

Lecture link : <https://www.scaler.com/meetings/i/dsml-advanced-ensemble-boosting-2/archive>

## Content

1. Recap (2:00-14:06)
2. Random Forest - introduction(14:06 - 26:00 )
3. Hyper parameter tuning (35:00 - 1:05:06)
4. SKlearn library (1:55:30 - 2:16:00)
5. Extra trees (2:16:00 - 2:29:22)
6. Summary (2:29:22 - 2:31:01)

### ▼ Recap

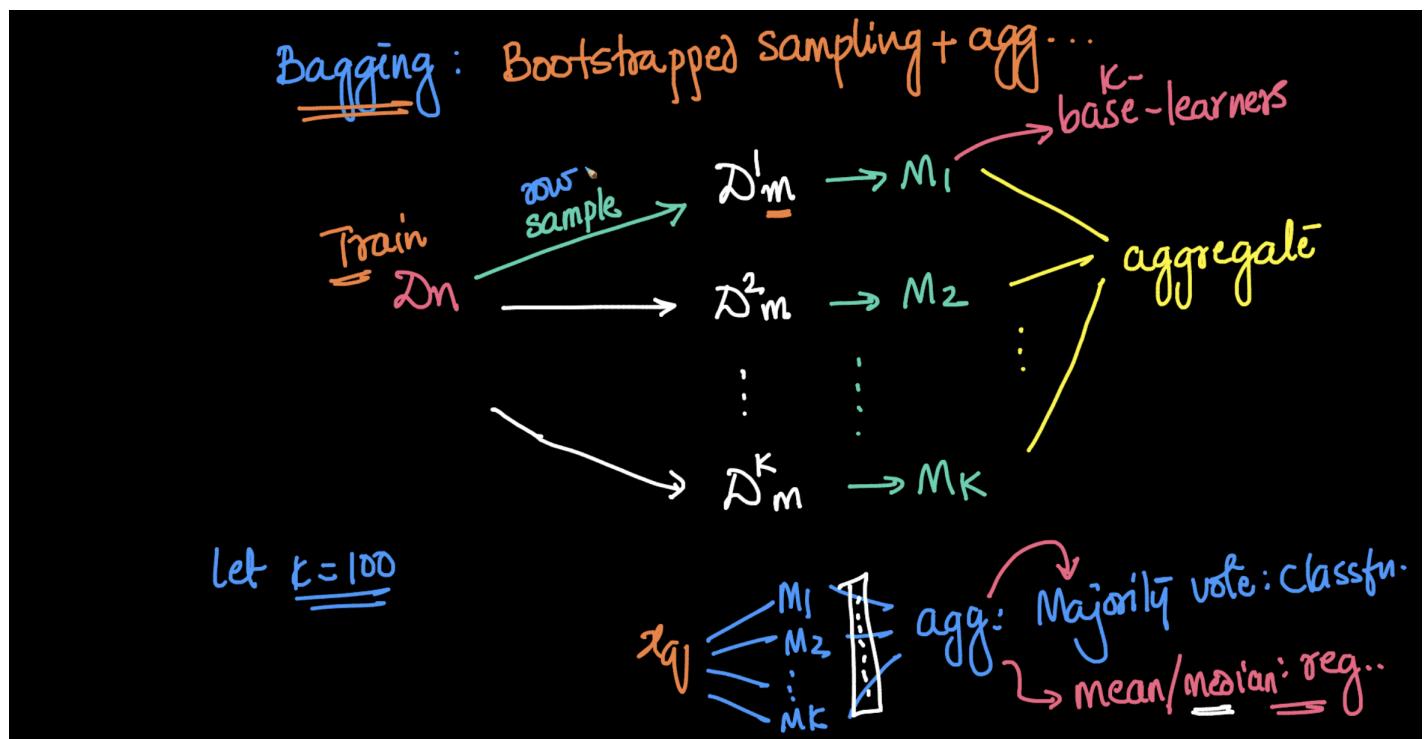
In previous lecture we saw what bagging is, lets take a quick recap of what it is and how it works

- Bagging is nothing but **Bootstrapped sampling + Aggregation**

#### How it works ?

Let us assume a train data set  $D$  with  $n$  data points i.e  $D_n$

- Now, we **sample  $m$  data points with replacement** to get  $D'_m$
- We do the sampling again for  $m$  points to get  $D'_2$ , **Repeat the same for  $k$  times** and we get  $D'_k$
- Now, we **train  $k$  different models**( $M_1, M_2, \dots, M_k$ ) basing on the  $k$  datasets obtained , there models are called **Base Learners**.
- After training we **cross validate each model with remaining  $n - m$  data points**
- Now, we do **Aggregation**
  - We use majority vote for Classification
  - We use Mean/Median for Regression



Let's now see what a Random forest is

## ▼ Random Forest

- As you can understand from the name a Random forest is nothing but a group of trees.
- It's basically a bagging based system with Decision tree
  - Base learners of Random Forest are Decision tree
- So, as a whole, Random forest is Applying Randomisation on the data and ensemble of Decision trees and later bagging these.

Let's see Random forest is built

- Let us assume we have a train data set ( $D_n$ ) with  $n$  rows and  $d$  columns
- And now we apply **row sampling and column sampling** on the data to get new data set ( $D_m^1$ ) with  $m$  rows and  $d'$  columns
  - where  $m < n$  and  $d' < d$

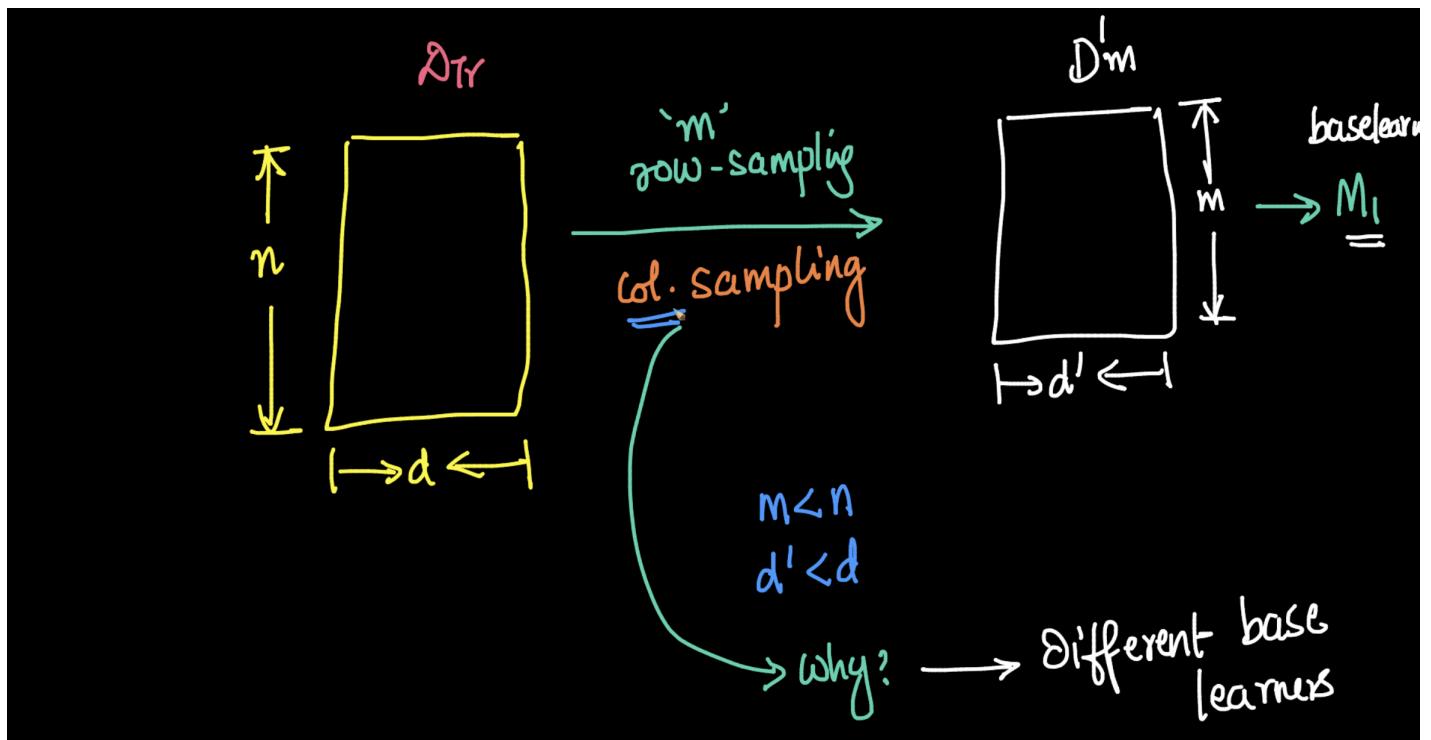
But why do we apply column sampling?

- Because we want all our **Base learners** to be different, choosing different set of features may help them to be different, for achieving this we do column sampling as well.

- Adding column sampling makes the models to be little more different

We do the same again and again to get  $k$  different data sets

- So we get  $D_m^1, D_m^2, \dots, D_m^k$  now the models  $M_1, M_2, \dots, M_k$  will be different
- Now we aggregate the models



