

### Decision Tree (DT)

- DTs are a type of supervised machine learning algorithms where the data is continuously split according to a certain parameter/feature
- As it is a type of supervised machine learning algorithm so we can use it for classification as well as for regression problems.
- DT can be explained by two entities,  
① Decision nodes.  
② Leaves nodes.

- Decision nodes represent conditions on the basis of these conditions data is splitted
- Leaves nodes are the decisions on the final outcome.
- Edges represent possible values of node.
- In case of classification leaves represent classes and in case of regression leaves represent continuous values
- Most important attribute select as root node.

### → ID3 Algorithm

- ID stands for interactive dichotomizer (3<sup>rd</sup> version)
- The first step is to find the root node
- It uses a special function GAIN, to evaluate the gain information of each attribute
- For example if there are 3 attributes it will calculates the gain information for each
- whichever attribute has the maximum gain information becomes the root node.
- The rest of the attributes then fight for the next slots.
- ID3 based on information gain

### Entropy

- In order to define information gain precisely, we use entropy, which characterizes the purity and impurity of an arbitrary collection of examples.
- The range of entropy is 0 to 1.
- Minimum impurity collection of examples belong to only 1 class.
- Maximum impurity collection of examples have equal ratio of classes.

- Given a collection ( $S$ ), containing +ve and -ve examples of some target concept, the entropy of  $S$  relative to this Boolean classification is:

$$\text{Entropy of } S = E(S) = -P_+ \log_2 P_+ - P_- \log_2 P_-$$

where  $P_+$  is the proportion of positive class/examples in  $S$  and  $P_-$  is the proportion of negative examples in  $S$ .

- In all calculations involving entropy we define  $\log 0$  to be 0.

- For example  $S$  is a collection of 14 examples of some Boolean concept/problem, including 9 positive and 5 negative examples, then the entropy of  $S$  relative to this Boolean classification is:

$$E(S) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.940$$

So if the collection contains unequal numbers of +ve and -ve examples, the entropy is between 0 and 1.

→ Formula of Entropy for multiple classes

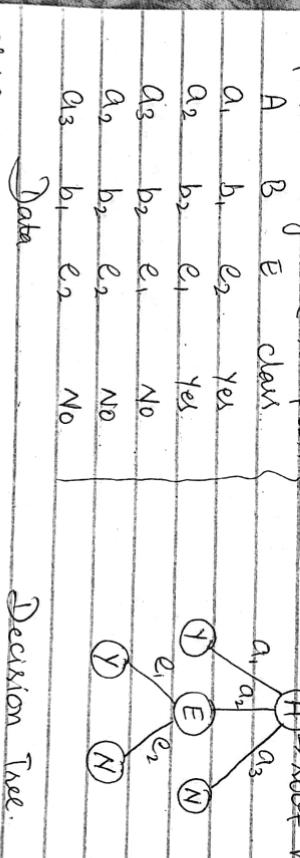
$$E(S) = -\sum_{i=1}^n P_i \log_2 P_i$$

where  $n$  represents the number of classes

→ Information Brain

- Given entropy as a measure of the impurity in a collection of training examples, we can now define a measure of the effectiveness of an attribute in

- The measure called information gain, is simply the expected reduction in entropy caused by partitioning the examples according to this attribute.
- That is, if we use the attribute with the maximum information gain as the node, then it will classify some of the instances as positive or negative with 100% accuracy, and this will reduce the entropy for the remaining instances.



Steph.

## Data Decision Tree

$$\begin{aligned}
 E(S) &= -\sum_{i=1}^n p_i \log_2 p_i \\
 &= -P_1 \log_2 P_1 - P_2 \log_2 P_2 \\
 &= -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \quad : P_1 = 2/5 \\
 &= 0.97
 \end{aligned}$$

(2) calculate the information gain for each attributes

- For attribute A

$$G_{(S,A)} = E(S) - \frac{|Sa_1|}{|S|} E(Sa_1) - \frac{|Sa_2|}{|S|} E(Sa_2) - \frac{|Sa_3|}{|S|} E(Sa_3)$$

