

Decision Trees

- Used for classification and Regression as well.

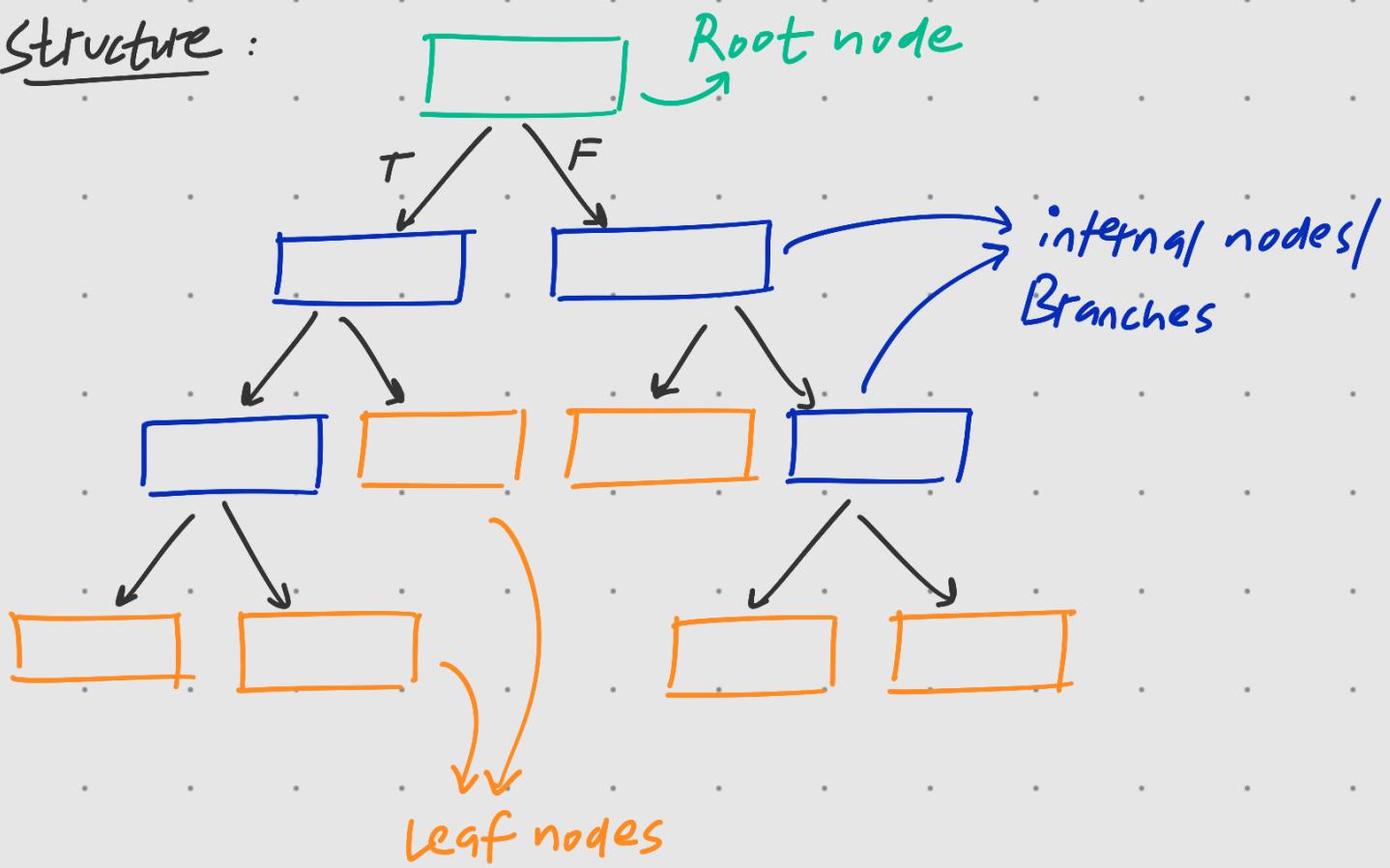
Attributes

→ Tree Structure
Decision nodes
Leaf nodes
Splitting
Entropy
Information Gain
Pruning

- Decision Tree makes a statement & then makes decision based on whether or not statement is True or False.
- Assumption: If, statement is True → go Left.
Statement is False → go Right.

