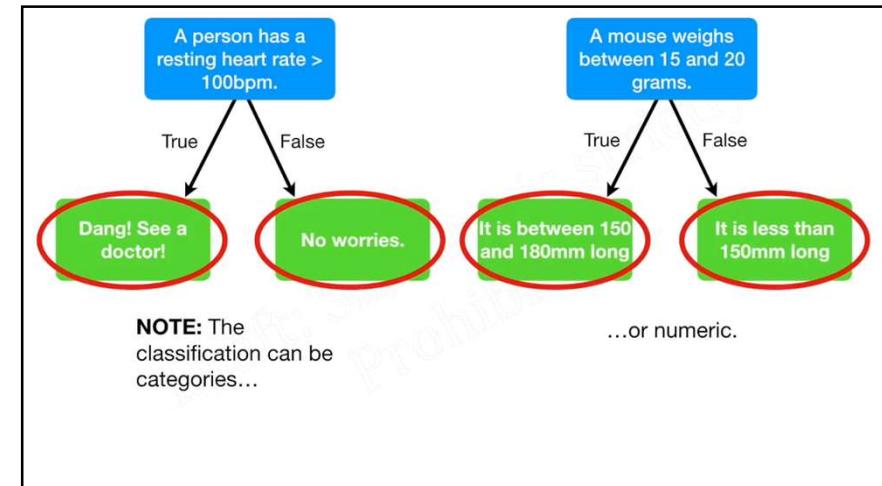
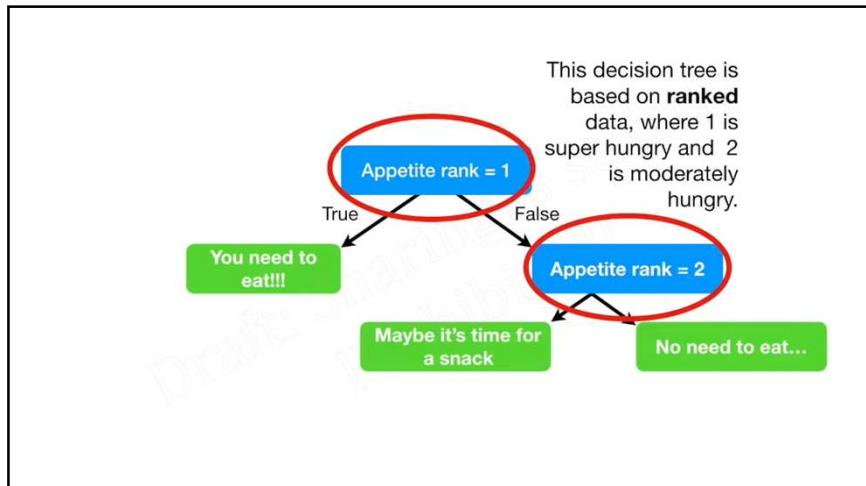
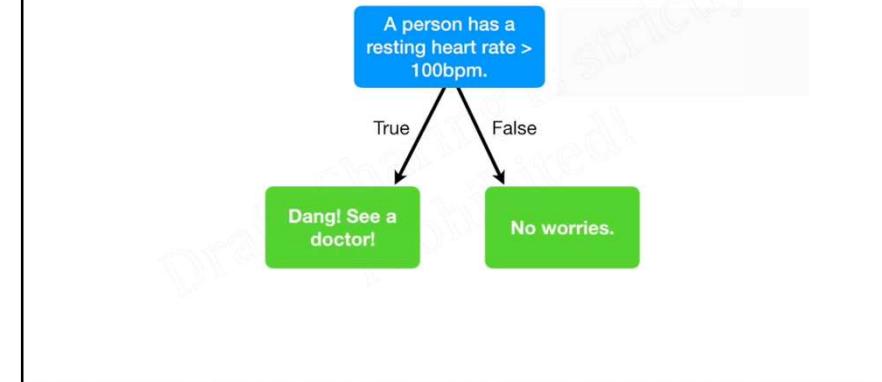
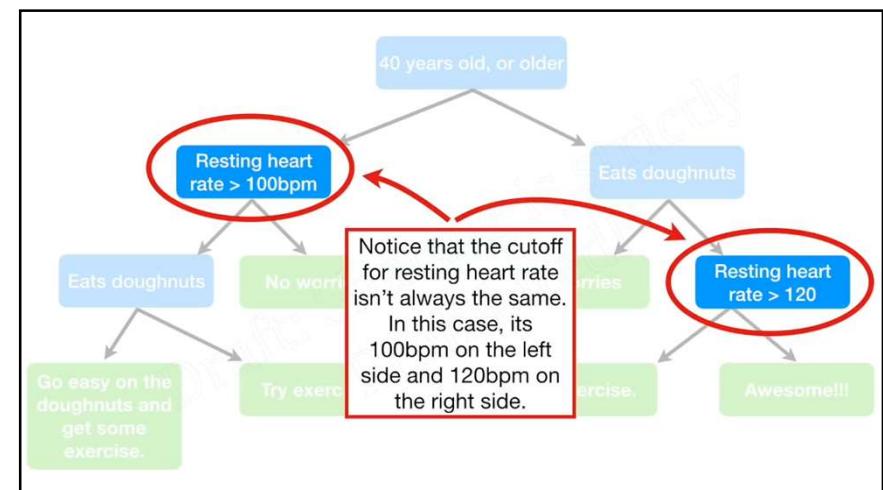
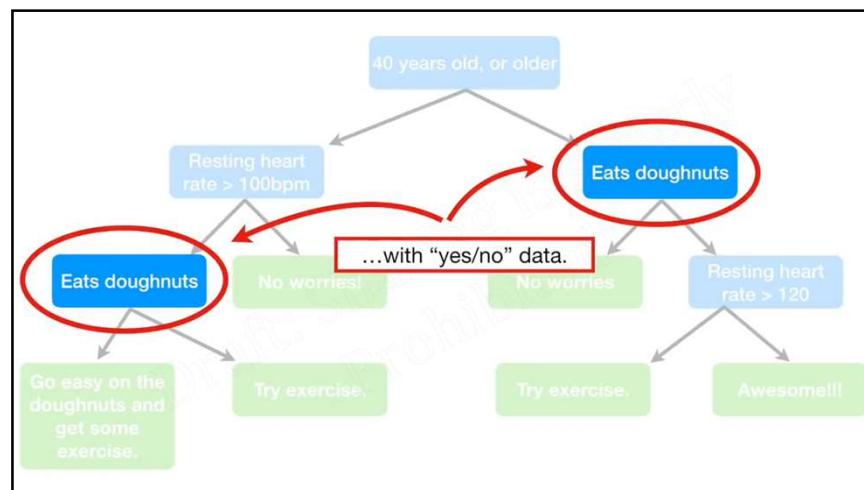
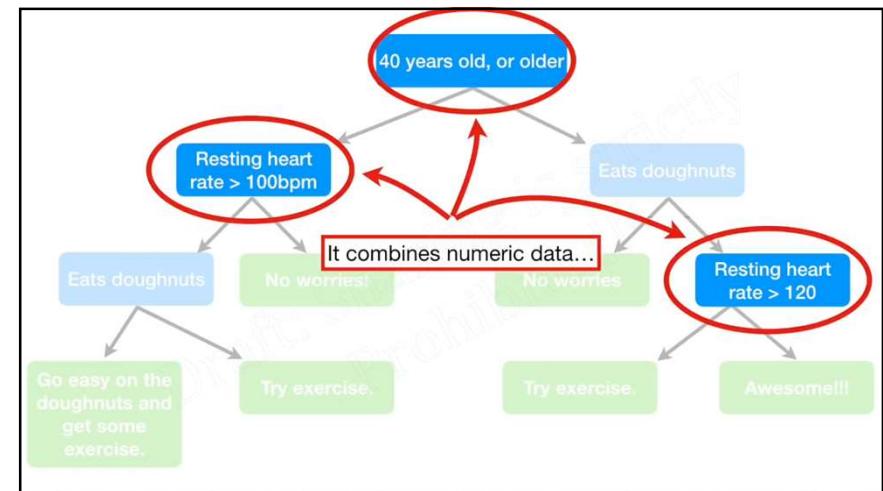
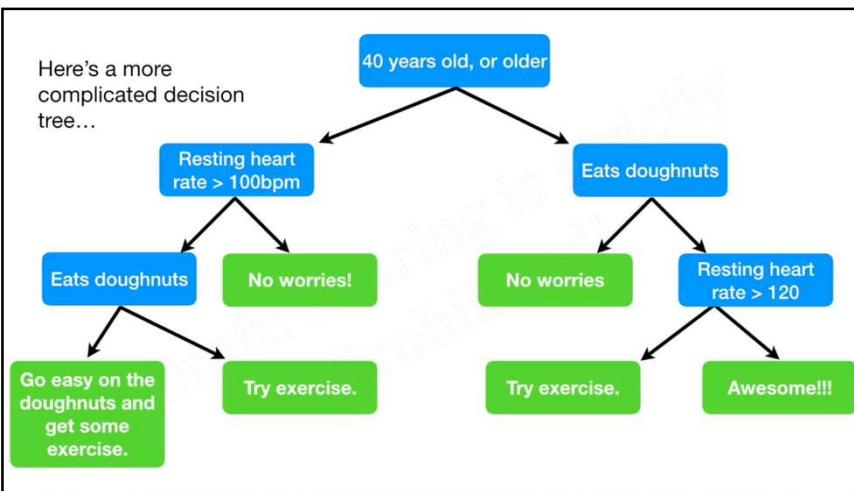
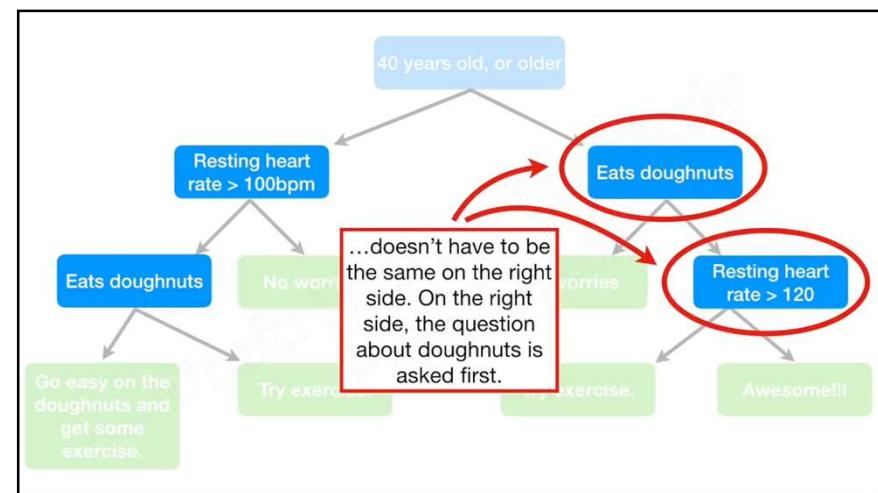
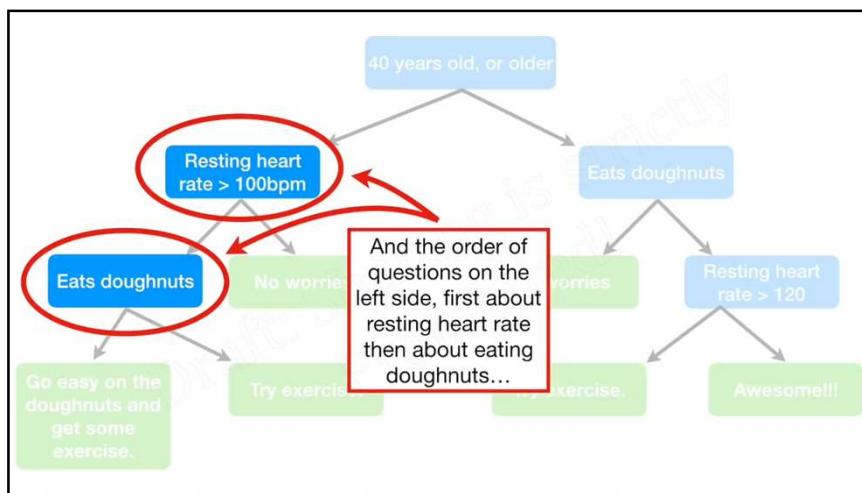


# DECISION TREES & RANDOM FORESTS



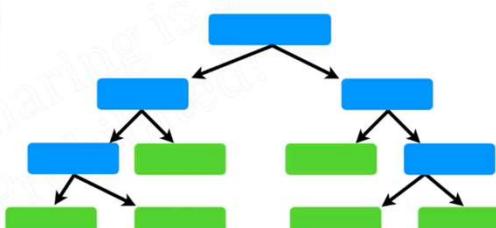




Now we are ready to talk about how to go from a raw table of data...

