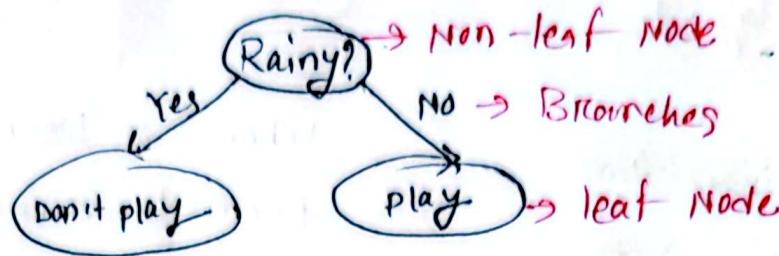


Classification Tree

Decision Tree -

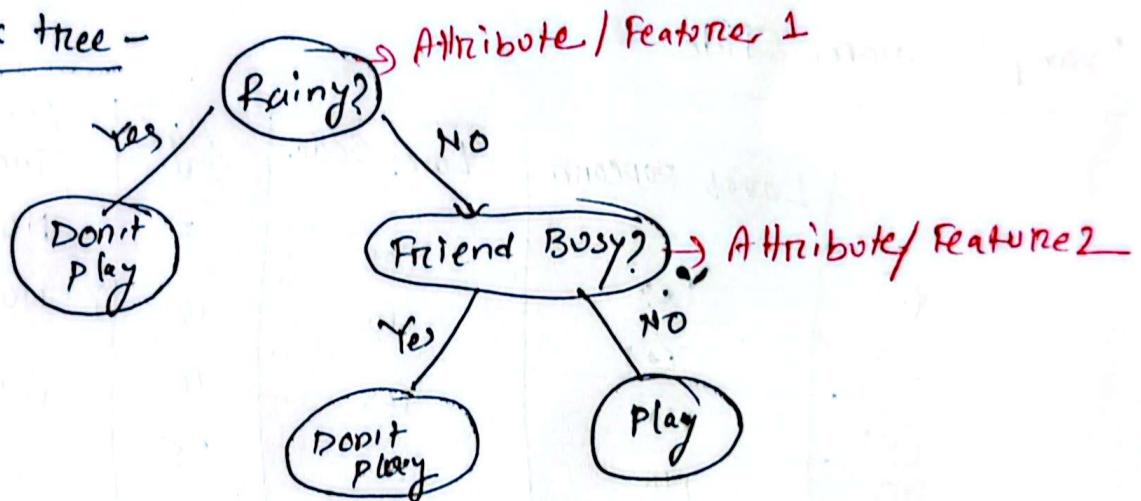


Non-leaf Nodes are Attributes / Features

leaf Nodes are Labels / classes

Branches are outcomes of a feature (can be more than 2)

more complex tree -

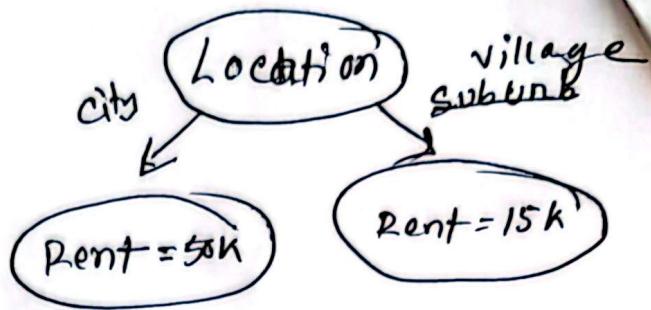
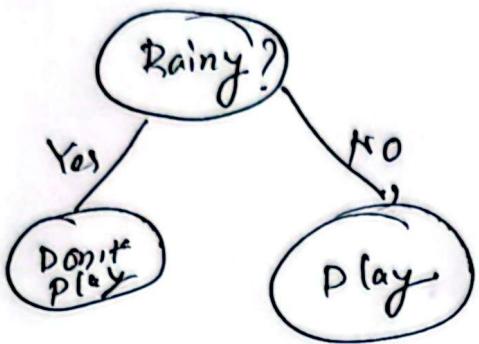


This tree can be represented as the following Data Set -

| | Rainy? | Friend Busy? | Play? |
|----------|--------|--------------|---------------|
| Sample 1 | Yes | Yes | No/Don't play |
| Sample 2 | Yes | No | No |
| Sample 3 | No | Yes | No |
| Sample 4 | No | No | Play/Yes |

Attributes / Features

We usually derive model/tree from this type of Dataset.



When a Decision tree classifies things into categories - **Classification Tree**

When a Decision tree predicts numeric values - **Regression Tree**

The features can be both numerical and categorical.