Classification Tree

Structure :



Eg

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
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

To be predicted using, previous 3 columns.

Sol'n

what should be the node?

Determining it... Hit & Trial (initially)



Loves Popcorn

Gini
Imp
0.375

True

False

Loves cool as Ice	
yes	no
1 ①	1+1+1 ③

Loves cool as Ice	
yes	no
1+1 ②	1 ①

Gini imp
0.444

Doing the same thing for "loves Soda"...

Loves Soda

True

False

Loves cool as Ice	
yes	no
1+1 ③	1 ①

Gini imp
0.375

Loves cool as Ice	
yes	no
0 ①	1+1+1 ③

Gini imp
0

An absolute
Answer.

No clear decision

, as mixture of i.e. Impure
people likes & dislikes "cool as Ice"

In "loves Popcorn" \rightarrow Both leaves are impure
(less information gain)

In "loves Soda" \rightarrow only one leaf is impure
(more information gain
compared to previous.)

Does better job
at predicting who will
"love cool as Ice"

A way to quantify the impurity of leaves...

Gini Impurity

Start by calculating gini Impurity of individual leaves...

$$\text{Gini Impurity} = 1 - (P(\text{Yes}))^2 - (P(\text{no}))^2$$

\rightarrow Loves Popcorn (leaf 1)

$$\begin{aligned}\text{Gini impurity} &= 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 \\ &= 1 - \frac{1}{16} - \frac{9}{16} \\ &= \frac{6}{16} = 0.375\end{aligned}$$

Loves Popcorn (leaf 2)

$$\begin{aligned}\text{Gini impurity} &= 1 - \left(\frac{2}{3}\right)^2 - \left(\frac{1}{3}\right)^2 \\ &= 1 - \frac{4}{9} - \frac{1}{9} \\ &= \frac{4}{9} = 0.444\end{aligned}$$

→ Loves Soda (leaf 1)

$$\begin{aligned}\text{Gini Impurity} &= 1 - \left(\frac{3}{4}\right)^2 - \left(\frac{1}{4}\right)^2 \\ &= 1 - \frac{9}{16} - \frac{1}{16} = \frac{6}{16} \\ &= 0.375\end{aligned}$$

Loves Soda (leaf 2)

$$\begin{aligned}\text{Gini impurity} &= 1 - \left(\frac{0}{3}\right)^2 - \left(\frac{3}{3}\right)^2 \\ &= 1 - 0 - 1 \\ &= 0\end{aligned}$$

Both leaves of "Loves Popcorn" doesn't have same no. of people...

Total Gini impurity = weighted Avg. of gini impurities of the leaves

calculating the weight...

$$\text{weight} = \frac{\text{Total people in 1st leaf}}{\text{Total people in both leaves}} \times \text{Gini impurity}_{(\text{leaf 1})}$$

$$\begin{aligned} \text{Total Gini impurity} &= \frac{4}{(4+3)} \cdot (0.375) + \frac{3}{(4+3)} \cdot (0.444) \\ (\text{lives popcorn}) &= \boxed{0.405} ! \end{aligned}$$

$$\begin{aligned} \text{Total Gini impurity} &= \frac{4}{(4+3)} \cdot (0.375) + \frac{3}{(4+3)} \cdot (0) \\ (\text{lives Soda}) &= \boxed{0.214} ! \end{aligned}$$

→ Gini Impurity by Age

Age
7

12

18

35

38

50

83

- First we sort the numeric data.
- Then, we calculate average for all adjacent people.
- Lastly, calculate Gini impurities for each average age.

