## Random Forest

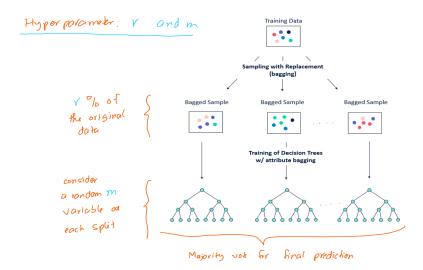
Darta set

X1 Y2 -- X10 +

anestan: How to creak 1,000 different trees from this data?

- 1) Crack treas on only a rondom subset of the data.
- 2) Decide the best split on only a few random variates

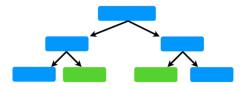
## Big Picture of Random Forest



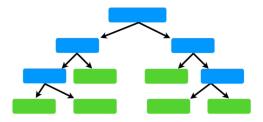
Random Forests are made out of decision trees



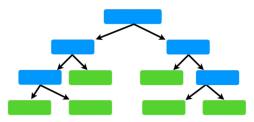
Random Forests are made out of decision trees



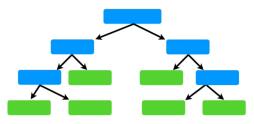
Random Forests are made out of decision trees



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

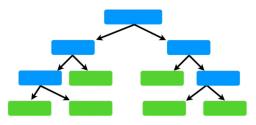


Decision Trees are easy to build, easy to use and easy to interpret...

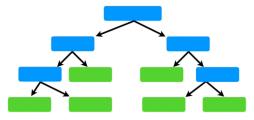


...but in practice they are not that awesome.

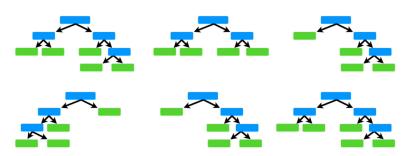
To quote from *The Elements of Statistical Learning* (aka The Bible of Machine Learning), "Trees have one aspect that prevents them from being the ideal tool for predictive learning, namely **inaccuracy**."



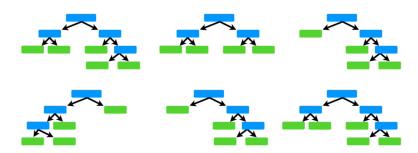
In other words, they work great with the data used to create them, but they are not flexible when it comes to classifying new samples.



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.

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...

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
-----------------------------	---------------------	--------	------------------

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

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.

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

167

Yes

Yes

No

#### Bootstrapped Dataset



This is the first sample that we randomly select...

Chest Pain     Good Blood Blood Poin     Blocked Arteries     Weight Disease       No     No     No     125     No       Yes     Yes     Yes     180     Yes					
		Blood		Weight	
Yes Yes Yes 180 Yes	No	No	No	125	No
	Yes	Yes	Yes	180	Yes

 Yes
 Yes
 No
 210
 No

 Yes
 No
 Yes
 167
 Yes

### **Bootstrapped Dataset**

Chest	Blood	Blocked	Weight	Heart
Pain	Circ.	Arteries		Disease
Yes	Yes	Yes	180	Yes

...so it's the first sample in our 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

#### **Bootstrapped Dataset**

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease					
Yes	Yes	Yes	180	Yes					

This is the second randomly selected sample from the original dataset...

#### Original Dataset **Bootstrapped Dataset** Good Good Blocked Heart Blocked Heart Chest Chest Blood Blood Arteries Pain Disease Pain Arteries Disease Circ. No 125 No No No Yes Yes Yes 180 Yes No Yes Yes Yes 180 Yes No No 125 No ...so it's the second Yes Yes No 210 No sample in our bootstrapped Yes No Yes 167 Yes dataset.

3							
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			
Voc	No	Vos	167	Vos			

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

Here's the third randomly selected sample...

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		
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		

...so here it is in the 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

#### **Bootstrapped Dataset**

	Boototiappou Butaoot					
	hest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease	
١	Yes	Yes	Yes	180	Yes	
-	No	No	No	125	No	
١	res .	No	Yes	167	Yes	

Here's the fourth randomly selected sample (**Note**: it's the same as the third)...



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.

#### We've created 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

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.

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.

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

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

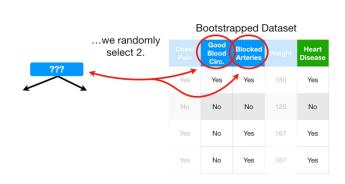
In this example, we will only consider 2 variables (columns) at each step.

