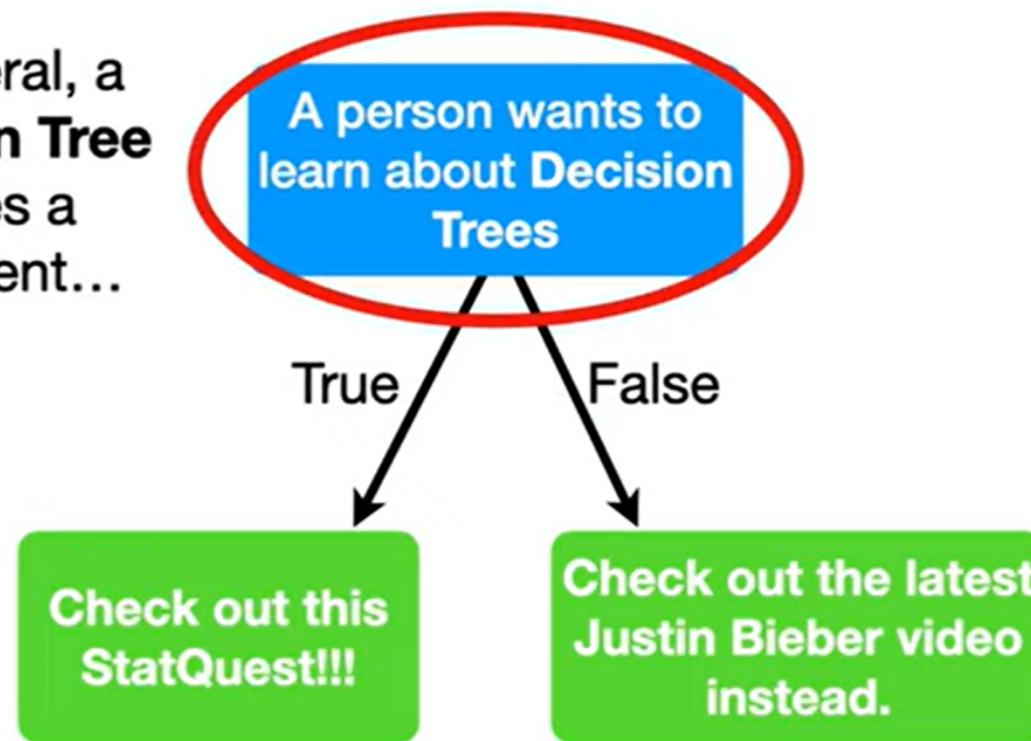




In general, a **Decision Tree** makes a statement...





A person wants to learn about **Decision Trees**

True

False

Check out this StatQuest!!!

Check out the latest Justin Bieber video instead.

When a **Decision Tree** classifies things into categories...

...it's called a **Classification Tree**.



A person wants to learn about **Decision Trees**

True False

Check out this StatQuest!!!

Check out the latest Justin Bieber video instead.

A mouse eats a special diet.