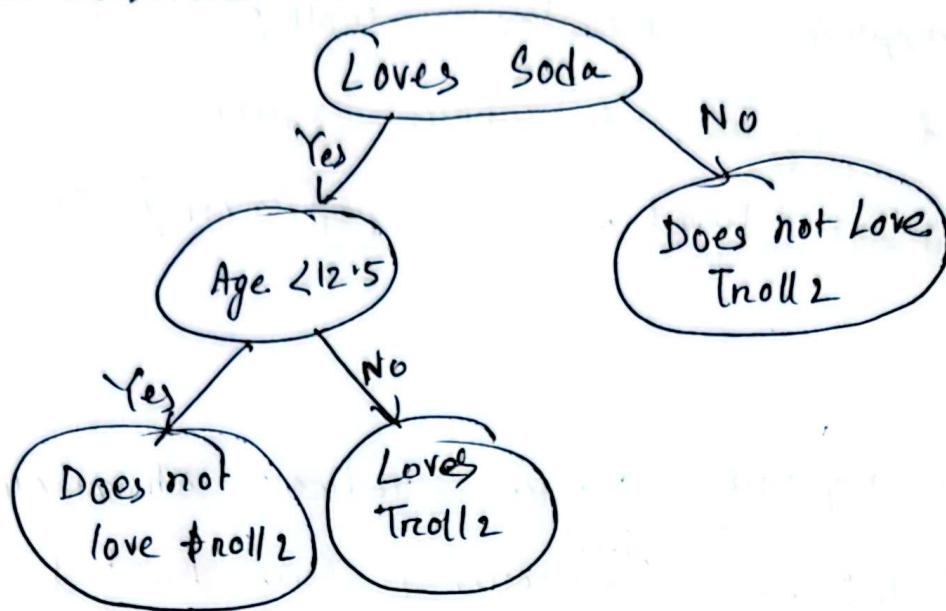
Example from slide -

| | Loves Popcorn | Loves Soda | Age | Loves Troll12 |
|---|---------------|------------|-----|---------------|
| 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 |

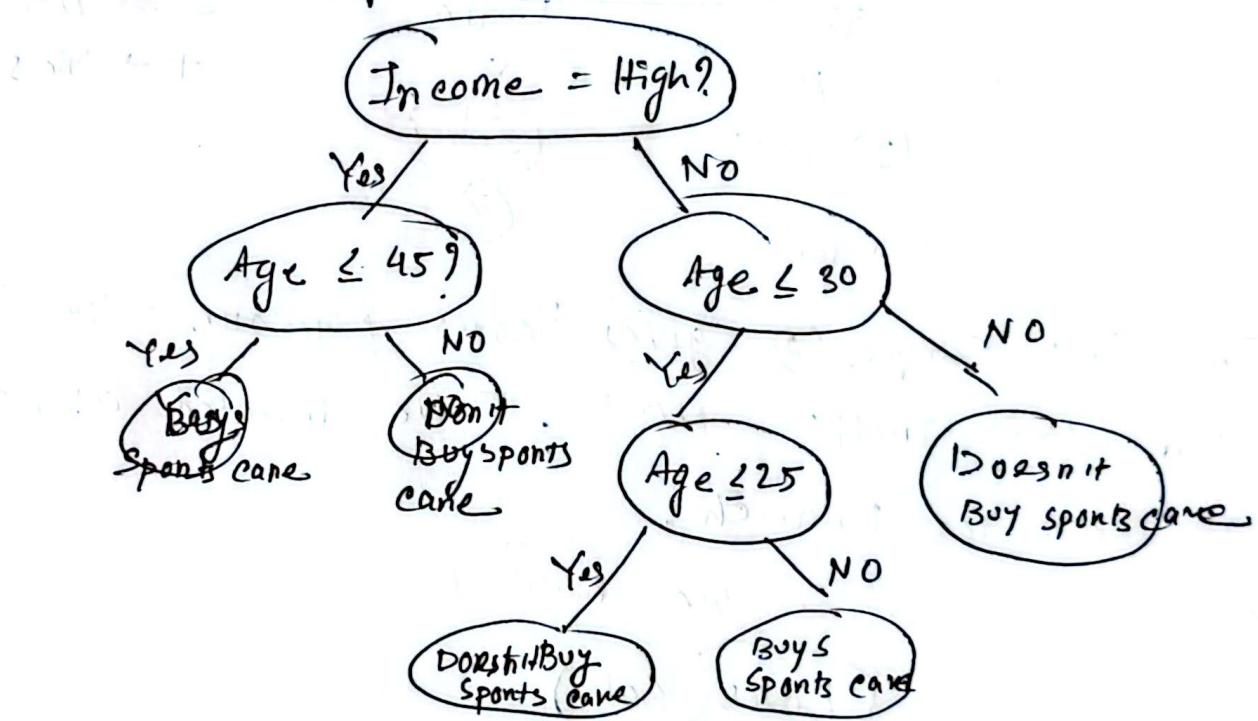
Categorical ↓ Numerical

We want to predict if someone loves troll12 or not using Classification Tree.

