| Subject : Decision Tree | Date : | | |
| --- | --- | --- | --- |

We try to form the Decision tree as concise as possible. A concise Decision tree is likely to be efficient and given a test data it can reach the decision easily.

## Information Purity:

| Yess | No | | Yes | No |
| --- | --- | --- | --- | --- |
| (-) | (+ +) | | + | (+ / -) Min of info outcome |

Sunny                                              Rainy

Take Umbrella → +ve                    More Pure = More Information
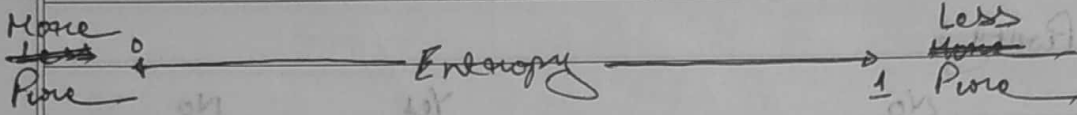Don't take Umbrella → -ve

Information wise, sunny is more pure than rain. Because under Yes and No there is no mixture of + and -. Under the class Yes and No they (+, -) are both getting conclusively separated.

Since, sunny is more pure so it has more information.
So, if our target is to make a tree as concise as possible then sunny should have more priority as the root node.

It is impossible to intuitively calculate information purity visually when the data set is very big like thousands of samples, with 30/40 attributes. In this case, the amount of decision with respect to attribute will also be greater.

The mathematical term to calculate information purity is called entropy.

The higher the value of entropy is, the attribute is less pure.

More
~~Less~~ : 0 ———————— Entropy ————————→ ~~Less~~ Less
Pure                                          1  Pure

$$\text{Entropy} = \sum_{i=1}^{n} -P \, Lg \, P \; ; \; i = \text{no. of labels}$$

Under a attribute, entropy is calculated for unique ~~vales~~ values.

| Sunny | Rainy | Decision |
|-------|-------|----------|
| Yes | No | Don't Take Umbrella |
| No | Yes | Take Umbrella |
| No | No | Take Umbrella |

Less Entropy = More Pure = More Information

$$\text{Entropy (Sunny = Yes)} = -P(\Delta T) Lg_2 P(\Delta T) - P(T) Lg_2 P(T)$$
$$= -\frac{1}{1} Log_2 \frac{1}{1} - \frac{0}{1} Lg_2 \frac{0}{1}$$
$$= 0$$
$$= \text{So this is very pure}$$

$$\text{Entropy (Sunny = No)} = -P(\Delta T) Lg_2 P(\Delta T) - P(T) Lg_2 P(T)$$
$$= -\frac{0}{2} Lg_2 P(0/2) - \frac{2}{2} Lg_2 \frac{2}{2}$$
$$= Lg_2 1$$
$$= 0 \quad \text{So the information is very pure}$$

Entropy $(Rainy = \overset{No}{Yes}) = -P(\Delta T) Lg_2 P(\Delta T) - P(T) Lg_2 P(T)$

$$= -\tfrac{1}{2} Lg_2 (\tfrac{1}{2}) - \tfrac{1}{2} Lg_2 (\tfrac{1}{2})$$

$$= -0.5 Lg_2 0.5 - 0.5 Lg_2 0.5$$

$= 1 \quad \leftarrow$ Max possible value. So the information
is very impure

Entropy $(Rain = Yes) = -P(\Delta T) Lg_2 P(\Delta T) - P(T) Lg_2 P(T)$

$$= -0/1 \, Lg_2 \, 0/1 - 1/1 \, Lg_2 \, 1/1$$

$\therefore$ Entropy (Sunny) $<$ Entropy (Rainy)

$\therefore$ Sunny is more pure.

Feature selection is based on purity. The more impure the
feature is, the less it is prioritized.

$\rightarrow$ So, for root node Sunny will be picked.

Based on the amount of information gain we select feature
for the root node. The decision tree creation algorithm based
on the information gain is called ID3. The more the information
gain, the more priority it gets to be the root node.

E.g.