True False

It is between 150 and 180mm long

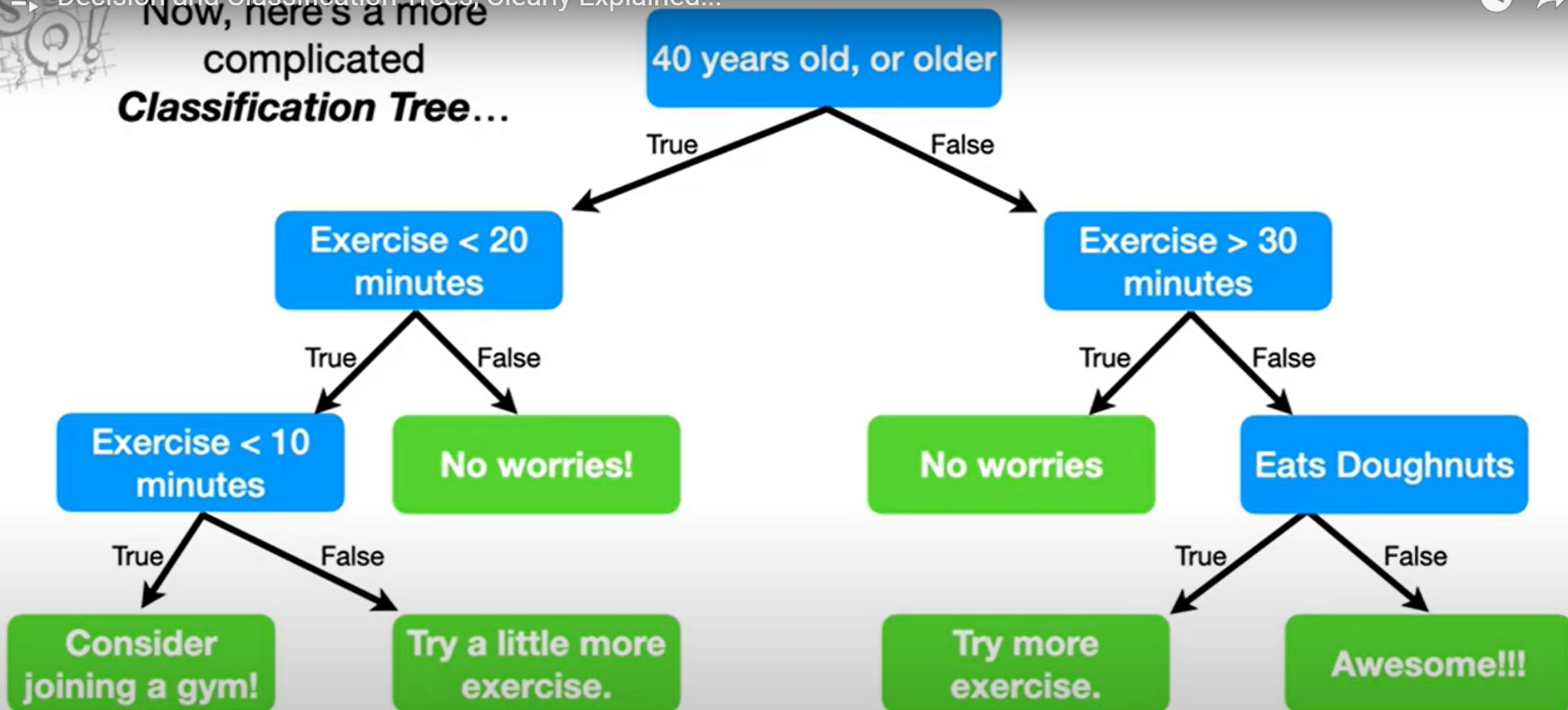
It is less than 150mm long

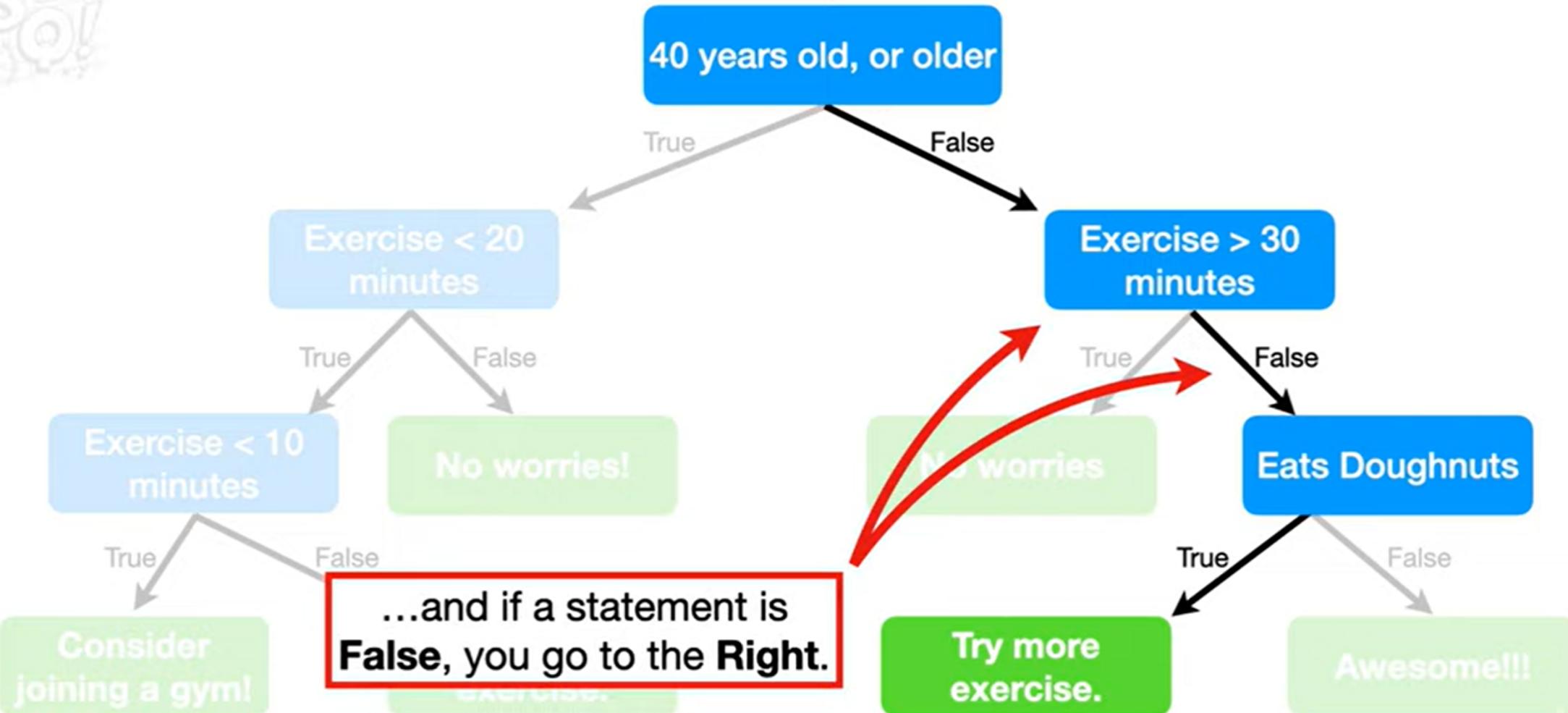
When a **Decision Tree** classifies things into categories...

...it's called a **Classification Tree**.

And when a **Decision Tree** predicts numeric values...