The classification tree —



Another tree Example —

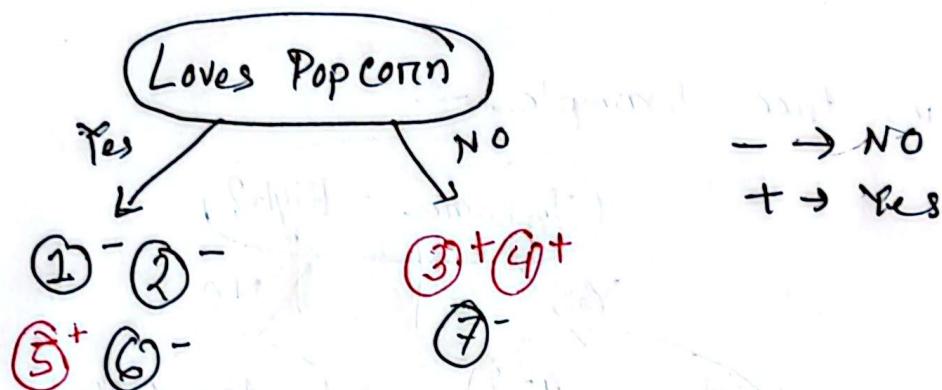


Numerical features' thresholds can be different for the same data at different levels of the tree.

From a classification tree we can predict a label / class for a new / unknown sample.

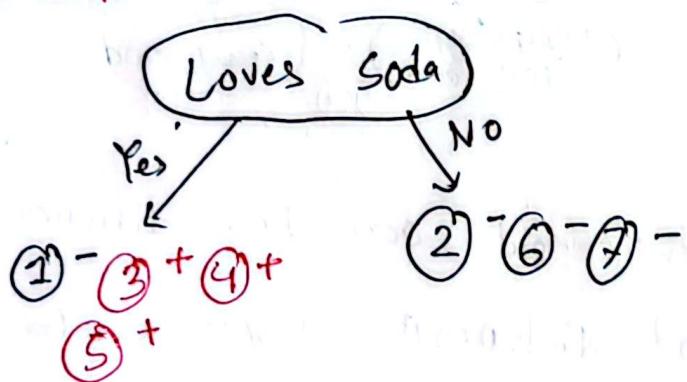
Now if we want to build a tree for the dataset to predict whether someone loves troll 2, we at first need to figure out which feature should be at the Root level → Loves popcorn / Loves soda / Age.

How well 'Loves Popcorn' feature predicts whether someone likes troll 2 → [if 'Loves Popcorn' is the Root]



The first sample (1) gives 'Loves popcorn = Yes' and 'Loves Troll 2 = NO'. We denoted as (1)⁻ on the Yes Branch.

[if 'Loves Soda' is the Root]

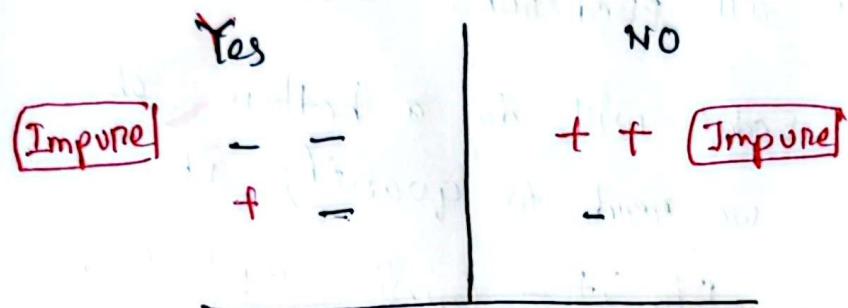


Neither of the categorical features does the perfect job to predict whether someone

Loves Troll2 .

