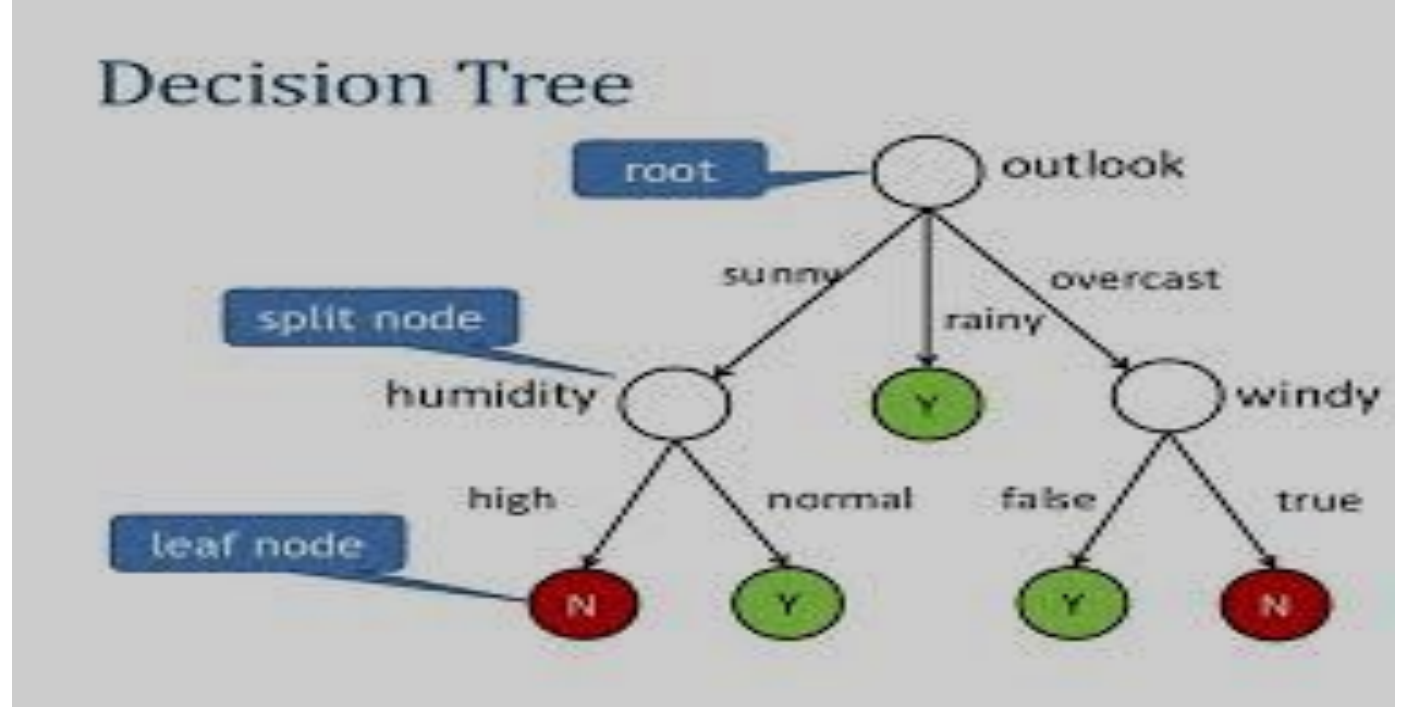


Decision Tree



Decision tree



- Decision tree builds classification or regression models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed.
- The final result is a tree with **decision nodes** and **leaf nodes**.
- A decision node (e.g., Outlook) has two or more branches (e.g., Sunny, Overcast and Rainy). Leaf node (e.g., Play) represents a classification or decision.
- The topmost decision node in a tree which corresponds to the best predictor called **root node**.

Decision tree: ID3 Algorithm

ID3 **Iterative Dichotomiser** is one of the most common decision tree algorithm.

Day	Outlook	Temp.	Humidity	Wind	Decision
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 can summarize the ID3 algorithm as illustrated below

$$\text{Entropy}(S) = \sum - p(I) \cdot \log_2 p(I)$$

$$\text{Gain}(S, A) = \text{Entropy}(S) - \sum [p(S|A) \cdot \text{Entropy}(S|A)]$$

ID3 Algorithm

$$E(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

Play Golf	
Yes	No
9	5



Entropy(PlayGolf) = Entropy (5,9)
= Entropy (0.36, 0.64)
= - (0.36 log₂ 0.36) - (0.64 log₂ 0.64)
= 0.94

Day	Outlook	Temp.	Humidity	Wind	Decision
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 need to calculate the entropy first. Decision column consists of 14 instances and includes two labels: yes and no. There are 9 decisions labeled yes, and 5 decisions labeled no.
- Entropy(Decision) = - p(Yes) . log₂p(Yes) - p(No) . log₂p(No)
- Entropy(Decision) = - (9/14) . log₂(9/14) - (5/14) . log₂(5/14) = 0.940
- Now, we need to find the most dominant factor for decision.