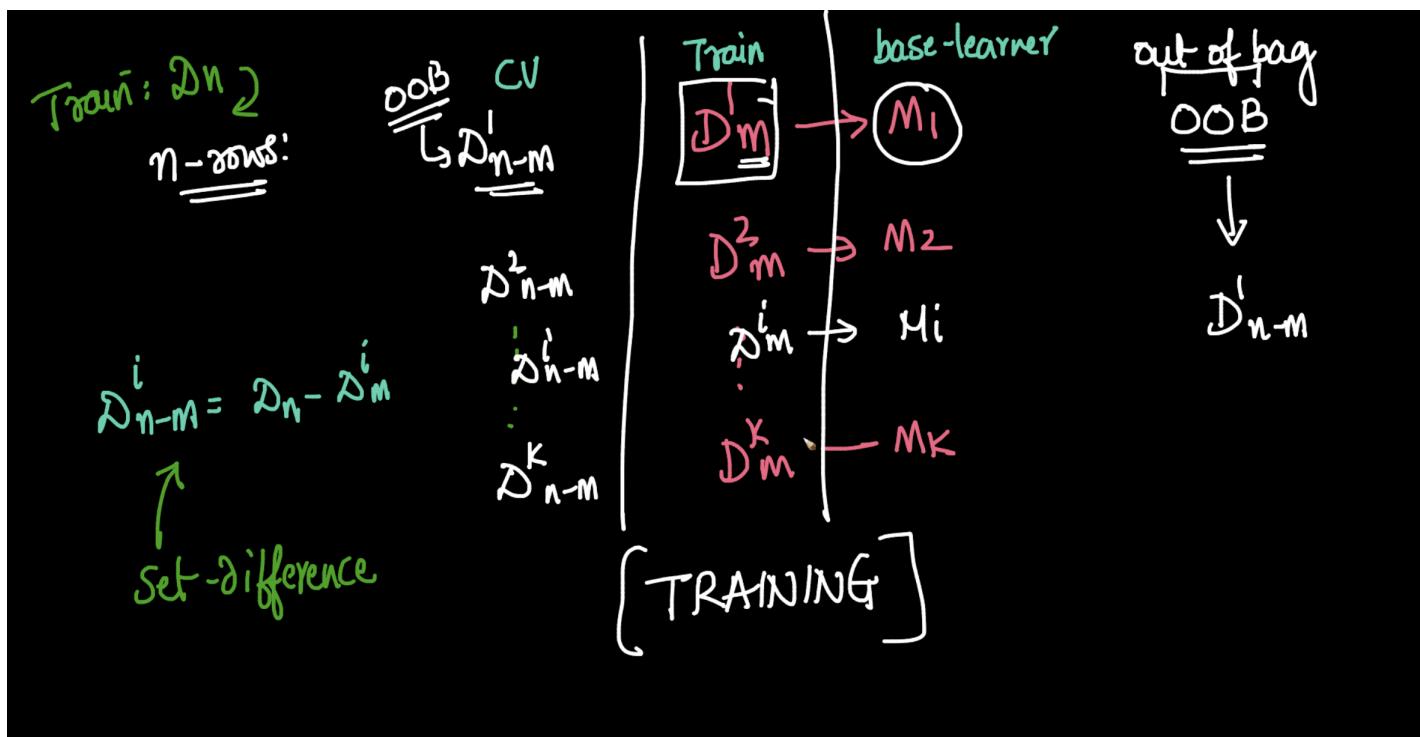
**Do you see a problem here?**

- We have  $k$  different train data sets ( $D_m^1, D_m^2, \dots, D_m^k$ ) to train  $k$  different models ( $M_1, M_2, \dots, M_k$ ) respectively

**Where is the cross validation data?**

The remaining  $n - m$  rows can be used as cross validation data

- for  $D_{n-m}^i = D_n - D_m^i$ , is Set Difference
  - This is referred to as Out Of Bag(OOB)
- So, we train the models on ( $D_m^1, D_m^2, \dots, D_m^k$ ) these data sets and cross validate them on ( $D_{n-m}^1, D_{n-m}^2, \dots, D_{n-m}^k$ ) this data which is OOB sample
- If some data points are not present in any of the train data sets, they will be in every cross validation data.



Will Random Forest always out perform Decision tree?

- Yes, in most cases of the cases they should ut perform

When will the performance of both Random Forest and Decision Tree be similar?

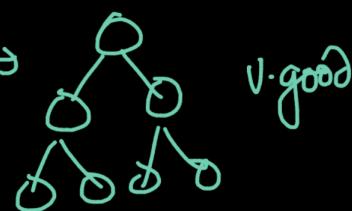
- This may happen when the data is very simple, that a simple Decision tree by using simple rules can get a very good performance.

Simpler model if it works:

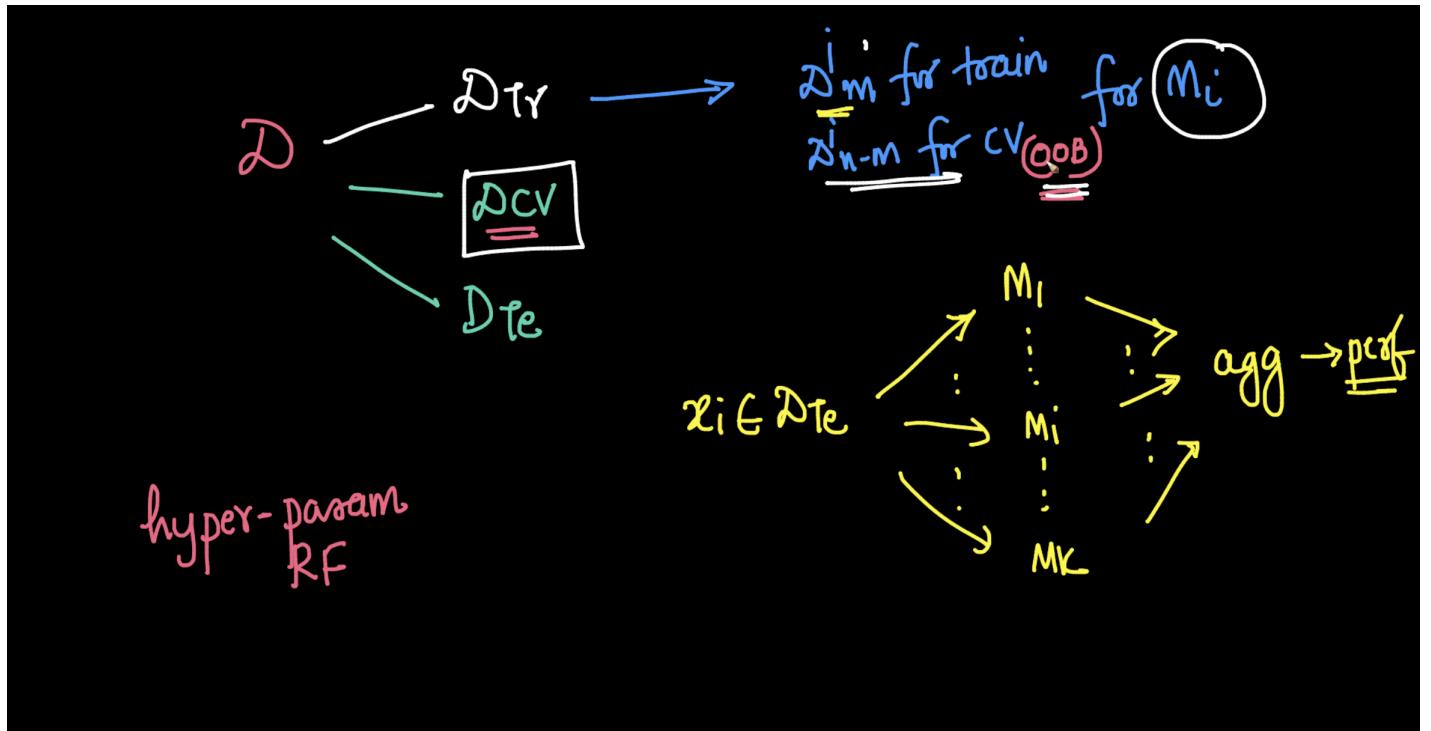
RF better than DT  
 (most-cases)

when so DT & RF have similar perf

Simple task:

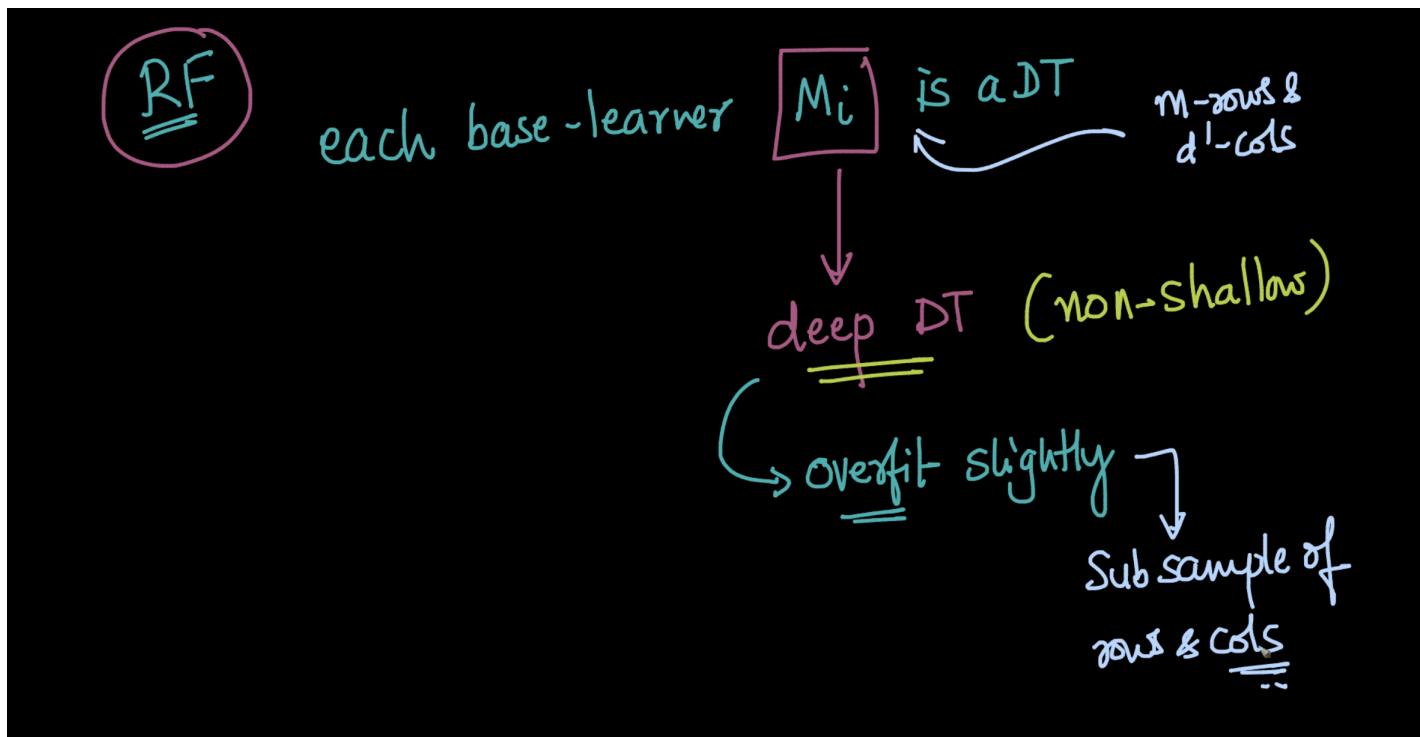


- In random forest The Base learners don't need separate cross validation data set because the remaining data after sampling(OOB) are used for cross validation
  - i.e from  $D_n$ ,  $D_m^i$  is sampled and  $D_{m-n}^i$  is used for cross validation
- But the random forest as a whole has the cross validation data and test data to tune the Hyper Parameters of random Forest

