

Ensemble - Collection / grp of things

Used in combination when we have multiple models being used together to build a more powerful model.

There are 4 types of models -

- i) Bagging (Bootstrapped Aggregation)
- ii) Boosting
- iii) Stacking
- iv) Cascading

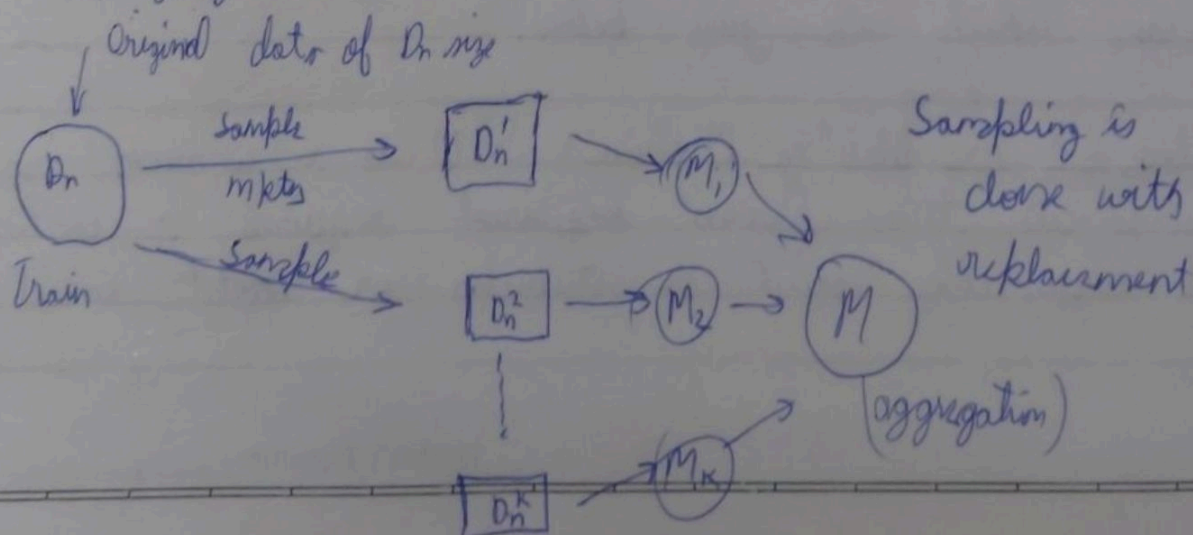
* Key aspect -

$M_1, M_2, M_3 \dots M_K$

The more different these models are,
the better you can combine them.

Bagging (Bootstrap Aggregation)

Popular Bagging model - Random Forest.



Each model M_i is built using D_n^i of size m ($m \leq n$)

\Rightarrow each model M_i has seen a different ~~subset~~ ^{sample} of data.

