

## Finding Missing Values Using Random Forest

Let's say we have the following dataset →

Chest pain	Good Blood circulation	Blocked Arteries	Weight	Heart Disease
NO	NO	NO	125	NO
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	NO			No

Missing value.

Dataset can consist of missing values. Before training our model using this training dataset, we need to fill up these missing values.

We do that by →

- making an (bad) initial guess for this missing values
- gradually refining these bad guesses until it's a good guess.

### 3) Finding Initial guess:

Chest pain	Good Blood circulation	Block Arteries	Weight	Heart Disease)
		NO	125	NO
		Yes	180	Yes
		NO	210	NO
		??	??	NO

4th sample has missing values

We will find initial guess for 'block Arteries'

target is NO

The 4th sample has target values 'NO'  $\rightarrow$  Do not have any heart diseases. So the most frequent/common value found in 'Block Arteries' column

that has target "NO" should be the initial guess for Block Arteries.

Chest pain	Good Blood circulation	Block Arteries	Weight	Heart Disease)
		NO	125	NO
		Yes	180	Yes
		NO	210	NO
		??	??	NO

The samples with 'No' target has 2 'NO' in the

"Block Arteries" feature and '0' 'Yes', so,  
 NO is the most common for "Heart Disease",  
 = NO "target".

So our initial guess for Block Arteries' 4th  
 sample will be 'NO'.

Chest pain	Good Blood circulation	Block Arteries	Weight	Heart Diseases
		NO	125	NO
		Yes	180	Yes
		NO	210	NO
		NO	??	NO

↓                            ↓

Initial guess              Numerical feature.

As weight is the numerical feature, the initial guess will be the median values of weight feature, samples that has Heart disease = NO.

Sample (1) and (3) has heart diseases = NO.

So Median of <sup>this</sup> two data =  $\frac{125 + 210}{2} = 167.5$ .

[You need to know how to find Median of odd and even number of data]

## 2) Refining the initial guess -

### Iteration

We will refine the values until the ~~value~~ guesses become good. This will take multiple iteration. We will see the first iteration; all iterations will be same.

### Iteration 1<sup>o</sup>

	Chest pain	Good Blood circulation	Blocked Arteries	Weight	Heart Disease?
①	No	No	No	125	No
②	Yes	Yes	Yes	180	Yes
③	Yes	Yes	No	210	No
④	Yes	Yes	NO	167.5	No

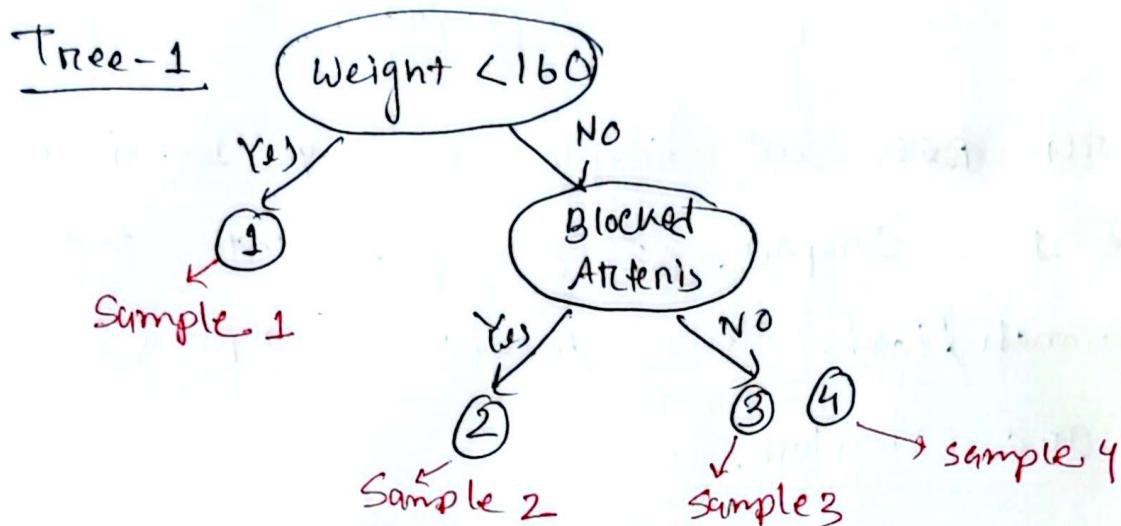
initial guesses.