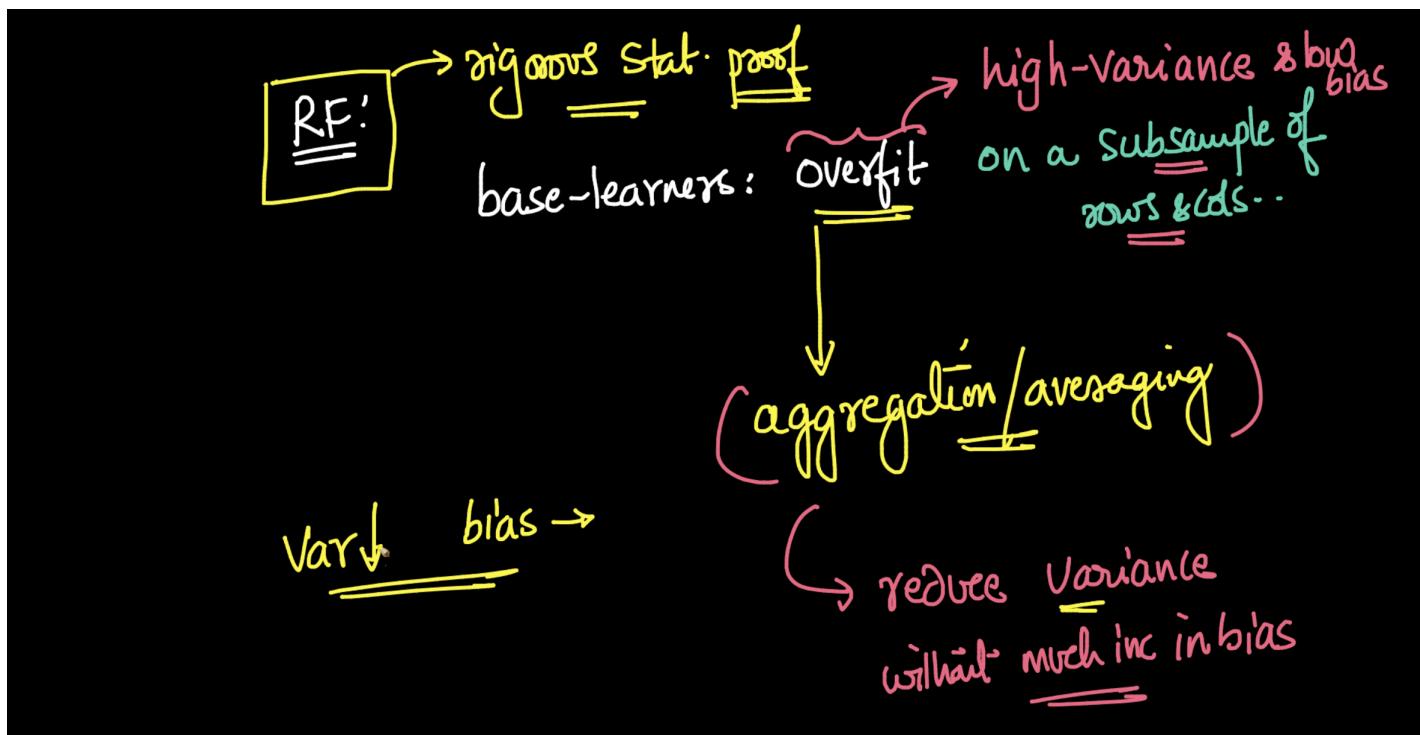


### What happens if $k$ increases?

- The Decision tree base learners in random forest are Deep Decision Trees (non-shallow).
- As they are deep, the models overfit slightly on a sub sample of data, as the base learners are made of  $m$  rows and  $d'$  columns only.
- Now, we perform aggregation\averaging on these slightly overfit models, on a sub-sample of data, which have high variance and low bias



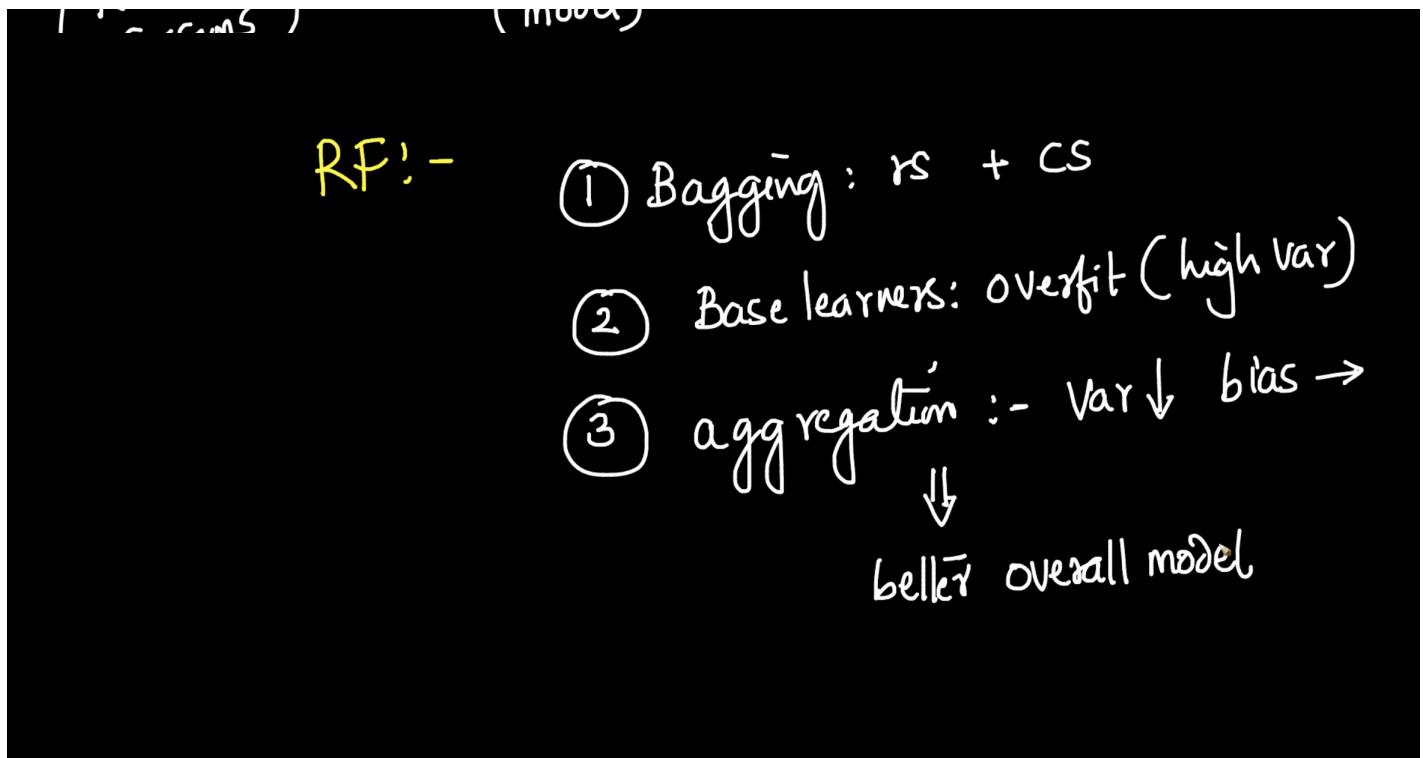
- But it is statistically proved that aggregating reduces the variance without much increase in the bias.
- In statistical Machine learning the error of a model can be represented as



- **Error = Bias<sup>2</sup> + Variance + Irreducible error.**
- Here we saw due to aggregation, the variance decreases without trading-off to bias, due to which the overall error of the model decreases.
  - Therefore the error of the random forest reduces.

Now, as we saw each step of random forest

- Let's see the core working of a random forest
  1. It uses Bagging, which is row sampling, and column sampling
  2. It takes Base learners which are slightly overfit i.e having high variance and low bias
  3. It reduces this variance by aggregation without much increase in its bias.



### Can we use Non Decision Tree models in Random forest?

- Yes we can you use slightly overfit models i.e having high variance models.
  - For example: We can use KNN with small  $k$
  - We can use Logistic Regression with higher order features + low value of  $\lambda$

## ▼ Hyper-Parameter tuning of Random Forest

Now lets see various hyper parameters of Random Forest

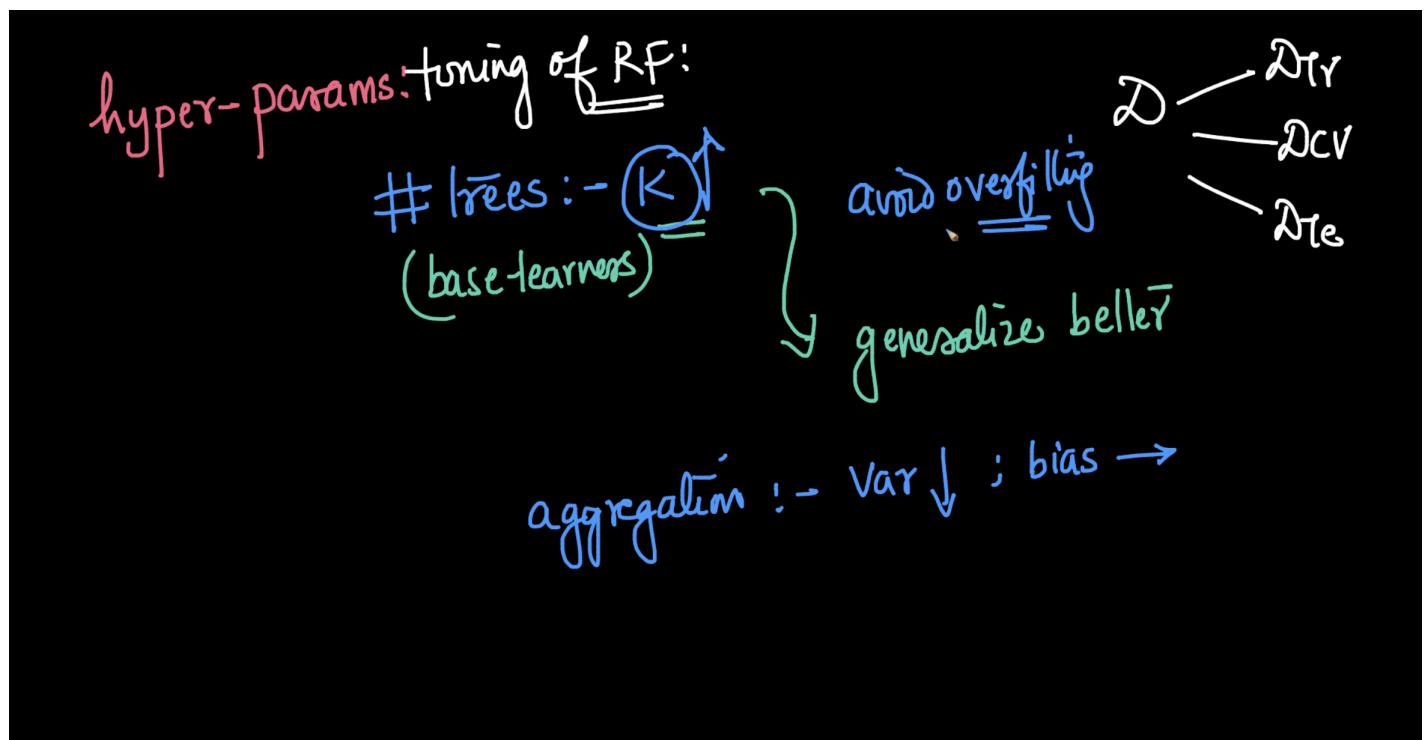
- Note that we consider the whole dataset  $D$  not only the  $D_{train}$  which is split into train data set for base learners

▼ What are the hyper parameters for random forest?