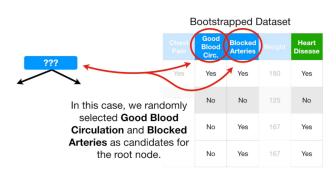
NOTE: We'll talk more about how to determine the optimal number of variables to consider later...

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



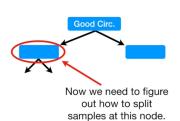




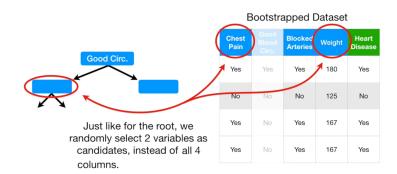
Just for the sake of the example, assume that **Good Blood Circulation** did the best job separating the samples.

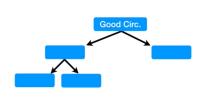


Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes



Boototi appoa Batacot				
Chest Pain		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



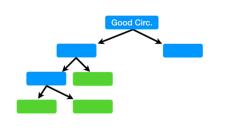


#### **Bootstrapped Dataset**

Yes	Yes	Yes	180	Yes

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

Yes	No	Yes	167	Yes

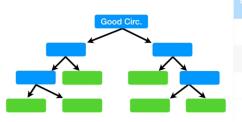


#### **Bootstrapped Dataset**

Yes	Yes	Yes	180	Yes

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

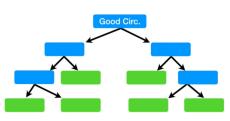
			407	
Yes	No	Yes	167	Yes



Boototrapped Bataset						
				Heart Disease		
Yes	Yes	Yes	180	Yes		
No	No	No	125	No		
Yes	No	Yes	167	Yes		
Yes	No	Yes	167	Yes		

#### We built a tree...

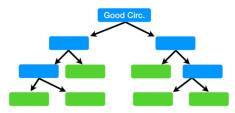
## 1) Using a bootstrapped dataset



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

We built a tree...

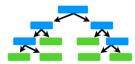
- 1) Using a bootstrapped dataset
- 2) Only considering a random a subset of variables at each step.

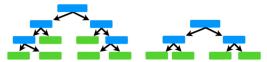


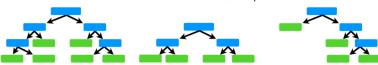
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease		
Yes	Yes	Yes	180	Yes		
No	No	No	125	No		
Yes	No	Yes	167	Yes		
Yes	No	Yes	167	Yes		

Here's the tree we just made...





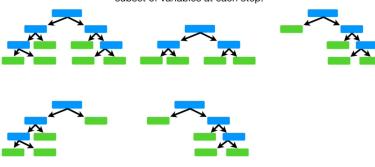


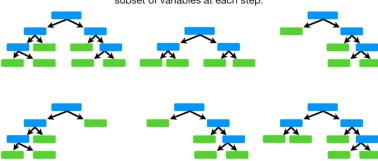




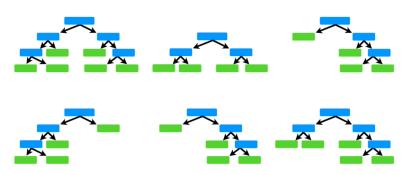




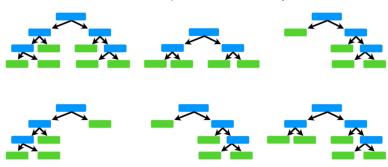




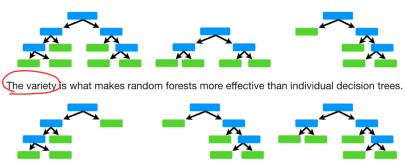
Ideally, you'd do this 100's of times, but we only have space to show 6... but you get the idea.



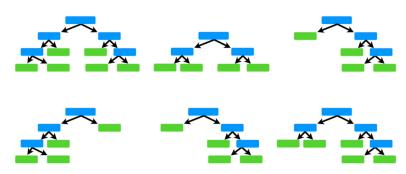
Using a bootstrapped sample and considering only a subset of the variables at each step results in a wide variety of trees.



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



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



## Well, first we get a new patient...

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

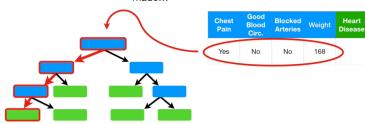


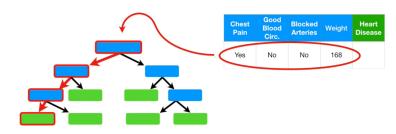
...we've got all the measurements...