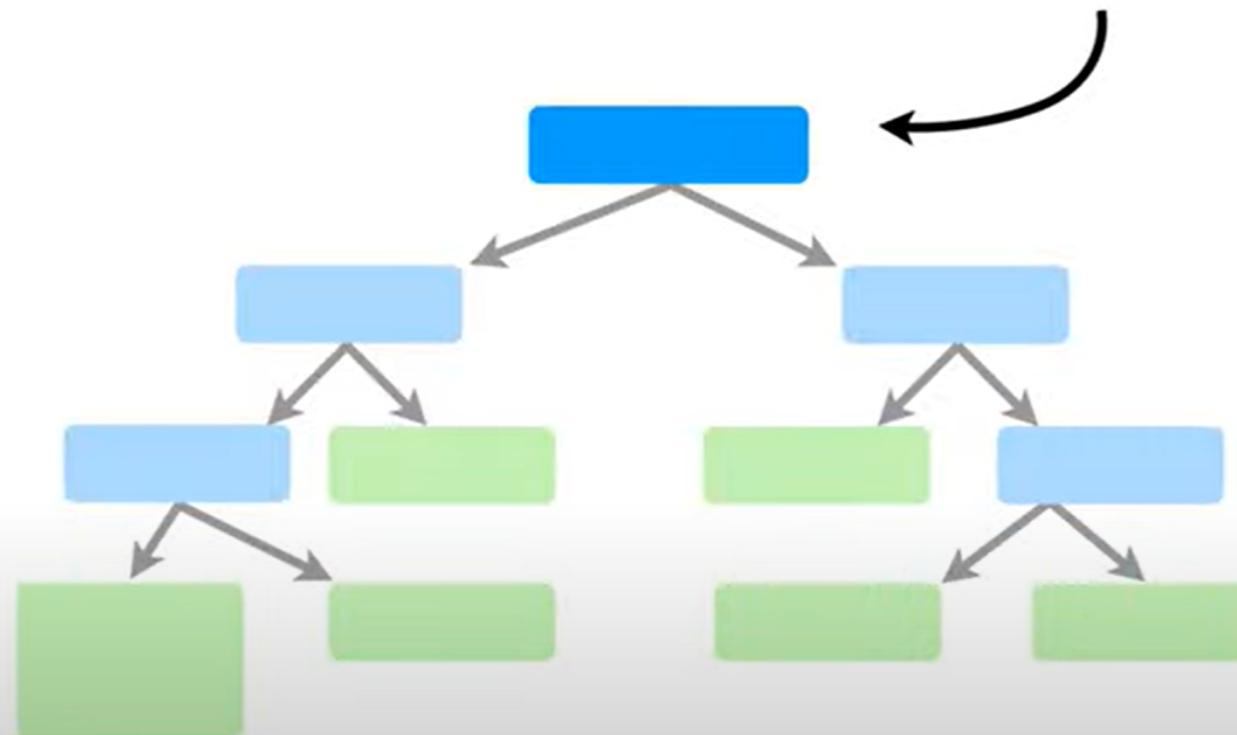
Now, here's a more complicated **Classification Tree...**