▼ 1. Number of Trees( $k$ )

What happens if number of base learners ( $k$ ) increases?

- As  $k$  increases, we Aggregate more number of base learners
- We already saw that, Aggregation decreases the variance without any significant increase in the bias.
- Hence, we can say that we avoid overfitting.



What if my base learners underfit instead being overfit?

- If the base learners underfit, the variance is already low
- Hence the aggregation which should reduce the variance doesn't work as expected.

▼ 2. Row sample size ( $m$ )

- In some libraries  $\frac{m}{n}$  is also considered where this value becomes maximum i.e 1 ,when  $m = n$

### What happens when $m = n$ ? ( $d = d'$ )

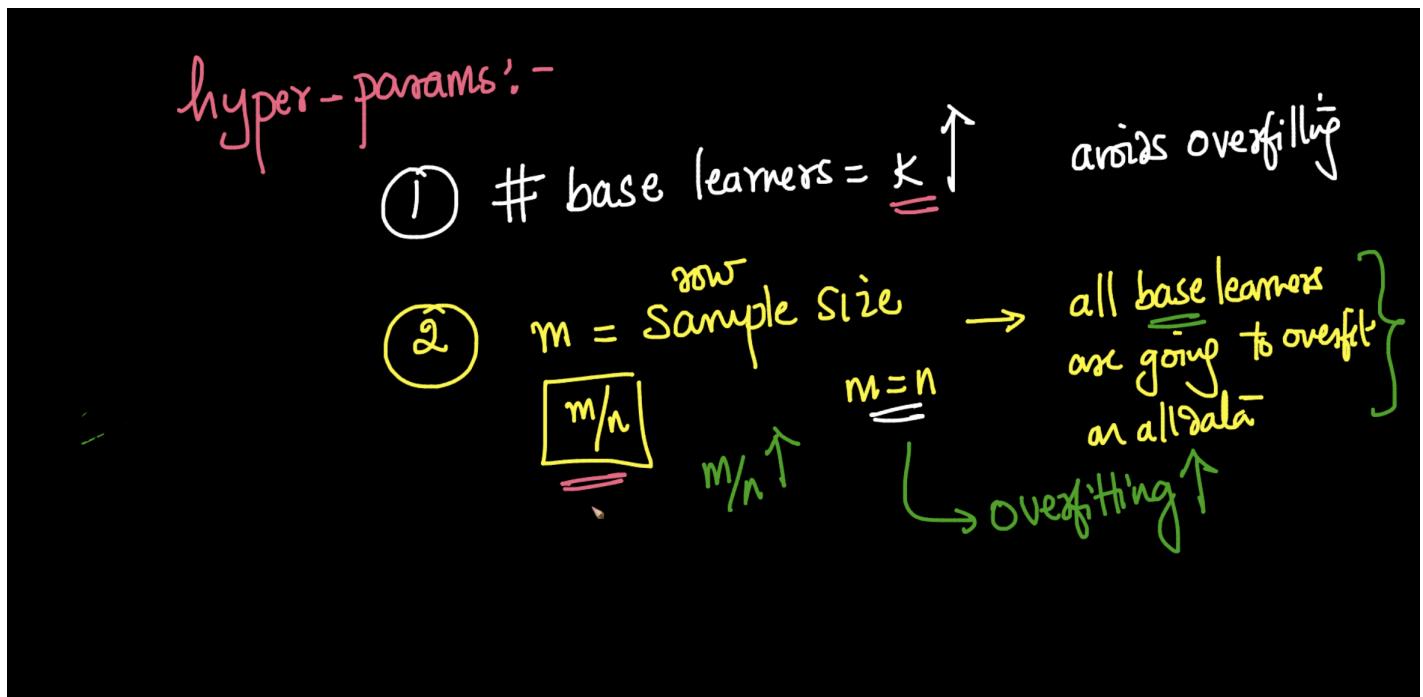
- All Base learners overfit on the whole data and they be similar.
- So, as the  $\frac{m}{n}$  increases the over-fitting of the model increases

### ▼ 3. Number of columns Sampled ( $d'$ )

- Consider  $m = n$
- As  $\frac{d'}{d}$  increases, then all the features are being used by each of the tree.
- This implies the overfitting increases i.e the variance increases

### ▼ 4. Depth of Base learners

- This is optional, some libraries let us use this when we dont want to use the Out Of Bag (OOB) data fro cross validation
- As depth increases, the overfitting increases
  - When the ratios  $\frac{m}{n}$  and  $\frac{d'}{d}$  are very less, that is when the base lerners are seeing very few columns and rows , and when the number of trees increases ( $k$ ) the overfitting chances are low.



(3)  $d' = \# \text{cols sampled}$

$$\frac{d'}{d} \uparrow \Rightarrow \text{overfitting} \uparrow$$

(4) depth of base learners  $\frac{\text{OOB Samples}}{cv}$

depth  $\uparrow$  overfit

depth  $\uparrow \rightarrow \frac{m}{n} \downarrow \frac{d'}{d} \downarrow \quad K \uparrow$  less chance of overfit

## ▼ Questions

What if the models have high bias in an ensemble?

- We know Random forest is an ensemble of high variance models, that is slightly overfitting models
- We see the Gradient Boosted Decision Tree which applies boosting, which deals with high bias or underfitting models

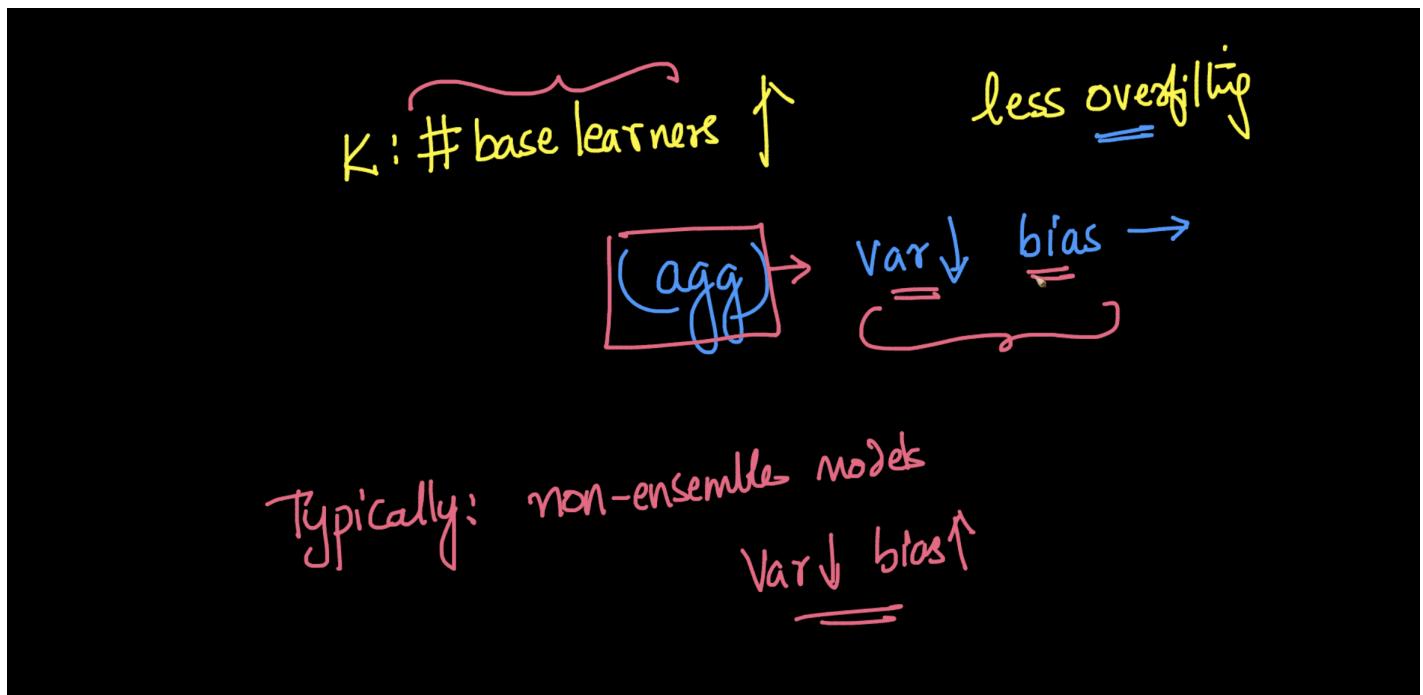
Bagging  $\rightarrow$  ensemble of high-var  $\text{agg} \circlearrowleft$ : Var  $\downarrow$  bias  $\Rightarrow$

RF: ensemble of high-var

GBT:  $\left. \begin{array}{l} \text{Boosting} \\ \text{...} \end{array} \right\} \rightarrow \text{ensemble of high bias models} \dots$

What if ensemble a deep decision tree and a simple linear classifier?

- This doesn't work as expected because the random forest expects highly variance models, as it does aggregation at the end which reduces the variance, but do not effect the bias. Due to which if a model of high bias is there, the high bias remains the same which we should avoid.



### What happens when $m$ is very close to $n$ ?

- When  $m$  is close to  $n$  the base learners overfit on the whole train data set. Due to which Aggregation doesn't work in reducing the variance
- So, when  $\frac{m}{n}$  increases, Random forest overfits.
- But when  $\frac{m}{n}$  decreases the base learners overfit, due to which we require more number of trees to reduce the variance.