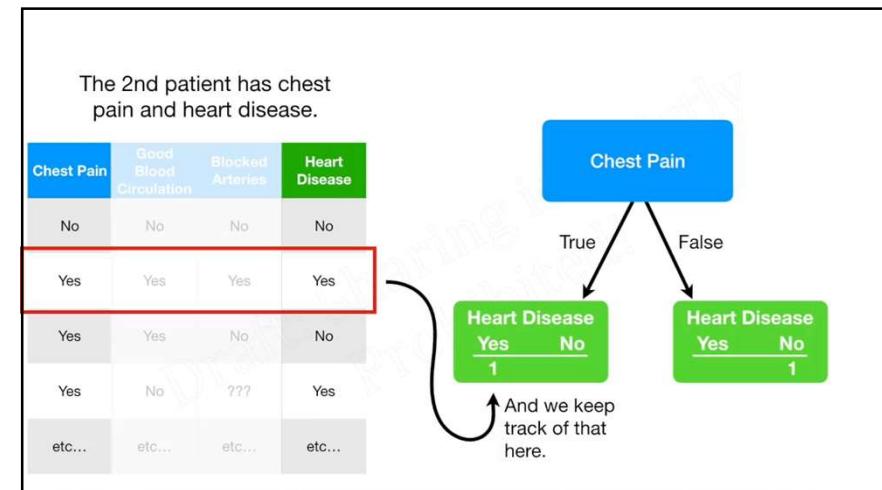
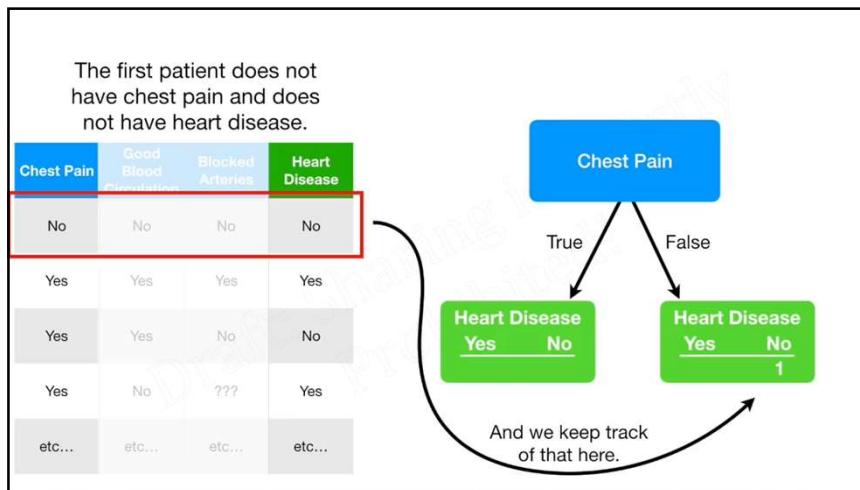
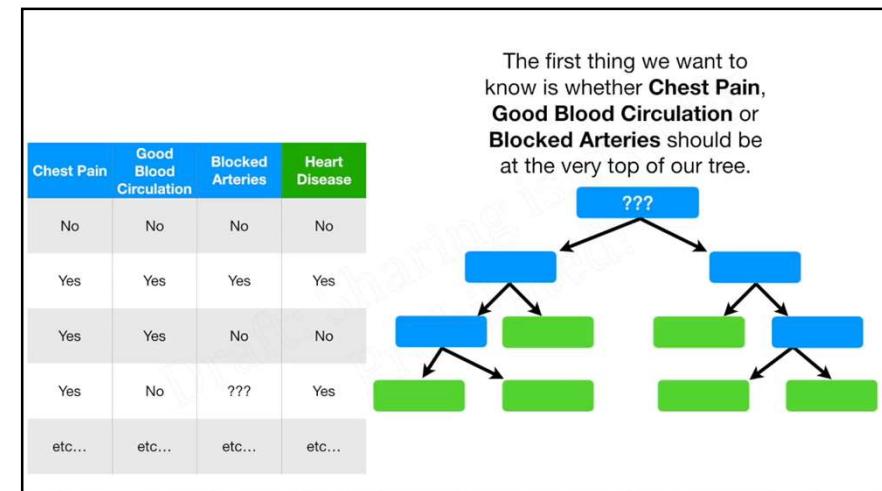
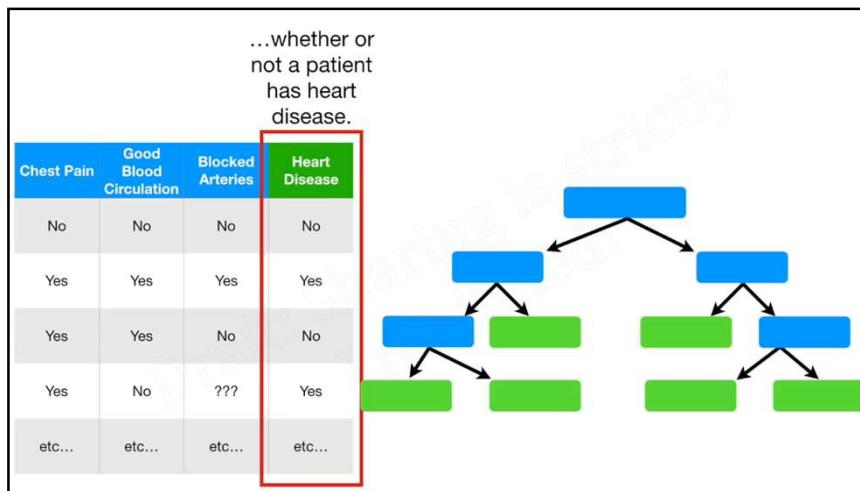
...to a decision tree!!!

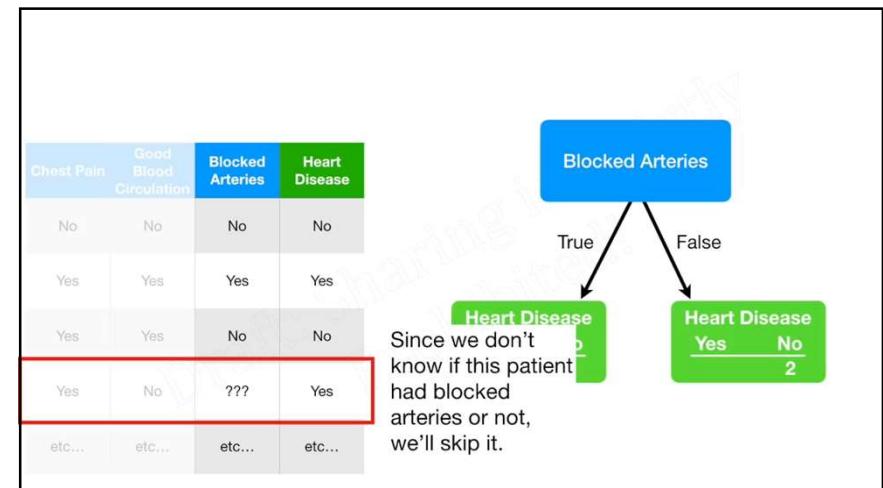
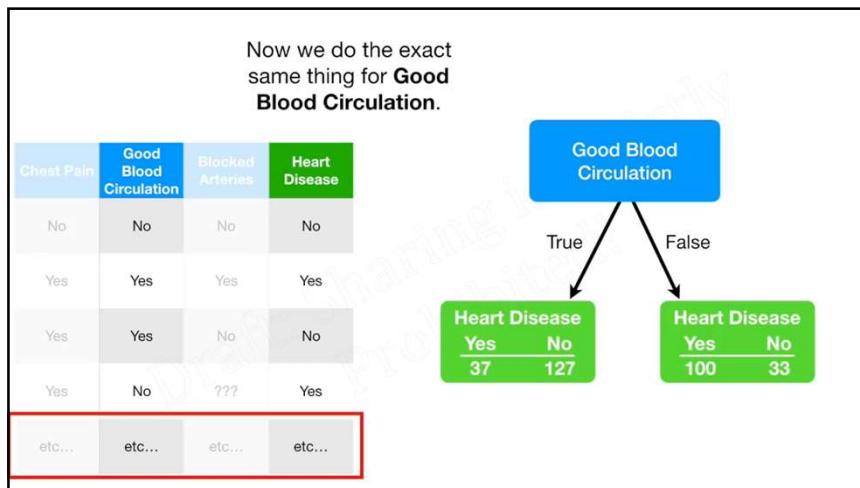
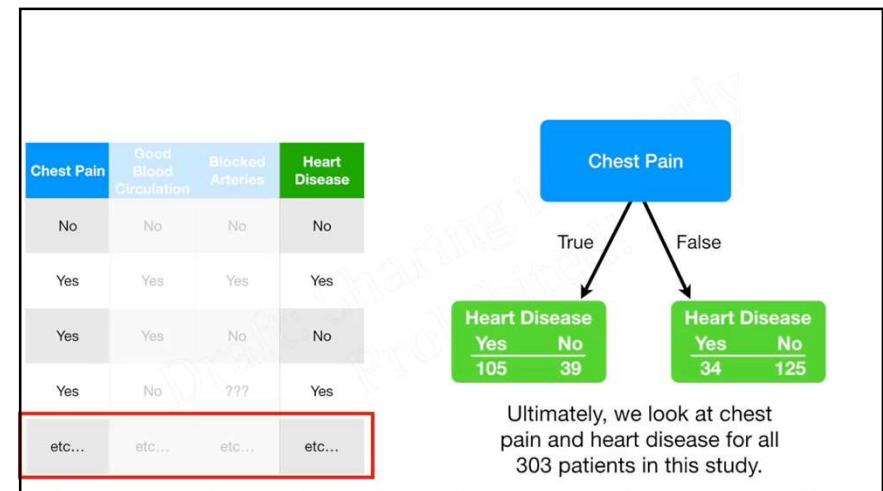
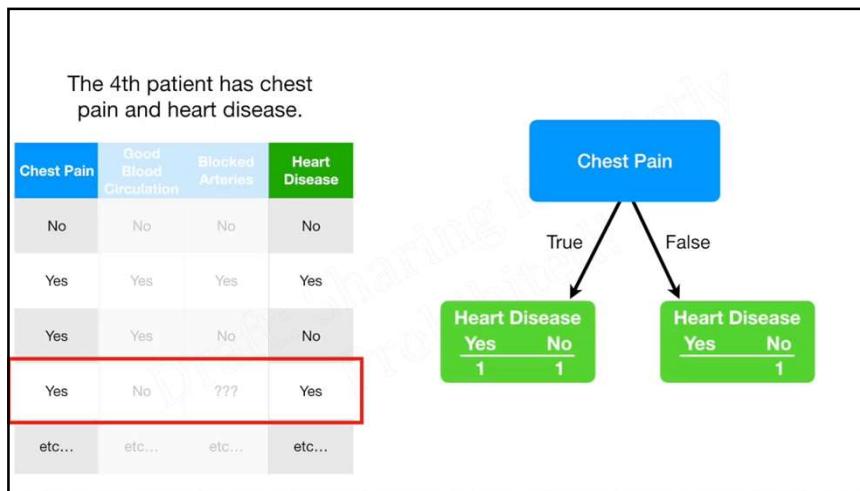
Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...

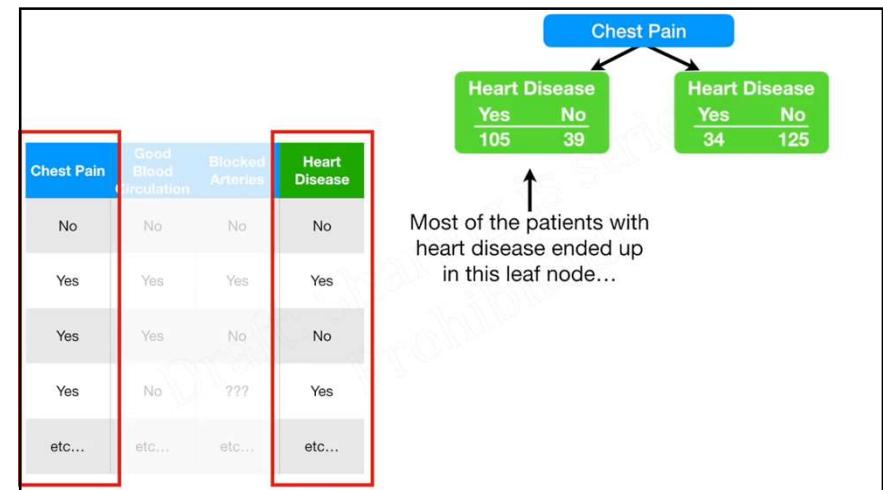
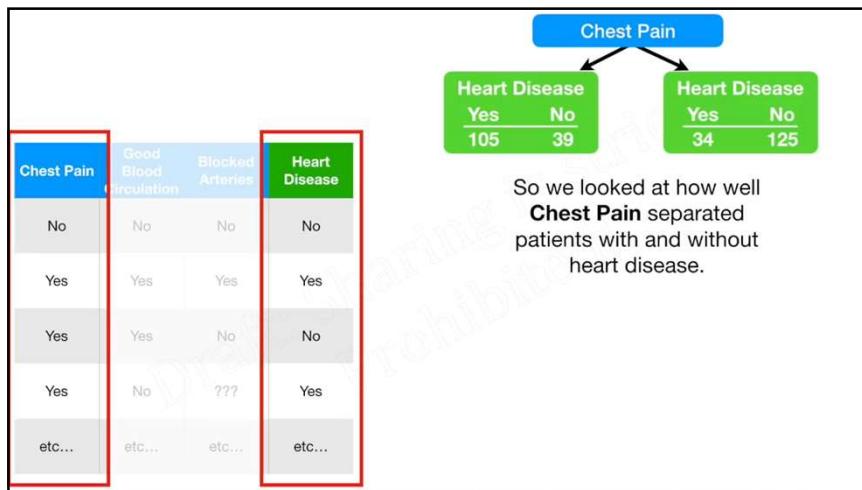
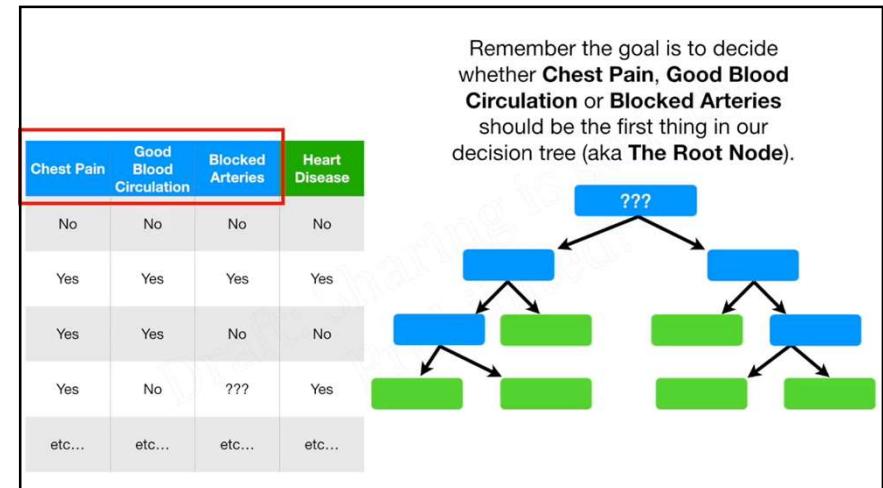
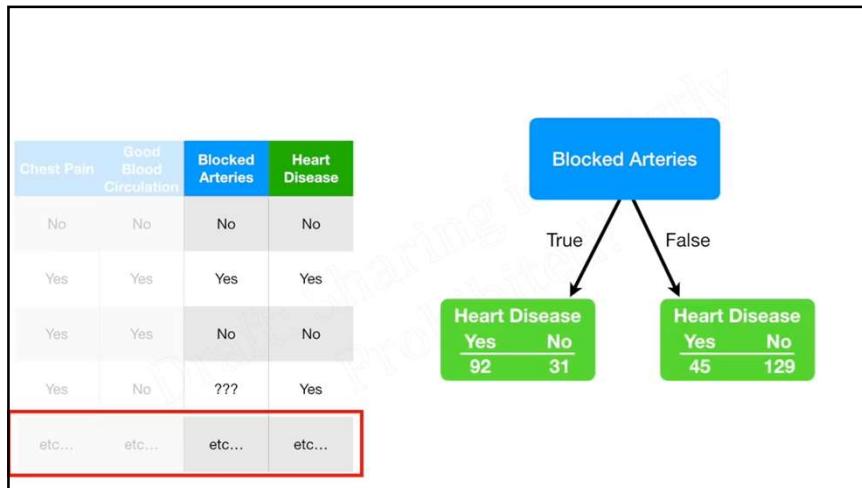


