# Random Forest - Part 1
## Building using, and Evaluating

- Decision trees are not flexing in order to classify new samples.

## Building Random Forest

Step1: create a "bootsraped" dataset.

original-dataset      (Target)      Bootstrapped Dataset

| Chest pain | Good Blood | F3 | F4 | F5 |
|---|---|---|---|---|
| No | No | No | 125 | Yes |
| Yes | Yes | Yes | 180 | No |
| Yes | Yes | No | 210 | No |
| Yes | No | Yes | 167 | Yes |
| | | | | |

| Chest Pain | Good blood | F3 | F4 | F5 |
|---|---|---|---|---|
| Yes | Yes | Yes. | 180 | No |
| No | No | No | 125 | Yes |
| Yes | No | Yes | 167 | Yes |
| Yes | Yes | Yes | 180 | No |
| | | | | |

→a) To create bootstrapped datset, ie of same size of original dataset, we just randomly select samples from the original dataset.

b) we are allowed to pick the same sample more than once.

Step-2: create a decision tree using the bootsrapped dataset, but only use a random subset of variables (features) at each step.
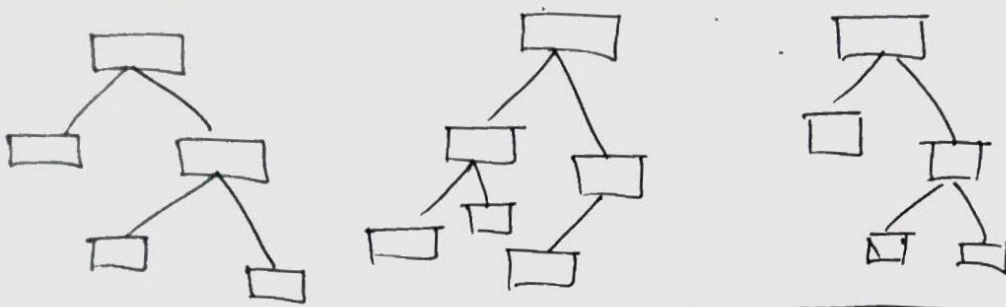
In this case, we will randomly select two variables, "chest" and "Good blood" as a Candidate for root node.

Assume, "Good blood" has lowest cross-entropy,

at this point also, randomly select two variables from "3" variables left.


Good blood

and build the tree as usual.

Now, go back to step 1 and repeat: Make a new bootstrapped dataset and build a tree Consider random subset of variables at each split.



This will result in wide variety of decision trees.

↓

Makes Random forest more effecting than individual

decision trees.

How to predict?

Run the sample on all tree we have created. and Calculates voting of Yes vs No

| Heart disease (Target) | |
|---|---|
| Yes | No |
| 5 | 1 |

⇒ class = Yes

# Random Forests- Part 2
## Missing data and sample clustering

Random forest Consider two types of missing data

1. Missing data in original dataset.
2. Missing data in New sample.

↳ we can make initial guess → refine our guess
↓
using mean, median or other imputation methods.

→ initial guess

| Chest pain | Good blood | Blocked Artery | Weight | Target |
|------------|-----------|----------------|--------|--------|
| Yes | Yes | No | 167.5 | No |

How to refine our guesses

Step-1 : Build Random forest...

Step-2 : Run- all data in all of distinct decision trees.

Find the samples which end up at same 'leaf Node' as our data-sample which we want to make a guess.

No. of Rows = No. of Column
= Total sample

Build Proximity matrix

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | | 2 | 1 | 1 |
| 2 | 2 | | 1 | 1 |
| 3 | 1 | 1 | | 8 |
| 4 | 1 | 1 | 8 | |

Suppose→ 3 and 4 end up in same leaf Node, we will add 1 in P[3][4] and P[4][3]

→ we divide each value by Total No of Trees.

| | 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 | |

# Random Forest

## Terminology:

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

→ Typically $1/3^{rd}$ samples does not end up in bootstrapped dataset.

| 33% of data | → Called the (out - of bag dataset)
$\downarrow$
doesn't appear in bootstrapped dataset.

\* Since out of bag-datasets were not used in ~~class~~ building decision tree, we will use it to evaluate our model.

