

# Ensembling Techniques

Till now we have learnt to use only one ML model at a time

Now can we make a model integrating more than one ML-algo or using one model to create multiple ml model.

Yes we can create

→ But why we need that?

\* Why we want multiple models?

Analogy behind it → I go to teacher asked i want a job in data science which  
Now here maximum teachers (3 of 5) saying  $t_1$  (Gen AI)  $t_2$  (ML)  $t_3$  (ML)  $t_4$  (data analysis)  
ML, saying  $t_5$  (MC)  
On, i choose as subfield  
→ Then ML is more confident field.

Now similar analogy → we use similar analogy and train multiple models.

$\{m_1, m_2, m_3, m_4, m_5, m_6\}$  → combined prediction and chose one that most of model predicted by this cost of error ↴

- Ensembles:
    - Combine multiple models
    - Prediction which is more stable and accurate
- Compared to individual models

either we combine multiple model

It can be of same type / different

of some algorithm

of different Algo

$DT_1$  (maxdepth = 3)

$DT_2$  (maxdepth = 5)

$DT_3$  (maxdepth = 7)

Logistic Regression ( $M_1$ )

S.V.M.C ( $M_2$ )

Decision Tree Classif. ( $M_3$ )

\* ensemble → not necessarily only one type of algorithm

Types of ensemble techniques →

Parallel  
Technique  
(Bagging)

↓  
Sequential  
Technique  
(Boosting)

→ Parallel + sequential

technique

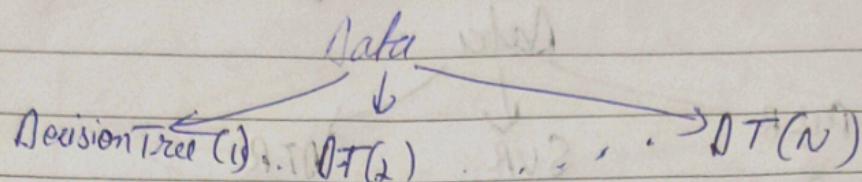
(stacking)

↓  
can be used multiple  
types of model

one type of model  
can be used

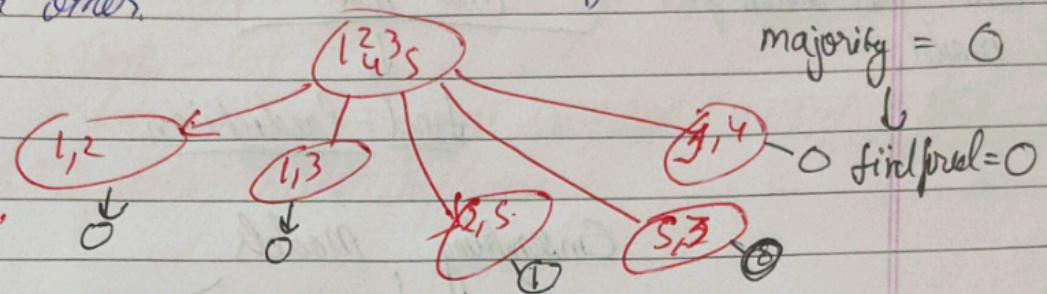
(1)

Parallel Technique of ensemble  $\rightarrow$  Bagging



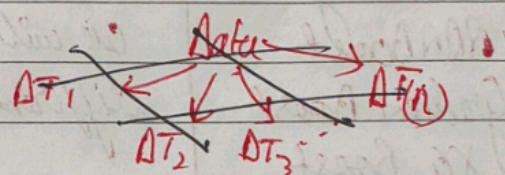
$\rightarrow$  All the model are built parallelly & independent of each other.

Data is divided in such a way all get different sets.

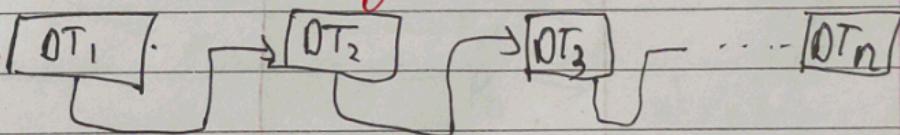


(2)