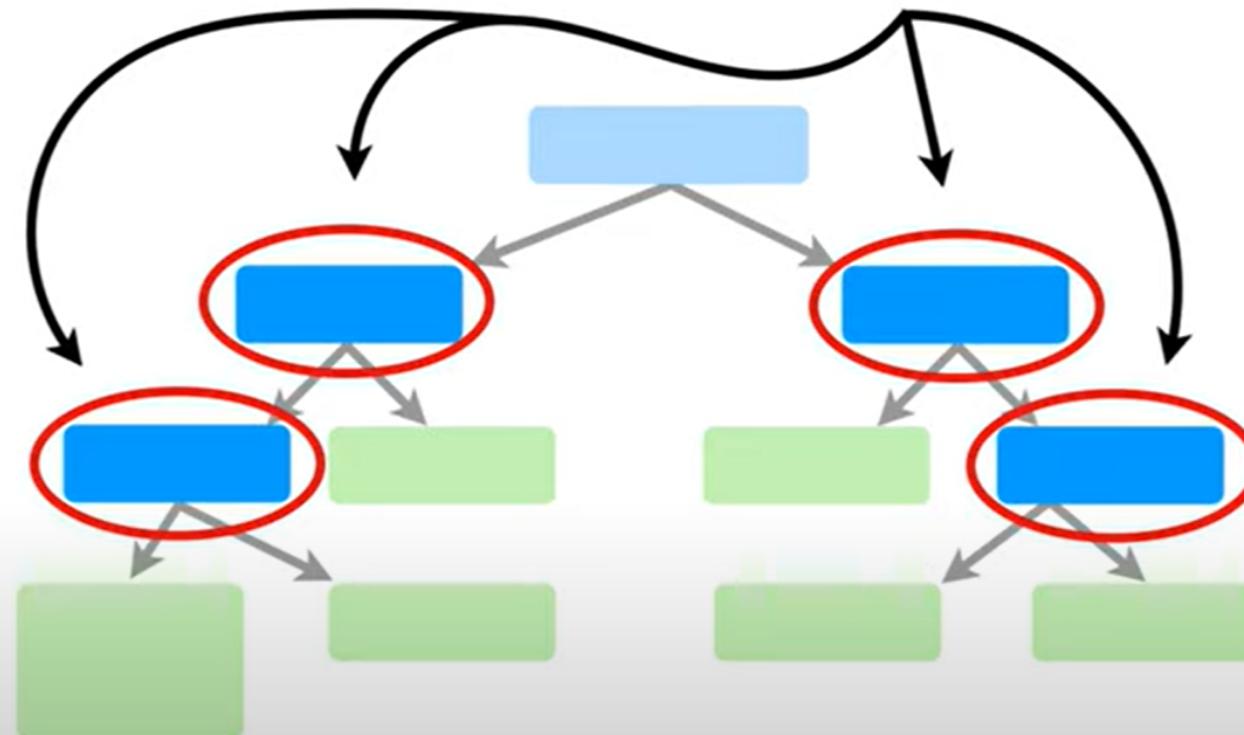


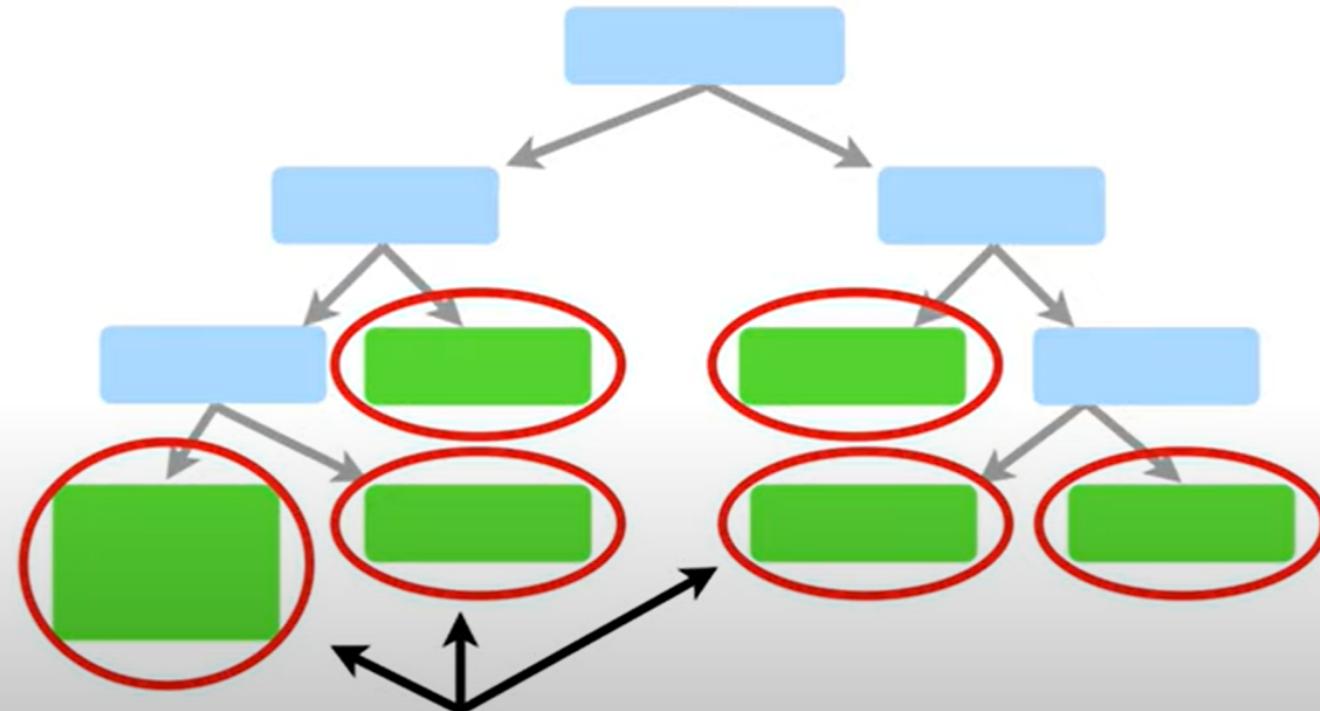
The very top of the tree is called the **Root Node** or just **The Root**.



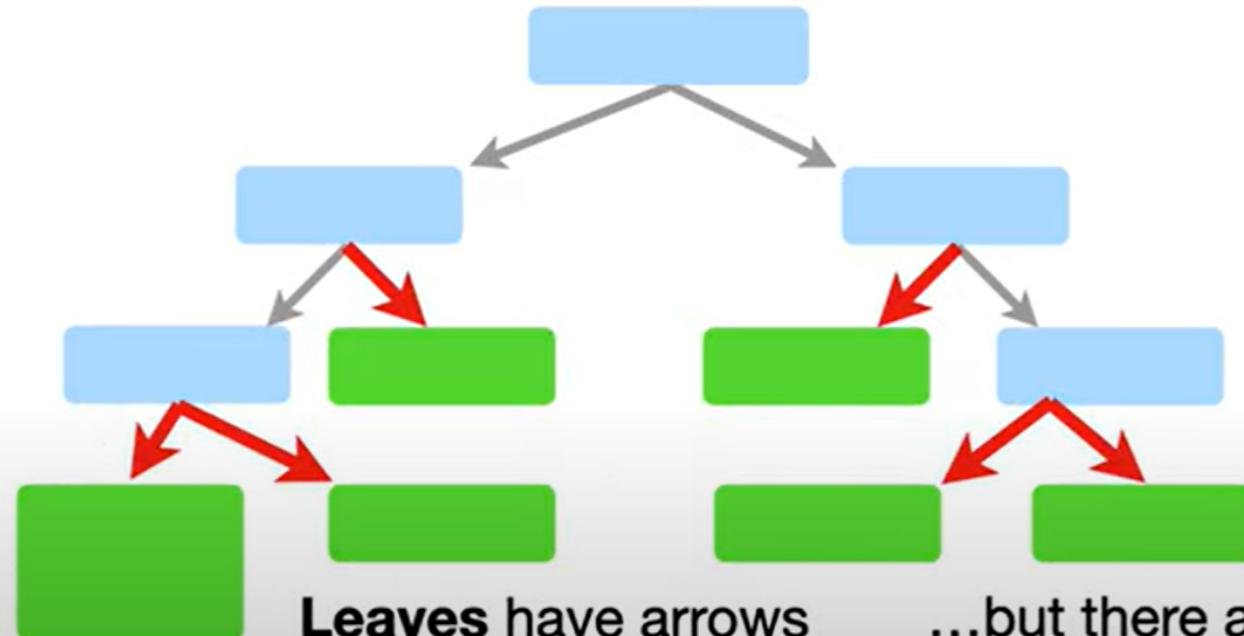


These are called **Internal Nodes**, or **Branches**.





Lastly, these are called **Leaf Nodes**, or just **Leaves**.