$m \approx n \Rightarrow$  base learners are similar

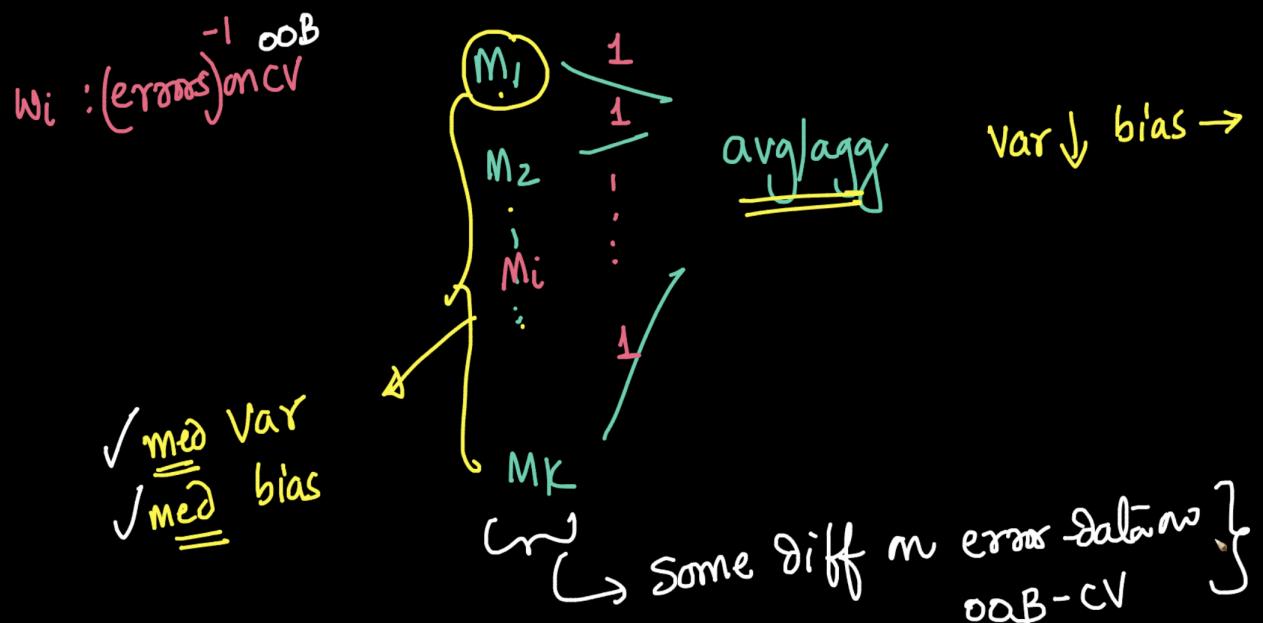
$\frac{m}{n} \uparrow$  RF will overfit

$\downarrow$  (agg) doesn't work well

Double-click (or enter) to edit

So here, Should we take weighted average of the base models in aggregation ?

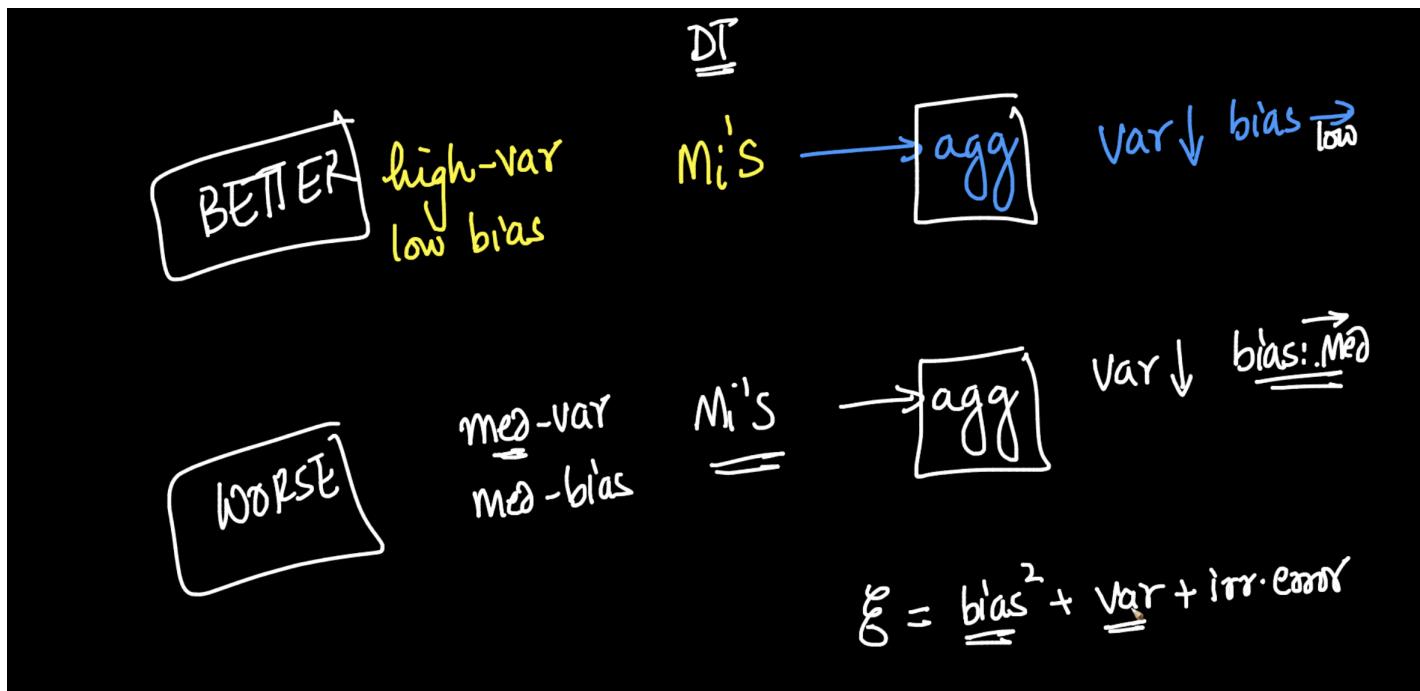
No, this is not required because we are taking Base learners with equal weight



- We saw that aggregation is the key in Random forest ,
- But we might have two different types of Base learners in the random forest

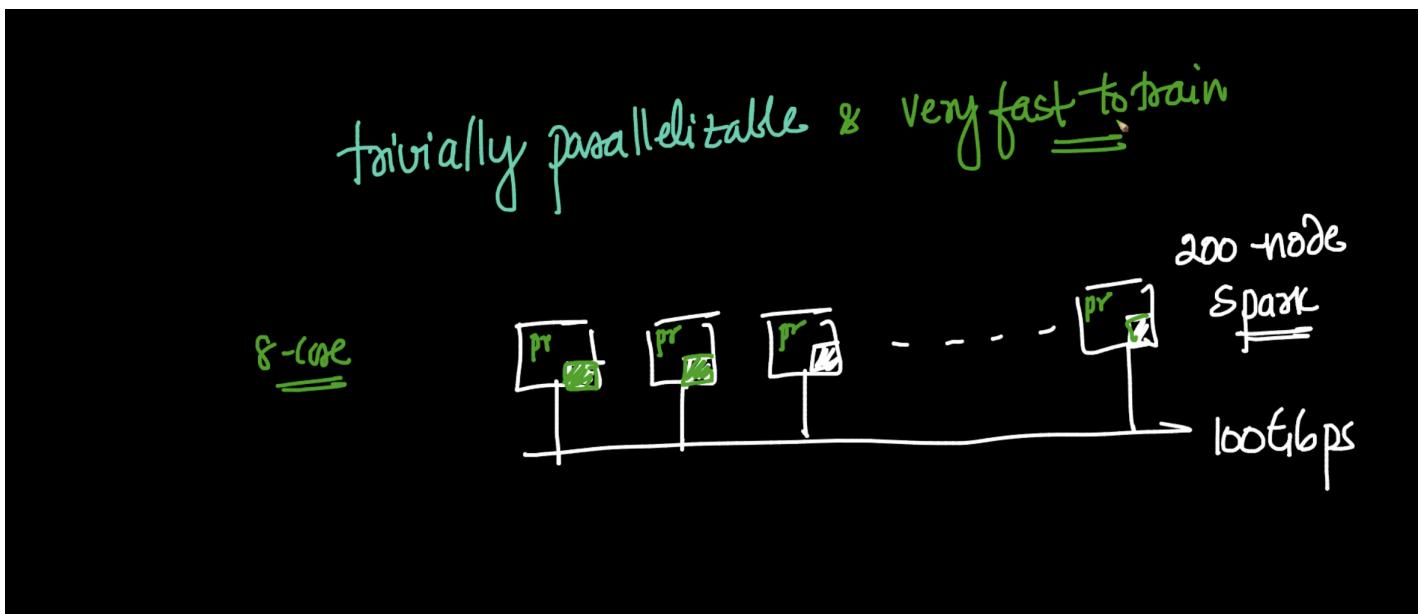
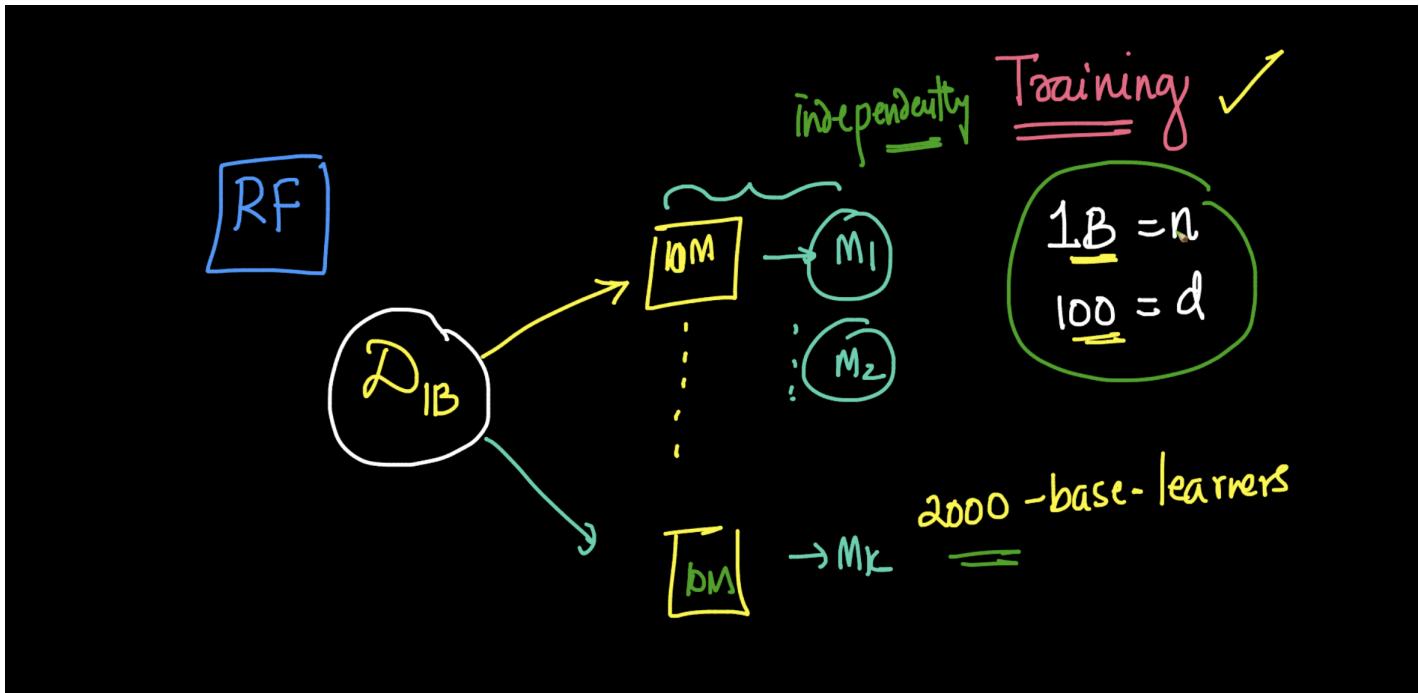
- With high variance and low bias or
- With medium variance and medium bias
- Here as we do aggregation both the models become
  - with low variance and low bias and
  - with low variance and medium bias respectively