Wind factor on decision

- $\text{Gain}(\text{Decision}, \text{Wind}) = \text{Entropy}(\text{Decision}) - \sum [p(\text{Decision} | \text{Wind}) \cdot \text{Entropy}(\text{Decision} | \text{Wind})]$
- Wind attribute has two labels: weak and strong. We would reflect it to the formula.
- $\text{Gain}(\text{Decision}, \text{Wind}) = \text{Entropy}(\text{Decision}) - [p(\text{Decision} | \text{Wind}=\text{Weak}) \cdot \text{Entropy}(\text{Decision} | \text{Wind}=\text{Weak})] - [p(\text{Decision} | \text{Wind}=\text{Strong}) \cdot \text{Entropy}(\text{Decision} | \text{Wind}=\text{Strong})]$
- Now, we need to calculate $(\text{Decision} | \text{Wind}=\text{Weak})$ and $(\text{Decision} | \text{Wind}=\text{Strong})$ respectively.

Weak wind factor on decision

Day	Outlook	Temp.	Humidity	Wind	Decision
1	Sunny	Hot	High	Weak	No
3	Overcast	Hot	High	Weak	Yes
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes
13	Overcast	Hot	Normal	Weak	Yes

There are 8 instances for weak wind. Decision of 2 items are no and 6 items are yes as illustrated below.

1- $\text{Entropy}(\text{Decision} | \text{Wind}=\text{Weak}) = -p(\text{No}) \cdot \log_2 p(\text{No}) - p(\text{Yes}) \cdot \log_2 p(\text{Yes})$

2- $\text{Entropy}(\text{Decision} | \text{Wind}=\text{Weak}) = - (2/8) \cdot \log_2 (2/8) - (6/8) \cdot \log_2 (6/8) = 0.811$

Strong wind factor on decision

Day	Outlook	Temp.	Humidity	Wind	Decision
2	Sunny	Hot	High	Strong	No
6	Rain	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
14	Rain	Mild	High	Strong	No

Here, there are 6 instances for strong wind. Decision is divided into two equal parts.

$$1- \text{Entropy}(\text{Decision} | \text{Wind}=\text{Strong}) = - p(\text{No}) \cdot \log_2 p(\text{No}) - p(\text{Yes}) \cdot \log_2 p(\text{Yes})$$

$$2- \text{Entropy}(\text{Decision} | \text{Wind}=\text{Strong}) = - (3/6) \cdot \log_2 (3/6) - (3/6) \cdot \log_2 (3/6) = 1$$

Now, we can turn back to Gain(Decision, Wind) equation.

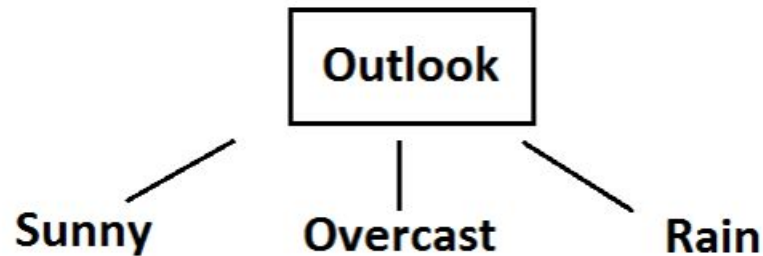
$$\begin{aligned} \text{Gain}(\text{Decision}, \text{Wind}) &= \text{Entropy}(\text{Decision}) - [p(\text{Decision} | \text{Wind}=\text{Weak}) \cdot \text{Entropy}(\text{Decision} | \text{Wind}=\text{Weak})] - p(\text{Decision} | \text{Wind}=\text{Strong}) \cdot \\ &\text{Entropy}(\text{Decision} | \text{Wind}=\text{Strong})] \\ &= 0.940 - [(8/14) \cdot 0.811] - [(6/14) \cdot 1] = 0.048 \end{aligned}$$

Other factors on decision

- Now, we can turn back to Gain(Decision, Wind) equation.
- $\text{Gain}(\text{Decision}, \text{Wind}) = \text{Entropy}(\text{Decision}) - [p(\text{Decision} | \text{Wind}=\text{Weak}) \cdot \text{Entropy}(\text{Decision} | \text{Wind}=\text{Weak})] - p(\text{Decision} | \text{Wind}=\text{Strong}) \cdot \text{Entropy}(\text{Decision} | \text{Wind}=\text{Strong})]$
- $= 0.940 - [(8/14) \cdot 0.811] - [(6/14) \cdot 1] = 0.048$
- Similar calculation on the other columns.
- 1- $\text{Gain}(\text{Decision}, \text{Outlook}) = 0.246$
- 2- $\text{Gain}(\text{Decision}, \text{Temperature}) = 0.029$
- 3- $\text{Gain}(\text{Decision}, \text{Humidity}) = 0.151$

Root node Selection

- As seen, outlook factor on decision produces the highest score. That's why, outlook decision will appear in the root node of the tree.



