

Classification

Tree



PZR

Loves Popcorn	Loves Soda	Age	Loves Troll 2
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

⇒ We will predict if someone loves troll 2 or not.

→ We will solve the problem using
classification tree.

\Rightarrow When it comes to tree we need to have a root node.

— So what will be the root node?

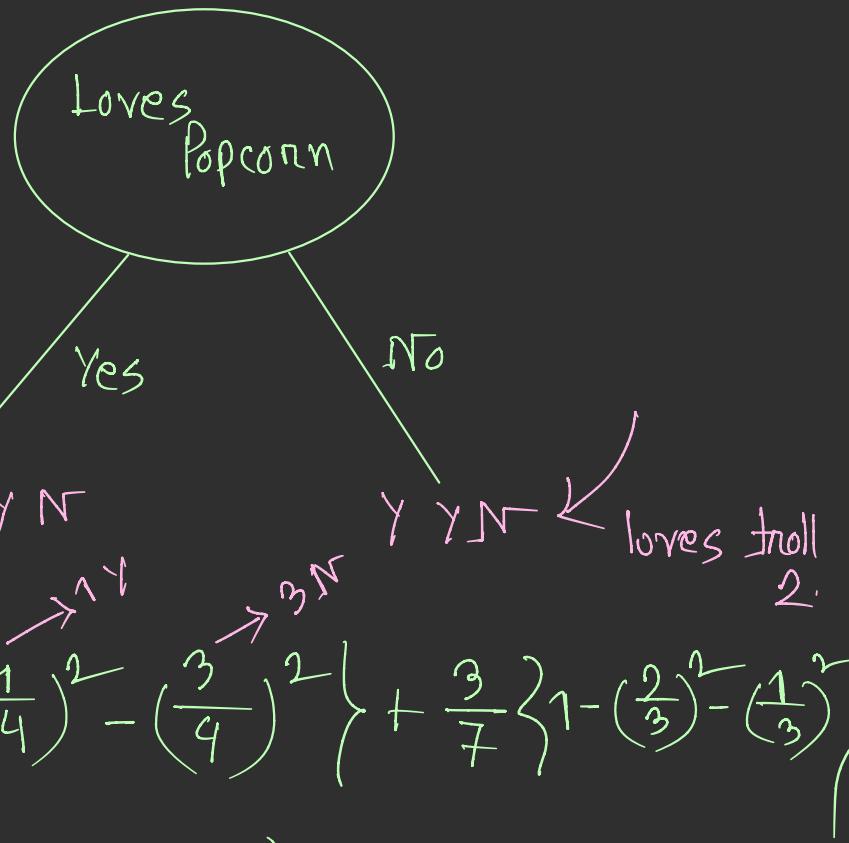
— How do we get this?

\Rightarrow There are multiple ways but we will use "Gini Impurity"

\rightarrow To calculate Gini Impurity

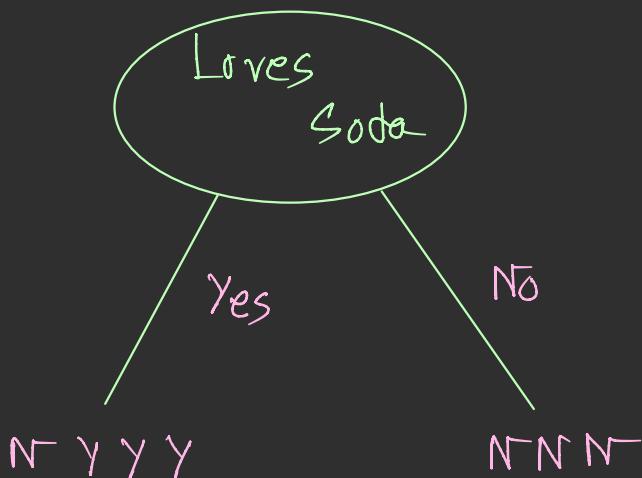
— Pick an attribute from the dataset

Let's say we pick "Loves Popcorn".



$$\Rightarrow \frac{4}{7} \left\{ 1 - \left(\frac{1}{4} \right)^2 - \left(\frac{3}{4} \right)^2 \right\} + \frac{3}{7} \left\{ 1 - \left(\frac{2}{3} \right)^2 - \left(\frac{1}{3} \right)^2 \right\}$$

$$\Rightarrow \text{Gini} (\text{Loves Popcorn}) = 0.405$$



$$\frac{4}{7} \left\{ 1 - \left(\frac{3}{4} \right)^2 - \left(\frac{1}{4} \right)^2 \right\} + \frac{3}{7} \left\{ 1 - \left(\frac{0}{3} \right)^2 - \left(\frac{3}{3} \right)^2 \right\}$$

