

Decision Tree

In Decision tree we measure Entropy and Information Gain of a feature to decide which feature is better to split the data.

Entropy is similar to Gini Impurity. So, Lower the entropy, better the feature. Information Gain is opposite of these two, which helps the decision tree algorithm to identify the most informative feature to make splits that lead to clearer and more predictable classification. So, the more the Information Gain, the better the feature.

Formula —

$$\text{Entropy } H(X) = \sum_{i=1}^c -P_i \log_2 P_i$$

where c = number of classes.

$$\text{Information Gain } (S, A) = H(S) - \sum P(S|A) \times H(S|A)$$

Here, S = Dataset and A is a feature of the dataset.

Dataset →

	Loves popcorn	Loves soda	Age	Loves Troll 2
①	Yes	Yes	7	NO
②	Yes	No	12	NO
③	NO	Yes	18	Yes
④	NO	Yes	35	Yes
⑤	Yes	Yes	38	Yes
⑥	Yes	No	50	NO
⑦	NO	No	83	No

features target column

Now we have to find Information Gain for Loves popcorn, Loves soda and Age Features to decide which should be the root node of the tree.

First let's find the Entropy of the dataset.

Entropy of Dataset -

→ [If there was another class we would add that as well]

$$\text{info}(D) / H(D) / H(S) = -P(\text{Yes}) \log_2 P(\text{Yes}) - P(\text{NO}) \log_2 P(\text{NO})$$

$$= -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7}$$

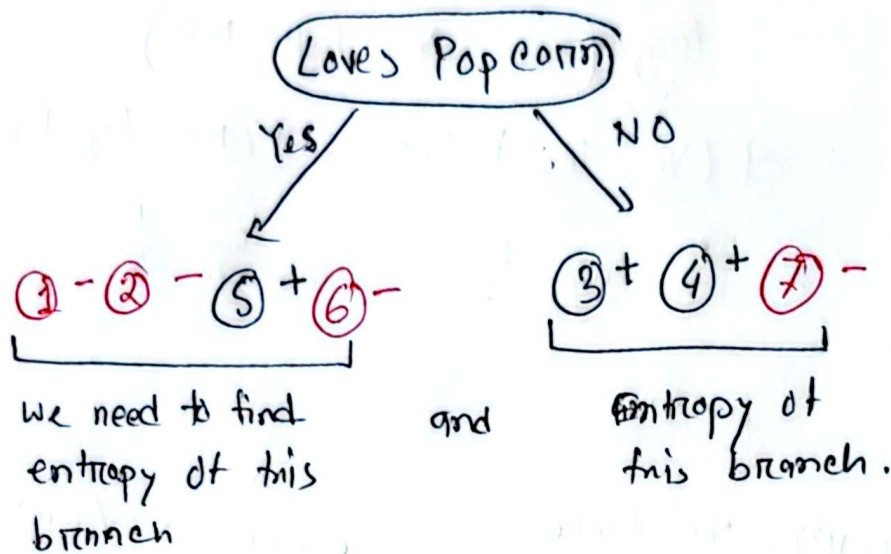
$$= 0.985$$

there are a total 7 data points
in the data set, among those
there are 3 'Yes' in the
target column

Among 7 data points,
there are 4 'No'
in the target column

'Loves popcorn feature' →

Let us find Entropies first and then we will find I_h of Loves popcorn.



Entropy of ~~left~~ 'Yes' branch →

$$\begin{aligned}
 H(\text{Loves popcorn} = \text{Yes}) &= -P(\text{Yes}) \log_2 P(\text{Yes}) - \boxed{P(\text{No})} \log_2 \boxed{P(\text{No})} \\
 &= -P(\text{Loves popcorn} = \text{Yes} | \text{Loves popcorn} = \text{Yes}) \log_2 (\text{Loves Troll 2} \\
 &= \text{Yes} | \text{Loves popcorn} = \text{Yes}) - \boxed{P(\text{Loves Troll 2} = \text{No} | \text{Loves} \\
 &\quad \text{popcorn} = \text{Yes})} \log_2 \boxed{P(\text{Loves Troll 2} = \text{No} | \text{Loves popcorn} = \text{Yes})}
 \end{aligned}$$

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

$$\approx 0.811$$

Among 4 Loves popcorn = Yes,
there is 1 Loves Troll 2 = No.