We can see / visualize this graphically ←

Loves Popcorn -

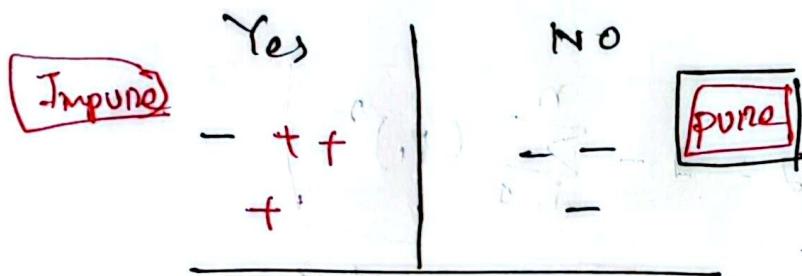


Loves Troll2

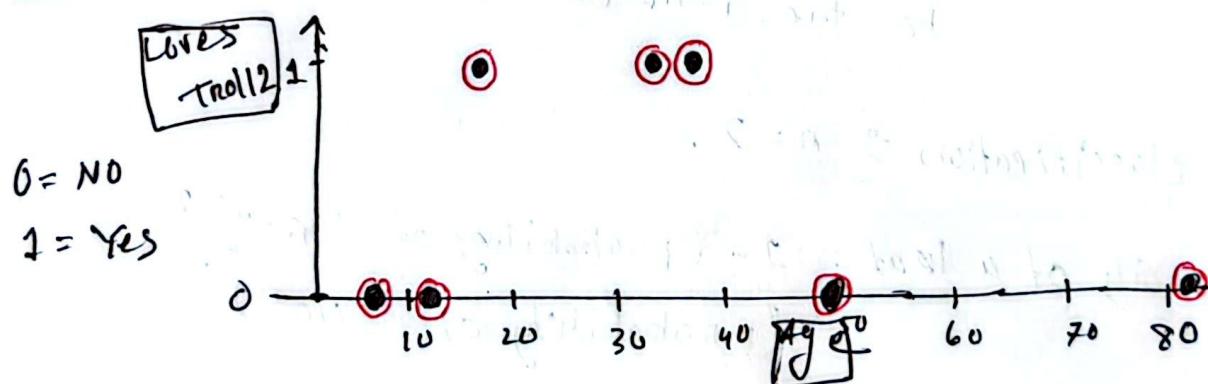
- NO

+ Yes

Loves Soda -



Age -



We cannot fit a best fit line for these points to perfectly classify Loves Troll2. So Age feature cannot perfectly predict as well.

As the leaves of 'Loves pop corn' and left leaf of 'Loves Soda' contains mixture of Yes, NO (+, -), they are called impure. Only the right leaf of 'Loves Soda' is pure as all predictions are NO / -.

It seems that 'Loves Soda' will do a better job to predict the Labels. But we need to quantify it. There are methods to quantify it - **Gini Impurity**, **Entropy**, **Information Gain**.

Gini Impurity:

Formula \rightarrow
$$\text{Gini Impurity} = 1 - \sum_{i=1}^n (P_i)^2$$

(of a leaf)

n = number of Label or class

⇒

P_i = Probability of sample being belonging to the i-th class.

For Binary Classification $\Rightarrow n = 2$.

$$\therefore \text{Gini Impurity of a leaf} = 1 - (\text{probability of "Yes"})^2 - (\text{probability of "No"})^2$$

Loves Popcorn

Yes

NO

① - ② - ③ + ④ -
Yes Yes No Yes

⑤ + ⑥ + ⑦ -

Gini impurity of Leaf

$$\begin{aligned} &= 1 - (P(\text{Yes}))^2 - (P(\text{No}))^2 \\ &= 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \\ &= 0.375 \end{aligned}$$

Gini impurity of leaf

$$\begin{aligned} &= 1 - (P(\text{Yes}))^2 - (P(\text{No}))^2 \\ &= 1 - \left(\frac{1}{3}\right)^2 - \left(\frac{2}{3}\right)^2 \\ &= 0.444 \end{aligned}$$

As there are different numbers of people at each leaf (4 at left leaf, 3 at right leaf), we will calculate Total Gini Impurity by taking weighted average of the Gini of each leaf.

Total Gini impurity = weighted average of Gini impurities for the leaves.

For binary classification

Total Gini impurity = $P(\text{Left leaf}) \times \text{Impurity of left leaf} + P(\text{Right leaf}) \times \text{impurity of right leaf.}$

∴ Total Gini impurity (Loves Popcorn)

$$\begin{aligned} &= \frac{4}{7} \times 0.375 + \frac{3}{7} \times 0.444 \\ &= 0.405 \end{aligned}$$

