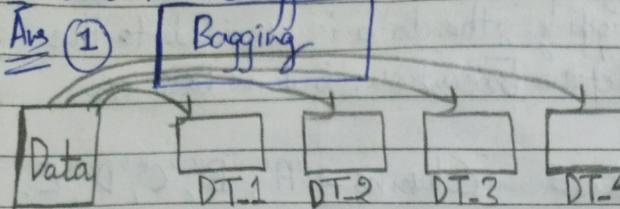


# Bagging .....

Date.....

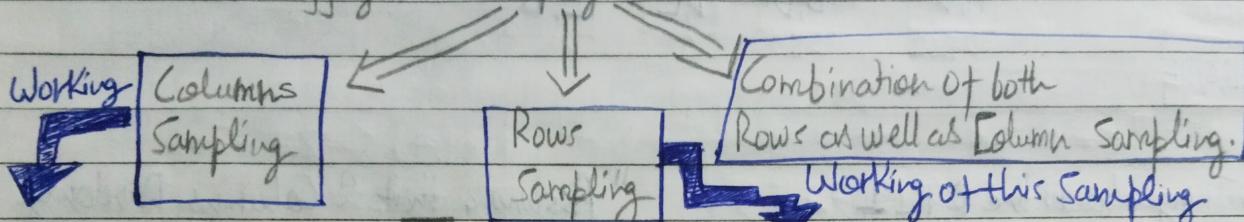
Q What is the difference between the Bagging and Random Forest...

Ans ①



As, per the Sampling, the data will be shared with all the Models.

In Bagging, the "Sampling" is done.



Let's say, My Data  $\Rightarrow$  "df"  
includes 100 Rows, 5 Columns.  
and We have 4 Decision-Tree

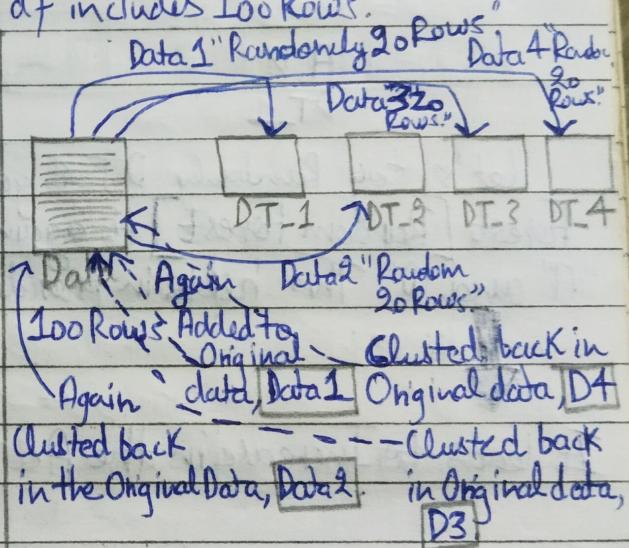
Models to train. Columns  $\Rightarrow$  A, B, C, D, E  
of data "df".

DT-1  $\Rightarrow$  Randomly (Any 2 Columns)

DT-2  $\Rightarrow$  Randomly (Any 2 Columns)

DT-3  $\Rightarrow$  Randomly (Any 2 Columns)

DT-4  $\Rightarrow$  Randomly (Any 2 Columns)

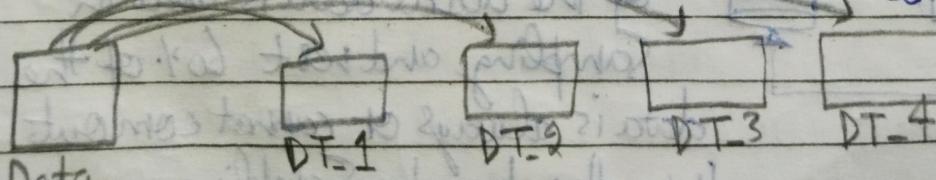


\*So, before Execution, the Columns were randomly provided to the Models [DT-1, DT-2, DT-3, DT-4, DT-5].

②

## Decision Forest

Sampling is same i.e. data is transferred either in Row Sampling Method or Column Sampling Method....