Among 4 Loves popcorn =
Yes, there are 3
Loves Troll 2 = No

Entropy of No branch \rightarrow

$$\begin{aligned} H(\text{Loves popcorn} = \text{NO}) &= -P(\text{Troll} = \text{Yes} | \text{pop} = \text{NO}) \log_2 P(\text{Troll} = \text{Yes} | \\ &\quad \text{pop} = \text{NO}) - P(\text{Troll} = \text{NO} | \text{pop} = \text{NO}) \\ &\quad \log_2 P(\text{Troll} = \text{NO} | \text{pop} = \text{NO}) \\ &= -P(\text{Yes}) \log_2 P(\text{Yes}) - P(\text{NO}) \log_2 P(\text{NO}) \\ &= -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \\ &= 0.918 \end{aligned}$$

Among Loves popcorn = NO, there
are 2 Loves Troll = Yes

Among 3 Loves popcorn
= NO, there are 1
Loves Troll = NO

\therefore Weighted Entropy of Loves popcorn \rightarrow

$$\text{Info}(\text{Loves popcorn}) / H(\text{Loves popcorn})$$

$$= P(\text{Loves popcorn} = \text{Yes}) \times H(\text{Loves popcorn} = \text{Yes})$$

$$+ P(\text{Loves popcorn} = \text{NO}) \times H(\text{Loves popcorn} = \text{NO})$$

$$= \frac{4}{7} \times 0.811 + \frac{3}{7} \times 0.918$$

$$= 0.857$$

\hookrightarrow [If there was another
outcome of Loves popcorn
we would add that
as well]

There are 4
Loves popcorn = Yes

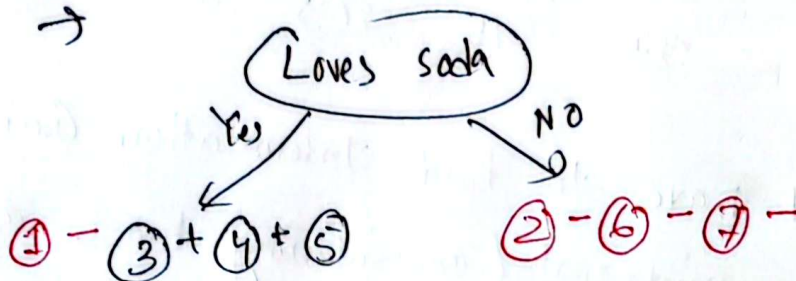
There are total 3
Loves popcorn = NO

Information Gain of Loves popcorn

$$\begin{aligned}\text{Info Gain}(\text{Loves popcorn}) &= \text{Info}(D) - \text{Info}(\text{Loves popcorn}) \\ &= 0.985 - 0.857 \\ &= 0.128\end{aligned}$$

'Loves soda feature':

We can directly find the $\text{Info}(\text{Loves soda})$ or weighted entropy \rightarrow



$$\begin{aligned}\text{Info}(\text{Loves Soda}) &= \frac{4}{7} \left[-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right] \\ &\quad + \frac{3}{7} \left[-\frac{0}{3} \log_2 \frac{0}{3} - \frac{3}{3} \log_2 \frac{3}{3} \right]\end{aligned}$$

$P(\text{Loves soda} = \text{Yes}) \rightarrow H(\text{Loves soda} = \text{Yes})$

$P(\text{Loves soda} = \text{No}) \rightarrow H(\text{Loves soda} = \text{No})$

$$= 0.464$$

$$\begin{aligned}\therefore \text{Info Gain}(\text{Loves Soda}) &= \text{Info}(D) - \text{Info}(\text{Loves soda}) \\ &= 0.985 - 0.464 \\ &= 0.521\end{aligned}$$

Numerical Feature 'Age' →

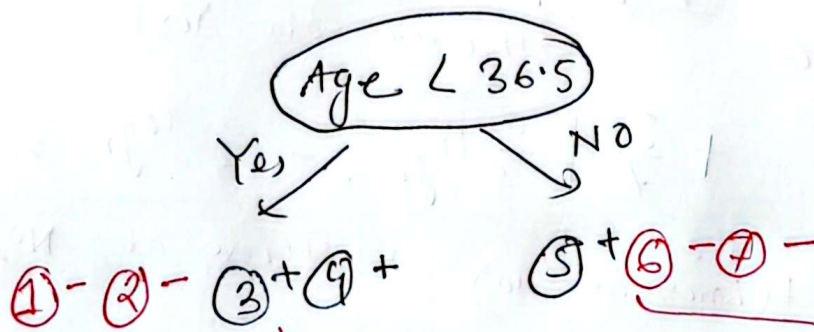
As we did in Mini Impurity, we will sort the numerical features/column first in ascending order and then find the midpoints / Averages of adjacent data points.