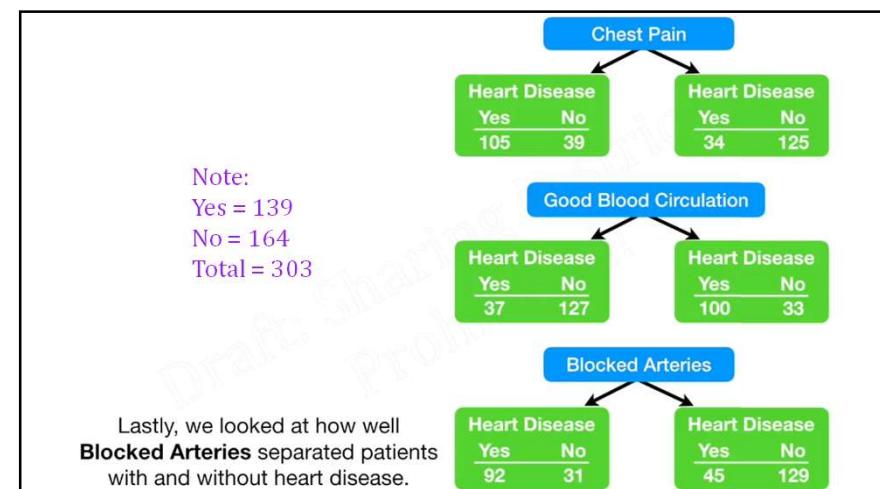
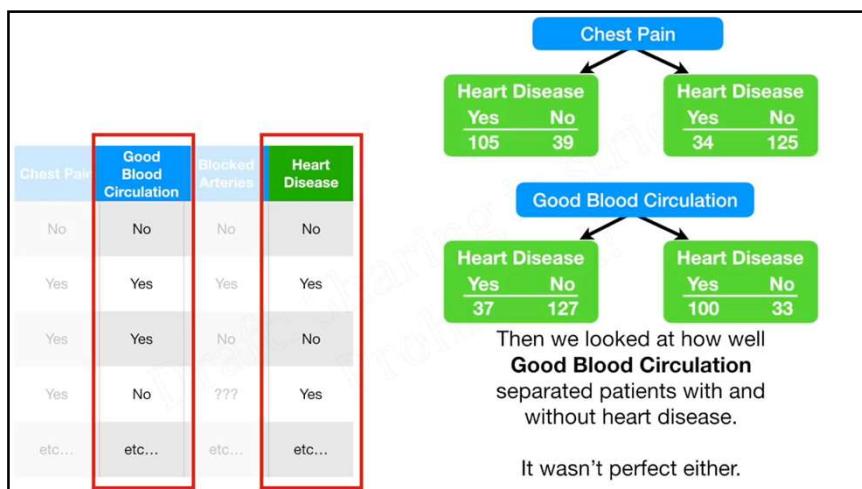
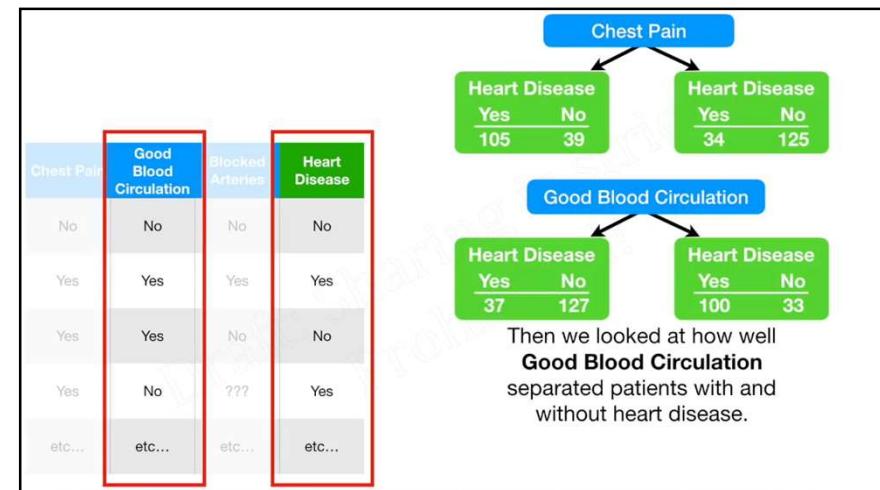
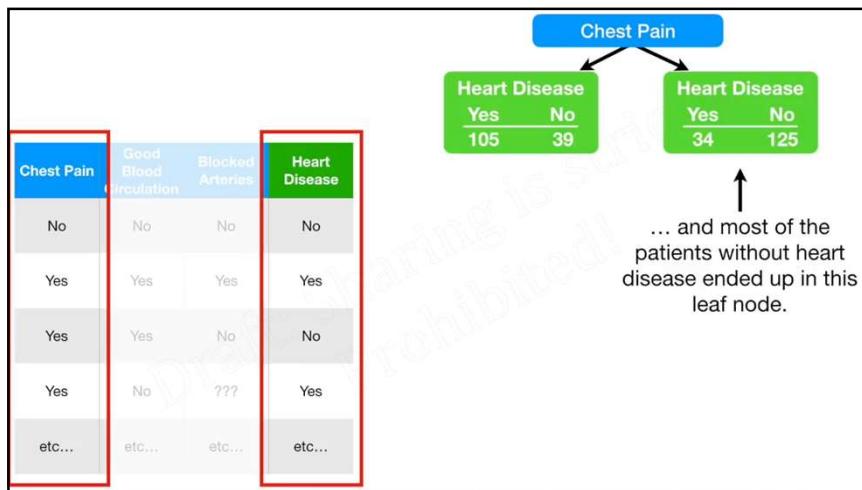
In this example, we want to create a tree that uses **chest pain**, **good blood circulation** and **blocked artery status** to predict...

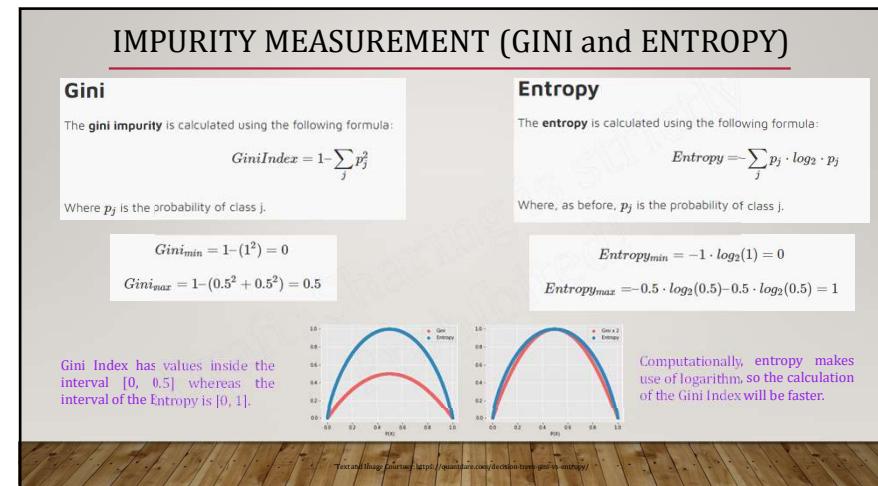
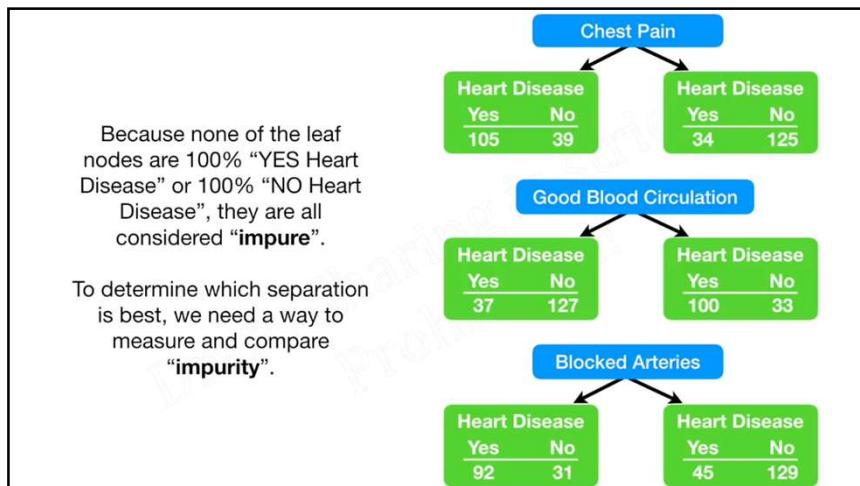
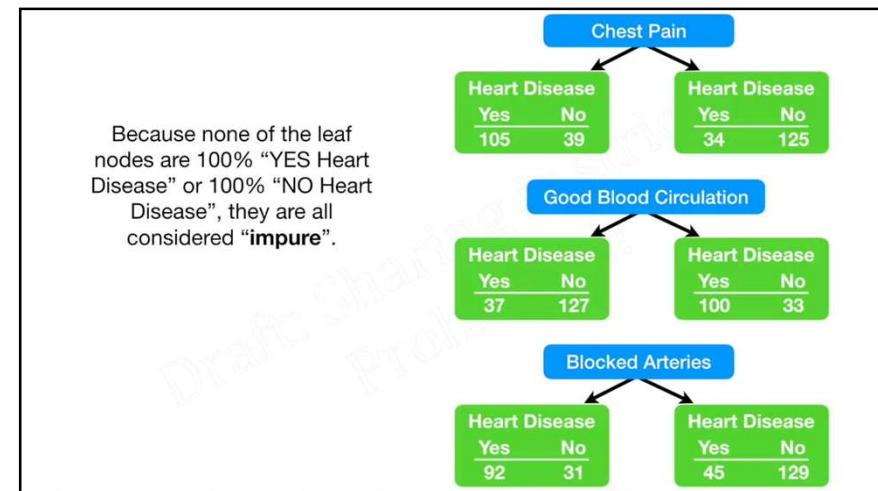
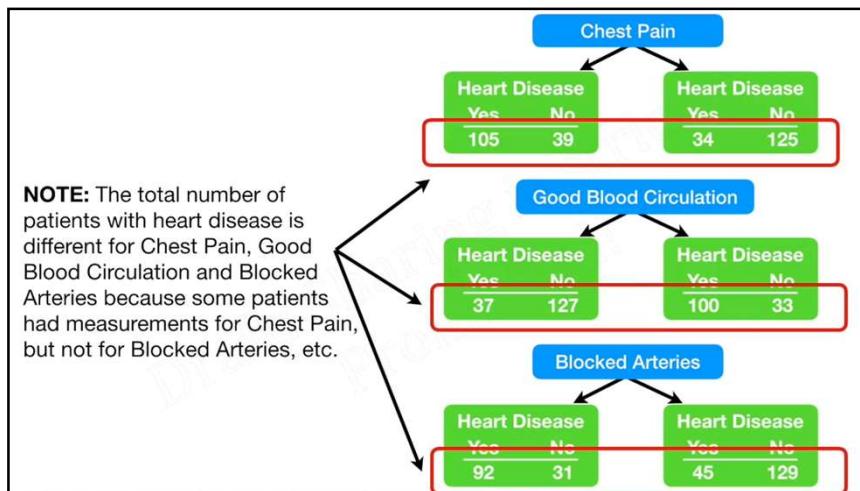
Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...



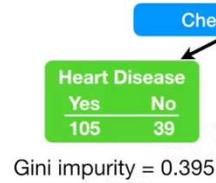
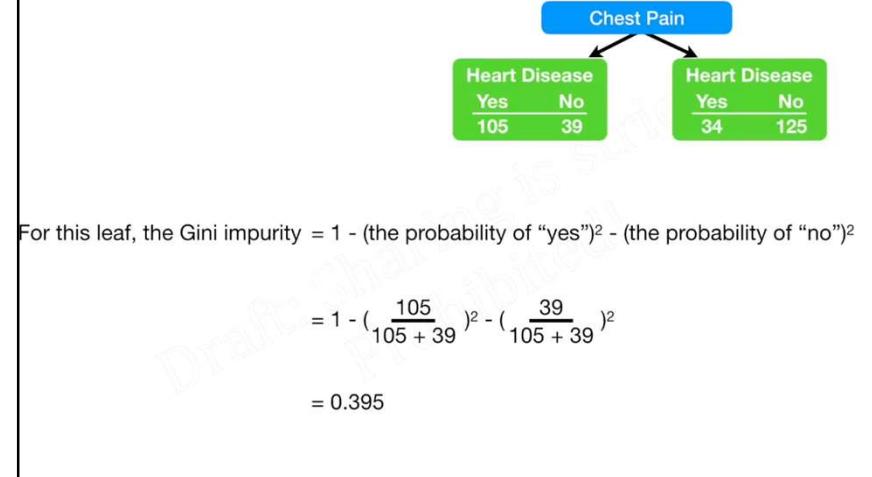
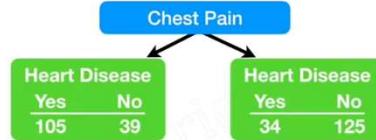