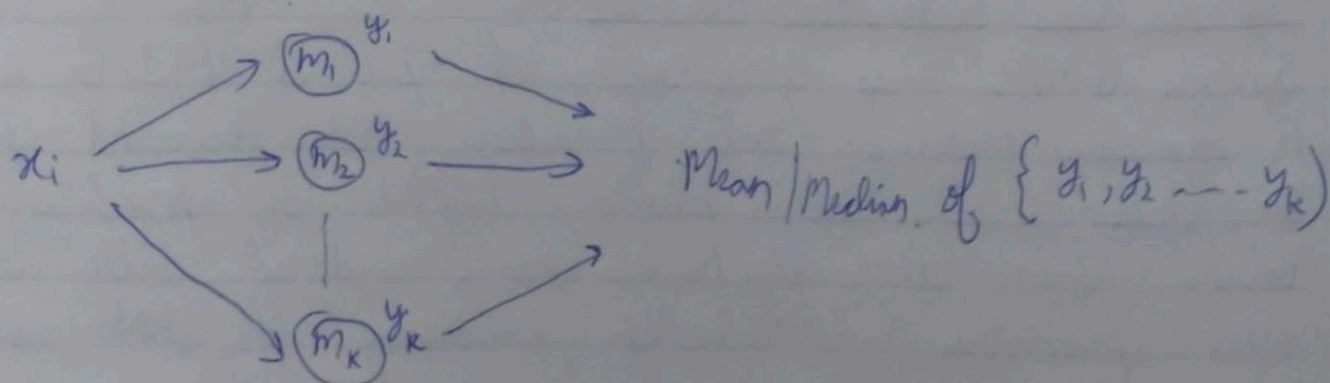
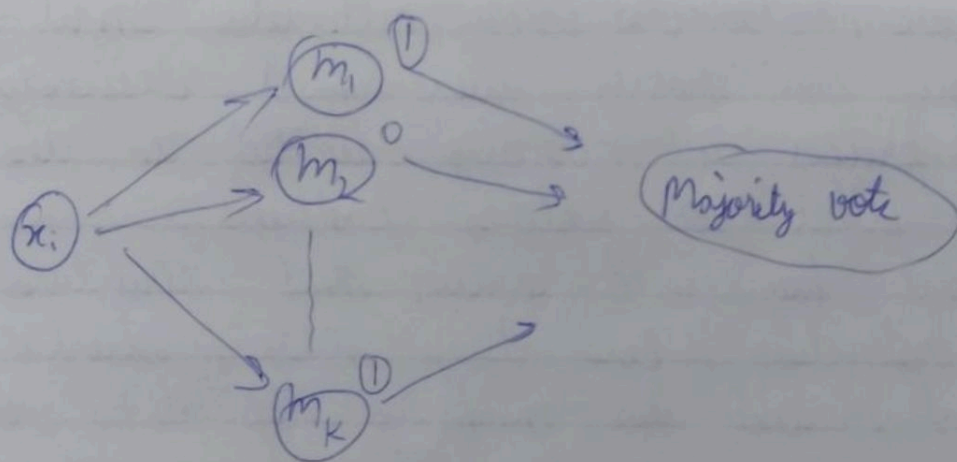
So, each model is Aggregated to create a one powerful model.

Aggregation :-

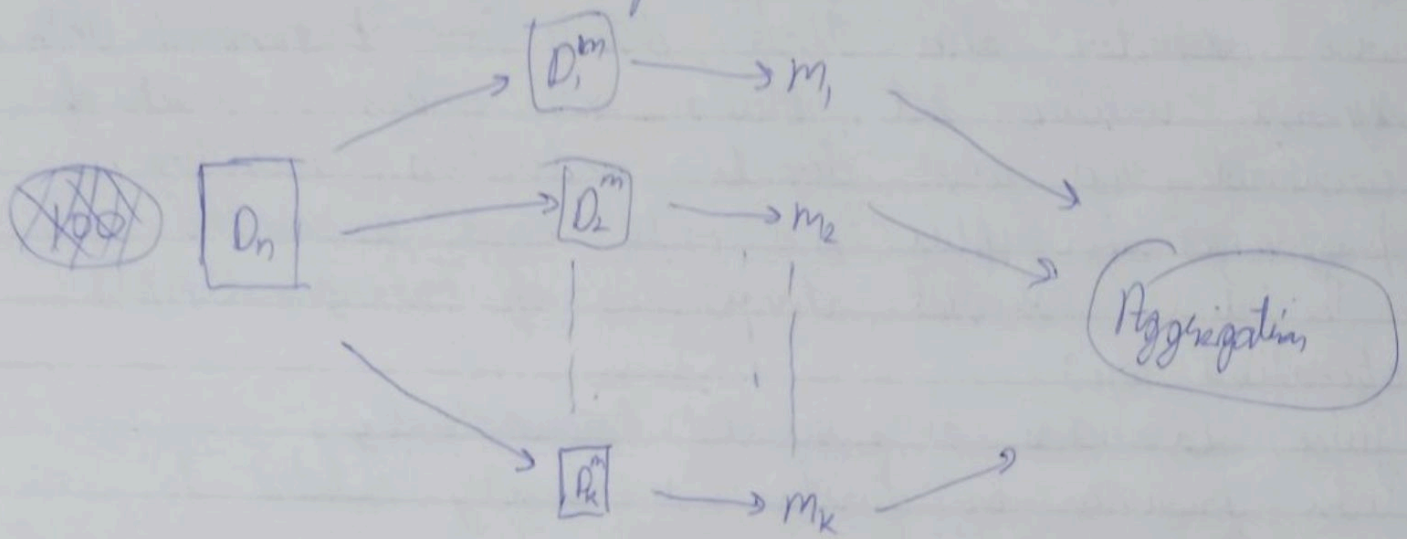
Classification :- Majority vote

Regression :- Mean/Median

So if we get a new datapoint x_i , it will go to all the base models.