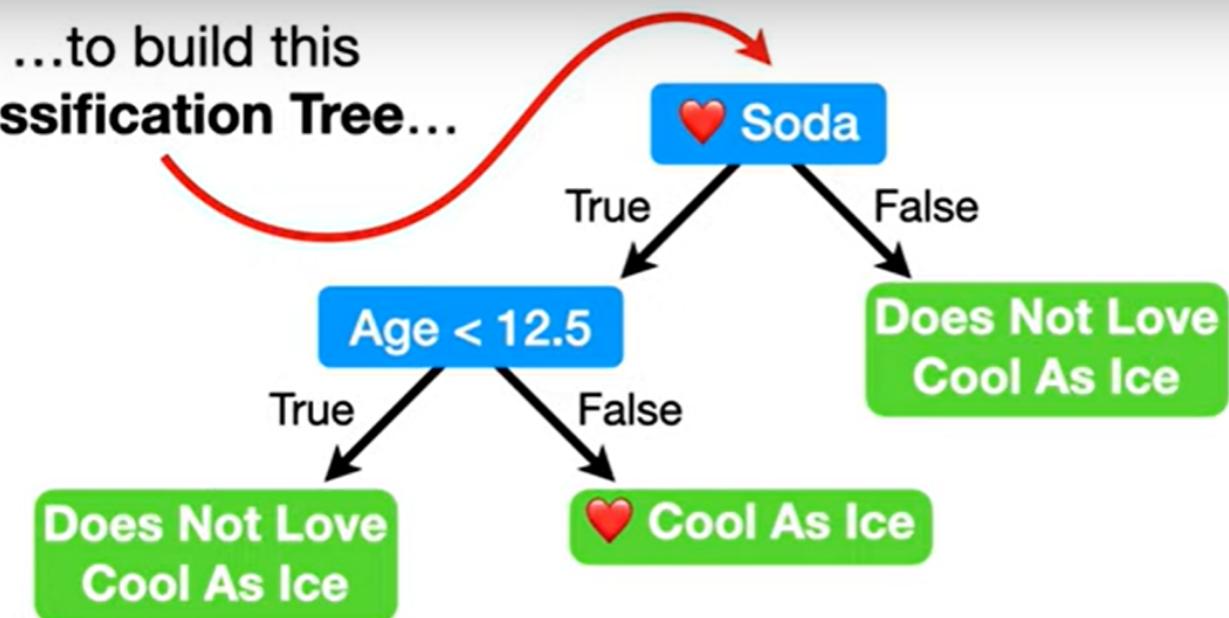
Leaves have arrows
pointing to them...

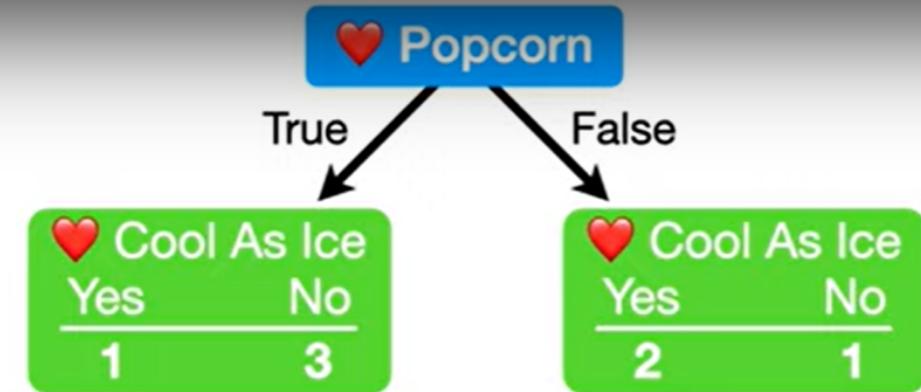
...but there are no
arrows pointing
away from them.



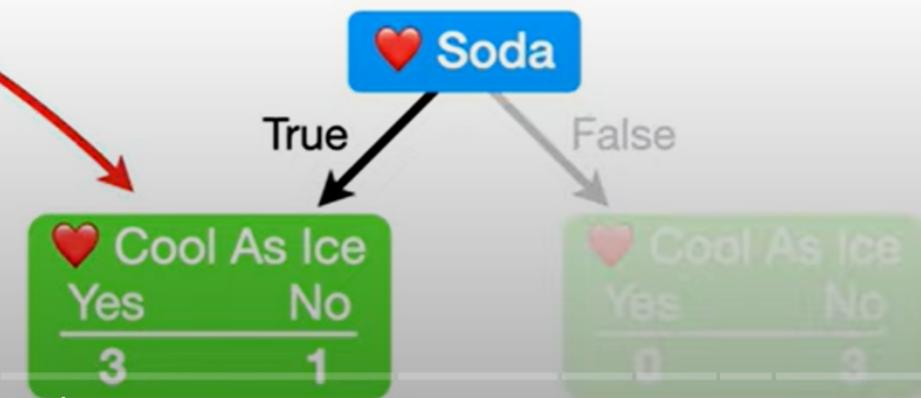
Loves Popcorn	Loves Soda	Age	Loves Cool As Ice
Yes	Yes	7	No
Yes	No	12	No
No	Yes	18	Yes
No	Yes	35	Yes
Yes	Yes	38	Yes
Yes	No	50	No
No	No	83	No

...to build this
Classification Tree...





Specifically, these three **Leaves** contain mixtures of people that *do and do not Love Cool As Ice.*





Popcorn

True

False

Cool As Ice	
Yes	No
1	3

Cool As Ice	
Yes	No
2	1

In contrast, this **Leaf** only contains people who *do not* Love Cool As Ice.