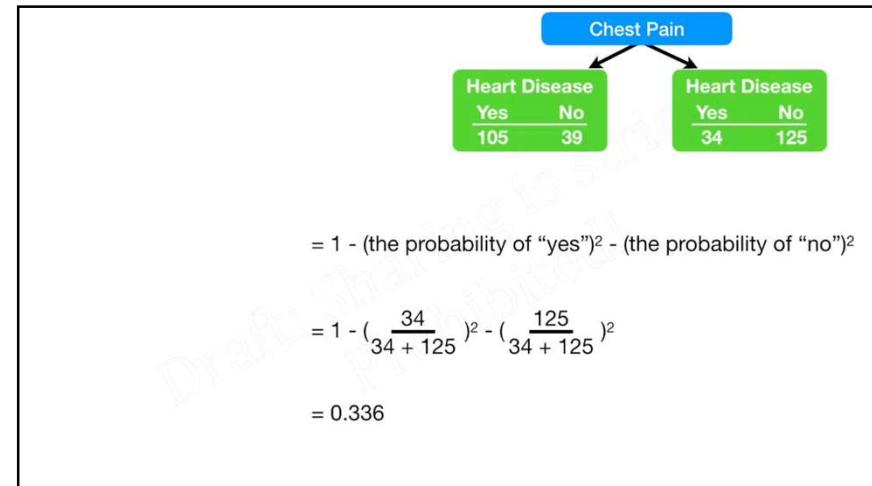


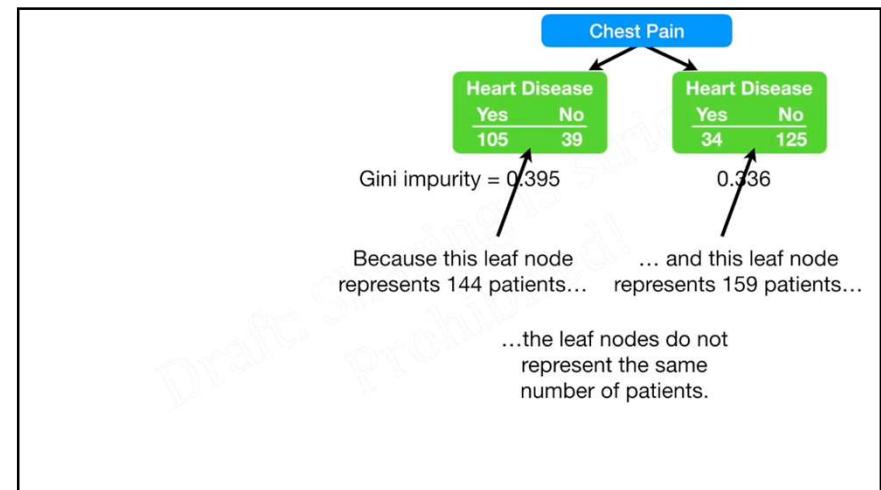
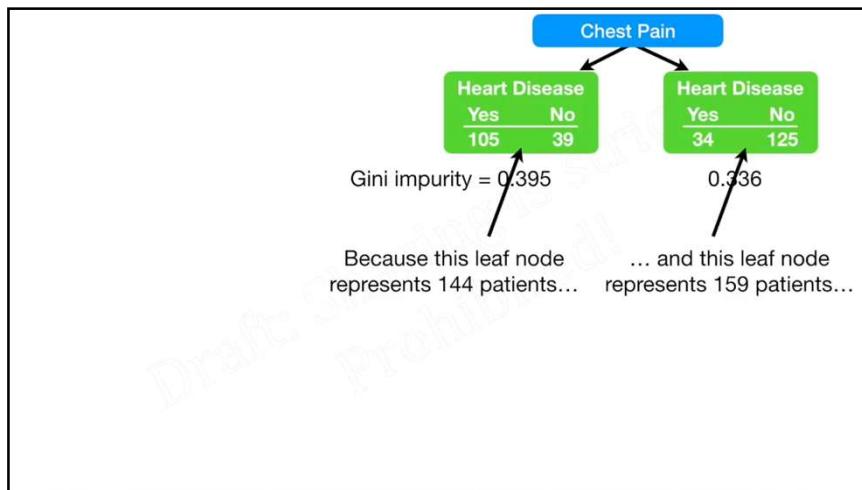
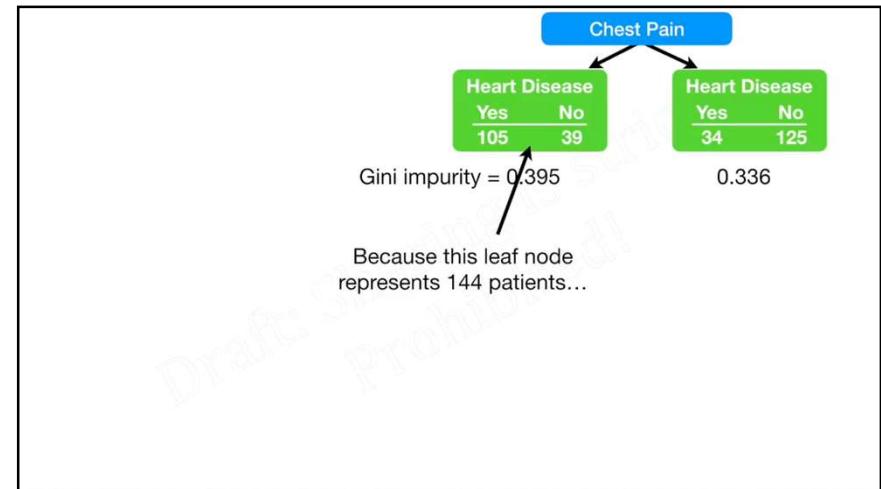
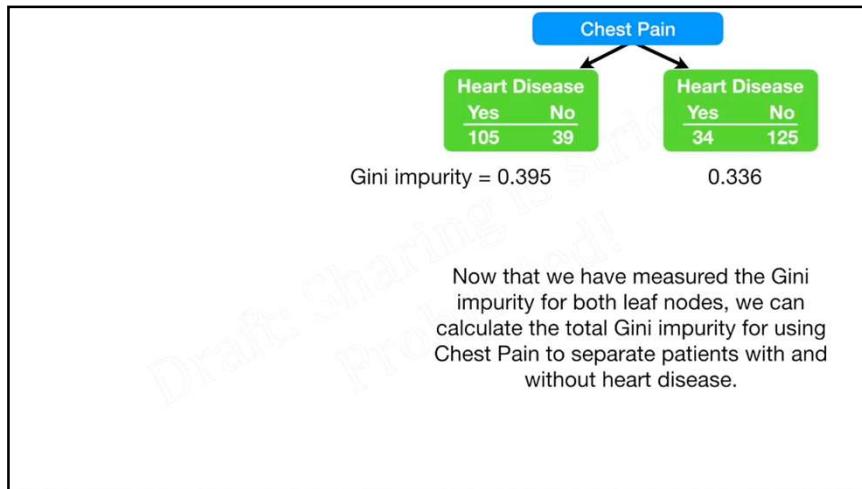


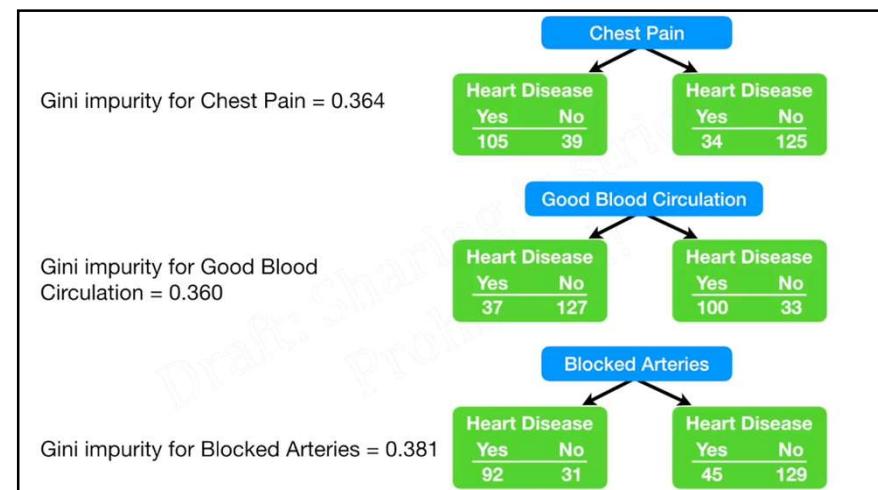
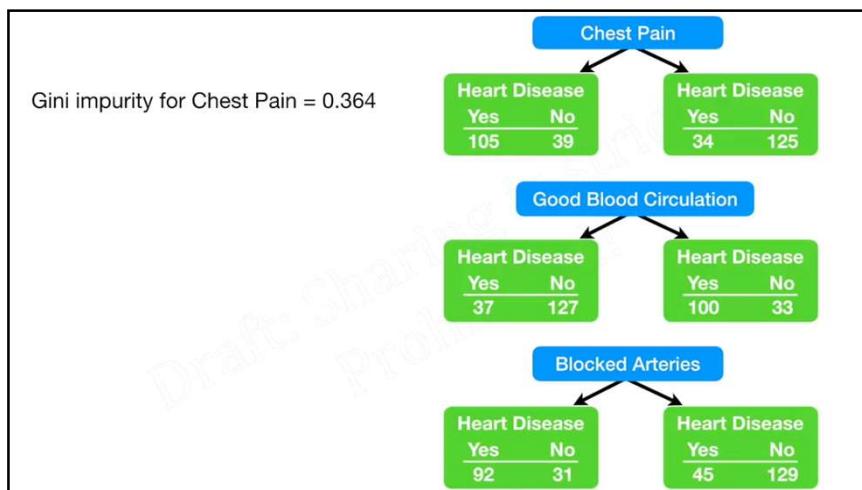
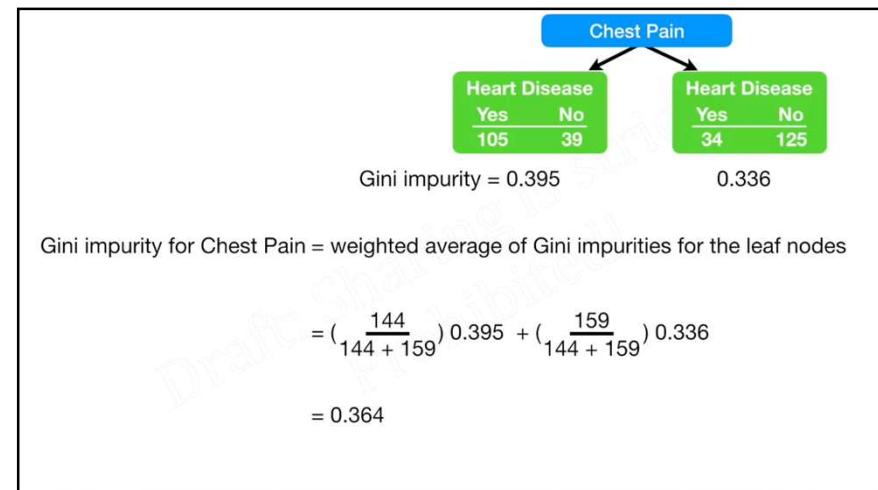
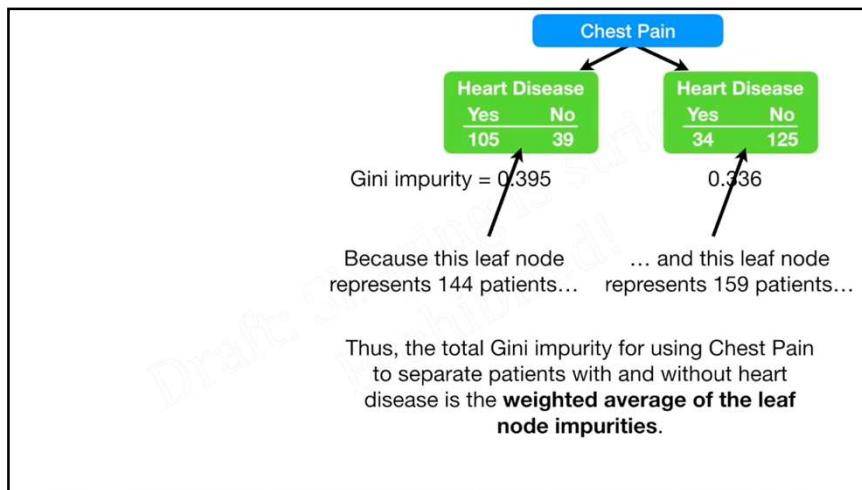
Let's start by calculating Gini impurity for Chest Pain...



$$\text{Gini impurity} = 0.395$$







Gini impurity for Chest Pain = 0.364

Gini impurity for Good Blood Circulation = 0.360

**Good Blood Circulation** has the lowest impurity (it separates patients with and without heart disease the best)...

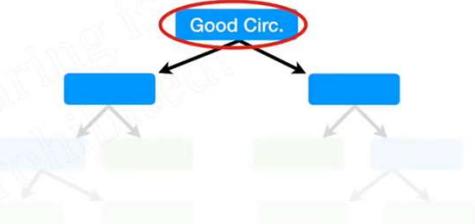
Gini impurity for Blocked Arteries = 0.381

Gini impurity for Chest Pain = 0.364

...so we will use it at the root of the tree.