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

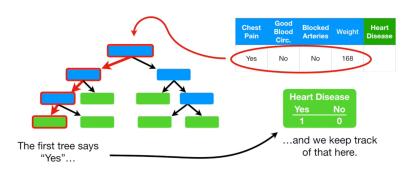
...and now we want to know if they have heart disease or not.

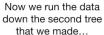
So we take the data and run it down the first tree that we made...

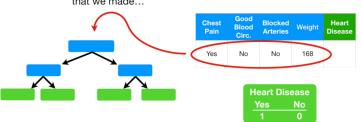


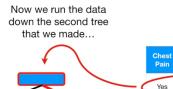


The first tree says "Yes"...

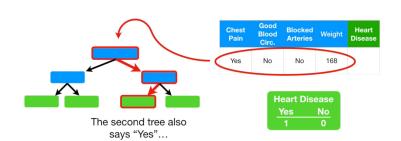


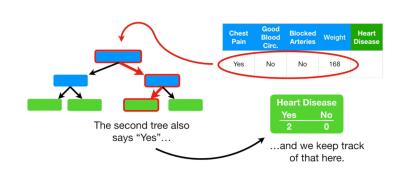


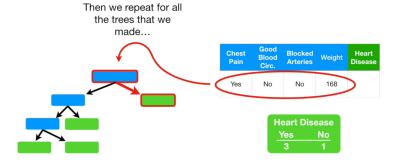


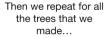




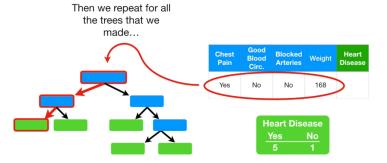












Chest Pain	Good Blood Circ.	Blood Blocked Weigh		Heart Disease
Yes	No	No	168	

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	YES

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



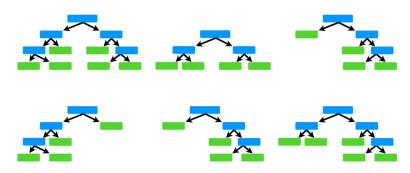
# **Terminology Alert!!!**

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

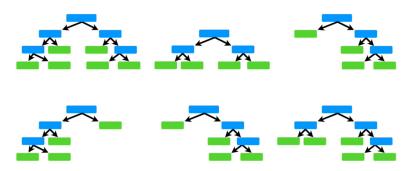
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	No	No	168	YES



# OK, now we've seen how to create and use a random forest...



# How do we know if it's any good?



	Tree 1	Tree 2	Tree 3	Tree 4	Tree 5	Foras 7	True value
065.1	V	1	$\checkmark$	l	6	Ĩ	1
06s. L	V	V	$\cup$	V	O	٥	0
085.3	D	0	V	$\checkmark$	0	0	1
065.4	1	V	l	0	V		0
065-5	O	V	0	l	$\checkmark$	6	0
·	1		,		)		

error mis. classification = 
$$\frac{2}{5}$$
 = 40% is an out of sos

# Out-of-bag errors

Random Forest do not need a validation dataset!

# Remember when we created 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

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease

# 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 Bat					asci	
	Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease	
\	Yes	Yes	Yes	180	Yes	
	<b>↓</b> 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.

	ualase					
Original Dataset						
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease		
No	No	No	125	No		
Yes	Yes	Yes	180	Yes	Į	
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		

# 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	$\downarrow$	
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

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

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

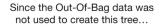
		•		
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	210	No

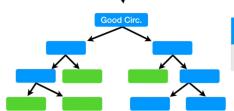
### Original Dataset

_							
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	Ī		

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

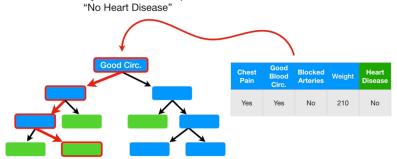
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	210	No

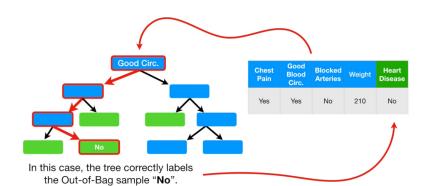


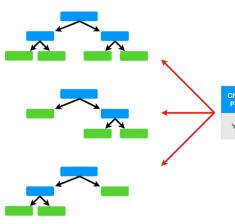


Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	210	No

...we can run it through and see if it correctly classifies the sample as

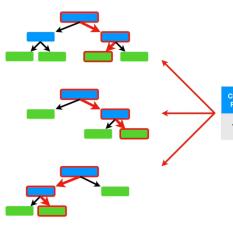






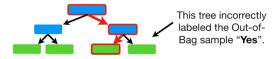
Then we run this Out-Of-Bag sample through all of the other trees that were built without it...