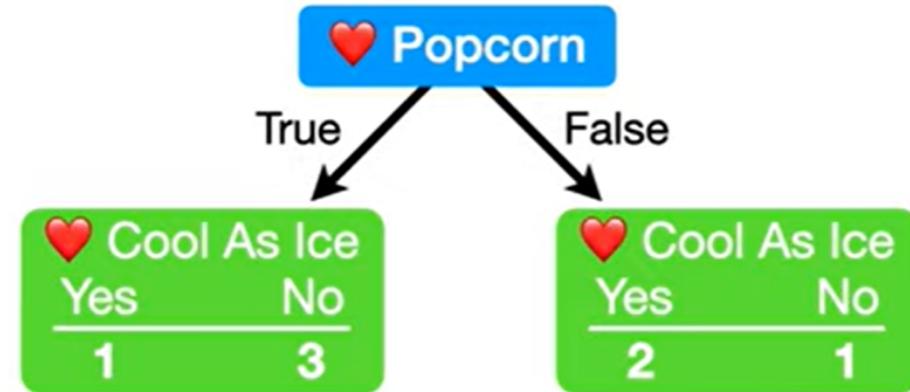
Soda

True

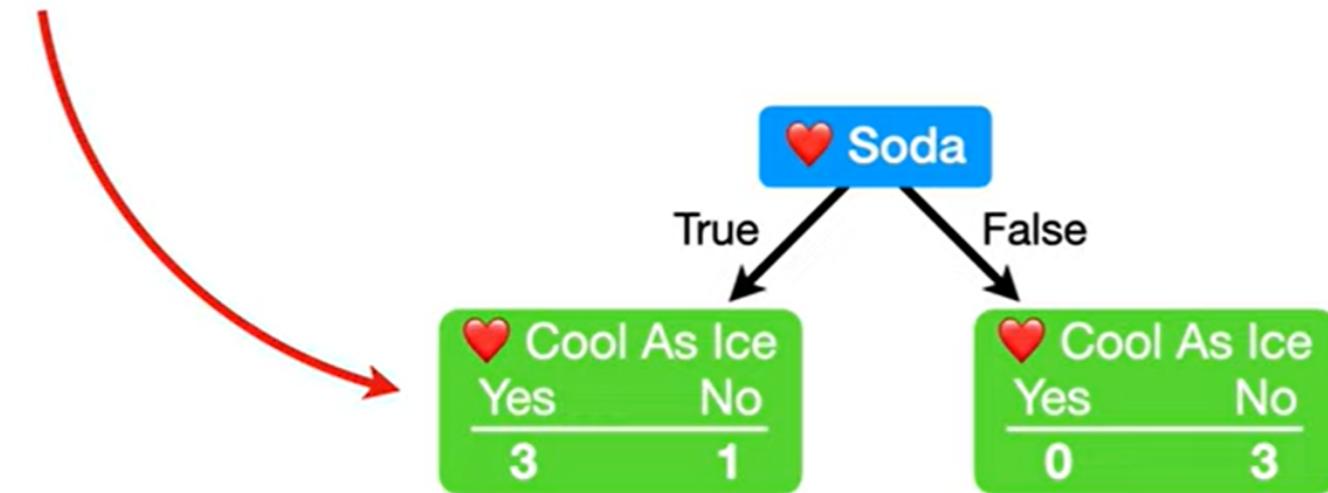
False

Cool As Ice	
Yes	No
0	3

Cool As Ice	
Yes	No
0	3

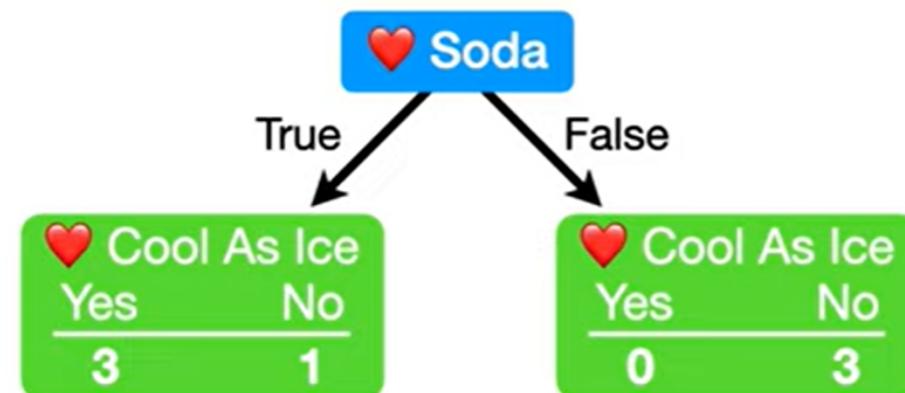
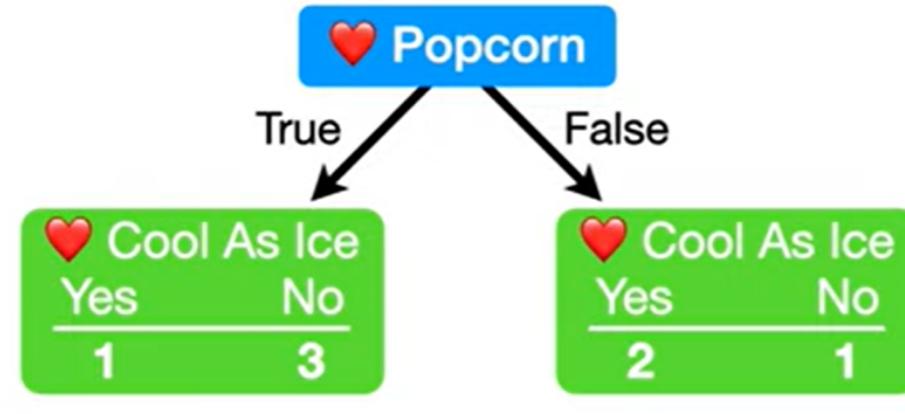


...and only one **Leaf** in the
Loves Soda tree is **Impure**...



