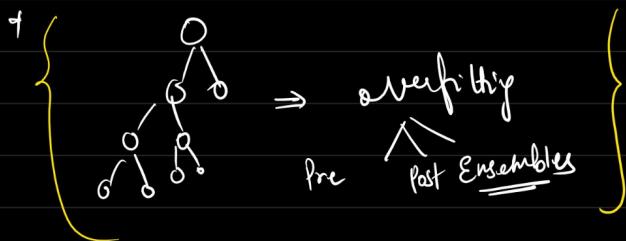
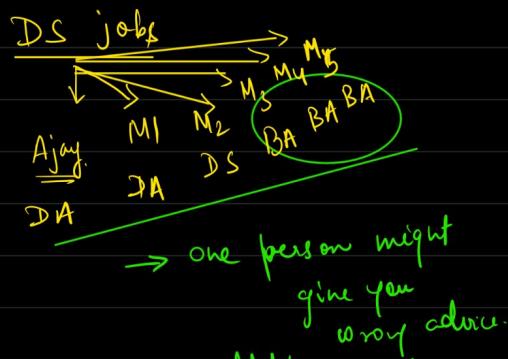


Till now

- Linear data  $\rightarrow$  MLR, Log Regression
- Non linear  $\rightarrow$  DT, SVM

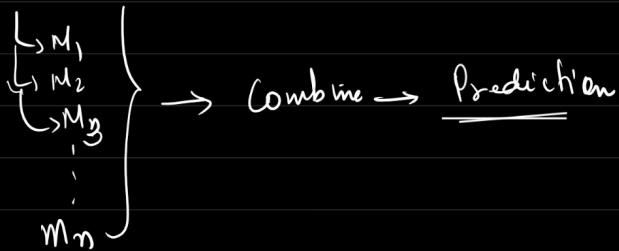


## \* Ensembles and its technique



\* data  $\rightarrow$  Model  $\rightarrow$  train  $\rightarrow$  prediction  $\rightarrow$  multiple mentors  $\rightarrow$  chance of getting wrong is minimized.

+ data



## \* Ensembles: combine multiple models

: Prediction : more stable and accurate as compared to individual model.

of same Algorithm

- |  |  |
|--|--|
| $DT_1 \rightarrow \text{max depth} - 5$  | $DT_2 \rightarrow \text{max depth} - 10$ |
| $DT_3 \rightarrow \text{max depth} - 12$ |  |

↓ different algorithms

- |                                    |                   |
|------------------------------------|-------------------|
| $\rightarrow$ Logistic Reg         | $\rightarrow$ SVC |
| $\rightarrow$ DTC <sub>1</sub> (S) |                   |

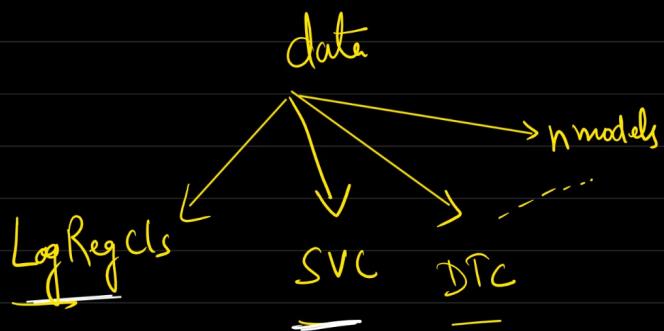
$$\left\{ \begin{array}{l} \\ n \end{array} \right. \rightarrow DT_1 (1^{\text{st}}) // \\ \rightarrow NB_C$$