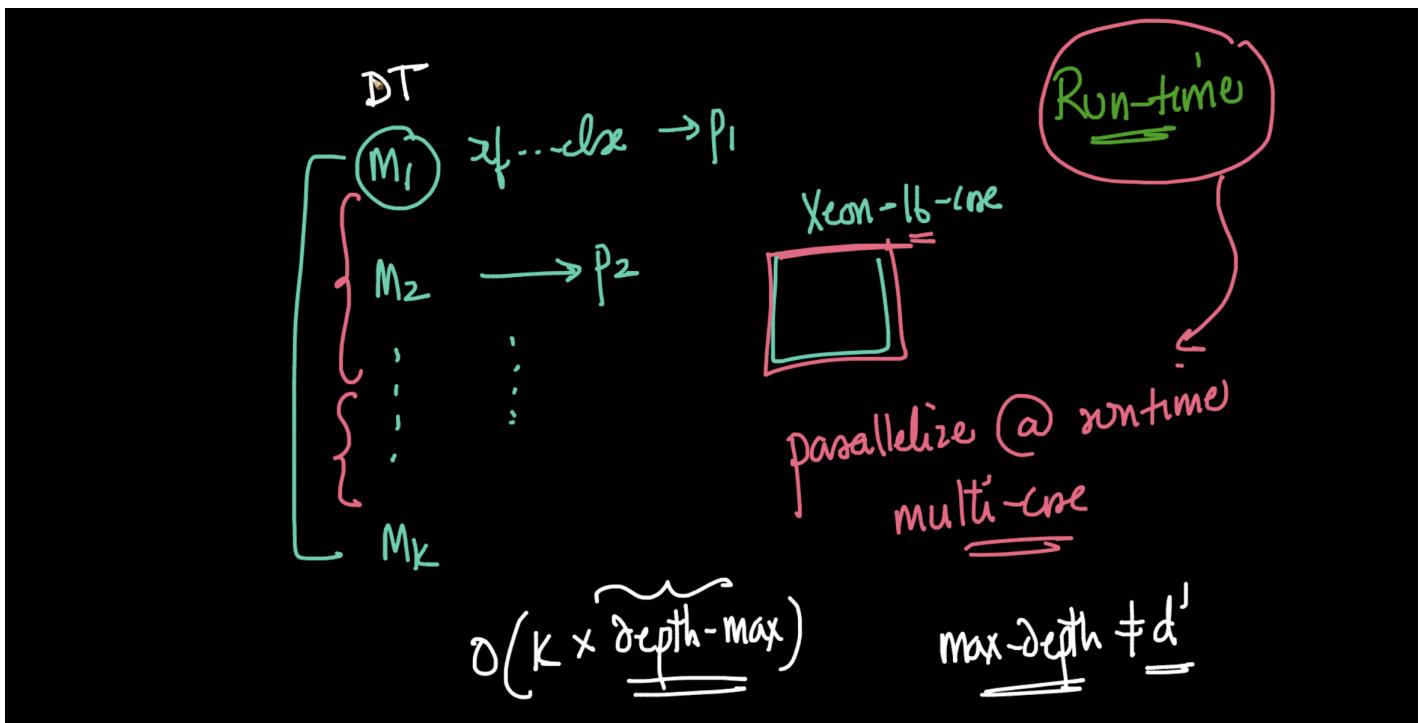
So it's better to have high variance and low bias base learners.



## ▼ Training a random forest

- With present day improvements we can easily run a Random forest with a billion data with 10 million sample size and with thousand base learners
- These base learners can be trivially parallelised.
- As each model is **trained independently** we can even take these models to distributed computing
- In a distributed computing system each processor is given with different model's data set , the system can work on multiple models parallelly, if the processor is a 8 core processor the datasets of 8 different trees can be worked on a system.
- So, we are parallelising the Random forest on multiple cores due to which the process becomes fast.
- The time complexity here is  $O(k * \text{max\_depth of tree})$

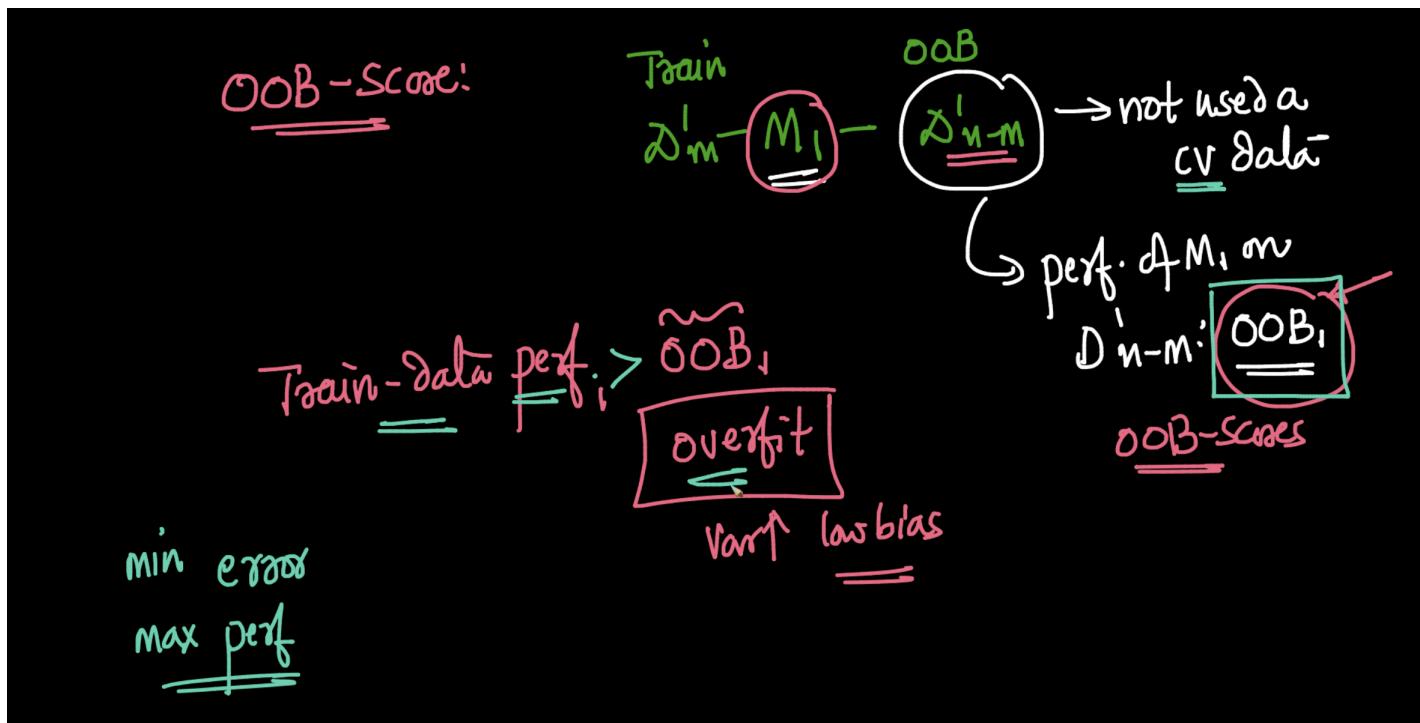




## ▼ OOB Score

As we already studied that after sampling the data for each model we have some data remaining which is the **Out Of Bag** data

- This OOB data ( $D'_{n-m}$ ) is not used as cross validation data, it is used to find the performance of the model on data, this is called the **OOB score** of the data
- When the trained data performance is more than the OOB data performance, we can conclude that the model is over-fitting, which we want upto some extend
- In the same way we can say that the model is overfitting if the trained data error is less than OOB data error



## ▼ Sklearn library (terms and usage):

### ▼ OOB score

- By default the OOB score is set to false, we can keep it true to get the OOB scores of each model
- Which can be compared by the performance of the trained data know the extent of overfitting in the base learner model



scikit-learn 1.0.2

Please cite us if you use the software.