If model changes a lot with ~~the~~ changes in training data it is said to have high variance.



Now, suppose I change or just remove 100 pts from my ~~training~~ training data so due to this out of K samples ~~only~~ only small no. of samples would be impacted due to which only small no. of Base models would be impacted due to which the overall aggregated model (with majority voting or mean) won't be impacted much, which shows that our final model won't change much. Hence, from here we can understand that ~~Bagging~~ with change in training data i.e., it has low variance.

Hence from here we can understand that Bagging can reduce the variance of a model without impacting the bias.

$$\text{model error} = \text{Bias}^2 + \text{Var}$$

Now, if we have a base model with low bias & high variance then using bagging on such multiple models would provide a final model with low bias and reduced variance.

↓
Core Idea of Bagging

Random Forest (Bagging)

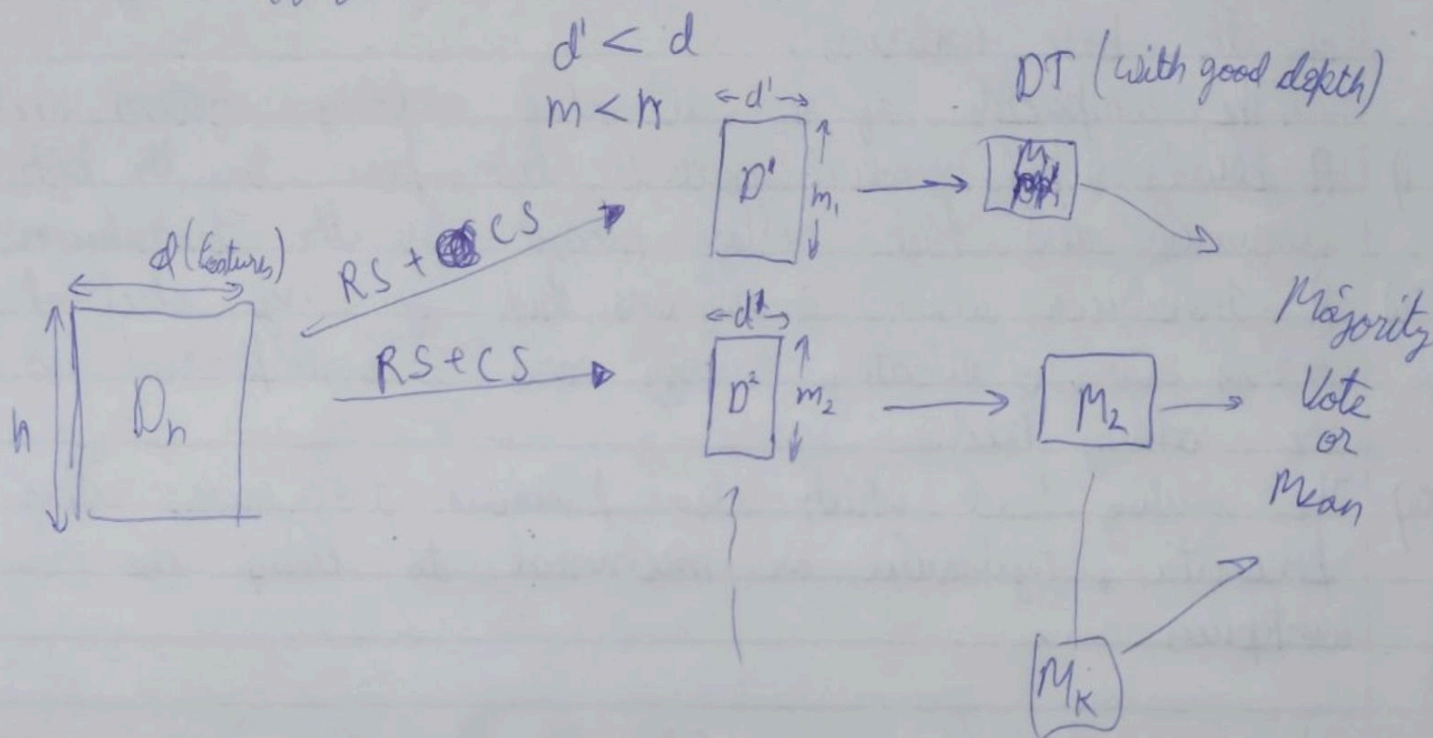
Several Decision Trees

Random Bootstrap Sampling

RF :- Decision Trees + Bagging + Col. Sampling
Base model Feature Bagging