As  $E(S) = 0.97$ ,  $|S| = 5$ ,  $|Sa_1| = 1$ ,  $|Sa_2| = 2$ ,  $|Sa_3| = 2$

$$E(Sa_1) = -P_Y \log_2 P_Y - P_N \log_2 P_N$$

$$= -\frac{1}{2} \log_2 \frac{1}{2} - 0 \log_2 0 = 0$$

Similarly

$$E(Sa_2) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

and

$$E(Sa_3) = -0 \log_2 0 - \frac{2}{2} \log_2 \frac{2}{2} = 0$$

So

$$G_{(S,A)} = 0.97 - \frac{1}{5} \times 0 - \frac{2}{5} \times 1 - \frac{2}{5} \times 0 = 0.51$$

- For attribute B

$$G_{(S,B)} = E(S) - \frac{|Sb_1|}{|S|} E(Sb_1) - \frac{|Sb_2|}{|S|} E(Sb_2)$$

$$= 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times E(Sb_2)$$

$$E(Sb_2) = -P_Y \log_2 P_Y - P_N \log_2 P_N$$
$$= -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.01$$

No

$$G_{(S,B)} = 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times 0.01 = 0.024$$

- For attribute C

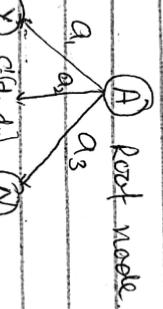
$$G_{(S,C)} = E(S) - \frac{|Sc_1|}{|S|} E(Sc_1) - \frac{|Sc_2|}{|S|} E(Sc_2)$$

$$= 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times E(Sc_2)$$

$$E(Sc_2) = -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = 0.91$$

$$= 0.97 - \frac{2}{5} \times 1 - \frac{3}{5} \times 0.91 = 0.024$$

(3) Select the attribute that has maximum Gain  
and define splitting criteria



④ Repeat all steps for sub-data set

- A B E class
- a<sub>2</sub> b<sub>2</sub> c<sub>1</sub> Yes
- a<sub>2</sub> b<sub>2</sub> c<sub>2</sub> No.

S'

$$E(S') = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{2} \log \frac{1}{2} = 1$$

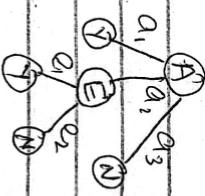
$$G_1(S', B) = E(S') - \frac{|S'|b_2|}{|S'|} E(S'|b_2|)$$

$$= 1 - \frac{2}{2} \times 1 = 0$$

$$G_2(S', E) = E(S') - \frac{|S'|e_1|E(S|c_1)}{|S'|} - \frac{|S'|e_2|E(S|c_2)}{|S'|}$$

$$= 1 - \frac{1}{2} \times 0 - \frac{1}{2} \times 0 = 1$$

So



Practice - 1.

Day(D)	Weather(W)	Temperature(T)	Humidity(H)	Wind(W)	Play(P)
1	Sunny	Hot	High	Weak	No
2	Cloudy	Hot	High	Weak	Yes
3	Sunny	Mild	Normal	Strong	Yes
4	Cloudy	Mild	High	Strong	Yes
5	Rainy	Mild	High	Strong	No
6	Rainy	Cool	Normal	Strong	No
7	Rainy	Mild	High	Weak	Yes
8	Sunny	Hot	High	Strong	No
9	Cloudy	Hot	Normal	Weak	Yes
10	Rainy	Mild	High	Strong	No

Dataset(S)

Solution.

- ① Find entropy of entire dataset(S).

$$E(S) = -P_1 \log_{\frac{1}{2}} P_1 - P_2 \log_{\frac{1}{2}} P_2 \\ = -\frac{5}{10} \log_{\frac{1}{2}} \frac{5}{10} - \frac{5}{10} \log_{\frac{1}{2}} \frac{5}{10} = 1.$$

- ② Select the attribute that has maximum gain.