\* Ensemble → No necessarily only one type of Algorithms

### Ensemble techniques

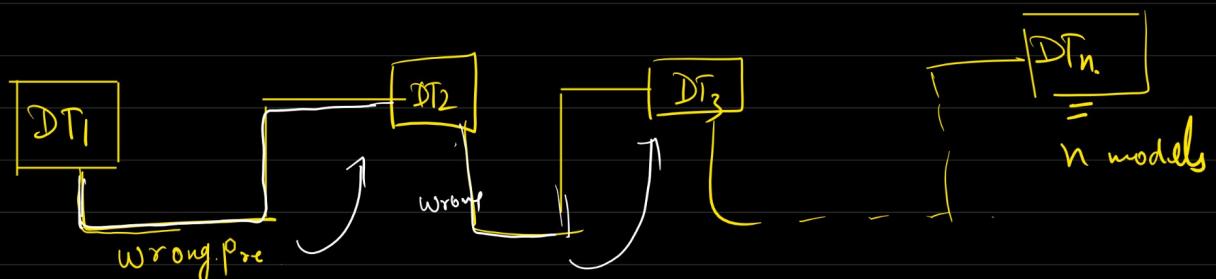


\* Parallel technique of Ensembles

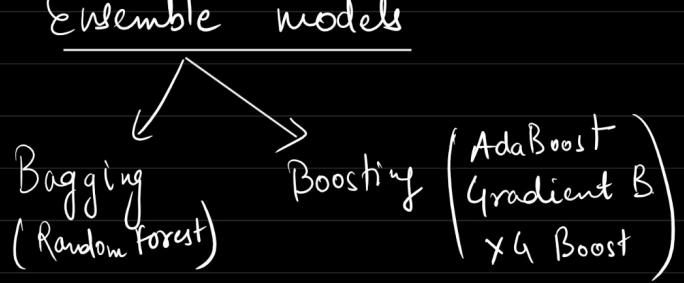


→ All the models here are built parallelly and independent of each other.