$$\text{Gini}(\text{Loves Soda}) = 0.214$$

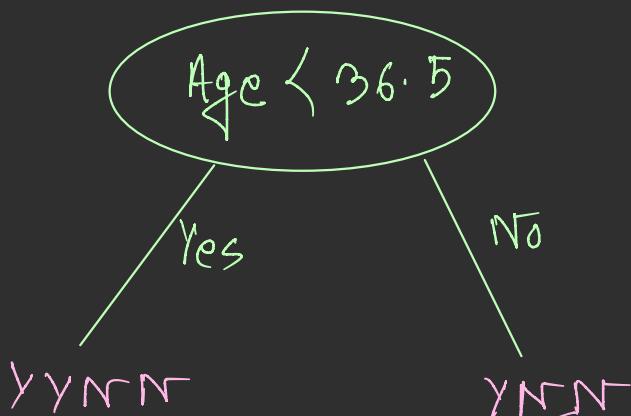
\Rightarrow Now we will do the same for the age attribute.

- But this attribute contains numeric value
So the process will be different.

→ first sort the numbers in ascending order, then compute average of two consecutive values.
 → also the label accordingly.

So,

$$\begin{array}{ccccccccc}
 7 & \downarrow & 12 & \downarrow & 18 & \downarrow & 35 & \downarrow & 38 & \downarrow & 50 & \downarrow & 83 \text{ (Age)} \\
 9.5 & & 15 & & 26.5 & & 36.5 & & 44 & & 66.5 \text{ (Avg)}
 \end{array}$$



$$G_1(\text{Age} < 36.5)$$

$$= \frac{4}{7} \left\{ 1 - \left(\frac{2}{4} \right)^2 = \left(\frac{2}{4} \right)^2 \right\} + \frac{3}{7} \left\{ 1 - \left(\frac{1}{3} \right)^2 = \left(\frac{1}{3} \right)^2 \right\} \\
 = 0.476$$

Similarly we will find the Gini impurity for all the average values of age.

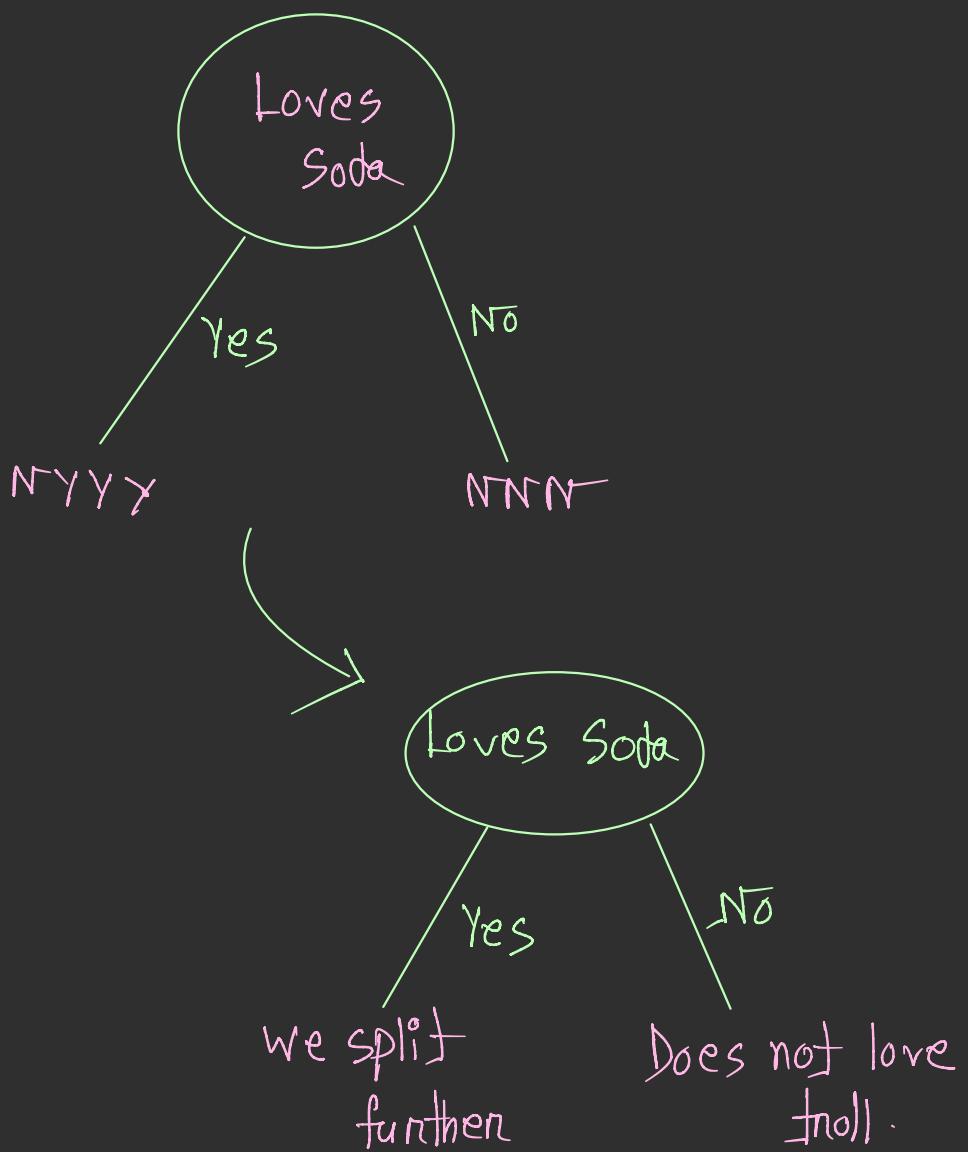
- Among all these $A < 15$ or $A < 44$ will have the lowest impurity. So we pick random one.
- Here we pick