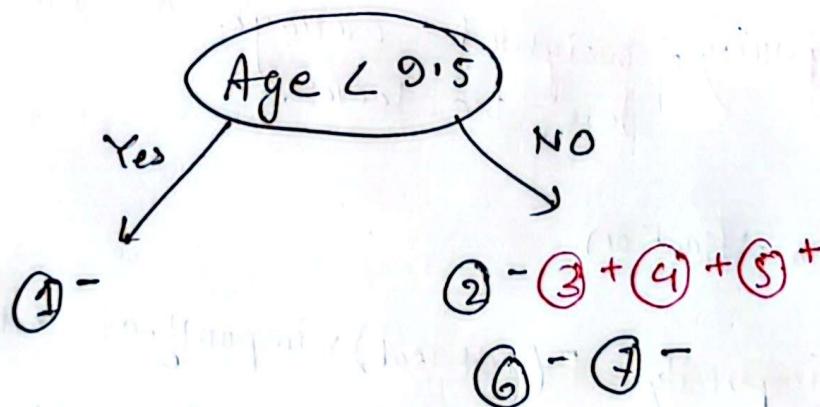
Similarly, total Gini impurity (Loves Soda) = 0.214.

Gini impurity of Numeric Data:

For the Age Feature - [Sort the column ~~Dataset by~~ feature Data set by Age. [Already sorted]. Then calculate average of all adjacent rows / Age of people]

$$\begin{array}{ccccccccc} 7 & 12 & 18 & 35 & 38 & 50 & 83 \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ 9.5 & 15 & 26.5 & 36.5 & 49 & 66.5 \end{array}$$

Now we have to calculate gini impurity for by setting each ^{average} Age as threshold.



$$\begin{aligned} \text{Total Gini impurity (Age < 9.5)} &= \frac{1}{7} \left\{ 1 - \left(\frac{0}{7}\right)^2 - \left(\frac{1}{7}\right)^2 \right\} \\ &\quad + \frac{6}{7} \left\{ 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2 \right\} \\ &= \frac{1}{7} \times 0 + \frac{6}{7} \times 0.5 \\ &= 0.429 \end{aligned}$$

Similarly Gini impurity of all thresholds \rightarrow

| | | | | | | |
|-----|----|------|------|----|------|----|
| 7 | 12 | 18 | 35 | 38 | 50 | 83 |
| ↓ | ↓ | ↓ | ↓ | ↓ | ↓ | ↓ |
| 9.5 | 15 | 26.5 | 36.5 | 44 | 66.5 | |

Total sum for Age < 15 and Age < 50 ~~am.~~ is

0.343 and this is the smallest.

Small impurity = Better feature to predict the class.

We can pick either anyone.

so calculated Gini impurity values

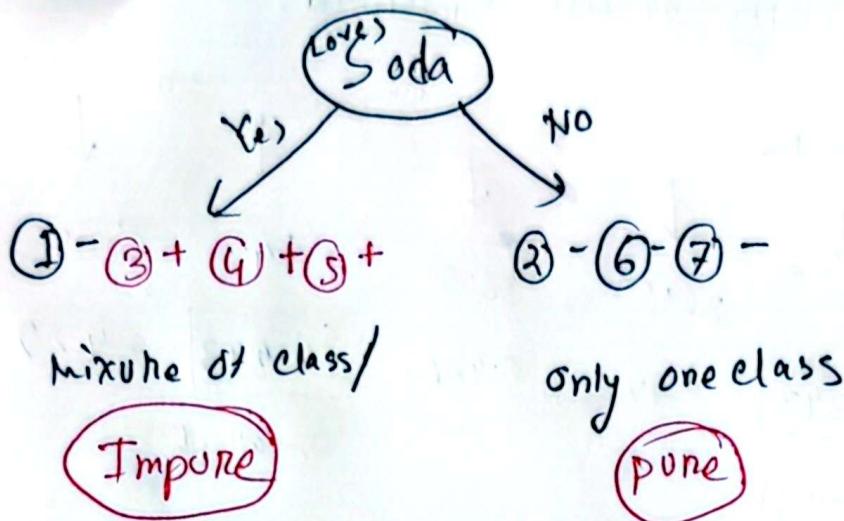
Total Gini (Loves popcorn) = 0.405

$$\text{Total Gini (Loves Soda)} = \cancel{0.343} 0.214$$

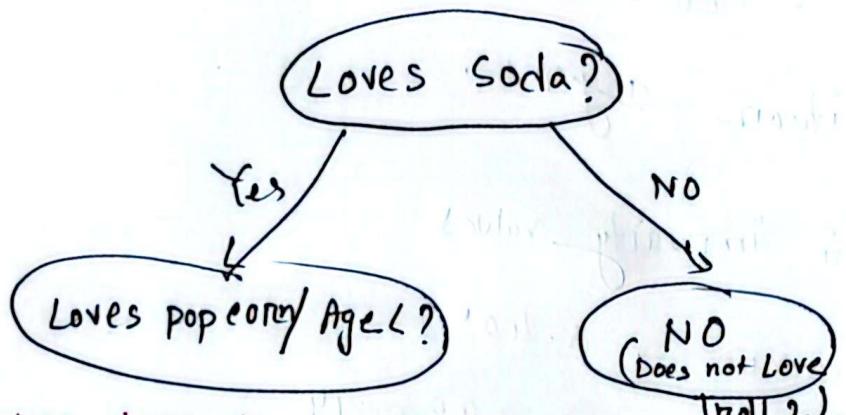
$$\text{Total } \text{C}_1\text{Ni} (\text{Age} < 15) = 0.343.$$

smallest aini impurity \rightarrow Loves Soda. [it's loves has lowest impurity]

so we put loves soda at the top. Root of the classification tree.