$E(\text{Decision}) = -P(Y)Lg_2 P(Y) - P(N)Lg_2 P(N)$

$= -9/14 \, Lg_2 \, 9/14 - 5/14 \, Lg_2 \, 5/14$

$= 0.940$

$E(\text{Outlook} = \text{Sunny}) = -P(Y|S)Lg_2 P(Y|S) - P(N|S)Lg_2 P(N|S)$

$= -2/5 \, Lg_2 \, 2/5 - 3/5 \, Lg_2 \, 3/5$

$= 0.971$

$E(\text{Outlook} = \text{Rain}) = -P(Y|R)Lg_2 P(Y|R) - P(N|R)Lg_2 P(N|R)$

$= -3/5 \, Lg_2 \, 3/5 - 2/5 \, Lg_2 \, 2/5$

$= 0.971$

$E(\text{Outlook} = \text{Overcast}) = -P(Y|O)Lg_2 P(Y|O) - P(N|O)Lg_2 P(N|O)$

$= -4/4 \, Lg_2 \, 4/4 - 0/4 \, Lg_2 \, 0/4$

$= 0$

$IG(\text{Outlook}) = E(\text{Decision}) - E(S) - E(R) - E(O).$

$= 0.940 \quad 5/14$

$= E(\text{Decision}) - P(S)E(S) - P(R)E(R) - P(O)E(O)$

$= 0.940 - 5/14 (0.971) - 5/14 (0.971) - 4/14 (0)$

$= 0.246$

For Humidity

$$E(D) = -P(Y) \, Lg_2 \, P(Y) - P(N) \, Lg_2 \, P(N)$$

$$= -9/14 \, Lg_2 \, 9/14 - 5/14 \, Lg_2 \, 5/14$$

$$= 0.940$$

$$E(Humidity = High) = -P(H|H) \, Lg_2 \, P(H|H) - P($$

$$E(Humidity = High) = -P(Y|H) \, Lg_2 \, P(Y|H) - P(N|H) \, Lg_2 \, P(N|H)$$

$$= -3/7 \, Lg_2 \, 3/7 - 4/7 \, Lg_2 \, 4/7$$

$$= 0.985$$

$$E(Humidity = Normal) = -P(Y|N) \, Lg_2 \, P(Y|N) - P(N|N) \, Lg_2 \, P(N|N)$$

$$= -6/7 \, Lg_2 \, 6/7 - 1/7 \, Lg_2 \, 1/7$$

$$= 0.592$$

$$IG(H) = E(D) - P(H) * E(H) - P(N) * E(N)$$

$$= 0.940 - [7/14 * 0.985] - [7/14 * 0.592]$$

$$= 0.151$$

$$IG(W) = 0.048$$             More IG = More Priority.

$$IG(T) = 0.029$$

The less the entropy, the greater the information gain.

Here, outlook has the most information gain. So, it will be the root node.

The rows of overcast have only the decision value of YES. So, there is no conflict of decision. We are. getting pure information.

② For Temperature under sunny.

$$E(T) = -\frac{2}{5} Lg_2 \frac{2}{5} - \frac{3}{5} Lg_2 \frac{3}{5}$$

$$= 0.970$$

$$E(T=H) = -P(Y|H) Lg_2 P(Y|H) - P(N|H) Lg_2 P(N|H)$$

$$= -\frac{2}{2} Lg_2 \frac{2}{2} - \frac{0}{2} Lg_2 \frac{0}{2}$$

$$= 0$$

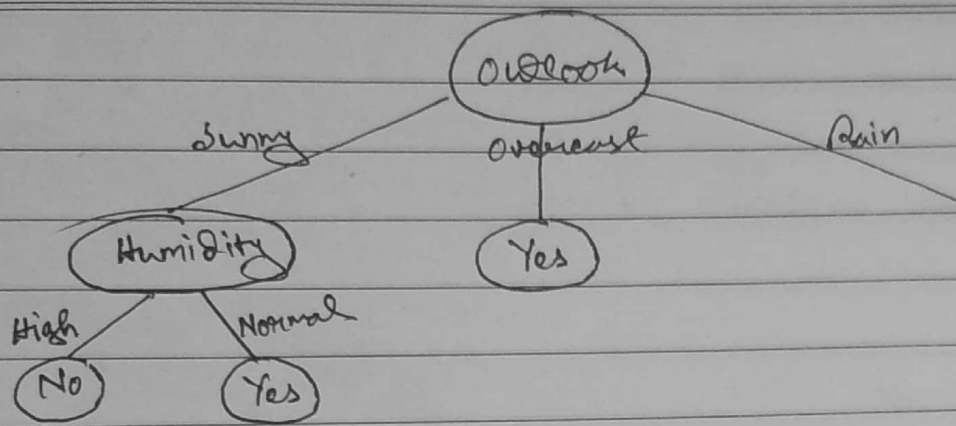$$E(M) = -\frac{1}{2} Lg_2 \frac{1}{2} - \frac{1}{2} Lg_2 \frac{1}{2} = 1$$

$$E(C) = 0$$

$$IG(T|0) = 0.970 - \frac{2}{5}(0) - \frac{2}{5}(1) - \frac{1}{5}(0)$$

$$= 0.570$$

$$IG(H|0) = 0.970$$

$$IG(W|0) = 0.019$$

∴ Under sunny, Humidity will sit.

Decision tree: Outlook → Sunny → Humidity (High → No, Normal → Yes); Overcast → Yes; Rain

20+  } Highest Info needed
20-  }  High uncertainly

20+  } no info needed
0-   } no uncertainty

~~80~~ 14+  } Needs less info
2-   } Less uncertainly.

The higher the information gain, the higher the separating ability.

The higher the info gain, the higher the no. of branches, i.e info. gain favours attributes having large no. of branches.

$$Gain\ Ratio = \frac{Info.\ Gain}{Split\ Information}$$