\* Sequential technique of Ensembles

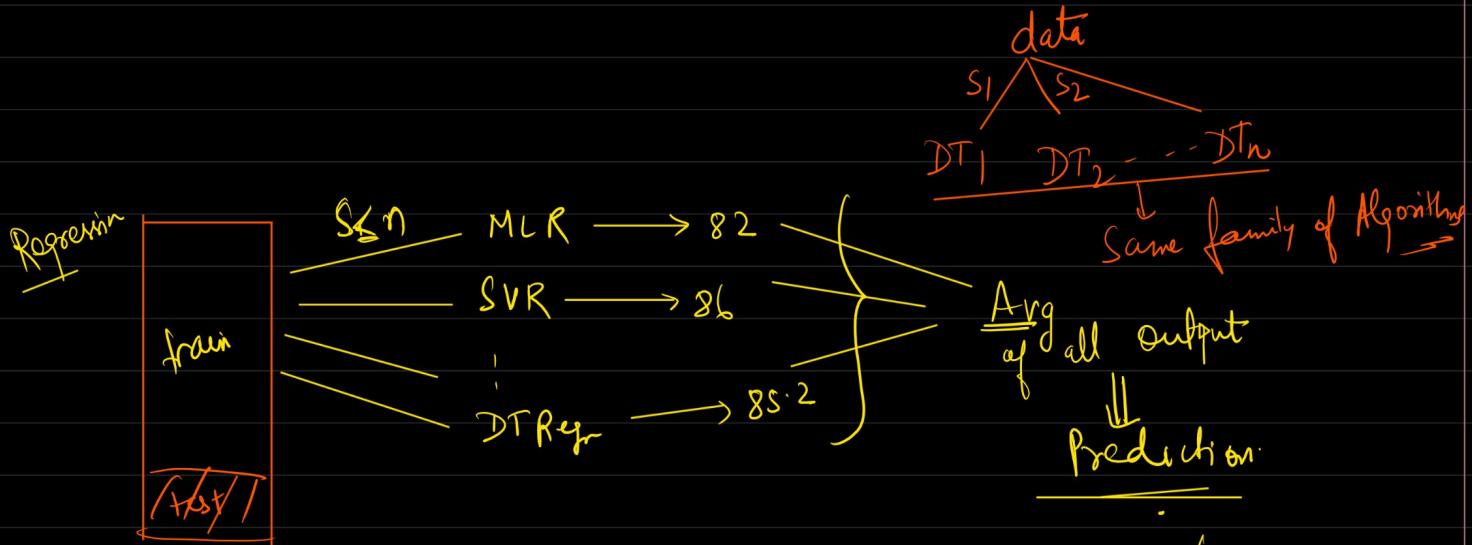
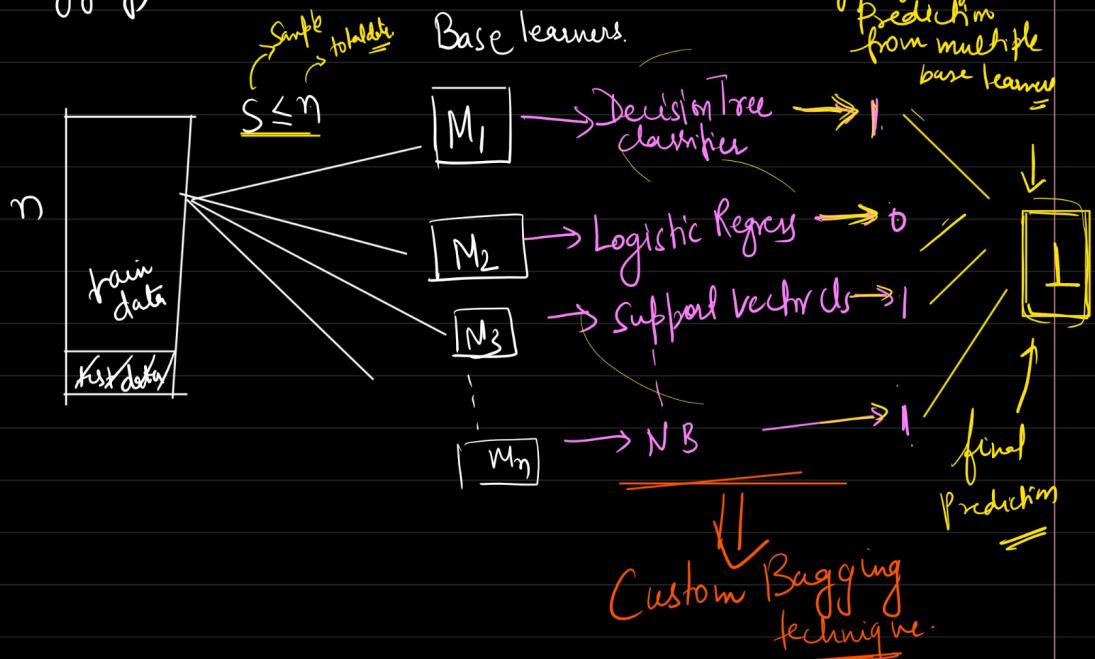


# Ensemble models

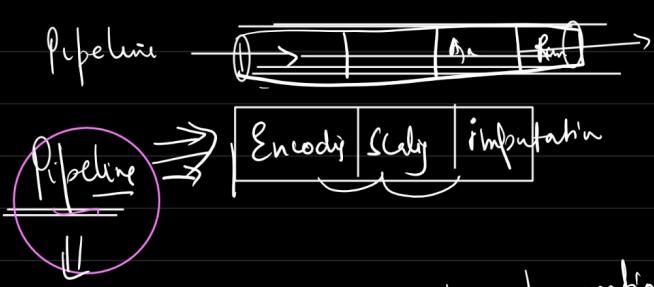


\* Bagging (custom Bagging)