Now we have to figure out what will be the left node = Loves popcorn? Age < ?



* We have to calculate

Initial impurity of Loves Soda and

Age where Loves Soda = Yes

→

| | Loves Pop Corn | 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 |

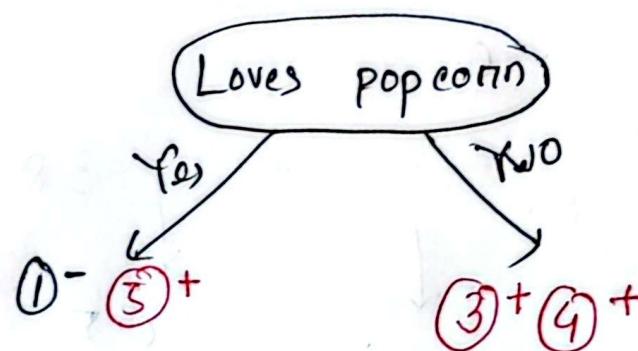
⑥

[NOTE] We will skip the categorical features if it has already been used as node, but can reuse numerical features with different threshold multiple times.

Data set with "Loves Soda = Yes" →

| | 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 |

Gini Impurity of Loves popcorn -



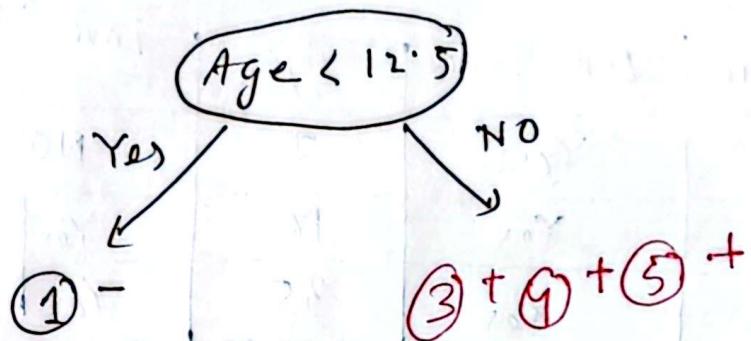
$$\begin{aligned} \text{Gini}(\text{Loves popcorn}) &= \frac{2}{4} \left\{ 1 - \left(\frac{1}{2}\right)^2 - \left(\frac{1}{2}\right)^2 \right\} + \frac{2}{4} \left\{ 1 - \left(\frac{2}{2}\right)^2 - \left(\frac{0}{2}\right)^2 \right\} \\ &= 0.25 \end{aligned}$$

Gini of Age -

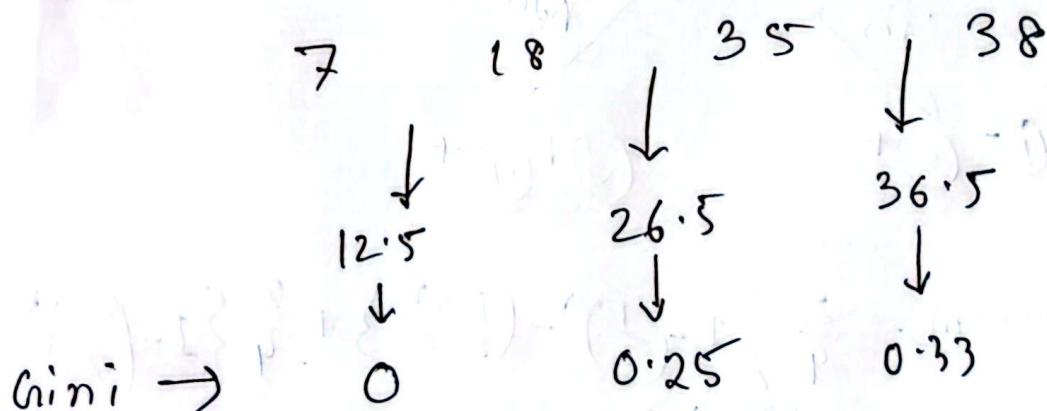
| | | | |
|---|----|----|----|
| 7 | 18 | 35 | 38 |
| ↓ | ↓ | ↓ | ↓ |