Sorted order \rightarrow 7 12 18 35 38 50 83

mid points 9.5 15 26.5 36.5 44 66.5

Now we will have to find Information Gain for all of this midpoints (considering those as thresholds)

Let us find for Age < 36.5 only \rightarrow



$$I_{nto}(Age < 36.5) = \underbrace{\frac{4}{7} \left[-\frac{2}{4} \log_2 \frac{2}{4} - \frac{2}{4} \log_2 \frac{2}{4} \right]}_{P(Age < 36.5)} + \underbrace{\frac{3}{7} \left[-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right]}_{P(Age \geq 36.5)}$$

features. So, it is a

$$\text{Info Gain}(\text{Age} < 36.5) = 0.985 - 0.956 = 0.02$$

Similarly if we find information gain for all of the mid points we get \rightarrow

$$\begin{array}{lll} \text{Ig}(\text{Age} < 12.5) & \text{Ig}(\text{Age} < 15) & \text{Ig}(\text{Age} < 26.5) \\ 0.128 & 0.292 \checkmark & 0.02 \end{array}$$

$$\begin{array}{lll} \text{Ig}(\text{Age} < 36.5) & \text{Ig}(\text{Age} < 44) & \text{Ig}(\text{Age} < 56.5) \\ 0.02 & 0.292 \checkmark & 0.128 \end{array}$$

Among these the highest Information Gain is $\text{Ig}(\text{Age} < 15) = 0.292$ and $\text{Ig}(\text{Age} < 44) = 0.292$.
we can pick any one. Let's pick $\text{Age} < 15$.

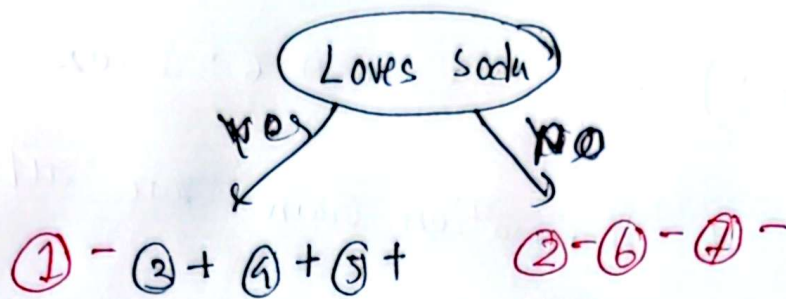
Info gain of ~~all~~ all of the features -

$$\text{Info Gain}(\text{Loves popcorn}) = 0.128$$

$$\text{Info Gain}(\text{Loves Soda}) = 0.521$$

$$\text{Info Gain}(\text{Age} < 15) = 0.292$$

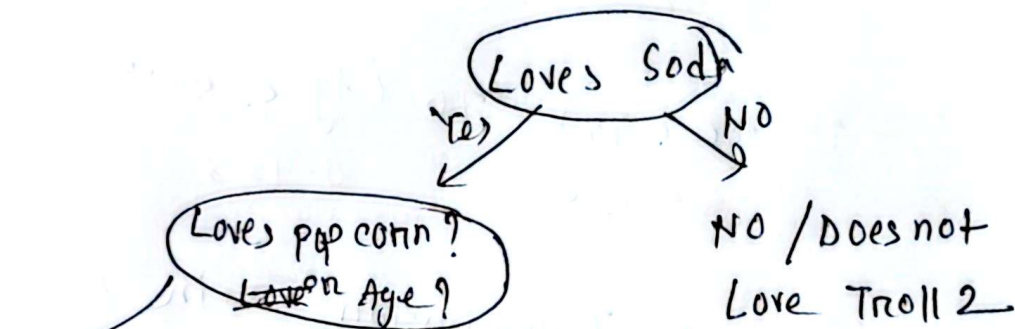
Among these three, Loves Soda has the highest Ig. So it will split the dataset better than the other features. So, it will be chosen as root of the tree.



impure leaf

We need to find another feature to split this branch.

pure leaf → we can assign class/ Label here



→ To decide this we will now consider the data points where Loves Soda = Yes only.

The new dataset →

	Loves popcorn	Loves Soda	Age	Loves Troll 2
①	Yes	Yes	7	NO
③	NO	Yes	18	Yes
④	No	Yes	35	Yes
⑤	Yes	Yes	38	Yes

Now we will find information gain for Loves popcorn and Age features for this dataset. Loves Soda has been used as node already, we will skip it.