learn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

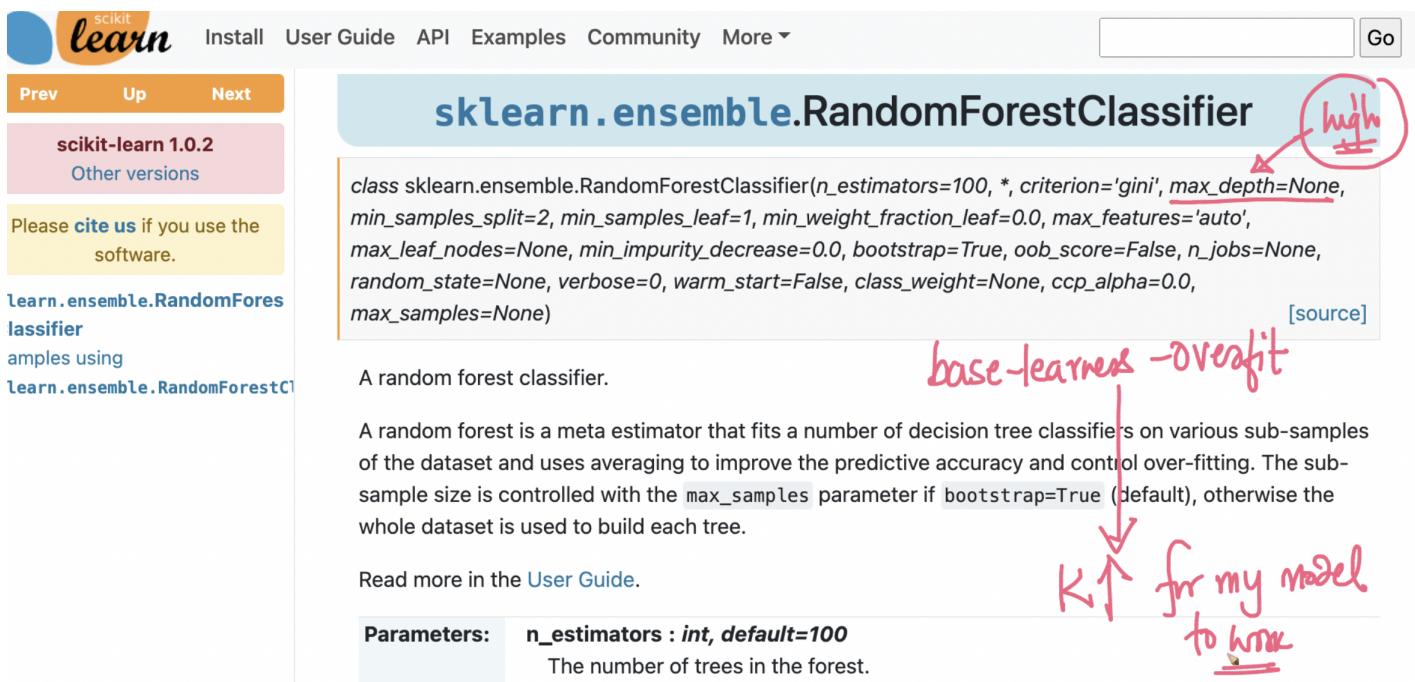
A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

**Parameters:** `n_estimators : int, default=100`  
The number of trees in the forest.

## ▼ max depth

- This is used to decide the depth of the each Decision tree in the model
- As we already discussed, as depth increases the model overfits more due to which we may need more base learners for model to work



scikit-learn 1.0.2

Please cite us if you use the software.

learn.ensemble.RandomForestClassifier

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the `max_samples` parameter if `bootstrap=True` (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

**Parameters:** `n_estimators : int, default=100`  
The number of trees in the forest.

## ▼ min\_samples\_split

- This is used to set the minimum number of data points required in a node to split further

## What happens if we increase the min\_sample\_split?

- If we increase the min\_sample\_split value the splitting stops if there are less number of datapoints than the set value, which controls the overfitting and may also result in underfitting and shallow trees.

**scikit-learn** Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2 Other versions

Please cite us if you use the software.

learn.ensemble.RandomForestClassifier

amples using learn.ensemble.RandomForestC

*base learners underfit*

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the max\_samples parameter if bootstrap=True (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

**Parameters:** **n\_estimators : int, default=100**  
The number of trees in the forest.

## ▼ min\_samples\_leaf

- This gives the minimum number of data points required in the leaf node.
- If the value of min\_samples\_leaf is increased the base learners become shallower

**scikit-learn** Install User Guide API Examples Community More ▾

Prev Up Next

scikit-learn 1.0.2 Other versions

Please cite us if you use the software.

learn.ensemble.RandomForestClassifier

amples using learn.ensemble.RandomForestC

*Shallow base-learners*

```
class sklearn.ensemble.RandomForestClassifier(n_estimators=100, *, criterion='gini', max_depth=None, min_samples_split=2, min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0, bootstrap=True, oob_score=False, n_jobs=None, random_state=None, verbose=0, warm_start=False, class_weight=None, ccp_alpha=0.0, max_samples=None)
```

[source]

A random forest classifier.

A random forest is a meta estimator that fits a number of decision tree classifiers on various sub-samples of the dataset and uses averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is controlled with the max\_samples parameter if bootstrap=True (default), otherwise the whole dataset is used to build each tree.

Read more in the [User Guide](#).

**Parameters:** **n\_estimators : int, default=100**  
The number of trees in the forest.

## ▼ n\_jobs

- This gives the number of jobs that can be run in parallel.
- This does the multi core parallelisation
- If the value of n\_jobs is set to 4, four cores of your processor will be used
- We can use all the cores of the processor by setting this value to -1.



### **n\_jobs : int, default=None**

The number of jobs to run in parallel. `fit`, `predict`, `decision_path` and `apply` are all parallelized over the trees. `None` means 1 unless in a `joblib.parallel_backend` context. `-1` means using all processors. See [Glossary](#) for more details.



### **random\_state : int, RandomState instance or None, default=None**

Controls both the randomness of the bootstrapping of the samples used when building trees (if `bootstrap=True`) and the sampling of the features to consider when looking for the best split at each node (if `max_features < n_features`). See [Glossary](#) for details.

### **verbose : int, default=0**

Controls the verbosity when fitting and predicting.

### **warm\_start : bool, default=False**

When set to `True`, reuse the solution of the previous call to fit and add more estimators to the ensemble, otherwise, just fit a whole new forest. See [the Glossary](#).

## ▼ max\_samples

- This is used to set the number of samples to draw from the dataset to train each of the base learners
- This is nothing but the row sampling which we studied earlier and represented as  $m$
- We can give both int and float values, if we give a int value it takes as the number of samples to be considered and float is for percentage of the rows to be considered



**100 D.L**

### max\_samples : int or float, default=None

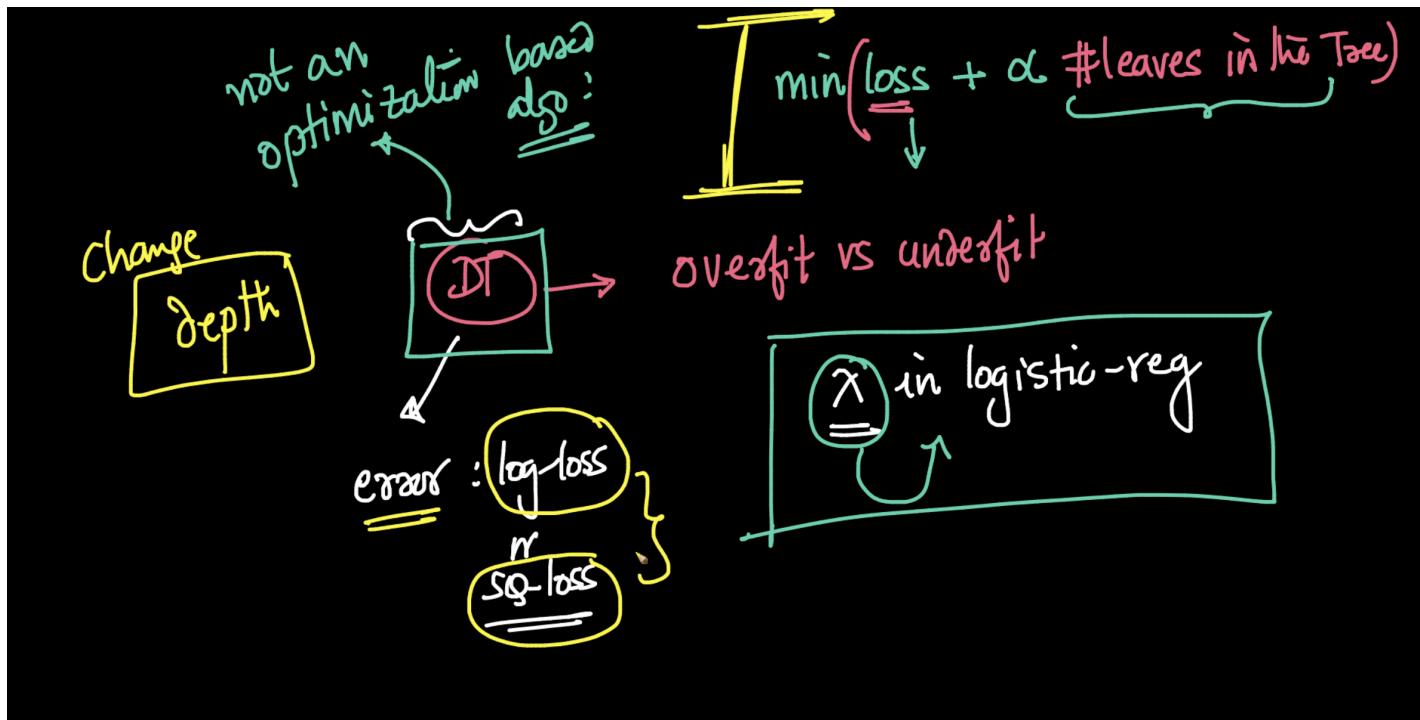
If bootstrap is True, the number of samples to draw from X to train each base estimator.