Bootstrap Sampling:- Row Sampling

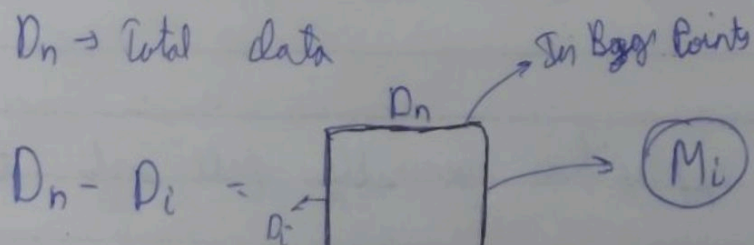
Feature Bagging:- Column Sampling



So considering my first base model as M_i

M_i gets trained on D' data having d' cols and m rows

$D_n \rightarrow$ Total data



Out of Bag points (OOB)

(can be used for Cross Validation of M_i)

RF : (DT) ^{Having reasonable depth} base learner

- + Row Sampling with replacement
- + Col. Sampling
- + Aggregation (Majority Vote / Mean or Median)

Bias - Variance Tradeoff

RF :- Reduce Variance

→ Low Bias Because base learners (M_i 's) are low-bias.
→ Variance

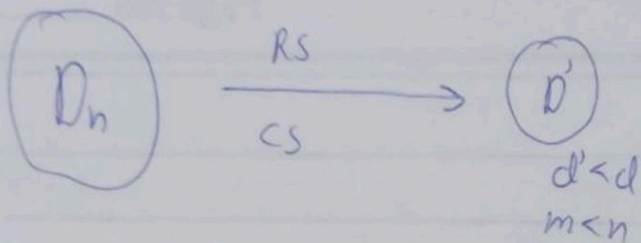
$$\overset{\text{final model}}{M} = \text{Agg}(M_1, M_2, M_3, \dots, M_K)$$

So, $K \uparrow$, Variance \downarrow

$K \downarrow$, Variance \uparrow

Bias of $M \approx$ Bias of M_i
 \uparrow \downarrow
Final model Base model

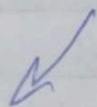
Other than K (No. of Base Learners) there is one more hyperparameter on which the variance depends :-



$\left(\frac{d'}{d} = \text{Col. Sampling Ratio} \right)$
 $\left(\frac{m}{n} = \text{Row Sampling Ratio} \right)$