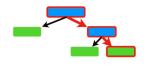
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	210	No



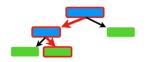
Then we run this Out-Of-Bag sample through all of the other trees that were built without it...

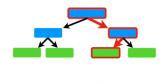
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	210	No





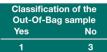
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	No	210	No











Since the label with the most votes wins, it is the label that we assign this Out-of-Bag sample.

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

Classification of the Out-Of-Bag sample Yes No 1 3

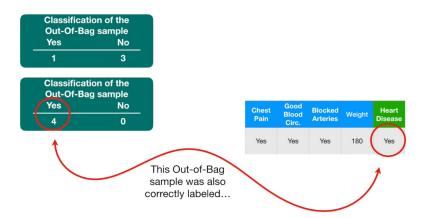
Since the label with the most votes wins, it is the label that we assign this Out-of-Bag sample.

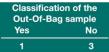
In this case, the Out-of-Bag sample is correctly labeled by the Random Forest.

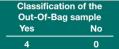


Classification of the Out-Of-Bag sample		
Yes	ag sample No	
1	3	

We then do the same thing for all of the other Out-Of-Bag samples for all of the trees.







Classification of the Out-Of-Bag sample Yes No 1

This Out-of-Bag sample was

Chest

Pain

No

Good

Blood

Circ.

No

Blocked

Arteries

No

Heart

Disease

No

Weight

125

incorrectly labeled...



Classification of the Out-Of-Bag sample Yes No

Classification of the Out-Of-Bag sample Yes No

etc... etc... etc...

Ultimately, we can measure how accurate our random forest is by the proportion of Out-Of-Bag samples that were correctly classified by the Random Forest.



Classification of the Out-Of-Bag sample Yes No

Classification of the Out-Of-Bag sample Yes No

etc... etc... etc...

Ultimately, we can measure how accurate our random forest is by the proportion of Out-Of-Bag samples that were correctly classified by the Random Forest.

The proportion of Out-Of-Bag samples that were *incorrectly* classified is the "Out-Of-Bag Error"

OK, we now know how to:	

3) Estimate the accuracy of a Random Forest.

Build a Random Forest

2) Use a Random Forest

2) 000 4 1 4 1 4 1 1 1 1 1 1 1 1 1 1

OK, we now know how to:

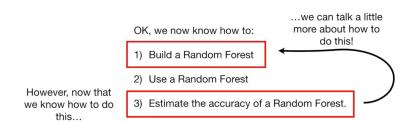
1) Build a Random Forest

2) Use a Random Forest

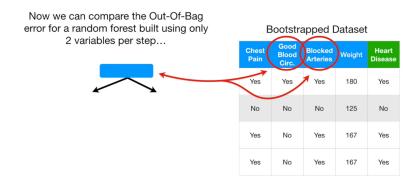
However, now that

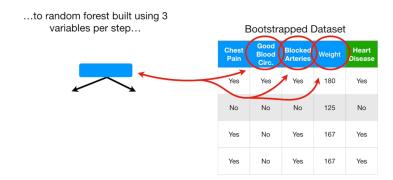
we know how to do this...

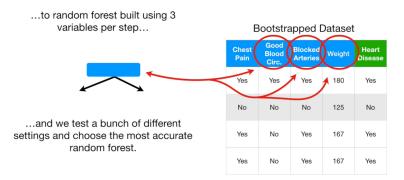
3) Estimate the accuracy of a Random Forest.



Remember when we built our first tree and we only used 2 variables (columns of **Bootstrapped Dataset** data) to make a decision at each step? Good Chest Blocked Heart Weight Blood Arteries Pain Disease Circ. 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

In other words

1) Build a Random Forest

2) Estimate the accuracy of a Random Forest.

In other words...
...change the number of variables used per step...

1) Build a Random Forest

2) Estimate the accuracy of a Random Forest.

In other words...

1) Build a Random Forest

2) Estimate the accuracy of a Random Forest.

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

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

In other words...
....change the number of variables used per step...

1) Build a Random Forest

2) Estimate the accuracy of a Random Forest.

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