Let us find the Entropy of this Small dataset first

Entropy of dataset \rightarrow

$$H(S) / \text{Info}(D) = -P(\text{Yes}) \log_2 P(\text{Yes}) - P(\text{No}) \log_2 P(\text{No})$$

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

$$= 0.811$$

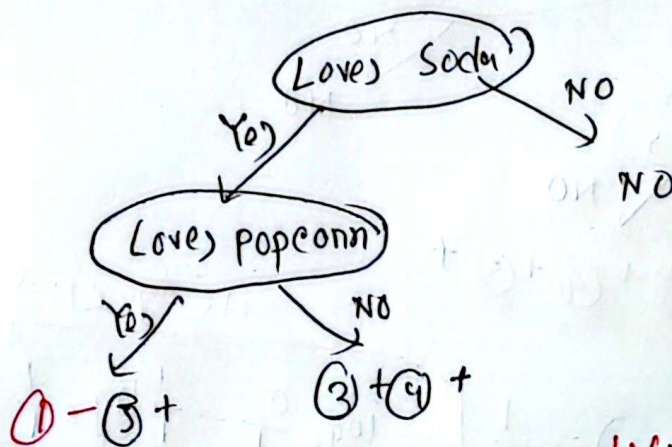
There are total 3

Yes in the Love Troll 2 /

target feature among
4 data points

There is only one
No in the target
feature.

Information Gain of Love popcorn \rightarrow



$$\text{Info}(\text{Love popcorn}) = \frac{2}{4} \left[-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right] +$$

$$P(\text{Love popcorn} = \text{Yes}) \frac{2}{4} \left[-\frac{2}{2} \log_2 \frac{2}{2} - \frac{0}{2} \log_2 \frac{0}{2} \right]$$

$$P(\text{Love popcorn} = \text{No}) H(\text{Love popcorn} = \text{No})$$

$$= 0.50$$

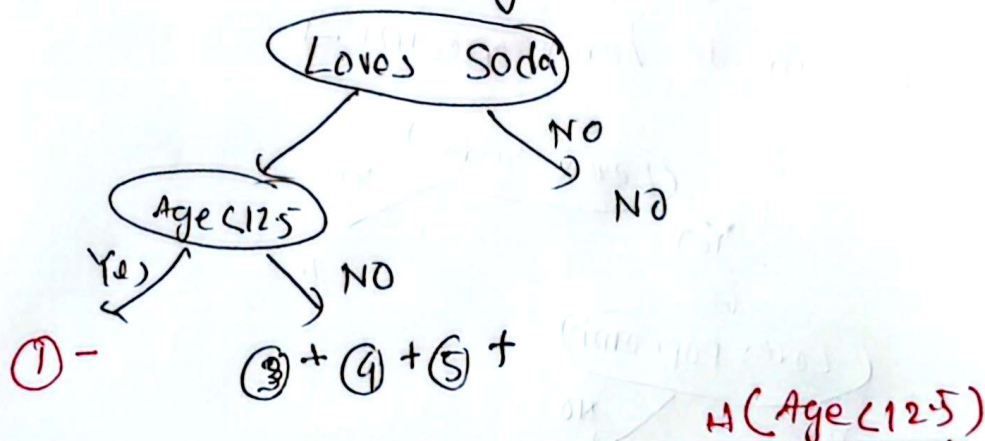
$$\begin{aligned}\therefore \text{Info Gain (Loves popcorn)} &= \text{Info (D)} - \text{Info (Loves popcorn)} \\ &= 0.811 - 0.50 \\ &= 0.311\end{aligned}$$

Information Gain of Age \rightarrow

Age in Sorted order \rightarrow 7 18 35 38

mid points \rightarrow 12.5 26.5 36.5

Let us find info gain of Age < 12.5



$$\begin{aligned}\text{Info (Loves Age } < 12.5) &= \frac{1}{4} \left[-\frac{0}{1} \log_2 \frac{0}{1} - \frac{1}{1} \log_2 \frac{1}{1} \right] + \\ &\quad \underbrace{\frac{3}{4} \left[-\frac{3}{3} \log_2 \frac{3}{3} - \frac{0}{0} \log_2 \frac{0}{0} \right]}_{\substack{P(\text{Age} < 12.5) \leftarrow H(\text{Age} < 12.5) \\ P(\text{Age} > 12.5) \leftarrow H(\text{Age} > 12.5)}} \\ &= 0.80\end{aligned}$$

$$\begin{aligned}\therefore \text{Info Gain (Age } < 12.5) &= \text{Info (D)} - \text{Info (Age } < 12.5) \\ &= 0.811 - 0.00 = 0.811\end{aligned}$$

Similarly In of other thresholds,

Age < 12.5

0.811

Age < 26.5

0.311

Age < 36.5

0.123

So Info Gain (Age < 12.5) is greater than other too.