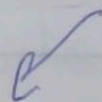
If Col. Sampling Ratio ↓

Row Sampling Ratio ↓



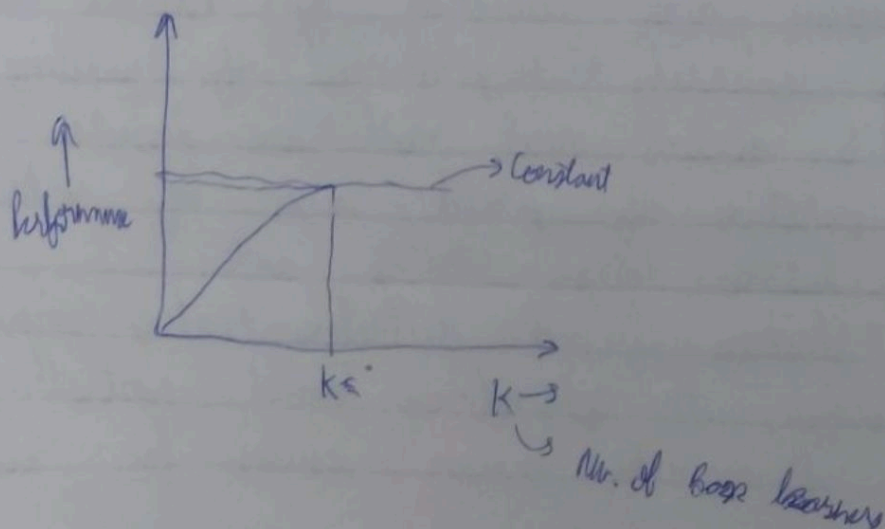
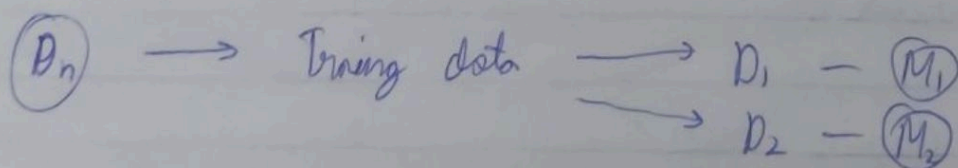
Each sample would be more different from others

→ ~~Change~~ Impact of change in training data would ↓



low variance

But we can't keep these ratios too small that we don't have any data left to train our model.



Train + Run Complexity

RF with ~~and~~ K base learners (DT)

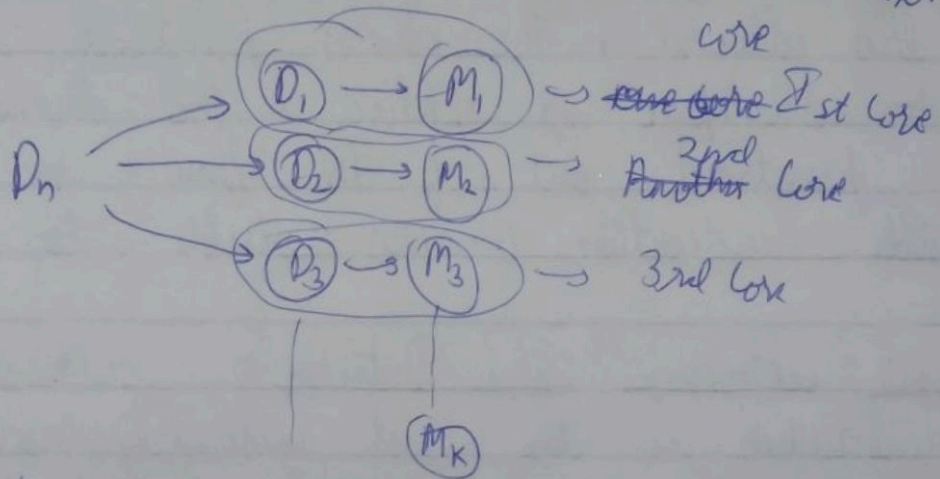
Train Time Complexity:- $O(n \log(n) * K)$

If we have a multi core devices then ~~then~~ each base model can be trained in ~~the~~ separate core

That's why it is also called

Trivially Parallelizable

Multiple models can be trained parallelly



Run time Complexity:- $O(\text{depth} * K)$

Space Complexity:- $O(\text{DT} * K)$

Random Forest Hyperparameters:-

$n_estimators = 10$ // no. of base learners

$max_depth = None$ // Putting no limit of depth.

$bootstrap = True$ // To ~~get~~ ~~take~~ bootstrap samples.

$oob_score = \text{False}$ // Will give Cross Validation Score of each DT if becomes True.

$n_jobs = 1$ // Scikit-learn will train each DT on all CPU core parallelly if becomes -1.

Extremely Randomized Trees

It is almost similar to RF except one more Randomization factor.

In case of RF, to check threshold for a column, all the values are first sorted and checked.

But in extremely randomized trees, only certain random samples are taken to check the threshold to split node of tree.

Extreme Trees = Col. Sampling + Row Sampling + Agg
+ Randomization while selecting threshold

↓
This led to more uniqueness in Base models.

↙
Reduce Variance Better than RF.

8.) Random Forest Cons

RF :- (DT) + RS + CS + Agg.