Sequential technique of ensemble  $\rightarrow$  Boosting



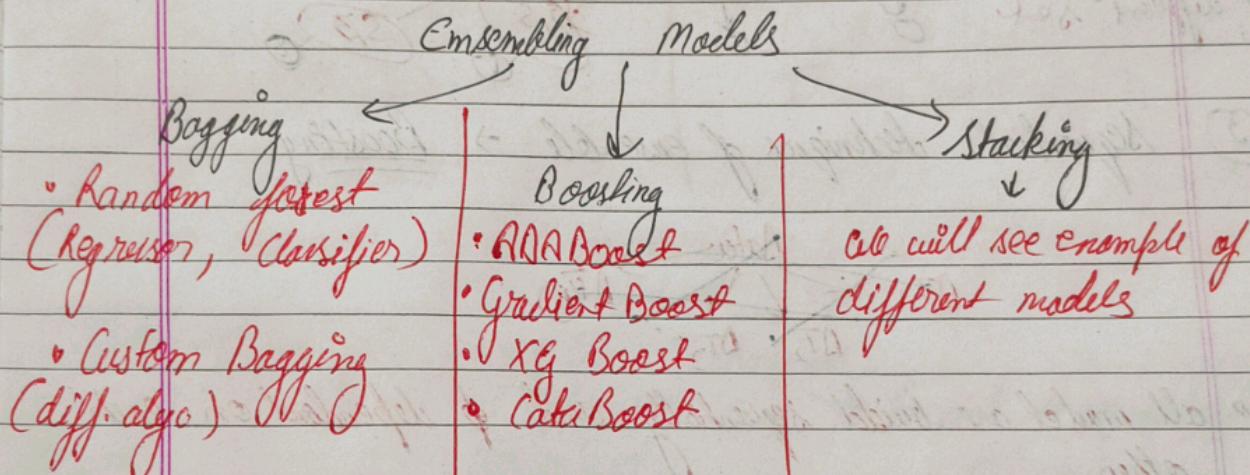
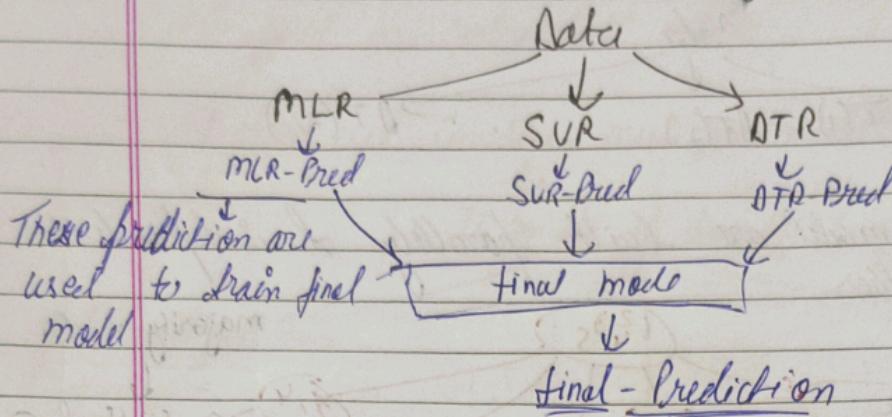
all model are build sequentially and  $\neq$  dependent on each other.



learning from mistake of previous model.

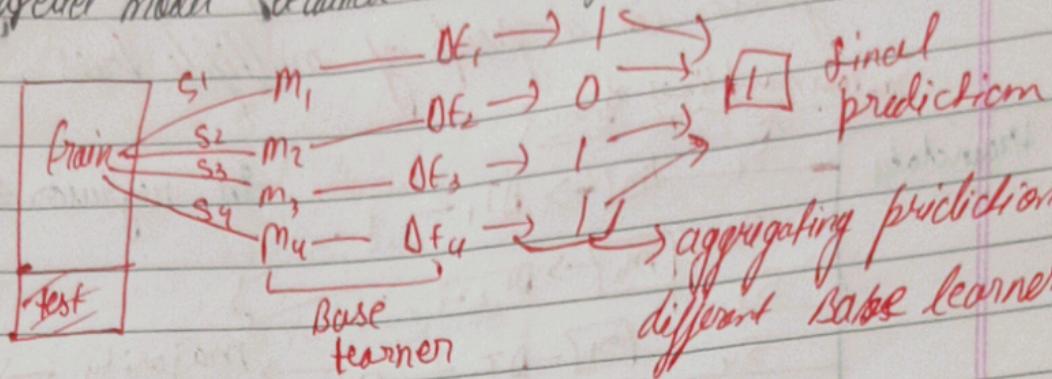
$\rightarrow$   $1 \rightarrow DT_1 \rightarrow$  trained  $\rightarrow$  then the dp's misclassified are trained on  $DT_2 \dots$  so upto  $n^{\text{th}}$  model