$$G(Age < 15) - 0.343$$

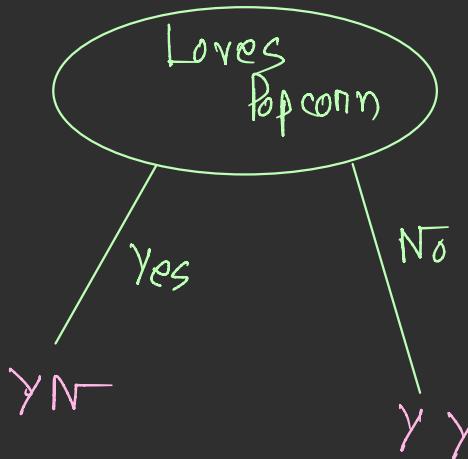
$$G(Loves \text{ popcorn}) - 0.405$$

$$G(Loves \text{ Soda}) - 0.214$$

\Rightarrow So, for our dataset Loves Soda will be the root node.



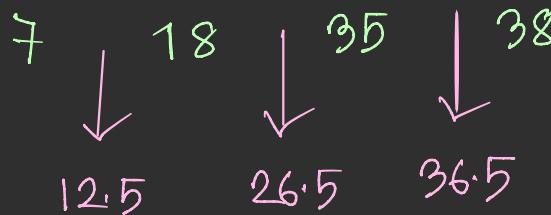
\Rightarrow We will continue for Love Soda = Yes
- Gini Impurity of Loves Popcorn



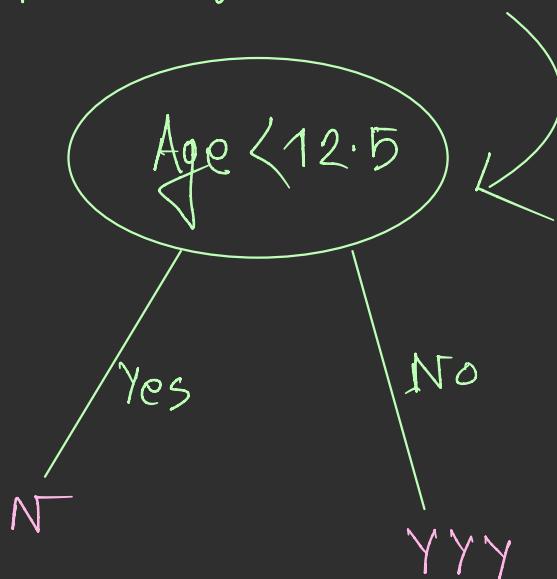
$$\text{Gini}(\text{Loves popcorn}) = 0.25$$

$$\frac{2}{4} \left\{ 1 - \left(\frac{1}{2} \right)^2 - \left(\frac{1}{2} \right)^2 + \frac{2}{4} \right\} 1 - \left(\frac{2}{2} \right)^2 - \left(\frac{0}{2} \right)^2 \} \\ = 0.5$$