Age

7

$$9.5 \xrightarrow{\text{Gini imp.}} 0.429$$

12

$$15 \rightarrow 0.343$$

18

$$26.5 \rightarrow 0.476$$

35

$$36.5 \rightarrow 0.476$$

38

$$44 \rightarrow 0.343$$

50

$$66.5 \rightarrow 0.429$$

83

averages

$\boxed{\text{Age} < 9.5}$

True

False

cool as ice	
yes	no
0	1

cool as ice	
yes	no
3	3

gini imp.

$$\text{gini impurity} = 1 - \left(\frac{0}{7}\right)^2 - \left(\frac{1}{7}\right)^2$$

$$= 0$$

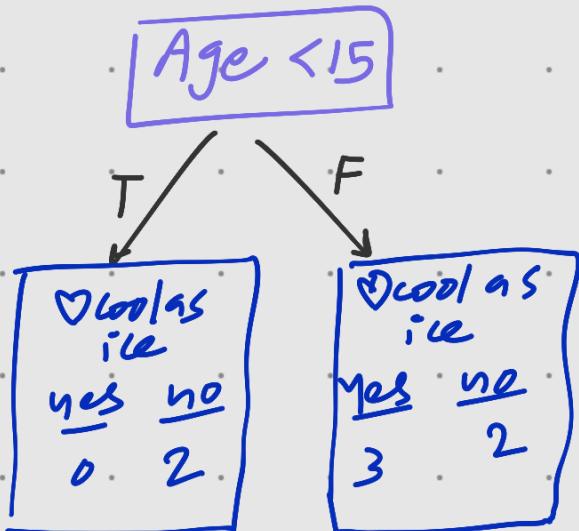
$$= 1 - \left(\frac{3}{6}\right)^2 - \left(\frac{3}{6}\right)^2$$

$$= 1 - \frac{7}{36} - \frac{9}{36}$$

$$= \frac{18}{36} = 0.5$$

$$\text{Total Gini impurity} = \frac{1}{1+6}(0) + \frac{6}{1+6}(0.5)$$

$$(Age) = \boxed{0.429}$$

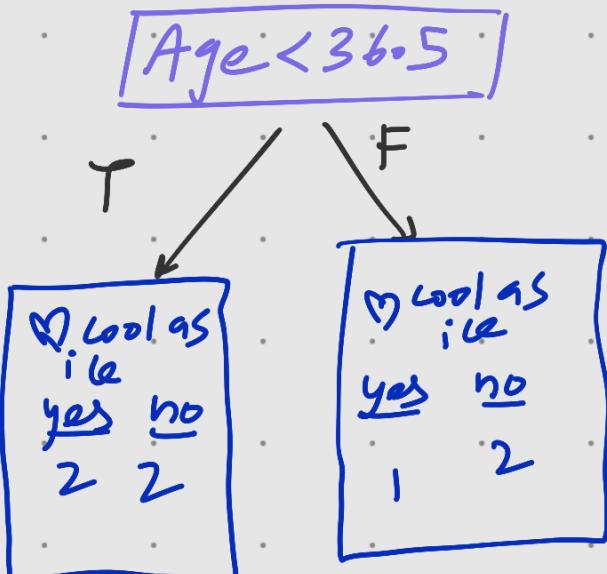
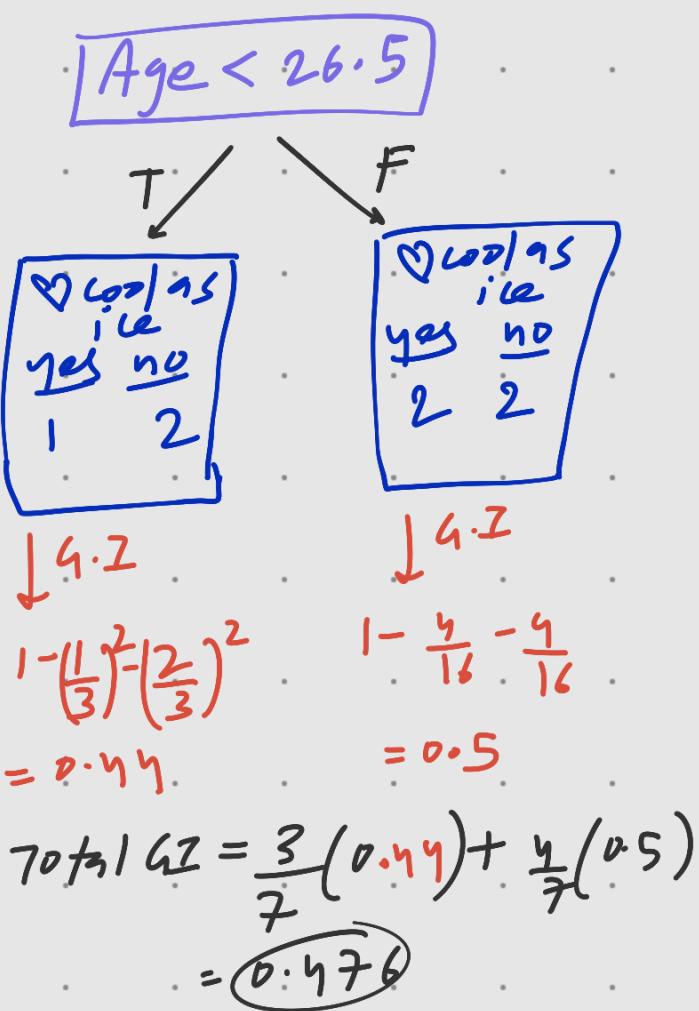