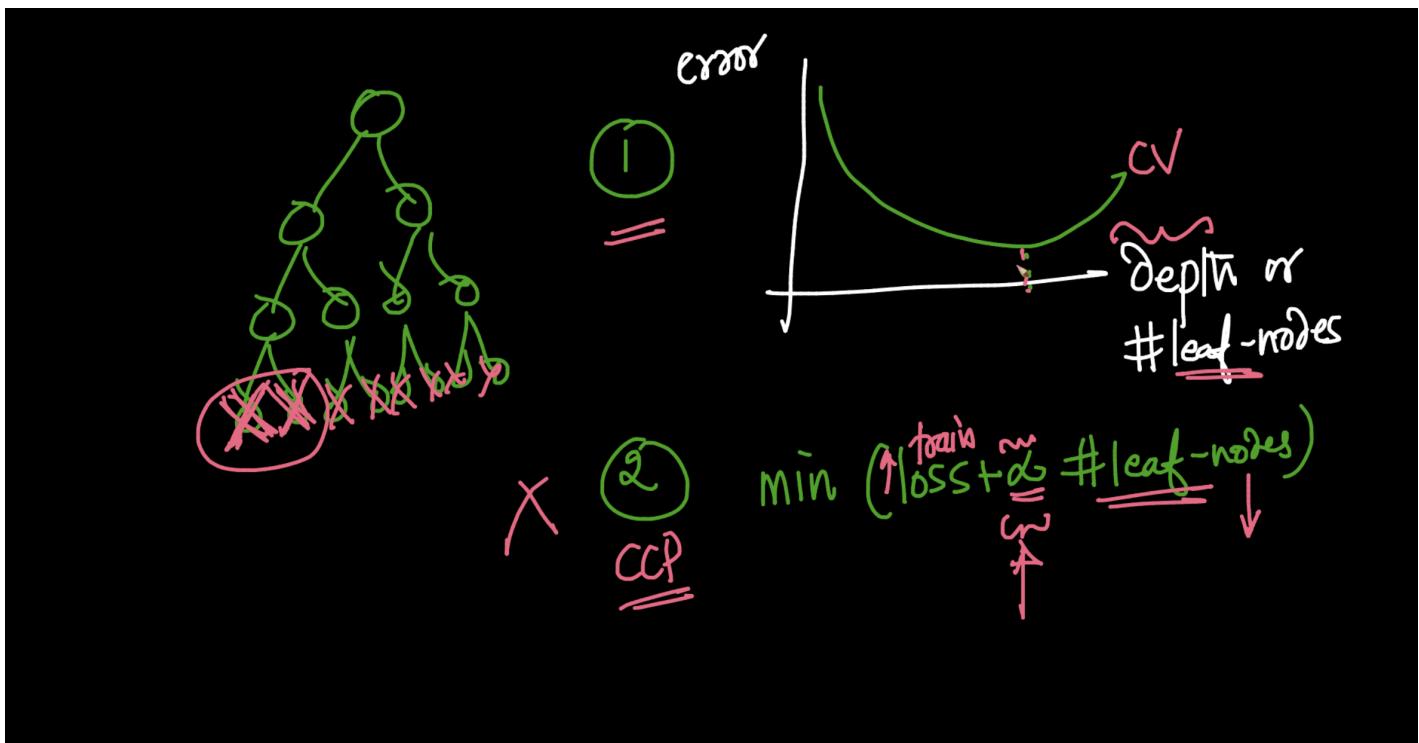
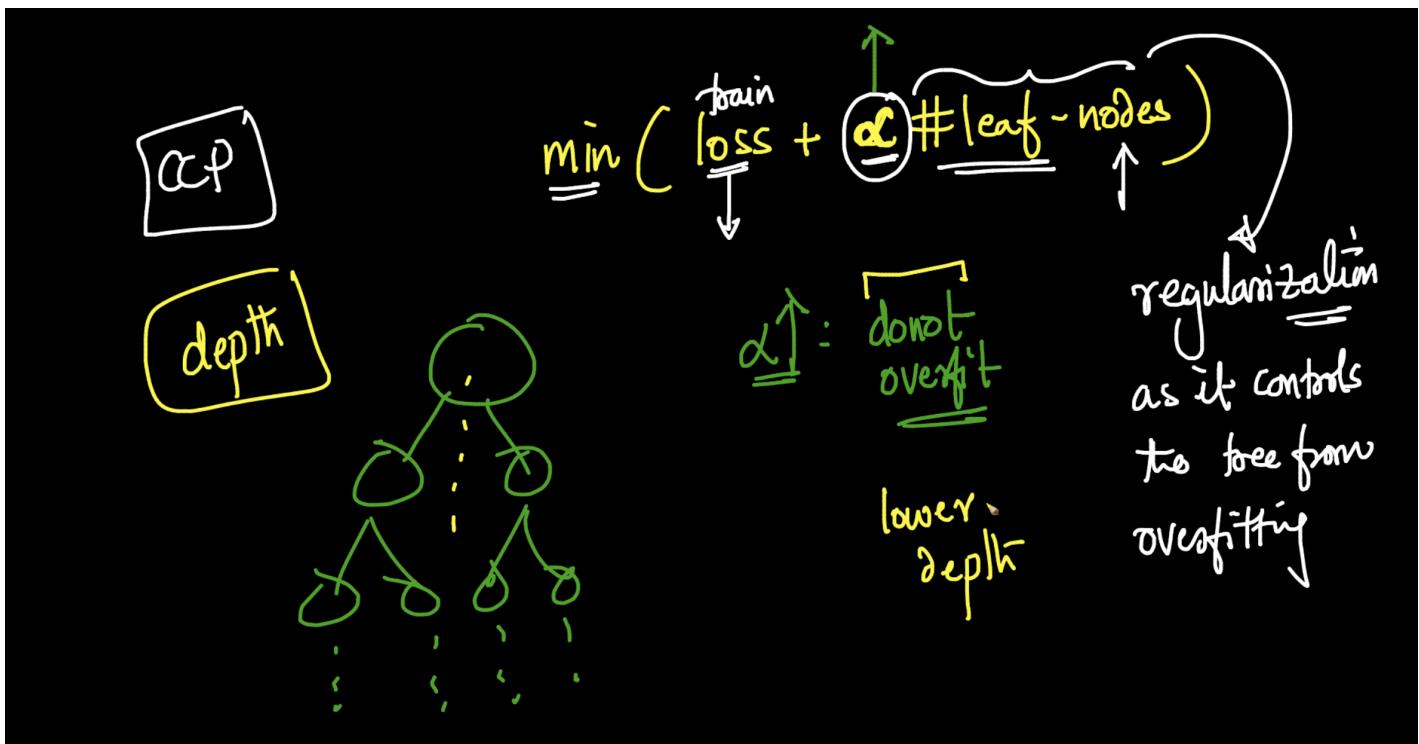
- If None (default), then draw `X.shape[0]` samples.
- If int, then draw `max_samples` samples.
- If float, then draw `max_samples * X.shape[0]` samples. Thus, `max_samples` should be in the interval  $(0.0, 1.0]$ .

New in version 0.22.

## ▼ ccp\_alpha

- CCP is cost complexity pruning
- This is basically used for pruning the base learners
- We can control the overfitting and underfitting of the base learners using the value  $\alpha$ , this is almost similar to  $\lambda$  which we used in linear and logistic regression
- so the idea here is to minimise the loss associated with the decision tree and  $\alpha$  times the number of terminal nodes (leaves) which controls overfitting
  - $\min(\text{loss} + \alpha * \text{number of leaves in the tree})$
- As the depth of tree increases we know the loss decreases, where the number of leaf nodes increases, this trade-off between the loss and number of leaves can be controlled using  $\alpha$  like regularisation.





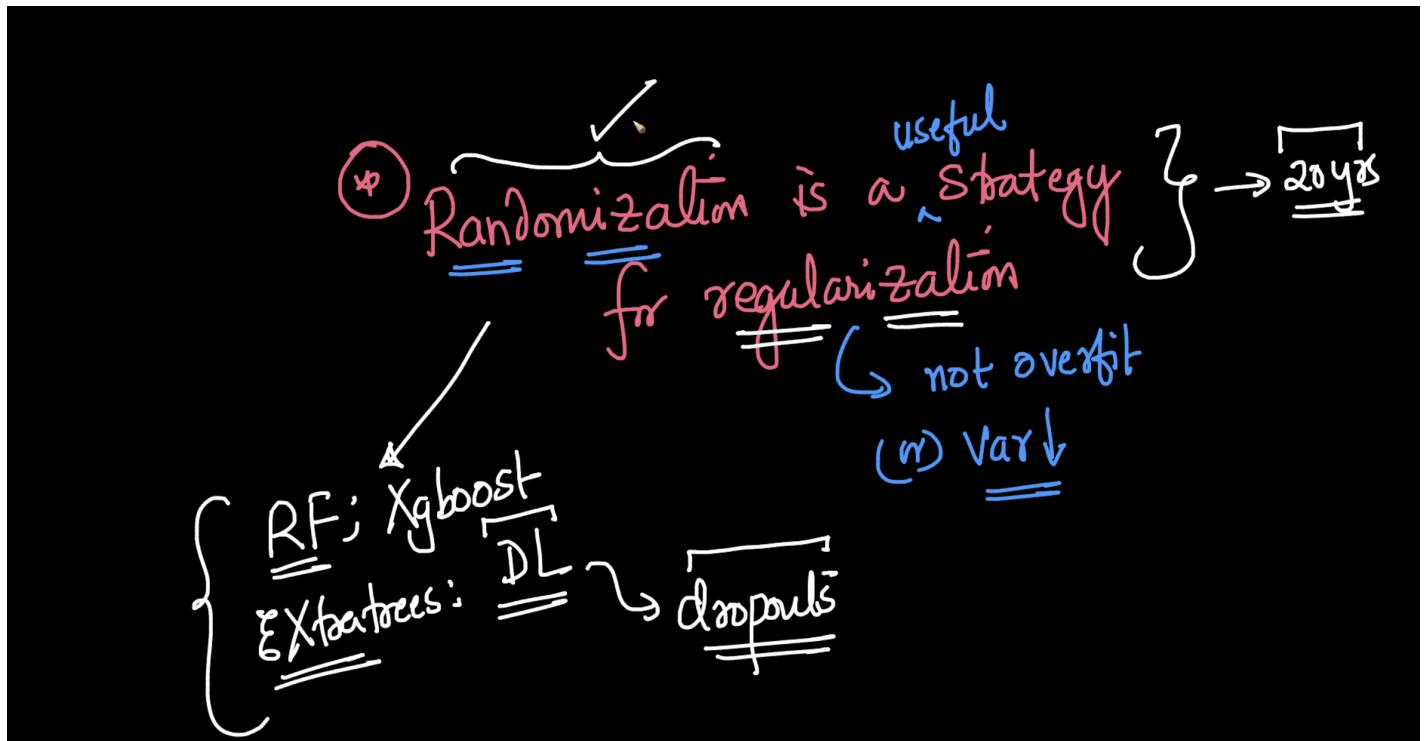
## ▼ Extra Trees or Extremely Randomised Trees

let's compare these with Random Forest

What do we do in random forest?

- We do random row sampling and column sampling and then we do aggregation, in which this randomisation and aggregation play a key role in reducing the variance keeping the bias similar, further avoiding overfitting

**Randomisation is a great, useful and very powerful strategy for Regularization**



### what do extra trees do?

In extra tree does random row and column sampling and aggregation just like decision trees but it also randomly picks the threshold ( $\tau$ ) to split for numerical features.

- In decision trees we saw using the feature values as threshold we calculate information gain and then choose the threshold of the split, basing on the values of information gain
- In an extra tree we choose this threshold ( $\tau$ ) also randomly, adding one more randomisation on top of random forest,
  - say there are  $m$  rows in a feature we randomly select few rows and set the threshold basing on their IG values
- This is useful when the  $m$  is large, but these require more base learners as the trees are not perfect
- Hence these are not widely used

$\varphi_1 - \text{loc} \rightarrow \text{Col. Sampling}$