We can also decide it by seeing the weighted

Entropy, Info (Age < 12.5) = 0.00. This is the

lowest Entropy possible. So, Age < 12.5 will be chosen

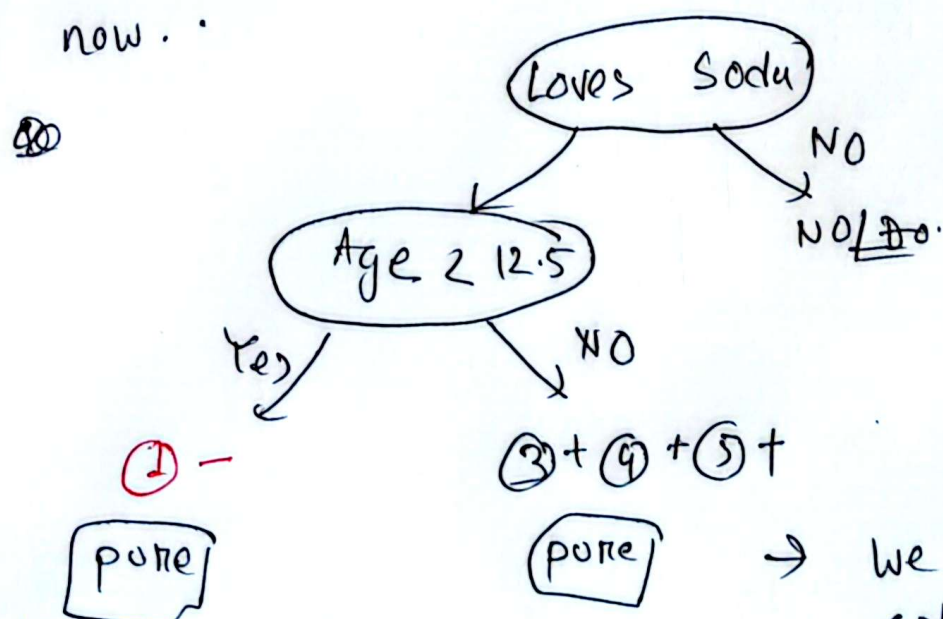
So, Now In of all features \rightarrow

Info Gain (Loves popcorn) = 0.311

Info Gain (Age < 12.5) = 0.811 **winner**

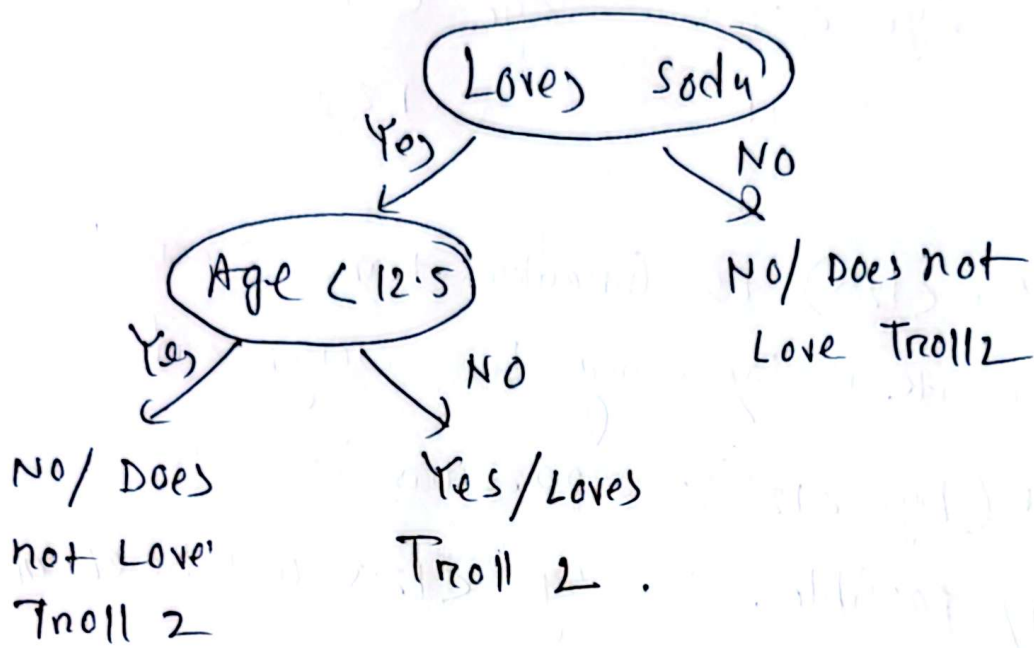
Hence Age < 12.5 has the highest IG. So

it will be chosen as the node of the tree now.



\rightarrow We do not need to split any further.

∴ The Final Decision Tree →



This is the same tree we found in Classification tree (Gini impurity) Process.