$\downarrow \text{gini imp}$

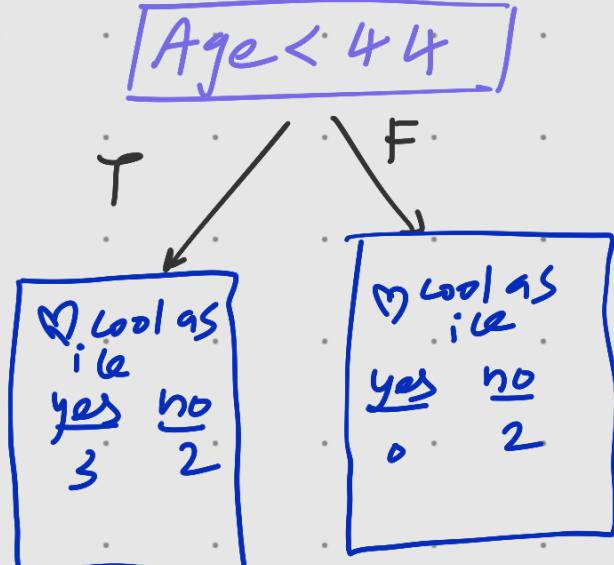
$$1 - \left(\frac{0}{2}\right)^2 - \left(\frac{2}{2}\right)^2 = 0$$

$$1 - \left(\frac{3}{5}\right)^2 - \left(\frac{2}{5}\right)^2 = 0.48$$

$$\begin{aligned} \text{Total Gini Impurity} &= \frac{2}{7}(0) + \frac{5}{7}(0.48) \\ &= 0.343 \end{aligned}$$

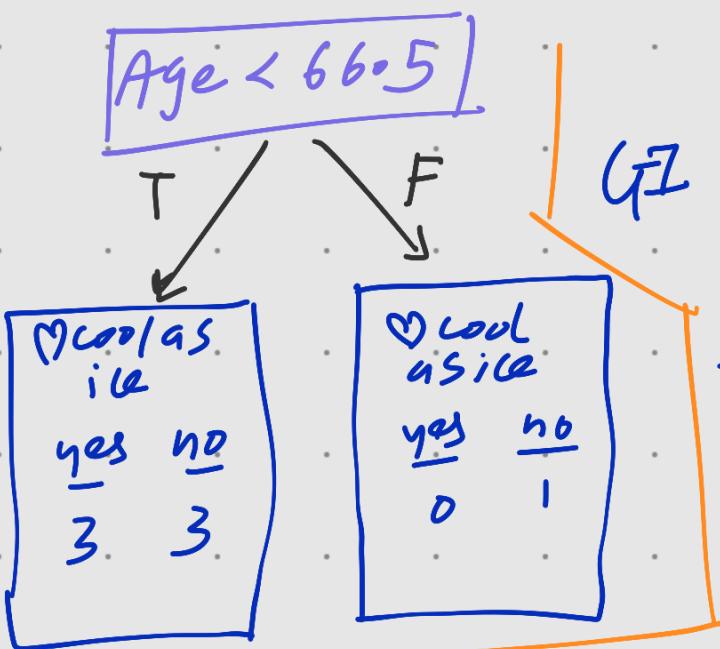


$$\begin{aligned} \text{Total GI} &= 0.476 \\ (\text{direct calculation}) \end{aligned}$$