Now we need to find which samples (among 1, 2, 3) are similar to sample ④ [the one with missing data].

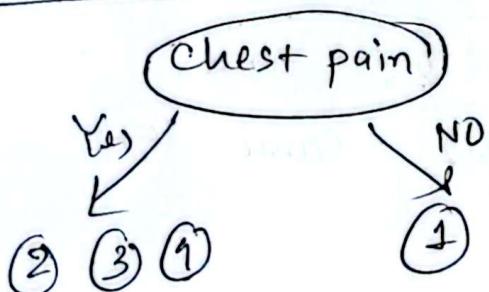
## Steps of finding Similarity:

### ① Build a random forest!

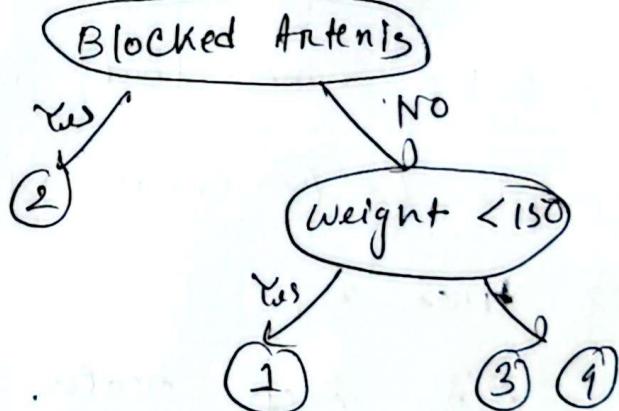
we need to build a random forest consisting of multiple trees. Let's say we have total 10 trees in our random forest. I will draw only tree of them here for your example.



### Tree-2



### Tree -3

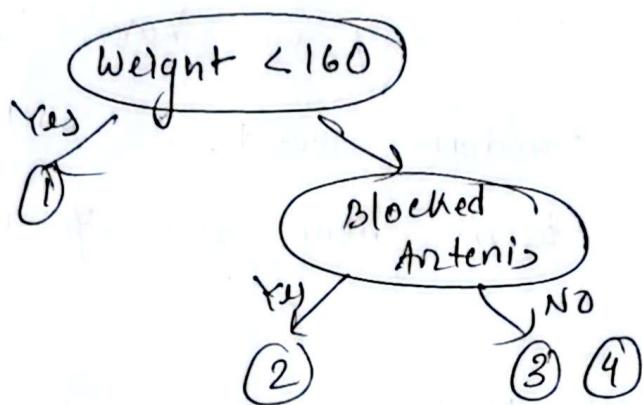


we will have total 10 trees like this for our example.

Step-2: Run all data down each of the tree.

We have already seen that in Step-1 trees.

for example in tree-1:



If we run ~~down~~ all sample 1, 2, 3, 4 ~~but~~ down the tree-1, Sample (3) (4) goes into the same branch/leaf. This means Sample (3) and (4) are similar.

Similarly we will run ~~at~~ the samples down ~~to~~ all of the trees. Sample (2), (3), (4) ends up in the same leaf of tree-2 and Sample (3) (4) ends up in the same leaf of tree-3.

We will keep track of this similarity using a proximity matrix.

### Step-3 : Building Proximity Matrix

As there are 4 samples in our data set, there will be 4 rows and columns in the proximity matrix - which looks like —

	1	2	3	4
1				
2				
3				
4				

These are sample numbers.

In tree-1 , sample ③ ④ ended up in same leaf so we put 1 ~~and~~ into the cells that corresponds to 3 and 4 in the matrix

①	2	3	3	4
1				
2				
3			1	
4		1		

row-3, column-4

row-4, column 3

If there were other pairs of samples ended in the same leaf node, we will keep track of that as well.

After running samples into the second tree, we saw sample ② ③ ④ ends up in the same leaf. so we will put 1 for each pair of (2,3) , (3,4) and (2,4) . As (3,4) pair cells already has 1, we will add 1 with that. So we are basically counting the similar pairs through the proximity matrix.