\* Parallel models



S-1 missing value treatment  
S-2 One hot encoding  
S-3 Scaling.  
\* Pipeline

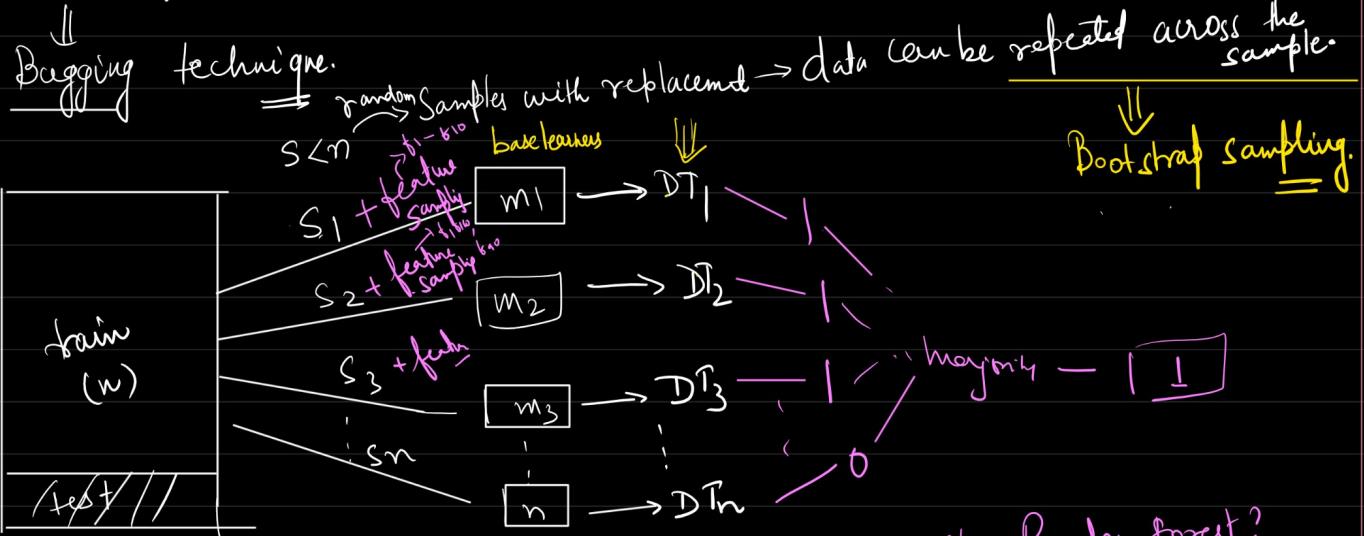


→ A sequence of data transformation  
 $f_1 \ f_2 \ f_3 \ \dots \ f_{100}$   
 Set of all features  
 ↓  
Column Transform.

Group all the pipeline steps for each of the column

\* Random Forest classifier and regressor

Random forest



\* Why Random forest?

- Multiple DT's in parallel
- row & features will be randomly selected.

Bagging ?

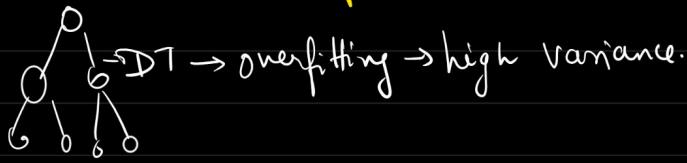
Bootstrap Aggregating

different samples with replacement

Bootstrap Sample

Aggregating the prediction

\* Random forest reduces the Variance



## Random forest

↳ random sampling of rows & features

data is splitted into small chunks (randomly)

due to random subset of rows  
and features, each of the  
subset is different representation  
of itself