Gini impurity for Good Blood Circulation = 0.360

Gini impurity for Blocked Arteries = 0.381



Gini impurity for Chest Pain = 0.364

Information gain = impurity reduction

Note:

Yes = 139

No = 164

Total = 303

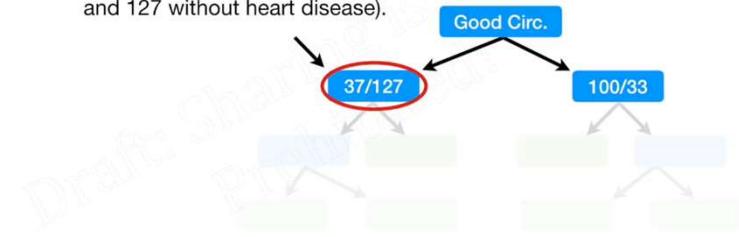
Gini impurity for Good Blood Circulation = 0.360

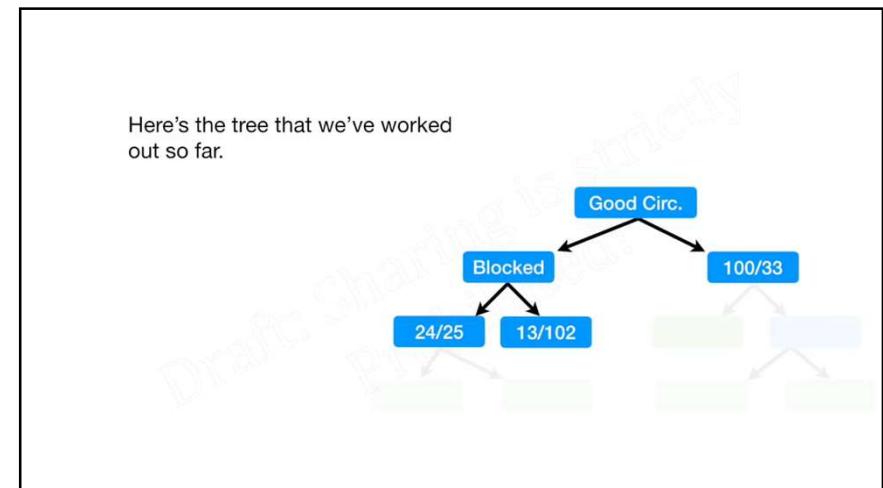
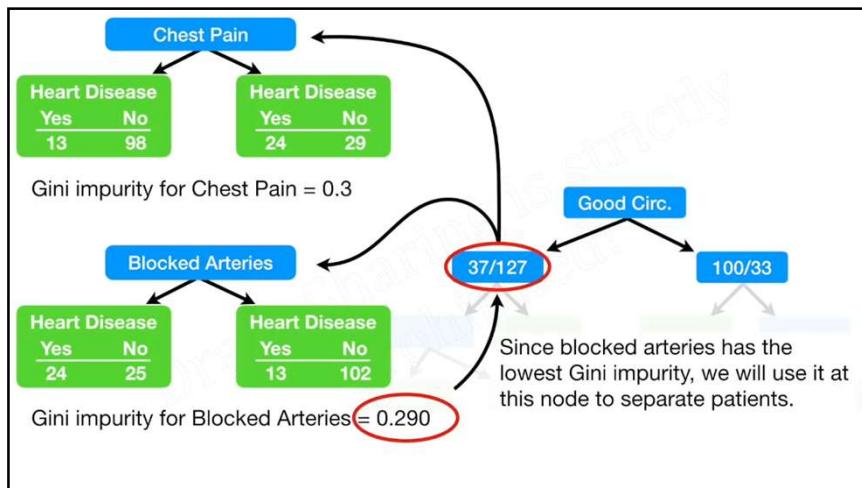
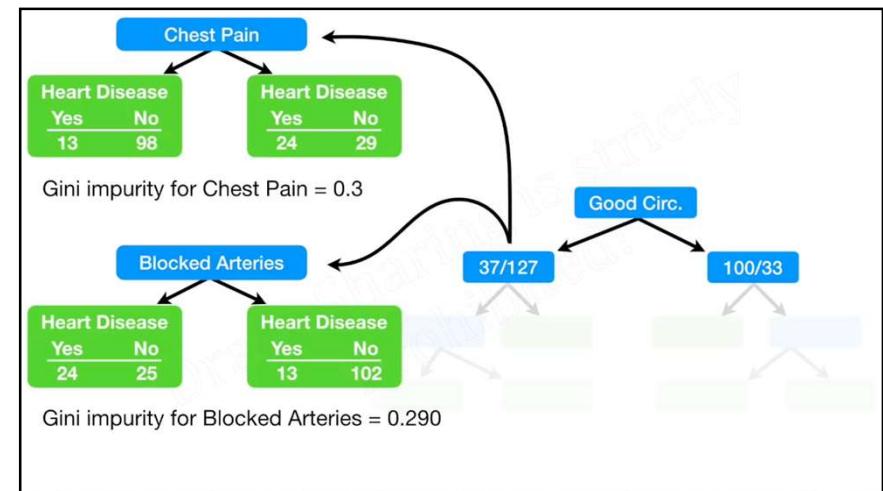
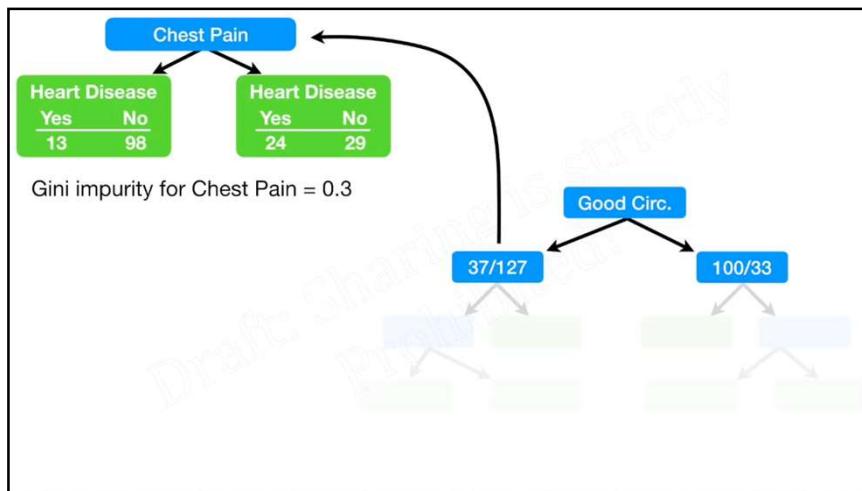
Initial Impurity in the data set was  
 $= 1 - ((139/(139+164))^2 - ((164/(139+164))^2)$   
 $= 1 - 0.210 - 0.293 = .497$

Information gain = .497 - .360 = .137

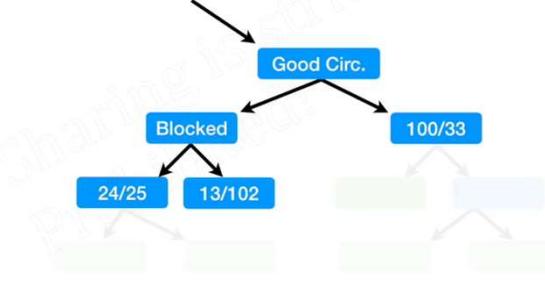
Gini impurity for Blocked Arteries = 0.381

Now we need to figure how well **chest pain** and **blocked arteries** separate these 164 patients (37 with heart disease and 127 without heart disease).

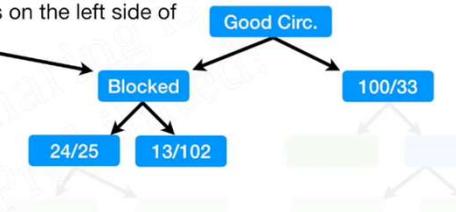




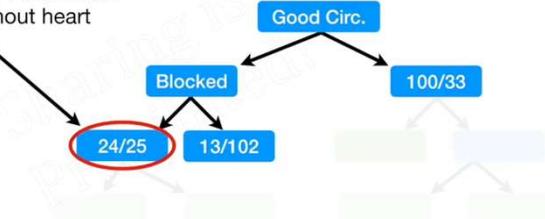
We started at the top by separating patients with Good Circulation...



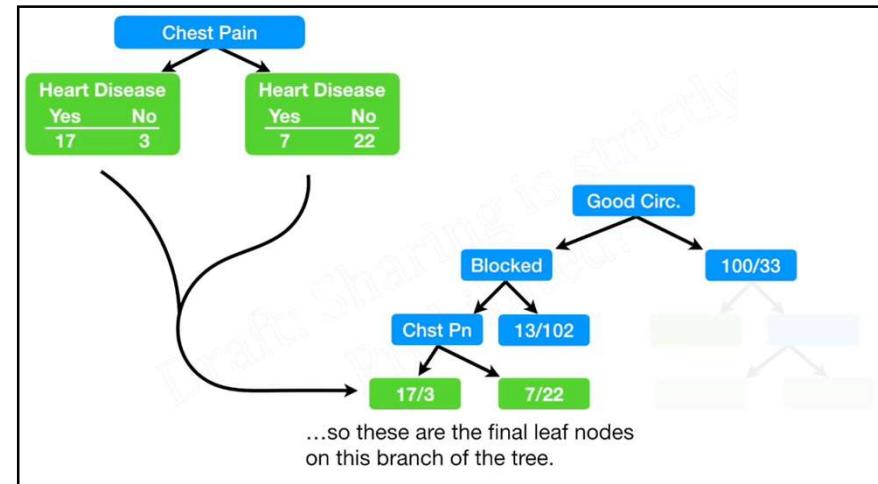
...then we used Blocked Arteries to separate patients on the left side of the tree.



All we have left is Chest Pain, so first we'll see how well it separates these 49 patients (24 with heart disease and 25 without heart disease).



...so these are the final leaf nodes on this branch of the tree.

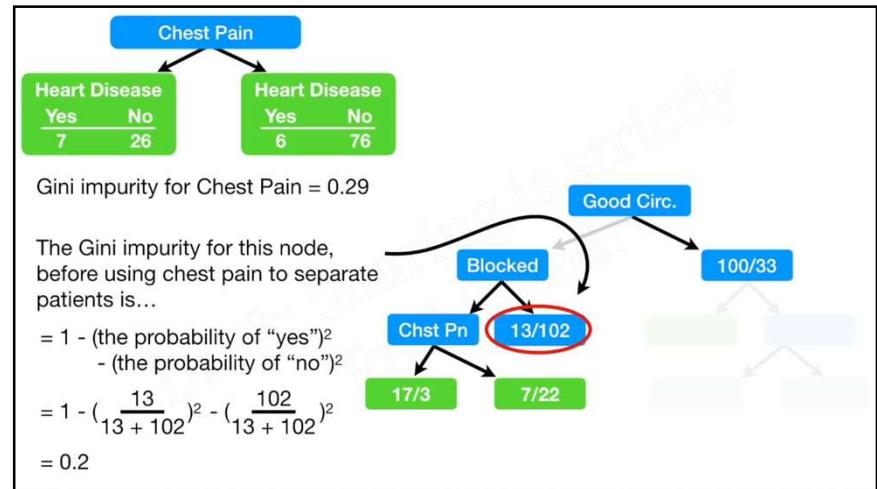
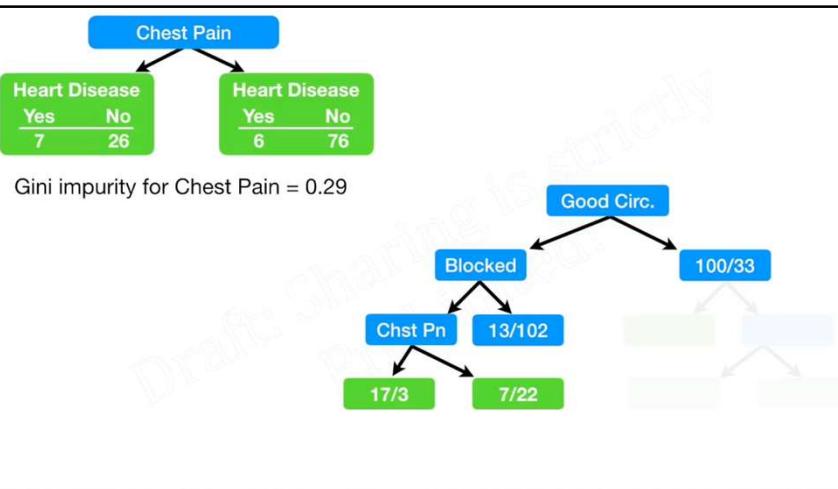
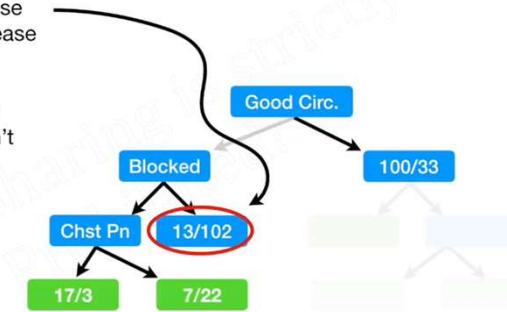


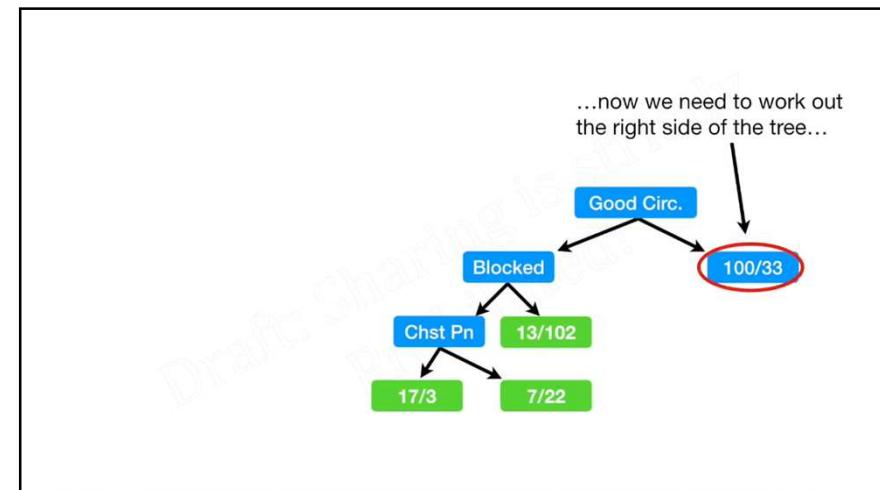
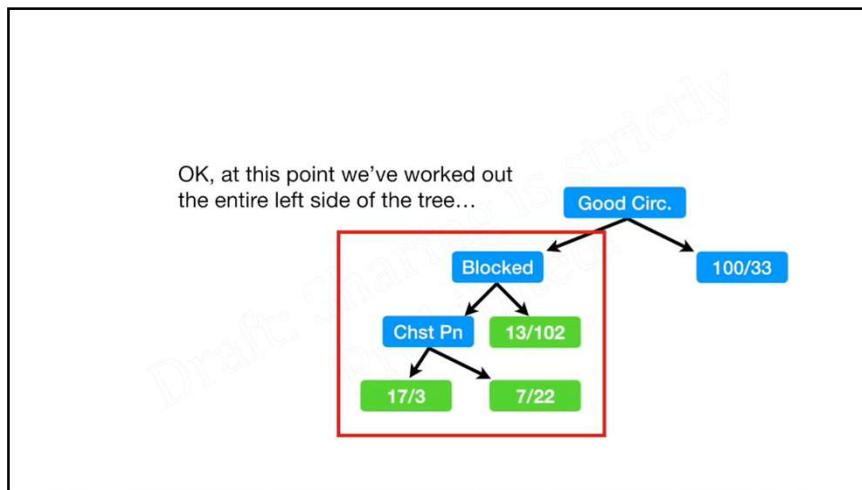
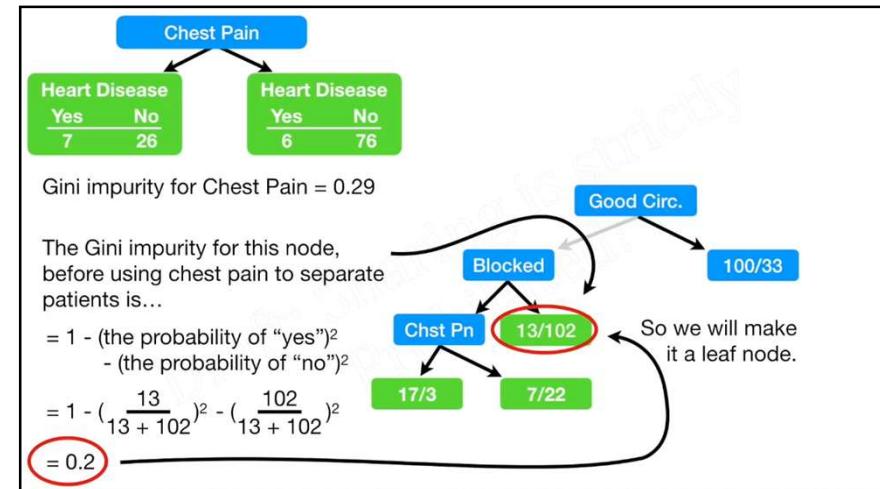
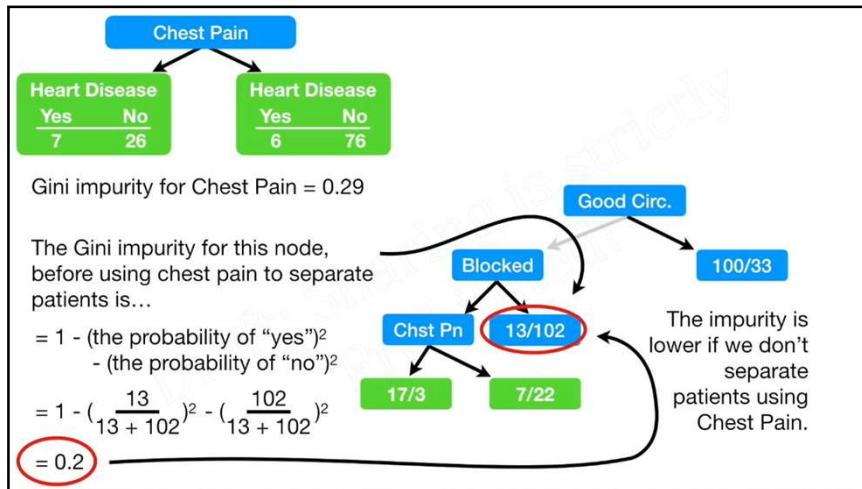
Now let's see what happens when we use chest pain to divide these 115 patients (13 with heart disease and 102 without).