\* Parallel + Sequential  $\rightarrow$  (Stacking)



## \* Bagging technique (with classification problem)

→ Parallel model trained → used for Bagging model

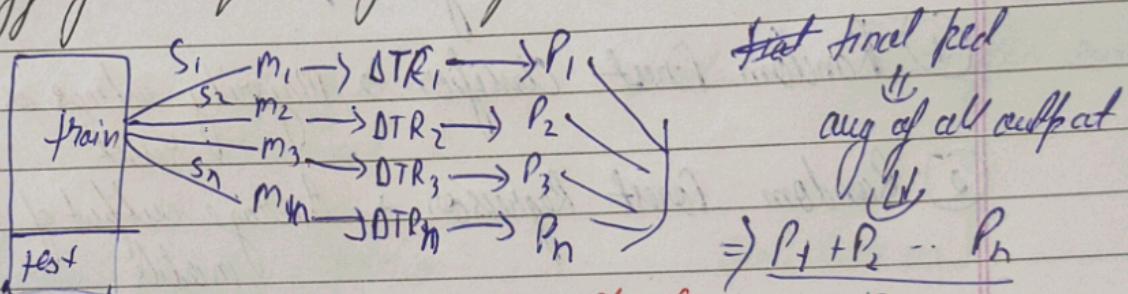


aggregating prediction of  
different base learner.

\* Data → Train - test  
    ↳ Samples → three samples (subset) for  
Each model is taken with replacement

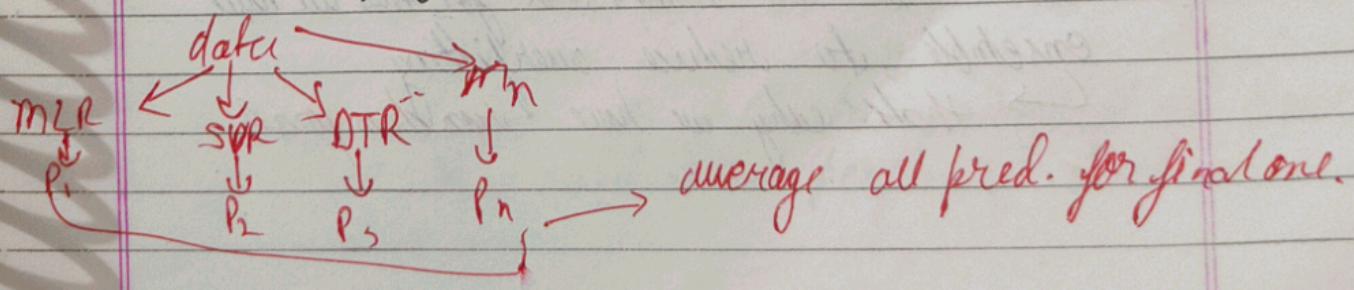
~~Bagging~~ → aggragation  
Bootstrapping → diff. samples with replacement.

\* Bagging technique for regression



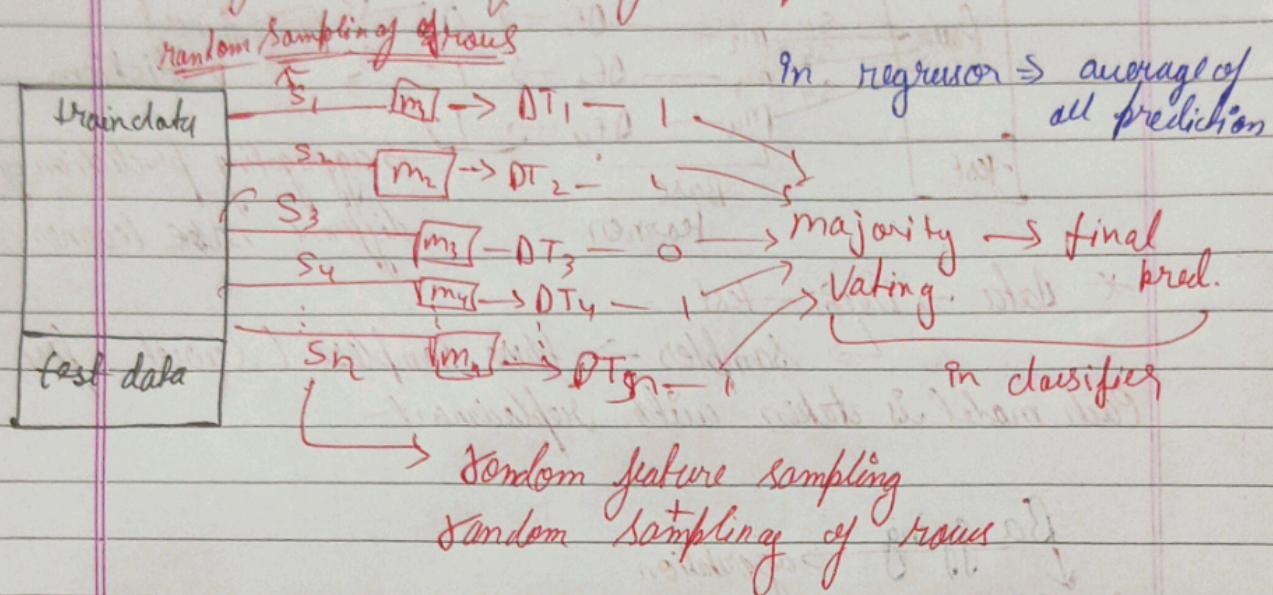
# every thing is same as of Classification problem