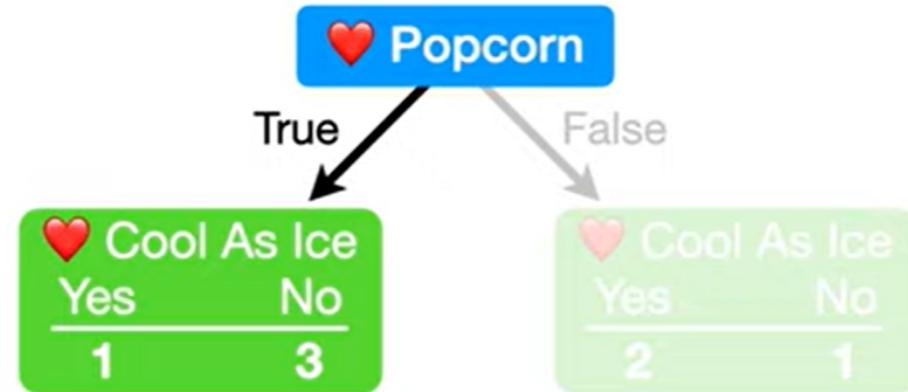


One of the most popular methods is called **Gini Impurity**, but there are also fancy sounding methods like **Entropy** and **Information Gain**.

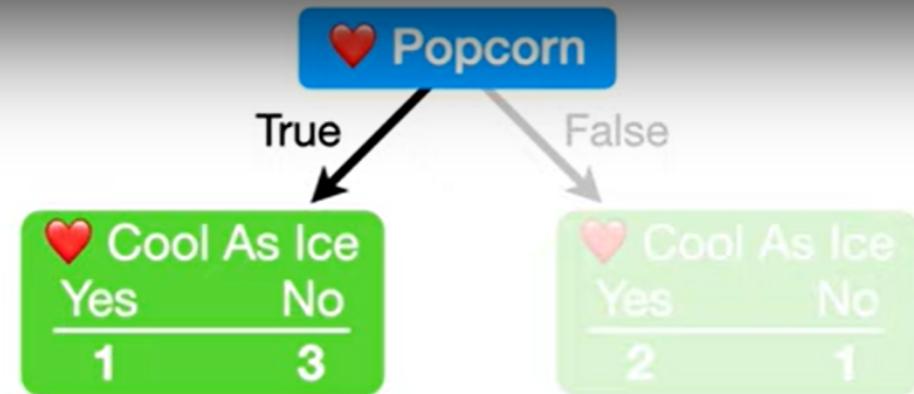




The **Gini Impurity** for the Leaf on the left is...



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



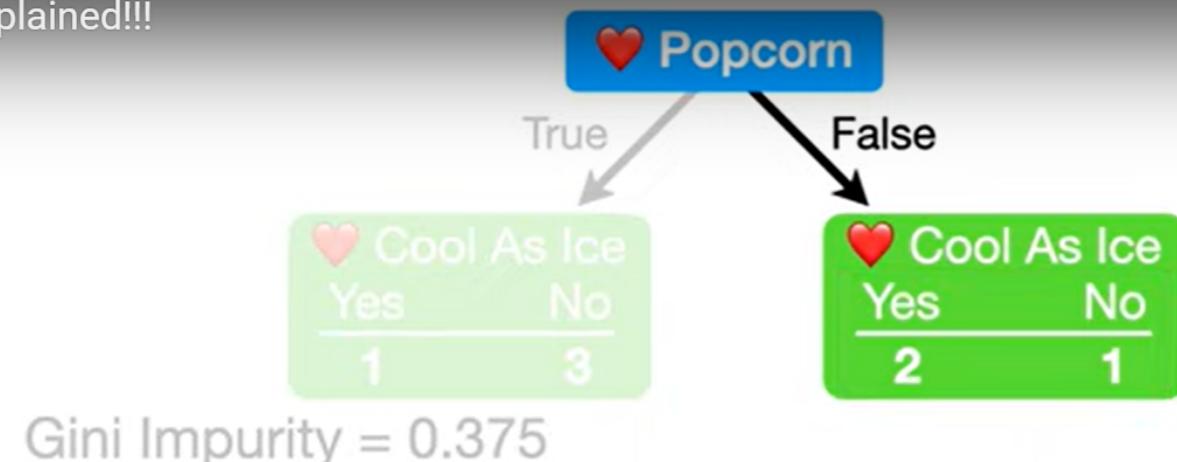
Gini Impurity = 0.375

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

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

$$= 0.375$$

So let's put **0.375** under the **Leaf** on the left so we don't forget it.