Now let's see what happens when we use chest pain to divide these 115 patients (13 with heart disease and 102 without).

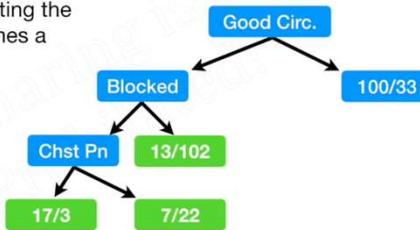
**NOTE:** The vast majority of the patients in this node (89%) don't have heart disease.



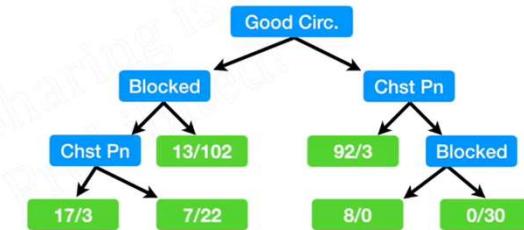


The good news is that we follow the exact same steps as we did on the left side:

- 1) Calculate all of the Gini impurity scores.
- 2) If the node itself has the lowest score, than there is no point in separating the patients any more and it becomes a leaf node.
- 3) If separating the data results in an improvement, than pick the separation with the lowest impurity value.

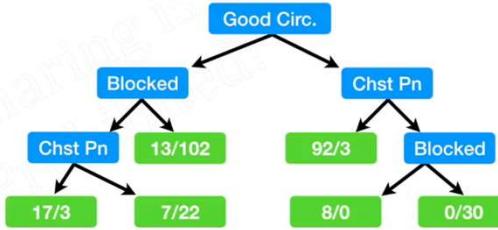


Hooray!!! We made a decision tree!!!



So far we've seen how to build a tree with "yes/no" questions at each step...

...but what if we have numeric data, like patient weight?



Imagine if this were our data...

Weight	Heart Disease
220	Yes
180	Yes
225	Yes
190	No
155	No

	Weight	Heart Disease
Lowest	155	No
	180	Yes
	190	No
	220	Yes
Highest	225	Yes

Step 1) Sort the patients by weight, lowest to highest.

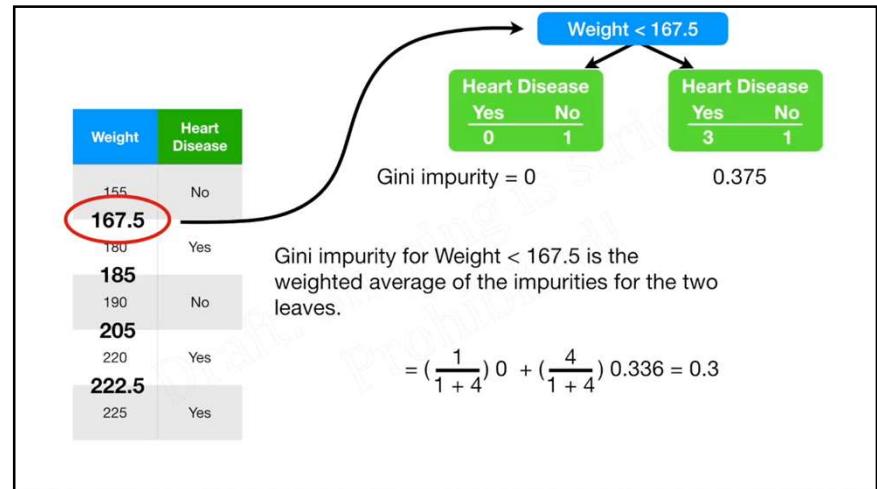
	Weight	Heart Disease
	155	No
	<b>167.5</b>	
	180	Yes
	<b>185</b>	
	190	No
	<b>205</b>	
	220	Yes
	<b>222.5</b>	
	225	Yes

Step 2) Calculate the average weight for all adjacent patients.

	Weight	Heart Disease
	155	No
	<b>167.5</b>	
	180	Yes
	<b>185</b>	
	190	No
	<b>205</b>	
	220	Yes
	<b>222.5</b>	
	225	Yes

Step 3) Calculate the impurity values for each average weight.

Gini impurity = ?



Weight	Heart Disease
155	No
<b>167.5</b>	Yes
180	No
<b>185</b>	Yes
190	No
<b>205</b>	Yes
220	No
<b>222.5</b>	Yes
225	Yes

Gini impurity = 0.3  
 Gini impurity = 0.47  
 Gini impurity = 0.27  
 Gini impurity = 0.4

Weight	Heart Disease
155	No
<b>167.5</b>	Yes
180	No
<b>185</b>	Yes
190	No
<b>205</b>	Yes
220	No
<b>222.5</b>	Yes
225	Yes

The lowest impurity occurs when we separate using **weight < 205**...

Weight	Heart Disease
155	No
<b>167.5</b>	Yes
180	No
<b>185</b>	Yes
190	No
<b>205</b>	Yes
220	No
<b>222.5</b>	Yes
225	Yes

The lowest impurity occurs when we separate using **weight < 205**...  
 ...so this is the cutoff and impurity value we will use when we compare weight to chest pain or blocked arteries.

Now we've seen how to build a tree with...

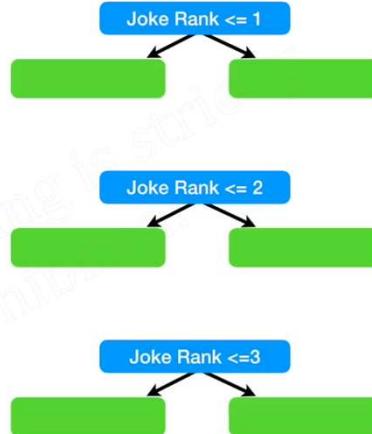
1) "yes/no" questions at each step...

2) Numeric data, like patient weight...

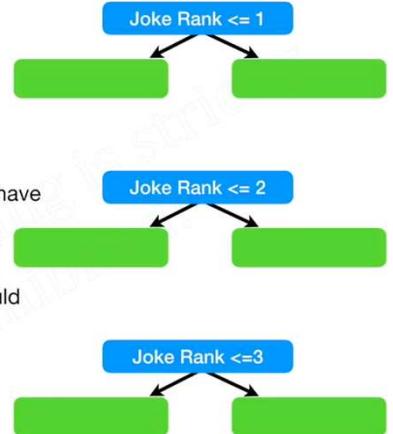
Now let's talk about **ranked data**, like "rank my jokes on a scale of 1 to 4", and **multiple choice data**, like "which color do you like, red, blue or green?"



Rank my jokes...	Likes StatQuest
1	Yes
1	No
3	Yes
1	Yes
etc...	etc...



Rank my jokes...	Likes StatQuestion
1	Yes
1	No
3	Yes
1	Yes
etc...	etc...



Color Choice	Likes StatQuest
Green	Yes
Blue	No
Red	Yes
Green	Yes
etc...	etc...

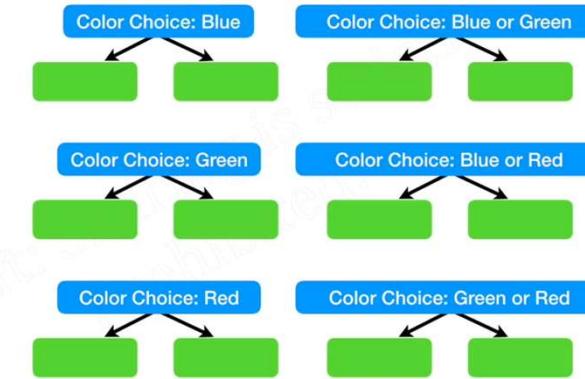
When there are **multiple choices**, like "**color choice can be blue, green or red**", you calculate an impurity score for each one as well as each possible combination.

Color Choice	Likes StatQuest
Green	Yes
Blue	No
Red	Yes
Green	Yes
etc...	etc...

When there are **multiple choices**, like "**color choice can be blue, green or red**", you calculate an impurity score for each one as well as each possible combination.

For this example, with three colors (blue, green and red) we get the following options...

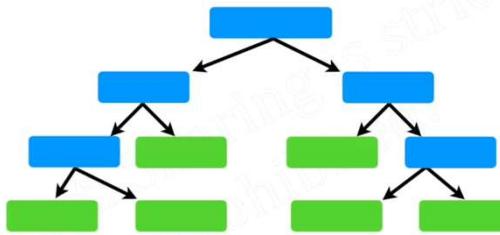
Color Choice	Likes StatQuest
Green	Yes
Blue	No
Red	Yes
Green	Yes
etc...	etc...



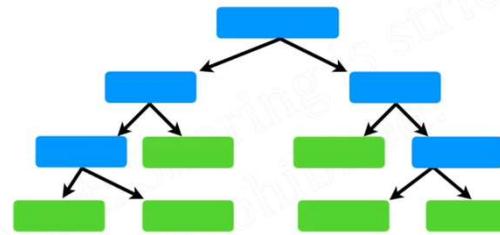
Color Choice	Likes StatQuest
Green	Yes
Blue	No
Red	Yes
Green	Yes
etc...	etc...



Decision Trees are easy to build, easy to use  
and easy to interpret...

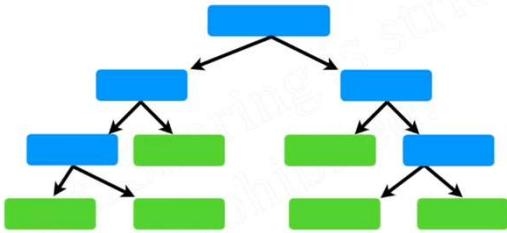


Decision Trees are easy to build, easy to use  
and easy to interpret...

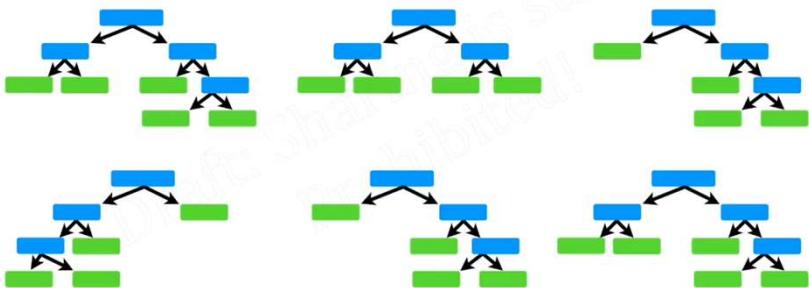


...but in practice they are not that awesome.

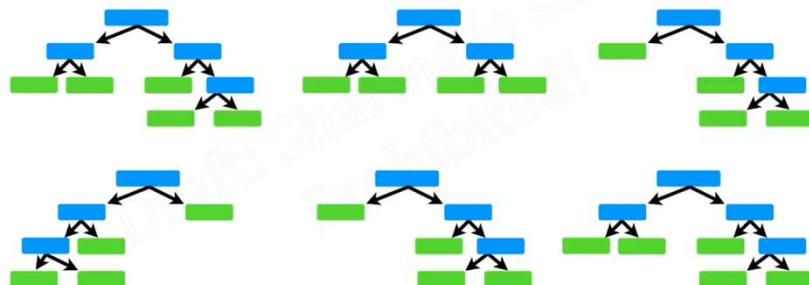
In other words, they work great with the data used to create them, but **they are not flexible when it comes to classifying new samples.**



The good news is that **Random Forests** combine the simplicity of decision trees with flexibility resulting in a vast improvement in accuracy.



So let's make a Random Forest!!!



**Step 1:** Create a “bootstrapped” dataset.

Original Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

Imagine that these 4 samples are the entire dataset that we are going to build a tree from...

Original Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

To create a bootstrapped dataset that is the same size as the original, we just randomly select samples from the original dataset.

The important detail is that we're allowed to pick the same sample more than once.

Original Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

...and here it is.

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

**Step 2:** Create a decision tree using the bootstrapped dataset, but only use a random subset of variables (or columns) at each step.

Bootstrapped Dataset

In this example, we will only consider 2 variables (columns) at each step.

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Thus, instead of considering all 4 variables to figure out how to split the root node...