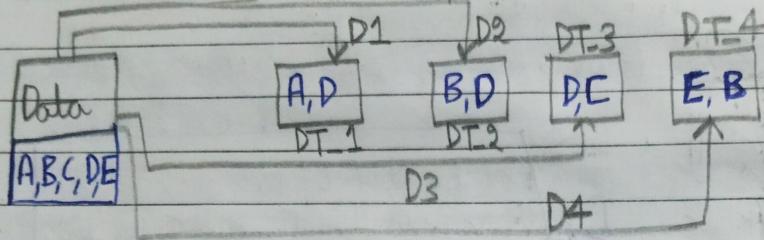
Spiral

## Random Forest....

Date.....

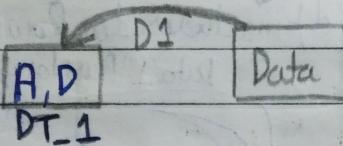
But, In Decision Forest, the Sampling of data is same as in the Case of Bagging but, we have seen that in Bagging, the data is provided to the Models [DT-1, DT-2, DT-3, DT-4] before Execution. But in ~~Decision~~ Random Forest this happens:-

"Column Sampling is done".... Columns = "A", "B", "C", "D", "E".



\* Talking about the DT-1

"Assuming that 2 Columns Randomly will be provided to each Column."



Let's say randomly, DT-1 got columns "A", "D". But, Here in decision Forest [Random Forest] it again randomly selects from the columns "A" and "D". This helps in providing more Randomness to the Models.

It leads to Increase in the Accuracy of the Prediction.....

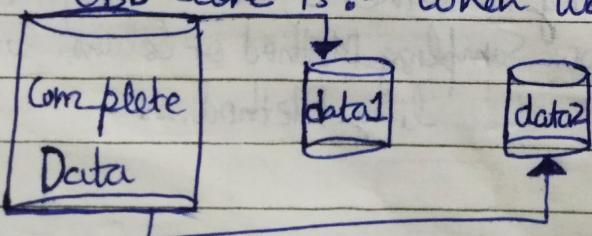
This happens with all the Models available in the Random Forest...

Q What is OOB in Bagging?

Ans Out of Bag Score or

Out of Bag Score [OOB\_Score]

OOB Score is :- When we take the data after Sampling i.e;



data 1 or data 2 only 30%. of the actual data is used in Sampling and rest 60% of the data is always or cannot come out from the bag for Sampling...

Spiral

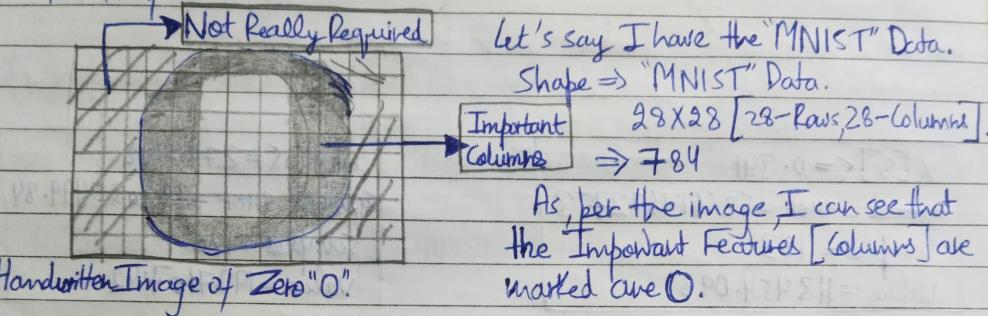
## Feature Importance....

Date.....

Q Why we use the Feature Importance?

A<sub>1</sub> Feature Importance  $\Rightarrow$  Feature Importance is used to take the important features [columns] from the complete data. These important features are important because their availability is majorly affecting the Prediction [Output].

For example:-



So, the "Random Forest" have the availability to calculate the Feature Importance within the data...

## Calculating Feature Importance

(i) Calculating the Feature Importance of the Decision Trees !!

### FORMULA....

Where:-  $N-t \Rightarrow$  No. of Rows in the Node.

