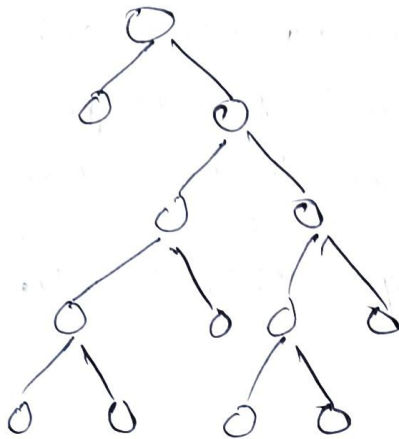


13-Nov-22

Decision Tree



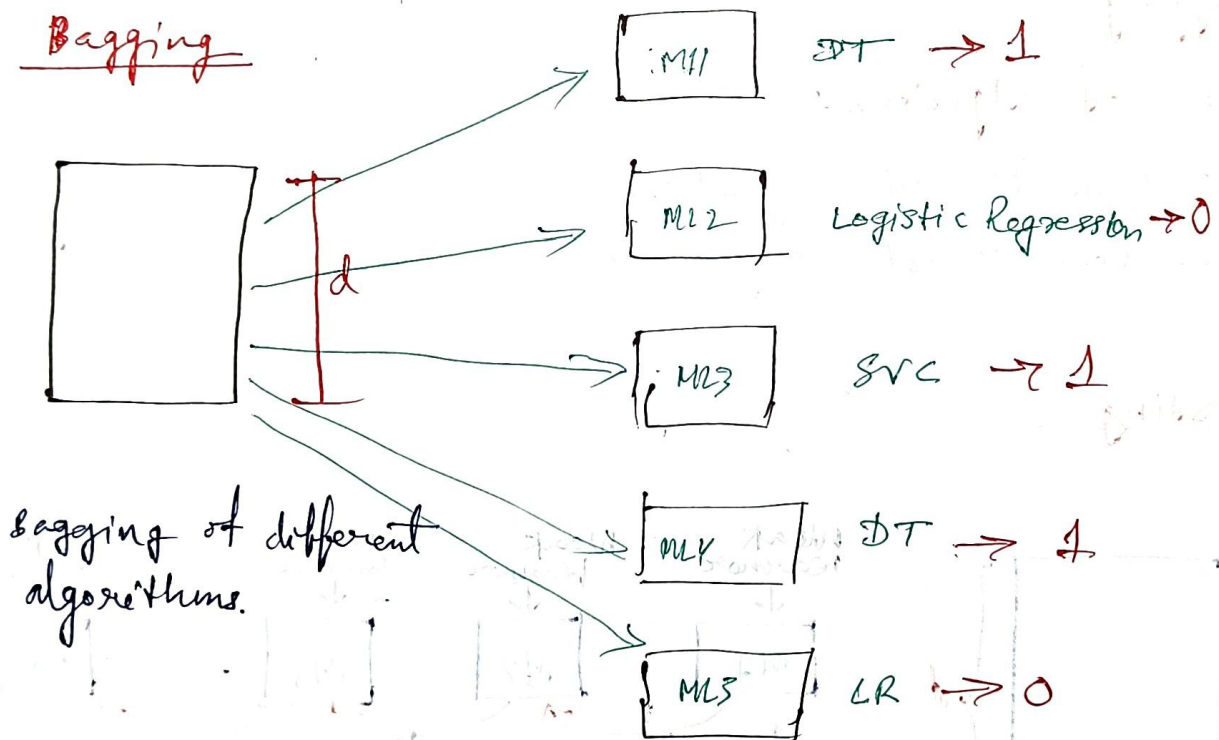
overfitting issue

Low Bias
High Variance

- Using Pre Pruning or Post Pruning we can avoid overfitting issue.
- Using bagging techniques also we can remove overfitting problem.

Bagging and Boosting

Bagging



> Final output will be based on majority voting classifier. Here it will be '1'.

> $\left\{ \begin{array}{l} \text{Bootstrap} \\ \downarrow \\ \text{dividing into diff algorithms} \end{array} \right\} \text{Aggregation} \left\{ \begin{array}{l} \downarrow \\ \text{aggregate to get final result} \end{array} \right\}$ Ensemble Technique

> Each algorithms can learn something unique from the problem statement.

Example:

KBC competition,

Subhash \rightarrow Data science domain knowledge.

X (UPSC) \rightarrow Different domain knowledge

So X (UPSC) candidate perform better than me at KBC competition. because he/she has different domain knowledge. This is called

Bagging.

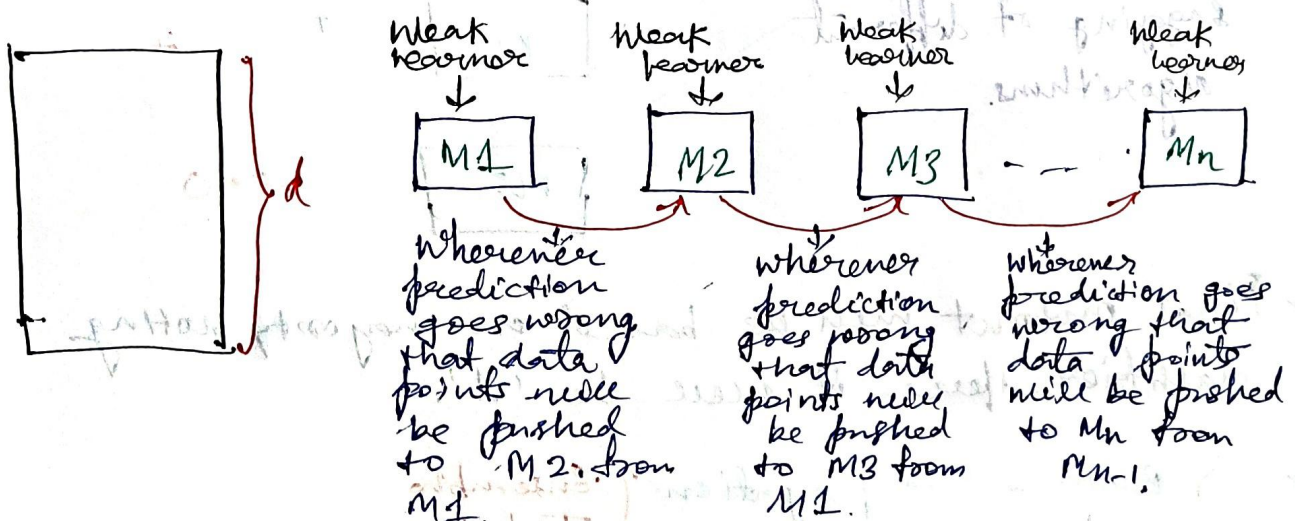
\rightarrow Most of the time, Ensemble Techniques is a good choice.