Not good in large dimensionality data,
Categorical features with large no. of categories

We have seen that almost all the cons that hold for DT also hold for RF except some.

Bias Variance Tradeoff

In case of DT, $K=1$ is the major factor of Bias Variance Tradeoff.

In case of RF, K is the factor of Bias Var. tradeoff

Boosting Intuition

In Bagging we do,

Reduce Var.

Bagging

high var + low bias + Randomization + Aggregation
Base models (LS, RS)

Boosting

Low var + High Bias

+

Additionally combine

{ Reduce Bias while
keeping var variance low }

Core idea on how actually boosting reduces bias.

In Stage 0, $D_{\text{train}} \rightarrow M_0$ ^{DT (Example only)}
High Bias & Low Variance (Shallow Tree)

A high bias & low-variance model will be created from D_{train} (whole data) having high bias & low variance.

$$-\log_{10}(2.3 \times 10^{-11})$$

High Bias :- Large training error

Stages

① a) $D_{\text{train}} \rightarrow M_0$
 $(x_i, y_i)_{i=1}^n$ $\left(y_i - h_0(x) \right)$
 \rightarrow Predictions

b) $y_i - h_0(x) \Rightarrow$ Errors

Now at the end of stage 0, we have $\{x_i, y_i, \text{error}_i\}_{i=1}^n$

① (M_1) , the model in stage 1 takes x_i as input and tries to predict error_i .

$(M_1) \leftarrow \{x_i, \text{error}_i\}$

$\downarrow \rightarrow$ Predicts

$h_1(x) \approx \text{error}_i (y_i - h_0(x))$

~~$h_1(x)$~~ tries to fit

$F_1(x)$ = model at the end of stage 1

$F_1(x) \approx \alpha_0 h_0(x) + \alpha_1 h_1(x)$

② $\{x_i, \text{error}_i\}$
 $\rightarrow y_i - F_1(x_i)$
 model at end of stage 2.

$F_2(x) \approx \alpha_0 h_0(x) + \alpha_2 h_2(x) + \alpha_1 h_1(x)$

At the end of stage K :-

$$F_K(x) = \sum_{i=0}^K \alpha_i h_i(x)$$

additive
weighted
model.

trained to fit the residual
error at the end of previous
stage.

$$F_K(x) = \sum_{i=0}^K \alpha_i h_i(x)$$

Final model ends up having low residual error.

so, at the end it is reducing Training error

Train error $\downarrow \rightarrow$ Bias \downarrow

Regularization + Shrinkage

$$F_M(x) = h_0(x) + \sum_{m=1}^M \gamma_m h_m(x)$$

$M \rightarrow$ # of base models

As $M \uparrow \Rightarrow$ overfit $\uparrow \Rightarrow$ variance \uparrow

As # Base model \uparrow , Bias \downarrow but variance \uparrow .

Shrinkage

$$F_m(x) = F_{m-1}(x) + \gamma \cdot h_m(x), \quad (0 < \gamma < 1)$$

If $\gamma = 1$, the eqⁿ becomes as previous

If γ is small, overfit \downarrow

~~For~~ To prevent overfitting we have 2 hyperparameters M and γ .

Train & Run Time Complexity

Train complexity :- $O(n \log n * M)$

RF is parallelizable but GBDT are sequential so it takes more time to train than RF.

Runtime

$$\text{GBDT} : O(\text{depth} * M)$$

Space complexity

$$O(\text{store each tree} + \gamma_m)$$

§ §

Scikit Learn GBDT \approx GBPT + Row sampling

PAGE 5

→ From RF

$$Xgboost = GBDT + (RS + CS)$$

booster = 'gltree' / 'glinear'

Both tree & linear model can be high bias & low variance model.

col sample by tree :- While constructing tree, check particular sample of columns

~~column sample~~

col sample by tree = 1, While constructing each node of a tree ~~only~~ only a random sample is considered.

reg - alpha = 0, L1

reg - lambda = 1, L2

$$F_m(x) = h_0(x) + \underbrace{(\alpha_1)}_{\text{weights}} h_1(x) + \alpha_2 h_2(x) + \dots$$

$$\alpha L1 + \lambda L2$$

↓
Regularizes

↓
Regularizes