$N \Rightarrow$  Total No. of Rows in the table.

impurity  $\Rightarrow$  gini impurity.

$$\text{Node Impurity} = \frac{N-t}{N} \left[ \text{impurity} - \left( \frac{N-t-r}{N} \times \text{right impurity} \right) - \right.$$

$$\left. \left( \frac{N-t-l}{N} \times \text{left impurity} \right) \right]$$

$N-t-r \Rightarrow$  No. of Rows in the right

Node.

$N-t-l \Rightarrow$  No. of rows in the left Node.

right impurity  $\Rightarrow$  gini impurity of right node. left impurity  $\Rightarrow$  gini impurity of left node.

Spiral

Date.....

For example:-

$x[7] \leq 154.519$
Squared_error = 2474808801.4
Samples = 311
Value = 148543.65

$x[8] \leq 0.341$
Squared_error = 526617672753
Samples = 183
Value = 113454.098

$x[7] \leq 278.708$
Squared_error = 983035499.84
Samples = 128
Value = 198710742

Squared_error = 198469644.133	Squared_error = 118663744.084	Squared_error = 812549382.84	Squared_error = 50218750.0
Samples = 105	Samples = 78	Samples = 125	Samples = 3
Value = 97051.730	Value = 135534.135	Value = 196584.6	Value = 267300.0

"7" Feature  $\Rightarrow 2 \Rightarrow$  Let's say the  $F_i$  will be " $X$ "  $\Rightarrow X_1 + X_2 / Y + X_1 + X_2$   
 "8" Feature  $\Rightarrow 1 \Rightarrow$  Let's say the  $F_i$  will be " $Y$ "  $\Rightarrow Y / X_1 + X_2 + Y$

Firstly Calculating for the "X"

$$X_1 \Rightarrow \frac{311}{389} / 2474808801.4 \left[ \left( \frac{128}{389} \times 983035499.84 \right) - \left( \frac{183}{389} \times 526617672.753 \right) \right]$$

$$\Rightarrow 0.799 / 2474808801.4 - \left[ \left( 0.329 \times 983035499.84 \right) - \left( 0.470 \times 526617672.753 \right) \right]$$

Date.....

$X_1$

$$\Rightarrow 0.799 / 2474808801.4 - [323418679.4 - 247510306.2]$$

$$\Rightarrow 0.799 / 2474808801.4 - 75908373.8$$

$$\Rightarrow 0.799 / 2398900428$$

$$\Rightarrow \boxed{1916721442} - ②$$

$$X_2 \Rightarrow 125 / 983035499.84 - \left[ \frac{3}{389} \times 50213750.0 \right] - \left[ \frac{125}{389} \times 812549382.84 \right]$$

$$\Rightarrow 0.329 / 983035499.84 - \left[ (7.71 \times 50213750.0) - (0.32 \times 812549382.84) \right]$$

$$\Rightarrow 0.329 / 983035499.84 - \left[ 387148012.5 - 260015802.5 \right]$$

$$\Rightarrow 0.329 / 983035499.84 - \boxed{127132210}$$

$$\Rightarrow 0.329 / 855903289.8$$

$$\Rightarrow \boxed{281592182.3} - ③$$

Calculating for the 4

$$Y \Rightarrow 183 / 526617672.753 - \left[ \frac{78}{329} \times 118663744.084 \right] - \left[ \frac{105}{329} \times 198469644.133 \right]$$

$$\Rightarrow 0.55 / 526617672.753 - \left[ (0.23 \times 118663744.084) - (0.31 \times 198469644.133) \right]$$

$$\Rightarrow 0.55 / 526617672.753 - \left[ 27292661.14 - 61525589.68 \right]$$

$$\Rightarrow 0.55 / 526617672.753 - \boxed{-34232928.54}$$

$$\Rightarrow 0.55 / 560850601.3$$

$$\Rightarrow \boxed{308467930.7} - ③$$

Date.....

Feature Importance of "7" Feature will be :-  $\frac{1916721442 + 281592182.3}{2198313624 + 4}$

$$\Rightarrow \underline{2198313624}$$

$$2506781455$$

$$\Rightarrow \boxed{0.87694}$$

Feature Importance of "8" Feature will be :-  $\frac{308467830.7}{2506781455}$

$$\Rightarrow \boxed{0.12305}$$