Age

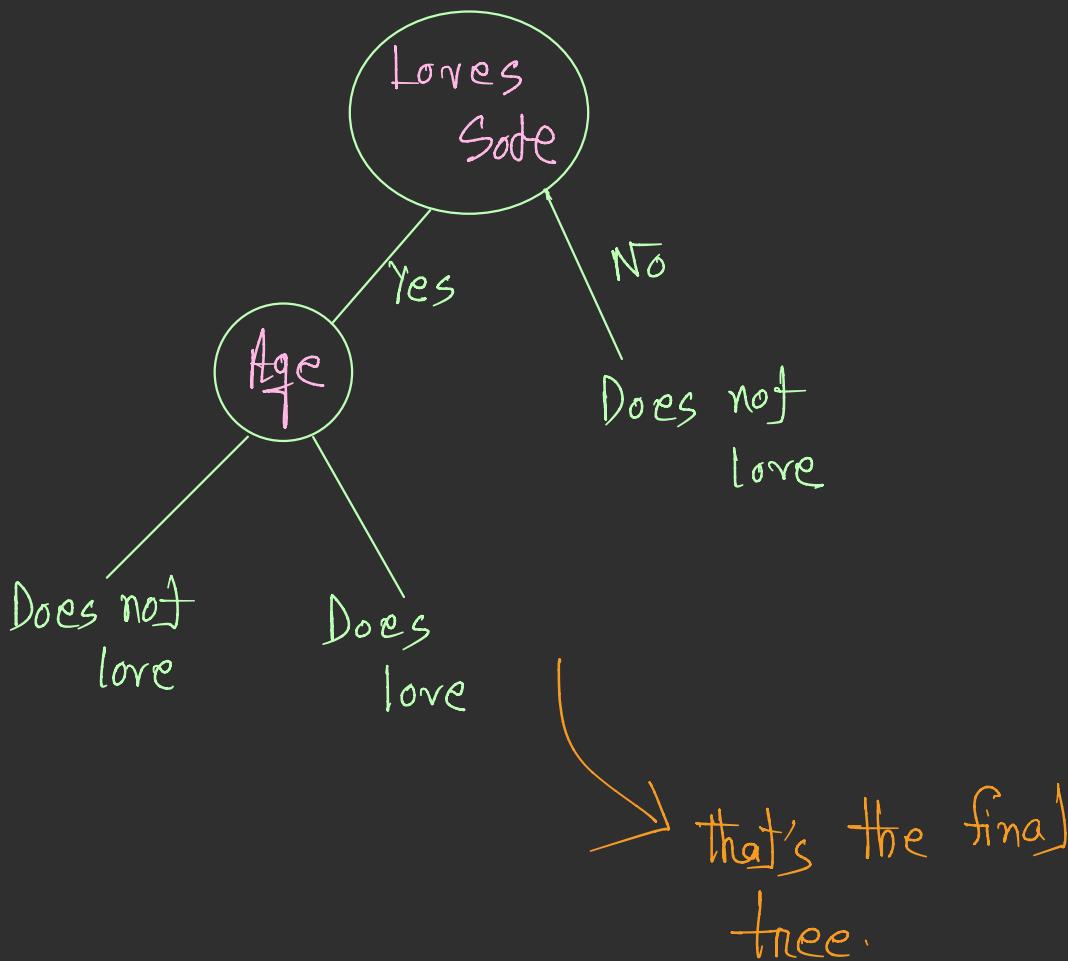


⇒ For $\text{Age} < 12.5 = 0$ is the minimum Gini impurity among all the age values.



→ Between Age < 12.5 and Loves Popcorn

Age will be picked.



⇒ Solve the same problem using
Information Gain.

$$\text{Entropy} = \sum_{i=1}^n -P \log_2 P$$

→ i is the number
 of labels.
 → Computes 'impurity'

→ Entropy (troll 2)

$$= -P(\text{Yes}) \log_2 P(\text{Y}) - P(\text{N}) \log_2 P(\text{N})$$

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

$$= 0.985$$

→ Entropy (Loves popcorn = Yes)

$$= -P(\text{Yes}) \log_2 P(\text{Y}) - P(\text{N}) \log_2 P(\text{N})$$

$$= -\frac{1}{4} \log_2 \frac{1}{4} - \frac{3}{4} \log_2 \frac{3}{4} = 0.811$$

\rightarrow Entropy ($\text{Loves popcorn} = \text{No}$)

$$= -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} = 0.918$$

\rightarrow Information Gain (Loves popcorn)

$$= \text{Entropy}(\text{froll 2}) - [P(\text{Loves popcorn} = \text{Yes}) *$$

$$\text{Entropy}(\text{Loves popcorn} = \text{Yes}) - P(\text{Loves popcorn} = \text{Yes}) * \text{Entropy}(\text{Loves popcorn} = \text{No})]$$

$$= 0.985 - \left(\frac{4}{7} * 0.811 \right) - \left(\frac{3}{7} * 0.918 \right)$$

$$= 0.128$$

→ Information Gain (Loves Soda)

$$= 0.985 - \left[\frac{4}{7} * \left(-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right) \right] -$$
$$\left[\frac{3}{7} * \left(-\underbrace{\frac{0}{3} \log_2 \frac{0}{3}}_0 - \underbrace{\frac{3}{3} \log_2 \frac{3}{3}}_0 \right) \right]$$

$$= 0.985 - 0.463 - 0$$

$$= 0.522$$

- Next calculate information gain for Age like before.
- Among these three attribute Loves Soda has the max information gain.

So, root will be Loves Soda.

→ Then continue the splitting until you meet the conditions.