\* Custom Bagging technique / pseudo Bagging tech.



# Random Forest Classifier and Regressor

Here forest is a group of multiple trees



Random Forest  $\rightarrow$  multiple decision tree models in parallel  
 $\rightarrow$  the rows and features will be randomly sampled

① Random Forest Classifier  $\rightarrow$  Majority voting as predicted result

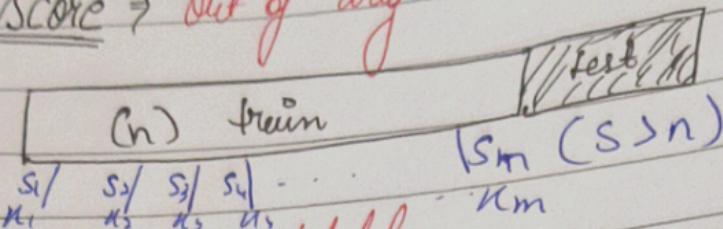
② Random Forest Regressor  $\rightarrow$  Average output of all the models

Why Random Forest if we have Decision Tree  $\rightarrow$  ?

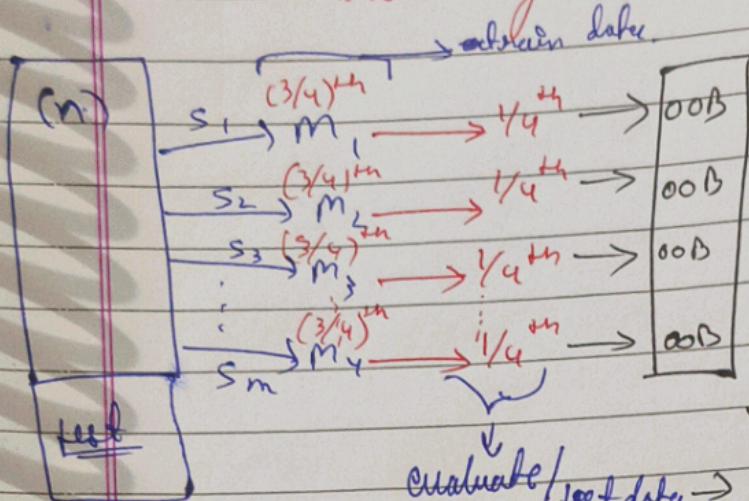
$\hookrightarrow$  DT is greedy algo, it will keep splitting until all are memorized. for this we have ensemble to reduce overfitting

$\hookrightarrow$  that's why we have random forest

OOB Score  $\rightarrow$  out of bag score



Samples: randomly selected rows and features (with replacement)



this  $1/4$  data is nothing but  
evaluate / test data  $\rightarrow$  out of bag sample

OOB sample  $\rightarrow$  part of train data, not used in model  
training for individual models/decision tree.

$\rightarrow$  now this act as an validation data for each  
specific individual AT.

$\rightarrow \{ \text{OOB score}_1, \text{OOB score}_2, \dots, \text{OOB score}_m \}$

$\rightarrow$  average of all OOB score  $\Rightarrow \frac{\text{OOB score}}{m}$

act as a validation score for  
random forest.

if in training itself if oob score  
is low, model is not performing well on  
train data