→ Each of random subset training a diff. model

Reduces the Variance

## Scen-1

high var

feature | rows

Scen-2

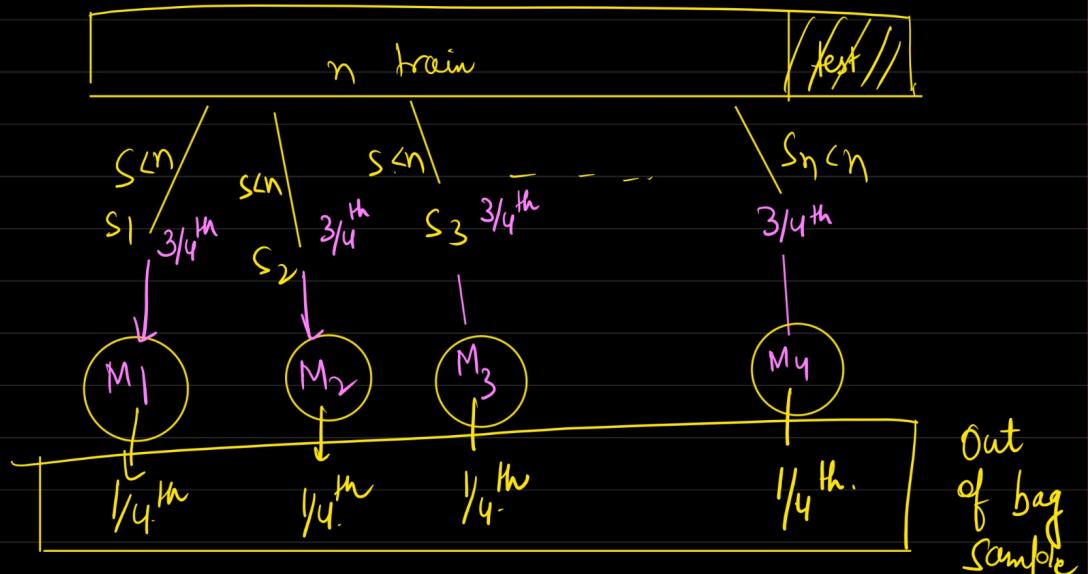
random

Sample

of rows and features

Final torchich no of trees  
overfitting.

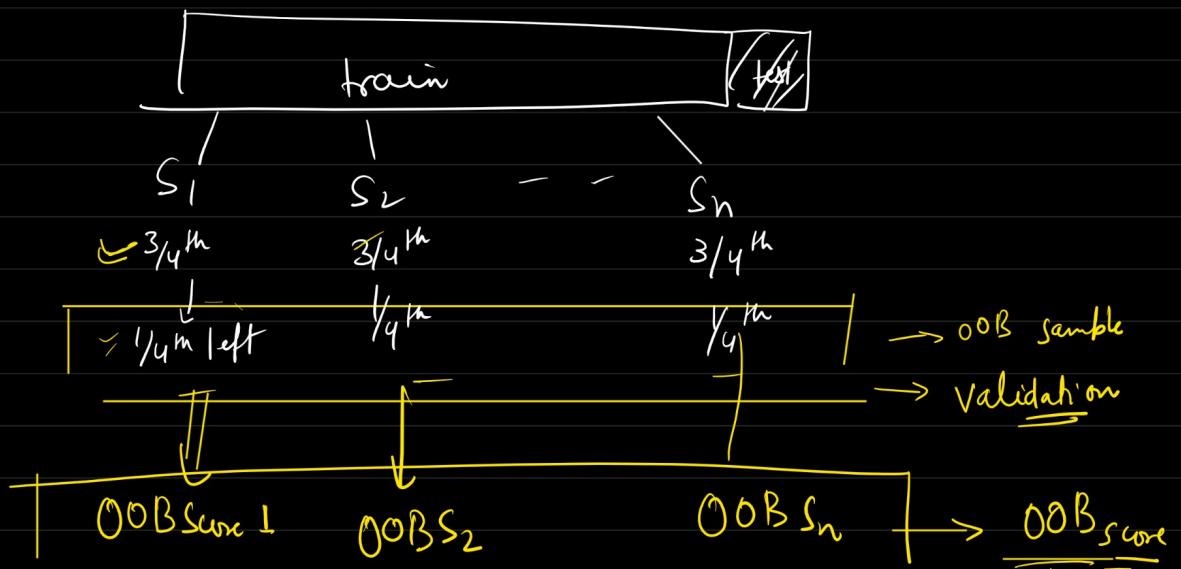
## OOB Score



→  $1/4^{\text{th}}$  acts as validation for each specific Individual DT.

Part of train  
data not used  
in model training  
for individual  
DT:

## OOB score



→ In training itself  
OOB is low, model  
is not performing well  
on train data.