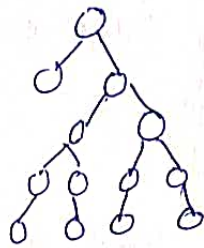


13) 11/22.

Bagging and Boosting

1. Random forest classifiers
2. " " regression

Decision tree :-



Overfitting :- low bias + high variance..

Pre Pruning

Post Pruning

In order to avoid that we use

while constructing the decision tree, we use specific hyperparameters like max-depth, min-split, max-features.

first of all we will construct the decision tree & we will cut the branches.

Bagging & Boosting :-

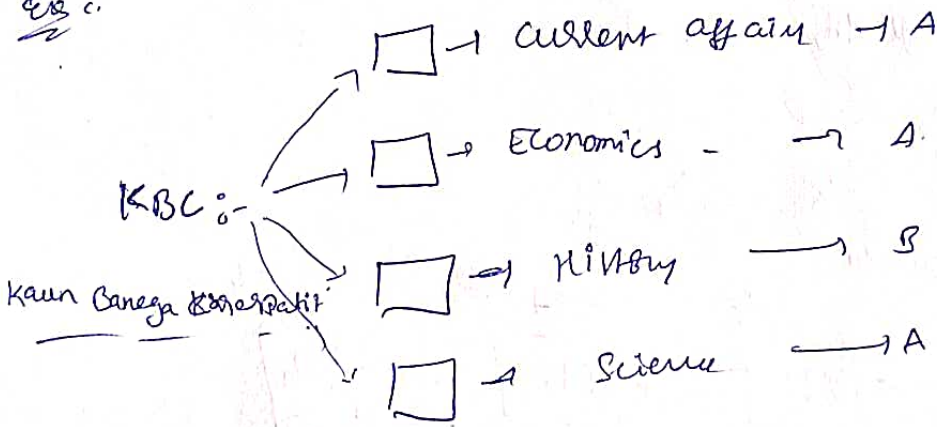
KBC :- (Amitabh Bachchan)

KRISH -> Data science

UPSC -> Diversified



Ex c.



'A'

Majority voting.

Bagging :- using ML algorithms parallel each other, the final O/P, we will

[Bootstrap Aggregation] consider majority voting classifiers.

ML1 -> DT -> 1 Ensemble Techniques.

ML2 -> logistic -> 0

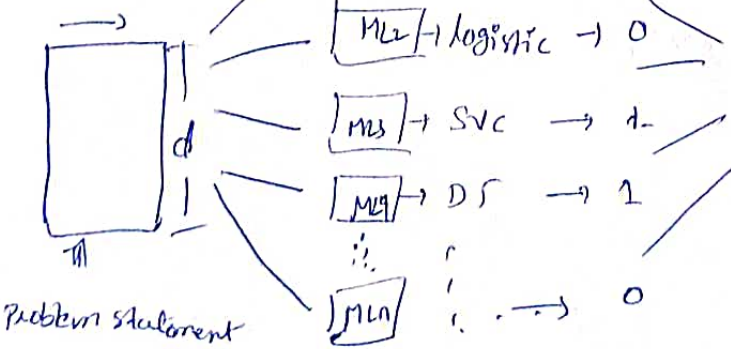
ML3 -> SVC -> 1

ML4 -> DS -> 1

MLn -> . -> 0

here '1' is m_{aj}.

bagging of algorithm.



Ensemble Techniques, when this technique came

||
Kaggle, Hacker Rank

||
Outperforming results.

bagging.

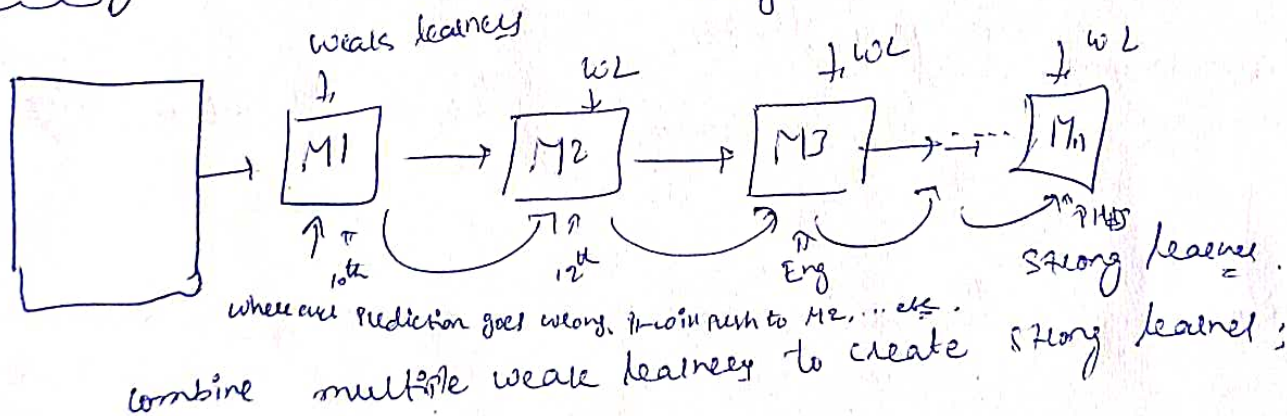
in regression we will take avg/
mean.

in classification :- Majority voting

Bagging :-

1. Random forest classifiers
2. Random forest regress.

Boosting :- we create models sequentially.