\* The proportion of out of bag samples incorrectly classified by Random Forest is known as "out of bag error".

1. Build Random Forest

(Change No. of variables)
$\leftarrow$

2. Estimate Accuracy of Random Forest

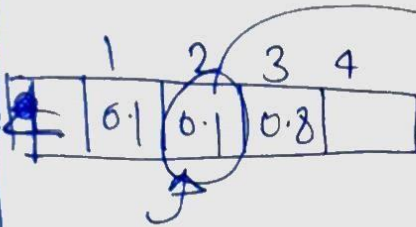Usually, we start with $\sqrt{No. of variables}$.

# Random Forest Part 2.2

Use proximity matrix to calculate missing values.

## 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 | Yes | ??? | ??? | No |

frequency $(Yes) = \frac{1}{3}$, $f(NO) = \frac{2}{3}$

$$weighted\_f(Yes) = \frac{1}{3} \times \underbrace{(The\ weight\ for\ Yes)}_{(O\ 1)} = \boxed{0.03}$$

The weight of Yes $= \left(\dfrac{Proximity\ of\ Yes}{All\ Proximities}\right) =$ Proximity Value for sample 2

```
   1   2   3   4
| 0.1 | 0.1 | 0.8 |   |
```

$$= \left(\frac{0.1}{0.1 + 0.1 + 0.8}\right) = \frac{0.1}{1} = 0.1$$

$$weighted\_f(No) = \frac{2}{3} * \underbrace{(The\ weight\ for\ No)}_{(0.9)} = \underline{0.6}$$

The weight for NO = sample 1 and sample 3

$$\frac{0.8 + 0.1}{0.1 + 0.8 + 0.1} = \frac{0.9}{1} = \boxed{0.9}$$

$$\underline{(No > Yes)}$$

$\longrightarrow$ New guess = No

weightage_average of (weight) = $125 *$ (proximity weight)

$+ (180) *$ (proximity_weight)

Proximity weight $(125) = \dfrac{0.1}{0.1 + 0.1 + 0.8} = 0.1$

$+ (210) *$ (proximity weight)

Proximity weight $(180) = \dfrac{0.1}{0.1 + 0.1 + 0.8} = 0.1$

Proximity_weight $(210) = \dfrac{0.8}{1} = 0.8$

$= 125 \times (0.1) + 180 \times (0.1) + 210 \times (0.8) = \underline{(198.5)}$

After refining guesses, we run random forest again to refine guesses, until missing values converge. $\underline{(6 \text{ to } 7 \text{ times})}$.

# Random Forest - Clustering

## Suppose:

$$P = \begin{array}{|c|c|c|} \hline & 2 & 1 & 1 \\ \hline 2 & & 1 & 1 \\ \hline 1 & 1 & & 10 \\ \hline 1 & 1 & 10 & \\ \hline \end{array} \rightarrow P/10 = \begin{array}{|c|c|c|} \hline & 0.2 & 0.1 & 0.1 \\ \hline 0.2 & & 0.1 & 0.1 \\ \hline 0.1 & 0.1 & & 1 \\ \hline 0.1 & & 1 & \\ \hline \end{array} \rightarrow$$

$\downarrow$ Total (No. of Trees)

$$1 - P/10 = \text{distance matrix} = $$

|   | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 |   | 0.8 | 0.9 | 0.9 |
| 2 | 0.8 |   | 0.9 | 0.9 |
| 3 | 0.9 | 0.9 |   | 0 |
| 4 | 0.9 |   | 0 |   |

$\Rightarrow$ sample '3' and '4' are closest.



## How to handle missing data from new sample

| Yes | No | 22.2 | 168 | Target→ |
|-----|-----|------|-----|---------|

**Step1** - create copies of this same with all target values.

→ | Yes | No | 222 | 168 | Yes |

→ | Yes | No | 222 | 168 | No |

Use iterative method to fill missing values.

suppose we got

| Yes | No | Yes | 168 | Yes | → Run through randomforest and see how many times it has been correctly classified

↓

| Yes | No | No | 168 | No |

↓

If (Yes, Yes) was classified more than (No, No), then prediction for missing data is (Yes) and target is (Yes)