①	2	2	3	4	
1					(2,3)
2			1	1	(2,4)
3		1		2	(3,4)
4		1	2		
	(2,3)	(2,4)	(3,4)		

Now if we run all the samples down the cell of the remaining 8 trees we end up with this proximity matrix →

	1	2	3	4	
1		2	1	1	
2	2		1	1	
3	1	1		8	
4	1	1	8		

Now we divide each proximity value (each cells) with total number of trees (which is 10 here) and end up with →

	1	2	3	4	
1	0.2	0.1	0.1		→ 1/10
2	0.2	0.1	0.1		
3	0.1	0.1	0.0	0.8	→ 8/10
4	0.1	0.1	0.8		

Now we will use this proximity values to refine our initial guesses.

#### Step-4: updating the missing data

##### ① updating Categorical Feature

Now we need to find the weight frequency for each category (Here Yes and No <sup>are</sup> the categories of Blocked arteries) and check which one is higher/ highest.

Weighted Frequency of a 'category' = Frequency of 'category'  $\times$  weight of 'category'.

We have to find weighted frequency of Yes

and weighted frequency of NO.

Weighted frequency of 'Yes' %

$\therefore$  Weighted frequency of 'Yes' = Frequency of 'Yes'  $\times$  weight of 'Yes'

↓

From dataset

from proximity matrix

## Frequency of Yes -

Chest pain	Blood circulation	Blocked Arteries	Weight	heart Disease)
		No		
		Yes		
		No		
		???		

↳ There is 2 NO and 1 Yes

$\therefore$  Frequency of 'Yes' =  $\frac{1}{3} \rightarrow 1$  Yes  
 $\rightarrow 3$  Samples without missing data

Similarly

Frequency of 'No' =  $\frac{2}{3}$

## Weight of 'Yes'

Weight of 'category' =  $\frac{\text{Proximity of 'category'}}{\text{All proximities}}$

$\therefore \text{Weight of Yes} = \frac{\text{Proximity of Yes}}{\text{All proximities}}$

~~Key~~,

	1	2	3	4
1		0.2	0.1	0.1
2	0.2		0.1	0.1
3	0.1	0.1		0.8
4	0.1	0.1	0.8	

→ All proximities

Proximity of Yes → Proximity w.r.t. value of the samples with only Yes (in Blocked antennas). Here only sample-2 has Yes in Blocked antenna. We will take the proximity from the row of sample

(Q) as sample 4 has the missing value.

$\therefore \text{Proximity of Yes} = 0.1$

All proximities  $\rightarrow$  All of the proximities of the sample that has missing value. The sample (g) has missing value here.

$$\therefore \text{All proximities} = 0.1 + 0.1 + 0.8$$

$$\therefore \text{The weight for Yes} = \frac{0.1}{0.1 + 0.1 + 0.8} = 0.1$$

The weight for NO  $\rightarrow$

$$\text{the weight of NO} = \frac{\text{Proximity of NO}}{\text{All proximities}}$$

	1	2	3	9	
1		0.2	0.1	0.1	
2	0.2		0.1	0.1	
3	0.1	0.1		0.8	
9	0.1	0.1	0.8		

$$\text{Proximity of NO} = 0.1 + 0.8$$

Sample (D) and (B) has NO in blocked antenna column.

$$\therefore \text{Weight of NO} = \frac{0.1 + 0.8}{0.1 + 0.1 + 0.8} = 0.9$$

$\therefore$  weighted frequency of Yes

$$= \frac{\text{Frequency of Yes}}{\text{Total}} \times \text{Weight of Yes}$$

$$= \frac{1}{3} \times 0.1 = 0.03$$

$\therefore$  weighted frequency of NO

$$= \frac{\text{Frequency of NO}}{\text{Total}} \times \text{Weight of NO}$$

$$= \frac{2}{3} \times 0.9 = \boxed{0.6} \rightarrow \text{large}$$

So, our now revised guess for Blocked Arteries will be NO

Chest pain (Good Blood circulation)	Blocked Arteries	Weight	Hazardous
	NO	125	
	Yes	180	
	NO	210	
	NO	??	