$$\text{Total GI} = 0.343$$

Total Gini = 0.429



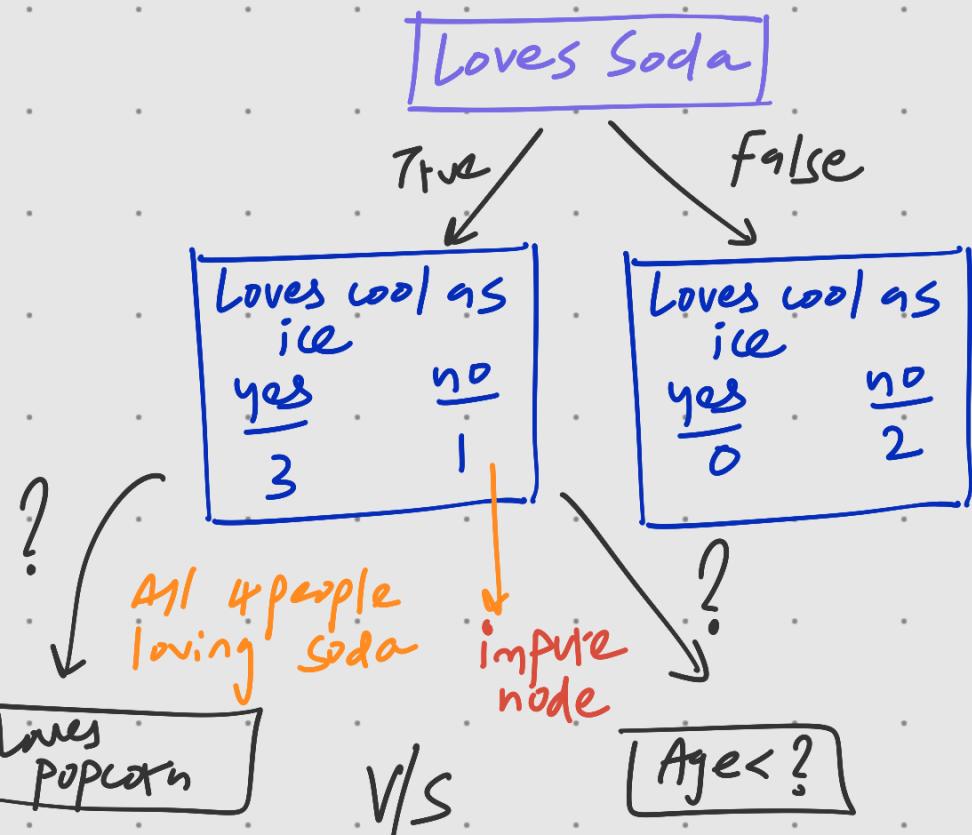
GI 0.343 is least, we'll pick that i.e (age < 15)

Total gini^o (loves popcorn) = 0.405
Impurity

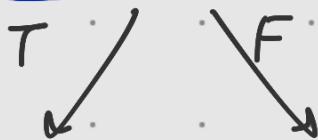
Total gini^o (loves soda) = 0.214
Impurity

↳ least impurity (i.e Root node)

Total gini^o (Age) = 0.343
Impurity



→ **Loves Popcorn**



we'll start by asking 4 people
that love Soda if they also
loves popcorn.

Loves Popcorn		Age		Loves Cool As Ice
Loves Popcorn	Age	Loves Soda	Age	Loves Cool As Ice
Yes	7	Yes	12.5	No
No	12	No	No	No

$$GI = 0.5$$

Loves Popcorn		Age		Loves Cool As Ice
Loves Popcorn	Age	Loves Soda	Age	Loves Cool As Ice
Yes	18	Yes	26.5	Yes
No	35	Yes	36.5	Yes

$$GI = 0$$

$$\frac{\text{Total } GI}{(\text{Popcorn})} = \frac{2}{4}(0.5)$$

$$= 0.25$$

Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
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

→ **Age <**

Consider age of
people loving Soda.