Bagging

Types of Algorithms

- ① Random Forest classifier
- ② Random Forest Regressor

Boosting



\rightarrow After combining all the models, it will be strong learner.

Types of Algorithms

- ① Adaboost Regressor and classifier
- ② Gradient Boost Regressor and classifier
- ③ Xgboost Regressor and classifier

↓
Extreme Gradient Boost

Random Forest Classification and Regression

→ It is an Ensemble Techniques.

↓
combination of many models.



→ In sklearn, we combine multiple decision trees

$$d' < d$$

$$M < M$$

For each model

{ Row sampling (d') + Feature Sampling + with Replacement (M) }

Q Why are we using Random Forest?

Ans

Pre-pruning and Post-pruning is a bigger task for large dataset. So, to avoid overfitting we use Random Forest. For large dataset, we cannot see complete tree.

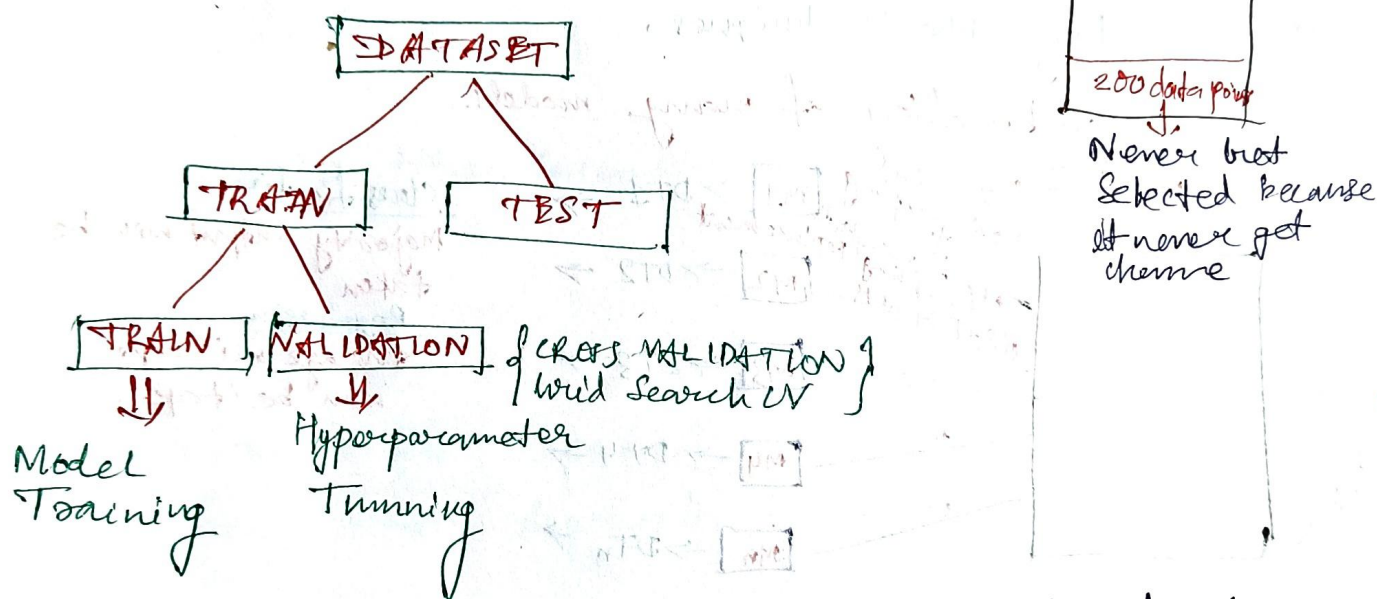
overfitting

Low Bias connect using Random Forest
High Variance → Low Variance

→ Using Random Forest, we can convert high variance to low variance.

Out of Bag Evaluation

- while doing Row sampling, feature sampling and with replacement for each model, so there is chance that some set of values never get selected.
- we can use unselected data points for testing and validation purpose.



- > If you set oob = true, then it will automatically handle missing data points for testing purpose.

7- out of

Out of Bag Error (Validation Error)

$$\text{Out of Bag Error} = 1 - \text{Out of Bag Score}$$

- oob score is the accuracy ^{score} over the validation data.

WRT Regression

- > R-2 score can be used.

- > We decide algorithm based on the dataset.
- > Random forest has more time and complexities.
- > Bagging can be heterogeneous as well as homogeneous.

Hyperparameters

$n_estimators$ = number of Decision Trees

b bootstrap = combination of models

oob-score = accuracy score, need test data for testing purpose.

n_jobs = How much CPU our model will take while running
 $verbose$ = can see messages about configuration while running
 $verbose$ = can see messages about configuration while running

Boosting classification

- We can use different algorithms.