↓  
revised

No we need to revise the guess of weight  
(numerical feature)

## ② Updating Numerical Feature

To find revised weight for sample ④, we need to calculate Weighted Average.

Weighted Average =  $(\text{Sample } ① \text{ 's value} \times \text{sample } ① \text{ 's weighted average}) + (\text{Sample } ② \text{ 's value} \times \text{sample } ② \text{ 's weighted average}) + (\text{sample } ③ \text{ 's value} \times \text{sample } ③ \text{ 's weighted average})$

→ we ~~add~~ <sup>sum</sup> all samples without ~~met~~ missing data.

Sample x's weighted average

$$w = \frac{\text{proximity of sample } x}{\text{sum of proximity}}$$

∴ sample ④'s weighted average

$$= \frac{\text{proximity of sample } 1}{\text{sum of proximity}}$$

of sample ④

here sample ④ has <sup>the</sup> missing data, so we will consider sample ④

	1	2	3	4
1		0.2	0.1	0.1
2	0.2		0.1	0.1
3	0.1	0.1		0.8
4	0.1	0.1	0.8	

All proximity

proximity of sample ①

$$\therefore \text{Sample ①'s weighted average} = \frac{0.1}{0.1 + 0.1 + 0.8} = 0.1$$

$$\text{Similarly " ②'s " " } = \frac{0.1}{0.1 + 0.1 + 0.8} = 0.1$$

$$\text{and " ③'s " " } = \frac{0.8}{0.1 + 0.1 + 0.8} = 0.8$$

sample ①'s value from dataset

$$\therefore \text{Weighted Average} = (125 \times 0.1) + (180 \times 0.1) + (210 \times 0.8)$$

Sample ①'s weighted average.

$$= 198.5$$

so this is our revised guess for sample ①'s weight.

So, after first iteration →

Chest Pain	Good Blood circulation	Blocked Arteries	Weight	Heart Diseases
		No	125	
		Yes	180	
		No	210	
		No	198.5	

Now if we do this <sup>steps</sup> over and over again until we reach, the missing value converges, when there is no changes happening in the missing values where each time we recalculate.

This process is for when the dataset has missing value. However, if the testing data/ new unknown data has missing value we will need to <sup>do a</sup> few steps extra.

## Finding missing value in new Sample.

chest pain	Good Blood circulation	Blocked arteries	Weight	Heart Disease
Yes	NO	??	168	We will predict the target

missing data.

We want to classify this new patient, But we have a missing value in Blocked arteries feature. So we need to make a guess about Blocked arteries.

Hence we already have the Random Forest Model. We are trying to predict for a new data. So we run <sup>this data</sup> down ~~this data~~ all the trees of the Random Forest.

Step-1: Create two copies of the data for two target.

copy-1	chest pain	Good Blood Circulation	Blocked Arteries	Weight	Heart Disease
	Yes	NO	??	168	Yes

## Copy - 2

chest pain	Blood Blood Cine	Blocked Arteries	weight	Heart Diseases
Yes	No	??	168	NO

Step-2: Fill the missing values.

NOW, we will run the previous method of

finding missing value for this two copies of

data, and fill the missing values.

Let's assume, we end up with →

## Copy-1

chest pain	Blood Blood cine	Blocked Arteries	weight	Heart Diseases
Yes	No	Yes	168	Yes

↓  
After refining

## Copy-2

chest pain	Blood Blood cine	Blocked Arteries	weight	Heart Diseases
Yes	No	No	168	NO

Step-3: Run the two data down the trees of

### Random forest

We will run down these two copy of new data ~~tree~~ down the trees of random forest. and count ~~win~~ how many times each ~~data~~ copy was correctly labelled (Yes, No)

- = Let's say copy-1 was correctly labelled by the trees 3 times. (Actual target is Yes and trees ~~gave~~ predicted Yes)
- = copy-2 was correctly labelled by the trees 1 time (Actual target is NO and trees predicted NO)

so, copy-1 wins as it was correctly labelled more time than copy-2.

so, our filled in missing data with prediction will be copy-1 →

Chest pain	Blood sugar	Blocked arteries	Weight	Heart Disease
Yes	No	Yes	168	Yes