Repeating the process
of GI of averages

Age < 12.5



Age < 12.5		Age		Loves Cool As Ice
Age < 12.5	Age	Loves Soda	Age	Loves Cool As Ice
Yes	0	No	1	
No	1			

$$GI = 0$$

$$GI = 0$$

$$\text{Total } GI = 0$$

Age < 26.5



Age < 26.5		Age		Loves Cool As Ice
Age < 26.5	Age	Loves Soda	Age	Loves Cool As Ice
Yes	1	No	1	
No	2		0	

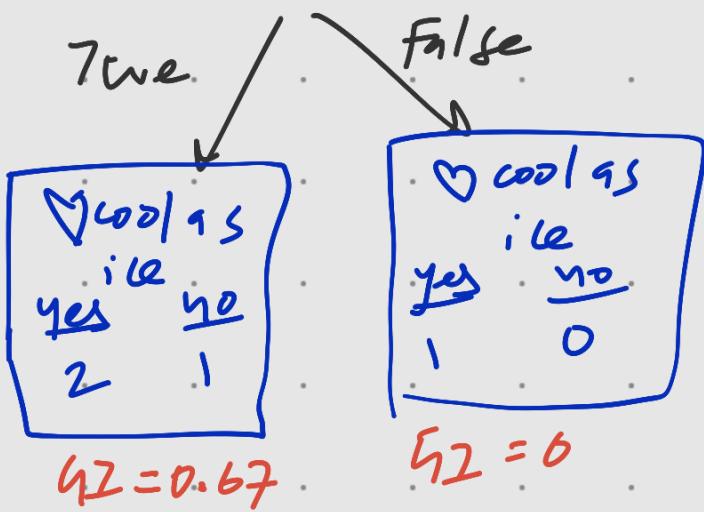
$$GI = 0.5$$

$$GI = 0$$

$$\text{Total } GI = \frac{1}{2}(0.5) = 0.25$$

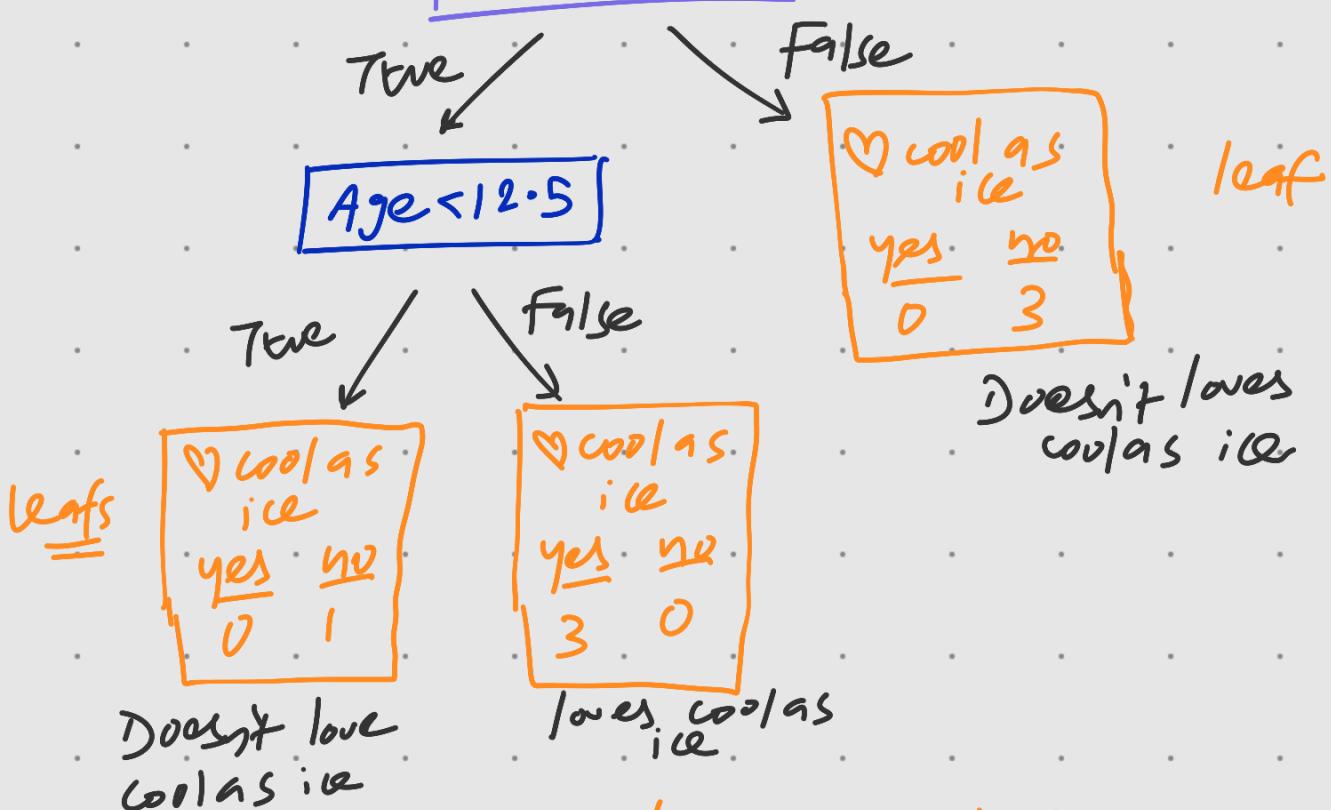
The least
among ages
& less than Popcorn too. i.e. the branch