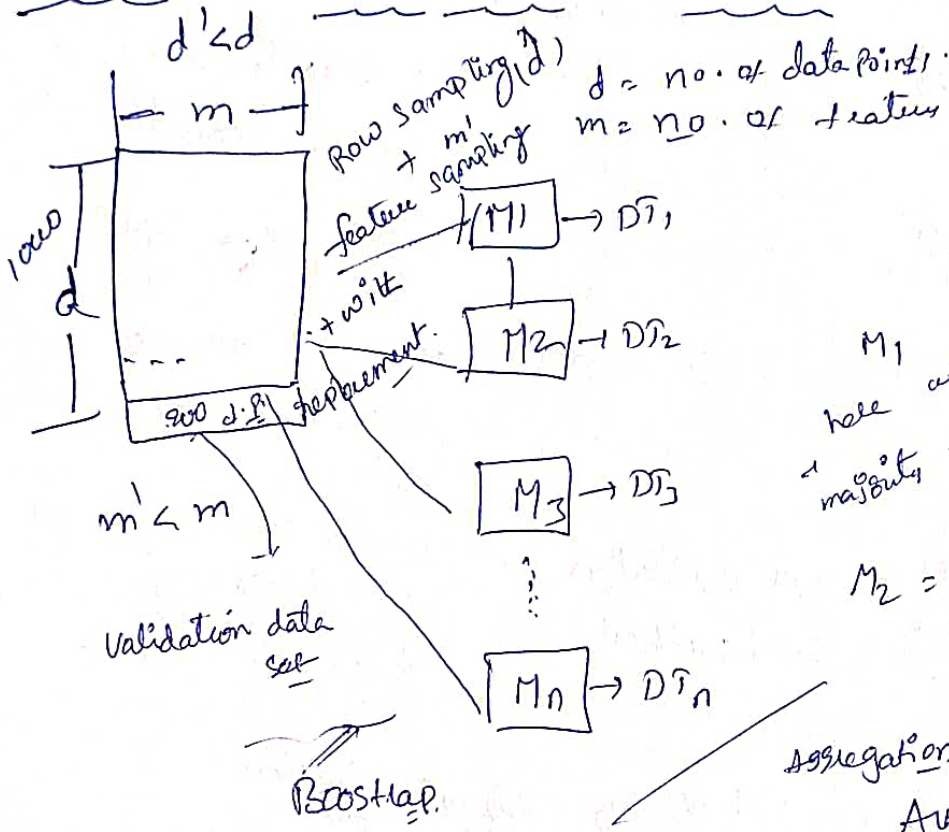
1. Adaboost Regress & classifier

2. Gradient-boost

3. XGBoost

4. Xstream Gradient boost

Random forest classification & Regression Ensemble techniques.



Combine many models
combine multiple
 DT'_1

M_1 we consider
"majority voting"

Row sampling is
Sample of row's

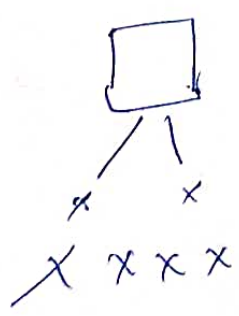
$M_2 = \text{Row sampling} + \text{feature sampling}$

Aggregation.
Average in case of regression.

New test data

$DT_1 \rightarrow \text{Overfitting} \rightarrow \text{Pruning (or) Post Pruning}$

1. low bias \rightarrow Training
 2. High variance \rightarrow here we need to get low variance.
- using pruning/post pruning process.

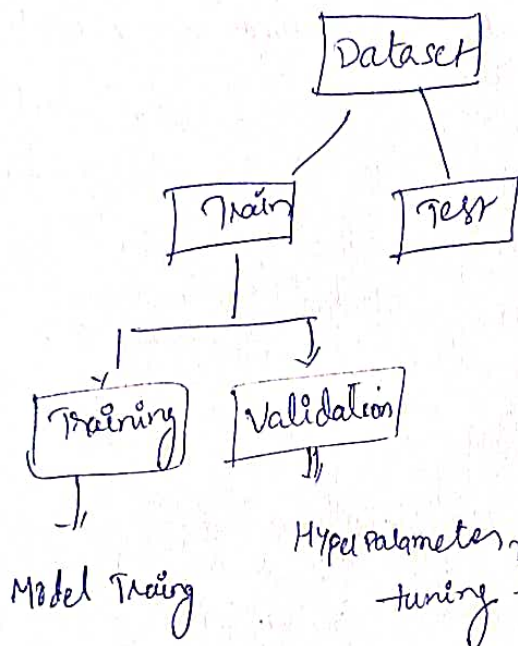


For huge data sets we can't see the DT properly.
generalised models.

classification DT's will get trained using Entropy, gini.

Regression : DT's " " by average, mean, threshold, MSE

Out of Bag Evaluation:



OOB \rightarrow There will be a scenario, some of the data points will never get selected.

those called as validation data points. It will produce some score, it's called

OOB score \approx TRUE

In random forest we can see

OOB score = TRUE

use this validation data

OOB score = 83%

out of Bag Error \Rightarrow Validation

$1 - \text{OOB score} = 1 - 83 = 0.17\%$

Assignment :-

Use bagging classifier and regressor, Voting classifier, regression and random forest classifiers, regressor on household consumption data + census data