For weather(W)

$$G(S, W) = E(S) - |S_{\text{Sunny}}|E(S_{\text{Sunny}}) - |S_{\text{Cloudy}}|E(S_{\text{Cloudy}})$$

$$|S| = |S_{\text{Sunny}}| + |S_{\text{Cloudy}}|$$

- $|S_{\text{Sunny}}|E(S_{\text{Sunny}})$ .

|S|

$$= 1 - \frac{3}{10} \times \left( -\frac{1}{3} \log_{\frac{1}{2}} \frac{1}{3} - \frac{2}{3} \log_{\frac{1}{2}} \frac{2}{3} \right) = \frac{3}{10} \times \left( \frac{3}{3} \log_{\frac{1}{2}} \frac{3}{3} - 0.0 \right)$$

$$= \frac{4}{10} \left( -\frac{1}{4} \log_{\frac{1}{2}} \frac{1}{4} - \frac{3}{4} \log_{\frac{1}{2}} \frac{3}{4} \right)$$

$$= 1 - 0.3 \times 0.91 - 0.3 \times 0 - 0.4 \times 0.81 = 0.40.$$

$G(S', w)$

For Temperature( $T$ )

$$G_1(S, T) = E(S) - \frac{|S|_{\text{Hot}}}{|S|} E(ST_{\text{hot}}) - \frac{|S|_{\text{Med}}}{|S|} E(ST_{\text{med}}) - \frac{|S|_{\text{Cold}}}{|S|} E(ST_{\text{cold}})$$

$$= 1 - \frac{4}{10} \left( -\frac{2}{4} \log_{2} \frac{2}{4} - \frac{2}{4} \log_{2} \frac{2}{4} \right) - \frac{5}{10} \left( -\frac{3}{5} \log_{2} \frac{3}{5} - \frac{2}{5} \log_{2} \frac{2}{5} \right)$$

$$= \frac{1}{10} (-0.4 \log_{2} 0 - \frac{1}{4} \log_{2} \frac{1}{4})$$

$$= 1 - 0.4 \times 1 - 0.5 \times 0.97 - 0.1 \times 0 = 0.11$$

For Humidity( $H$ )

$$G_1(S, H) = E(S) - \frac{|S|_{\text{High}}}{|S|} E(SH_{\text{high}}) - \frac{|S|_{\text{Normal}}}{|S|} E(SH_{\text{normal}})$$

$$= 1 - \frac{7}{10} \left( -\frac{3}{7} \log_{2} \frac{3}{7} - \frac{4}{7} \log_{2} \frac{4}{7} \right) - \frac{3}{10} \left( -\frac{2}{3} \log_{2} \frac{2}{3} - \frac{1}{3} \log_{2} \frac{1}{3} \right)$$

$$= 1 - 0.7 \times 0.98 - 0.3 \times 0.91 = 0.04$$

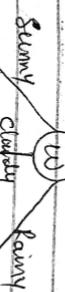
For Wind( $WD$ )

$$G(S, WD) = E(S) - \frac{|S|_{\text{Weak}}}{|S|} E(SWD_{\text{weak}}) - \frac{|S|_{\text{Strong}}}{|S|} E(SWD_{\text{strong}})$$

$$= 1 - \frac{4}{10} \left( -\frac{3}{4} \log_{2} \frac{3}{4} - \frac{1}{4} \log_{2} \frac{1}{4} \right) - \frac{6}{10} \left( -\frac{2}{6} \log_{2} \frac{2}{6} - \frac{4}{6} \log_{2} \frac{4}{6} \right)$$

$$= 1 - 0.4 \times 0.81 - 0.6 \times 0.91 = 0.13$$

③ Select Attribute that has maximum gain as root node.



$$S'_1 = [d_1, d_3, d_6]$$

④

$$S''_1 = [d_5, d_6, d_7, d_{10}]$$

④ Repeat 1 to 3 steps for  $S'$  and  $S''$

$$E(S') = -P_H \log_{2} P_H - P_M \log_{2} P_M$$

$$E(S'') = -\frac{1}{3} \log_{2} \frac{1}{3} - \frac{2}{3} \log_{2} \frac{2}{3} = 0.91$$

$$G_1(S', T) = E(S') - \frac{|S'_{\text{Total}}|}{|S'|} E(S'_{\text{Total}}) - \frac{|S'_{\text{Initial}}|}{|S'|} E(S'_{\text{Initial}})$$