Average of
Adjacent rows.

| | | |
|------|------|------|
| 12.5 | 26.5 | 36.5 |
|------|------|------|



$$\text{Gini}(\text{Age} < 12.5) = \frac{1}{4} \left\{ 1 - \left(\frac{1}{1} \right)^2 - \left(\frac{1}{1} \right)^2 \right\} + \\ \frac{3}{4} \left\{ 1 - \left(\frac{3}{3} \right)^2 - \left(\frac{0}{3} \right)^2 \right\} \\ = 0$$



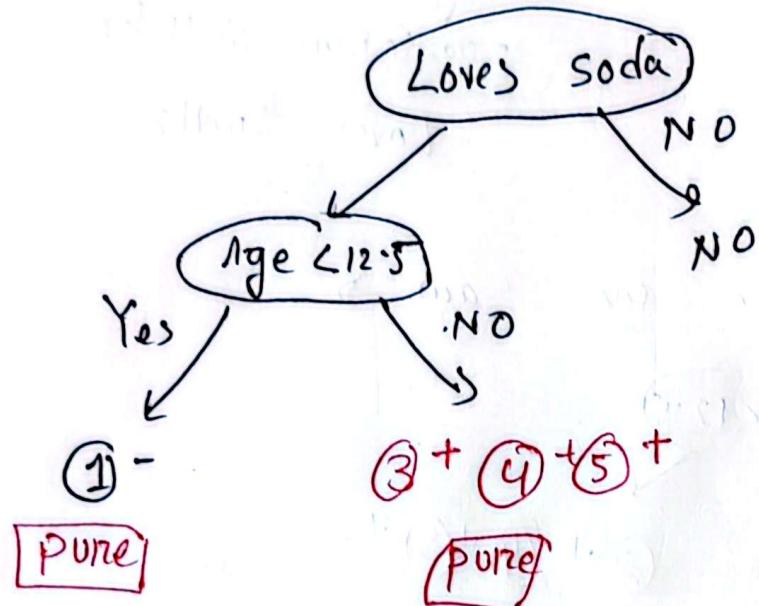
so smallest is $\text{Gini}(\text{Age} < 12.5) = 0$

calculated Gini \rightarrow

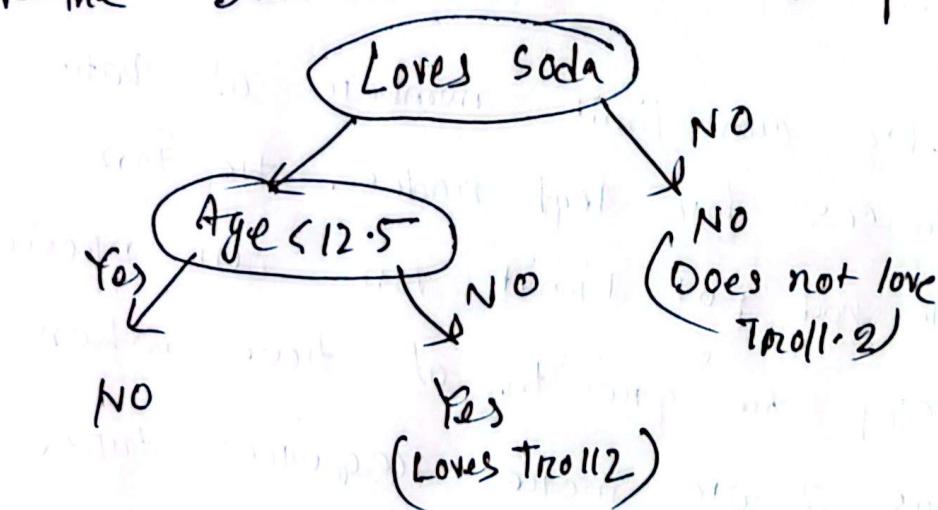
$$\text{Gini}(\text{Popcorn}) = 0.25$$

$$\text{Gini}(\text{Age} < 12.5) = 0$$

$\therefore \text{Gini}(\text{Age} < 12.5)$ is smallest



As we got pure leaves, we can stop building the tree now. So the Classification Tree for the Data set after assigning output values \rightarrow