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes



Bootstrapped Dataset

In this case, we randomly selected **Good Blood Circulation** and **Blocked Arteries** as candidates for the root node.

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
No	Yes	Yes	167	Yes
No	Yes	Yes	167	Yes

Just for the sake of the example, assume that **Good Blood Circulation** did the best job separating the samples.



Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Now we need to figure out how to split samples at this node.

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Just like for the root, we randomly select 2 variables as candidates, instead of all 3 remaining columns.

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
Yes	No	No	167	Yes
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

And we just build the tree as usual, but only considering a random subset of variables at each step.

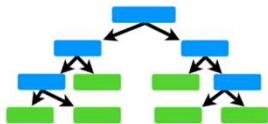
We built a tree...

- 1) Using a bootstrapped dataset
- 2) Only considering a random a subset of variables at each step.

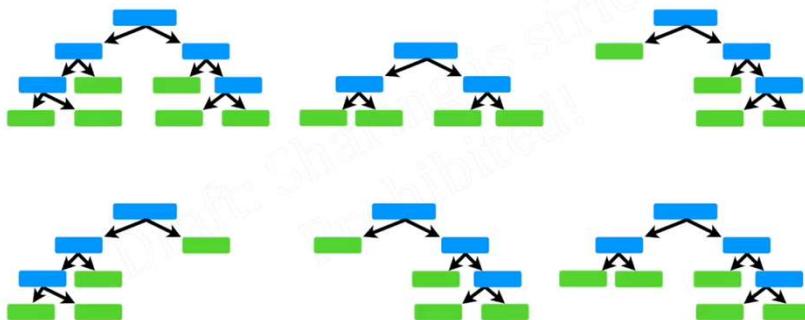
Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

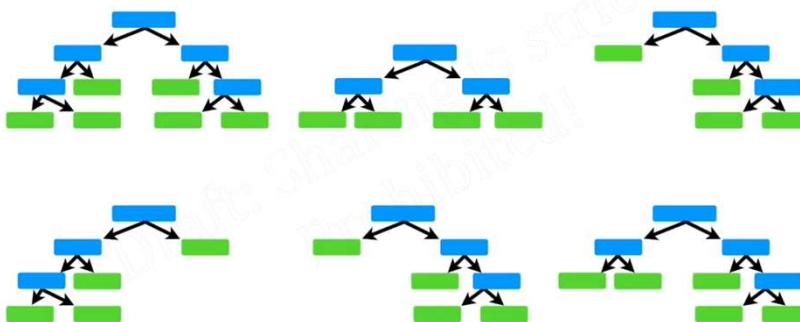
Here's the tree we just made...



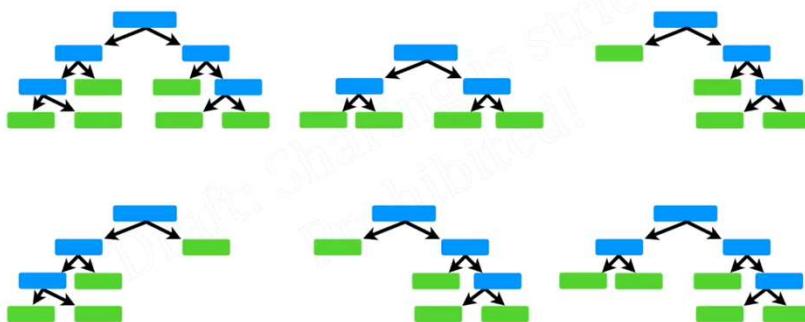
**Now go back to Step 1 and repeat:** Make a new bootstrapped dataset and build a tree considering a subset of variables at each step.



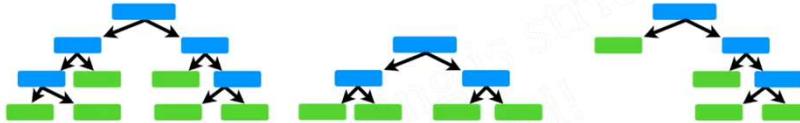
Ideally, you'd do this 100's of times, but we only have space to show 6... but you get the idea.



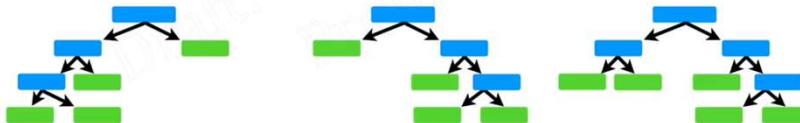
Using a bootstrapped sample and considering only a subset of the variables at each step results in a wide variety of trees.



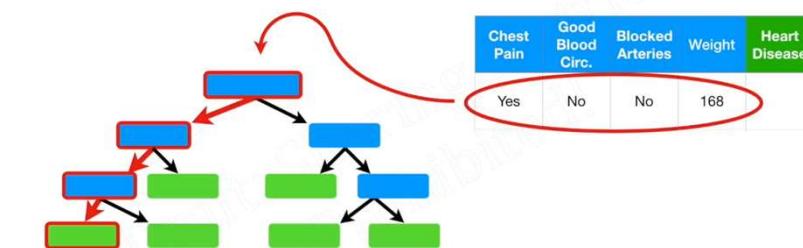
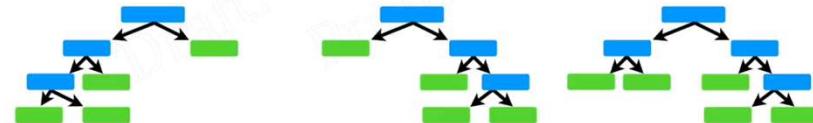
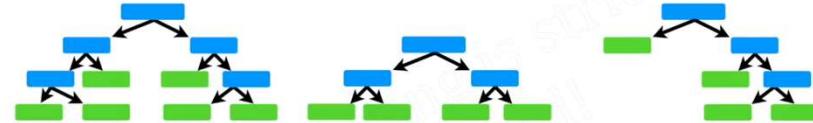
Using a bootstrapped sample and considering only a subset of the variables at each step results in a wide variety of trees.



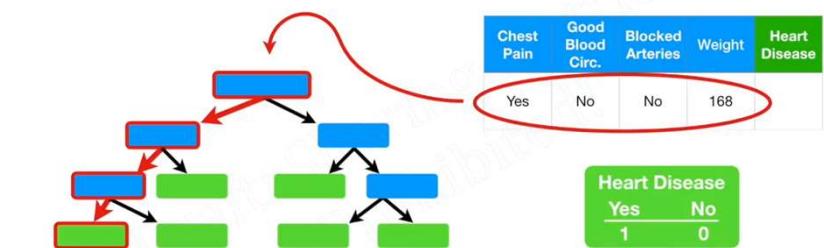
The variety is what makes random forests more effective than individual decision trees.



Sweet!!! Now that we've created a random forest, how do we use it?



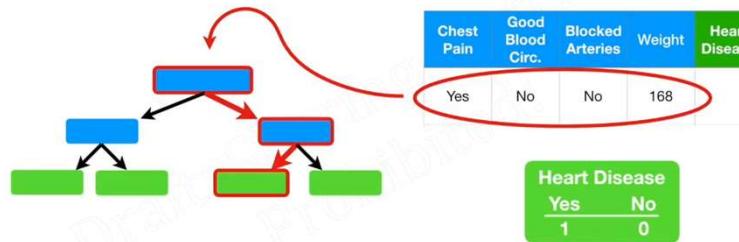
The first tree says  
"Yes"...



The first tree says  
"Yes"...

...and we keep track  
of that here.

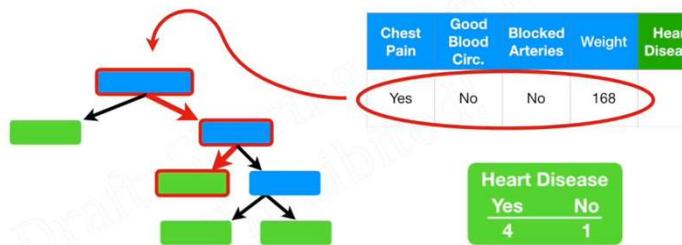
Now we run the data down the second tree that we made...



The second tree also says "Yes"...

...and we keep track of that here.

Then we repeat for all the trees that we made...



After running the data down all of the trees in the random forest, we see which option received more votes.

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	No	No	168	

Heart Disease	
Yes	No
5	1

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	No	No	168	<b>YES</b>

In this case, "Yes" received the most votes, so we will conclude that this patient has heart disease.

Heart Disease	
Yes	No
5	1

### Terminology Alert!!!

Bootstrapping the data plus using the aggregate to make a decision is called "Bagging"

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	No	No	168	<b>YES</b>

Heart Disease	
Yes	No
5	1

How do we know if it's any good?

Remember when we created the bootstrapped dataset?

Original Dataset					Bootstrapped Dataset				
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease	Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No	Yes	Yes	Yes	180	Yes
Yes	Yes	Yes	180	Yes	No	No	No	125	No
Yes	Yes	No	210	No	Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes					

A red box highlights the last row of the original dataset, and a black arrow points from it to the corresponding row in the bootstrapped dataset.

We allowed duplicate entries in the bootstrapped dataset...

**Original Dataset**

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

**Bootstrapped Dataset**

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

As a result, this entry was not included in the bootstrapped dataset.

**Original Dataset**

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes

**Bootstrapped Dataset**

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Typically, about 1/3 of the original data does not end up in the bootstrapped dataset.

**Original Dataset**

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

**Bootstrapped Dataset**

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Here is the entry that didn't end up in the bootstrapped dataset..

**Original Dataset**

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

**Bootstrapped Dataset**

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Original Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

(psst... if the original dataset were larger, we'd have more than just 1 entry over here...)



Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	210	No

Original Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

This is called the "Out-Of-Bag Dataset"

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	210	No

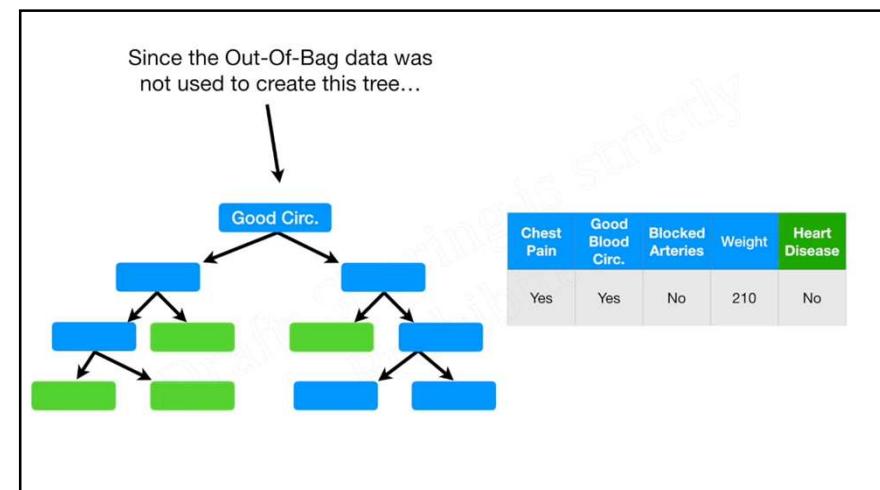
Original Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

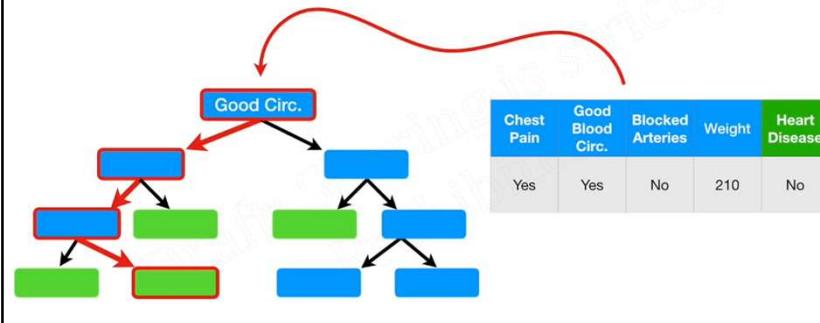
This is called the "Out-Of-Bag Dataset"

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	210	No

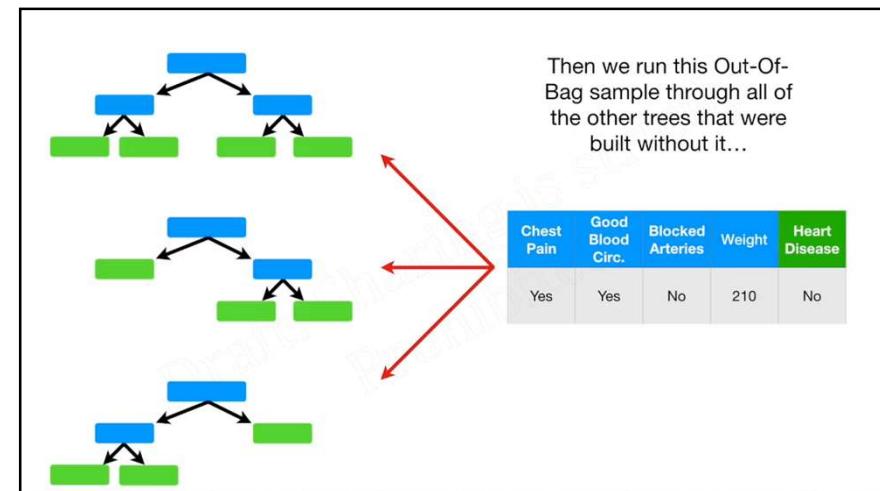
(If it were up to me, I would have named it the "Out-Of-Boot Dataset", since it's the entries that didn't make it into the bootstrap dataset.)