Age < 36.5



$$\text{Total GI} = \frac{3}{4} (0.67) = 0.5$$

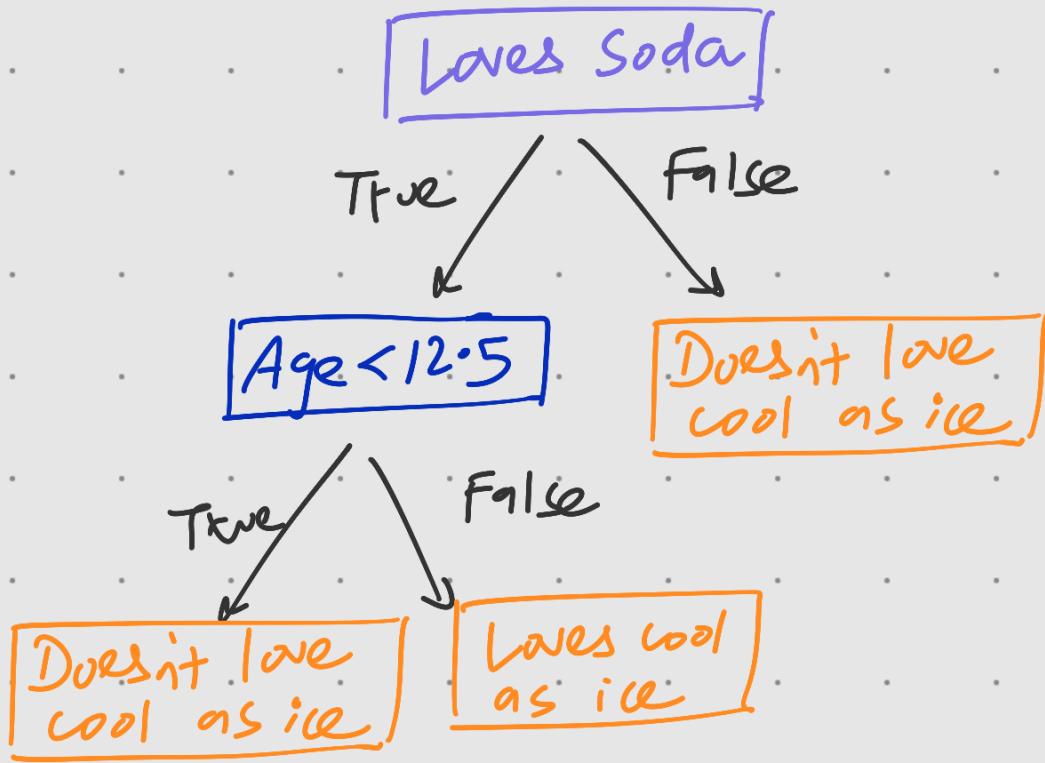
Loves Soda



we need to assign output

values for each leaf. i.e category

Final Tree



→ overfitting, because only 1 person reaches this node | Prune

Check for no. of values per node
by cross validation.