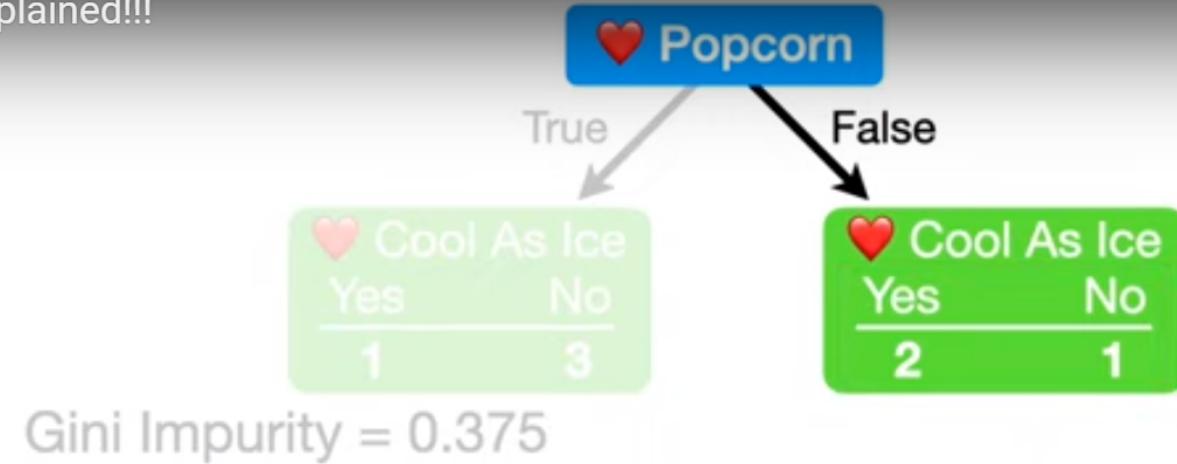


Gini Impurity for a Leaf = $1 - (\text{the probability of "Yes"})^2 - (\text{the probability of "No"})^2$

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

$$= 0.444$$

And when we do the math we get **0.444**.



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

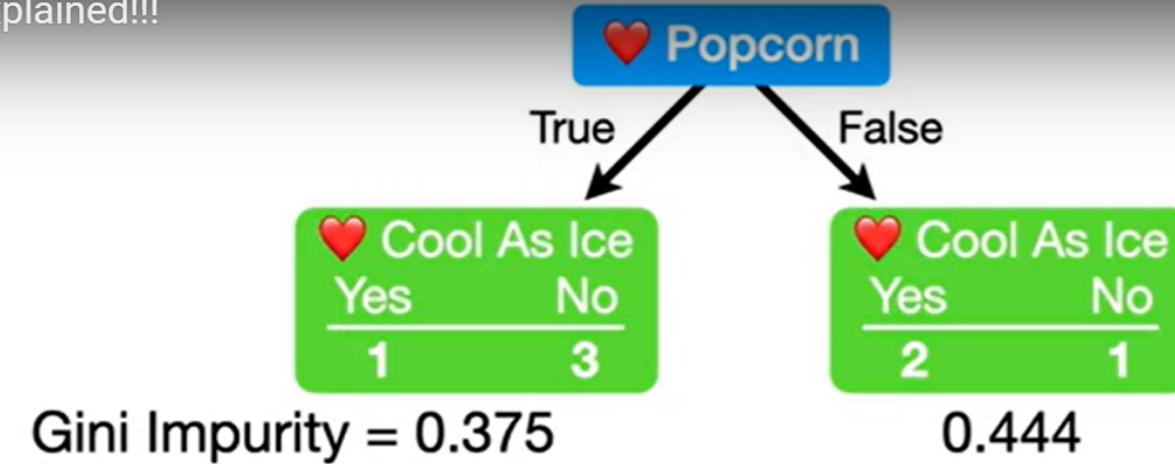
$$= 1 - \left(\frac{2}{2+1}\right)^2 - \left(\frac{1}{2+1}\right)^2$$

$$= 0.444$$

And when we do the
math we get **0.444**.



And when we do the math, we get **0.405**.



Total **Gini Impurity** = weighted average of **Gini Impurities** for the **Leaves**

$$= \left(\frac{4}{4+3} \right) 0.375 + \left(\frac{3}{4+3} \right) 0.444$$

$$= 0.405$$



Lastly, we calculate the **Gini Impurity** values for each average age.

Age	Loves Cool As Ice
9.5	No
12	No
15	
18	Yes
26.5	
35	Yes
36.5	
38	Yes
44	
50	No
66.5	
83	No

Gini Impurity = 0.429

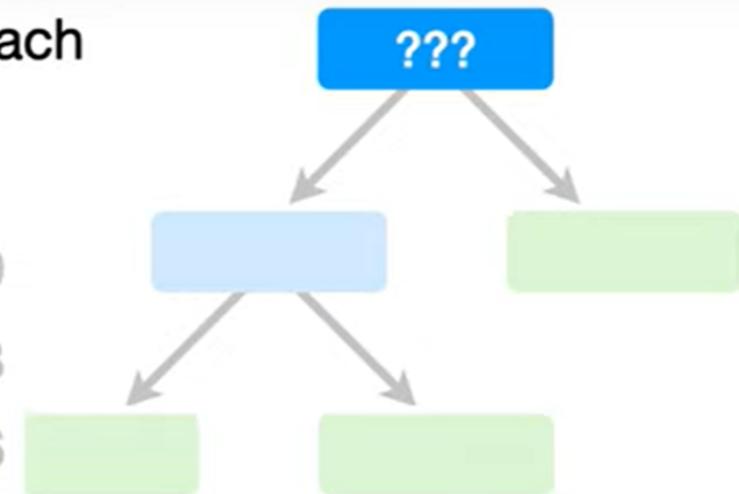
Gini Impurity = 0.343

Gini Impurity = 0.476

Gini Impurity = 0.476

Gini Impurity = 0.343

Gini Impurity = 0.429





Lastly, we calculate the **Gini Impurity** values for each average age.

Age	Loves Cool As Ice
9.5	No
12	No
15	
18	Yes
26.5	
35	Yes
36.5	
38	Yes
44	
50	No
66.5	
83	No

Gini Impurity = 0.429

Gini Impurity = 0.343

Gini Impurity = 0.476

Gini Impurity = 0.476

Gini Impurity = 0.343

Gini Impurity = 0.429