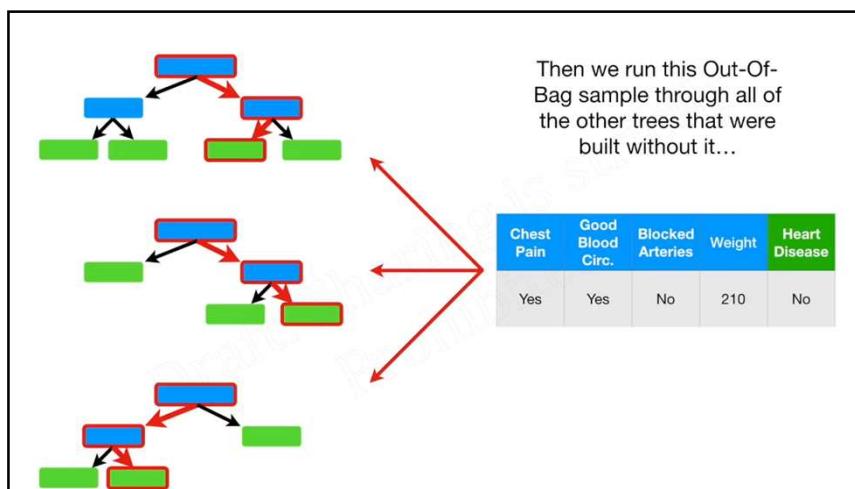
...we can run it through and see if it correctly classifies the sample as "No Heart Disease"



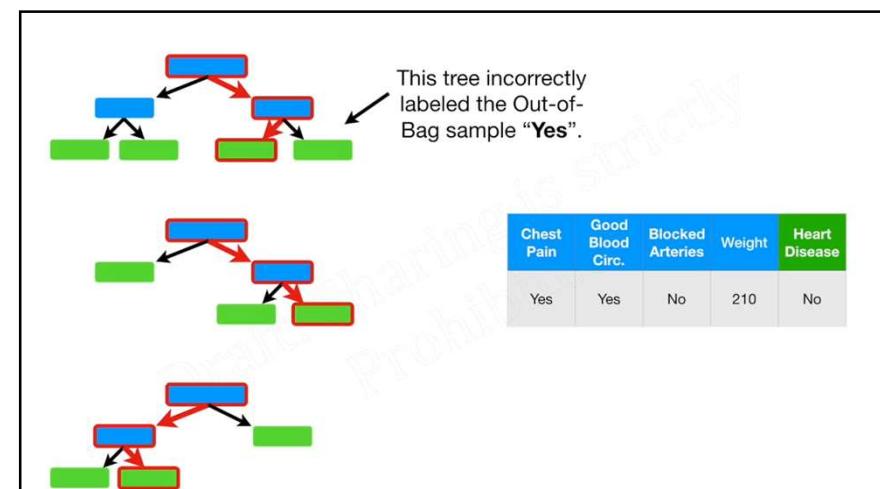
Then we run this Out-Of-Bag sample through all of the other trees that were built without it...

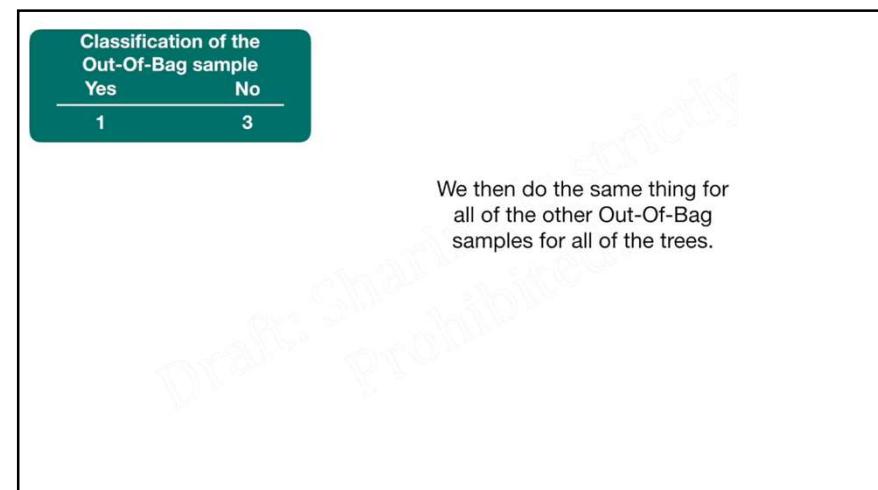
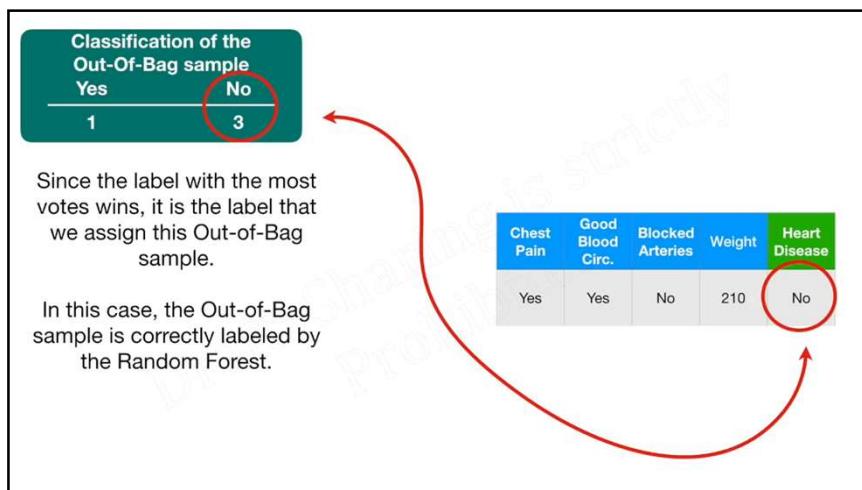
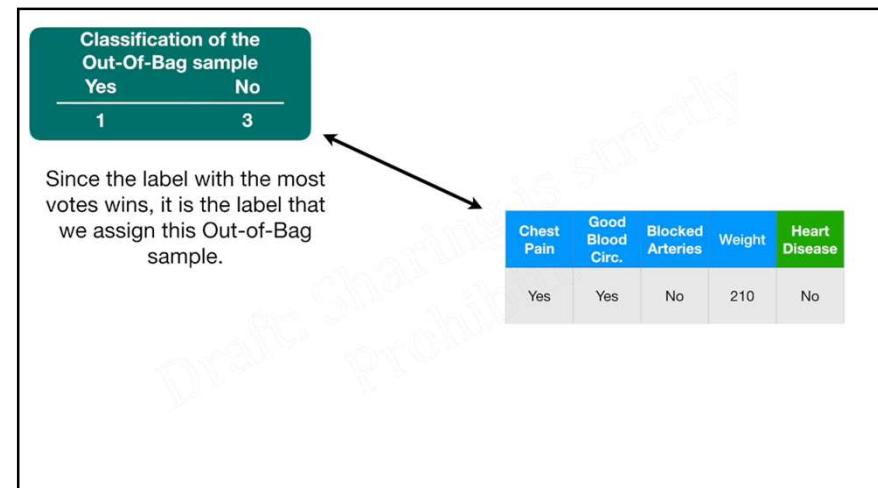
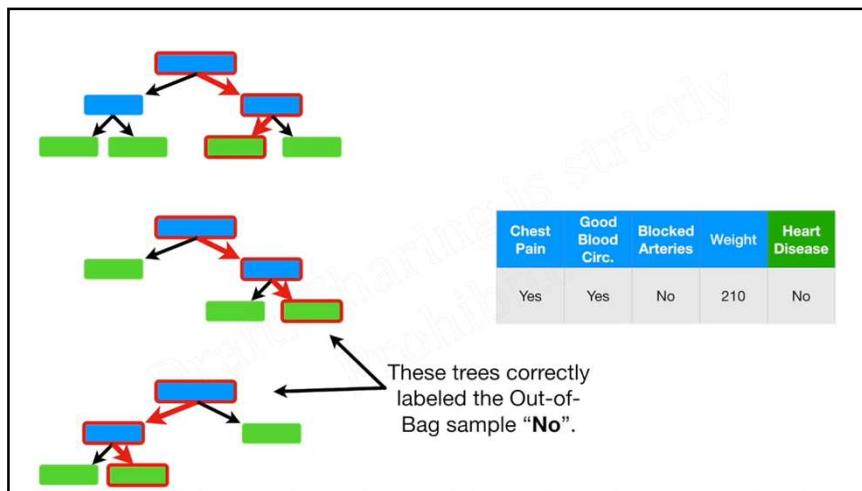


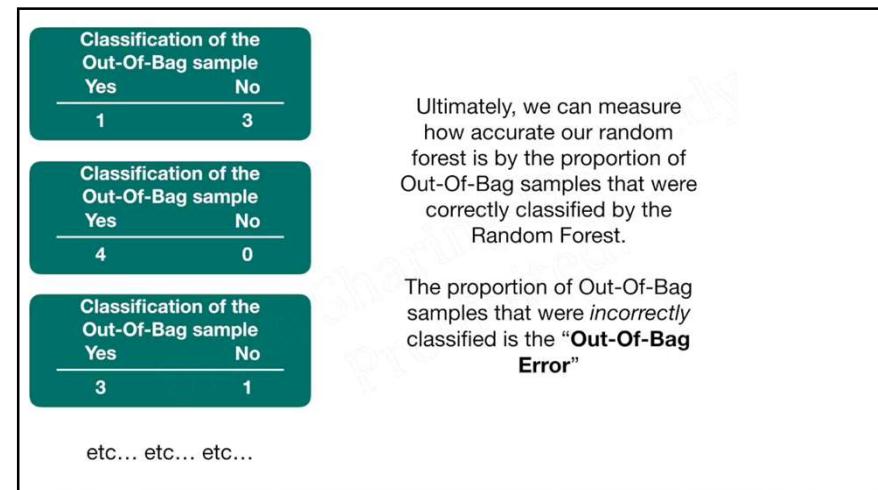
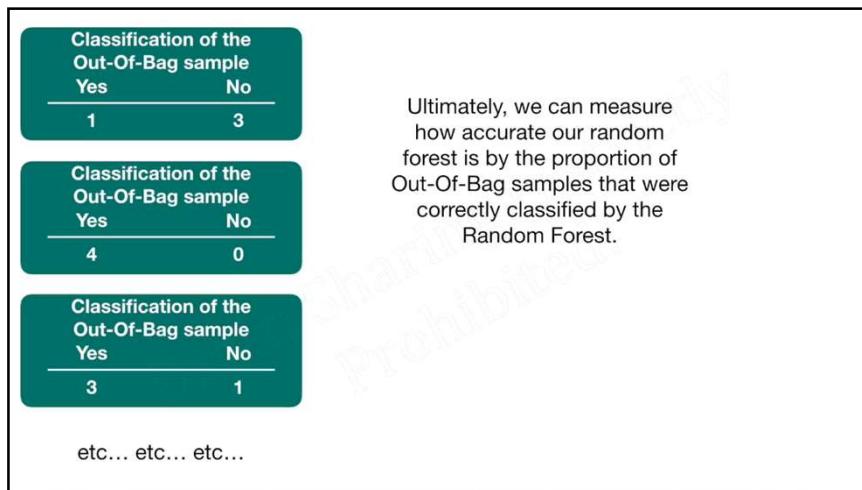
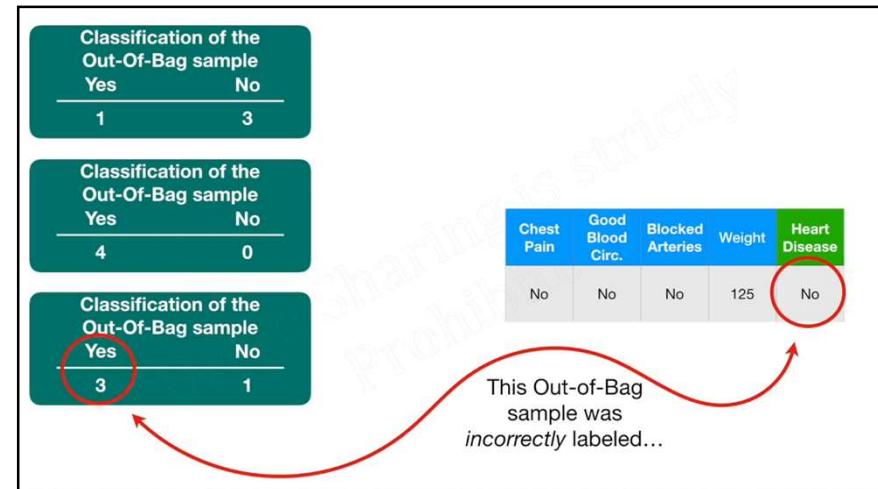
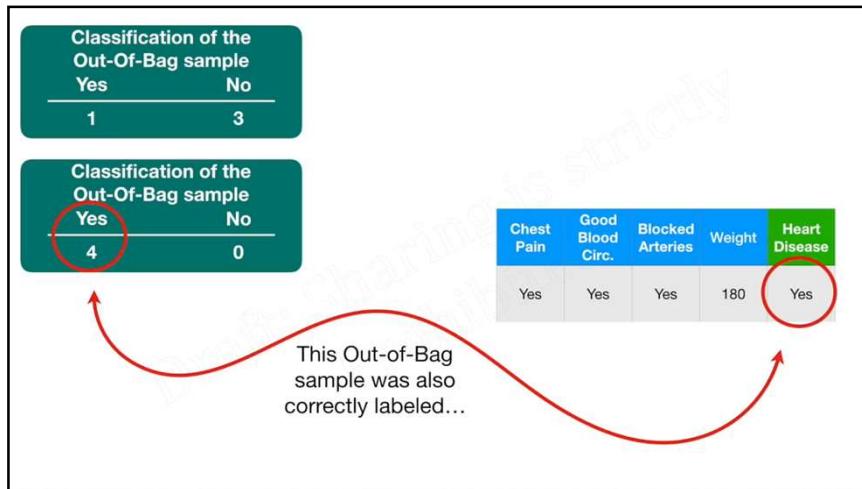
Then we run this Out-Of-Bag sample through all of the other trees that were built without it...



This tree incorrectly labeled the Out-of-Bag sample "Yes".







OK, we now know how to:

- 1) Build a Random Forest
- 2) Use a Random Forest
- 3) Estimate the accuracy of a Random Forest.

Remember when we built our first tree and we only used 2 variables (columns of data) to make a decision at each step?

Bootstrapped Dataset					
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease	
Yes	Yes	Yes	180	Yes	
No	No	No	125	No	
Yes	No	Yes	167	Yes	
Yes	No	Yes	167	Yes	

Now we can compare the Out-Of-Bag error for a random forest built using only 2 variables per step...



...to random forest built using 3 variables per step...



...to random forest built using 3 variables per step...



...and we test a bunch of different settings and choose the most accurate random forest.

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

In other words...

1) Build a Random Forest

2) Estimate the accuracy of a Random Forest.

...change the number of variables used per step...

In other words...

1) Build a Random Forest

2) Estimate the accuracy of a Random Forest.

...change the number of variables used per step...

Do this for a bunch of times and then choose the one that is most accurate.

In other words...

1) Build a Random Forest

2) Estimate the accuracy of a Random Forest.

...change the number of variables used per step...

Typically, we start by using the **SQRT** square of the number of variables and then try a few settings above and below that value.