$$= 0.91 - \frac{2}{3} (-0.1 \log_2 0 - \frac{2}{2} \log_2 \frac{2}{2}) - \frac{1}{3} (-\frac{1}{1} \log_2 \frac{1}{1} - 0 \log_2 0)$$

$$= 0.91 - 0.67 \times 0 - \frac{1}{3} \times 0 = 0.91$$

$$G_1(S', H) = E(S') - \frac{|S'_{\text{Hazard}}|}{|S'|} E(S'_{\text{Hazard}}) - \frac{|S'_{\text{Normal}}|}{|S'|} E(S'_{\text{Normal}})$$

$$= 0.91 - \frac{2}{3} (-0.1 \log_2 0 - \frac{2}{2} \log_2 \frac{2}{2}) - \frac{1}{3} (-\frac{1}{1} \log_2 \frac{1}{1} - 0 \log_2 0)$$

$$G_1(S', WD) = E(S') - \frac{|S'_{\text{WDweak}}|}{|S'|} E(S'_{\text{WDweak}}) - \frac{|S'_{\text{WDSound}}|}{|S'|} E(S'_{\text{WDSound}})$$

$$= 0.91 - \frac{1}{3} (-0.1 \log_2 0 - \frac{1}{1} \log_2 \frac{1}{1}) - \frac{2}{3} (-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2})$$

$$= 0.91 - 0.33 \times 0 - 0.67 \times 1 = 0.24.$$

For  $S''$

$$E(S'') = -P_A \log_2 P_A - P_B \log_2 P_B$$

$$= -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4}$$

$$= 0.81$$

$$G_1(S'', T) = E(S'') - \frac{|S''_{\text{Total}}|}{|S''|} E(S''_{\text{Total}}) - \frac{|S''_{\text{Initial}}|}{|S''|} E(S''_{\text{Initial}}).$$

$$= 0.81 - \frac{3}{4} \left( -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right) - \frac{1}{4} (-0.1 \log_2 0 - \frac{1}{1} \log_2 \frac{1}{1})$$

$$= 0.81 - 0.75 \times 0.91 - 0.25 \times 0 = 0.12$$

$$G_1(S'', H) = E(S'') - \frac{|S''_{\text{Hazard}}|}{|S''|} E(S''_{\text{Hazard}}) - \frac{|S''_{\text{Normal}}|}{|S''|} E(S''_{\text{Normal}})$$

$$= 0.81 - \frac{3}{4} \left( -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right) - \frac{1}{4} (-0.1 \log_2 0 - \frac{1}{1} \log_2 \frac{1}{1})$$

$$= 0.81 - 0.75 \times 0.91 - 0.25 \times 0 = 0.12.$$

G(S', WD)

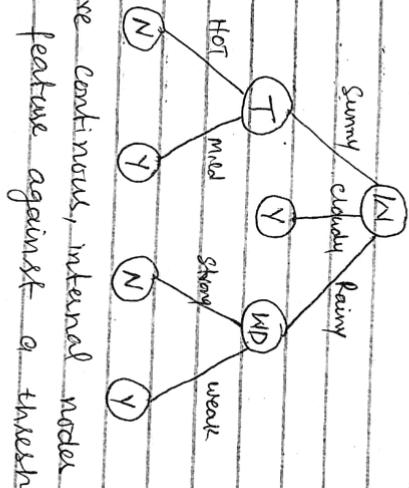
$$G(S', WD) = E(S') - |S'| \log_2 |S'|$$

$$= 0.81 - \frac{3}{4} \left( -0.69 \log_2 \frac{3}{2} - \frac{3}{3} \log_2 \frac{3}{2} \right) - \frac{1}{4} \left( -\frac{1}{1} \log_2 \frac{1}{1} - 0.69 \log_2 \frac{1}{2} \right)$$

$$= 0.81 - 0.75 \times 0 - 0.25 \times 0 = 0.81$$

As we see that.

$G(S', T)$  and  $G(S', H)$  have maximum equal gain so can select randomly.



Note:- If features are continuous, internal nodes can test the value of a feature against a threshold.