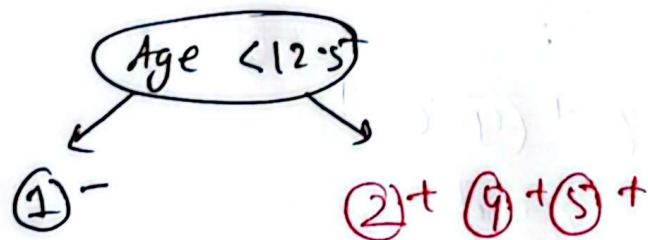


Predicting for unknown data -

| <u>Loves Popcorn</u> | <u>Loves Soda</u> | <u>Age</u> | <u>Loves Trollz</u> |
|----------------------|-------------------|------------|---------------------|
| NO | Yes | 20 | ? |

using the tree, the prediction will be
'Loves Trollz'

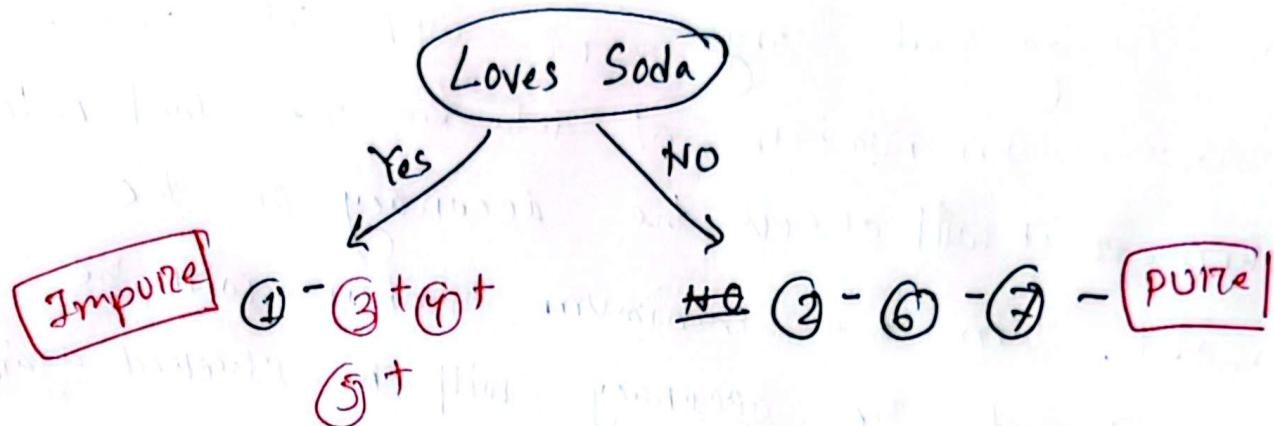
While building the tree we saw →



only one person / data could make it to this leaf.
This could possibly cause **Overshooting** and might not work ~~so~~ well for unknown data.

To solve this, we can limit number of data points or instances for leaf nodes. For example if we had put limit for our previous example as ~~stop the growth of tree when~~ leaf node has 3 or more requires three or more data/instance, then the tree

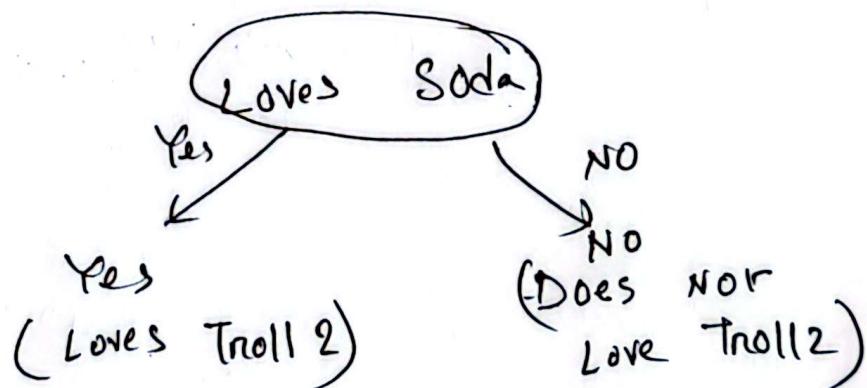
would stop at →



cause splitting these leaves ~~will~~ ^{can} cause leave us with leaves less than three data per leaves.

For pure leaves we can easily decide what will be the output value. But for impure leaves in this ~~case~~ we will count which output / label is maximum.

Hence positive/Yes occurs 3 times in the impure leaf, so we will put Yes.



The limit can be chosen by doing Cross Validation. We don't know which limit will give us better accuracy.

Cross Validation — After splitting the data set for training and testing, we will set the minimum number of data for the leaf node. Then it will check the accuracy of the model. Then the minimum number will be reduced and the accuracy will be checked again. We will continue doing it until we are getting high accuracy. If we get low accuracy we will stop.

For example we can start with limit = 10, then calculate model accuracy then limit = 9, calculate accuracy and continue until the accuracy is decreasing.