- Now, we need to test dataset for custom subsets of outlook attribute.

Overcast outlook on decision

Day	Outlook	Temp.	Humidity	Wind	Decision
3	Overcast	Hot	High	Weak	Yes
7	Overcast	Cool	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes

Decision will always be yes if outlook is overcast.

Sunny outlook on decision

Day	Outlook	Temp.	Humidity	Wind	Decision
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes

Here, there are 5 instances for sunny outlook. Decision would be probably 3/5 percent no, 2/5 percent yes.

1- $\text{Gain}(\text{Outlook}=\text{Sunny} \mid \text{Temperature}) = 0.570$

2- $\text{Gain}(\text{Outlook}=\text{Sunny} \mid \text{Humidity}) = 0.970$

3- $\text{Gain}(\text{Outlook}=\text{Sunny} \mid \text{Wind}) = 0.019$

Now, humidity is the decision

Decision

Day	Outlook	Temp.	Humidity	Wind	Decision
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
8	Sunny	Mild	High	Weak	No

Now, humidity is the decision because it produces the highest score if outlook were sunny.

At this point, decision will always be no if humidity were high.

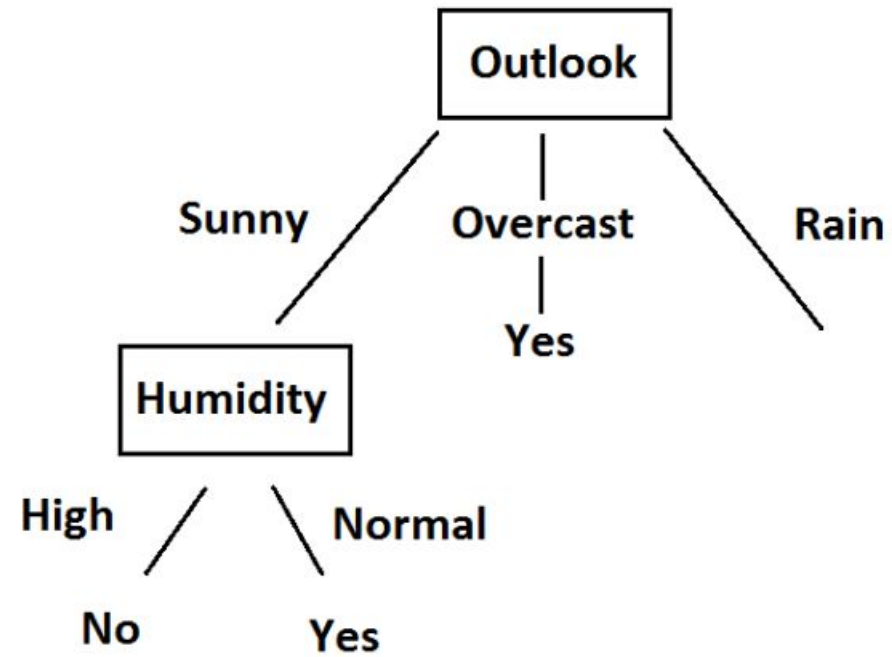
Day	Outlook	Temp.	Humidity	Wind	Decision
9	Sunny	Cool	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes

On the other hand, decision will always be yes if humidity were normal

Finally, it means that we need to check the humidity and decide if outlook were sunny.

Decision

Finally, it means that we need to check the humidity and decide if outlook were sunny.



Rain outlook on decision

Day	Outlook	Temp.	Humidity	Wind	Decision
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
6	Rain	Cool	Normal	Strong	No
10	Rain	Mild	Normal	Weak	Yes
14	Rain	Mild	High	Strong	No

1-Gain(Outlook=Rain | Temperature)

2- Gain(Outlook=Rain | Humidity)

3- Gain(Outlook=Rain | Wind)

Here, wind produces the highest score if outlook were rain. That's why, we need to check wind attribute in 2nd level if outlook were rain.

Decision

Day	Outlook	Temp.	Humidity	Wind	Decision
4	Rain	Mild	High	Weak	Yes
5	Rain	Cool	Normal	Weak	Yes
10	Rain	Mild	Normal	Weak	Yes

So, Decision will always be yes if wind were weak and outlook were rain.

Day	Outlook	Temp.	Humidity	Wind	Decision
6	Rain	Cool	Normal	Strong	No
14	Rain	Mild	High	Strong	No

And decision will be always no if wind were strong and outlook were rain.

Decision tree Construction

- So, decision tree construction is over. We can use the following rules for decisioning.

